




## Article

# My-Trac: System for Recommendation of Points of Interest on the Basis of Twitter Profiles

Alberto Rivas <sup>1,2,\*</sup> , Alfonso González-Briones <sup>1,2,3</sup> , Juan J. Cea-Morán <sup>1</sup>, Arnau Prat-Pérez <sup>4</sup> and Juan M. Corchado <sup>1,2</sup> 

<sup>1</sup> BISITE Research Group, University of Salamanca, Edificio I+D+i, Calle Espejo 2, 37007 Salamanca, Spain; alfonso@usal.es (A.G.-B.); juanju\_97@usal.es (J.J.C.-M.); corchado@usal.es (J.M.C.)

<sup>2</sup> Air Institute, IoT Digital Innovation Hub, Carbajosa de la Sagrada, 37188 Salamanca, Spain

<sup>3</sup> Research Group on Agent-Based, Social and Interdisciplinary Applications (GRASIA), Complutense University of Madrid, 28040 Madrid, Spain

<sup>4</sup> Sparsity-Technologies, 08034 Barcelona, Spain; arnau@sparsity-technologies.com

\* Correspondence: rivis@usal.es

**Abstract:** New mapping and location applications focus on offering improved usability and services based on multi-modal door to door passenger experiences. This helps citizens develop greater confidence in and adherence to multi-modal transport services. These applications adapt to the needs of the user during their journey through the data, statistics and trends extracted from their previous uses of the application. The My-Trac application is dedicated to the research and development of these user-centered services to improve the multi-modal experience using various techniques. Among these techniques are preference extraction systems, which extract user information from social networks, such as Twitter. In this article, we present a system that allows to develop a profile of the preferences of each user, on the basis of the tweets published on their Twitter account. The system extracts the tweets from the profile and analyzes them using the proposed algorithms and returns the result in a document containing the categories and the degree of affinity that the user has with each category. In this way, the My-Trac application includes a recommender system where the user receives preference-based suggestions about activities or services on the route to be taken.

**Keywords:** users' profiling; data extraction; natural language processing; recommender system; mapping application

**Citation:** Rivas, A.; González-Briones, A.; Cea-Morán, J.J.; Prat-Pérez, A.; Corchado, J.M. My-Trac: System for Recommendation of Points of Interest on the Basis of Twitter Profiles.

*Electronics* **2021**, *10*, 1263. <https://doi.org/10.3390/electronics10111263>

Academic Editors: Dimitris Apostolou and Osvaldo Gervasi

Received: 13 April 2021

Accepted: 20 May 2021

Published: 25 May 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Humans are social beings; we always seek to be in contact with other people and to have as much information as possible about the world around us. The philosopher Aristotle (384–322 B.C.) in his phrase “Man is a social being by nature” states that human beings are born with the social characteristic and develop it throughout their lives, as they need others in order to survive. Socialization is a learning process; the ability to socialize means we are capable of relating with other members of the society with autonomy, self-realization and self-regulation. For example, the incorporation of rules associated with behavior, language, and culture improves our communication skills and the ability to establish relationships within a community.

In the search for improvement, communication, and relationships, human beings seek to get in contact with other people and to obtain as much information as possible about the environment in order to achieve the above objectives. The emergence of the Internet has made it possible to define new forms of communication between people. It has also made it possible to make a large amount of information on any subject available to the average user at any time. This is materialized in the development of social networks. The concept of social networking emerged in the 2000s as a place that allows for interconnection between people, and, very soon, the first social networking platforms appeared on the Internet that

served to bring people together. Among the first platforms that emerged were Fotolog, MySpace, Hi5, Buzz, and SecondLife; however, most of them have declined in popularity or disappeared. Today, Facebook, Twitter, and Instagram stand out; they are used by millions of people all over the world. Thanks to these technologies, people from different parts of the world can engage in conversation, post photos of their latest trip, or keep their followers updated by sharing their opinions or experiences.

With regard to writing opinions or experiences, Twitter is the social network par excellence. Twitter is based on the concept of microblogging, i.e., users can post messages about their opinions, preferences, experiences, etc., with a maximum of 280 characters. Twitter allows its users to follow other accounts that interest them, or to comment on events in real time using hashtags. All this translates into one word: information. The information that users provide on social networks can be used in a variety of ways, many of them negative. Exposure on the Internet means that anyone can access the users' data and use it for financial gain. However, it can also be used to make life easier for users who choose to do so, always bearing in mind that there must be express consent on their part. This is precisely the case of the work presented here.

Since the emergence of the first social networks much progress has been made towards the current state of maturity of the social network life cycle. As presented above, they are of vital importance to society, as they fulfill the innate communication function of human beings. In addition to their use as a means of communication, they have begun to be exploited for business purposes in order to profit from the enormous amount of information that is generated on a daily basis. This information, which is generated by the society, is of great value once analyzed and processed correctly.

The data generated by users on social networks allows for the development of commercial and advertising actions that are much more effective than with traditional formats. Advertisement platforms have developed the ability to segment advertising on the basis of the behavior of each user, to show products and services to those who are really interested in those products and services. This increases the effectiveness of advertisements.

Social network users publish practically everything that happens in their daily lives, their opinions, where they are, what they eat, what they would like to buy, where they are going on holiday, and a long list of behaviors that are transformed into valuable information for analysis. The information obtained through the analysis is very interesting as it allows for the elaboration of demographic, socio-economic, and consumer trend profiles. The companies that own these social networks sell high-value information to other companies to enable them to carry out much more powerful and effective marketing and advertising strategies. Another remarkable aspect of data analytics on social networks is the ability to perform real-time analysis of the information to offer products and services according to their characteristics.

Information is a very precious commodity, and, as presented above, Twitter is a great source of data when analyzing human behavior and interactions or when learning about the opinion of certain users on certain topics. This information can be used to improve the multi-modal experience of users when they use the My-Trac application. Therefore, an adaptation of these systems for adoption in mapping applications is proposed.

The European My-TRAC project focuses on providing user-centered services to improve the multi-modal experience of passengers from door-to-door. This helps citizens develop greater confidence in and adherence to multi-modal transport services. In addition, My-TRAC improves customization to users' needs through data, statistics, and trends provided by passengers' experiences when using the proposed platform. Part of the tailoring of services and recommendations to users is determined by the knowledge obtained from their Twitter posts through the use of NLP techniques to classify and understand users.

There are other services that offer similar functionalities to My-Trac, such as ROSE [1], CTRR, and CTRR+ based systems for city-based tourism [2], or to participate in solidarity projects in rural environments [3]. From among the above, ROSE (ROuting Service) stands out, which is a mobile phone application that suggests events and places to the user and

guides them via public transport. There are many different systems that incorporate both recommendation and navigation. However, there is no system that combines event recommendation and pedestrian navigation with (real-time) public transport. However, it does not employ multi-modal navigation between different public transport modes (bus, train, carpooling, plane, etc.) in different countries and that would use information from the user's social network profile. Instead, current systems utilize a set of information initially entered into the application which is not updated afterwards. Finally, Tables 1 and 2 present a review of similar works.

**Table 1.** Review of similar works: Part I.

Title / Publication	Functionality	Advantages	Shortcomings
ROSE (ROuting SErvice) [1]	Mobile phone application that suggests events and places to the user and guides them via public transport.	The current systems utilize a set of information initially entered into the application which is not updated afterwards.	There is no system that combines event recommendation and pedestrian navigation with (real-time) public transport. It does not employ multi-modal navigation between different public transport modes (bus, train, carpooling, plane, etc.) in different countries and that would use information from the user's social network profile.
Systems for city-based tourism [2]	A personalized travel route recommendation based on the road networks and users' travel preferences.	The experimental results show that the proposed methods achieve better results for travel route recommendations compared with the shortest distance path method.	It does not use information from public transport services in route recommendations.
Tourism routes as a tool for the economic development of rural areas—vibrant hope or impossible dream? [3]	This paper argues that the clustering of activities and attractions, and the development of rural tourism routes, stimulates co-operation and partnerships between local areas. The paper further discusses the development of rural tourism routes in South Africa and highlights the factors critical to its success.	The article analyzes the realization of routes that include activities and attractions in a way that encourages and enhances rural development in Africa.	Preliminary project that requires public cooperation (institutions, transport, services) for a comprehensive improvement of the proposal.

This article improves on the previous system for the extraction of information regarding Twitter users [4]. The system is capable of obtaining information about a particular user and of elaborating a profile with the user's preferences in a series of pre-established categories. A review of existing reputation systems is presented in Section 2. Section 3 describes the proposal. Section 4 presents the assessment made with synthetic data. Section 5 shows how the system is integrated in My-Trac app. Finally, Section 6 presents the conclusions.

Table 2. Review of similar works: Part II.

Title / Publication	Functionality	Advantages	Shortcomings
Social Recommendations for Events [5]	Outlife recommender assists in finding the ideal event by providing recommendations based on the user's personal preferences.	In addition to the user's preferences, the recommender uses information from the user's group of friends to make event recommendations more satisfactory.	Although it uses information from the user's groups of friends, no use is made of information from the user's social networks to complement the analysis and recommendation.
Smart Discovery of Cultural and Natural Tourist Routes [6]	This paper presents a system designed to utilize innovative spatial interconnection technologies for sites and events of environmental, cultural and tourist interests. The system discover and consolidate semantic information from multiple sources, providing the end-user the ability to organize and implement integrated and enhanced tours.	The system adapts the services offered to meet the needs of specific individuals, or groups of users who share similar characteristics, such as visual, acoustic, or motor disabilities. Personalization is done in a dynamic way that takes place at the time and place of the service.	The very comprehensive system that uses external services, scraping, crawling, geo-positioning but does not include information from social networks to complement the analysis and recommendation of events.
Enhancing cultural recommendations through social and linked open data [7]	Hybrid recommender system (RS) in the artistic and cultural heritage area, which takes into account the activities on social media performed by the target user and her friends	The system integrates collaborative filtering and community-based algorithms with semantic technologies to exploit linked open data sources in the recommendation process. Furthermore, the proposed recommender provides the active user with personalized and context-aware itineraries among cultural points of interest.	The main drawback is the absence of extensive control over the semantics that are not taken into account. It generates difficulties in justifying, explaining, and hence analyzing the resulting scores.
Personalized Tourist Route Generation [8]	Intelligent routing system able to generate and customize personalized tourist routes in real-time and taking into account public transportation.	We have modeled the tourist planning problem, integrating public transportation, as the Time Dependent Team Orienteering Problem with Time Windows (TDTOPTW). We have designed a heuristic able to solve it in real time, precalculating the average travel times between each pair of POIs in a preprocessing step.	Future works consists on extending the system to more cities with a different public transport network topology. The next one consists on integrating an advanced recommendation system in a wholly functional PET. The systema don't use social network capabilities, that allows to store, share and add travel experiences to better help tourists on the destination.

## 2. Natural Language Processing Techniques Applied to Twitter Profiles

In this section, we review the main techniques applied in the analysis that make it possible to get to know the users preferences through their tweets. This allows for recommendations to be made according to the user profile.

### 2.1. Word Embedding Techniques

NLP techniques allow computers to analyze human language, interpret it, and derive its meaning so that it can be used in practical ways. These techniques allow for tasks, such as automatic text summarization, language translation, relation extraction, sentiment

analysis, speech recognition, and item classification, to be carried out. Currently, NLP is considered to be one of the great challenges of artificial intelligence as it is one of the fields with the highest development activity since it presents tasks of great complexity: how to really understand the meaning of a text, how to intuit neologisms, ironies, jokes, or poetry? It is a challenge to apply the techniques and algorithms that allow us to obtain the expected results.

One of the most commonly used NLP techniques is Topic Modeling. This technique is a type of statistical modeling that is used to discover the abstract “topics” that appear in a series of input texts. Topic modeling is a very useful text mining tool for discovering hidden semantic structures in texts. Generally, the text of a document deals with a particular topic, and the words related to that topic are likely to appear more frequently in the document than those that are unrelated to the text. Topic Modeling collects the set of more frequent words in a mathematical framework, which allows one to examine a set of text documents and discover, on the basis of the statistics of the words in each one, what the topics may be and what the balance is between the topics in each document.

The input of topic modeling is a document-term matrix. The order of words does not matter. In a document-term matrix, each row is a question (or document), each column is a term (or word), we label “0” if that document does not contain that term, “1” if that document contains that term once, “2” if that document contains that term twice, and so on.

Algorithms, such as Bag-of-words or TF-IDF, among others, make it possible to represent the words used by the models and create the matrix defined above, representing a token in each column and counting the number of times that token appears in each sentence (represented in each row).

- **Bag-of-words.** This model allows to extract the characteristics of texts (also images, audios, etc.). It is, therefore, a feature extraction model. The model consists of two parts: a representation of all the words in the text and a vector representing the number of occurrences of each word throughout the text. That is why it is called Bag-of-words. This model completely ignores the structure of the text, it simply counts the number of times words appear in it. It has been implemented through the Genism library [9].
- **Term Frequency - Inverse Document Frequency (TF-IDF).** This is the product of two measures that indicate, numerically, the degree of relevance that a word has in a document within a collection of documents [10]. It is broken down into two parts:
  - *Term frequency:* Measures the frequency with which certain terms appear in a document. There are several measurement options, the simplest being the gross frequency, i.e., the number of times a term  $t$  appears in a document  $d$ . However, in order to avoid a predisposition towards long documents, the normalized frequency is used:

$$tf(t, d) = \frac{f(t, d)}{\max\{f(t, d) : t \in d\}}. \quad (1)$$

As shown in Equation (1), the frequency of the term is divided by the maximum frequency of the terms in the document.

- *Inverse document frequency:* If a term appears very frequently in all of the analyzed documents, its weight is reduced. If it appears infrequently, it is increased.

$$idf(t, D) = \log \frac{|D|}{|\{d \in D : t \in d\}|}. \quad (2)$$

As shown in Equation (2), the total number of documents is divided by the number of documents containing the term. Term frequency—Inverse document frequency: The entire formula is as shown in Equation (3).

$$tf - idf(t, d, D) = tf(t, d) \times idf(t, D). \quad (3)$$

Word embedding is a term used for the representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning [11]. Word embeddings can be obtained using a set of language modeling and feature learning techniques where words or phrases from the vocabulary are mapped to vectors of real numbers [12].

- **Word2vec.** This technique uses huge amounts of text as input and is able to identify which words appear to be similar in various contexts [13–15]. Once trained on a sufficiently big dataset, 300-dimensional vectors are generated for each word, forming a new vocabulary where "similar" words are placed close to each other. Pre-trained vectors are used, achieving a wealth of information from which to understand the semantic meaning of the texts.
- **Doc2vec.** This technique is an extension of Word2Vec and is applied to a document as a whole instead of individual words, it uses an unsupervised learning approach to better understand documents as a whole [16]. Doc2Vec model, as opposed to Word2Vec model [17], is used to create a vectorized representation of a group of words taken collectively as a single unit. It does not only give the simple average of the words in the sentence.

## 2.2. Topic Modeling

As already presented in the previous section, topic modeling is a tool that takes an individual text (or corpus) as input and looks for patterns in word usage; it is an attempt to find semantic meaning in the vocabulary of that text (or corpus).

This set of tools enables the extraction of topics from texts; a topic is a list of words that is presented in a way that is statistically significant. Topic modeling programs do not know anything about the meaning of the words in a text. Instead, they assume that each text fragment is composed (by an author) through the selection of words from possible word baskets, where each basket corresponds to a topic. If that is true, then it is possible to mathematically decompose a text into the baskets from which the words that compose it are most likely to come. The tool repeats the process over and over again until the most probable distribution of words within the baskets, the so-called topics, is established.

The techniques executed by the proposed system are used to discover word usage patterns of each user on Twitter, and they make it possible to group users into different categories. To this end, a thorough review of the main tools for topic modeling has been carried out. Most of the algorithms are based on the paradigm of unsupervised learning. These algorithms return a set of topics, as many as indicated in the training. Each topic represents a cluster of terms that must be related to one of those categories. Precisely for this reason, a large number of tweets have been retrieved as training data. Keywords have been searched for for each category. As part of this research, a total of three algorithms have been evaluated: LDA, LSI, and NMF. In the NMF experiment, the best results were obtained, although the techniques applied in other works have been reviewed in order to contrast their results with this method.

Apart from the comparison itself, there are numerous studies that have made similar comparisons between these techniques so that the decision is supported by similar studies. In the work of Tunazzina Islam, in 2019 a similar experiment was carried out to the one proposed in this paper [18]. In this paper, Apache Kafka is employed to handle the big streaming data from Twitter. Tweets on yoga and veganism are extracted and processed in parallel with data mining by integrating Apache Kafka and Spark Streaming. Topic modeling is then used to obtain the semantic structure of the unstructured data (i.e. Tweets). They then perform a comparison of the three different algorithms LSA, NMF, and LDA, with NMF being the best performing model.

Another noteworthy work is that carried out by Chen et al. [19], in which an experiment is carried out to detect topics in small text fragments.

This is similar to the proposal made in this paper, since tweets can be considered small texts. In this work a comparison is made between the LDA and NMF methods, the latter being the one that provided the best results.

- **Latent Dirichlet allocation (LDA).** Is a generative statistical model that allows sets of observations to be explained by unobserved groups that explain why some parts of the data are similar [20–22]. For example, if observations are collections of words in documents, each document is a mixture of a small number of topics and each word's presence is attributable to one of the document's topics. LDA is an example of a topic model and belongs to the machine learning toolbox and in wider sense to the artificial intelligence toolbox.
- **Nonnegative Matrix Factorization (NMF).** Is an unsupervised learning algorithm belonging to the field of linear algebra. NMF reduces the dimensionality of an input matrix by factoring it in two and approximating it to another of a smaller range. The formula is  $V \approx WH$ . Let us suppose, observing Equation (4), a vectorization of  $P$  documents with an associated dictionary of  $N$  terms (weight). That is, each document is represented as a vector of  $N$  dimensions. All documents, therefore, correspond to a  $V$  matrix

$$V \in \mathbb{R}^{N \times P} = \begin{pmatrix} \square & \cdots & \cdots & \square \\ \square & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \square & \cdots & \square & \square \end{pmatrix}, \quad (4)$$

where  $N$  is the number of rows in the matrix, and each of them represents a term, while  $P$  is the number of columns in the matrix and each of them represents a document. Equations (5) and (6) shows matrices  $W$  and  $H$ . The value  $r$  marks the number of topics to be extracted from the texts.

Matrix  $W$  contains the characteristic vectors that make up these topics. The number of characteristics (dimensionality) of these vectors is identical to that of the data in the input matrix  $V$ . Since only a few topic vectors are used to represent many data vectors, it is ensured that these topic vectors discover latent structures in the text.

The  $H$ -matrix indicates how to reconstruct an approximation of the  $V$ -matrix by means of a linear combination with the  $W$ -columns.

$$W \in \mathbb{R}^{N \times r} = \begin{pmatrix} \square & \cdots & \cdots & \square \\ \square & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \square & \cdots & \square & \square \end{pmatrix}, \quad (5)$$

where  $N$  is the number of rows in matrix  $W$ , and each of them represents a term (weight), and  $r$  is the number of columns in matrix  $W$ , where  $r$  is the number of characteristics to be extracted.

$$H \in \mathbb{R}^{r \times P} = \begin{pmatrix} \square & \cdots & \cdots & \square \\ \square & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \square & \cdots & \square & \square \end{pmatrix}, \quad (6)$$

where  $r$  is the number of rows in matrix  $H$ ,  $r$  is the number of characteristics to be extracted, and  $P$  is the number of columns, with one column for each document. The result of the matrix product between  $W$  and  $H$  is, therefore, a matrix of dimensions  $N \times P$  corresponding to a compressed version of  $V$ .

The use of Machine Learning techniques for the analysis of information extracted from Twitter is a very common case study today. It is convenient to study what kind of

research is being carried out on this subject. One of the main applications is the use of Twitter and Natural Language Processing techniques in order to extract a user's opinion about what is being tweeted at a given time. The article "A system for real-time Twitter sentiment analysis of 2012 U.S. presidential election cycle", written by Hao Wang et al. [23], presents a system for real-time polarity analysis of tweets related to candidates for the 2012 U.S. elections.

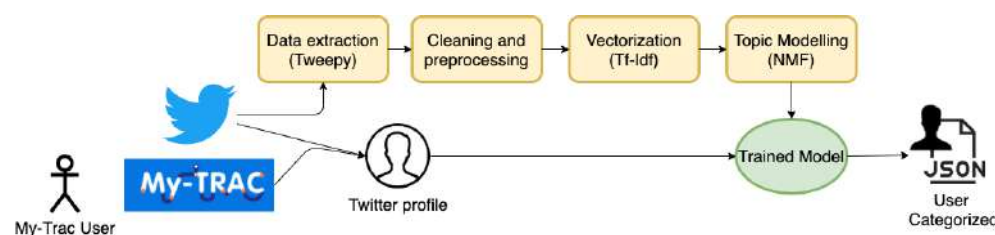
The system collects tweets in real time, tokens and cleans them, identifies which user is being talked about in the tweet, and analyzes the polarity. For training, it applies Naïve Bayes, a statistical classifier. It uses hand-categorized tweets as input. Another study similar to this one is the one proposed by J.M.Cotelo et al. from the University of Seville: "Tweet Categorization by combining content and structural knowledge" [24]. It proposes a method to extract the users' opinion about the two main Spanish parties in the 2013 elections. It uses two processing pipelines, one based on the structural analysis of the tweets, and the other based on the analysis of their content.

Another possible line of research is based on categorizing Twitter content. This is the case of the article "Twitter Trending Topic Classification" written by Kathy Lee et al. [25]. It studies the way to classify trending topics (hashtags highlighted) in 18 different categories. To this end, Topic Modeling techniques were used. The key point lies in providing a solution based on the analysis of the network underlying the hashtags and not only the text: "our main contribution lies in the use of the social network structure instead of using only textual information, which can often be noisy considering the social network context".

As it can be seen, there are many studies currently oriented to the analysis of Twitter using Machine Learning tools. The challenge to be faced in this work is to find the optimal way of classifying users according to their tweets. The sections that follow describe the objectives of the project and detail the research and testing that led to the construction of a stable system fit for the purpose for which it has been designed.

### 3. Proposal

This section proposes a system for the extraction of information about Twitter users. The system is capable of obtaining information about a particular user and of elaborating a profile with the user's preferences in a series of pre-established categories. From an abstract point of view, the proposal could be seen as a processing pipeline, as shown in Figure 1. The different phases of this pipeline contribute to the achievement of the main objective: user classification.



**Figure 1.** Pipeline representing the system processing steps.

#### 3.1. Category Definition

Matching a given profile to a specific category or topic is one of the objectives of NLP algorithms. As a starting point, it is necessary to prepare the training dataset that is used when investigating the algorithmic model. The strategy followed is based on the model of the Interactive Advertising Bureau (IAB) association [26]. Today, IAB is a benchmark standard for the classification of digital content. In particular, the IAB Tech Lab has developed and released a content taxonomy on which the present categorization is based. This taxonomy proposes a total of 23 categories with their corresponding subcategories covering the main topics of interest. In this way, 8000 tweets from each of these categories have been ingested. As a result, 23 datasets with examples of tweets related to each category



were obtained, these datasets have been used to train the system at a later stage. Specifically, the list of topics is shown in Table 3.

**Table 3.** Categories taxonomy.

Topics
Arts & Entertainment
Automotive
Business
Careers
Education
Family & Parenting
Health & Fitness
Food & Drink
Hobbies & Interests
Home & Garden
Law, Gov't & Politics
News
Personal Finance
Society
Science
Pets
Sports
Style & Fashion
Technology & Computing
Travel
Real Estate
Shopping
Religion & Spirituality

### 3.2. Twitter Data Extraction

The Twitter data extraction mechanism is a fundamental element of the system. The goal of this mechanism is to recover two types of data.

On the one hand, the system extracts a set of anonymous tweets related to each of the defined preference categories; these tweets are used to train the data classification algorithms.

On the other hand, the mechanism extracts information about the given user for the analysis of their preferences.

Twitter's API enables developers to perform all kinds of operations on the social network. It is, therefore, necessary for our system to use this powerful API. This API could be used by elaborating a module that would make HTTP requests to the API so that the endpoints of interest are executed. However, this involves a remarkably high development cost.

Another option would be to make use of one of the multiple Python libraries that encapsulate this logic and offer a simple interface to developers. The latter option has been chosen for the development of this system, more specifically, library Tweepy [27].

### 3.3. Preprocessing of Tweets

Once the data has been extracted, it must be prepared for the classification algorithms. Cleaning and preprocessing techniques must be applied, so that the text is prepared for topic modeling algorithms. Libraries, such as NLTK and Spacy, have been used, as can be observed in Listing 1.

The first step involves cleaning tweets, by removing content that does not provide information for language processing. More specifically, this task consists in eliminating URLs, hashtags, mentions, punctuation marks, etc.

Another of the techniques applied to obtain more information from tweets is the transformation of the emojis contained in the text into a format from which it is possible to extract information. To do this, a dictionary of emojis is used as a starting point for the

conversion of the data. This dictionary contains a series of values that interpret each of the existing emojis when applying the corresponding analysis. In this way, it has been possible to identify and give a certain value to each emoji for its treatment.

The key activity performed during the preprocessing consist of eliminating stopwords and tokenization. Whether it is a paragraph, an entire document or a simple tweet, every text contains a set of empty words or stopwords. This set of words is characterized by its continuous repetition in the document and its low value within the analysis. These words are mainly articles, determiners, synonyms, conjunctions, and others.

---

**Listing 1.** Preprocessing step pseudocode.

---

```
from nltk.tokenize import word_tokenize
import~spacy

sp = spacy.load('en_core_web_sm')
stopwords_dict = sp.Defaults.stop_words

def tweet_preprocessing(tweet):
    tweet = hashtag_removal(tweet)
    tweet = mentions_removal(tweet)
    tweet = url_removal(tweet)
    tweet = html_removal(tweet)
    tweet = punctuation_removal(tweet)
    tweet = emojis_removal(tweet)
    tweet = word_tokenize(tweet)
    tweet = [word for word in tweet if not word in stopwords_dict]
    return tweet
```

---

Table 4 shows the results obtained after the tweets have gone through the preprocessing and preparation process which had been carried out using the tools listed above.

### 3.4. Vectorization

Vectorization is the application of models that convert texts into numerical vectors so that the algorithms can work with the data. Two algorithms have been considered for the performance of this task, “Bag-Of-Words” and “Tf-Idf”. Both are widely used in the field of NLP, but, in general, creation of tf-idf weights from text works properly and is not very expensive computationally. Moreover, NMF expects as input a Term-Document matrix, typically a “Tf-Idf” normalized.

The vectorizer have been tuned manually with some parameters according to the dataset, as can be observed in Listing 2. *Min\_df* was set to 100 to ignore words that appear in less than 100 tweets. In the same way, *max\_df* was set to 0.85 to ignore words that appear in more than 85% of the tweets. Thanks to that feature, it is possible to remove words that introduce noise in the model. Finally, the algorithm only takes into account single words, so, in order to include bigrams, the parameter *ngram\_range* was set to (1, 2)

**Table 4.** Preprocessing results using NLTK tokenization.

	Text	Nltk_tokenized
0	Read This Before Taking a Road Trip with a Pet	[read, taking, road, trip, pet]
1	@kenwardskorner @Senators @Canucks In addition, does ...	[also, name, imply, take, acid, road, trips, I...]
2	Our Art is our Passion \n#apnatruckart #truck	[art, passion, apnatruckart, truck, art, uniqu...]
3	Lelang drop acc budget 40–50 k dong?	[lelang, drop, acc, budget, dong]
4	We agree...and want everyone to know that ou	[agree, want, everyone, know, tours, relaxed, ...]
5	Choosing a hotel for a break away with the fam...	[choosing, hotel, break, away, family, special...]
6	@_JassyJass Are you,camping?	[camping]
7	How to Pack Your Electronics for Air Travel ht	[pack, electronics, air, travel]
8	Dasar low budget! <a href="https://t.co/2YUmUrGjj5">https://t.co/2YUmUrGjj5</a>	[dasar, low, budget]

**Listing 2.** Vectorization step pseudocode.

```

from sklearn.feature_extraction.text import TfidfVectorizer

def tfidf_vectorization(tweets):
    vectorizer = TfidfVectorizer(
        min_df=100,
        max_df=0.85,
        ngram_range=(1, 2),
        preprocessor='_'.join,
        use_idf=True
    )
    vectorized_tweets = vectorizer.fit_transform(tweets)
    return vectorized_tweets

```

### 3.5. Topic Modeling

Topic Modeling is a typical NLP task that aims to discover abstract topics in texts. It is widely used to discover hidden semantic structures. In the present work, this technique has been used to discover the main topics of interest of the My-Trac application users based on their Twitter profiles, which should correspond to some of the previously defined categories.

Regarding the features of the model, in Section 3.4, the training tweets were vectorized to create a Term-Document matrix which has been the input of the NMF model. In addition, NMF needs one important parameter, the number of topics to be discovered “*n\_components*”. In this case, *n\_components* was set manually to 23, which is the number of topics that were defined initially in the categories taxonomy. Following this approach,

the algorithm is trained with 184,000 tweets (8000 per category) with the aim of obtaining as many topics as categories were defined in the taxonomy. Once the model has been trained, it has been possible to determine in which topics a user's profile fits on the basis of their tweets. The implementation of the topic modeling algorithm has been carried out on the basis of NMF using SKLearn library, as is deailed in Listing 3.

**Listing 3.** NMF Sklearn implementation.

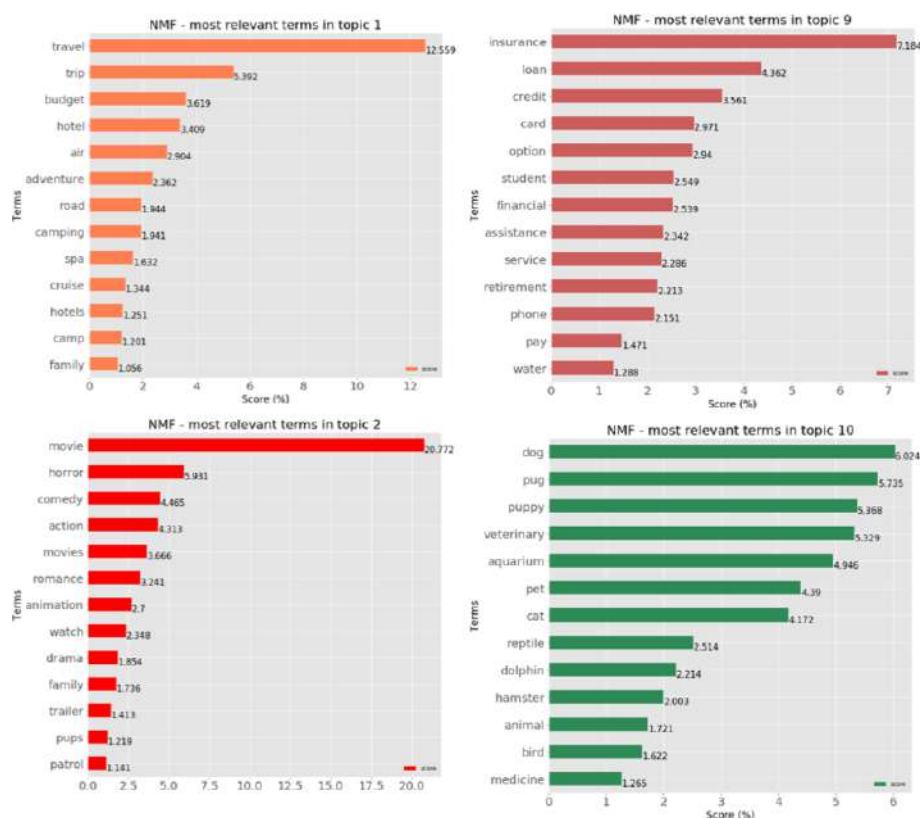
```
from sklearn.decomposition import NMF

def train_model(vectorized_tweets):
    nmf_model = NMF(n_components=23, alpha=.1,
                    l1_ratio=.5, init='nndsvda')
    nmf_model.fit(vectorized_tweets)
    return nmf_model
```

Finally, it is worth mentioning the use of some extra parameters which were set in the implementation of the model. The method used to initialize the procedure was set to "NNDsva" which works better with the tweet dataset since this kind of data it is not sparse. *Alpha* and *l1\_ratio* both are parameters which helps to define regularization.

#### 4. Evaluation and Results

In order to evaluate the results of the algorithm, the most relevant terms have been identified for each resulting topic. Then, by reviewing the main terms for each topic, it is possible to determine if that words really represent the content of the topic. An example is shown in Figure 2, where the most relevant terms have been identified for 4 different topics, proving how well the algorithm identifies the terms associated with each one. As it can be seen, all of them are unambiguously related to their defined categories. Topic 1: Travel. Topic 2: Arts & Entertainment. Topic 9: Personal Finance. Topic 10: Pets.



**Figure 2.** Example topics generated by NMF.

The full list of topics and their top 10 related keywords identified by the algorithm can be seen in Table 5. It should be noted that some of the previously defined categories in Table 3 have been removed during the evaluation of this model. This fact is due to the lack of tweets that would fit into those categories, as well as some topics were quite overlapped amongst them. The initially defined categories that have been removed during training process and evaluation are: “Home & Garden”, “Real State”, “Society”, and “News”. In the same way, the algorithm has been able to discover new categories related to the original ones, such as: “Movies”, “Videogames”, “Music”, “Events”, and “Medicine & Health”, leaving a total of 23 categories in the system.

**Table 5.** Topics obtained by the algorithm.

	Topic	Top 10 Words
1	Travel	travel, trip, budget, air, hotel, adventure, road, camp, family, day
2	Movies	movie, horror, action, comedy, movies, romance, watch, animation, family, drama
3	Videogames	game, xbox, pc, mmo, nintendo, videogame, esports, rpg, play, console
4	Careers	apprenticeship, internship, job, search, career, interview, vocational, training, remote, advice
5	Events	amusement, concert, cinema, restaurant, birthday, match, holiday, football, funeral, park
6	Health & Fitness	health, nutrition, therapy, physical, fitness, workout, exercise, wellness, medicine, weight
7	Religion & Spirituality	islam, christianity, hinduism, judaism, buddhism, spirituality, religion, astrology, atheism, sikhism
8	Shopping	grocery, lotto, shopping, gift, discount, sale, card, coupon, sales, code
9	Personal Finance	insurance, loan, credit, option, card, student, financial, service, assistance, phone
10	Pets	veterinary, dog, pug, puppy, aquarium, pet, cat, reptile, dolphin, hamster
11	Automotive	truck, car, auto, motorcycle, tesla, van, scooter, pickup, luxury, minivan
12	Science	chemistry, geography, geology, biology, physics, genetic, astronomy, environment, math, science
13	Law, Gov't & Politics	election, political, law, issue, vote, news, state, people, country, trump
14	Education	preschool, college, university, exam, electoral, education, homework, student, school, language
15	Food & Drink	coffee, vegetarian, beer, eat, vegan, cook, drink, wine, tea, dining
16	Family & Parenting	marriage, parent, single, baby, daycare, teen, date, life, toddler, adopt
17	Style & Fashion	wear, beauty, clothing, perfume, deodorant, wallet, casual, fashion, shave, trainer
18	Technology & Computing	software, app, developer, mongodb, database, email, android, internet, computer, ai
19	Hobbies & Interests	meme, draw, puzzle, collect, comic_strip, antique, guitar, art, woodwork, painting
20	Sports	martial, rugby, golf, sport, climb, pool, racing, cricket, skating, basketball
21	Business	industry, agriculture, construction, startup, recall, economy, business, automotive, butterfly, turkey
22	Medicine & Health	vaccine, menopause, pregnancy, health, mental, surgery, injury, disease, psychology, substance
23	Music	music, radio, rock, funk, pop, soul, classic, songwriter, listen, classical

Once the resulting model has been evaluated and verified, the next step is to check the effectiveness of the model with real Twitter profiles. The tests have been performed extracting 1200 tweets from different users and predicting for each user the most related topics based on their tweets. The final test results are shown in Table 6, where it can be observed how each profile name match with related topics according to the profile.

As an example, the main topics for the profile “Tesla” are “Automotive”, “Technology and computing”, and “Travel”.

Finally, in order to suggest the main topics of a specific user in the My-Trac app, for each user, the model returns the associated categories, along with the percentage of weight that each category has on the user. The lower the percentage, the less relation the user has with the category. The results of the final classification using some known Twitter accounts are given in Table 7. It should be noted that only the three main categories are shown in the table (together with their associated percentage), as they are the most accurate for categorizing the user.

**Table 6.** NMF evaluation with data from real Twitter profiles.

	Cat 1	Cat 2	Cat 2
Pontifex	Religion & Spirituality	Family & Parenting	Tech & Computing
Tesla	Automotive	Tech & Computing	Travel
BBCNews	Law, Gov’t & Politics	Sports	Family & Parenting
NintendoAmerica	Videogames	Hobbies & Interests	Sports
Theresa_may	Law, Gov’t & Politics	Business	Personal Finance
Oprah	Events	Family & Parenting	Sports
SkyFootball	Sports	Events	Hobbies & Interests
ScuderiaFerrari	Sports	Automotive	Law, Gov’t & Politics
IMDb	Events	Movies	Hobbies & Interests
ScienceMagazine	Science	Medicine & Health	Tech & Computing
Spotify	Music	Hobbies & Interests	Family & Parenting
Airbnb	Travel	Events	Careers

**Table 7.** Final results with different accounts.

	First Category	Second Category	Third Category
Tesla	Automotive (70.92%)	Tech & Computing (10.07%)	Travel (3.86%)
RealDonaldTrump	Law, Gov’t & Politics (28.02%)	Business (11.36%)	Sports (8.02%)
ScienceMagazine	Science (24.41%)	Medicine & Health (16.73%)	Tech & Computing (12.03%)
NintendoAmerica	Videogames (41.65%)	Hobbies & Interests (12.19%)	Sports (9.74%)

## 5. Final System Integration in My-Trac Application

Having passed the entire research and evaluation process, a trained algorithm has been obtained capable of classifying different Twitter accounts according to defined and discovered categories. In addition, a reliable data extraction method has been developed. Therefore, the next step consists of applying the algorithm to the My-Trac app to create a system that allows recommendations to My-Trac users based on their Twitter profiles, which is the objective of the present work.

The final system for My-Trac app consists of a mobile app where the user logs in, as it can be seen in Figure 3, and is asked to grant access to their Twitter data. Once the user signs in to the application, My-Trac seeks for the optimal means of transport to reach a specific destination given by the user and suggests the best conveyance for the trip, as Figures 4 and 5 show.

Finally, when the user chooses the route and mean of transport that best fits his trip, based on the present work, My-Trac app recommends some activities and points of interest for the user during the way based on its Twitter information, as can be observed in Figures 6 and 7.

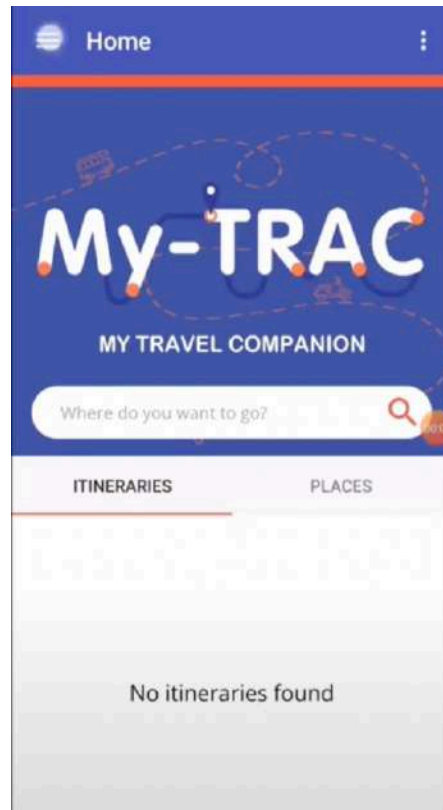


Figure 3. My-Trac application.

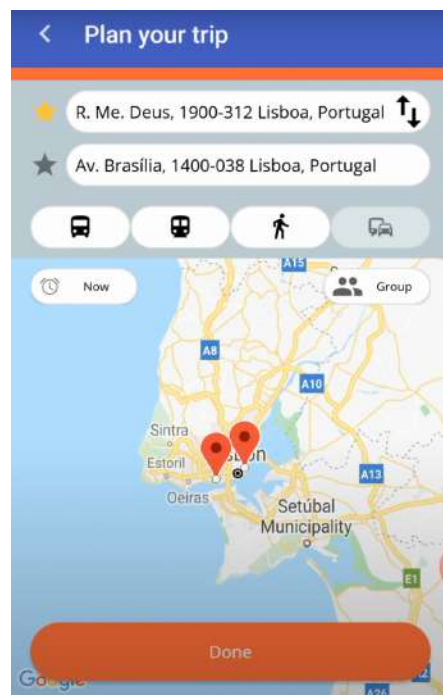


Figure 4. Trip planification using My-Trac app.

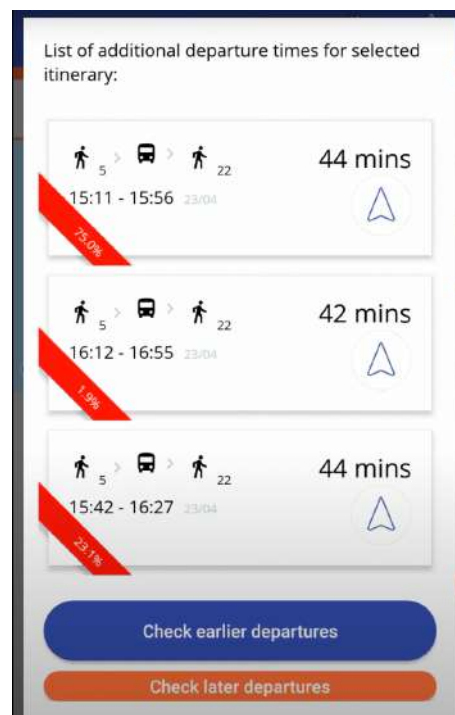


Figure 5. My-Trac suggests optimal means of transport for the destination.

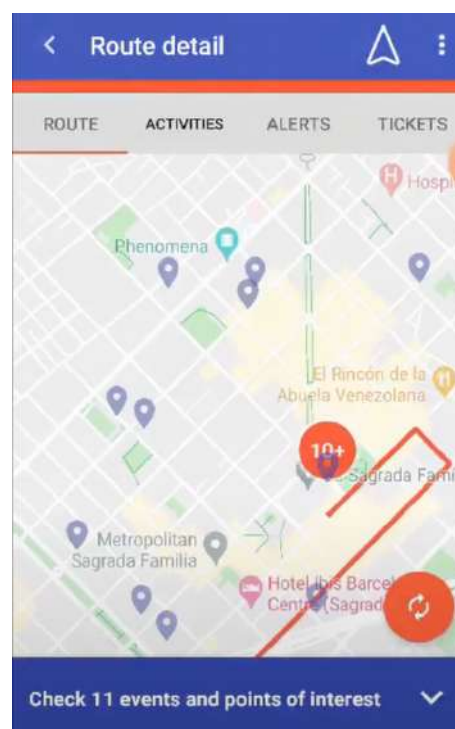


Figure 6. My-Trac recommends activities and point of interest for the user using its Twitter information.





**Figure 7.** My-Trac suggestions.

Moreover, it is possible to get some detailed information for each activity recommended, as Figure 8 shows. In this way, thanks to My-Trac app, the user can improve his experience not only by receiving suggestions for the best conveyance for the trip but also receiving customized activity recommendations and points of interest.



**Figure 8.** Detailed information about a suggested point of interest.

## 6. Conclusions and Future Work

This article presents a novel approach to extracting preferences from a Twitter profile by analyzing the tweets published by the user for use in mapping applications. This approach has successfully defined a consistent and representative list of categories, and the mechanisms needed for information extraction have been developed, both for model training and end-user analysis. It is a unique system, with which it has been possible to

develop an important feature in the My-Trac app, whereby it is possible to recommend relevant point of interest to the end users.

Regarding future work on this system, many areas of improvement and development have been identified. Tweets are not the only source of information that allows to discern the interests of a profile. It may be the case that a user only writes about football but is the follows many news-related and political accounts. The current system would only be able to extract the sports category. Therefore, one of the improvements would be the implementation of a model that would analyze followed users. This has been started, by extracting the followers and creating wordclouds with the most relevant ones. Similarly, hashtags also provide additional information suitable for analysis. Another line of research is the training of a model that allows to analyze the tweets individually. This would open the doors to performing a polarity analysis that would allow us to know if a user who writes about a certain category does it in a positive, negative, or neutral way.

As for the limitations of the system, it is possible that, in some regions, there may be restrictive regulations on the use of information published on social networks for this type of analysis. Therefore, the user should carry out a study of data protection and the legal framework adapted to each region where the service is to be provided. Furthermore, in terms of performance, it is possible that specific context-dependent systems training an algorithm for each individual user may perform slightly better than the proposed solution.

**Author Contributions:** Conceptualization, A.R. and A.G.-B.; methodology, A.R. and A.G.-B.; software, A.R. and J.J.C.-M.; validation, A.R., A.G.-B., J.J.C.-M. and A.P.-P.; formal analysis, A.R. and J.J.C.-M.; investigation, A.R., A.G.-B. and A.P.-P.; resources, A.R., A.G.-B., J.J.C.-M.; data curation, J.J.C.-M.; writing—original draft preparation, A.R., A.G.-B., J.J.C.-M.; writing—review and editing, A.R., A.G.-B., J.J.C.-M., A.P.-P. and J.M.C.; visualization, A.R. and J.J.C.-M.; supervision, A.G.-B., A.P.-P. and J.M.C.; project administration, A.G.-B., A.P.-P. and J.M.C.; funding acquisition, J.M.C. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by the Spanish Ministerio de Ciencia e Innovación under grant number TIN2017-89314-P.

**Acknowledgments:** This research has been supported by the European Union's Horizon 2020 research and innovation programme under grant agreement No 777,640.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Ludwig, B.; Zenker, B.; Schrader, J. Recommendation of Personalized Routes with Public Transport Connections. In Proceedings of the International Conference on Intelligent Interactive Assistance and Mobile Multimedia Computing, Rostock-Warnemünde, Germany, 9–11 November 2009; Springer: Berlin/Heidelberg, Germany, 2009; pp. 97–107.
2. Cui, G.; Luo, J.; Wang, X. Personalized travel route recommendation using collaborative filtering based on GPS trajectories. *Int. J. Digit. Earth* **2018**, *11*, 284–307. [CrossRef]
3. Briedenhann, J.; Wickens, E. Tourism routes as a tool for the economic development of rural areas—Vibrant hope or impossible dream? *Tour. Manag.* **2004**, *25*, 71–79. [CrossRef]
4. Cea-Morán, J.J.; González-Briones, A.; De La Prieta, F.; Prat-Pérez, A.; Prieto, J. Extraction of Travellers' Preferences Using Their Tweets. In Proceedings of the International Symposium on Ambient Intelligence, L Aquila, Italy, 17–19 June 2020; pp. 224–235.
5. De Pessemier, T.; Minnaert, J.; Vanhecke, K.; Dooms, S.; Martens, L. Social recommendations for events. In Proceedings of the CEUR workshop Proceedings, Miami, FL, USA, 1 October 2013; Volume 1066.
6. Stathopoulos, E.A.; Paliokas, I.; Meditskos, G.; Diplaris, S.; Tsafaras, S.; Valkouma, E.; Pehlivanides, G.; Riggas, C.; Vrochidis, S.; Votis, K.; et al. Smart discovery of cultural and natural tourist routes. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence-Companion, Thessaloniki, Greece, 14–17 October 2019; pp. 208–214.
7. Sansonetti, G.; Gasparetti, F.; Micarelli, A.; Cena, F.; Gena, C. Enhancing cultural recommendations through social and linked open data. *User Model. User-Adapt. Interact.* **2019**, *29*, 121–159. [CrossRef]
8. Garcia, A.; Arbelaitz, O.; Linaza, M.T.; Vansteenwegen, P.; Souffriau, W. Personalized tourist route generation. In Proceedings of the International Conference on Web Engineering, Vienna, Austria, 5–9 July 2010; Springer: Berlin/Heidelberg, Germany, 2010; pp. 486–497.
9. University, Y. About Yale: Yale Facts. 2017. Available online: <https://www.yale.edu/about-yale/yale-facts> (accessed on 24 May 2021).

10. Demestichas, K.; Kosmides, P. An offline, statistical method for cost efficient design of experiments and field trials involving electric vehicles. In Proceedings of the 11th ITS European Congress, Glasgow, Scotland, 6–9 June 2016.
11. Ferrari, A.; Donati, B.; Gnesi, S. Detecting domain-specific ambiguities: An NLP approach based on Wikipedia crawling and word embeddings. In Proceedings of the 017 IEEE 25th International Requirements Engineering Conference Workshops (REW), Lisbon, Portugal, 4–8 September 2017; pp. 393–399.
12. Wallace, E.; Wang, Y.; Li, S.; Singh, S.; Gardner, M. Do nlp models know numbers? Probing numeracy in embeddings. *arXiv* **2019**, arXiv:1909.07940.
13. Church, K.W. Word2Vec. *Nat. Lang. Eng.* **2017**, *23*, 155–162. [CrossRef]
14. Rong, X. word2vec parameter learning explained. *arXiv* **2014**, arXiv:1411.2738.
15. Lilleberg, J.; Zhu, Y.; Zhang, Y. Support vector machines and word2vec for text classification with semantic features. In Proceedings of the 2015 IEEE 14th International Conference on Cognitive Informatics & Cognitive Computing (ICCI\* CC), Beijing, China, 6–8 July 2015; pp. 136–140.
16. Bilgin, M.; Şentürk, İ.F. Sentiment analysis on Twitter data with semi-supervised Doc2Vec. In Proceedings of the 2017 International Conference on Computer Science and Engineering (UBMK), London, UK, 5–7 July 2017; pp. 661–666.
17. Chen, Q.; Sokolova, M. Word2Vec and Doc2Vec in unsupervised sentiment analysis of clinical discharge summaries. *arXiv* **2018**, arXiv:1805.00352.
18. Islam, T. Yoga-veganism: Correlation mining of twitter health data. *arXiv* **2019**, arXiv:1906.07668.
19. Chen, Y.; Zhang, H.; Liu, R.; Ye, Z.; Lin, J. Experimental explorations on short text topic mining between LDA and NMF based Schemes. *Knowl.-Based Syst.* **2019**, *163*, 1–13. [CrossRef]
20. Resnik, P.; Armstrong, W.; Claudino, L.; Nguyen, T.; Nguyen, V.A.; Boyd-Graber, J. Beyond LDA: Exploring supervised topic modeling for depression-related language in Twitter. In Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality, Denver, CO, USA, 5 June 2015; pp. 99–107.
21. Mehrotra, R.; Sanner, S.; Buntine, W.; Xie, L. Improving lda topic models for microblogs via tweet pooling and automatic labeling. In Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval, Dublin, Ireland, 28 July–1 August 2013; pp. 889–892.
22. Tajbakhsh, M.S.; Bagherzadeh, J. Semantic knowledge LDA with topic vector for recommending hashtags: Twitter use case. *Intell. Data Anal.* **2019**, *23*, 609–622. [CrossRef]
23. Wang, H.; Can, D.; Kazemzadeh, A.; Bar, F.; Narayanan, S. A system for real-time twitter sentiment analysis of 2012 us presidential election cycle. In Proceedings of the ACL 2012 System Demonstrations. Association for Computational Linguistics, Jeju Island, Korea, 8–14 July 2012; pp. 115–120.
24. Coteló, J.M.; Cruz, F.L.; Enríquez, F.; Troyano, J. Tweet categorization by combining content and structural knowledge. *Inf. Fusion* **2016**, *31*, 54–64. [CrossRef]
25. Lee, K.; Palsetia, D.; Narayanan, R.; Patwary, M.M.A.; Agrawal, A.; Choudhary, A. Twitter trending topic classification. In Proceedings of the 2011 IEEE 11th International Conference on Data Mining Workshops, Vancouver, BC, Canada, 11 December 2011; pp. 251–258.
26. IAB Categories | MoPub Publisher UI | MoPub Developers. Available online: <https://developers.mopub.com/publishers/ui/iab-category-blocking/> (accessed on 4 May 2021).
27. Tweepy. Available online: <https://www.tweepy.org/> (accessed on 4 May 2021).



## Article

# A Generic Data-Driven Recommendation System for Large-Scale Regular and Ride-Hailing Taxi Services <sup>†</sup>

Xiangpeng Wan, Hakim Ghazzai \*  and Yehia Massoud

School of Systems and Enterprises, Stevens Institute of Technology, Hoboken, NJ 07030, USA;  
xwan6@stevens.edu (X.W.); ymassoud@stevens.edu (Y.M.)

\* Correspondence: hghazzai@stevens.edu

<sup>†</sup> This paper is an extended version of our paper published in IEEE Conference on Vehicular Electronics and Safety (ICVES'19), Cairo, Egypt, 4–6 September 2019.

Received: 6 March 2020; Accepted: 9 April 2020; Published: 15 April 2020

**Abstract:** Modern taxi services are usually classified into two major categories: traditional taxicabs and ride-hailing services. For both services, it is required to design highly efficient recommendation systems to satisfy passengers' quality of experience and drivers' benefits. Customers desire to minimize their waiting time before rides, while drivers aim to speed up their customer hunting. In this paper, we propose to leverage taxi service efficiency by designing a generic and smart recommendation system that exploits the benefits of Vehicular Social Networks (VSNs). Aiming at optimizing three key performance metrics, number of pick-ups, customer waiting time, and vacant traveled distance for both taxi services, the proposed recommendation system starts by efficiently estimating the future customer demands in different clusters of the area of interest. Then, it proposes an optimal taxi-to-region matching according to the location of each taxi and the future requested demand of each region. Finally, an optimized geo-routing algorithm is developed to minimize the navigation time spent by drivers. Our simulation model is applied to the borough of Manhattan and is validated with realistic data. Selected results show that significant performance gains are achieved thanks to the additional cooperation among taxi drivers enabled by VSN, as compared to traditional cases.

**Keywords:** intelligent transportation systems; demand prediction; taxi recommendation; vehicle social network; ride-hailing

## 1. Introduction

Modern urbanization has significantly changed people's living arrangements, making public transportation, particularly taxi services, a convenient and affordable means of travel for most people, especially when owning a car and paying parking fees is exorbitant. In New York city, 80% of the residents do not own a car [1]. This leads to an explosive growth of the taxi fleet size (e.g., regular yellow taxis in New York city), and ride-hailing service demand, which results in increasing congestion and inefficient exploitation of the resources. For regular taxi services, like yellow taxis in New York city, the taxi drivers do not know the exact locations of potential customers, while for the ride-hailing taxi services, such as Uber, Lyft, and Didi, customers send requests with their locations to nearby ride-hailing vehicles. In both taxi services, and independently of the level of knowledge about the customers' demand, users experience long waiting time periods before getting a ride. At the same time, taxi drivers are engaged in a tedious customer hunting search, traveling long distances. Indeed, even with the ride-hailing service, customers may find out that the nearest available vehicle needs a long time to pick them up. Therefore, there is a pressing need to improve the utilization of such

a means of transportation and enhance the efficiency of both services for the benefits of both customers and drivers.

In regular taxi services, traditional ways for taxi drivers to find potential customers include driving around the city and waiting at some 'hot spots', e.g., taxicab stands. For the first option, taxi drivers usually follow an intuition-based trajectory hoping to find customers as soon as possible, while for the second option, most of the drivers will target the same hot spots since based on their personnel experience, they know when and where customers will be gathered. In the latter case, regular taxi drivers may be subject to an unfair competition since the number of taxis is higher than the demand or vice versa. Hence, traditional solutions for customer hunting are usually exhaustive and inaccurate. On the other hand, for the ride-hailing taxi services, although a central server is dedicated to manage the requests of customers and allocate them to drivers, similar problems that face regular taxi services still exist. Customers' requests might still be raised far away from drivers' locations and high vacant distances are accumulated, resulting in huge and redundant fuel consumption. In Portland, the average waiting times are estimated to be around six and ten minutes for regular and ride-hailing taxi services, respectively, according KGW News [2]. Therefore, it is recommended to enhance the efficiency of such transportation services by tackling the offer/demand problem in both taxi categories.

Thanks to the spread of on-board and infrastructure-based sensors [3], collecting and sharing data have become very common, especially in urban areas, where several novel data-driven applications exist, including Google Navigation, Waze, and parking localization service. This is additionally boosted by the emerging concept of vehicular social network (VSN), which effectively exploits the data availability in transportation networks [4,5]. With the installation and spread of on-board sensors, the data sharing ability has dramatically increased [3]. Mobile apps like Google Navigation and Waze utilize the historical traffic data and human-report accidents to improve the navigation services. The emerging concept of vehicular social network (VSN) has been proposed to better exploit the data availability among road users and transportation networks. A variety of applications and use cases have been discussed in [4–7]. VSN enables interactions between different participants, including human-to-vehicle and vehicle-to-vehicle interactions [8,9]. As an example, the connected vehicle technology in NYC is developed to leverage the safety of road users. It relies on vehicle-to-vehicle (V2V), vehicle-to-infrastructure (V2I) and infrastructure-to-pedestrian (IVP) to share information among them and better assess the transportation network (<https://www.cvp.nyc/>). Hence applied to our context, VSN can be utilized for improving the communication among taxi drivers and exploit their information to revamp the operation of taxi drivers by enabling efficient and real-time identification and sharing of their locations, as well as knowledge about customers such as pick up time, pick up places, drop off time, and drop off places, as well as accurate and relevant data about the traffic situation. Such real-time data sharing can provide a clearer vision about the current customers' requests and help continuously predict the future demand at different regions of the navigation map [10,11]. This technological advance significantly contributes to designing novel taxi recommendation solutions for customer hunting [12] or involving highly connected autonomous taxis [13].

In this paper, we propose a combination of data-driven solutions that jointly improve the taxi service efficiency by recommending the operation of both regular taxicabs and ride-hailing taxis [14]. The proposed recommendation systems consists of three phases: (i) a demand prediction phase, (ii) a taxi-to-region matching phase, and (iii) a route planning phase. The proposed system divides the geographical urban area into several sub-regions and predicts the future demand during the next time periods for each region. Afterwards, it assigns taxis to this region based on the predicted demand. The number of taxis associated to each region is determined such that redundant taxi travel is avoided/reduced. This is performed by taking into account the current locations of the taxis and the predicted demand of each region. The problem is modeled as a bipartite graph which is designed such that the total expected traveled distance for taxis during the transition phase (i.e., taxis moving to their assigned locations) is minimized. Finally, the taxi recommendation system employed for realistic maps provide to drivers optimized trajectories to follow given real-time traffic data.

The realistic map is converted into a graph and the Dijkstra's algorithm is applied to determine the fastest paths for each member of the taxi fleet when needed. Three key principal performance indicators, namely total number of pick-ups, total customer waiting time, and total traveled distance for vacant taxis, are evaluated for both regular and ride-hailing taxi services and employed to compare our proposed system versus traditional solutions using realistic data of the area of Manhattan, Borough of New York city.

The main contributions of this paper are summarized as follows:

- We develop demand prediction models to precisely estimate the future demands in each region of the area of interest. One is online learning time series and the other is Long Short Term Memory model (LSTM). Their accuracies are validated using realistic data.
- We convert the taxi-to-region matching problem into a bipartite perfect matching graph, where we evenly assign taxis to different regions based on their future demands and the current locations of the taxis.
- We optimize the routing of each taxi by minimizing the expected time spent from its current location to the guided destination using the Dijkstra's algorithm by considering the real-time traffic data and geodesic distances of the road network.
- We develop a real-time simulated taxi operation based on the recommendation system using realistic maps, which provides evidence that significant performance gains can be achieved as compared to the traditional case.

The rest of the paper is organized as follows. Section 2 provides a literature review. Section 3 presents the system model and the adopted methodology. Section 4 develops the taxi recommendation system. Section 5 describes the proposed simulated taxi operation. Section 6 presents and discusses selected simulation results. Finally, concluding remarks and future directions are drawn in Section 7.

## 2. Related Work

Over the last few years, researchers have focused on designing solutions to support taxi drivers in enhancing their services. One of the main research directions is the identification of hot spot areas and the prediction of the demand, e.g., using Gaussian process regression [15] or reinforcement learning [16]. The objective is to identify regions with high likelihood of finding potential customers by predicting the spatial distribution of taxi passengers for a short-term time horizon [17,18]. The recommendation system assigns hot spot areas to vacant taxi drivers in order to shorten the waiting time for customers [19]. In [20], the authors proposed a mutual recommendation system that assigns hot spots for both taxi and passengers based on the trajectory of taxis. In [21], the authors developed a route recommendation engine to minimize vacant traveled distance through Monte Carlo tree search algorithm. These studies mainly focus on a single taxi and do not consider the situations where some hot spot areas are attracting a number of taxis larger than the needed demands or the opposite case. Some researchers focused on designing algorithms for ride-sharing services while addressing different research questions including taxi-to-customers assignment, demand and pricing, competition impacts, etc. [22]. In [23,24], the authors proposed Integer Linear Programs (ILP) that can match large groups of riders to a fleet of shared vehicles in real-time with certain capacity size. The algorithms are designed to address the current situation without considering future demands. Moreover, their computational complexity remains high. The adopted routing methods are based on the shortest path algorithm which does not consider traffic data and congestion level.

Spatial-demand prediction was one of the essential topics that are investigated in the context of taxi recommendation systems. In [25], the authors proposed Multi-View Spatial-Temporal Network (DMVST-Net) approach to predict the taxi demand. It is shown that the proposed method achieves a Mean Absolute Percentage Error (MAPE) of  $\approx 16\%$ . However, the predicted results are daily instead of hourly, which is not suitable for assisting drivers. Moreover, the running time to generate hourly results is also high. In [26], the authors predicted the short-term supply-demand gap of taxis by partitioning the city area into various regular Hexagon lattices-based Convolutional Neural Networks

(H-CNN). However, the proposed model is also computationally expensive compared to traditional methods while achieving slightly better performance. More importantly, it is not necessarily true that cities have uniform partitioning of their area, such as the case of Manhattan, NYC. Therefore, in this study, we use the cities' own region partition to predict the future demand using a faster algorithm in a real-time manner.

Recently, taxi recommendation studies consider more generalized scenarios and are not limited to a specific task. For instance, in [27], the authors developed a recommendation system for taxis by jointly considering the benefits of both drivers and passengers. The driver's utility includes expected revenue, searching time for next passenger, traveled distance, while the passenger's utility includes the waiting time. The authors grouped pick-up locations into clusters and defined them as the hot spot areas, to which it assigns taxis according to their scores. The recommendation system efficiently works for selected hot spot areas but ignores other areas with lower demand. Also, the speed of the vacant taxis is assumed to be constant which is not very practical. In [28], the authors presented a receding horizon control framework to dispatch taxis, with the demand prediction based on the estimated demand distribution. The system is evaluated on a square region without practical road network. In [29], the authors analyzed the dynamic spatial equilibrium of taxis and provided efficient regulation for taxi services in different regions. In [30], the authors presented a two-stage stochastic optimization formulation to consider expected future demand to solve the spatio-temporal matching problem, i.e., taxi matching. Generally, most of the studies discussed earlier do not consider the real-time locations of the taxis.

Furthermore, some other researchers focused on the cruising and matching for the taxi drivers. In [31], the authors provided a data-driven simulation framework for ride-sharing taxis simulated in a simplistic grid map. The proposed approach provides a path for a taxi while optimizing a certain cost function, such as traveled distance or gasoline consumption. In [32], the authors found out that driver's cruising choice is learned from his/her previous experience and his/her interactions with other drivers. In [33], the authors proposed pCruise system to reduce the taxi's cruising miles by providing the shortest cruising route with at least one expected available passengers for this route. In [34], the authors developed efficient algorithms for non-myopic adaptive routing to minimize the collective travel time of all vehicles in the system. In [35], the authors proposed solutions to reduce the number of cruising miles while increasing the number of live miles of taxis by suggesting profitable locations to taxicab drivers. Other research directions have investigated dynamic models to arrange ride-sharing vehicles with discrete simulation environment [36,37]. The authors of [38,39] have proposed data-driven vehicle re-balancing across regions but lack future demands prediction. Some researchers provided a graph partitioning methodology to partition the bipartite graph with lower computational complexity and implemented it in the one-to-one ride-matching problems [40]. Another study has modeled the matching problem as a competition strategy between different ride-hailing companies [41]. Despite the previous studies providing solutions for taxi cruising and matching problem, most of the methods are built in simplistic maps without convincing evidence to show the practicality of their methods. Moreover, they did not take the demand prediction, taxi dispatch, and route selection together into consideration. To the best of our knowledge, the recommendation system that we propose is the first one which jointly takes into account the prediction of future demands, taxi dispatch, and cruising routes selection for both regular and ride-hailing taxi services and is validated using realistic data and map.



### 3. System Model and Methodology

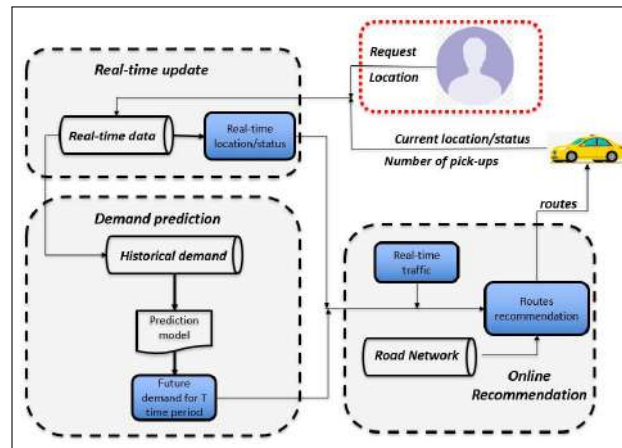
We propose to design a novel recommendation system for taxis cruising on a large geographical area. The latter is sub-divided into multiple regions for which we aim to predict the demand based on their respective historical data. The demand in the area of interest is estimated and updated in every time period  $T$ . In this paper, we focus on both the regular and ride-hailing taxi services. The difference is that regular taxi drivers are supposed to not know the exact locations of the customers as the ride-hailing vehicles, which are informed by the exact locations once they receive the request. Hence, we assume that for regular taxis, the pick-up happens when a taxi driver sees a customer waving his/her hand (e.g., when the distance between the customer and taxi is less than 100 m). In Table 1, we present the summary of the different taxi services managed by the proposed recommendation system.

**Table 1.** The three taxi services managed by the proposed recommendation system.

Service	Trad. Regular	Smart Regular	Ride-Hailing
Taxi Call	Waving	Waving	Online Request
Information Sharing	Without VSN	With VSN	With VSN
Knowledge level	Taxis locations only	Taxis locations only	Taxis and customers locations
	Future demand every $T$	Future demand in real-time	Future demand in real-time
Recommendation	Once every $T$	Continuously	Continuously

Note that the regular taxi services can be split into two categories: the traditional regular taxi services without VSN and the smart regular taxi services with VSN. In both services, taxi drivers are not aware of the locations of customers, but with the smart regular taxi services, when a pick-up happens, all other taxis via the recommendation system are aware about it. Hence, they are instantaneously updated about the changes in the area of interest. In other words, the system can adjust the hunting search locations for vacant taxis during the period  $T$  in a real-time manner instead of waiting until the end of the time period  $T$ , as it is the case with traditional services. For the ride-hailing taxi service, the taxis are aware of the locations of both users' demand and taxis in real-time and hence, it continuously provides recommendations to vacant taxis.

In Figure 1, we present the overview of the proposed framework for regular and ride-hailing taxi services. There are three major phases: the first phase is the real-time data update phase where information is collected from customers and taxi drivers. The data includes the current locations of customers and taxis in addition to the statuses of taxis (vacant or occupied) and the number of pick-ups already done. The second phase is the demand prediction phase that is executed every period  $T$ . In this phase, the historical data is used to predict the demand of the area of interest. Note that for every  $T$  time period, the system would predict the demand only once, set at the beginning of that time period  $T$ . Then, the demand would be updated by considering the number of pick-ups happening during the entire time period  $T$ . The third phase encompasses the process of taxi-to-region matching and taxi routing. For the taxi-to-region matching, the recommendation system assigns vacant taxis to the different regions based on their locations and the potential future demand on that region, e.g., if the system recommends several taxis to some regions, it will only send them to the nearby ones. For the route selection, the system determines the routes for all taxis to reach their destinations by minimizing the expected time spent on their trips by considering the collected real-time traffic data.



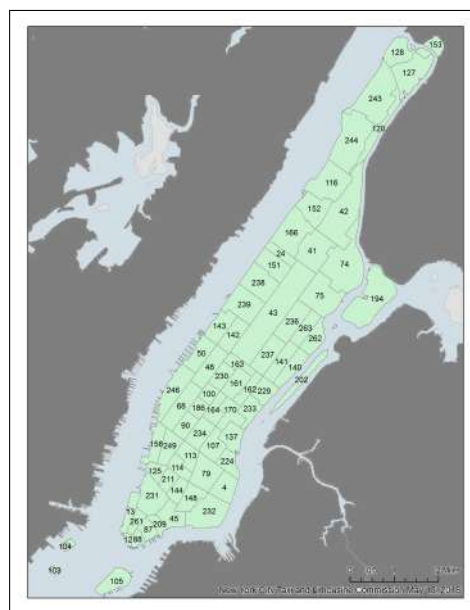
**Figure 1.** Recommendation framework for regular (without red curved rectangle) and ride-hailing (with red curved rectangle) taxi services.

#### 4. Proposed Taxi Recommendation System

In this section, we introduce the different components and steps of the proposed taxi recommendation system: (1) the taxi demand predictor, (2) the taxi-to-region matching component, and (3) the taxi routing optimizer.

##### 4.1. Taxi Demand Predictor

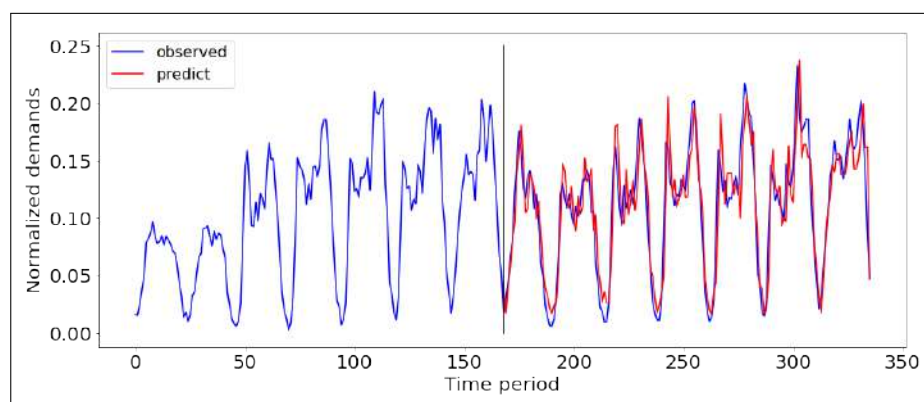
The first step is to predict the customer demand in the area of interest. We introduce and compare two models that fit the scope of this task. One is Long Short Term Memory (LSTM) model and the other is Autoregressive integrated moving average (ARIMA) model. To illustrate their accuracy, we collect the data about the operation of yellow taxis from the Taxi Limousine Commission (TLC) (<https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page>), which contains the taxi operation information in New York City including the pick-up instants, pick-up locations, drop-off time, drop-off region, trip fare, and trip distance. We then explore the historical demands on the borough of Manhattan which is split into 69 regions as shown in Figure 2. Before feeding the data into our models, we normalize the demands over  $T$  ( $T = 60$  min in this case) at first.



**Figure 2.** The borough of Manhattan and its taxi region subdivision.

The next step is to predict the future demand for the next period  $T$  on each region with ARIMA. In our case, we are using the demand of the previous 168 h to predict the demand of the next hour. In other words, we are using historical data for the previous week to predict the demands in the next hour, which automatically takes into account the weekday and weekends by assessing the trend of its consecutive features (the model could figure out if the date of prediction is a weekend or weekday). To prevent inputting extra information as weekday or weekends, we proceed by predicting the next hour of demands based on previous one-week data inputs. That is, using 168 previous inputs, we predict the next time period's demand, which would take the holidays, weekends, and weekdays into consideration by assessing the trend of its consecutive features. The choice of the demand prediction period is not arbitrary. It considers the objective of the next phase of the proposed recommendation system. Indeed, every hour, the taxi-to-region matching is provided after precisely predicting the hourly demand. Hence, choosing longer values of  $T$  may be unadapted with the demand variation in the region and may lead to taxi operation delay. Selecting lower values of  $T$  will increase the frequency of executing the taxi-to-region matching algorithm, which may lead to either an excessive re-assignment of taxis, which is not practical, redundant results similar to the ones of the previous time period, as well as extra computational complexity. More importantly, when predicting the traffic demand over the time period  $T$ , we aim to estimate the demand of each region at each instant of that period and not only a constant demand. With the help of VSN, the traffic and demand information are shared among the fleet instantaneously. For instance, the predicted remaining customers during the time period are estimated while considering the number of already picked customers.

We set the ARIMA parameter  $d$  to zero. In order to find the best model, we test different ARIMA models with different parameter combinations  $(p, d, q)$  where we pick the parameters with the lowest Akaike information criteria (AIC) value in the end. We find out that ARIMA with  $(p = 5, d = 0, q = 3)$  provides the lowest AIC where  $AIC = -2.9 \times 10^4$ . The ARIMA model fitting is based on the latest updated 168 data time periods before predicting the next time period. In this case, the predicting process is converted into an online learning where the model parameters are updated continuously. The prediction result from the ARIMA model is shown in Figure 3 where the red and blue series correspond to the predicted and actual values, respectively. The resulting mean square error (MSE) is  $4.7 \times 10^{-4}$ . Hence, we can conclude that the online ARIMA model can provide accurate prediction of the future demand, which can be effectively used to provide precise prediction for the taxi recommendation system.



**Figure 3.** Hourly predicted demand for region 114 using the Autoregressive integrated moving average (ARIMA).

We have compared the performance of the ARIMA model to the LSTM one, which is trained on the historical data first and then employed to predict the future taxi demand. The LSTM model contains two hidden layers and one output neuron. The input shape is 168, which contains previous one-week hourly demand data. It achieves an MSE equal to  $6.9 \times 10^{-4}$  as illustrated in Figure 4.

Unlike the ARIMA model, the LSTM is not trained in an incremental/online manner, which requires a more important amount of data compared to ARIMA. From the comparison results, we find out that the online ARIMA model is more accurate, hence we adopt it in our system.

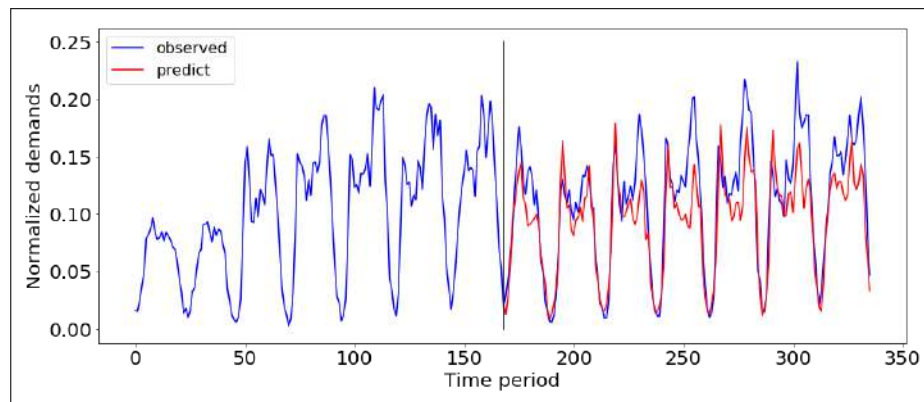


Figure 4. Hourly predicted demand for region 114 using Long Short Term Memory model (LSTM).

#### 4.2. Taxi-to-Region Matching Component

Once accurate future taxi demand is determined for each region, we proceed by assigning vacant taxis to these regions according to the region demands and the taxi current locations. The first metric is used to ensure that the taxi assignment is made proportionally to the demand. In this way, when the expected demand is high, more taxis will be sent to that region and vice versa. The second metric (taxi current locations) is considered in order to minimize the transition phase during which vacant taxis need to reach their assigned regions. This helps in reducing the waiting time of the customers looking for rides at the beginning of the time period. An example is shown in Figure 5 where four regions and eight taxis exist. Obviously, given the demand, we should assign one taxi to region A, two taxis to region B, four taxis to region C, and one taxi to region D based on their respective demand ratios (10, 20, 40, 10).

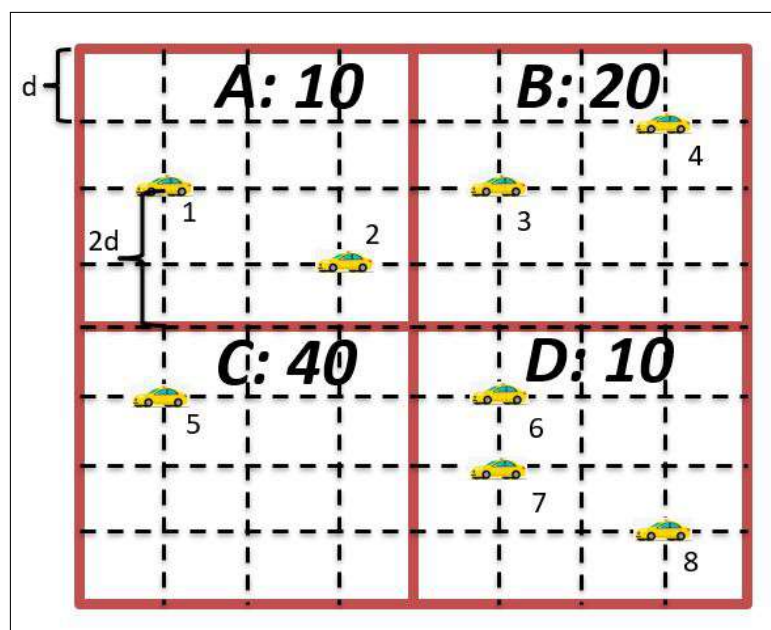
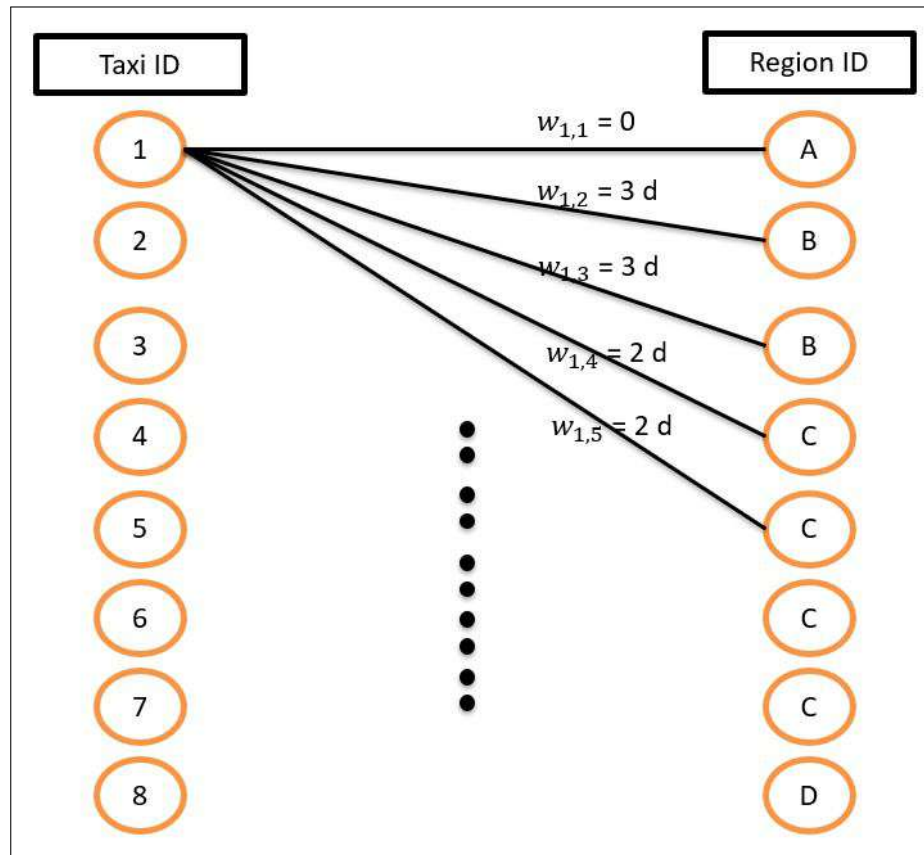


Figure 5. An illustrative example of an  $8 \times 8 d^2$  area composed of four regions A, B, C, and D with different demands 10, 20, 40, and 10, respectively. Eight vacant taxis are circulating around the region. The black dashed lines represent roads.

To ensure an efficient taxi-to-region matching for large-scale problems, we propose to model it by a bipartite weighted graph presented in Figure 6. The weights associated to the graph edges are computed based on the shortest distance needed by the taxi to reach the closest border of the region. To reflect the demand of each region in the graph, we duplicate the ones having higher demand multiple times according to their normalized demand levels with respect to the total number of taxis and total demand in the Borough of Manhattan during that time period. Consequently, the objective of the taxi-to-region matching component is to minimize the sum of the weights while maintaining the perfect matching. In other words, each taxi is assigned to one region. Note that, in practice, the number of taxis is usually higher than the number of regions. We refer to the taxi ID by the index  $i$  and the region ID after duplication by the index  $j$ . Hence, as shown in Figure 6,  $j = 2$  and  $j = 3$  refer to the same region B.



**Figure 6.** An illustrative example of the bipartite matching graph applied to the scenario presented in Figure 5.

The aforementioned matching procedure can be optimized using an ILP problem formulated as follows:

$$(P): \underset{x_{i,j} \in \{0,1\}}{\text{minimize}} \quad \sum_i \sum_j w_{i,j} x_{i,j} \quad (1)$$

subject to:

$$\sum_j x_{i,j} = 1, \quad \forall i, \text{ and } \sum_i x_{i,j} = 1, \quad \forall j, \quad (2)$$

where  $w_{i,j}$  represents the weights of the edges linking the taxis and the regions while  $x_{i,j}$  is a decision variable indicating whether a taxi  $i$  is assigned to region  $j$  or not. It is equal to 1 if this is the case. In (P), constraints (2) ensure the perfect matching, which forces a taxi to be assigned to only one region.

The matching problem can be also solved optimally using the heuristic minimum weight perfect matching algorithm: the Hungarian method. This algorithm solves the problem in a polynomial time  $\approx \mathcal{O}(N)$ , much faster than the NP-complete ILP-based solution that adopts the branch-and-bound algorithm, where  $N$  is the number of taxis.

#### 4.3. Taxi Routing Optimizer

The road network can be constructed in many ways, such as simple search techniques and complex fuzzy logic theory [42,43]. In this paper, we transform the traffic network of the area of interest into a complex graph composed of intersections and roads. Each road  $r$ , where  $r \in \{1, \dots, R\}$  connecting at most two intersections, is divided into multiple small segments with the same length  $l_r$ . The graph then has as vertices the connections of different segments and as edges the segments themselves. We define the current location of taxi  $i$  as  $(S_i, Sg_i)$  and its destination as  $(D_i, Dg_i)$ , here  $S_i$  and  $D_i$  represent the ID of the streets and  $Sg_i, Dg_i$  represent the ID of their segments. In [44], we propose an optimal solution for route planning problem that takes the real-time traffic into consideration. Integer linear programs are formulated to determine the fastest route given the current locations of vacant taxis and their assigned regions. The fastest paths can then be determined as the real-time traffic feed-back is obtained by the system. With the recurrent updates, ILP is solved regularly to determine the best routes according to the recent data, in other words, the route keeps updated as the new data is received. In order to reduce the complexity of the routing optimizer, we employ the recurrent Dijkstra's algorithm using the metrics evaluating the traffic level at each segment defined in [44,45]. In our approach, routes might be updated every 1 min. The detailed process is provided in Procedure 1. Note that the same routing approach is adopted to determine the trips of occupied taxis after pick-ups.

---

#### Procedure 1 Routing Optimizer for Taxi $i$

---

- 1: Inputs =  $\{(S_i, Sg_i), (D_i, Dg_i)\}$ , time instant  $t$ .
  - 2: **while** vehicle does not reach the destination **do**
  - 3:   Obtain the latest update traffic data based on the collected information.
  - 4:   Update the weights of the road network graph.
  - 5:   Run Dijkstra's algorithm to find the fastest route from  $(S_i, Sg_i)$  to  $(D_i, Dg_i)$ .
  - 6:   Vehicle follows the proposed route for one minute.
  - 7:   Update  $(S_i, Sg_i)$ .
  - 8: **end while**
- 

### 5. Simulated Taxi Operation and Validation

In this section, we introduce our framework to simulate the operation of taxis in the area of interest. Then, we validate the proposed model with realistic data to ensure that our simulations after determining routes are close to real-world situations.

#### 5.1. Simulation Model

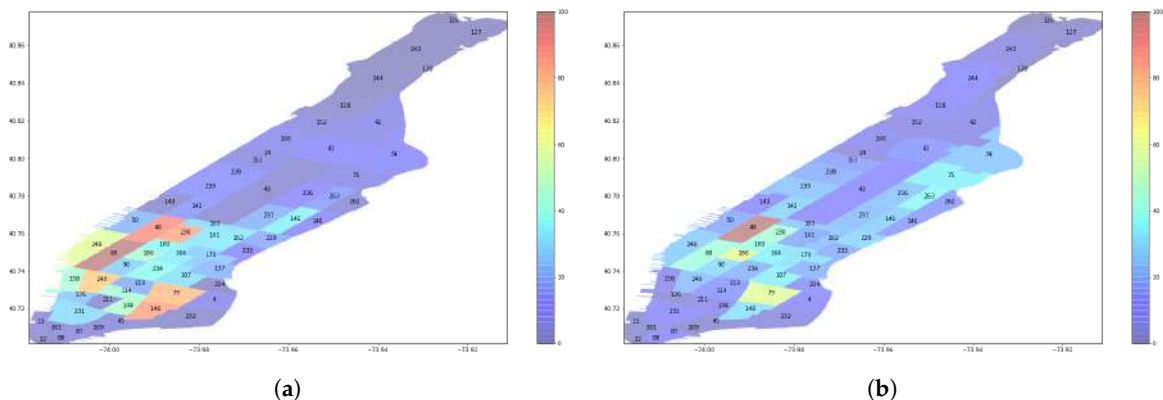
In our simulations, we consider the area of Manhattan, New York city, which is divided into 69 regions. We assume there are  $N$  taxis circulating in the area of interest. If it is vacant, we assume that the taxi picks up a customer when the distance separating them is less than 100 m. In the traditional system, where data exchange and knowledge about the customers' demand are absent, we consider that the  $N$  taxis move randomly in the whole area when they are vacant, while with the recommendation system, the taxis are always assigned to different regions at the beginning of the time period and will move randomly only within that region. Once a pick-up is made, the status of the taxi is changed to occupied until the customer is dropped off.

For the taxi routing optimization, we extract the parameters of the off-line map from Open Street Map [46]. In total, there are 9070 roads and 4146 intersections in the area of interest. We split each



road into segments having length at maximum 100 m. Thus, we obtain a graph of 11,760 edges and 6393 nodes.

Two scenarios are provided to strengthen the persuasive of the model. We consider the one hour demand information on 1 June 2018 from 3 am to 4 am that contains 1813 pick-ups in total (time instants and GPS locations) as the first scenario while the demand information on 1 January 2018 from 5 am to 6 am that contains 2027 pick-ups as the second scenario. We choose these two periods instead of rush hours for tractability and clarity reasons. Indeed, over rush hours, the number of pick-ups is huge and it will be difficult to visualize the results. This also impacts the simulation time, which is expected to be very expensive. Although we have developed low complexity algorithms for both the taxi-to-region matching component and the taxi routing optimizer, simulating the instantaneous operation of a huge number of taxis remains time consuming. It is worth noting that in our simulation results (Section 6) where we compare the different scenarios after simulating the taxi operations, we have investigated the same time periods where identical traffic conditions are experimented with. Since the customer arrival time and waiting time are missing in the dataset, without loss of generality, we assume that their arrival instants are the taxi pick-up times. Three key performance indicators are evaluated in our simulations: (1) the total number of pick-ups, (2) the waiting time of each customer corresponding to the difference between its pick-up time and its arrival time instants, and (3) the vacant traveled distance where no passengers are in the taxis. Precisely, the waiting time of customers corresponds to the period starting from the time instant when the customer arrives on the road for regular taxis or requests the service for ride-hailing taxis. The deadheading or idle distance of drivers is defined as the distance travelled by a taxi without serving any customers either before finding or after dropping a customer. All of these metrics are measured after simulating the taxi operation, as indicated in Section 5. The demands on 69 regions for both scenarios are presented in Figure 7. We notice that the demands mainly exist in mid and lower Manhattan. Although the two scenarios have similar total requests, their distributions in the regions are different. Customers in Scenario 2 are mainly located in regions 48, 68, 246, 230, 249, 79, 148, unlike Scenario 1 where most of them are gathered in regions 48, 186 and 79.



**Figure 7.** Demands percentage heat-map for 69 different regions on 1 January 2018 from 5 am to 6 am and 1 June 2018 from 3 am to 4 am separately. (a) Scenario 1. (b) Scenario 2.

The detailed algorithm to perform the simulations for regular taxi services without VSN is provided in Algorithm 1. Note that the recommendation occurs at the beginning of every time period  $T$  where  $T = 1$  h. Hence, the regions assigned to different taxis remain unchanged during this hour. For the next time period, the recommendation system updates its matching procedure for the vacant taxis according to their latest locations and the new demand.

The detailed algorithm to perform the simulations for regular taxi services with VSN is provided in Algorithm 2, where the recommendation occurs on the fly during the time period whenever a pickup is reported to the system. Here,  $N_{idle}(t)$  represents the number of vacant taxis at time instant  $t$ . In this algorithm, the system continuously provides recommendation during the time period  $T$  as

the number of pick-ups changes over time. Every  $\bar{t}$  minutes, the system sends the vacant vehicles to different regions considering the distance as well as the potential demand for the rest of the time period  $T$ . Note that within  $T$  the demand in the current step is highly correlated with the one of the next step. Hence, very few re-assignments will occur for vacant taxis.

---

**Algorithm 1** Simulated Taxi Operation for Regular Taxi Services Without VSN
 

---

```

1: Inputs =  $(S_i, Sg_i)$   $i \in \{1, \dots, N\}$ .
2: Determine the best assignment of taxi-to-region using the Hungarian method.
3: Send vacant taxis to recommended regions using the Routing Optimizer given in Procedure 1.
4:  $t = 0$ .
5: while  $t \leq T$  do
6:   for each Taxi  $i \in \{1, \dots, N\}$  do
7:     if Taxi  $i$  is vacant then
8:       Taxi  $i$  circulates towards or within the assigned region during this minute.
9:       Calculate the shortest distance  $d_{ik}$  between taxi  $i$  and potential nearby customers  $k$ 's.
10:      if  $\exists k$  such that  $d_{ik} < 100$  then
11:        Taxi  $i$  notices customer  $k$  waiving his/her hand and then heads to him/her.
12:        Record the waiting time of customer  $k$ .
13:        Change the status of taxi  $i$  to occupied.
14:      end if
15:    else
16:      Taxi  $i$  drives towards its destination as per customer request using the Routing Optimizer given in Procedure 1.
17:    end if
18:    Update  $(S_i, Sg_i)$ .
19:  end for
20:   $t = t + 1$ .
21: end while

```

---



---

**Algorithm 2** Simulated Taxi Operation for Regular Taxi Services With VSN
 

---

```

1: Inputs =  $(S_i, Sg_i)$ ,  $i \in \{1, \dots, N\}$ .
2:  $t = 0$ .
3: while  $t \leq T$  do
4:   if  $\text{mod}(t, \bar{t}) == 0$  then
5:     Update the demand by subtracting the pick-ups happened already.
6:     Find the vacant taxis  $i$ ,  $i \in \{1, \dots, N_{idle}\}$  at  $t$ .
7:     Determine the best assignment of taxi-to-region using the Hungarian method.
8:     Send vacant taxis to recommended regions the Routing Optimizer given in Procedure 1.
9:   end if
10:  for each Taxi  $i \in \{1, \dots, N\}$  do
11:    if Taxi  $i$  is vacant then
12:      Taxi  $i$  circulates towards or within the assigned region during this minute.
13:      Calculate the shortest distance  $d_{ik}$  between taxi  $i$  and potential nearby customers  $k$ 's.
14:      if  $\exists k$  such that  $d_{ik} < 100$  then
15:        Taxi  $i$  notices customer  $k$  waiving his/her hand and then heads to him/her.
16:        Record the waiting time of customer  $k$ .
17:        Change the status of taxi  $i$  to occupied.
18:      end if
19:    else
20:      Taxi  $i$  drives towards its destination as per customer request using the Routing Optimizer given in Procedure 1.
21:    end if
22:    Update  $(S_i, Sg_i)$ .
23:  end for
24:   $t = t + 1$ .
25: end while

```

---



Unlike the regular taxi services, the system for ride-hailing taxi services is aware of the locations for both taxis and customers' requests. Once a customer sends the request to the system, one of the nearby vacant taxis would head to him/her directly. Note that only the vacant vehicles that are within the search range  $R_g$  receive the request. In our simulation, we set the minimum search range  $R_g$  to 2 km. The detailed algorithm to perform the simulations is provided in Algorithm 3. Here, similar to the regular service with VSN, the system assigns the regions to vacant taxis every  $\bar{t}$  minutes as the demand for the rest of  $T$  is changing. However, taxis in ride-hailing services do not need to find customers waiving their hands on the street, in other words, the search range of taxis in ride-hailing service  $R_g$  is much larger than the regular taxi services. In our algorithm, we collect the location information of every customer and calculate their distance to all vacant vehicles within the search range. The closest available vehicle within that search region will be assigned to the customer. In our simulation, we set  $\bar{t} = 5$  min.

---

**Algorithm 3** Simulated Taxi Operation for Ride-Hailing Taxi Services
 

---

```

1: Inputs =  $(S_i, Sg_i)$ ,  $i \in \{1, \dots, N\}$ .
2:  $t = 0$ .
3: while  $t \leq T$  do
4:   if  $\text{mod}(t, \bar{t}) == 0$  then
5:     Update the demand by subtracting the pick-ups happened already.
6:     Find the vacant taxis  $n_i$ ,  $i \in \{1, \dots, N_{idle}\}$  at  $t$ .
7:     Determine the best assignment of taxi-to-region using the Hungarian method.
8:     Send vacant taxis to recommended regions using the Routing Optimizer given in Procedure 1.
9:   end if
10:  for each Customer  $k \in \{1, \dots, \mathcal{K}\}$  who shows up before  $t$  do
11:    Customer  $k$  sends its request and its location to the system.
12:    Calculate the shortest distance  $d_{ik}$  for customer  $k$  with all nearby taxis  $i$ ,  $i \in \{1, \dots, N_{idle}\}$ .
13:    Find the closest taxi  $i'$  and the shortest distance  $d_{i'k}$ .
14:    if  $d_{i'k} < R_g$  then
15:      Taxi  $i'$  heads to the customer  $k$  to pick him/her up.
16:    end if
17:    Record the waiting time of customer  $k$ .
18:    Change the status of taxi  $i'$  to occupied.
19:  end for
20:  for each Taxi  $i \in \{1, \dots, N\}$  do
21:    if Taxi  $i$  is vacant then
22:      Taxi  $i$  cruises towards or within the assigned region during this minutes.
23:    else
24:      Taxi  $i$  drives towards its destination as per customer request using the Routing Optimizer given in Procedure 1.
25:    end if
26:    Update  $(S_i, Sg_i)$ .
27:  end for
28:   $t = t + 1$ .
29: end while

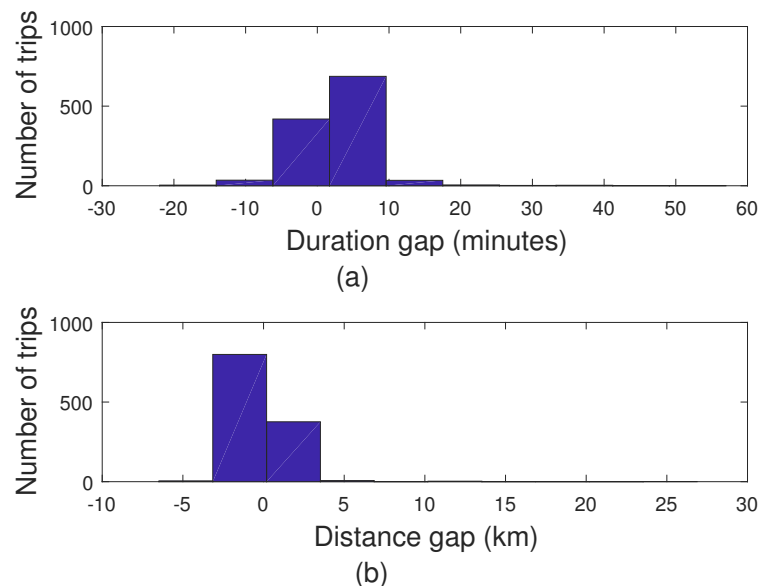
```

---

## 5.2. Model Validation

Figure 8, we propose to compare the simulation results with the current data to prove the efficiency of the model with respect to real-world scenarios. In the figure, we provide two histograms comparing the gap in terms of trip duration and traveled distance between actual data and simulated data for the different trips. From Figure 8a, we find that the majority of the simulated trips have duration close to the real data with a difference ranging from  $-3$  to  $5$  min. This is due to the difference between the true traffic status and the simulated one, as well as different drivers' routing preferences, that do not necessarily follow the obtained paths using the approach presented in Section 4. The difference is not huge since the average is close to 0. Moreover, from Figure 8b, we find out that distance differences

of the majority of trips are plus or minus 2.5 km from the realistic data since the available dataset only contains the pick-up and drop-off region ID without specifying the exact geographical points. In Figure 8 shows that the simulated model is very close to the real-world case and validates the system model and routing optimization algorithms that we developed.



**Figure 8.** Comparison between the actual and simulated data: (a) histogram representing the trip duration and (b) histogram representing the distance gap.

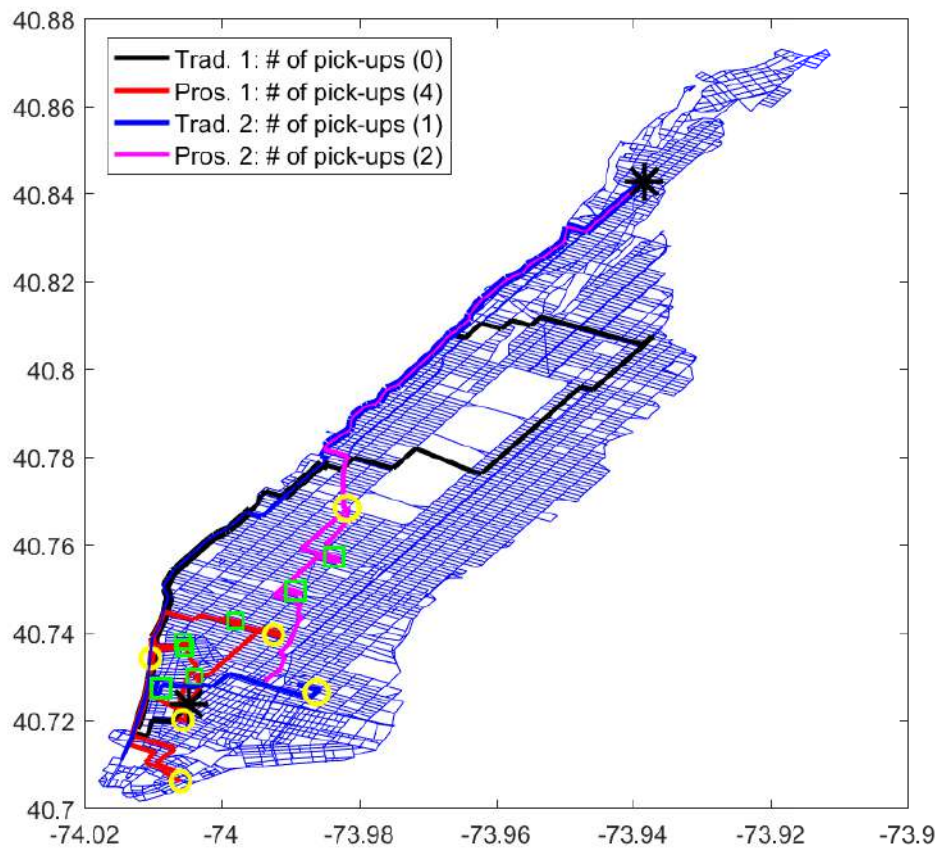
## 6. Performance Evaluation of the Proposed Recommendation System

In this section, we evaluate the performance of our proposed recommendation system and compare it to the traditional case where taxi drivers work individually and based on their own experience for both regular and ride-hailing taxi services. To sum-up, in our simulations, we compare five cases: Two traditional taxi services without recommendations (Regular Trad. and Ride-hailing Trad.) and three taxi services based on our proposed recommendation system (Regular Recom. (w/o VSN), Regular Recom. (VSN), and Ride-hailing Recom.). We start by providing a detailed analysis of the key performance metrics for Scenario 1, followed by a short discussion about Scenario 2.

### 6.1. Taxi Operation Visualization (Scenario 1)

In Figure 9, we illustrate an example of two selected taxis circulating in the area of interest while considering Scenario 1 (Figure 7a) for regular taxi services without VSN. Two of them, i.e., black and blue trajectories, are moving randomly looking for customers based on their own experience while two other taxis, colored in red and pink, follow the recommendations of the proposed system using Algorithm 1. The starting positions of the black and red taxis as well as the blue and pink taxis are the same, and by comparing the number of pick-ups between those two groups of taxis, we find out that the number of pick-ups increases when the recommendation system is applied. On the other hand, we can notice that the red vehicle spends most of its time cruising within the same region compared to the black vehicle and thus has a greater chance to find customers with lower vacant traveled distance. The starting position of the pink vehicle has lower number of potential customers so it is assigned to other regions that have higher probability to find customers.

In the sequel, we evaluate the performances of the proposed taxi recommendation systems for both regular and ride-hailing taxi services and compare them to the ones of the traditional cases.

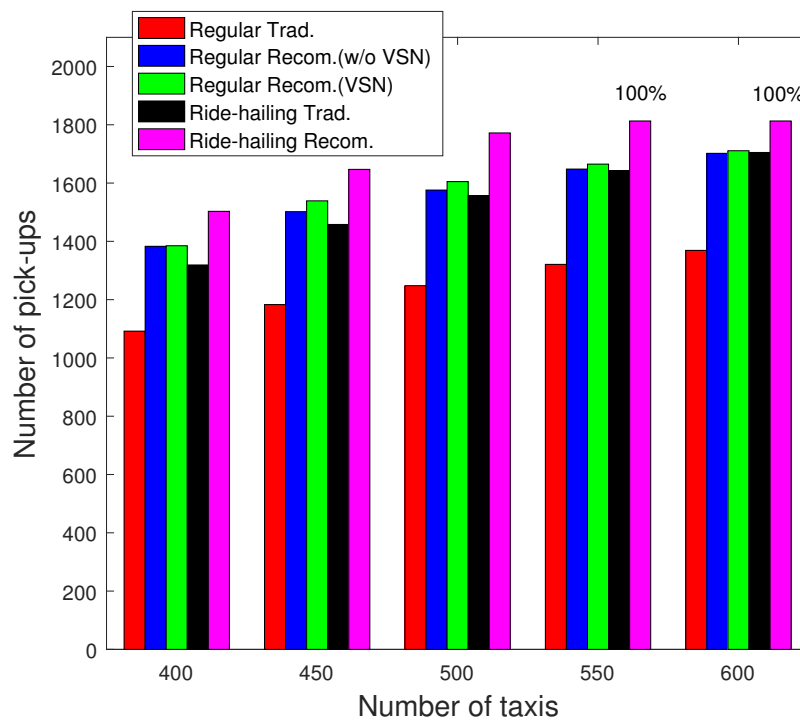


**Figure 9.** Example of two selected taxis circulating in the area of interest using the traditional and recommendation-based modes for regular taxi services without Vehicular Social Network (VSN). The 'black' and 'blue' trajectories correspond to two taxis moving in a traditional manner. The 'red' and 'pink' trajectories are of the same taxis following the recommendation system instructions (Circles (yellow) = drop off, squares (green) = pick-up locations).

## 6.2. Number of Pick-Ups (Scenario 1)

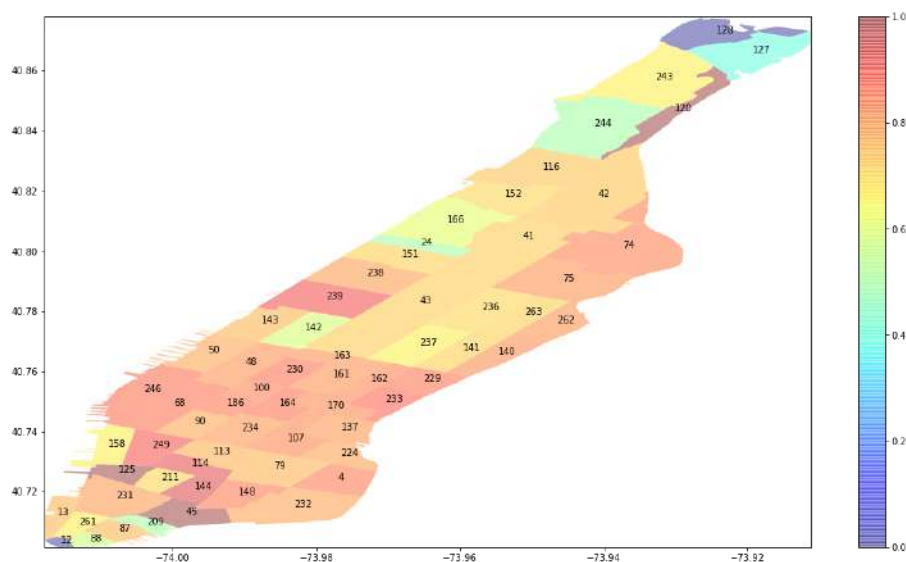
In Figure 10, we depict the number of pick-ups achieved by recommendation system in regular (without and with VSN) and ride-hailing taxi services (blue, green, pink) and compare them to the traditional cases of both services (red, black) with different taxi fleet sizes ( $N = \{400, 450, 500, 550, 600\}$ ) for Scenario 1. We can clearly notice that higher performance are achieved with the recommendation system regardless of the taxi fleet size.

For instance, the number of pick-ups with regular taxis increases by around 20% with a fleet size  $N = 450$ . Adding VSN option also helps in slightly improving the performance. On the other hand, the number of pick-ups in ride-hailing service is higher than those of regular service. For instance, when  $N = 600$ , with the recommendation system, 100% of the customers have been picked up using ride-hailing taxi service while 95% of the customers have been picked up using regular taxi service. Also, we notice that when  $N = 550$  and  $N = 600$ , the ride-hailing taxis are able to pick up all the customers. In other words, an excess supply is obtained with a taxi fleet of  $N = 600$ , which corresponds to an unnecessary wasting of fuel and may cause redundant congestion.



**Figure 10.** Number of pick-ups with different taxi fleet sizes for Scenario 1.

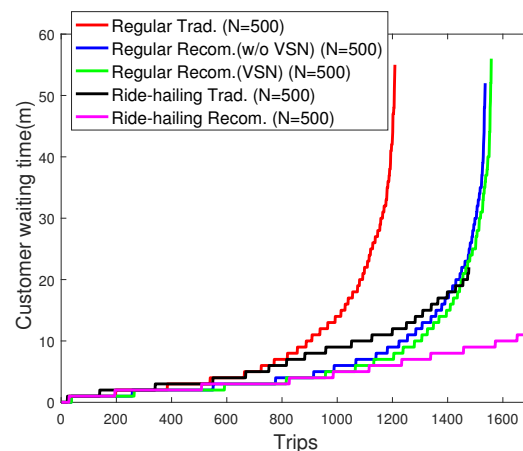
In order to deeply visualize the number of pick-ups for each region, we plot, in Figure 11, the ratio of number of pick-ups to the total customer's requests with  $N = 500$ . We notice that the ratio of pick-ups is small in the upper Manhattan since the customers' demands are mainly located in the lower Manhattan. Although we evenly assign the vehicles to different regions according to the expected customers' requests, there is a high probability that a vehicle heading to upper Manhattan from lower Manhattan ends up finding customers on the route before arriving.



**Figure 11.** Ratio of pickups with  $N = 500$  using the Ride-hailing recommendation system.

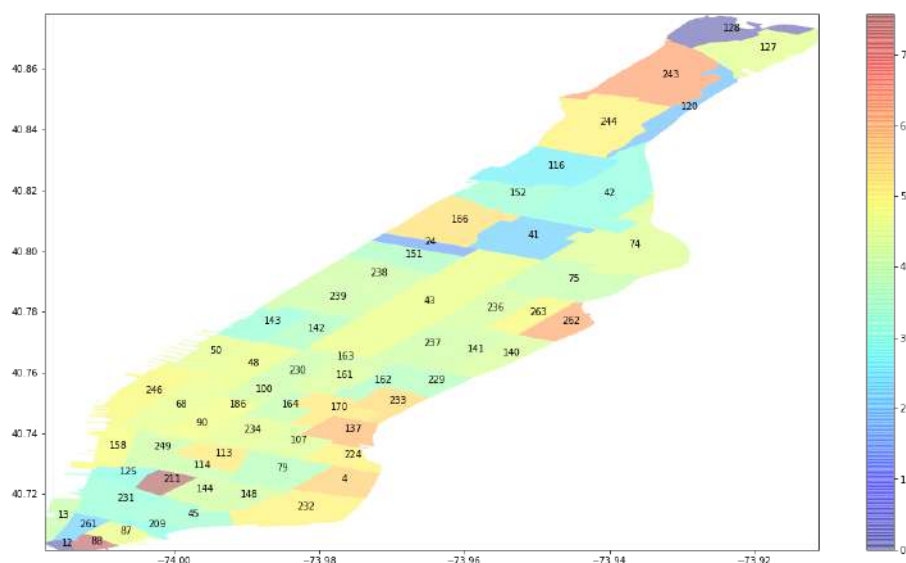
### 6.3. Customer Waiting Time (Scenario 1)

In Figure 12, we evaluate the satisfaction of customers (waiting time) for all the recorded trips during the time period  $T$  with  $N = 500$  for Scenario 1. We notice that, with the same fleet size of taxis cruising throughout the map, customers wait much less time with the recommendation system. With VSN, the performance of the recommendation system in regular taxi service is improved slightly. It is worth noting that 98% of the customers wait less than 10 min before finding a vacant ride-hailing taxi thanks to the proposed recommendation system compared to 70% with the traditional case. We also notice that without recommendation system, the ride-hailing service provides customers with shorter average waiting time compared to the regular taxi services, which is true in practice. If we apply the recommendation system for both services, then the average waiting time in ride-hailing is much lower than the one obtained with regular taxis.



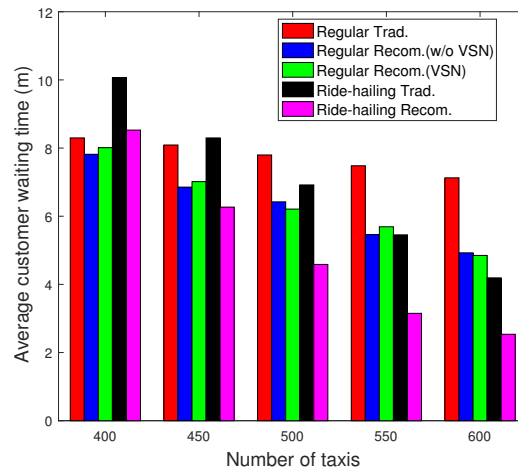
**Figure 12.** Traditional versus recommendation systems: sorted customer waiting time.

In addition, we present the average customer waiting time per region with  $N = 500$  for Scenario 1, as shown in Figure 13. We find out the average waiting time in upper Manhattan is lower than that in lower Manhattan, since the customers are gathered in lower Manhattan and there is competition among customers to find available taxis.



**Figure 13.** Average customer waiting time with  $N = 500$  using Ride-hailing recom.

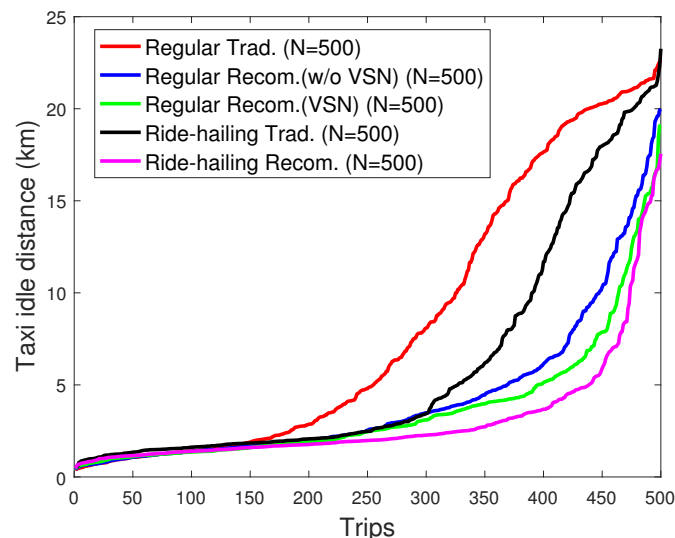
Finally, in Figure 14, we depict the average customer waiting time including recommendation system in regular and ride-hailing taxi services (blue, green, pink) and traditional case using both services (red, black) with different taxi fleet size ( $N = \{400, 450, 500, 550, 600\}$ ) for Scenario 1. We notice that higher performance is achieved with the recommendation system regardless of the taxi fleet size. For instance, when  $N = 600$ , on average, with the recommendation system, customers wait 1.66 min lower compared to the traditional case in ride-hailing taxi service and 2.28 min lower for regular taxi service.



**Figure 14.** Average customer waiting time with different taxi fleet sizes for Scenario 1.

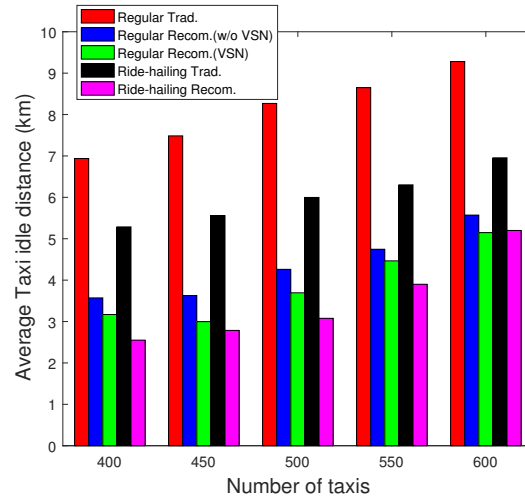
#### 6.4. Vacant Traveled Distance (Scenario 1)

Similarly, in Figure 15, we evaluate the satisfaction of taxi drivers represented by their idle traveled distance. We notice that with the proposed recommendation system, the taxi drivers have less idle traveled distance, and as expected, drivers in ride-hailing service have less idle traveled distance than those of the regular taxi service. We also notice that with VSN, the performance slightly increases in regular taxi services. It is worth noting that 92% of the taxis idly travel for less than 10 km during one hour when  $N = 500$  thanks to the proposed recommendation system. However, with the traditional techniques, only 78% of the fleet achieves a similar result.



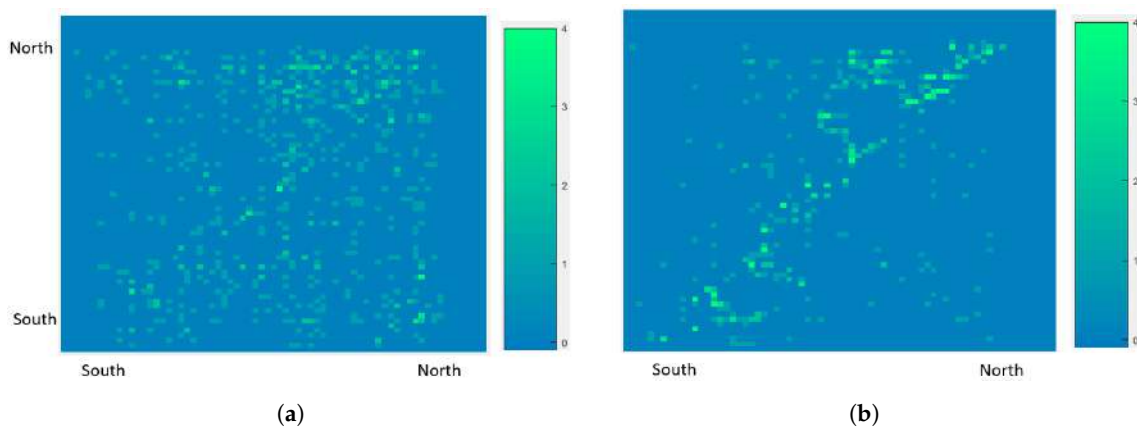
**Figure 15.** Traditional versus recommendation systems: sorted idle traveled distance.

In addition, we present, in Figure 16, the average idle traveled distance using the recommendation system for regular and ride-hailing taxi services, as well as the traditional cases for Scenario 1. Again, higher performances are achieved with the recommendation system regardless of the taxi fleet size. For instance, when  $N = 600$ , with the recommendation system, taxis travel 1.7 km less in vacant status compared to the traditional case in ride-hailing taxi service and 4.1 km less compared to the traditional case in regular taxi service. Close performances are achieved with the other fleet sizes.



**Figure 16.** Average idle traveled distance with different taxi fleet sizes for Scenario 1.

In Figure 17, we provide heatmaps for both traditional and proposed schemes illustrating the regions crossed by  $N = 500$  regular taxicabs during idle periods when looking for customers. In this figure, we sort the regions according to their geographical locations and place the regions next to each other in both axes where the horizontal axis is the origin region (last dropoff region) and the vertical axis is the destination region (the region where the next customer is found). The figure shows that the taxis in the traditional case are moving from a region to another in a near-uniform pattern where taxi drivers search for customers following their own intuition, while thanks to the recommendation system, taxi drivers are able to reduce their travelling idle distance by searching for customers within the same or nearby regions as it is corroborated by the diagonal pattern given in Figure 17b. In this way, the hunting time is minimized, which allows taxi drivers to save additional time and fuel.



**Figure 17.** Movement of vacant taxis from a region to another before finding new customers for regular taxicabs. The x-axis represents the origin regions while the y-axis represents the destination. Regions are sorted according to their geographical locations, in Manhattan area, from South to North. (a) Traditional taxi service. (b) Proposed recommendation system.



### 6.5. Taxi Re-Assignment Frequency (Scenario 1)

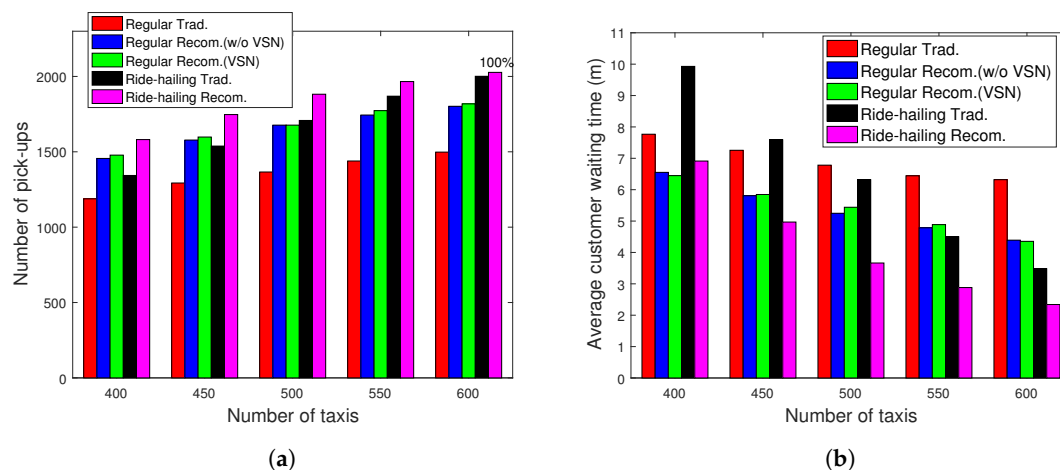
We have shown that our recommendation system could achieve outstanding progress for the different key metrics. We then explore whether the system (with VSN) provides excessive re-assignments to taxis during time period  $T$  or not and, hence, check the practicality of the system. In other words, we need to pay attention to the number of re-assignments since taxi drivers do not prefer such recommendations in practice. In Table 2, we provide the average number of re-assignments for Scenario 1 based on our simulations. On average, the number of re-assignments is lower than 2 during one hour for vacant taxis. On the other hand, ride-hailing taxis have less re-assignments compared to the those of regular taxi services since the locations of customers are known in ride-hailing taxi services. Also, we notice that when the number of taxis increase, taxi drivers are more likely re-assigned before finding customers since the supply is exceeding the demand.

**Table 2.** Average taxis re-assignment frequency using Regular Recom. (with VSN) and Ride-hailing Recom.

	Regular Recom. (VSN)	Ride-Hailing Recom.
N = 400	1.62	0.59
N = 450	1.64	1.08
N = 500	1.77	1.09
N = 550	1.83	1.19
N = 600	1.96	1.37

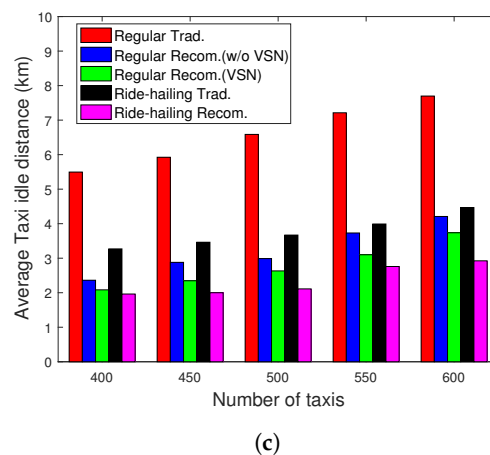
### 6.6. Summary and Discussion for Scenario 2

Finally, a comprehensive summary of the results for Scenario 2 is shown in Figure 18, which presents similar performance to Scenario 1. For instance, when  $N = 600$ , customers wait 1.62 min lower compared to the traditional case in ride-hailing taxi service and 2.79 min lower for regular taxi service. Also, with the recommendation system, taxis travel 1.8 km less in vacant status compared to the traditional case in ride-hailing taxi service and 3.9 km less compared to the traditional case in regular taxi service. On the other hand, by comparing the performance between ride-hailing and regular taxi services, we find out the average customer waiting time and the idle traveled distance of drivers are improved when customers' locations are sent to the system. Furthermore, it is worth noting that when the number of taxis increases, the customers' waiting time decreases while the idle traveled distance of taxi drivers increase. It is important to determine the appropriate size of taxi fleet for each time period of the day so that both customer and taxi drivers are satisfied without overloading the region with redundant taxis.



**Figure 18.** Cont.





**Figure 18.** System performance for the three major taxi services for Scenario 2 with 2027 customers in total. (a) Number of pick-ups. (b) Average customer waiting time. (c) Average idle traveled distance.

## 7. Conclusions

In this paper, we have designed and validated an effective recommendation system for three main taxi services: regular (without VSN), regular (with VSN) and ride-hailing taxi services. The system includes three major components: an incremental predictor of future demands, taxi-to-region matching component, and taxi routing optimizer. By comparing the performance of the proposed recommendation services to the ones of the traditional cases, we found that the proposed approach achieves significant gains in terms of pick-ups efficiency, time, and energy saving for both customers and taxis. The proposed framework can be used as an effective tool for different taxi services by exploiting the power of vehicular social networks and data sharing taxi drivers. Enabling timely and erroneous information exchange of the automatic sensing crowd-sourcing framework will be the scope of our future work in order to ensure efficient operation of the proposed recommendation system.

**Author Contributions:** Conceptualization, Methodology, Validation, Data Analysis, X.W. and H.G.; Supervision, Project Administration, Funding Acquisition and Writing—review and editing by H.G. and Y.M. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research received no external funding.

**Acknowledgments:** The authors would like to acknowledge the NYC DOT government for their open-access traffic data.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. NYCEDC's Blog. New Yorkers and Their Cars. Technical Report. April 2018. Available online: <https://edc.nyc/article/new-yorkers-and-their-cars> (accessed on 14 April 2020).
2. Roth, S. Report: Taxis Have Longer Wait Times than Uber, Lyft. Technical Report. July 2015. Available online: [https://www.google.com.hk/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=2ahUKEwj0i9mxzOnoAhVKQd4KHTjMBowQFjAAegQIAhAB&url=https%3A%2F%2Fwww.its.dot.gov%2Fitspac%2Fdec2014%2Fridesourcingwhitepaper\\_nov2014.pdf&usg=AOvVaw2BIBGvWSCTY3plQN BaxA4d](https://www.google.com.hk/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=2ahUKEwj0i9mxzOnoAhVKQd4KHTjMBowQFjAAegQIAhAB&url=https%3A%2F%2Fwww.its.dot.gov%2Fitspac%2Fdec2014%2Fridesourcingwhitepaper_nov2014.pdf&usg=AOvVaw2BIBGvWSCTY3plQN BaxA4d) (accessed on 14 April 2020).
3. Kong, X.; Xia, F.; Ning, Z.; Rahim, A.; Cai, Y.; Gao, Z.; Ma, J. Mobility Dataset Generation for Vehicular Social Networks Based on Floating Car Data. *IEEE Trans. Veh. Technol.* **2018**, *67*, 3874–3886, doi:10.1109/TVT.2017.2788441. [CrossRef]
4. Zhao, Y.; Han, Q. Spatial crowdsourcing: Current state and future directions. *IEEE Commun. Mag.* **2016**, *54*, 102–107. [CrossRef]

5. Ning, Z.; Xia, F.; Ullah, N.; Kong, X.; Hu, X. Vehicular Social Networks: Enabling Smart Mobility. *IEEE Commun. Mag.* **2017**, *55*, 16–55. [CrossRef]
6. Vegni, A.M.; Loscrí, V. A Survey on Vehicular Social Networks. *IEEE Commun. Surv. Tutor.* **2015**, *17*, 2397–2419. [CrossRef]
7. Uhlemann, E. Connected-Vehicles Applications Are Emerging [Connected Vehicles]. *IEEE Veh. Technol. Mag.* **2016**, *11*, 25–96. [CrossRef]
8. Ning, Z.; Hu, X.; Chen, Z.; Zhou, M.; Hu, B.; Cheng, J.; Obaidat, M.S. A Cooperative Quality-Aware Service Access System for Social Internet of Vehicles. *IEEE Internet Things J.* **2018**, *5*, 2506–2517, doi:10.1109/JIOT.2017.2764259. [CrossRef]
9. Lucic, M.C.; Ghazzai, H.; Khattab, A.; Massoud, Y. Rapid Management of Unexpected Events in Urban V2I Communications Systems. In Proceedings of the 2019 IEEE International Conference on Vehicular Electronics and Safety (ICVES), Cairo, Egypt, 4–6 September 2019; pp. 1–6, doi:10.1109/ICVES.2019.8906309. [CrossRef]
10. Niu, B.; Li, Q.; Zhu, X.; Cao, G.; Li, H. Achieving K-anonymity in privacy-aware location-based services. In Proceedings of the IEEE Conference on Computer Communications (INFOCOM 2014), Toronto, ON, Canada, 27 April–2 May 2014.
11. Shao, J.; Lu, R.; Lin, X. FINE: A fine-grained privacy-preserving location-based service framework for mobile devices. In Proceedings of the IEEE Conference on Computer Communications (INFOCOM 2014), Toronto, ON, Canada, 27 April–2 May 2014.
12. Yuan, N.J.; Zheng, Y.; Zhang, L.; Xie, X. T-Finder: A Recommender System for Finding Passengers and Vacant Taxis. *IEEE Trans. Knowl. Data Eng.* **2013**, *25*, 2390–2403, doi:10.1109/TKDE.2012.153. [CrossRef]
13. Wan, X.; Ghazzai, H.; Massoud, Y. Online Recommendation System for Autonomous and Human-driven Ride-hailing Taxi Services. In Proceedings of the 2019 31st International Conference on Microelectronics (ICM), Cairo, Egypt, 15–18 December 2019; pp. 351–354.
14. Wan, X.; Ghazzai, H.; Massoud, Y. Incremental Recommendation System for Large-scale Taxi Fleet in Smart Cities. In Proceedings of the 2019 IEEE International Conference on Vehicular Electronics and Safety (ICVES), Cairo, Egypt, 4–6 September 2019; pp. 1–6, doi:10.1109/ICVES.2019.8906451. [CrossRef]
15. Kong, X.; Xia, F.; Wang, J.; Rahim, A.; Das, S.K. Time-Location-Relationship Combined Service Recommendation Based on Taxi Trajectory Data. *IEEE Trans. Ind. Inform.* **2017**, *13*, 1202–1212, doi:10.1109/TII.2017.2684163. [CrossRef]
16. Verma, T.; Varakantham, P.; Kraus, S.; Lau, H.C. Augmenting decisions of taxi drivers through reinforcement learning for improving revenues. In Proceedings of the International Conference on Automated Planning Scheduling (ICAPS'17), Pittsburgh, PA, USA, 18–23 June 2017.
17. Moreira-Matias, L.; Gama, J.; Ferreira, M.; Mendes-Moreira, J.; Damas, L. Predicting Taxi–Passenger Demand Using Streaming Data. *IEEE Trans. Intell. Transp. Syst.* **2013**, *14*, 1393–1402, doi:10.1109/TITS.2013.2262376. [CrossRef]
18. Li, X.; Pan, G.; Wu, Z.; Qi, G.; Li, S.; Zhang, D.; Zhang, W.; Wang, Z. Prediction of urban human mobility using large-scale taxi traces and its applications. *Front. Comput. Sci.* **2012**, *6*, 111–121.
19. Xu, X.; Zhou, J.; Liu, Y.; Xu, Z.; Zhao, X. Taxi-RS: Taxi-Hunting Recommendation System Based on Taxi GPS Data. *IEEE Trans. Intell. Transp. Syst.* **2015**, *16*, 1716–1727, doi:10.1109/TITS.2014.2371815. [CrossRef]
20. Yuan, J.; Zheng, Y.; Zhang, L.; Xie, X.; Sun, G. Where to find my next passenger. In Proceedings of the ACM International Conference on Ubiquitous Computing, Beijing, China, 17–21 September 2011.
21. Garg, N.; Ranu, S. Route Recommendations for Idle Taxi Drivers: Find Me the Shortest Route to a Customer! In Proceedings of the ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD'19), London, UK, 19–23 August 2018.
22. Wang, H.; Yang, H. Ridesourcing systems: A framework and review. *Transp. Res. Part B Methodol.* **2019**, *129*, 122–155. [CrossRef]
23. Alonso-Mora, J.; Samaranayake, S.; Wallar, A.; Frazzoli, E.; Rus, D. On-demand high-capacity ride-sharing via dynamic trip-vehicle assignment. *Proc. Natl. Acad. Sci. USA* **2017**, *114*, 462–467. [CrossRef]
24. Simonetto, A.; Monteil, J.; Gambella, C. Real-time city-scale ridesharing via linear assignment problems. *Transp. Res. Part C Emerg. Technol.* **2019**, *101*, 208–232. [CrossRef]

25. Yao, H.; Wu, F.; Ke, J.; Tang, X.; Jia, Y.; Lu, S.; Gong, P.; Ye, J.; Li, Z. Deep multi-view spatial-temporal network for taxi demand prediction. In Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, LA, USA, 2–7 February 2018.
26. Ke, J.; Yang, H.; Zheng, H.; Chen, X.; Jia, Y.; Gong, P.; Ye, J. Hexagon-Based Convolutional Neural Network for Supply-Demand Forecasting of Ride-Sourcing Services. *IEEE Trans. Intell. Transp. Syst.* **2019**, *20*, 4160–4173, doi:10.1109/TITS.2018.2882861. [CrossRef]
27. Wang, X.; Zhang, H.; Wang, L.; Ning, Z. A Demand-Supply Oriented Taxi Recommendation System for Vehicular Social Networks. *IEEE Access* **2018**, *6*, 41529–41538, doi:10.1109/ACCESS.2018.2857002. [CrossRef]
28. Miao, F.; Han, S.; Lin, S.; Stankovic, J.A.; Zhang, D.; Munir, S.; Huang, H.; He, T.; Pappas, G.J. Taxi Dispatch With Real-Time Sensing Data in Metropolitan Areas: A Receding Horizon Control Approach. *IEEE Trans. Autom. Sci. Eng.* **2016**, *13*, 463–478. doi:10.1109/TASE.2016.2529580. [CrossRef]
29. Buchholz, N. Spatial Equilibrium, Search Frictions and Efficient Regulation in the Taxi Industry. Working Paper; Technical Report. December 2017. Available online: [https://www.google.com.hk/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=2ahUKEwixg-DrzOnoAhXSP3AKHdq2Ct8QFjAAegQIAhAB&url=https%3A%2F%2Fscholar.princeton.edu%2Fsites%2Fdefault%2Ffiles%2Fnbuchholz%2Ffiles%2Ftaxi\\_draft.pdf&usg=AOvVaw3fPzZJAw6mn13xqbuotuLA](https://www.google.com.hk/url?sa=t&rct=j&q=&esrc=s&source=web&cd=1&cad=rja&uact=8&ved=2ahUKEwixg-DrzOnoAhXSP3AKHdq2Ct8QFjAAegQIAhAB&url=https%3A%2F%2Fscholar.princeton.edu%2Fsites%2Fdefault%2Ffiles%2Fnbuchholz%2Ffiles%2Ftaxi_draft.pdf&usg=AOvVaw3fPzZJAw6mn13xqbuotuLA) (accessed on 14 April 2020).
30. Lowalekar, M.; Varakantham, P.; Jaillet, P. Online spatio-temporal matching in stochastic and dynamic domains. *Artif. Intell.* **2018**, *261*, 71–112. [CrossRef]
31. Ota, M.; Vo, H.; Silva, C.; Freire, J. A scalable approach for data-driven taxi ride-sharing simulation. In Proceedings of the 2015 IEEE International Conference on Big Data (Big Data), Santa Clara, CA, USA, 29 October–1 November 2015; pp. 888–897, doi:10.1109/BigData.2015.7363837. [CrossRef]
32. Liu, S.; Wang, S.; Liu, C.; Krishnan, R. Understanding taxi drivers’ routing choices from spatial and social traces. *Front. Comput. Sci.* **2015**, *9*, 200–209. [CrossRef]
33. Zhang, D.; He, T. pCruise: Reducing Cruising Miles for Taxicab Networks. In Proceedings of the 2012 IEEE 33rd Real-Time Systems Symposium, San Juan, PR, USA, 4–7 December 2012.
34. Liu, S.; Yue, Y.; Krishnan, R. Non-Myopic Adaptive Route Planning in Uncertain Congestion Environments. *IEEE Trans. Knowl. Data Eng.* **2015**, *27*, 2438–2451, doi:10.1109/TKDE.2015.2411278. [CrossRef]
35. Powell, J.W.; Huang, Y.; Bastani, F.; Ji, M. Towards reducing taxicab cruising time using spatio-temporal profitability maps. In Proceedings of the International Symposium on Spatial and Temporal Databases, Minneapolis, MN, USA, 24–26 August 2011; Springer: Berlin/Heidelberg, Germany, 2011; pp. 242–260.
36. Sayarshad, H.R.; Chow, J.Y. Non-myopic relocation of idle mobility-on-demand vehicles as a dynamic location-allocation-queueing problem. *Transp. Res. Part E Logist. Transp. Rev.* **2017**, *106*, 60–77. [CrossRef]
37. Nourinejad, M.; Ramezani, M. Ride-Sourcing modeling and pricing in non-equilibrium two-sided markets. *Transp. Res. Part B Methodol.* **2020**, *132*, 340–357. [CrossRef]
38. Chen, X.; Miao, F.; Pappas, G.J.; Preciado, V. Hierarchical data-driven vehicle dispatch and ride-sharing. In Proceedings of the 2017 IEEE 56th Annual Conference on Decision and Control (CDC), Melbourne, Australia, 12–15 December 2017; pp. 4458–4463.
39. Ma, T.Y.; Rasulkhani, S.; Chow, J.Y.; Klein, S. A dynamic ridesharing dispatch and idle vehicle repositioning strategy with integrated transit transfers. *Transp. Res. Part E Logist. Transp. Rev.* **2019**, *128*, 417–442. [CrossRef]
40. Tafreshian, A.; Masoud, N. Trip-based graph partitioning in dynamic ridesharing. *Transp. Res. Part C Emerg. Technol.* **2020**, *114*, 532–553. [CrossRef]
41. Pandey, V.; Monteil, J.; Gambella, C.; Simonetto, A. On the needs for MaaS platforms to handle competition in ridesharing mobility. *Transp. Res. Part C Emerg. Technol.* **2019**, *108*, 269–288. [CrossRef]
42. Quddus, M.A.; Ochieng, W.Y.; Noland, R.B. Current map-matching algorithms for transport applications: State-of-the art and future research directions. *Transp. Res. Part C Emerg. Technol.* **2007**, *15*, 312–328. [CrossRef]
43. Quddus, M.A. High Integrity Map Matching Algorithms for Advanced Transport Telematics Applications. Ph.D. Thesis, Imperial College London, London, UK, January 2006.
44. Wan, X.; Ghazzai, H.; Massoud, Y. Real-Time Navigation in Urban Areas Using Mobile Crowd-Sourced Data. In Proceedings of the IEEE International Systems Conference (SYSCON’19), Orlando, FL, USA, 8–11 April 2019.

45. Wan, X.; Ghazzai, H.; Massoud, Y. Mobile Crowdsourcing for Intelligent Transportation Systems: Real-Time Navigation in Urban Areas. *IEEE Access* **2019**, *7*, 136995–137009, doi:10.1109/ACCESS.2019.2942282. [CrossRef]
46. Boeing, G. OSMnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks. *Comput. Environ. Urban Syst.* **2017**, *65*, 126–139. [CrossRef]



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).

## Article

# Using a Hybrid Recommending System for Learning Videos in Flipped Classrooms and MOOCs

Jaume Jordán <sup>1,\*</sup>, Soledad Valero <sup>1,†</sup>, Carlos Turró <sup>2,†</sup> and Vicent Botti <sup>1,†</sup>

<sup>1</sup> Valencian Research Institute for Artificial Intelligence (VRAIN), Universitat Politècnica de València, Camino de Vera s/n, 46022 Valencia, Spain; svalero@dsic.upv.es (S.V.); vbotti@dsic.upv.es (V.B.)

<sup>2</sup> Área de Sistemas de Información y Comunicaciones, Universitat Politècnica de València, Camino de Vera s/n, 46022 Valencia, Spain; turro@cc.upv.es

\* Correspondence: jjordan@dsic.upv.es; Tel.: +34-963-877-000

† These authors contributed equally to this work.

**Abstract:** New challenges in education require new ways of education. Higher education has adapted to these new challenges by means of offering new types of training like massive online open courses and by updating their teaching methodology using novel approaches as flipped classrooms. These types of training have enabled universities to better adapt to the challenges posed by the pandemic. In addition, high quality learning objects are necessary for these new forms of education to be successful, with learning videos being the most common learning objects to provide theoretical concepts. This paper describes a new approach of a previously presented hybrid learning recommender system based on content-based techniques, which was capable of recommend useful videos to learners and lecturers from a learning video repository. In this new approach, the content-based techniques are also combined with a collaborative filtering module, which increases the probability of recommending relevant videos. This hybrid technique has been successfully applied to a real scenario in the central video repository of the Universitat Politècnica de València.

**Citation:** Jordán, J.; Valero, S.; Turró, C.; Botti, V. Using a Hybrid Recommending System for Learning Videos in Flipped Classrooms and MOOCs. *Electronics* **2021**, *10*, 1226. <https://doi.org/10.3390/electronics10111226>

Received: 15 April 2021

Accepted: 17 May 2021

Published: 21 May 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

**Keywords:** learning recommender system; learning object; learning videos; content-based; collaborative filtering

## 1. Introduction

New challenges in education has raised due to the students' profile changes in the last decade. They demand new ways of learning, better adapted to their way of life and moving away from classical teaching. Academic institutions must be agile in adapting their teaching methodology to the new forms required, taking into account the opportunities offered by the global world [1,2]. In this way, there has been a great increase in the supply of massive online open courses (MOOCs) by academic institutions, as well as in the number of students opting for this type of training [3]. MOOCs mainly relay on learning objects (LOs). As IEEE proposes, a LO is "any entity, digital or non-digital, which can be used, re-used or referenced during technology supported learning" [4]. Thus, MOOCs use different types of LOs, with videos being one of the most commonly used to teach theoretical concepts. In addition, new teaching approaches are emerging in higher education, such as flipped teaching [5–7], in which the theoretical content is studied at home by the students, while the face-to-face sessions are eminently practical, where the knowledge acquired is put into practice by solving problems. To this end, the lecturer instructs the students which LOs they should work on at home before the next face-to-face session. In this way, students are encouraged to acquire the theoretical concepts not only through books or specialised articles, but also through audio–visual material.

Furthermore, in this pandemic context, face-to-face classes have been replaced by online classes in many institutions, increasing the adoption of flipped learning. Lecturers need to plan subjects taking into account possible connectivity problems, as the possibility

of students being unable to attend online classes due to health problems (quarantine, hospitalisation) increases. In addition, students may have difficulties in accessing devices during working hours, because they share them with their parents, etc. All these new circumstances make it even more important to provide educators with useful tools to search for and recommend good LOs, which can be accessed by students at any time.

Universitat Politècnica de València (<http://www.upv.es>, accessed on 14 May 2021) (UPV) is a Spanish Public university which offers undergraduate degrees, dual degrees, masters and doctoral programs. UPV has more than 28,000 students. UPV has been promoting new pedagogical methodologies in their degrees in the last decade, such as flipped teaching [8]. It has also made a great effort in developing MOOCs within the edX platform (<https://www.edx.org/>, accessed on 14 May 2021), with more than 2 million enrollments and having three courses in Class Central all time top 100 MOOCs in 2019 (<https://www.classcentral.com/report/top-moocs-2019-edition/>, accessed on 14 May 2021), and two courses top 30 in 2020 (<https://www.classcentral.com/report/best-free-online-courses-2021/>, accessed on 14 May 2021). In addition, UPV participated in the movement that arose during the first months of the pandemic to offer free certificates for some of the MOOCs offered. This fact, together with the need for new training channels, has led to a significant increase in the number of students on this type of courses (<https://www.classcentral.com/report/mooc-stats-2020/>, accessed on 14 May 2021).

Students need access to a variety of resources to understand the theoretical concepts required in blended and flipped classroom environments. To facilitate this difficult work, UPV has had a long-standing digital resources project with the aim of producing video content as LOs. This video content is handled in the university central video repository, called mediaUPV (<https://media.upv.es/>, accessed on 14 May 2021). This portal is not only used in the field of MOOCs, but also in other educational projects.

mediaUPV allows UPV lecturers to upload and manage video content for students. Students access mediaUPV usually through suggestions made by their lecturers through the learning management System (LMS), but they also access the video portal and browse through the content on their own. A relevant feature of mediaUPV against other alternatives (e.g., YouTube) is that the videos have been prepared and recorded by lecturers from the institution, so the quality of the content is guaranteed.

mediaUPV portal has become an essential tool to the institution during this pandemic; however, the size of the mediaUPV content is a growing problem, and hence, it is increasingly difficult to find the most relevant content for both students to view and lecturers to suggest. Thus, UPV determined that both students and lecturers would benefit from the development of a new Learning Recommendation System (LRS), that could recommend relevant and related videos. Therefore, in our previous work [9], we presented the first recommender engine proposed to carry out that purpose, which combined two content-based techniques to recommend useful learning videos to learners and lecturers. However, with that engine it was only possible to recommend videos labelled as high quality, that is, with transcripts. Therefore, it is necessary to apply other techniques to cover the entire mediaUPV catalogue.

In this paper we present an enhancement of the previous work, which extends the proposed recommendation engine also using collaborative filtering techniques. Thus, we describe how we have designed and developed a hybrid recommender system based on both content-based techniques and collaborative filtering. Furthermore, a complete analysis of the results in production of the initial content-based LRS since its application in October 2019 until March 2021 is provided.

This article is structured as follows. In Section 2 related works are described. Following, Section 3 specifies a description about the LOs to recommend and the potential users of the system. In Section 4, the proposed recommender system is explained. Then, Section 5 shows the experimental results of the proposed recommender using the data of the mediaUPV portal. In addition, Section 6 provides an analysis of the results of the LRS used in production. Finally, Section 7 draws the conclusions and future work of this paper.

## 2. Related Work

Learning Recommender Systems (LRS) should assist learners in discovering relevant LO than keep them motivated and enable them to complete their learning activities [10]. Most of the LRS adopt the same techniques than regular recommender systems [10–13], such as: content-based, in which recommendations are determined considering user profiles and content analysis of the learning objects already visited by the user; collaborative filtering, in which recommendations are based on the choices of other similar user profiles; knowledge-based, in which it is inferred whether a LO satisfies a particular learning need of the user to recommend it; and hybrid, in which recommendations are computed by combining more than one of the above techniques.

In recent years, different approaches have been proposed in order to improve the efficiency and accuracy of the recommendations and retrieval of useful LOs. In this way, in [14], authors provide new metrics for applying collaborative filtering in a learning domain, so users with better academic results have greater weight in the calculation of the recommendations. However, the experiments did not carried out in a learning environment. In another proposal, Zapata et al. provide a tool for filtering the retrieved results from a user query, which uses a combination of different filtering techniques, such as content comparison, and collaborative and demographic searches [15].

Other proposals focus on recommending to students those LOs that can be most useful to them, providing solid arguments. This is the case of [13], which combine content-based, collaborative and knowledge-based recommenders using an argumentation-based module to recommend LOs inside a LMS. In this case, information on student profiles and learning styles is also available. An item-based collaborative filtering method is combined with a sequential pattern mining algorithm to recommend LOs to learners in [16]. In this case, LOs are ranked by the students and it is also possible to obtain the browsing sequences made by them. In a similar way, ref. [17] proposes a hybrid knowledge-based recommender system based on ontology and sequential pattern mining for recommendation of LO. Authors can adequately characterize learners and LOs using an ontology, since they have detailed information about them. Personalized learning paths (sequence of LOs) that maximizes the performance of the learner and effectiveness of learning are provided in [18], hybridizing ant colony optimization with genetic algorithm. In this case, authors use both learners and LOs attributes to determine the appropriate learning paths. In addition, in [19], Dwivedi et al. recommend learning paths using a variable length genetic algorithm. This approach considers learners' learning styles and knowledge levels extracted from the learners' registration process. Finally, in [20], authors propose a method, based in a collaborative filtering approach, for building a unified learner profile which is used to recommend LOs to a group of individuals.

Besides, other works have been done to improve the accuracy of the searches in mediaUPV. For example, in [21] a semi-supervised method is applied to cluster and classify the LOs of mediaUPV, obtaining specific keywords that represent each cluster. In [22], authors applied a custom approach for indexing and retrieving educational videos using their transcripts, which are available in mediaUPV. Videos are classified in different domains using the method described in [21]. In addition, they applied a Latent Dirichlet Allocation algorithm [23] to get a list of topics and their score. User queries are classified in one of the domains, recovering from that cluster those videos whose transcripts are the closest to the query.

As can be seen, most previous work on LRS adopts a hybrid strategy, seeking to harness the strength of each particular technique, overcoming its limitations by using them together. Furthermore, the different strategies that can be applied depend on the data available to describe LOs and users. In our case, a hybrid strategy will also be applied, combining content-based methods. On the other hand, previous experiences on the improvement of searches in the mediaUPV repository show us the usefulness of characterizing LOs using available transcripts and titles.

### 3. Problem Description

mediaUPV portal started in 2011 and by the end of 2019 it had 55,600 different videos mainly from STEM topics, with more than 10 million views. During 2020 until the end of March 2021, the mediaUPV catalogue increased almost a 44.8%, reaching 80,500. This significant growth may be due to the increased need for the use of online learning resources, not only for new forms of learning, but also for the institution's traditional courses that had to be converted to purely online teaching during the pandemic. From this database, only 13,232 by the end of 2019, and 20,135 by the end of March 2021 are certified as high quality LOs (an increase of 52.2%). All these high quality videos have a transcript (with more than 100 characters), a title, and an author, however the videos are not classified using any taxonomy and no useful keywords are associated to them for all cases. mediaUPV platform generates the transcripts of what is said in the videos using the poli[Trans] service (<https://politrans.upv.es>, accessed on 14 May 2021), an online platform for automated and assisted multilingual media subtitling offered by UPV. poli[Trans] service is based on transLectures-UPV Platform [24].

Hence, a 75% of the mediaUPV catalogue is composed by videos without transcripts. It was detected that a huge percentage of that videos have a medium quality and could be also interesting to lecturers and students. In fact, some of them are recordings of classes from some subjects, which can be used for reinforcing the learning of students of similar subjects in other degrees. By default, mediaUPV does not obtain the transcripts of these recordings, because is the lecturer who has the permissions to demand them.

mediaUPV portal is mainly used by students and lecturers. The students are mainly formal students of the UPV, but it also receives many visits from anonymous users, who can register in some MOOC offered by the UPV. Moreover, mediaUPV is not connected to the LMS of the UPV, so even though the user is authenticated, it is not possible to know his student profile (e.g., enrolled subjects).

Therefore, the aim of this proposal is to be able to offer recommendations not only to authenticated users but also to anonymous users of the system. In addition, the system should be able to recommend not only students, but also lecturers who want to find quality videos to suggest to their students. Therefore, the main objective is to be able to offer recommendations on learning videos (that is LOs) that meet the quality standards set by the UPV.

It is important to note that all recommendations made on the mediaUPV portal are always through a video. That is, when a user logged into the system is watching a video (which can suggest videos similar to the current video or content that the user has watched previously), or an anonymous user watching a video (which can only suggest content similar to the current video). This implies that there are no recommendations only associated with a logged-in user. Thus, a hybrid system needs to be used in which the current video, and (optionally) the logged-in user are considered.

### 4. Learning Recommender System Proposal

Collaborative recommendation has some well-known difficulties, such as the necessity of a huge quantity of data for making accurate recommendations. In the same way, content-based approaches also have some problems, such as the lack of serendipity. For this reason, a hybrid approach which combines the previous techniques can improve the accuracy in the provided recommendations and reduce the cold start impact.

In this way, our LRS is based on two different approaches:

- First, a content-based recommender module based on two components. One component recommends taking into account the activity of the user identified in the system, the profile-based module (PB). The other component offers recommendations based on the content of the video being watched at the time, the item-based module (IB).
- Second, a collaborative filtering module (CF), which offers suggestions based on similar other users with a viewing history similar to current user.



Thus, the computed recommendations are based on these two approaches, getting a hybrid recommendation system, so that the user receives recommendations of videos similar to the one she is watching at that moment, or also of videos that may interest her due to her viewing history, and additionally, content considered interesting due to the opinions of other users with a viewing history similar to her own. This fact increases the serendipity, discovering contents to the user that are of her interest, although they are not so similar, a priori, to learning videos already watched by the user so far.

In our proposal, LOs/videos are characterized by their title and their transcript (when available). This characterization is used for the calculation of the similarity between the LOs. As transcripts are in different languages, it is possible to recommend videos from different languages.

Because there is a large set of words in the transcript and title of the videos, it is necessary to have an algorithm that filters out the too common words, which do not serve to differentiate the content. Thus, it is possible to focus on the particular words in the entire collection, which serve to identify the content of a video. We need to use an unsupervised term weighting approach, as our video repository is not categorized. Therefore, the algorithm chosen is the well-known term frequency–inverse document frequency (TF-IDF) [25], as it is a common term weighting scheme used to represent documents. In order to improve the performance of the TF-IDF algorithm, also the stop words from the *nlTK Python* (<https://www.nltk.org/>, accessed on 14 May 2021) package are used. Furthermore, some ad hoc words have been added to this package.

The item-based module takes only the information from the content of the videos, i.e., there is no information about the users. So, we can consider this module as a recommendation by item-item similarity to be used when a (maybe anonymous) user is watching a video. To do this, we take the characteristics (transcript and title) of each of the videos to calculate the item-item matrix with the TF-IDF algorithm. Then, the cosine similarity of two instances of this item-item matrix returns the similarity among the different videos. Cosine similarity calculates the similarity between two n-dimensional vectors by the angle between them in the vector space:

$$\text{cosine\_sim}(\vec{p}, \vec{q}) = \frac{\vec{p} \cdot \vec{q}}{|\vec{p}| * |\vec{q}|} \quad (1)$$

The profile-based module considers the content information (the terms extracted by the TF-IDF algorithm from the transcript and title) of the videos viewed by users. In this way, the recommendations to a user are made based on the similarity among the viewed videos of the user in mediaUPV. In this case, the similarity is calculated using the cosine profile-item matrix.

The collaborative filtering module considers the similarity between the users based on their viewed videos. Thus, two users will be more similar depending on how many videos they have both watched. In the opposite way, a user will have no similarity with another user if the sets of videos watched by each of the two users are completely disjoint. Therefore, the similarity is calculated using the cosine similarity between the users based on their watched videos. Then, the prediction values of recommendation for each user are calculated taking the dot product between the similarity of users and the viewed videos matrix, normalizing the data properly. With these process, we have the prediction to recommend each user in the system considering all users similarity (we will call this approach all users). However, we can also obtain a different prediction to recommend if we only consider the top  $k$  most similar users to the user to be recommended. To do that, we calculate the top  $k$  nearest neighbours (NN) for each user in the system, and then, the prediction values to recommend videos are also calculated (we will call this approach top  $k$  NN). We note that the computation cost of this last approach considering the top  $k$  NN is higher than the all users approach. The computation cost is also increased as  $k$  grows.

Finally, our *hybrid recommendation (HR)* system consists of two components, one content-based (with two modules) and another collaborative-based (with only one module).

Thus, the recommendations are made taking into account the intersection of the videos recommended by the modules (see Equation (2)). The rest of the recommendations are obtained considering first the weight applied to each component, i.e.,  $w_{CB}$  for the content-based component and  $w_{CF}$  for the collaborative filtering module. Furthermore, it is possible to balance the importance given to the two modules of the content-based component by means of other two weights,  $w_{IB}$  for the item-based module and  $w_{PB}$  for the profile-based one. In this way, in cases where there is no user (anonymous) the  $w_{PB}$  and  $w_{CF}$  are set to 0. Likewise, when an authenticated user is not yet watching a video, the  $w_{IB}$  is set to 0 (however, this case is not currently applicable to mediaUPV portal).

$$HR = (IB \cap PB \cap CF) \cup (w_{CB} \cdot ((w_{IB} \cdot IB) \cup (w_{PB} \cdot PB)) \cup (w_{CF} \cdot CF)) \quad (2)$$

In the next section, we explain the experiments carried out to determine the combination of weights which offers the best performance.

## 5. Experimental Results

In this section, we explain the experiments carried out with the proposed LRS. We have made tests with data from the videos watched by the users from September 2018 to July 2019, dividing this data into a training set and a test set. We used this data to be able to compare the new collaborative filtering module with the experiments already presented in our previous work [9]. In this section, we first present the general experimental setup in Section 5.1. Then, the experiments of both content-based modules are explained in Section 5.2. Section 5.3 presents the experiments with the collaborative filtering module. The experiments with the hybrid approach, which is the combination of the content-based component and the collaborative filtering module, are explained in Section 5.4. Finally, Section 5.5 contains a discussion of the results of all these experimental results.

### 5.1. Experimental Setup

As mentioned above, the database is comprised of all the videos available on the mediaUPV as items to be recommended, as well as the usage data (views) of the users logged into the system. Although the platform had more than 55,600 videos by the end of 2019, in the following tests we have filtered out the hidden videos (only available with a direct link), the videos that do not have a transcript, and the videos in which the transcript has less than 100 characters, leaving a total of 13,232 videos suitable for recommendation.

The data set used for these experiments is formed by of learning videos viewed by users during an academic year at the UPV, from September 2018 to July 2019. The training data set consists of the data from September 2018 to April 2019, while the testing data set consists of the data from May to July 2019, i.e., 8 months for training and 3 months for testing. The videos considered are those present on the platform until July 2019, after filtering them as described above. Thus, we try to simulate a real scenario, in which the training data represent the past activity of the users, while the test data is formed by the activity of the following 3 months ("future"). Therefore, any recommendation from the recommender that is among the videos that users have actually watched in the test set is considered a success. Additionally, in these tests we only consider five recommendations since it is the number required by the mediaUPV portal.

We use the well-known precision and recall measures to evaluate the success of the recommender. Precision can be defined as the successful recommendations made (videos that have been viewed by the user in the test set) divided by the number of recommendations made:

$$precision = \frac{success\_recommendation}{recommendations\_made} \quad (3)$$

Recall is defined as the successful recommendations made divided by the number of watched videos in the test set:

$$recall = \frac{success\_recommendation}{watched\_videos} \quad (4)$$

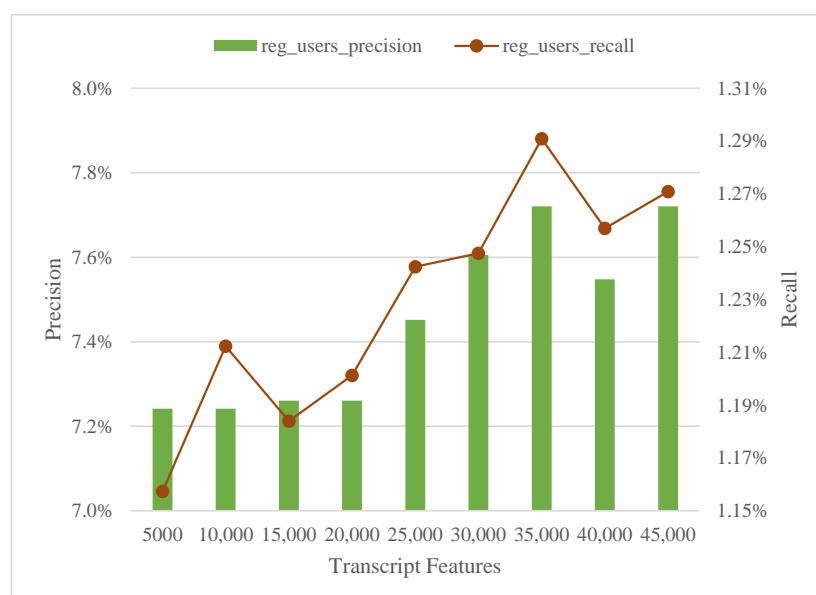
For our tests we considered a set of regular users (*reg\_users*) that we define as those who have watched between 10 and 150 videos both in the training and test periods, having 1044 users in this set. We also considered a set of new users that have watched between 1 and 9 videos both in the training and test periods (there are 815 users in this set). If nothing is specified, the regular users set is the one used.

## 5.2. Content-Based Component

To test the content-based modules we use both of them together to combine their efforts. So this is like using the hybrid recommender engine but without considering the collaborative filtering module. Thus, we focus on the content-based module centred on the video being watched, that is, the item-based module (*IB*); and on the content-based module that considers the user's views, namely the profile-based module (*PB*). The weights for each module in the hybrid recommender engine are set to  $w_{IB} = 50\%$ ,  $w_{PB} = 50\%$  and  $w_{CF} = 0\%$ .

### 5.2.1. Setting Transcript and Title Features

In this first test, we analyze the success of recommendations for the set of regular users considering different amount of features for the transcript and the title of the video to train the LRS, in order to establish the better amount of both. The graph in Figure 1 shows the precision and recall for different values of the number of features considered for the transcript, while the number of features for the title is kept at zero.



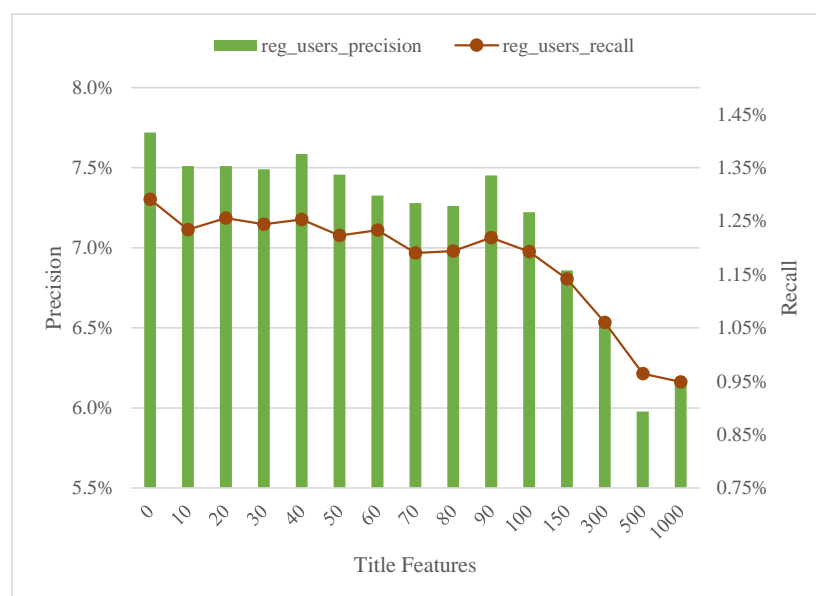
**Figure 1.** Precision and recall results for the regular users set with different amount of transcript features.

In general, precision and recall increase slightly as the value of the features for the transcript increases (from 7.2% to 7.7% for precision, and from 1.15% to 1.29% for recall), as would be expected when more information is available from the transcript. The best precision values are with 35,000 and 45,000 features for the transcript. However, the best recall value is in the case of 35,000 transcript features. Therefore, the best configuration for the recommender would be to use 35,000 transcript features, since this value achieves higher recall than any other and equals the precision obtained with 45,000 transcript

features. In addition, the computation of 35,000 transcript features is computationally less expensive and, in particular, implies a lower memory cost.

It should be noted that, in the experiment made with 35,000 transcript features, the number of users who have been recommended successfully is 254 of 1044, i.e., 24.33%.

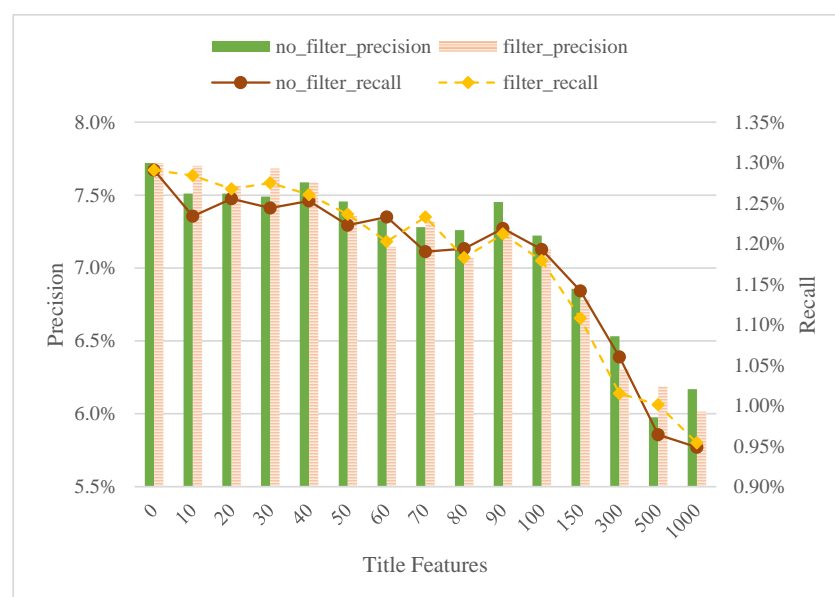
Figure 2 shows the precision and recall for different values in the number of features considered for the title having the transcript features fixed to 35,000. The best values of precision and recall are obtained with the number of features of the title at 0. In addition, both values decrease slightly and in a relatively uniform way as the number of title features increases. So, apparently it is better to skip the title features. However, it is interesting to analyze this considering the nature of the LOs of mediaUPV. In this way, we have analyzed the set of words that determines the TF-IDF algorithm. As we mentioned before, the videos on the platform correspond mainly to university courses, so there is a set of terms that are certainly repetitive in the titles of the videos and do not provide any differentiating information with respect to their content. Among these terms, we find the following: {'analysis', 'calculation', 'control', 'creation', 'data', 'design', 'exercise', 'engineering', 'introduction', 'management', 'mechanism', 'model', 'module', 'practical', 'practice', 'presentation', 'simulation', 'system', 'systems', 'theme', 'unit', 'virtual'}.



**Figure 2.** Precision and recall results for the regular users set with different amount of title features.

### 5.2.2. Filtering Title Features

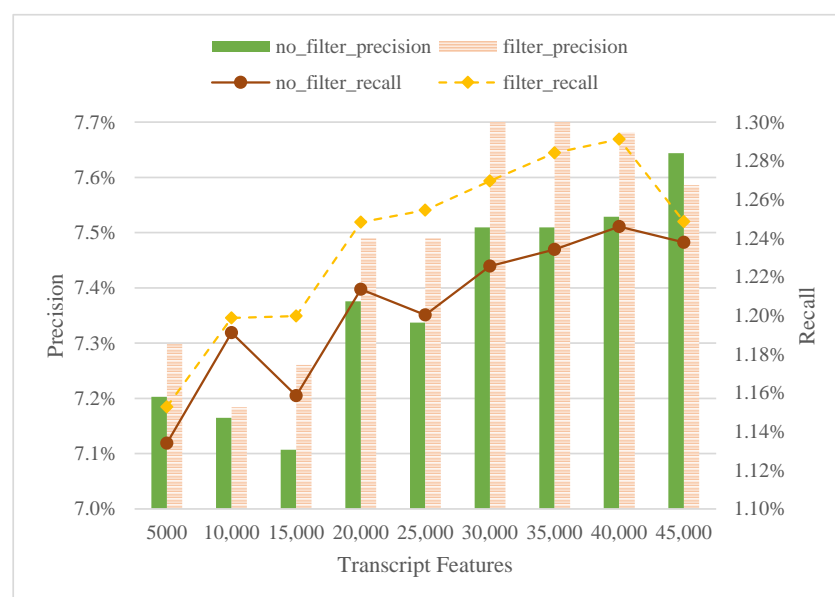
In order to increase the precision of the recommendations, we decided to filter the terms of the previous list from the titles of the videos by considering them as stop words. The results of applying this correction can be seen in the Figure 3, which shows the precision and recall for 35,000 transcript features and different amount of title features. In this case, it can be seen that the best global values of precision and recall are still obtained with 0 title features. However, it should be noted that with 10 title features, for the case where the specified title terms have been filtered out, the precision and recall almost reach this base case, slightly surpassing the case without filtering (with 10 title features). In addition, for values up to 30 title features, the precision and recall are better for the case with filtered terms. However, from 40 title features on, the effect of filtering is diluted and the results are generally slightly worse, and most of the cases slightly worse than the unfiltered case.



**Figure 3.** Precision and recall comparison with and without filtering title terms for different amount of title features.

Consequently, we can say that filtering has significantly improved the results but it is not enough to make the title relevant. Perhaps it would be necessary a still greater filtering of terms that we have not considered ‘commo’ and that the TF-IDF algorithm has not identified as such either. However, since we are considering 35,000 transcript terms, the inclusion of 10 to 50 terms from the title can be considered irrelevant after the analysis. Furthermore, it can also be interpreted as a video being better characterized by its own transcript than by its title.

Although we have already seen that it is better to skip the title in all cases, we will analyze in detail the difference between filtering the title terms and not filtering them for 10 title terms and different amount of transcript features. Figure 4 shows this comparison (precision and recall) with and without filtering. Precision and recall are significantly better if filtering of title terms is performed in all cases, with the only exception of 45,000 transcript features where precision is slightly higher for the unfiltered case (being also the best result for the different transcript values for the unfiltered case of title terms). In this particular case, it could be that by considering only 10 features of the title, but 45,000 for the transcript, the effect of the filtering of the title is diluted. However, this is not a very significant difference. On the other hand, as previously observed, the best results, both in terms of precision and recall, are obtained with 35,000 transcript features in the case of filtering the title terms.



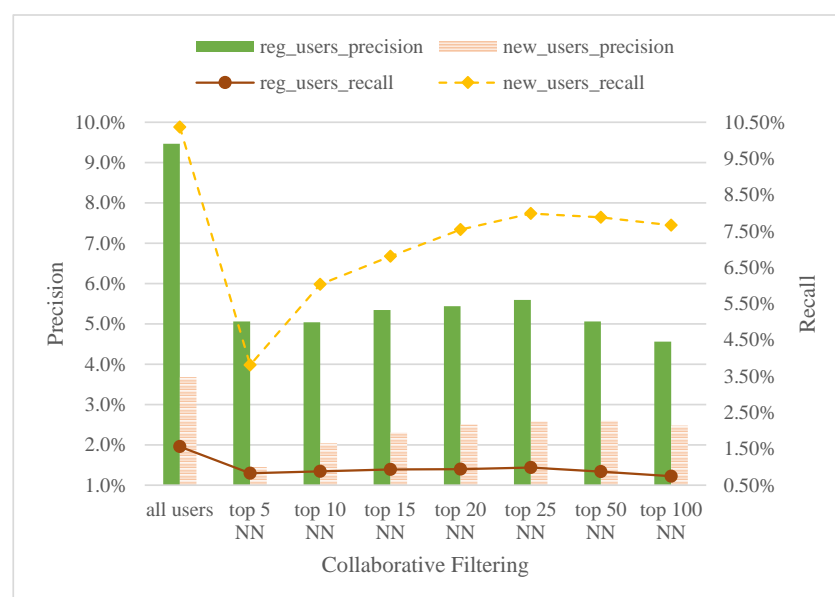
**Figure 4.** Precision and recall comparison with and without filtering title terms for different amount of transcript features with 10 title features.

### 5.3. Collaborative Filtering Component

In this subsection we conducted experiments to test the performance of the collaborative filtering module in its two variants, i.e., the one that considers the similarity with all users (which we refer to it as all users) to make recommendations (taking the videos with the highest recommendation values) versus the one that only considers the top  $k$  nearest neighbour (NN) users to make the relevant recommendations.

Figure 5 presents the results of precision and recall for the sets of regular users (watched from 10 to 150 videos) and new users (watched from 1 to 9 videos). The values shown correspond to the collaborative filtering module considering all users recommendations in the first column, and the subsequent columns consider the top  $k$  NN users for recommending. Generally, for both the regular and new users sets, the best results in precision and recall are obtained by the CF all users approach, which in fact doubles the values for almost all the top  $k$  NN approaches. In this way, it is clear that the best approach for collaborative filtering with mediaUPV data is the one that uses CF all users instead of any number of the top  $k$  NN. Additionally, this approach is less costly computationally. If we compare the different  $k$  values of the top NN, it seems that the best are 25 for the regular users set, and 50 for the new users set, slightly decreasing in both cases as  $k$  grows. All in all, the difference with the CF all users approach is significant enough to avoid any top  $k$  NN approach.

In the case of the regular users, the CF all users approach achieves almost 9.5% precision and 1.5% of recall. However, with the new users set, the precision is almost 3.7% and the recall over 10%. This difference between both users sets is mainly due to the amount of historical data of the users. In the case of the new users set, data about them is limited, as they have only watched 1 to 9 videos. This causes the precision of the recommendation to be significantly lower than the precision obtained in the regular users set, in which the amount of data is larger to build a more solid recommender. This difference is commonly known in literature as the cold start problem. However, since our recommender engine has different modules, we can leverage them to obtain satisfactory results in any case.

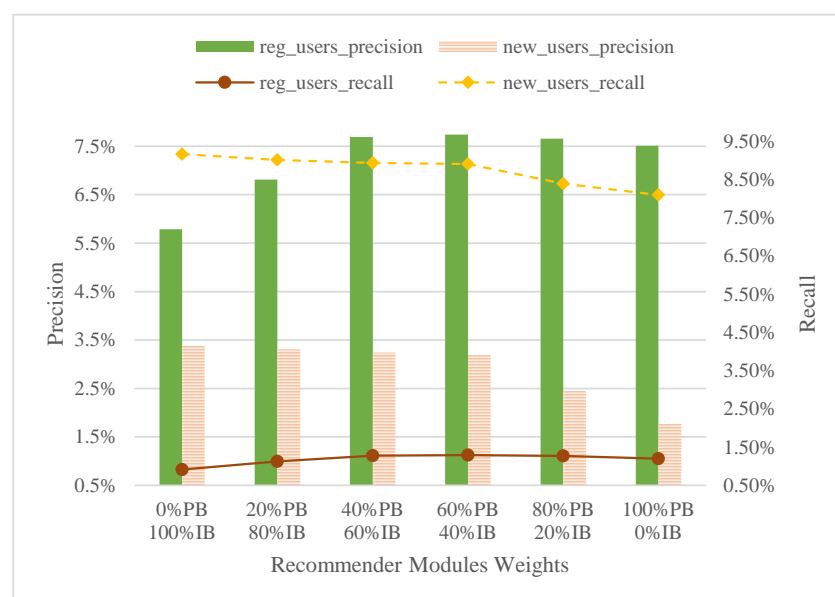


**Figure 5.** Precision and recall comparison of regular users and new users for collaborative filtering module with different values of top  $N$  users obtained by nearest neighbour (NN) technique.

#### 5.4. Hybrid Weights Setting

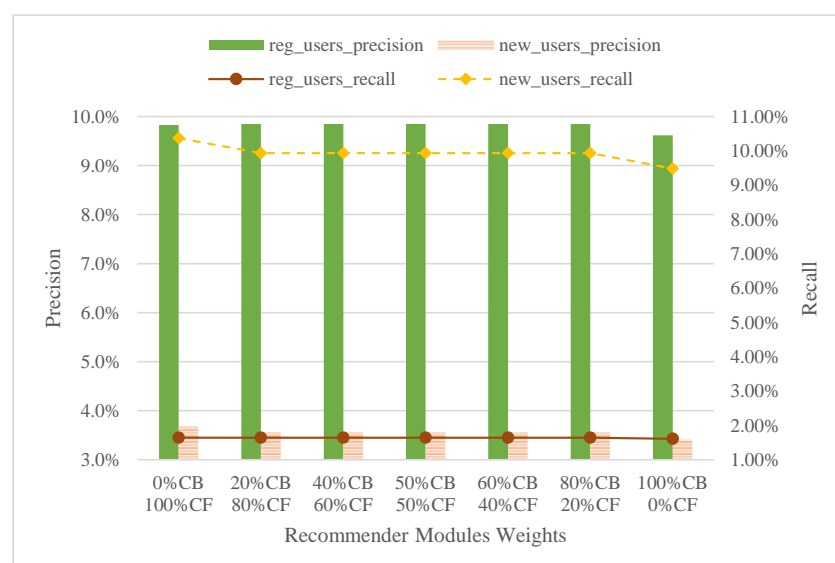
Our LRS is a hybrid approach that considers, on the one hand, a content-based component with two modules named profile-based module and item-based module, and on the other hand, a collaborative filtering component (with a unique module) as it is specified in Equation (2). In this subsection, we first analyse the best combination of weights for the content-based modules. Once found them, we study how to tune the weights of the content-based component and the collaborative filtering component.

In Figure 6 we make a comparison of precision and recall of two different sets of users, using different weights for the content-based modules of our hybrid LRS. We first consider the set of new users of which we have little knowledge as they have only watched 1 to 9 videos in the data set (815 users). The second set are the regular users, formed by users that have already watched between 10 and 150 videos in the data set (1044 users). For regular users, the best precision and recall is obtained with balanced weights, i.e.,  $\{w_{PB} = 40\%; w_{IB} = 60\%\}$  and  $\{w_{PB} = 60\%; w_{IB} = 40\%\}$ . However, for new users, the best precision and recall values are obtained with low  $w_{PB}$ , with 15.83% of new users receiving successful recommendations.



**Figure 6.** Precision and recall comparison of regular users and new users, using different content-based modules weights.

Having set the weights of the content-based modules to  $w_{PB} = 60\%$  and  $w_{IB} = 40\%$ , since they obtain the best results, we can now set the weights of the combination of the content-based component ( $w_{CB}$ ) with the collaborative filtering component ( $w_{CF}$ ) as specified in Equation (2). For this, Figure 7 shows the results of precision and recall for both groups of regular and new users with different weights for the recommender modules, the content-based modules (CB) and the collaborative filtering module (CF). In this case, the results for the different values of weights are the same except of the extreme values of  $\{w_{CB} = 0\%; w_{CF} = 100\%\}$ , and  $\{w_{CB} = 100\%; w_{CF} = 0\%\}$ , that are slightly lower. The reason behind this is that the intersection of recommendations of both the content-based modules and the collaborative filtering module already gives the best results, and hence, any combination of the weights (except the extremes) can be considered.



**Figure 7.** Precision and recall comparison of regular users and new users, using different content-based and collaborative filtering modules weights.



The incorporation of the collaborative filtering module improves the precision from 7.7% to 9.85% for regular users (for new users is increased from 3.3% to 3.56%) with respect to the version of the LRS with only the content-based modules proposed in [9]. Recall is also slightly improved from 1.3% to 1.6% for regular users (9.1% to 9.9% for new users). In this case, 28.5% of the regular users and 16.4% of the new users received a useful recommendation. We also note that the results of the collaborative filtering module alone (see Figure 5) when compared with the combined version with the content-based modules also improve from 9.5% to 9.85% of precision for the regular users. Therefore, all of these are positive results that justify the need of applying the collaborative filtering component to our LRS to improve its accuracy.

Since there is no significant difference between the weights of the content-based component and the collaborative filtering component, we propose balanced weights to apply in production  $\{w_{CB} = 50\%; w_{CF} = 50\%\}$ .

### 5.5. Discussion

For new users from which the system has few information is harder to make successful recommendations; however, this could improve in the future during the application of the LRS as the users get more engaged to mediaUPV portal. In the case of the users from which there is more historical data available, the accuracy of our LRS improves significantly, specially when using the new developed collaborative filtering module combined with the content-based modules. Furthermore, the collaborative filtering module will be able to recommend videos that do not have transcript (we remind that only 20,135 from the 80,500 videos have transcript), which will suppose a lot of more possibilities to recommend to the logged-in users when applied to production.

We emphasize that even though our hybrid LRS obtains a precision of 9.85% simulating a real environment (we improved previous results of 7.7% precision without the collaborative filtering module of [9]), 28.5% regular users and 16.4% of new users received some good recommendation. In addition, the precision and recall obtained are significantly better than a random recommendation.

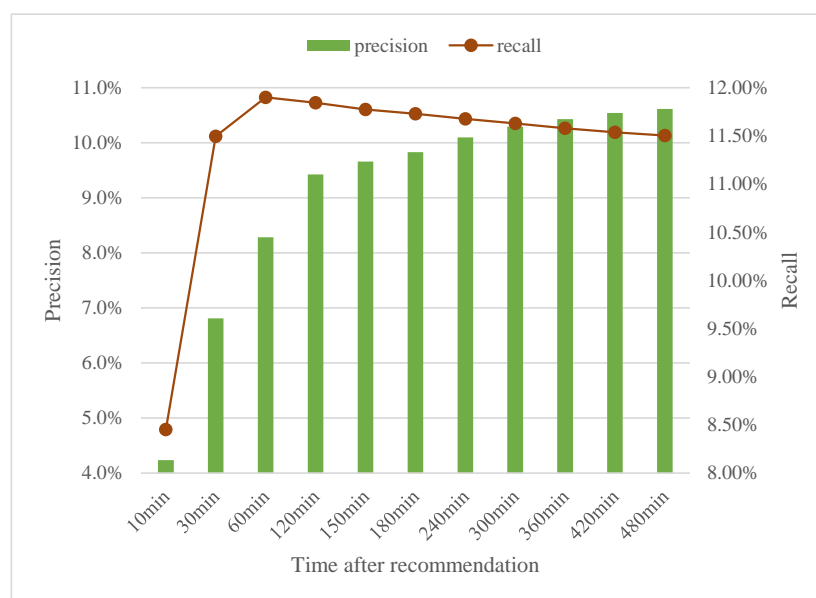
In conclusion, after this analysis we established the parameters of the LRS to be applied in production in the mediaUPV portal as 35,000 transcript features, 0 title features,  $w_{CB} = 50\% \cdot (w_{PB} = 60\%, w_{IB} = 40\%), w_{CF} = 50\%$ .

## 6. Production Results

In this section we show the results of the application of the LRS described in [9] to the mediaUPV portal. That proposal did not contain a collaborative filtering module. It should be noted that mediaUPV did not have a LRS until the application in October 2019. The parameter configuration of the recommender system applied to production was the one that worked best in the experiments of the previous work [9], that is, 35,000 transcript features, 0 title features,  $w_{PB} = 60\%$  and  $w_{IB} = 40\%$  (without the collaborative filtering module which has been developed for the present work).

The graph of Figure 8 presents the global precision and recall results of the recommendations made to the users for the period from October 2019 until March 2021. We compare different time ranges, from 10 to 480 min, in which the user can watch the recommended videos. Hence, the column of 10 min represents the precision and recall of the video recommendations that are watched within the 10 min after the recommendation is made. So the precision increases as the time range grows. As it can be seen in Figure 8, the precision is significantly lower for the cases below 120 min, and then, it only increases slightly. The reason behind that might be that the users of mediaUPV portal usually watch long videos, or they even do other tasks (like homework in the case of students, or preparing other classes or new material in the case of lecturers) between watching a video and the next. This may explain why the precision increases significantly if we consider a time after recommendation of at least 120 min instead of 10 or 30 min.

Overall, a precision of around 9.5% for 120 min after the recommendation, or higher than 10% if we consider more minutes is a suitable result for our recommender system. We note that most of the users that enter in mediaUPV portal do not seek for recommendations, since they only watch the video that they need to (i.e., a student who must watch the corresponding lesson). Additionally, the production results are significantly better (around 2–3% higher precision) than the experimental results with the training and test sets of September 2018 to July 2019 originally used in our previous work [9].



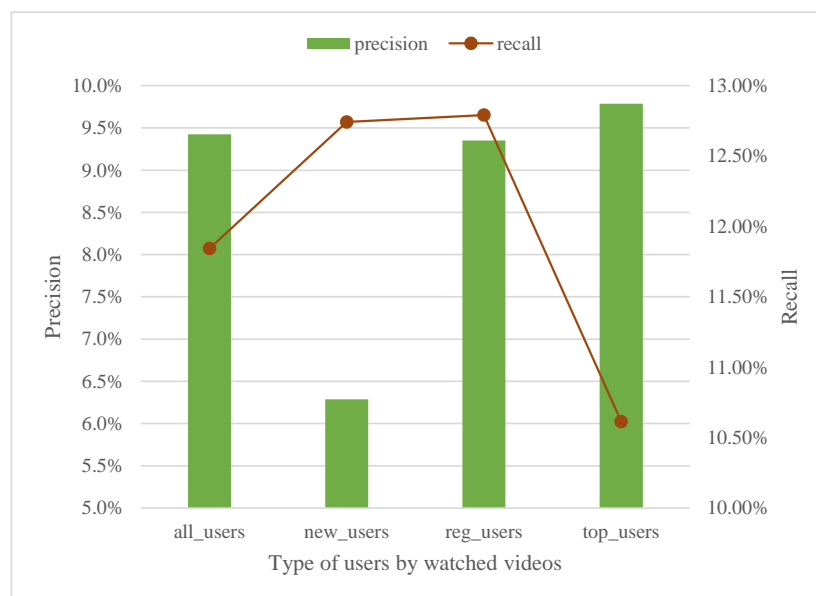
**Figure 8.** Production results of precision and recall considering different time range for the user to watch the recommended video after the recommendation.

Figure 9 presents the results of precision and recall for different type of users considering at most 120 min after the recommendation to watch the video. All users (all\_users) is the set that includes the total amount of users that watched any video and received any recommendation in our system from October 2019 to March 2021. The set of all users is divided in other three sets, namely: new users (new\_users), which are those that watched between 1 and 9 videos in the period of application of our recommender system; regular users (reg\_users), the ones that watched 10 to 150 videos in the period; and top users (top\_users), those who watched more than 150 videos in the period. According to this classification, we have 2297 new users, 8202 regular users, and 634 top users. The results in precision and recall for all users and regular users are reasonably similar since this last set is the larger, so it influences more the all users set, but also it is the middle point between the new users and top users set.

The main results of Figure 9 show that the precision of our recommender engine for new users is 6.3% with a recall of 12.74%. This precision is significantly lower than the precision for both the regular users and the top users, that is 9.23% and 9.8%, respectively, with a lower recall for the top users of 10.6%. From these results, we can confirm that as the known data from the user increases, the accuracy of the recommendations also increase. However, it is even hard to increase precision with top users due to the nature of the recommendations of mediaUPV portal, in which most users only come to watch the specific videos they have to. This scenario differs significantly from platforms like YouTube where most users enter looking to spend some leisure time.

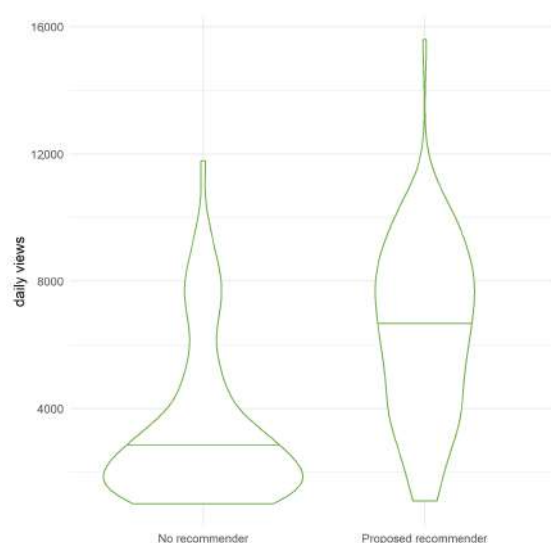
Globally, we make some successful recommendation to 35% of all users. Particularly, we have successful recommendations for 24% of new users, 51% in the case of regular users, and 78.6% for the top users. This results clearly show that a successful or useful recommendation for a user is totally related to the amount of (historical) data that the system has of the user. It is important to note that almost the 80% of the top users (the

most informed set of users we defined) received any successful recommendation, which is a very positive result for the recommender system. We can assume that our recommender would be more accurate as the users are more engaged with it.



**Figure 9.** Production results of precision and recall by different type of users within the 120 min after the recommendation. All users include the total amount of users that received any recommendation. New users are those that watched between 1 and 9 videos, regular users watched 10 to 150 videos, and top users watched more than 150 videos.

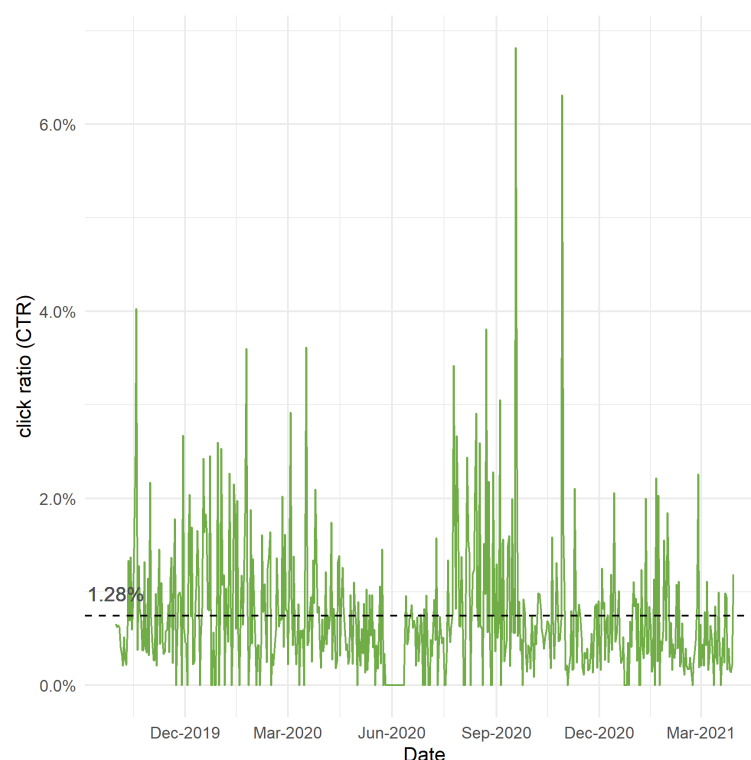
In Figure 10, we show a violin plot (<https://towardsdatascience.com/violin-plots-explained-fb1d115e023d>, accessed on 14 May 2021) depicting the distribution of daily accesses to the videos of the mediaUPV platform during the year 2019, with no recommender (before October) and with the proposed one. As can be seen clearly, the existence of the LRS is correlated with an increase of the number of accesses to the videos by the users.



**Figure 10.** Video access with and without a recommender in production.

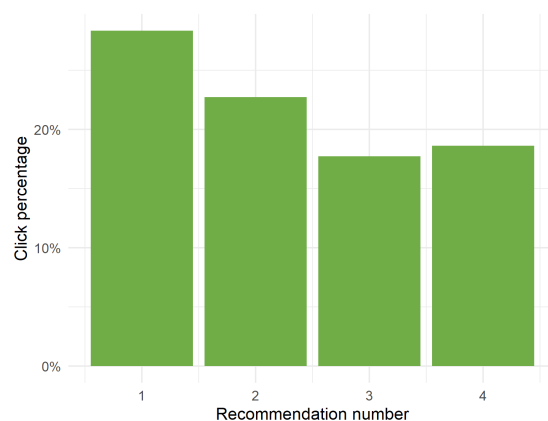
A common way in industry to measure relative quality of a recommender system is the Click-Through Rate (CTR) (<https://en.wikipedia.org/wiki/Click-throughrate>, accessed on 14 May 2021), that measures the percentage of clicks in the recommender per number

of views. As the CTR is used by the ads industry, there is an ongoing interest in CTR prediction techniques [26]. In the case of a generic recommender system, anything above 0.35% means you are doing a good job ([www.acquisio.com/blog/agency/what-is-a-good-click-through-rate-ctr/](http://www.acquisio.com/blog/agency/what-is-a-good-click-through-rate-ctr/), accessed on 14 May 2021). As can be seen in Figure 11, the CTR is 1.28% on average, with notable peaks over 4%. These results can be considered quite satisfactory, since they imply that the recommendations made to users generate interest in them.



**Figure 11.** Clicks in recommender (CTR) in production.

Finally, we point out the percentage of clicks on each of the recommended videos according to their order in the list, demonstrating the relevance of this order to users (see Figure 12). About 28% of users click on the first video, and more than half of the users click on the first two recommendations with a distribution that seems to be heavy tailed. So the most relevant video by far for users is the first one and then the rest of the videos follow a decreasing order of importance, which also points to a reasonable work of the presented LRS.



**Figure 12.** Percentage of clicks in recommendations per list position.

## 7. Conclusions

This work proposes a new hybrid LRS based on collaborative filtering and content-based components capable of recommend learning videos based on viewing history and current video content. Thus, the LRS proposed is able to recommend not only to authenticated users but also to anonymous users from the mediaUPV portal, independently if they are lecturers or students. In fact, mediaUPV portal has not information about learners' profiles or needs, as it is not connected with any LMS.

The hybrid LRS has been applied to a simulated environment, using a data set of learning videos and user profiles from the 2018-2019 academic year at UPV. The best hybrid LRS configuration obtained 9.85% of precision and 1.6% of recall, where 28.5% of regular users received some useful recommendation.

Furthermore, the content-based component of the approach has been applied in a real scenario, the mediaUPV portal of the UPV from October 2019 to March 2021. This portal is mainly used by learners and trainers to access to useful LOs for their MOOCs and flipped classrooms. We can state that the application of this LRS to the mediaUPV portal was positive as it improved the precision of the original experimental results of [9] (from 7.7% to 9.85%), it brought an increase in visits to the videos, and it had a significant CTR of 1.28% on average, with notable peaks of over 4%.

The results of our LRS must be seen in the context of its application. In this respect, it should be remembered that we are dealing with videos of university lectures or subject-specific lessons. This means that users of the system do not usually enter for leisure purposes as on YouTube or Netflix, or to make purchases as on platforms such as Amazon. Thus, users of this system usually enter to watch a specific lesson of the subjects they are studying, or to search for a specific video about a particular topic. Therefore, it is difficult for a recommender to obtain better results than the ones we show. In addition, it is important to highlight what has been observed with the results in production with respect to the time that elapses between the recommendation and the moment in which the users watch one of the recommended videos. In this sense, it has been shown (Figure 8) that the precision of the recommender almost doubles if instead of considering the 30 min after the recommendation we consider 120 min or more.

As future work, we want to analyse the results in production of both components working together, the collaborative filtering and content-based one. Thus, the collaborative filtering component presented in this work is able to recommend videos that have no transcript, which opens up more possibilities to increase the serendipity of the recommendations. In addition, it would be interesting to evaluate whether to add the classification of the mediaUPV videos obtained by [21] to the current characterization of the videos used by our proposal, which is currently based on the video transcript and collaborative filtering.

We also want to test if we achieve better results by changing the term frequency-inverse document frequency (TF-IDF) algorithm in the content-based module using other techniques such as delta TF-IDF [27] or TF.IDF.ICF [28], which try to avoid the TF-IDF problem of not considering intraclass or interclass distributions. In addition, we need to further study if it will be possible to apply other variant weighting approaches of TF-IDF, such as the presented in [29], in which the number of occurrences of a term, the number of documents that include the term, and the number of classes in which the term appears, are used to obtain a more accurate set of characteristic features.

Finally, we would like to conduct random user surveys to analyse the precision of the recommender in a less automatic way and to get direct feedback from users. In this way, we could know better the precision of the recommendations as the users could answer if these recommendations are useful instead of basing them on whether they have seen the video or not. This would avoid the uncertainty of the current analysis in which we do not know for sure if users do not find the recommendations useful or if they only enter the platform to watch the content they need without paying attention to any recommendation, whether it is useful or not.

**Author Contributions:** Conceptualization, J.J., S.V. and V.B.; Data curation, J.J. and C.T.; Formal analysis, J.J. and S.V.; Funding acquisition, V.B.; Investigation, J.J., S.V. and C.T.; Methodology, V.B.; Project administration, V.B.; Resources, C.T.; Software, J.J.; Supervision, C.T. and V.B.; Validation, J.J., S.V. and C.T.; Writing—original draft, J.J. and S.V.; Writing—review & editing, J.J., S.V., C.T. and V.B. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was partially supported by MINECO/FEDER RTI2018-095390-B-C31 and TIN2017-89156-R projects of the Spanish government, and PROMETEO/2018/002 project of Generalitat Valenciana.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

CB	Content-Based Recommender Modules
CF	Collaborative Filtering Module
CTR	Click-Through Rate
HR	Hybrid Recommender
IB	Item-Based Recommender Module
LMS	Learning Management System
LRS	Learning Recommender System
LO	Learning Object
MOOC	Massive Online Open Course
PB	Profile-Based Recommender Module
TF-IDF	term frequency—inverse document frequency
UPV	Universitat Politècnica de València

## References

1. Maassen, P.; Nerland, M.; Yates, L. (Eds.) *Reconfiguring Knowledge in Higher Education*; Higher Education Dynamics, 50; Springer International Publishing: Cham, Switzerland, 2018.
2. Zajda, J.; Rust, V. *Globalisation and Higher Education Reforms*; Springer: Cham, Switzerland, 2016; Volume 15.
3. van Dijck, J.; Poell, T. Higher Education in a Networked World: European Responses to U.S. MOOCs. *Int. J. Commun. IJoC* **2015**, *9*, 2674–2692.
4. Institute and Committee of Electrical and Electronics Engineers; Learning Technology Standards. *IEEE Standard for Learning Object Metadata*; IEEE Standard 1484.12.1; IEEE: Piscataway, NJ, USA, 2002.
5. O’Flaherty, J.; Phillips, C. The use of flipped classrooms in higher education: A scoping review. *Internet High. Educ.* **2015**, *25*, 85–95. [CrossRef]
6. Roehl, A.; Reddy, S.L.; Shannon, G.J. The flipped classroom: An opportunity to engage millennial students through active learning strategies. *J. Fam. Consum. Sci.* **2013**, *105*, 44–49. [CrossRef]
7. Tucker, B. The Flipped Classroom. Online instruction at home frees class time for learning. *Educ. Next* **2012**, *2012*, 82–83.
8. Turró, C.; Morales, J.C.; Busquets-Mataix, J. A study on assessment results in a large scale Flipped Teaching Experience. In Proceedings of the 4th International Conference on Higher Education Advances (HEAD’18), Valencia, Spain, 20–22 June 2018; pp. 1039–1048.
9. Jordán, J.; Valero, S.; Turró, C.; Botti, V. Recommending Learning Videos for MOOCs and Flipped Classrooms. In *Advances in Practical Applications of Agents, Multi-Agent Systems, and Trustworthiness*; Demazeau, Y., Holvoet, T., Corchado, J.M., Costantini, S., Eds.; The PAAMS Collection; Springer International Publishing: Cham, Switzerland, 2020; pp. 146–157.
10. Klačnja-Milićević, A.; Ivanović, M.; Nanopoulos, A. Recommender systems in e-learning environments: A survey of the state-of-the-art and possible extensions. *Artif. Intell. Rev.* **2015**, *44*, 571–604. [CrossRef]
11. Burke, R. Hybrid Recommender Systems: Survey and Experiments. *User Model. User-Adapt. Interact.* **2002**, *12*, 331–370. [CrossRef]
12. Herlocker, J.; Konstan, J.; Terveen, L.; Riedl, J. Evaluating collaborative filtering recommender systems. *ACM Trans. Inf. Syst.* **2004**, *22*, 5–53. [CrossRef]
13. Rodríguez, P.; Heras, S.; Palanca, J.; Duque, N.; Julián, V. Argumentation-Based Hybrid Recommender System for Recommending Learning Objects. In *Multi-Agent Systems and Agreement Technologies*; Springer International Publishing: Cham, Switzerland, 2016; pp. 234–248.
14. Bobadilla, J.; Serradilla, F.; Hernando, A. Collaborative filtering adapted to recommender systems of e-learning. *Knowl.-Based Syst.* **2009**, *22*, 261–265. [CrossRef]
15. Zapata, A.; Menéndez Domínguez, V.; Prieto, M.; Romero, C. A Hybrid Recommender Method for Learning Objects. *Int. J. Comput. Appl. (DEDCE)* **2011**, *1*, 1–7.

16. Chen, W.; Niu, Z.; Zhao, X.; Li, Y. A Hybrid Recommendation Algorithm Adapted in E-Learning Environments. *World Wide Web* **2014**, *17*, 271–284. [CrossRef]
17. Tarus, J.K.; Niu, Z.; Yousif, A. A hybrid knowledge-based recommender system for e-learning based on ontology and sequential pattern mining. *Future Gener. Comput. Syst.* **2017**, *72*, 37–48. [CrossRef]
18. Vanitha, V.; Krishnan, P.; Elakkiya, R. Collaborative optimization algorithm for learning path construction in E-learning. *Comput. Electr. Eng.* **2019**, *77*, 325–338. [CrossRef]
19. Dwivedi, P.; Kant, V.; Bharadwaj, K.K. Learning Path Recommendation Based on Modified Variable Length Genetic Algorithm. *Educ. Inf. Technol.* **2018**, *23*, 819–836. [CrossRef]
20. Dwivedi, P.; Bharadwaj, K.K. e-Learning recommender system for a group of learners based on the unified learner profile approach. *Expert Syst.* **2015**, *32*, 264–276. [CrossRef]
21. Stoica, A.S.; Heras, S.; Palanca, J.; Julian, V.; Mihaescu, M.C. A Semi-supervised Method to Classify Educational Videos. In *Hybrid Artificial Intelligent Systems*; Springer International Publishing: Cham, Switzerland, 2019; pp. 218–228.
22. Turcu, G.; Heras, S.; Palanca, J.; Julian, V.; Mihaescu, M.C. Towards a Custom Designed Mechanism for Indexing and Retrieving Video Transcripts. In *Hybrid Artificial Intelligent Systems*; Springer International Publishing: Cham, Switzerland, 2019; pp. 299–309.
23. Blei, D.M.; Ng, A.Y.; Jordan, M.I. Latent Dirichlet Allocation. *J. Mach. Learn. Res.* **2003**, *3*, 993–1022.
24. MLLP Research Group (Universitat Politècnica de València). TLP: The transLectures-UPV Platform. Available online: <https://www.mllp.upv.es/tlp/> (accessed on 14 May 2021).
25. Salton, G.; Buckley, C. Term-weighting approaches in automatic text retrieval. *Inf. Process. Manag.* **1988**, *24*, 513–523. [CrossRef]
26. Richardson, M.; Dominowska, E.; Ragno, R. Predicting clicks: Estimating the click-through rate for new ads. In Proceedings of the 16th International Conference on World Wide Web, Banff, AB, Canada, 8–12 May 2007; pp. 521–530.
27. Martineau, J.; Finin, T. Delta tfidf: An improved feature space for sentiment analysis. In Proceedings of the International AAAI Conference on Web and Social Media, San Jose, CA, USA, 17–20 May 2009; Volume 3.
28. Ren, F.; Sohrab, M.G. Class-indexing-based term weighting for automatic text classification. *Inf. Sci.* **2013**, *236*, 109–125. [CrossRef]
29. Kansheng, S.; Jie, H.; Liu, H.T.; Zhang, N.T.; Song, W.T. Efficient text classification method based on improved term reduction and term weighting. *J. China Univ. Posts Telecommun.* **2011**, *18*, 131–135.





## Article

# Forecasting Energy Consumption of Wastewater Treatment Plants with a Transfer Learning Approach for Sustainable Cities

Pedro Oliveira , Bruno Fernandes , Cesar Analide  and Paulo Novais 

ALGORITMI Centre, Department of Informatics, University of Minho, 4710-057 Braga, Portugal; bruno.fmf.8@gmail.com (B.F.); analide@di.uminho.pt (C.A.); pjon@di.uminho.pt (P.N.)

\* Correspondence: poliveira199208@gmail.com

**Abstract:** A major challenge of today's society is to make large urban centres more sustainable. Improving the energy efficiency of the various infrastructures that make up cities is one aspect being considered when improving their sustainability, with Wastewater Treatment Plants (WWTPs) being one of them. Consequently, this study aims to conceive, tune, and evaluate a set of candidate deep learning models with the goal being to forecast the energy consumption of a WWTP, following a recursive multi-step approach. Three distinct types of models were experimented, in particular, Long Short-Term Memory networks (LSTMs), Gated Recurrent Units (GRUs), and uni-dimensional Convolutional Neural Networks (CNNs). Uni- and multi-variate settings were evaluated, as well as different methods for handling outliers. Promising forecasting results were obtained by CNN-based models, being this difference statistically significant when compared to LSTMs and GRUs, with the best model presenting an approximate overall error of 630 kWh when on a multi-variate setting. Finally, to overcome the problem of data scarcity in WWTPs, transfer learning processes were implemented, with promising results being achieved when using a pre-trained uni-variate CNN model, with the overall error reducing to 325 kWh.

**Keywords:** deep learning; energy consumption; sustainable cities; transfer learning; wastewater treatment plants

**Citation:** Oliveira, P.; Fernandes, B.; Analide, C.; Novais, P. Forecasting Energy Consumption of Wastewater Treatment Plants with a Transfer Learning Approach for Sustainable Cities. *Electronics* **2021**, *10*, 1149. <https://doi.org/10.3390/electronics10101149>

Academic Editors: Juan M. Corchado, Josep L. Larriba-Pey, Pablo Chamoso and Fernando De la Prieta

Received: 31 March 2021

Accepted: 3 May 2021

Published: 12 May 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Over the years, there has been an increase in global urbanisation through a greater concentration of people in small spaces. According to the World Urbanisation Perspectives report carried out in 2017 by the United Nations on the number of people living in urban and rural areas worldwide, it was found that 4.1 billion people already lived in urban areas [1]. In fact, cities have a fundamental role in sustainable development, namely related to economic and environmental concerns.

With the increase in energy consumption, concerns about the energy sector have expanded substantially. Although there has been a greater awareness of the impact of non-renewable energy sources on the planet and the high emission of greenhouse gases, if concrete and imperative measures are not applied, this problem will only worsen. Thus, over the years, the term energy efficiency has become increasingly important and indispensable. Energy efficiency can help reduce energy production and, consequently, reduce greenhouse gas emissions and preserve fossil fuel resources, ensuring a notable contribution to reducing environmental problems on our planet [2].

There are several infrastructures where energy consumption is high in a city, with Wastewater Treatment Plants (WWTPs) being one of them. In a WWTP, achieving a high energy efficiency level has become an increasingly important topic [3]. WWTPs, with the execution of their functions, demand high levels of energy, reflecting about 7% of all energy consumed worldwide [4]. In Portugal, about 4% of the consumed electricity is urban water cycle's responsibility, with approximately 25% of that energy being used in WWTPs [5].

Reducing energy consumption, emission of greenhouse gases and operating costs has been one of the main concerns of WWTP managers, who have been adopting more efficient equipment and technologies [6,7]. Hence, a WWTP must always consider the efficient management of all its resources, including energy.

Currently, in most WWTPs, low levels of energy efficiency performance are found. In fact, several factors influence the consumed energy in this type of facilities, depending on their characteristics and the types of treatments being applied. In general, the lack of energy efficiency is due to [8]:

- A growing need for water recycling due to the scarcity of this resource;
- Types of motors and pumps that are used;
- Higher requirements on discharge parameters in the treated effluent;
- Water pumping processes, which require high energy consumption;
- Absence of energy recovery mechanisms;
- Low efficiency in operations, mixing, and aeration systems;
- Influent flow.

### 1.1. State of the Art

A study carried out by Li et al. [9] aimed at predicting energy consumption in a WWTP through the use of a Radial Basis Function (RBF) neural network. To evaluate the conceived models, they compared these with a Multi-variate Linear Regression (MLR) model. The data were based on a WWTP located in China, with daily periodicity. The data collected corresponded to 360 records, between December 2015 and December 2016, with six invalid records removed. To decide which features were given as input to the model, the authors used the Fuzzy C-Means (FCM) method. This method identified three indicators: the influential charge, the Chemical Oxygen Demand (COD), and the total nitrogen removed. The authors defined the FCM hyperparameters without any search for the best value for each of them, such as the number of iterations or clusters. Each of these selected indicators was used, one at a time, as input to the RBF model. The authors used min relative error, max relative error, and mean absolute percentage error (MAPE) for performance measurement metrics. In total, the authors developed four models with different inputs, three of them for each set of selected indicators and another with the total data. Using only data from each indicator's subset, the RBF model performed better than the MLR model. On the contrary, the MLRM model performed better when using the total dataset as input. Overall, both models performed better when using only the data subset of the indicators.

Harrou et al. [10] conducted a study to make short-term forecasts of energy consumption in a WWTP, using statistical and Deep Learning (DL) models. The data used in this study are between 2010 and 2017, belonging to a WWTP in Saudi Arabia. In total, the authors used six statistical models, such as the Auto-regressive Integrated Moving Average (ARIMA) or the Ordinary Least Square (OLS). Two types of networks were based on DL models, Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRU). The models conceived by the authors used a uni-variate approach, where only the feature they intend to forecast, the energy consumption, is given as input to the different candidate models. The data were normalised between 0 and 1 for all conceived models. There was no particular attention to the case of LSTM networks working internally with a hyperbolic tangent. Throughout the manuscript, no cross-validation or overfitting control techniques are mentioned in the conceived models. Regarding the evaluation metrics of the models, the authors used four, i.e., MAPE, Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Root Mean Squared Log Error (RMSLE). By observing the obtained results, the authors verified that the statistical-based models slightly outperformed the DL models, with ARIMA getting a MAPE of 2.29%, while the best DL model, LSTMs, presented a MAPE of 2.42%. The authors also verified that the models' parameters were updated recursively, given a better performance than the models with no updates. However, they concluded

that the DL models could provide forecast results with more significant performance when applying more data. No reference was made regarding fitting times.

The study carried out by Huang et al. [11] had as objective the construction of an energy consumption model in a WWTP based on Elman Neural Network-Energy Consumption Model (ENN-ECM) to identify the relationship between energy consumption and the quality of the effluent. The benchmark simulation model (BSM1) was used to compare the authors' model results. Both models were based on data related to an activated sludge model, being obtained from BSM1, which provided data for a period of two weeks in 15-minutes time intervals. Firstly, the authors used the energy consumption model to verify which effluent characteristics had a more significant relationship with the characteristics related to energy consumption. Then, they implemented the ENN-ECM with five characteristics of the effluent obtained from the energy consumption model to forecast four energy consumption parameters. The network architecture, namely the number of layers, was obtained through empirical formulas and the Kolmogorov theorem. The authors concluded that the ENN-ECM model obtained better performance concerning energy consumption with the analysis of the obtained results.

Ramli et al. [12] conducted a study to forecast energy consumption in a WWTP in Malaysia using an ARIMA model. To compare the obtained results, the authors used a linear regression method. The data used in this study were based on four years of active power in the WWTP. To achieve the best ARIMA model, the authors used the Time Series Modeler, incorporated in the SPSS software, obtaining the values (0, 1, 0) for ARIMA's parameters. The results allowed the authors to verify that the ARIMA model obtained better performance than the linear regression with an RMSE of 55.59, compared to 67.51, respectively. The authors further concluded that it was possible to increase energy efficiency by 10% of energy recovery, which could reduce the cost of electricity in the studied WWTP.

Another study carried out by Maki et al. [13] aimed to forecast the total energy consumption of a WWTP and the consumption in different processes, using a Markov switching model. The data were collected by applying several sensors connected to a WWTP energy distribution network in Japan and transmitted over a 3G line. The data collection was carried out between March 2015, and March 2017, with a 1-min periodicity. The authors then grouped the data into an hourly periodicity. In addition to the forecast of total energy consumption at the WWTP, the authors also forecast the energy consumption in the water treatment, sludge treatment, and auxiliary facilities processes. Additionally, as the sum of the three identified processes' energy consumed did not coincide with the total energy consumed in the WWTP, they made the forecast for the remaining operations, marked as "others". An analysis was made of energy consumption over time, where it was possible to verify that there is greater energy consumption in summer than winter. In addition to the data collected by the sensors, the authors added six more features to be used in the conceived model: holidays, office hours, temperature, humidity, wind speed, and the previous five hours of energy consumption. Only 1 week was considered as input. With the obtained results, the authors found that, except for the sludge treatment and auxiliary facilities, the values were below 10%. Besides, the relationships between the variables that affect the energy consumption forecast equation were verified in each process. The authors then concluded that an increase in the WWTP's energy consumption, together with the increase in seasonal temperatures, leads to a rise between 0.1% and 0.2% for each 1 °C in temperature.

Oulebsir et al. [14] conducted a study where they conceived an Artificial Neural Network (ANN) to create an energy consumption model in a WWTP using the active sludge process. The authors used data provided by a WWTP in Algeria between January 2006 and March 2016. In this study, the authors use four parameters: (1) the Biological Oxygen Demand ( $BOD_5$ ), (2) the COD, (3) suspended solids and (4) ammonium. In addition, they also use the water temperature, and flow of the influent, the flow of recirculated sludge, and the total consumed energy. The authors applied a set of methods to clean the dataset, keeping 318 days of observations even though the original dataset had 10 years of data. The

different ANNs had six hidden layers with a total of 200 neurons each. The architecture of the models was established using the trial-and-error method. In each conceived model, data were divided into 80% for training and 20% for testing, without using a time series cross-validator. The authors confirmed that the pollution load contributes more significantly to forecasting energy consumption than the removal efficiency. The authors also applied the k-means algorithm, observing three clusters. The authors were thus able to verify three classes of energy consumption: under-consumption, over-consumption, and optimal consumption.

As an overall conclusion, it can be said that some studies have already considered the use of DL models to forecast energy consumption in a WWTP. Typically, studies follow a single-step approach, i.e., they only forecast consumption value for the next day. Furthermore, it is usual to find studies that do not consider certain aspects of time series problems, such as using an appropriate cross-validator, not breaking the time series when removing missing values or missing timesteps, or even when searching the best hyperparameters. In addition, it is not easy to understand the existence of overfitting as learning curves are not analysed. All this may lead to significant problems when deploying the best candidate model in a real-life scenario.

### *1.2. Goals, Research Questions, and Paper Structure*

This work aims to conceive, tune, and evaluate a set of candidate DL models to forecast energy consumption in a WWTP, going from recurrent to convolutional candidates. In addition, the goal is to implement a recursive multi-step approach to forecast the next two days, providing a stronger understanding of future patterns. We also aim to experiment two different methods for outliers' handling and the performance of the candidates in uni- and multi-variate settings. Then, as last goal, we aim to evaluate the best candidate model in a WWTP with a low volume of data. For that, we are required to apply transfer learning processes, overcoming the problem of data scarcity.

This study uses data provided by a Portuguese water company. The elicited goals can be translated into the following research questions:

1. Do recurrent neural networks have a better performance when forecasting energy consumption in a WWTP than convolutional ones?
2. Which features facilitate the process of forecasting energy consumption in a WWTP?
3. Is it possible to apply transfer learning processes, with the goal being to use a pre-trained model to forecast the energy consumption of a WWTP with low volumes of data?

The remainder of this manuscript is structured in three more sections. Section 2 describes the materials and methods, namely the collection, exploration, and pre-processing of data, the developed DL models, and the conducted experiments. Section 3 is responsible for summarising the obtained results, as well as their interpretation. Finally, Section 4 discusses the obtained results and gathers the conclusions drawn from this study.

## **2. Materials and Methods**

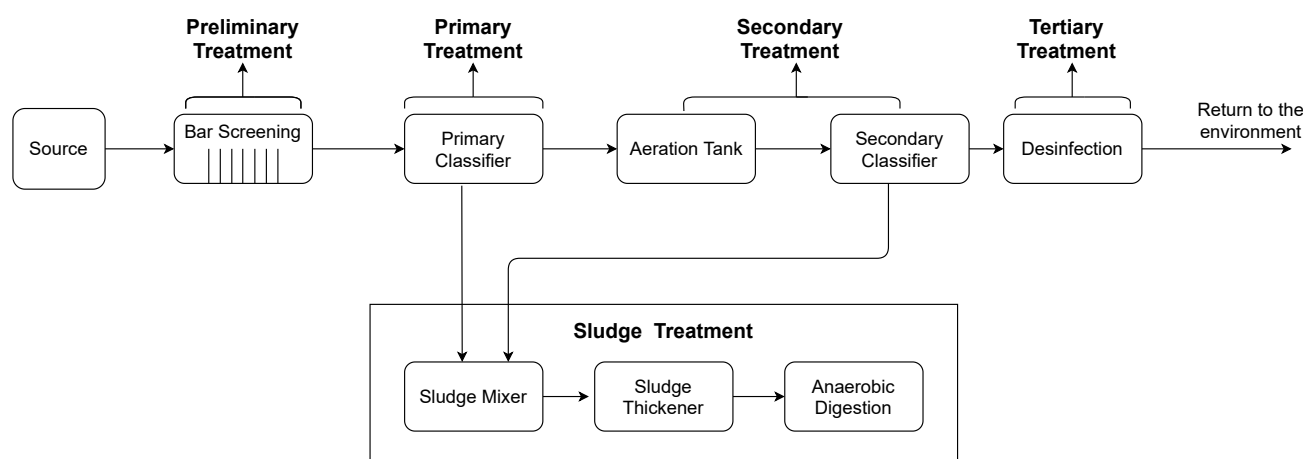
The following lines describe the materials and methods used throughout this study, including collecting, exploring, and treating data. Additionally, the models used throughout the work are described, as well as the evaluation metrics, the used technologies, and the designed experiments.

### *2.1. Dataset*

The data used in this study took into account three different datasets. Dataset one was related to energy consumption while the second dataset described the volume of the flow of water at the entrance of a WWTP. The third dataset described the climatological conditions. The first two datasets were made available by a Portuguese wastewater company and were related to a single WWTP. Regarding the energy consumption value, which is the target feature, there is an intrinsic relationship between the different processes present in a WWTP and the required energy (typically, the larger the WWTP, the greater its energy

consumption). However, this relation was captured and described in the time series in itself as the values were a snapshot of the state of the WWTP. The third dataset was collected using the Open Weather Map API, and contains climatological data regarding the same city where the WWTP was located. All datasets contained observations belonging to the period between January 2016 to May 2020.

Figure 1 illustrates the WWTP layout used in this study. This WWTP was based on four main stages: preliminary, primary, secondary and tertiary treatments. In addition, there was also a line responsible for the sludge treatment. The preliminary treatment, which included bar screening, was accountable for removing solids and materials of greater volume, an essential step in the WWTP process since some of these objects could damage some equipment in the following steps. The primary treatment, which included the primary classifier, aimed to remove the smaller volume solids, namely the suspended solids, from the previous stage and the organic matter present. In the secondary treatment, two processes were included, the aeration tank and the secondary classifier. This stage aimed to remove biodegradable organic matter from wastewater, in addition to suspended solids and nutrients, such as nitrogen. Finally, the tertiary treatment was responsible for removing the remaining suspended solids resulting from the previous stages. The sludge produced in the primary and secondary treatment was inserted in the sludge treatment line. This line was responsible for dewatering and disinfecting the sludge, reusing it as an energy source.



**Figure 1.** WWTP (Wastewater Treatment Plants) layout.

### 2.1.1. Data Exploration

The energy consumption dataset comprised two features: the energy consumption value (in kWh) and the corresponding timestamp, making 1522 records with a daily periodicity. The influent flow dataset also contained two features, i.e., the value of the influent flow (in m<sup>3</sup>) and the timestamp, with a total of 1535 records, again with a daily periodicity. Finally, the climatological dataset had a total of 25 features, including the timestamp, air temperature, and humidity, among others, with a total of 38,651 hourly timesteps. Table 1 presents the different features available in the three datasets, detailing its characteristics and presenting the corresponding units of measure.

None of the three datasets had missing values. However, as in its genesis the problem identified in this study was based on a time series problem, it was essential to pay attention to missing timesteps. In the case of the climatological dataset, there were no missing timesteps. On the contrary, both the energy consumption and the influent inflow datasets contained missing timesteps. In the former, there were 88 missing timesteps, while in the latter 75 missing timesteps were identified. In a subsequent section, it is explained how to overcome the missing timesteps problem.

As the main goal of this study was to forecast energy consumption, data exploration emphasized the *value\_energy* feature of the energy consumption dataset. Firstly, it is worth mentioning that this feature presented an accumulated value. Hence, it was necessary to subtract, from each observation, the value of the previous one, in order to obtain its real value. Since the first observation had no previous one, it was removed. A box plot analysis allowed us to identify the existence of some extreme outliers that were derived from an incorrect insertion of values by the operators of the WWTP.

**Table 1.** Features available in the used datasets. Only the main features of the climatological dataset are presented.

#	Features	Description	Unit
<i>Energy Consumption Dataset</i>			
1	<i>date</i>	Timestamp	date and time
2	<i>value_energy</i>	Total energy consumption	kWh
<i>Influent Flow Dataset</i>			
1	<i>date</i>	Timestamp	date and time
2	<i>flow</i>	Accumulated influent flow value	m <sup>3</sup>
<i>Climatological Dataset</i>			
1	<i>dt_iso</i>	Timestamp	date and time
2	<i>temp</i>	Temperature	°C
3	<i>feels_like</i>	Human perception of climate	°C
4	<i>temp_min</i>	Minimum temperature	°C
5	<i>temp_max</i>	Maximum temperature	°C
6	<i>pressure</i>	Atmospheric pressure	hPa
7	<i>humidity</i>	Humidity percentage	%
8	<i>wind_speed</i>	Wind speed value	m/s
9	<i>wind_deg</i>	Wind direction	Degrees
10	<i>rain</i>	Rain volume	mm
11	<i>clouds_all</i>	Cloudiness percentage	%

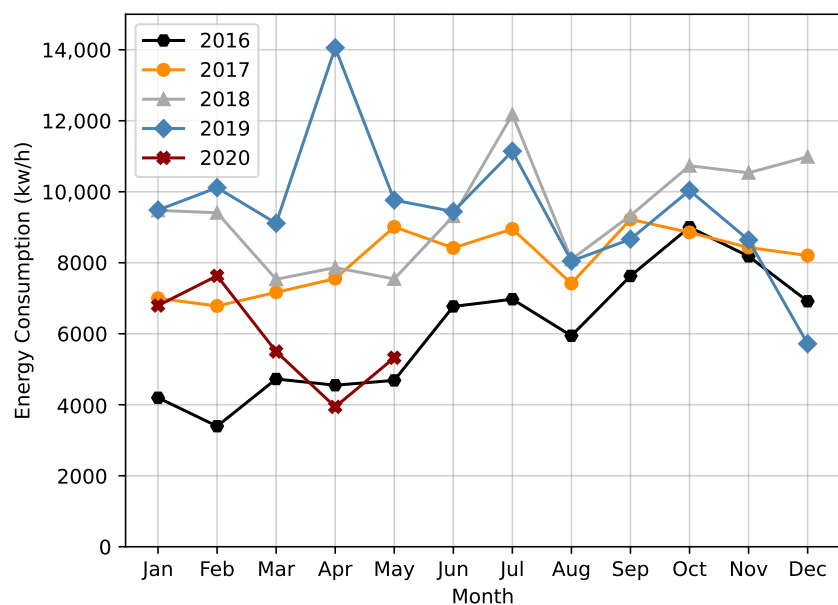
A statistical analysis of the energy consumption values was performed, being described in Table 2. It was possible to verify that the mean energy consumption value in the dataset presents a value of 8050.96 kWh, with a standard deviation of 3736.359 kWh. The skewness was 3.172, representing an asymmetric distribution, i.e., the positive value indicates a positive inclination in the distribution of the data, in which the tail size of the right hand is larger than that of the left. Regarding the kurtosis value, it was 28.101. A kurtosis value greater than 1 indicates that the distribution of energy consumption has a very high peak (a leptokurtic distribution).

**Table 2.** Descriptive statistics for energy consumption.

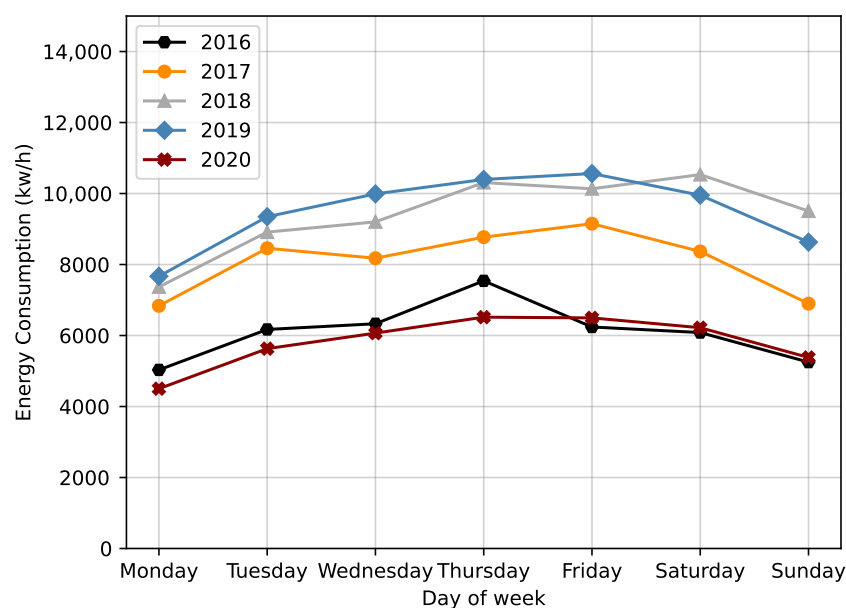
Number of Items	Mean	Median	Std. Deviation	Skewness	Kurtosis
1522	8050.960	7689	3736.359	3.172	28.101

We then explored the energy consumption over the months of a year, during the 5 years present in the dataset. In Figure 2 it is possible to verify a pattern in all the explored years, with a constant drop in energy consumption between July and August.

Another analysis took into account the variation in energy consumption over the different days of the week. This analysis was based on the mean value of the days of the week for each year. As shown in Figure 3, it is possible to verify that Sunday and Monday were the days when there was less energy consumption in the WWTP. In conclusion, it appears that the traditional working days had a higher energy consumption on average, while on weekends there was a decrease.

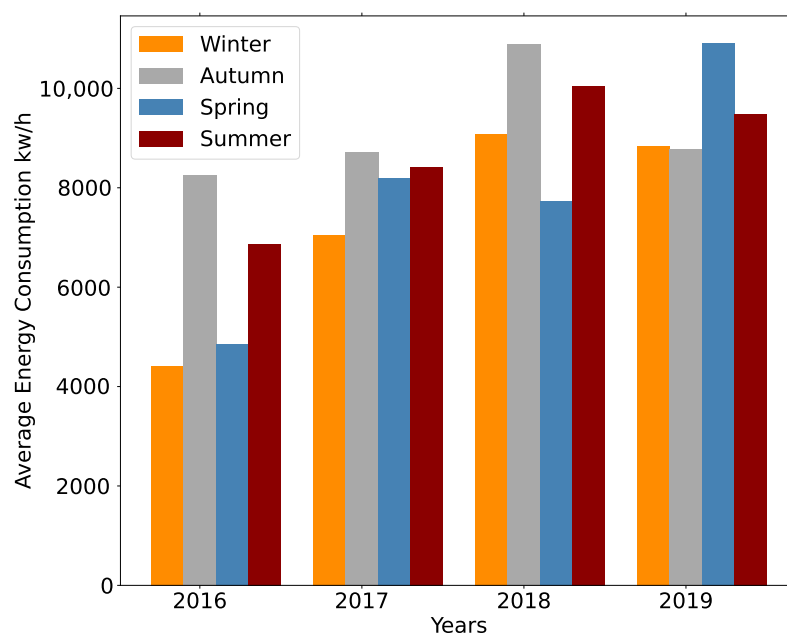


**Figure 2.** Monthly variation of energy consumption over the years present in the dataset.



**Figure 3.** Day of the week variation of energy consumption over the years.

To understand seasonality, we performed two different analyses on the energy consumption data between 2016 and 2019: the first relative to the average consumption by season and the second related to the energy consumption per trimester. Figure 4 depicts the first analysis, being possible to verify that, typically, more energy was consumed during the autumn. Interestingly, in 2019, autumn was the season with the lowest average energy consumption value. In general, it was also possible to see that over the years, energy consumption was rising in different seasons. Despite a higher number of average consumption values, it was not in the autumn that the highest average peak was reached, but in the spring of 2019 with a value of 10,912 kWh. Regarding the lowest peak, it occurred in the winter of 2016, with a value of 4398 kWh. Additionally, it was possible to verify that, in general, winter was the season with less consumption of energy.



**Figure 4.** Mean energy consumption per seasons of the year.

The trimesters analysis showed that the fourth trimester had the highest energy consumption values over the first three years. Despite this, the highest value was verified in the second trimester of 2019, with 11,072 kWh. As demonstrated in the seasons' analysis, in general, the average values increased during the first three years. In 2019, there was an increase in the first and second trimester and a decrease in the third and fourth ones.

Regarding the influent flow, an analysis was carried out considering the average for each year, described in Table 3. As can be seen, 2019 was the year with the highest volume of influent flow on the WWTP (1155.33 m<sup>3</sup>). Interestingly, checking the year of 2019 concerning the energy consumption (Figure 2), we verified that this year also obtained, in general, the highest average of energy. On the other hand, looking at 2016, excluding the incomplete year of 2020, this was where the lowest average influent flow value occurred, this being, in general, the year with the lowest energy consumption value.

**Table 3.** Average influent flow per year.

Year	Value (m <sup>3</sup> )
2016	910.69
2017	1025.23
2018	981.26
2019	1155.33
2020	849.19

#### 2.1.2. Data Preparation

The first step to prepare the data were to carry out a feature engineering process in the three datasets, thus creating three new features from the timestamps (i.e., *year*, *month*, and *day*). The dataset related to climatological data, as mentioned, had an hourly periodicity, so to match the same periodicity as the other datasets, these were grouped by day, month and year, aggregating the mean value per feature.

As referred above, as both the energy consumption and influent flow datasets presented accumulated values, a method was applied to obtain the value that would correspond to each specific day. The identified extreme outliers, which corresponded to miss



insertions of values by the operators of the WWTP (for example, extra digits), were also solved. The remainder of the data treatment is specified in the following lines.

#### Handling Missing Timesteps

To deal with the missing timesteps verified in the energy consumption and the influent flow datasets, a dataset was created comprising all days (i.e., timesteps) that should have been present in the dataset. In both cases, the start date was 2nd January 2016 and the end date 28 May 2020. The datasets were joined, with missing timesteps being added and having its features filled with the  $-99$  value. Solving the missing timesteps problem created a new one, missing values, i.e., timesteps that were missing were now present but all their features had the  $-99$  value.

#### Handling Missing Values

To fill the missing values, a queue-based approach was followed. Each record was read for each of the two datasets with missing values, saving its value (energy consumption or influent flow) in the mentioned structure, with a maximum size of eight values. Whenever reading a record, if the queue was full, a push operation would be performed at the beginning of the queue. When a timestep had a feature with the  $-99$  value, its value would be computed based on the average of the last eight records, i.e., the previous 8 days, present in the queue. Once calculated, this value would then be pushed to the queue, eliminating the oldest record. By the end of this process, no dataset had missing values neither missing timesteps.

#### Joining Datasets

When reaching this point, each one of the three datasets was made of 1609 observations. However, we were required to join the three datasets into a single one. This was performed using the features *year*, *month*, and *day*. In the end, a single dataset was created, having 1609 observations with 30 features each.

#### Correlation Analysis

To verify which features had a more significant correlation with the target feature (*value\_energy*), it was first necessary to check whether the data followed a normal distribution. Using a  $p < 0.05$  and the Kolmogorov–Smirnov test, it was possible to verify that all features assumed a non-Gaussian distribution. Hence, it was necessary to use the non-parametric Spearman’s rank correlation coefficient, being possible to verify that the features that had a more significant correlation with the target were the *year*, *month*, *temperature*, and *flow\_value*. Since the other features had a low correlation with the target, they were removed. After this treatment, the final dataset had 1609 observations with a shape (1609, 5). Table 4 shows an example of a record in the final dataset.

**Table 4.** Features present in the final dataset.

#	Features	Observation Example
1	<i>year</i>	2018
2	<i>month</i>	5
3	<i>temperature</i>	11.96
4	<i>flow_value</i>	829
5	<i>value_energy</i>	5155

#### Handling Outliers

Extreme outliers were above 14,000 kWh. Only six observations were below 2000 kWh. Since the range between the maximum and minimum values for the feature *value\_energy* was large, and considering the reduced amount of observations that were causing it, two different methods were experimented to handle outliers. These two methods provided a

comparative term for the different experiments, causing slight modifications to the input data that were fed to the models. The two methods were as follows:

- Method 1—to further reduce the amplitude of the target feature, the few timesteps with *value\_energy* greater than 10,000 kWh or lower than 2000 kWh had their value updated, using the queue-based approach described above. The goal was to use interpolation to replace the outliers;
- Method 2—to further reduce the amplitude of the target feature, the few timesteps with *value\_energy* greater than 10,000 kWh or lower than 2000 kWh had their value truncated. The goal was not to use interpolation to update the target value.

#### Normalisation

With the data prepared, the next step was to normalize them. Since LSTMs work internally with the hyperbolic tangent, we decided that the applied normalization would be in the range  $[-1, 1]$ , according to the following equation:

$$\frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (1)$$

#### Supervised Problem

The final step was to go from an unsupervised problem to a supervised one, with the respective inputs (X) and corresponding labels (y). Thus, it was necessary to create sequences of data, which depend on the number of timesteps used as input for the models. A sliding window was used over the initial dataset to create the different sequences and the respective labels, thus creating a set of sequences that can be fed to the models. As an example, if the shape of a model's input was (1601, 7, 5), the first element set the number of samples, the second the number of input timesteps, and the last the number of features. In this example, the labels would have the shape (1601, 1). A similar algorithm can be seen in the work of Fernandes et al. [15].

#### 2.2. Model Conception

To achieve the objective of forecasting energy consumption in a WWTP, three different DL models were conceived and evaluated, namely LSTMs, GRUs, and uni-dimensional Convolutional Neural Networks (CNNs). Regarding the choice of models, concerning the LSTM and GRU models, these were selected since they belong to the set of Recurrent Neural Networks (RNNs), which has shown an outstanding performance in time series problems. While traditional ANNs cannot remember what they learned in previous iterations, RNNs can learn from earlier timesteps [16–19]. Regarding the choice of CNNs as the third model to be used, despite its greater use in image processing, it has shown promising results in terms of time series problems when using uni-dimensional convolutions [20–23].

To find the best combination of hyperparameters, two error metrics were used. The RMSE is an error measure, as it measures the difference between the values predicted by the model ( $\hat{y}$ ) and the true values observed ( $y$ ). RMSE equation is as follows:

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}} \quad (2)$$

The second metric, the MAE, is the mean of the differences between predicted and observed values. Its use is mainly to complement and strengthen the confidence on the obtained values. Its equation is as follows:

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (3)$$

### 2.2.1. LSTMs

One of the models used in this study was based on a particular RNN, i.e., LSTMs. RNNs are a type of network that, unlike ANNs, can have as input the current input and pay attention to past inputs [24,25]. In other words, the decision taken on the timestep  $t - 1$  will affect the timestep  $t$ . LSTMs, introduced in 1997 by Hochreiter and Schmidhuber [26], can learn temporal dependencies over a long period, in addition to the short term. These networks came to fill an existing problem in RNNs, where there was an exponential drop in the backpropagated error in long periods. Nowadays, LSTMs are widely used in forecasting problems, such as in road traffic or weather, and their use in detecting anomalies in time series problems [27–31].

Regarding the architecture of LSTMs, it consists of multiple memory cells. There are two states in each of these memory cells: the hidden state and the cell state. The hidden state, already present in RNNs, is responsible for short-term memory, while on the other hand, the cell state (not present in RNNs) has the capacity for long-term memory. Additionally, each memory cell has internal gates, which allow a LSTM to forget ( $f_t$ ), include ( $i_t$ ), and output ( $o_t$ ) information [26]. The following equations describe the calculation performed on each of the gates.

$$i_t = \sigma(w_i[h_{t-1}, x_t] + b_i) \quad (4)$$

$$f_t = \sigma(w_f[h_{t-1}, x_t] + b_f) \quad (5)$$

$$o_t = \sigma(w_o[h_{t-1}, x_t] + b_o) \quad (6)$$

where  $\sigma$  represents the sigmoid function,  $w_x$  the weight for the respective gate,  $h_{t-1}$  the output of the previous block,  $x_t$  input at current timestep and  $b_x$  the biases for the respective gate.

First, through the sigmoid layer, it is necessary to decide which information will leave the cell state (forget gate) and remain the same. The action on what will keep information is divided into two stages, the first deciding which values should be updated through another sigmoid layer (input gate) and the second creating a vector of new deals that can add to the state through a hyperbolic tangent layer. The next cell state update is obtained through a point multiplication operation on the two previous steps results. Finally, the output is decided using a sigmoid layer (output gate) followed by a hyperbolic tangent one [26]. Figure 5 provides a graphical view of such a memory cell. The following equations describe the calculation of the cell state, the candidate cell state and the final output.

$$\tilde{c}_t = \tanh(w_c[h_{t-1}, x_t] + b_c) \quad (7)$$

$$c_t = f_t \times c_{t-1} + i_t \times \tilde{c}_t \quad (8)$$

$$h_t = o_t \times \tanh(c_t) \quad (9)$$

where  $c_t$  represents the cell state at timestep  $t$  and  $\tilde{c}_t$  represents the candidate for cell state at timestep  $t$ .

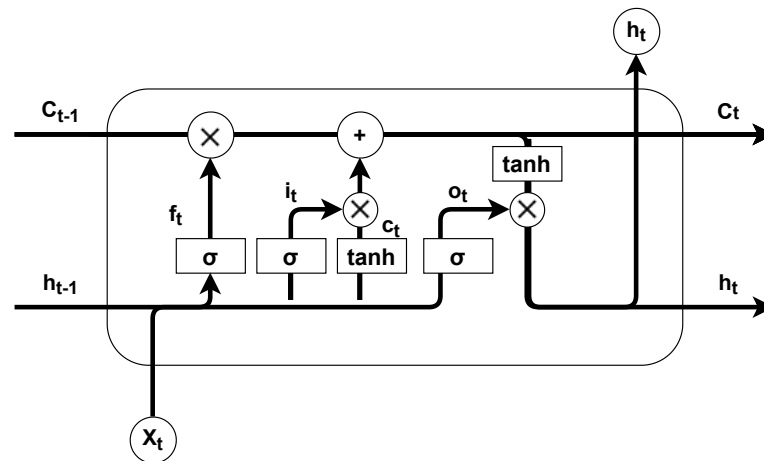


Figure 5. Architecture of a LSTM (Long Short-Term Memory) cell.

### 2.2.2. GRUs

Another model used in this study was the GRU. These networks are a subtype of RNNs, introduced in 2014 by Kyunghyun Cho [32]. Like LSTMs, GRUs were developed to solve the vanishing gradient problem of RNNs. GRUs are a simpler version of LSTMs, and they can be faster than these, obtaining similar performance. Unlike LSTMs, GRU cells only have the hidden state, which can maintain long and short term dependencies, thus eliminating the LSTM cell state. Another difference is that GRUs only have two layers of neural networks and have only two gates: reset ( $r_t$ ) and update ( $z_t$ ) [33]. The following equations describe the calculation performed on each of the gates.

$$z_t = \sigma(w_z \cdot [h_{t-1}, x_t]) \quad (10)$$

$$r_t = \sigma(w_r \cdot [h_{t-1}, x_t]) \quad (11)$$

The first step performed in a GRU cell is to represent the information removed by a sigmoid layer, from the previous hidden states, through the reset gate, working in a very similar way to the LSTM forget gate. Then, through the update gate, the amount of information from the previous timesteps is decided to be transmitted to the next state through a sigmoid layer. The next step uses the reset gate, applying a hyperbolic tangent layer, to introduce a new memory content, called the hidden state candidate. Finally, the update gate effect is incorporated to create the new hidden state [33]. GRUs are, like LSTM, widely used in forecasting problems in time series [34–36]. Figure 6 provides a graphical view of a GRU cell. The following equations describe the calculation of the current memory content and the final memory at current time step.

$$\tilde{h}_t = \tanh(w \cdot [r_t \times h_{t-1}, x_t]) \quad (12)$$

$$h_t = (1 - z_t) \times h_{t-1} + z_t \times \tilde{h}_t \quad (13)$$

where  $\tilde{h}_t$  represents the current memory cell and  $h_t$  the vector which holds information for the current unit.

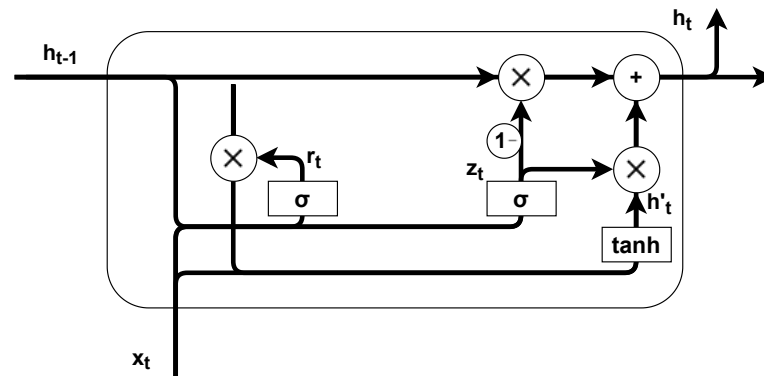


Figure 6. Architecture of a GRU (Gated Recurrent Units) cell.

### 2.2.3. CNNs

The last model used in this study was a CNN, a type of neural network developed a few decades ago [37,38]. Its appearance was based on a survey carried out by Hubel and Wiesel, in 1962, on the visual cortex of cats [39]. Over the past few years, CNNs has been closely linked to the classification of images and object detection [40,41]. In general, CNNs have a set of essential aspects: the convolutional layer, the pooling layer, and the fully connected one. Based on an image as an input, the convolutional layer is responsible for dividing the image's features, while the fully connected layer uses the output of the convolutional layer to classify. The pooling layer is used to reduce the amount of information coming from the convolutional one.

Recent times came with the use of CNNs for time series problems, mainly using uni-dimensional ones [21–23]. In the context of a time series problem, a significant aspect that needs to be taken into account is the approach being followed in terms of the data format, i.e., whether channels' last or channels' first. Concerning channels' last, this approach aims to reduce the number of timesteps while keeping the number of filters intact. On the other hand, the channels' first approach does just the opposite, i.e., reduces the number of filters and keeps the number of timesteps intact. Depending on the followed approach, this will always cause differences in the convolutional layer, which has the format (*timesteps, filters*). The kernel size is yet another parameter responsible for defining the timesteps window length that is affected by each filter. An illustrative example of a channels' last approach can be seen in work of Oliveira et al. [23]. Finally, the form of calculating the shape of the output follows the following equation:

$$(\text{Timesteps} - \text{KernelSize}) + 1 \quad (14)$$

## 2.3. Experiments

Several experiments were carried out, taking into account different scenarios as shown in the next lines. The same random seed (91195003) was used in all conducted experiments.

### 2.3.1. Technologies

For data exploration, the *Knime* platform was used as well as the Python programming language, version 3.7. Python was also used for data pre-processing and for the development and evaluation of the DL models. *Pandas*, *NumPy*, *scikit-learn*, and *matplotlib* were the used libraries. In addition to these, *TensorFlow v2.0.0* was used to develop the models. Regarding the hardware, all of it was made available by Google's Colaboratory.

### 2.3.2. Experimental Setup

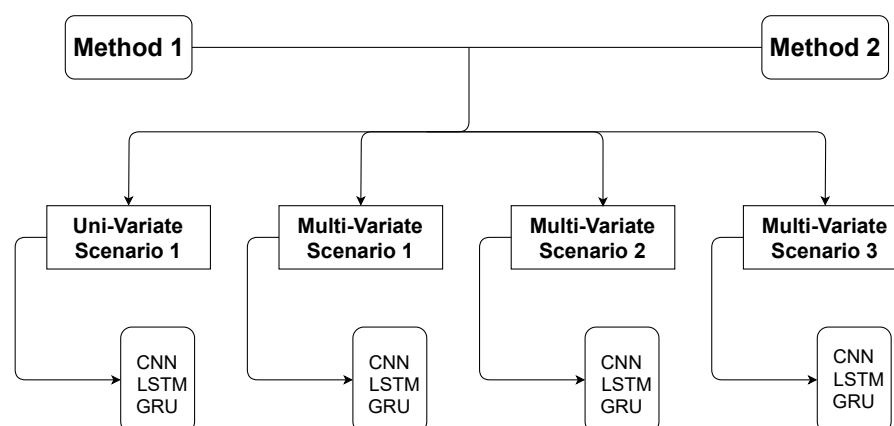
To achieve the goal of forecasting the energy consumption of a WWTP, it was necessary to evaluate multiple candidate models. All candidates were designed to follow a recursive multi-step approach, i.e., to forecast energy consumption for the next 2 days.

For each type of DL model used in this study, candidate models were designed based on an uni-variate and multi-variate approach. In the case of being uni-variate, the models would only receive, as input, the *value\_energy* feature. In the multi-variate approach, three distinct scenarios were defined, with each scenario consisting in a different set of features. Table 5 summarises the features that each scenario contains. These scenarios are useful to understand the importance of temporal and climatological context data in the energy consumption of WWTPs. The influent flow is included in all multi-variate scenarios, since it had the highest correlation coefficient with the target feature.

**Table 5.** Uni- and multi-variate data scenarios.

Uni-Variate	
Scenario 1	<i>value_energy</i>
Multi-Variate	
Scenario 1	<i>value_energy, year, month, temperature, flow_value</i>
Scenario 2	<i>value_energy, temperature, flow_value</i>
Scenario 3	<i>value_energy, flow_value</i>

Two distinct datasets were built, one for each outliers' method. For each method, two approaches were followed: uni- and multi-variate. Then, for each approach, a set of scenarios were defined. Figure 7 sets the different combinations of data used to fit and evaluate the candidate models.



**Figure 7.** Different combinations for the conception of the candidate models.

The search for the best hyperparameters' configuration was performed using grid search. This method was applied to tune parameters such as the model architecture, batch size, or the number of timesteps that make an input sequence. Table 6 describes the hyperparameters' searching space considered for each model type. Besides, two callbacks were defined over the validation's loss. One aimed to automatically reduce the learning rate, while the other stopped the training when the RMSE stopped improving.

To prevent overfitting and underfitting situations, learning curves were plotted, stored, and analyzed. It should also be noted, taking into account that we were facing a time series problem, that a time series cross-validator was used ( $k = 3$ ), namely the *TimeSeriesSplit* API of scikit-learn. This cross-validator, unlike traditional ones, had successive training sets as supersets of those that came before. Each of these training sets was further split into training and validation sets.

**Table 6.** Hyperparameters searching space.

Parameter	LSTM and GRU	CNN
Layers	[3, 4, 5]	[3, 4, 5]
Neurons	[32, 64, 128]	-
Activation	[ReLU, Tanh]	[ReLU, Tanh]
Timesteps	[14, 21, 28]	[14, 21, 28]
Batch Size	[5, 10, 20]	[5, 10, 20]
Dropout	[0.0, 0.5]	[0.0, 0.5]
Kernel Size	-	[3, 4, 5]
Filters	-	[16, 32]
Pool Size	-	[2, 3]

### 3. Results

Several hundred experiments were run in order to evaluate all possible candidate models. The candidates were evaluated considering their RMSE and MAE.

#### 3.1. Method 1

The first method had the outliers updated as per the conceived queue-based approach. Table 7 presents the best hyperparameter configurations for each combination in this method. Within these combinations, it was possible to verify that the best one concerned CNNs for the third multi-variate scenario, with a MAE of 630 and a RMSE of 690 kWh.

**Table 7.** Best results, per scenario, for Method 1. The letters stand as follows: a. timesteps; b. batch size; c. number of layers; d. number of neurons/filters; e. pool size; f. kernel size; g. dropout; h. activation; i. RMSE; j. MAE; k. time (s).

Model	a.	b.	c.	d.	e.	f.	g.	h.	i.	j.	k.
<i>Uni-Variate-Scenario 1</i>											
CNN	14	10	3	32	3	3	0.0	tanh	702.71	645.70	33
LSTM	21	20	3	32	-	-	0.5	tanh	779.13	714.58	36
GRU	21	20	5	64	-	-	0.0	ReLU	715.42	653.75	78
<i>Multi-Variate-Scenario 1</i>											
CNN	21	20	5	32	3	3	0.5	ReLU	737.47	677.48	32
LSTM	21	5	4	64	-	-	0.0	tanh	788.46	720.77	175
GRU	14	10	3	32	-	-	0.5	tanh	755.56	693.78	73
<i>Multi-Variate-Scenario 2</i>											
CNN	21	20	4	16	3	3	0.0	ReLU	742.35	684.92	19
LSTM	21	5	3	64	-	-	0.0	ReLU	760.75	699.45	136
GRU	28	20	3	64	-	-	0.0	ReLU	727.07	670.53	50
<i>Multi-Variate-Scenario 3</i>											
CNN	28	20	4	32	3	3	0.5	ReLU	690.00	630.63	27
LSTM	21	20	4	128	-	-	0.5	ReLU	729.73	668.62	91
GRU	21	20	3	32	-	-	0.5	ReLU	746.98	683.02	38

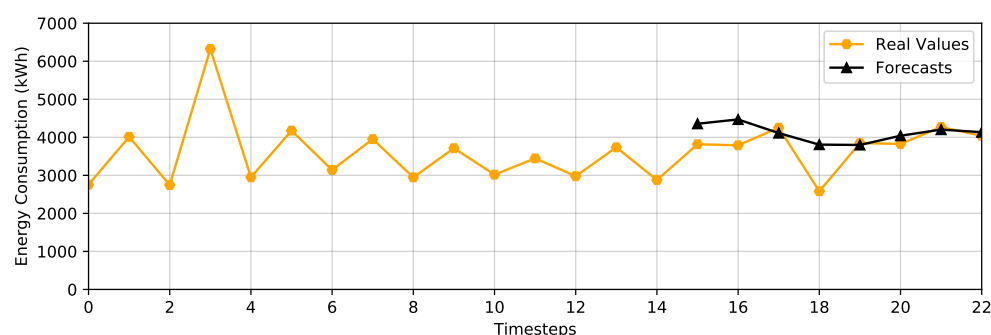
Regarding the uni-variate approach, it was possible to verify some differences between some hyperparameters between the RNN-based models and the CNN-based model. Concerning the number of timesteps and the value of the batch size, it appeared that the CNN-based had the lowest value of both models, 14 and 10, respectively. Regarding the number of layers and the number of neurons/filters, the GRU-based model presented the highest values of 3, 5 and 64, respectively.

Overall, CNN candidates showed better results in all uni- and multi-variate scenarios, except for the second scenario, where GRUs presented a better performance. Regarding

the training times, CNNs candidate models demonstrated lower values than the other two, with LSTM-based models being the ones taking more time to fit. It was also notable that the number of timesteps given as input increased, in general, with the number of features provided to the model. Concerning the activation function, it was also possible to verify that there was a tendency to use *tanh* in the uni-variate approach, while in the multi-variate, the best candidate models tended to use *ReLU*.

The best multi-variate scenario is the one that added, to the *value\_energy* feature, the *flow\_value*, i.e., the influent flow value combined with energy consumption value. In this approach, it was possible to verify that approach in terms of the number of timesteps, in Scenario 1 and Scenario 2, the model based on LSTM and the model based on CNN presented the same value in both cases, 21 timesteps (3 weeks). Regarding the batch size, note that the LSTM-based model in Scenarios 1 and 2 had a lower value than the others, while in Scenario 3, both models had the same value (20). It was also possible to verify that the best candidate models had a better performance with climatological context and without temporal context, except for CNN-based models. On the other hand, GRU-based models had their the best performance in the uni-variate approach, while the other two models presented their best performance in the multi-variate approach, more specifically in Scenario 3 (*value\_energy* and *flow\_value* features).

Figure 8 plots eight multi-step forecasts for the best candidate model in this method (the best CNN candidate in the third multi-variate scenario). These forecasts describe a set of 28 timesteps (i.e., days) given as input, making a successive two-day forecast for a total of 8 days.



**Figure 8.** Eight multi-step forecasts for the best candidate model in Method 1.

### 3.2. Method 2

The second method used a dataset that had the outliers truncated. Table 8 depicts the best hyperparameter configuration for each combination of this method, with the best candidate, a CNN, following a uni-variate approach and presenting a MAE of 784 and a RMSE of 869 kWh. This meant that when truncating the outliers, a uni-variate approach presented better results than a multi-variate one.

As in Method 1, the CNN-based models presented a training time shorter than the others. It was also possible to verify that the CNN-based models had better performance. These models show an interesting uniformity in the cardinality of timesteps, while in the other models there was a higher fluctuation. Regarding the number of layers, it was possible to verify a constant value in most models (three layers), except for two CNN-based models.

In the uni-variate approach, it was possible to verify that the model based on CNN presented a lower value of timesteps given as input to the model (14) than models based on RNN. On the other hand, regarding the batch size value, the CNN-based model presented a higher value than the others (30).

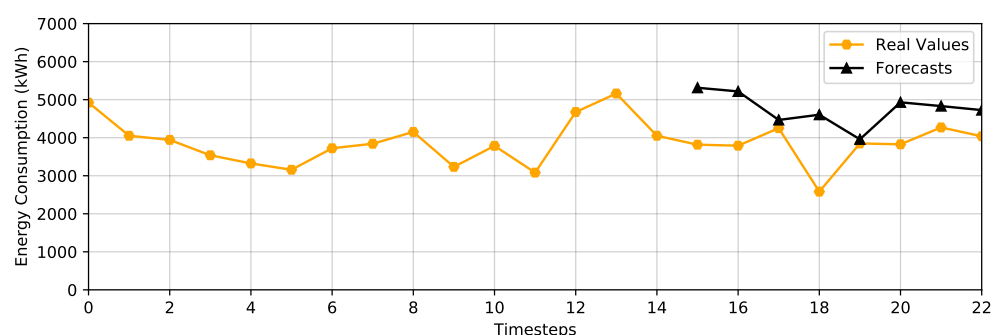


**Table 8.** Best results, per scenario, for Method 2. The letters stand as follows: a. timesteps; b. batch size; c. number of layers; d. number of neurons/filters; e. pool size; f. kernel size; g. dropout; h. activation; i. RMSE; j. MAE; k. time (s).

Model	a.	b.	c.	d.	e.	f.	g.	h.	i.	j.	k.
<i>Uni-Variate-Scenario 1</i>											
CNN	14	30	4	32	2	4	0.5	tanh	869.78	784.23	13
LSTM	21	20	3	32	-	-	0.5	tanh	913.90	828.48	72
GRU	28	20	3	64	-	-	0.5	ReLU	869.85	798.94	83
<i>Multi-Variate-Scenario 1</i>											
CNN	21	10	5	16	3	4	0.0	ReLU	926.23	845.11	27
LSTM	21	5	3	64	-	-	0.0	ReLU	961.34	881.92	186
GRU	21	10	3	32	-	-	0.5	ReLU	950.09	863.89	46
<i>Multi-Variate-Scenario 2</i>											
CNN	14	20	3	16	3	4	0.0	tanh	885.90	796.07	21
LSTM	28	20	3	128	-	-	0.0	ReLU	913.90	845.28	51
GRU	21	20	3	64	-	-	0.0	ReLU	887.41	808.75	43
<i>Multi-Variate-Scenario 3</i>											
CNN	14	20	3	16	3	3	0.0	tanh	916.27	831.90	12
LSTM	14	30	3	128	-	-	0.5	tanh	946.63	854.81	31
GRU	21	20	3	32	-	-	0.5	ReLU	898.67	816.27	39

Regarding the models conceived over the multi-variate approach, it was possible to verify that the best performance was again obtained by a CNN-based model but now in the second scenario. This scenario had, as input features, the *value\_energy*, *temperature*, and *flow\_value*. In this approach, it was possible to verify that in the scenario with the most significant number of features given with input to the models, all three models presented the same value of timesteps (21). In the remaining scenarios, where there was a decrease in the number of features, in general, the CNN-based model requires a lower timestamp value than the rest. It should also be noted that, for the most part, all DL models required an equal value of layers in each of the scenarios. It is also interesting to note that this scenario held the best multi-variate candidates for CNNs, LSTMs, and GRUs.

Figure 9 illustrates several multi-step forecasts made by the best candidate model in this method. Here, the input sequence was made of 14 timesteps (i.e., days).



**Figure 9.** Eight multi-step forecasts for the best candidate model in Method 2.

### 3.3. Transfer Learning

It is usual to find situations where an WWTP has insufficient data. Hence, a goal of this study was to understand the applicability of transfer learning processes in this domain. To achieve such a goal, data were obtained from a second WWTP. However, no influent flow data were available. Hence, we were limited to apply transfer learning

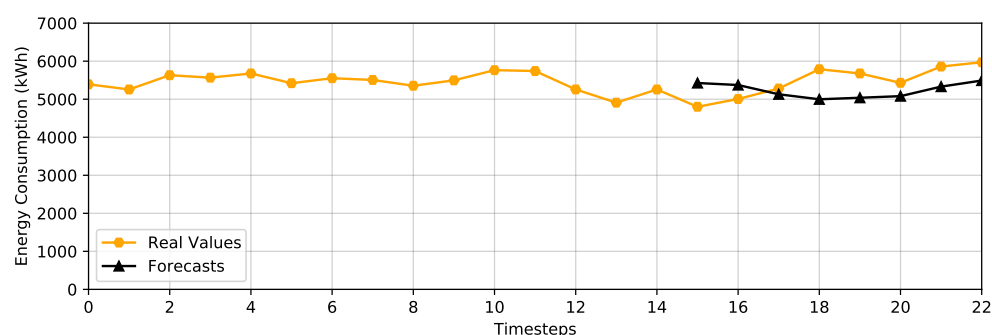
processes over the uni-variate approach since it only considers the *value\_energy* feature, which was only available in a daily periodicity for the years of 2016 and 2017. The best uni-variate candidate model, a CNN, was conceived over the first method, i.e., the one that had the outliers interpolated. Hence, the data from the second WWTP were treated similarly. Finally, 2016 data were used for training and 2017 for testing.

To carry out the transfer learning process, it was necessary to store several parameters of the best uni-variate CNN including its architecture, hyperparameters, and weights (the pre-trained model). Two different settings were tried. The first one re-trained the entire pre-trained CNN model, while the second one only re-trained the layers after the last *Conv1D/AveragePooling1D* pair, inclusive. This is achieved by enabling, or disabling, the *trainable* property of each layer. Table 9 describes the results achieved by the pre-trained uni-variate CNN model, in each setting.

**Table 9.** Results of the pre-trained CNN (Convolutional Neural Networks) model on the second WWTP (Wastewater Treatment Plants).

Setting	RMSE	MAE
1-All model re-trained	357.98	324.69
2-Re-train after the last pair, inclusive	367.72	334.18

It was possible to verify that the method with better performance was the one that re-trained the entire model. This method had a MAE of 324 and a RMSE of 357 kWh. Figure 10 illustrate eight multi-step forecasts for the best model. A total of 14 timesteps were used as input, with successive two-day forecasts encompassing the next 8 days.



**Figure 10.** Eight multi-step forecasts when re-train the entire model.

#### 4. Discussion and Conclusions

Energy consumption forecasting in a WWTP can significantly impact these installations, making them increasingly sustainable, obtaining greater energy efficiency, and reducing costs. After a diversity of experiments being carried out, from all the candidate models, the one achieving a better performance was a multi-variate CNN over the dataset created by Method 1, with a RMSE and MAE of 690 and 630 kWh, respectively.

Another interesting result was the differences in performance concerning the uni- and multi-variate approaches, for the two methods. If in Method 1 the best candidate model was a multi-variate one, in Method 2 it was uni-variate. Regarding both methods, it can be said that the method in which interpolations are made (Method 1) allowed all candidate models to achieve better performances when compared to the method that truncated the outliers (Method 2). Overall, CNN models presented a better performance than the remaining models. Table 10 summarises the obtained results.

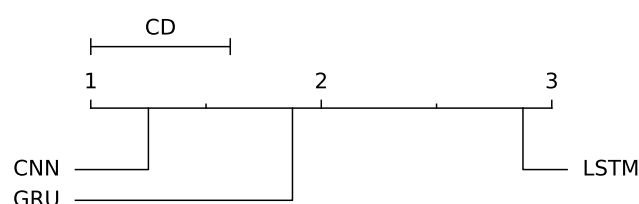
**Table 10.** Ordered list of best candidate models.

Method	Approach	Scenario	Best Candidate	RMSE	MAE
1	Multi-Variate	3	CNN	690.00	630.63
1	Uni-Variate	1	CNN	702.71	645.70
1	Multi-Variate	2	GRU	727.07	670.53
1	Multi-Variate	1	CNN	737.47	677.48
2	Uni-Variate	1	CNN	869.78	784.23
2	Multi-Variate	2	CNN	885.90	796.07
2	Multi-Variate	3	GRU	898.67	816.27
2	Multi-Variate	1	CNN	926.23	845.11

Within the different scenarios in the multi-variate approach, there were some differences between both methods. In Method 1, the best multi-variate scenario was found when combining the influent flow with the energy consumption values (Scenario 3). However, in Method 2, the best multi-variate scenario was found when adding the climatological context to the influent flow and energy consumption values (Scenario 2). In both methods, it was possible to verify that the temporal context (*year* and *month*) worsened the energy consumption forecasts.

Regarding the cardinality of timesteps required as input by the models, in CNN-based models, the increase in the number of features usually led to an increase in the number of timesteps. On the other hand, GRU-based models showed that more features led to a lower number of timesteps. LSTM candidates had their results varying significantly.

Finally, an analysis was carried out to compare the three models' performance. A critical difference diagram was developed to represent the results of a two-tailed Nemenyi post-hoc test, with a  $p < 0.05$ , as depicted in Figure 11. When the average ratings of two models differ by, at least, the critical difference, we can say that the performance between the two is statistically significant. Considering the mean MAE as measure, it is possible to verify that CNNs have better performance than LSTMs and GRUs, being this difference statistically significant.

**Figure 11.** Critical difference diagram showing pairwise comparison of the average ranks in terms of MAE (Mean Absolute Error) ( $p < 0.05$ ).

In regard to the applied transfer learning processes, promising results were achieved using a pre-trained uni-variate CNN model. The best performance was achieved when re-training the whole model. To answer the research questions raised at the beginning of the study, it can be said that (RQ1) CNNs performed better than RNNs, with CNN-based models being the best in practically the whole set of experiments; (RQ2) that the feature that most facilitated the process of forecasting energy consumption in a WWTP was the influent flow; and (RQ3) it was found that it is viable to use transfer learning processes in WWTP with a low volume of data and still present promising results.

However, it is known that other factors can be correlated with energy consumption in a WWTP, such as the concentration of certain pollutants in water like  $BOD_5$ . Nevertheless, to obtain this data, laboratory analysis of WWTP waters is required. Thus, it can take us several days to know the  $BOD_5$  value, among many others. Hence, from a data exploration perspective, it is interesting to understand the impact of such pollutants on energy consumption. Although, from an engineering point of view, this is a significant limitation as the goal of this study is to deploy the best DL model to have real-time forecasts of energy

consumption. If we were expected to include the concentration of such pollutants, it would only be possible to predict the value of energy consumption in the WWTP for tomorrow after obtaining the results from the laboratory, and this would only be available the day after tomorrow. In this way, we would not be able to implement the model to predict the value of energy consumption for tomorrow due to some input parameters of the model would be unknown and would only be available in a few days.

Considering that we are handling a real-life scenario and that the goal is to deploy the best candidate model in a WWTP, future work and research will focus on the use of more extensive sets of data, as well as the conception and evaluation of hybrid models to forecast energy consumption. An additional goal is to conceive a dashboarding platform for Machine Learning Operations (MLOps) to improve the process of monitoring the execution and performance of the deployed models.

**Author Contributions:** Conceptualization, P.O. and B.F.; methodology, P.O. and B.F.; software, P.O. and B.F.; validation, P.O. and B.F.; formal analysis, P.O. and B.F.; investigation, P.O. and B.F.; resources, P.N. and C.A.; data curation, P.O. and B.F.; writing—original draft preparation, P.O. and B.F.; writing—review and editing, P.N. and C.A.; visualization, P.O. and B.F.; supervision, P.N.; project administration, P.N. and C.A.; funding acquisition, P.N. and C.A. All authors have read and agreed to the published version of the manuscript.

**Funding:** The work of Paulo Novais and Cesar Analide has been supported by FCT—Fundação para a Ciência e Tecnologia within the R&D Units Project Scope: UIDB/00319/2020. The work of Pedro Oliveria and Bruno Fernandes is also supported by National Funds through the Portuguese funding agency, FCT—Fundação para a Ciência e a Tecnologia within project DSAIPA/AI/0099/2019.

**Data Availability Statement:** The data presented in this study are available on request from the corresponding author. The data are not publicly available due to having been made available by a multi-municipal water systems company.

**Conflicts of Interest:** The authors declare no conflict of interest.

## Abbreviations

The following abbreviations are used in this manuscript:

ANN	Artificial Neural Networks
ARIMA	Autoregressive Integrated Moving Average
BSM1	Benchmark Simulation Model
$BOD_5$	Biological Oxygen Demand
COD	Chemical Oxygen Demand
CNN	Convolutional Neural Network
DL	Deep Learning
ENN-ECM	Elman Neural Network-Energy Consumption Model
FCM	Fuzzy C-Means
GRU	Gated Recurrent Units
LSTM	Long Short-Term Memory
MAE	Mean Absolute Error
MAPE	Mean Absolute Percentage
MLOps	Machine Learning Operations
MLR	Multi-variable Linear Regression
OLS	Ordinary Least Square
RBF	Radial Basis Function
RMSE	Root Mean Square Error
RMSLE	Root Mean Squared Log Error
RQ	Research Question
WWTP	Wastewater Treatment Plant

## References

- World Urbanization Prospects–Population Division–United Nations. 2018. Available online: <https://population.un.org/wup/> (accessed on 21 January 2021).
- Omer, A.M. Energy, environment and sustainable development. *Renew. Sustain. Energy Rev.* **2018**, *12*, 2265–2300. [CrossRef]
- Daw, J.; Hallett, K.; DeWolfe, J.; Venner, I. *Energy Efficiency Strategies for Municipal Wastewater Treatment Facilities*; National Renewable Energy Lab.(NREL): Golden, CO, USA, 2012. [CrossRef]
- Liu, F.; Ouedraogo, A.; Manghee, S.; Danilenko, A. *A Primer on Energy Efficiency for Municipal Water and Wastewater Utilities*; World Bank: Washington, DC, USA, 2012.
- Frade, J.; Lacasta, N.; Mendes, P.; Cardoso, P.; Trindade, I.; Newton, F.; Franco, P.; Serra, A.; Póvoa, C.; Narciso, F. PensaAR 2020–Uma Estratégia ao Serviço da População: Serviços de Qualidade a um Preço Sustentável. Available online: <https://www.apambiente.pt/index.php?ref=16&subref=7&sub2ref=9&sub3ref=1098> (accessed on 22 January 2021).
- Rajaeifar, M.; Ghanavati, H.; Dashti, B.; Heijungs, R.; Aghbashlo, M.; Tabatabaei, M. Electricity generation and GHG emission reduction potentials through different municipal solid waste management technologies: A comparative review. *Renew. Sustain. Energy Rev.* **2017**, *79*, 414–439. [CrossRef]
- Zeng, S.; Chen, X.; Dong, X.; Liu, Y. Efficiency assessment of urban wastewater treatment plants in China: Considering greenhouse gas emissions. *Resour. Conserv. Recycl.* **2017**, *120*, 157–165. [CrossRef]
- De Haas, D.; Foley, J.; Marshall, B.; Dancey, M.; Vierboom, S.; Bartle-Smith, J. Benchmarking Wastewater Treatment Plant Energy Use in Australia. Available online: [https://www.researchgate.net/profile/David-De-Haas-2/publication/276921977\\_Benchmarking\\_Wastewater\\_Treatment\\_Plant\\_Energy\\_Use\\_in\\_Australia/links/5599093e08ae793d137e2735/Benchmarking-Wastewater-Treatment-Plant-Energy-Use-in-Australia.pdf](https://www.researchgate.net/profile/David-De-Haas-2/publication/276921977_Benchmarking_Wastewater_Treatment_Plant_Energy_Use_in_Australia/links/5599093e08ae793d137e2735/Benchmarking-Wastewater-Treatment-Plant-Energy-Use-in-Australia.pdf) (accessed on 25 January 2021).
- Li, Z.; Zou, Z.; Wang, L. Analysis and forecasting of the energy consumption in wastewater treatment plant. *Math. Probl. Eng.* **2019**, *2019*. [CrossRef]
- Harrou, F.; Cheng, T.; Sun, Y.; Leiknes, T.O.; Ghaffour, N. A Data-Driven Soft Sensor to Forecast Energy Consumption in Wastewater Treatment Plants: A Case Study. *IEEE Sens. J.* **2020**, *21*, 4908–4917. [CrossRef]
- Huang, X.; Han, H.; Qiao, J. Energy consumption model for wastewater treatment process control. *Water Sci. Technol.* **2013**, *67*, 667–674. [CrossRef] [PubMed]
- Ramli, N.A.; Abdul, M.F. Analysis of energy efficiency and energy consumption costs: A case study for regional wastewater treatment plant in Malaysia. *J. Water Reuse Desalin.* **2017**, *7*, 103–110. [CrossRef]
- Maki, S.; Chandran, R.; Fujii, M.; Fujita, T.; Shiraishi, Y.; Ashina, S.; Yabe, N. Innovative information and communication technology (ICT) system for energy management of public utilities in a post-disaster region: Case study of a wastewater treatment plant in Fukushima. *J. Clean. Prod.* **2019**, *233*, 1425–1436. [CrossRef]
- Oulebsir, R.; Lefkir, A.; Safri, A.; Bermad, A. Optimization of the energy consumption in activated sludge process using deep learning selective modeling. *Biomass Bioenergy* **2020**, *132*, 105420. [CrossRef]
- Fernandes, B.; Silva, F.; Alaiz-Moreton, H.; Novais, P.; Neves, J.; Analide, C. Long Short-Term Memory Networks for Traffic Flow Forecasting: Exploring Input Variables, Time Frames and Multi-Step Approaches. *Informatica* **2020**, *31*, 723–749. [CrossRef]
- Jin, X.; Yang, N.; Wang, X.; Bai, Y.; Su, T.; Kong, J. Integrated predictor based on decomposition mechanism for PM2.5 long-term prediction. *Appl. Sci.* **2019**, *9*, 4533. [CrossRef]
- Zhang, T.; Song, S.; Li, S.; Ma, L.; Pan, S.; Han, L. Research on gas concentration prediction models based on LSTM multidimensional time series. *Energies* **2019**, *12*, 161. [CrossRef]
- Mbatha, N.; Bencherif, H. Time series analysis and forecasting using a novel hybrid LSTM data-driven model based on empirical wavelet transform applied to total column of ozone at Buenos aires, Argentina (1966–2017). *Atmosphere* **2020**, *11*, 457. [CrossRef]
- Chatterjee, A.; Gerdes, M.W.; Martinez, S.G. Statistical explorations and univariate timeseries analysis on covid-19 datasets to understand the trend of disease spreading and death. *Sensors* **2020**, *20*, 3089. [CrossRef]
- Zhang, W.; Yu, Y.; Qi, Y.; Shu, F.; Wang, Y. Short-term traffic flow prediction based on spatio-temporal analysis and CNN deep learning. *Transp. Transp. Sci.* **2019**, *15*, 1688–1711. [CrossRef]
- Dong, X.; Qian, L.; Huang, L. Short-term load forecasting in smart grid: A combined CNN and K-means clustering approach. In Proceedings of the International Conference on Big Data and Smart Computing (BigComp), Jeju, Korea, 13–16 February 2017; pp. 119–125. [CrossRef]
- Hussain, D.; Hussain, T.; Khan, A.A.; Naqvi, S.A.A.; Jamil, A. A deep learning approach for hydrological time-series prediction: A case study of Gilgit river basin. *Earth Sci. Informatics* **2020**, *13*, 915–927. [CrossRef]
- Oliveira, P.; Fernandes, B.; Aguiar, F.; Pereira, M.A.; Analide, C.; Novais, P. A Deep Learning Approach to Forecast the Influent Flow in Wastewater Treatment Plants. In Proceedings of the International Conference on Intelligent Data Engineering and Automated Learning, Guimarães, Portugal, 4–6 November 2020; pp. 362–373. [CrossRef]
- Schmidhuber, J. Deep learning in neural networks: An overview. *Neural Netw.* **2015**, *61*, 85–117. [CrossRef] [PubMed]
- Medsker, L.; Jain, L. *Recurrent Neural Networks: Design and Applications*; CRC Press: Boca Raton, FL, USA, 1999.
- Hochreiter, S.; Schmidhuber, J. Long short-term memory. *Neural Comput.* **1997**, *9*, 1735–1780. [CrossRef]
- Kang, D.; Lv, Y.; Chen, Y. Short-term traffic flow prediction with LSTM recurrent neural network. In Proceedings of the 20th International Conference on Intelligent Transportation Systems (ITSC), Yokohama, Japan, 16–19 October 2017; pp. 1–6. [CrossRef]

28. Yang, B.; Sun, S.; Li, J.; Lin, X.; Tian, Y. Traffic flow prediction using LSTM with feature enhancement. *Neurocomputing* **2019**, 320–327. [CrossRef]
29. Fente, D.N.; Singh, D.K. Weather forecasting using artificial neural network. In Proceedings of the 2018 Second International Conference on Inventive Communication and Computational Technologies (ICICCT), Coimbatore, India, 20–21 April 2018; pp. 1757–1761. [CrossRef]
30. Kim, T.; Cho, S. Web traffic anomaly detection using C-LSTM neural networks. *Expert Syst. Appl.* **2018**, 106, 66–76. [CrossRef]
31. Feng, C.; Li, T.; Chana, D. Multi-level anomaly detection in industrial control systems via package signatures and LSTM networks. In Proceedings of the 47th Annual IEEE/IFIP International Conference on Dependable Systems and Networks (DSN), Denver, CO, USA, 26–29 June 2017; pp. 261–272. [CrossRef]
32. Cho, K.; Van Merriënboer, B.; Bahdanau, D.; Bengio, Y. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv* **2014**, arXiv:1409.1259.
33. Chung, J.; Gulcehre, C.; Cho, K.; Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv* **2014**, arXiv:1412.3555.
34. Niu, Z.; Yu, Z.; Tang, W.; Wu, Q.; Reformat, M. Wind power forecasting using attention-based gated recurrent unit network. *Energy* **2020**, 196, 117081. [CrossRef]
35. Wang, R.; Li, C.; Fu, W.; Tang, G. Deep learning method based on gated recurrent unit and variational mode decomposition for short-term wind power interval prediction. *IEEE Trans. Neural Netw. Learn. Syst.* **2019**, 31, 3814–3827. [CrossRef] [PubMed]
36. Wang, Y.; Liao, W.; Chang, Y. Gated recurrent unit network-based short-term photovoltaic forecasting. *Energies* **2018**, 11, 2163. [CrossRef]
37. Fukushima, K.; Miyake, S. *Neocognitron: A Self-Organizing Neural Network Model for a Mechanism of Visual Pattern Recognition*; Springer: Berlin/Heidelberg, Germany, 1982; pp. 267–285. [CrossRef]
38. LeCun, Y.; Bottou, L.; Bengio, Y.; Haffner, P. Gradient-based learning applied to document recognition. *Proc. IEEE* **1998**, 86, 2278–2324. [CrossRef]
39. Hubel, D.H.; Wiesel, T.N. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *J. Physiol.* **1962**, 160, 106–154. [CrossRef] [PubMed]
40. Li, Q.; Cai, W.; Wang, X.; Zhou, Y.; Feng, D.D.; Chen, M. Medical image classification with convolutional neural network. In Proceedings of the 13th International Conference on Control Automation Robotics & Vision (ICARCV), Singapore, 10–12 December 2014; pp. 844–848. [CrossRef]
41. Chen, C.; Liu, M.; Tuzel, O.; Xiao, J. R-CNN for small object detection. In Proceedings of the 13th Asian Conference on Computer Vision, Taipei, Taiwan, 20–24 November 2016; pp. 214–230. [CrossRef]

## Article

# Extending a Trust model for Energy Trading with Cyber-Attack Detection

Rui Andrade , Sinan Wannous , Tiago Pinto  and Isabel Praça 

GECAD—Knowledge Engineering and Decision Support Research Centre, School of Engineering, Polytechnic of Porto (ISEP/IPP), 4050-535 Porto, Portugal; [sinai@isep.ipp.pt](mailto:sinai@isep.ipp.pt) (S.W.); [tcp@isep.ipp.pt](mailto:tcp@isep.ipp.pt) (T.P.)

\* Correspondence: [rfaar@isep.ipp.pt](mailto:rfaar@isep.ipp.pt) (R.A.); [icp@isep.ipp.pt](mailto:icp@isep.ipp.pt) (I.P.)

**Abstract:** This paper explores the concept of the local energy markets and, in particular, the need for trust and security in the negotiations necessary for this type of market. A multi-agent system is implemented to simulate the local energy market, and a trust model is proposed to evaluate the proposals sent by the participants, based on forecasting mechanisms that try to predict their expected behavior. A cyber-attack detection model is also implemented using several supervised classification techniques. Two case studies were carried out, one to evaluate the performance of the various classification methods using the IoT-23 cyber-attack dataset; and another one to evaluate the performance of the developed trust mode.

**Keywords:** cyber-attack detection; IoT; trust; energy trading; trusted negotiations

**Citation:** Andrade, R.; Wannous, S.; Pinto, T.; Praça, I. Extending a Trust model for Energy Trading with Cyber-Attack Detection. *Electronics* **2021**, *10*, 1975. <https://doi.org/10.3390/electronics10161975>

Academic Editor: Myung-Sup Kim

Received: 15 July 2021

Accepted: 6 August 2021

Published: 17 August 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

The energy market and electric grid play a major role in everyday life. Most areas in modern society require electric energy to operate properly. The electric grid has become indispensable for life in modern society. Due to these reasons, it is important to maintain and improve the stability and reliability of the energy grid.

Currently, energy grids tend to follow a very strict and somewhat inefficient structure. A high number of entities that desire to consume energy are connected to a single centralized energy supplier entity. Traditional energy markets, such as wholesale or retail markets, were not designed to support the rising in distributed energy generation coming from Renewable Energy Sources (RES) in households, small commerce and small industry. Such facts raise questions about different ways of structuring energy markets to deal with these challenges.

One of the possible proposals to answer to this problem is the creation and implementation of local energy markets (LEMs). LEMs are structured in such a way as to enable small-scale negotiations and energy exchanges between participants who traditionally would only be final consumers. These markets are designed to operate within a regional area, such as a neighborhood or a city. Participants in this market are the local households, small commerce and small industry, that may be regular consumers or consumers with some type of local energy generation, being referred to as prosumers. Furthermore, local small-scale power plants can also participate in the LEM. The LEM is better designed to deal with distributed energy generation from RES because the surplus in generation from local energy producers and prosumers can be purchased and utilized by local consumers. This flexibility of response makes LEM an attractive proposition for the future of energy markets.

In order to guarantee the success and desired operation of the LEM, it is necessary to ensure security and trust in negotiations. While security is focused on the traditional measures of cyber-security, such as security in network communications, trust is focused on ensuring that the LEM participants and their proposals in the negotiations are viable and trustworthy.

The objective in this work is to create a LEM simulation, and incorporate a trust model. The trust model should be able to score participants' trust level during negotiations, allowing the untrustworthy participants (low trust score) to be prohibited from participating in the LEM. Furthermore, the goal is to also create a cyber-attack detection model utilizing supervised classification techniques.

A MAS is a system that combines several agents, which are software entities that have the capacity to interact among themselves. For this reason it is ideal to simulate the LEM, as each market participant can be simulated individually, and by the means of their interactions, it is possible to simulate a far more complex environment, as is the case of the LEM.

After this introductory Section, the document is organized as follows: Section 2 contextualizes the work, describing the concepts of the LEM, and of trust models for MAS. Section 3 presents possible approaches to obtain the cyber-security guaranties needed for the safe operation of the LEM. Section 4 describes the LEMMAS system. Section 5 describes the developed cyber-attack detection system. Section 6 presents the analyses done to Cyber-Attack detection models. Lastly, Section 7 presents the conclusions of this work.

## 2. Local Energy Market

The local energy market (LEM) is a novel energy market model. There is no unique definition on what a LEM is; however, many authors have addressed this issue, and among their work a general idea of the LEM begins to emerge. Authors exploring this topic tend to define three key aspects: (i) market structure; (ii) advantages; and (iii) challenges.

The structure of the LEM is generally defined as a group of local participants (such as a neighborhood) [1–3], which are capable of trading energy among themselves. Participants in the local market are separated into three kinds [3–5]:

- Consumers: who wish to buy energy;
- Producers: who wish to sell energy;
- Prosumers: (consumers with some source of energy generation) who wish to buy and sell energy.

Both the market participants and the underling electrical grid, which serves as a basis for the LEM, are defined as having monitoring sensors for consumption, generation, energy storage and other data sources; and network communication technologies to share this information [3,4]. Such an energy grid is referred to as a Smart-Grid [1].

The LEM brings several potential advantages when compared to traditional energy markets. Some authors [1,5] claim that the LEM would make a more efficient use of electrical grids. Simultaneously, it is believed that the shift to local energy markets (LEMs) could reduce the greenhouse effect [1] and create a more sustainable environment [4,5]. Participants in the LEM (especially traditional consumers) take a much more involved role in the market when compared to traditional markets. These participants gain the ability to directly negotiate and can achieve cost reductions or even profits with their participation [1,4,5]. Lastly, the versatility of the LEM makes it possible for the coexistence with traditional markets [4,5], that being the case, the local market can adapt to the needs of each specific community.

Currently, the LEM is facing some challenges that prevent its adoption at a large scale. Abidin et al. [1] identify security concerns as one of these challenges. The local market, and consequently the underling Smart-Grid, deal with a lot of sensitive information that needs to be properly secured from unauthorized access; and from malicious entities who may tamper with data in order to have some financial gain. The former also emphasizes the need for trust in negotiations in the local energy market (LEM). Interest from the community and an economic upfront investment by investors are also seen as one of the current challenges to the LEM adoption [2,5]. From a technical point of view, the implementation of Smart-Grids capable of providing the support needed for the LEM is still a challenge that needs further research [5]. Lastly, the support from governments and creation of adequate legislation is a must for the success of the LEM [5].



### 2.1. Trust in Multi-Agent Systems

Trust and reputation systems (TRS) are designed with the objective of predicting the reliability in the behavior of an entity by analyzing data from past interactions [6]. By performing such analyzes, TRS are able to associate a reputation to each user. Good reputation indicates that the user is trustworthy in its negotiations, and vice versa.

In [6], several trust models are identified, some of which are specific for the marketplace area of applicability. Two of these trust models seem interesting for this project since they are targeted at a marketplace, but apply different strategies. These models are the e-commerce model and ReGreT [7].

Reference [8] views the e-commerce trust model from the perspective of eBay. eBay operates as an online auction web site. Users of this platform can propose their sale offers and/or place bids on other users' offers. In online auction web sites such as eBay, the participants in the transactions are humans, and these platforms implement mechanisms for participants to review their experience in the transaction. This feedback provided by the users is then used to feed the TRS with the data necessary to access the reputation of the users [6].

ReGreT is a trust and reputation model proposed by [7]. This model is different from the eBay model because it does not consider trust as a global value. ReGreT has a focus towards modularity [9]. Modules might be used or not depending on the needs of each context. ReGreT considers three kinds of information for trust: the agent's own experiences, information from other agents and the social structure among agents. These types of information coincide with the three dimensions used in ReGreT to calculate trust, which are the following:

- Individual dimension: Considers the outcomes observed directly by the agent when in negotiation with another agent.
- Social dimension: Considers information provided by other entities. Something that can be useful when direct information is not available.
- Ontological dimension: Considers the contextual information that can be gained by the reputation.

In another work [10], the authors identify three distinct approaches that can be followed when developing a trust mechanism for a MAS. Each approach considers a different dimension of trust. These approaches are:

- Security Approach: Is focused on the traditional security mechanisms [11]: confidentiality, availability, authentication, integrity, non-repudiation. This dimension aims to prevent cyber-security threats;
- Institutional Approach: Considers the idea of a centralized entity that acts as an overseer in the MAS. This centralized entity, takes the place of an institution that must evaluate all agents and ensure that each one of them is trustworthy. This is the case for the e-commerce trust model;
- Social Approach: Is similar to the way humans interact in the real world. With this trust model, each agent decides who it considers trustworthy. Agents can make this decision based on their interactions with other agents, and/or by considering others' opinions. This is the case for the ReGreT trust and reputation model.

### 2.2. Security Risks

In a LEM, part of the physical layer corresponds to the network and sensor infrastructure that is necessary for collecting data and allowing communication, and this is what makes the grid be called a smart-grid. However, this infrastructure can be a vector of cyber-attacks [12]. The sensor infrastructure in the smart-grid is often composed of IoT devices. IoT devices can have a potential risk of being tied to a company or cloud network and having access to the data collected by the sensor. A security breach in the cloud network would also expose the data related to the sensor and intern the LEM where this sensor is [13].

Traditional cyber-attacks to the LEM's smart-grid are also a security risk [13]. These can be attacks that aim at gathering private information, such as Man In The Middle attacks, or can even be attacks that try to tamper with the communications in the network. There might even be financial incentives to try such an attack by an ill intending participant in the LEM since he might be able to change the final market price in order to have a financial gain. Unknown sensor hardware malfunctions can also be problematic since they can leave the system working with incorrect data. In this work, the aim is to provide tools to help detect both malicious data tampering and hardware malfunctions.

In this paper, we present a first approach towards including mechanisms to detect attacks that can be made on devices of market participants, with the aim of making LEMMAS a system that provides trusted and secured negotiation.

### 3. Cyber-Attack Detection

In order to obtain a secure environment for local energy market negotiations, traditional cyber-security cannot be forgotten. It is as important to trust in the participants as it is to have a secure network and computer systems. One option to create this environment is to combine intrusion detection systems with artificial intelligence algorithms as an anomaly/attack detection tool by analyzing network data.

In [14], the authors developed a system to perform intrusion detection of smart meters. They combined support vector machine (SVM) and temporal failure propagation graph (TFPG) techniques with a pattern recognition algorithm. The study showed that the system provided good results.

The authors in [15] tried to detect false data injection in a smart grid using deep learning techniques. Their approach combined a Convolutional Neural Network (CNN) with a Long Short Term Memory (LSTM) network. The system was able to achieve an accuracy result above 90% for certain kinds of attacks. They conclude that their approach can be combined with a different technique to obtain a highly accurate attack detection system for all kinds of attacks.

In [16], the authors also tackled the problem of false data injection in power systems. Their approach used an autoencoder network with 4 hidden layers. A case study was performed, and their system was able to detect the kinds of attacks the study was focusing on, and the system also outperformed the techniques currently used in that scenario.

In [17], the authors researched the problem of face spoofing attacks. Instead of using the traditional methods for such a problem, the authors opted to use ensemble based technique by combining multiple one-class classifiers. A case study was conducted to evaluate the performance of their approach using three face anti-spoofing datasets. Their proposed solution showed a good performance for the problem.

In [18], the authors developed a cyber-attack detection system for network based attacks. In this work, the methods of random forest, multi-layer perceptron and long-short term memory were implemented and experimented using the CIDDs-001 dataset [19,20]. The results of this study showed that the long-short term memory technique was the best, achieving an accuracy score above 99%.

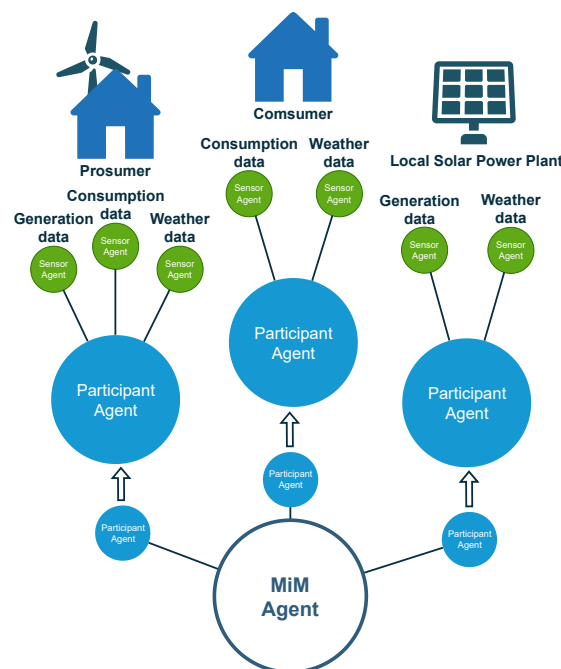
The authors in [21] studied the application of unsupervised learning techniques in order to perform cyber-attack anomaly detection. The authors experimented with six different techniques: Isolation Forest, K-Means, 1-Nearest Neighbor, Autoencoder, Scaled Convex Hull, Support Vector Machines; combined with the best pre-processing steps for each. A case study was performed with the NSL-KDD [22] and the ISCX [23] datasets in order to evaluate the algorithms. Based on the results, the authors concluded that all detection algorithms showed a good performance for the cyber-attack anomaly detection problem.

### 4. LEMMAS System

The developed MAS follows the agent structure as proposed in [24]. In that work, a computational model of a LEM is separated into three kinds of agents:

- **Market Interactions Manager:** Is at the center of every LEM and is responsible for managing all negotiations within the market. The task of ensuring trust in the negotiations is also a responsibility of the market interaction manager (MIM).
- **Participant agent:** Acts on behalf of consumers, producers or prosumers. This agent assumes the role of a negotiator that seeks to best satisfy the needs of the respective market participants (home owners, local commerce and small industry owners). The Participant Agent will have Sensor Agents that report to them the information needed for the negotiations. The Participant Agent is then able to make proposals to buy or sell energy in the market;
- **Sensor agent:** Has the single responsibility of acquiring one type of data and reporting said data to their respective Participant Agent. A Sensor agent can be, for example, connected to a meter measuring energy consumption in a household, while another sensor agent can be connected to a web service in order to obtain the weather forecast.

With these three kinds of agents, it is possible to create a reasonably complete representation of a LEM, which includes: consumers, producers and prosumers. The Sensor Agents allow the cyber-physical system, such as the ones of smart houses and other connected environments. A complete representation of the proposed LEM model is presented in Figure 1.



**Figure 1.** Proposed LEM model diagram [24].

The LEM is composed of several participants, represented by their respective Participant Agent, and all of these agents are connected to the MIM. In Figure 1, three participants are further detailed as examples of how real participants might be structured in a realistic scenario. These participants are the following:

- **Consumer:** Represents a household without self-generation that participates in the market.
- **Prosumer:** Represents an household that participates in the market and has its own energy generation, with a small wind generator.
- **Local Solar Power Plant:** Exemplifies a small photovoltaic power plant that is part of the LEM.

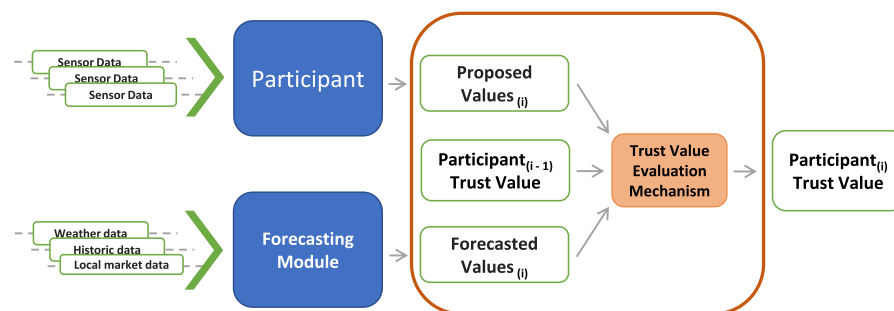
#### 4.1. Trust Model

To support the market, an institutional based trust model is proposed to be used by the MIM, capable of evaluating the behavior of participants and detecting faulty or malicious activities. This trust model was chosen over the social model because with a social model, participants might need access to sensitive (consumption, generation, etc.) data from other participants in order to make their own trust evaluation.

The idea for the trust mechanism is that with information such as weather, historical consumption and generation data, and other contextual data, it is possible to use forecasting methods to try to predict what the participant's consumption, generation or proposals should be in the coming market negotiation period.

Using such forecasted values, it is possible to obtain an idea if the participant is trustworthy over time. Since forecasting methods always have a certain degree of uncertainty, a single proposed value that does not match the forecasted value does not provide a reliable metric. So, by using an evaluation over time, it is thought that incorrectly forecasted values become negligible.

Figure 2 presents a diagram of the proposed trust evaluation process. As shown, the trust evaluation process takes three values as input: the participant's proposed values for the current market negotiation period, the participant's trust value from the previous negotiation period, and the forecasted value based on the participants historical and contextual data.



**Figure 2.** Trust evaluation process.

The definition of the proposed trust model is the following:

The trust value ranges from  $[0, 1]$  where 1 is the highest trust and 0 lowest trust value. The trust value for a participant  $p$  in negotiation period  $i$  is represented as  $t_{pi}$ .

The evaluation formula takes several variables into consideration that can be configured to obtain the best possible results, these variables are:

- $vr$ : represents the *variable acceptance range*, which is a percentage value;
- $fr$ : represents the *fixed acceptance range*, which is a static value;
- $tiv$ : represents the *trust increase value*, a value by which the participant's trust is increased;
- $tdv$ : represents the *trust decrease value*, a value by which the participant's trust is decreased;
- $sv_{pin}$ : represents the *submitted value* by participant  $p$ , from the data source of sensor  $n$  in market negotiation period  $i$ ;
- $fv_{pin}$ : represents the *forecasted value* for participant  $p$ , for sensor  $n$  in negotiation period  $i$ ;
- $t_{p0}$ : represents the *default trust value*, the trust value which all participants start with.

Equation (1) shows how the trust evaluation is calculated by being combined with either the Equation (2) for the asymmetric acceptance range or the Equation (3) for the symmetric acceptance range.

The difference between the asymmetric and the symmetric acceptance range is that the asymmetric has a higher acceptance range when the forecasting mechanism overestimates the value, since a percentage from a higher value results in a higher range.

$$t_{pi} = t_{p(i-1)} + trust\_eval(sv_{pi}, fv_{pi}) \quad (1)$$

$$trust\_eval_{asym}(sv_{pi}, fv_{pi}) = \begin{cases} tiv & \text{if } sv_{pi} > fv_{pi} * (1 - vr) \text{ AND } sv_{pi} < fv_{pi} * (1 + vr) \\ tiv & \text{if } sv_{pi} > fv_{pi} - fr \text{ AND } sv_{pi} < fv_{pi} + fr \\ tdv & \text{otherwise} \end{cases} \quad (2)$$

$$trust\_eval_{sym}(sv_{pi}, fv_{pi}) = \begin{cases} tiv & \text{if } fv_{pi} > sv_{pi} * (1 - vr) \text{ AND } fv_{pi} < sv_{pi} * (1 + vr) \\ tiv & \text{if } fv_{pi} > sv_{pi} - fr \text{ AND } fv_{pi} < sv_{pi} + fr \\ tdv & \text{otherwise} \end{cases} \quad (3)$$

There needs to be some consideration of how each participant's trust value is interpreted. Two things need to be taken into account: a participant that always submits real and true values should be fully trustworthy and so should be evaluated with a 1.0 trust value; on the other hand, a participant that always submits false values should not be trusted and should have a trust evaluation of 0.0.

There is, however, some subjectivity in considering these trust evaluations. For example, a participant that always submits real and true values and is evaluated with a 0.9 trust value, or a participant that always submits false values and is evaluated with a 0.1 trust value, also seem like acceptable evaluations. Given this subjective nature of the trust evaluation, three trust ranges are proposed:

- **Trustworthy:** range where the trust value is  $[h_t, 1]$ , and any participant in this range is fully trusted;
- **Unsure:** range where the trust value is  $[m_t, h_t]$ , and any participant in this range is considered to be a possible malicious or faulty participant and should, for example, be further evaluated by the market authority;
- **Untrustworthy:** range where the trust value is  $[0, m_t]$ , and any participant in this range is considered a malicious or faulty participant and should be prevented from participating in the market.

The values of  $h_t$  minimum threshold for high trust and  $m_t$  minimum threshold for medium trust are variable values that can be configured accordingly to the needs of the LEM.

#### 4.2. Cyber-Security Model

Having a trust model capable of correctly analyzing the trust evaluation of participants supports the LEM negotiations; however, this kind of analysis leaves an important aspect neglected, the origin of a malicious proposal.

To fully understand the safety of negotiations, it is also necessary to consider the traditional security aspects. The idea behind the security model is to analyze the data coming from the sensor agents to discover potential security intrusions.

Discovering a security intrusion would also influence the ability of a participant to negotiate in the market. Regarding the cyber-security aspect we consider a simple binary classification for participants: *Secure*, *Insecure*. This leaves us with the final possible participants classifications in Table 1.

**Table 1.** Possible participant classifications.

	Trustworthy	Unsure	Untrustworthy
Secure	Secure Trustworthy	Secure Unsure	Secure Untrustworthy
Insecure	Insecure Trustworthy	Insecure Unsure	Insecure Untrustworthy

Considering the classifications from Table 1 the participants with classifications of: *Insecure Trustworthy*, *Insecure Unsure* and *Insecure Untrustworthy* should be prevented from negotiating because they can be under a cyber-attack. Participants with a classification of *Secure Untrustworthy* will also be prevented from negotiating as their trust score does not allow it. Participants classified as *Secure Unsure* are in a grey area where they can be allowed to participant, but further investigation is required in order to ensure if the negotiations are at risk. Furthermore, lastly, the participants classified as *Secure Trustworthy* will be allowed to negotiate.

## 5. LEMMAS Case Study

The idea for this case study is to simulate a LEM with several participants that vary in the amount and intensity of false proposals and observing how the proposed trust model evaluates these participants. Since the trust model is based on forecasting, forecasting methods are simulated as a normal standard distribution based on what the real proposal value should be, this way forecasting methods with distinct levels of accuracy and precision can be estimated, and it is possible to see how the performance of the forecasting method influences the trust model performance.

The LEM was simulated for a 24 h period and with 15 min market negotiation period duration, which results in a total of 96 market negotiation periods. Each simulation was performed 10 times, and its results were averaged. The 24 h simulated were of a Monday, simulated from hour 00:00 to hour 24:00. The LEM aggregates 4 participants using real consumption data from private homes publicly available in [25]. Each participant has their own bias in the proposals it submits:

- TP—True Proposer: Is the only who does not have a bias and always sends the real value;
- LUaUP—Low Under and Over Proposer: Sends a real value 80% of the times and a value between 30% under to 30% over the real value the rest of the time;
- MUaOP—Medium Under and Over Proposer: Sends a real value 50% of the times and a value between 60% under to 60% over the real value the rest of the time;
- HUaOP—High Under and Over Proposer: Sends a real value 30% of the times and a value between 90% under to 90% over the real value the rest of the time.

With these participants configurations, the expected result is a correlation between the trust value of the participant and the amount of false submissions. The True Proposer acts as a base line showing if trustworthy participants are being correctly identified.

As for the estimated forecasting methods, four were simulated, in decreasing levels of accuracy and precision. The estimated forecasting methods have the following mean  $\bar{x}$  and standard deviation  $\sigma$ :

- Perfect Predictor:  $\bar{x} = 1.0$   $\sigma = 0.0$ ;
- Low Center Predictor:  $\bar{x} = 1.0$   $\sigma = 0.2$ ;
- High Center Predictor:  $\bar{x} = 1.0$   $\sigma = 0.4$ .

The simulations are preformed with both the symmetric and asymmetric acceptance methods. Lastly, the trust formula variables are configured as such:  $ar = 0.5$ ,  $tiv = 0.01$ ,  $tdv = -0.08$  and  $t_{p0} = 0.8$ ; and the trust ranges are: *Trustworthy*  $[0.8, 1]$ , *Unsure*  $[0.5, 0.8[$  and *Untrustworthy*  $[0, 0.5[$ .

These values were chosen after some experimentation, as they proved to be adequate values for the specific scenario in study.

### Case Study Results and Discussion

To present these results in a clear way, each simulation was divided into 2 graphs showing the trust value for each participant over time, separated by the forecasting method and acceptance method.

Looking at Figure 3, there is a clear distinction in the trust evaluation of each participant. Analyzing Figure 4, the results have some changes from the asymmetric model. All partic-

ipants obtained a higher trust value compared to the results of the asymmetric acceptance.

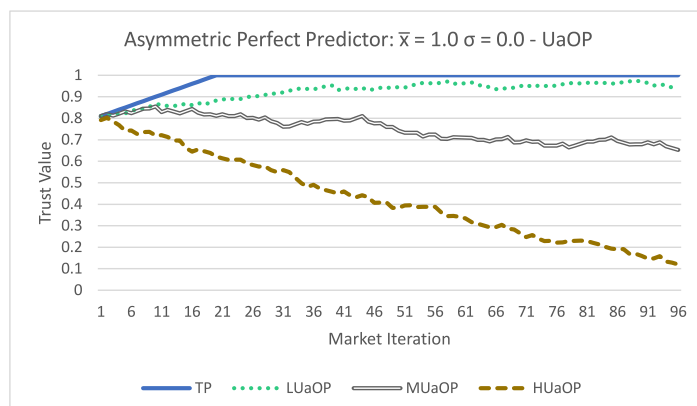


Figure 3. Asymmetric perfect predictor.

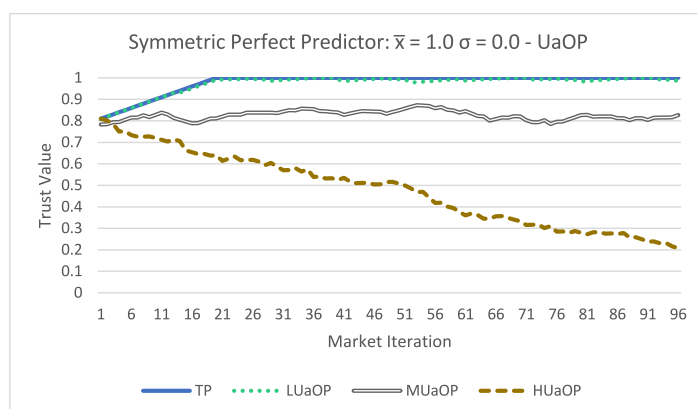


Figure 4. Symmetric Perfect Predictor.

Figures 5 and 6 show different result. In the previous estimator the trust value for the TP participant was always 1.0, but now with the uncertainty in the estimated forecasting method the trust value oscillates; however, it remains close to 1.0. In Figure 6, the results are very similar to the ones obtained with the Perfect Predictor. The biggest difference is in Figure 5 where the TP and LUaOP participants obtained trust values very similar to the ones obtained with the Perfect Predictor, and the MUaOP and HUaOP participants obtained evaluations significantly lower. This demonstrates that the acceptance formula used can make a big difference in the results.

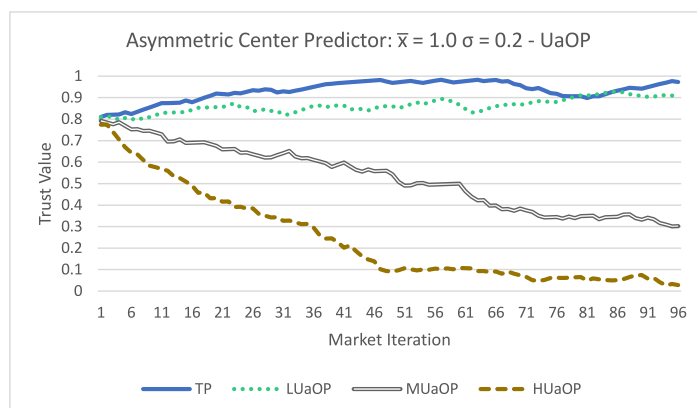
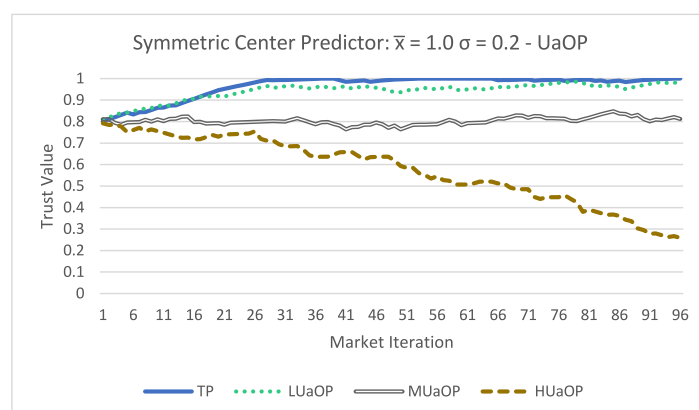
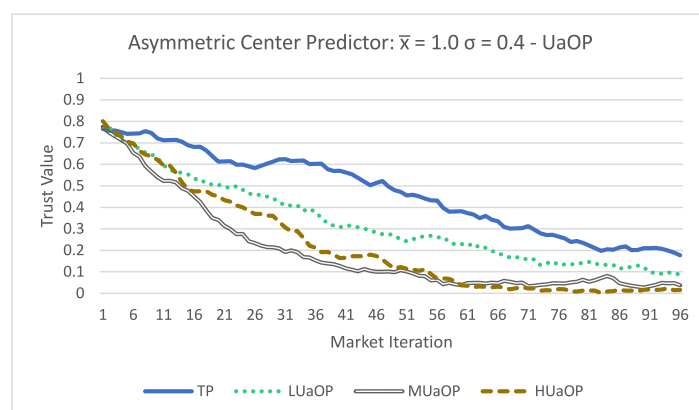


Figure 5. Asymmetric Low Predictor.

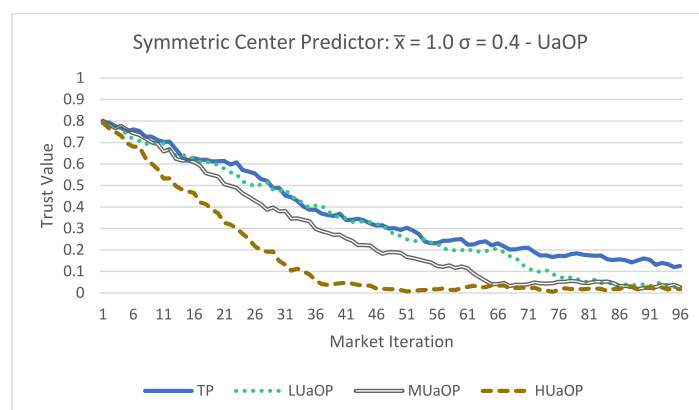


**Figure 6.** Symmetric Low Predictor.

Lastly both Figures 7 and 8 show low trust evaluations for all participants. There are some differences in the way the trust value changed over time between the asymmetric and symmetric acceptance formulas; however, at the end, the values are very similar (all below 0.2). Even the TP participants obtained a low trust evaluation, and this result shows that with a low performing forecasting method the trust evaluation is also low performing.



**Figure 7.** Asymmetric High Predictor.



**Figure 8.** Symmetric High Predictor.

Finally these results lead us to conclude that:

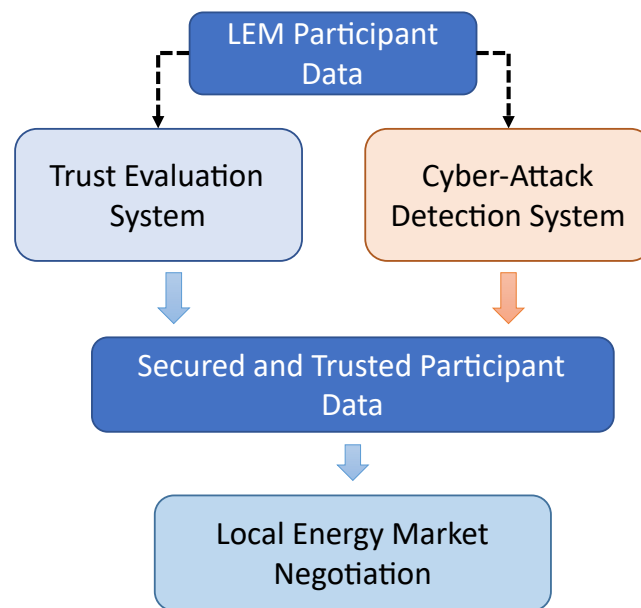
1. Using the proposed trust methodology, it is possible to dynamically update the trust value of a participant;
2. The MIM agent is able to use the proposed trust methodology to access the trust value of a participant;



3. The performance of forecasting methods has a direct impact on the trust evaluation;
4. The acceptance formula can have an impact on the trust evaluation;
5. The higher the amount of false values a participant submits, the lower their trust value will be.

## 6. Developed Security Model and Analysis

As our proposed LEM architecture is based on IoT sensors and their representation as sensor agents, it is fundamental that the information and network communications coming from the sensor agents is secured. With this goal in mind, a cyber-security module is needed to classify participants as described in Table 1. A cyber-attack detection system was developed to complement LEMMAS, as is shown in Figure 9, and the objective is to achieve a negotiation environment with only secured and trusted data.



**Figure 9.** Trust and security modules.

In order to create the necessary security model the python library, Scikit-Learn was used. The goal is to train an artificial intelligent supervised classification model that can analyze the sensor data and classify it as malicious or not.

Six classification models were selected and implemented in order to evaluate which ones were the best for this application. The models are the following.

- Nearest Neighbors;
- Decision Tree;
- Random Forest;
- Neural Network;
- AdaBoost;
- Naive Bayes.

### 6.1. Dataset

Aposemat IoT-23 [26] is a publicly available dataset containing Internet of Things (IoT) network traffic data. This dataset is labeled, including both benign and malicious data entries, and subdivided into 23 sub datasets, 20 containing malicious cyber-attack samples and 3 containing only benign data samples. The data was collected between 2018 and 2019 in three kinds of IoT network devices, namely a Philips HUE smart LED lamp, an Amazon Echo home intelligent personal assistant and a Somfy smart doorlock. These kinds reflect some of the devices that would be part of a smart-home in a smart-grid and as so are aligned with the data generated and collected by LEM participants.

Table 2 presents the different datasets showing which ones include malicious samples, the IoT devices involved, the duration of the attack, the number of packets recorded, the information flows and the size in (GB).

The dataset structure contains the following fields:

- ts—Time of the connection;
- uid—Unique identifier for the connection;
- id.orig\_h, id.orig\_p, id.resp\_h, id.resp\_p—Ports and addresses used in the communication;
- proto—Communication protocol used;
- service—Service protocol used;
- duration—Duration of the communication;
- orig\_bytes, resp\_bytes—Bytes sent in the communication;
- conn\_state—State of the connection;
- local\_orig, local\_resp—Origin of the communication;
- missed\_bytes—Bytes missed in the communication;
- history—History of the state of the connection;
- orig\_pkts—Number of packets sent in the original message;
- orig\_ip\_bytes—Number of IP level bytes sent in the original message;
- resp\_pkts—Number of packets sent in the response message;
- resp\_ip\_bytes—Number of IP level bytes sent in the response message;
- tunnel\_parents—UID of tunnel parents connections;
- label—Label of the sample (malicious or benign);
- detailed-label—Detailed label of the sample (specific attack used in case of a malicious sample).

**Table 2.** IoT-23 dataset description.

#	Type	Capture Name	Malware Device	Duration	Number of Packets	Total Flows	Total Size
1	Malicious	Capture-34-1	Mirai	24	233,000	23,146	121 MB
2	Malicious	Capture-43-1	Mirai	1	82,000,000	67,321,810	6 GB
3	Malicious	Capture-44-1	Mirai	2	1,309,000	238	1.7 GB
4	Malicious	Capture-49-1	Mirai	8	18,000,000	5,410,562	1.3 GB
5	Malicious	Capture-52-1	Mirai	24	64,000,000	19,781,379	4.6 GB
6	Malicious	Capture-20-1	Torii	24	50,000	3210	4 MB
7	Malicious	Capture-21-1	Torii	24	50,000	3287	4 MB
8	Malicious	Capture-42-1	Trojan	8	24,000	4427	3 MB
9	Malicious	Capture-60-1	Gagfyt	24	271,000,000	3,581,029	21 GB
10	Malicious	Capture-17-1	Kenjiro	24	109,000,000	54,659,864	7.8 GB
11	Malicious	Capture-36-1	Okiru	24	13,000,000	13,645,107	992 MB
12	Malicious	Capture-33-1	Kenjiro	24	54,000,000	54,454,592	3.9 GB
13	Malicious	Capture-8-1	Hakai	24	23,000	10,404	2 MB
14	Malicious	Capture-35-1	Mirai	24	46,000,000	10,447,796	3.6 GB
15	Malicious	Capture-48-1	Mirai	24	13,000,000	3,394,347	1.2 GB
16	Malicious	Capture-39-1	IRCBot	7	73,000,000	73,568,982	5.3 GB
17	Malicious	Capture-7-1	Linux Mirai	24	11,000,000	11,454,723	897 MB
18	Malicious	Capture-9-1	Linux Hajime	24	6,437,000	6,378,294	472 MB
19	Malicious	Capture-3-1	Muhstik	36	496,000	156,104	56 MB
20	Malicious	Capture-1-1	Hide & Seek	112	1,686,000	1,008,749	140 MB
21	Benign	Capture-7-1	Soomfy Doorlock	1.4	8276	139	2 MB
22	Benign	Capture-4-1	Phillips HUE	24	21,000	461	4 MB
23	Benign	Capture-5-1	Amazon Echo	5.400	398,000	1383	364 MB

## 6.2. Dataset Pre-Processing

In order to utilize this dataset to train and evaluate models, first a pre-processing step was needed.

The dataset was divided into X and Y, with Y being the target column “label” and X being the remaining data. The columns of “UID” and “ts” were dropped as they do not provide any valuable information. The column “detailed-label” was also dropped since the current objective is only to classify as “Malicious” or “Benign”, meaning a binary classification. All columns containing IPs were converted to the corresponding integer number. The columns of “proto”, “service”, “conn\_state” and “history” were also converted to a numeric value. Regarding missing values, all are imputed and replaced by the median corresponding value. Lastly, the data was randomly split in 80% train data and 20% test data.

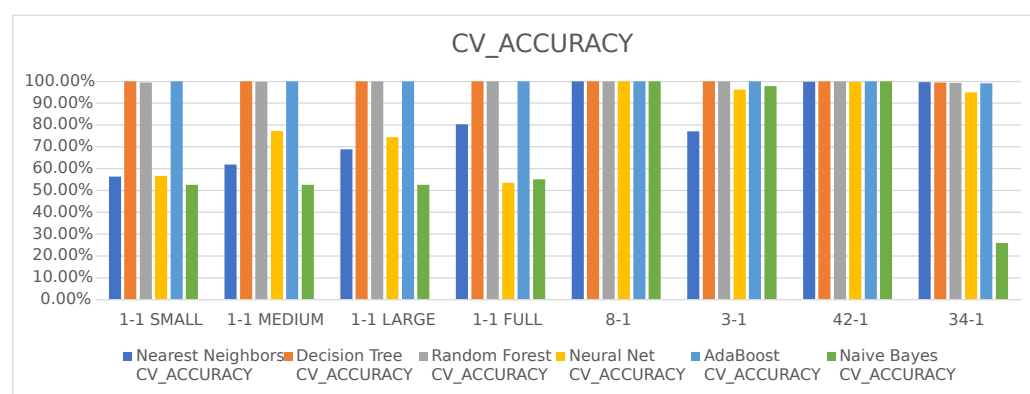
### 6.3. Train, Test and Results

Due to the specifications of our test machine, 64 GB of RAM, 20 cores CPU and a GPU with 8 GB, and the large size of parts of the datasets only some datasets were used, namely: Capture-1-1, Capture-8-1, Capture-3-1, Capture-42-1 and Capture-34-1. Within these datasets, only Capture-1-1 was balanced in the percentage of malicious and benign samples, so we decided to create 3 sub datasets: Capture-1-1 SMALL, Capture-1-1 MEDIUM and Capture-1-1 LARGE, created by randomly selecting samples of each category in a balanced way. The original Capture-1-1 dataset was also used with the name Capture-1-1 FULL. We decided to use this approach to analyze how the performance of the algorithms changes with more data. Table 3 presents in detail the information about the datasets used, including the time it took for each one to be processed. To train and test the model, we used a 80/20 data split, 80% for training and 20% for testing. The analyses were performed with a 5 fold cross validation.

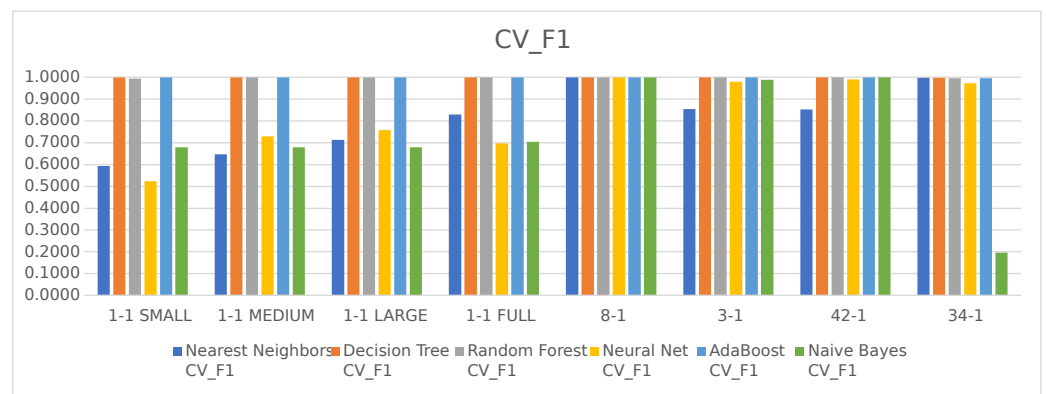
**Table 3.** Datasets used.

File	Total Time	Total Samples	Benign%	Malicious%
1-1 SMALL	0:00:22	20,000	50.00%	50.00%
1-1 MEDIUM	0:23:09	200,000	50.00%	50.00%
1-1 LARGE	1:33:50	400,000	50.00%	50.00%
1-1 FULL	5:31:42	1,008,748	46.52%	53.47%
8-1	0:00:10	10,403	20.96%	79.03%
3-1	0:06:36	156,103	2.90%	97.09%
42-1	0:00:07	4426	99.86%	0.13%
34-1	0:00:29	23,145	8.30%	91.69%

Looking at the results from training and testing presented in Figures 10 and 11, we can see how each technique performed with each dataset.



**Figure 10.** Accuracy of techniques per dataset.



**Figure 11.** F1 score of techniques per dataset.

Regarding the unbalanced datasets, Capture-8-1, Capture-3-1, Capture-42-1 and Capture-34-1, almost all techniques achieved great results with the exception of Naive Bayes when using the dataset Capture-34-1, where the results were low performing. However, these datasets are unbalanced, and the technique might just be over fitting one of the results.

When looking at the results for the balanced datasets, Capture-1-1 SMALL, Capture-1-1 MEDIUM, Capture-1-1 LARGE and Capture-1-1 FULL we see different results, only Decision Tree, Random Forest and AdaBoost were capable of maintaining the strong results of both Accuracy and F1 score. The different amounts of data did not change the results, it only increased the processing time, with the 1-1\_FULL dataset taking more than five and a half hours to process.

Lastly, it is necessary to analyze the percentage of false positives in these algorithms, this is because a model that generates a high percentage of false positives is impractical and will generate more confusion rather than help find and stop cyber-attacks. This metric is presented in Table 4, where each algorithm's false positive rate is shown for each sub-dataset tested. Looking at the results, we can see that once again Decision Tree, Random Forest and AdaBoost are the best options since they obtained a false positive scores below 1% on all sub-datasets, while the other algorithms reached more than 20% on some occasions.

**Table 4.** False positive percentage of each algorithm per sub-dataset.

File	Nearest Neighbors	Decision Tree	Random Forest	Neural Net	AdaBoost	Naive Bayes
1-1 SMALL	25.43%	0.09%	0.59%	25.46%	0.00%	47.34%
1-1 MEDIUM	23.24%	0.00%	0.13%	8.13%	0.00%	47.39%
1-1 LARGE	20.09%	0.01%	0.08%	16.32%	0.00%	47.38%
1-1 FULL	13.99%	0.00%	0.04%	46.52%	0.00%	44.93%
8-1	0.02%	0.02%	0.02%	3.35%	0.02%	0.02%
3-1	1.12%	0.06%	0.00%	2.32%	0.00%	2.13%
42-1	0.25%	0.00%	0.00%	0.27%	0.00%	0.00%
34-1	0.29%	0.46%	0.38%	4.24%	0.28%	0.34%

## 7. Conclusions

The local energy market (LEM) is an emergent market model that is aimed towards solving the challenges currently faced in the energy landscape. One of the requirements for the success of LEM is trust in its negotiations. The main goals in this work are the development of a multi-agent system (MAS) for simulation and modeling LEM; and the proposal of a trust model capable of ensuring trust the LEM negotiations.

A MAS was developed with three types of agents, namely: (i) the Sensor Agent, (ii) the Participant Agent and (iii) the Market Interaction Manager (MIM) Agent, each with their own responsibilities, thus facilitating, the process of modeling the market.

To give a response to the needs of trust in the LEM, a formulation was proposed to calculate a trust value for each participant based on the analysis of the participant's historical data, contextual data, such as weather data, and by using forecasting methods to predict the participants expected behavior. The trust value given to participants evolves over time and takes into consideration its market submissions to the LEM, the forecasting of those submissions and considers the disparity between those values.

A case study was carried out in which several simulations were made with four participants using realistic consumption data and with different biases towards submitting false values. Each simulation used a different estimated forecasting mechanism with distinct levels of accuracy and precision.

The LEM was simulated for a 24 h period and 15 min market negotiation period duration, which resulted in a total of 96 market negotiation periods. This case study's aim was to evaluate the ability of the proposed trust formulation to respond to market needs by evaluating each participant with an appropriate trust value. The realization of the case study made it possible to conclude that: (i) The forecasting methodology used has a big impact on the performance of the trust formulation, but the acceptance formula also needs to be considered; (ii) a bad forecasting method, will provide a bad trust evaluation; and (iii) the higher the amount of false values a participant submits the lower their trust value will be, which is the desired outcome.

A study was carried out to evaluate the six supervised classifications techniques implemented. The training and testing of these classifications techniques were done using the IoT-23, a dataset containing IoT device data under malicious cyber-attacks. The classifications techniques were evaluated using the Accuracy and F1 score metrics. The results showed that the techniques of Decision Tree, Random Forest and AdaBoost provided excellent results. With these results in mind we believe that further studying is necessary with larger datasets and using multi-class classification in order to improve this cyber-attacks detection model. None the less, these results show that such an approach seems viable for the local energy market.

Lastly, one aspect we want to further improve is to develop the LEMMAS system in order to make use of the security and trust models at the same time, and developing a case study to evaluate how these models perform when working together.

**Author Contributions:** Conceptualization, R.A. and I.P.; methodology, R.A., T.P. and I.P.; software, R.A. and S.W.; validation, R.A., S.W., I.P., and T.P.; investigation, R.A.; resources, I.P.; data curation, R.A. and S.W.; writing—original draft preparation, R.A. and S.W.; writing—review and editing, I.P. and T.P.; supervision, I.P.; project administration, I.P.; funding acquisition, I.P. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work has received funding from FEDER Funds through COMPETE program and from National Funds through FCT under the project SPET-PTDC/EEI-EEE/029165/2017 and UIDB/00760/2020.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Abidin, A.; Aly, A.; Cleemput, S.; Mustafa, M.A. Towards a Local Electricity Trading Market Based on Secure Multiparty Computation. 2016. Available online: <https://www.esat.kuleuven.be/cosic/publications/article-2664.pdf> (accessed on 6 August 2021).
2. Bremdal, B.A.; Olivella, P.; Rajasekharan, J. EMPOWER: A network market approach for local energy trade. In Proceedings of the 2017 IEEE Manchester PowerTech, Manchester, UK, 18–22 June 2017; pp. 1–6. [CrossRef]
3. Ampatzis, M.; Nguyen, P.H.; Kling, W. Local electricity market design for the coordination of distributed energy resources at district level. In Proceedings of the 2014 IEEE PES Innovative Smart Grid Technologies Conference Europe (ISGT-Europe), Istanbul, Turkey, 12–15 October 2014; pp. 1–6.
4. Teotia, F.; Bhakar, R. Local energy markets: Concept, design and operation. In Proceedings of the 2016 National Power Systems Conference (NPSC), Bhubaneswar, India, 19–21 December 2016; pp. 1–6.
5. Mendes, G.; Nylund, J.; Annala, S.; Honkapuro, S.; Kilkki, O.; Segerstam, J. Local Energy Markets: Opportunities, Benefits, and Barriers. 2018. Available online: <https://www.cired-repository.org/handle/20.500.12455/1265> (accessed on 6 August 2021).

6. Rahimi, H.; Bekkali, H.E. State of the art of Trust and Reputation Systems in E-Commerce Context. *arXiv* **2017**, arXiv:1710.10061.
7. Sabater, J.; Sierra, C. Regret: A reputation model for gregarious societies. In Proceedings of the Fourth Workshop on Deception Fraud and Trust in Agent Societies, Barcelona, Spain, 4 June 2001; Volume 70, pp. 61–69.
8. Houser, D.; Wooders, J. Reputation in auctions: Theory, and evidence from eBay. *J. Econ. Manag. Strategy* **2006**, *15*, 353–369. [CrossRef]
9. Sabater, J.; Sierra, C. Reputation and Social Network Analysis in Multi-agent Systems. In *Proceedings of the First International Joint Conference on Autonomous Agents and Multiagent Systems: Part 1*; ACM: New York, NY, USA, 2002; pp. 475–482. [CrossRef]
10. Pinyol, I.; Sabater-Mir, J. Computational trust and reputation models for open multi-agent systems: A review. *Artif. Intell. Rev.* **2013**, *40*, 1–25. [CrossRef]
11. Yan, Y.; Qian, Y.; Sharif, H.; Tipper, D. A Survey on Cyber Security for Smart Grid Communications. *IEEE Commun. Surv. Tutor.* **2012**, *14*, 998–1010. [CrossRef]
12. Ghirardello, K.; Maple, C.; Ng, D.; Kearney, P. Cyber security of smart homes: Development of a reference architecture for attack surface analysis. In Proceedings of the Living in the Internet of Things: Cybersecurity of the IoT, London, UK, 28–29 March 2018.
13. Hall, F.; Maglaras, L.; Aivaliotis, T.; Xagoraris, L.; Kantzavelou, I. Smart Homes: Security Challenges and Privacy Concerns. *arXiv* **2020**, arXiv:2010.15394.
14. Sun, C.C.; Cardenas, D.J.S.; Hahn, A.; Liu, C.C. Intrusion Detection for Cybersecurity of Smart Meters. *IEEE Trans. Smart Grid* **2020**, *12*, 612–622. [CrossRef]
15. Niu, X.; Li, J.; Sun, J.; Tomsovic, K. Dynamic detection of false data injection attack in smart grid using deep learning. In Proceedings of the 2019 IEEE Power & Energy Society Innovative Smart Grid Technologies Conference (ISGT), Washington, DC, USA, 18–21 February 2019; pp. 1–6.
16. Wang, C.; Tindemans, S.; Pan, K.; Palensky, P. Detection of false data injection attacks using the autoencoder approach. In Proceedings of the 2020 International Conference on Probabilistic Methods Applied to Power Systems (PMAPS), Liege, Belgium, 18–21 August 2020; pp. 1–6.
17. Fatemifar, S.; Awais, M.; Arashloo, S.R.; Kittler, J. Combining multiple one-class classifiers for anomaly based face spoofing attack detection. In Proceedings of the 2019 International Conference on Biometrics (ICB), Crete, Greece, 4–7 June 2019; pp. 1–7.
18. Oliveira, N.; Praça, I.; Maia, E.; Sousa, O. Intelligent cyber attack detection and classification for network-based intrusion detection systems. *Appl. Sci.* **2021**, *11*, 1674. [CrossRef]
19. Ring, M.; Wunderlich, S.; Grödl, D.; Landes, D.; Hotho, A. Flow-based benchmark data sets for intrusion detection. In Proceedings of the 16th European Conference on Cyber Warfare and Security; Dublin, Ireland, 29–30 June 2017; pp. 361–369.
20. Ring, M.; Wunderlich, S.; Grödl, D.; Landes, D.; Hotho, A. Creation of flow-based data sets for intrusion detection. *J. Inf. Warf.* **2017**, *16*, 41–54.
21. Meira, J.; Andrade, R.; Praça, I.; Carneiro, J.; Bolón-Canedo, V.; Alonso-Betanzos, A.; Marreiros, G. Performance evaluation of unsupervised techniques in cyber-attack anomaly detection. *J. Ambient Intell. Humaniz. Comput.* **2020**, *11*, 4477–4489. [CrossRef]
22. Tavallaee, M.; Bagheri, E.; Lu, W.; Ghorbani, A.A. A detailed analysis of the KDD CUP 99 data set. In Proceedings of the 2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications, Ottawa, ON, Canada, 8–10 July 2009; pp. 1–6.
23. Shiravi, A.; Shiravi, H.; Tavallaee, M.; Ghorbani, A.A. Toward developing a systematic approach to generate benchmark datasets for intrusion detection. *Comput. Secur.* **2012**, *31*, 357–374. [CrossRef]
24. Andrade, R.; Pinto, T.; Praça, I. Trust Model for a Multi-agent Based Simulation of Local Energy Markets. In *International Conference on Practical Applications of Agents and Multi-Agent Systems*; Springer: Berlin, Germany, 2020; pp. 183–194.
25. Open Data Sets IEEE PES Intelligent Systems Subcommittee. Available online: <https://site.ieee.org/pes-iss/data-sets/> (accessed on 8 May 2012).
26. Garcia, S.; Parmisano, A.; Erquiaga, M.J. IoT-23: A Labeled Dataset with Malicious and Benign IoT Network Traffic (Version 1.0.0). Available online: <https://zenodo.org/record/4743746#.YRscFt8RVPY> (accessed on 6 August 2021).

## Article

# SPD-Safe: Secure Administration of Railway Intelligent Transportation Systems

George Hatzivasilis <sup>1,2,\*</sup>, Konstantinos Fysarakis <sup>3</sup> , Sotiris Ioannidis <sup>4</sup>, Ilias Hatzakis <sup>2</sup>, George Vardakis <sup>2</sup>, Nikos Papadakis <sup>2</sup> and George Spanoudakis <sup>3</sup>

<sup>1</sup> Institute of Computer Science, Foundation for Research and Technology—Hellas, Vassilika Vouton, GR-70013 Heraklion, Greece

<sup>2</sup> Electrical and Computer Engineering, Hellenic Mediterranean University (HMU), Estavromenos, GR-71410 Heraklion, Greece; hatzakis@cs.teicrete.gr (I.H.); gvardakis@cs.hmu.gr (G.V.); npapadak@cs.hmu.gr (N.P.)

<sup>3</sup> Sphynx Technology Solutions AG, Innovation Department, 6300 Zug, Switzerland; fysarakis@sphynx.ch (K.F.); spanoudakis@sphynx.ch (G.S.)

<sup>4</sup> Electrical and Computer Engineering, Akrotiri Campus, Technical University of Crete, GR-73100 Chania, Greece; sotiris@ece.tuc.gr

\* Correspondence: hatzivas@ics.forth.gr; Tel.: +30-2810-391600

**Abstract:** The railway transport system is critical infrastructure that is exposed to numerous man-made and natural threats, thus protecting this physical asset is imperative. Cyber security, privacy, and dependability (SPD) are also important, as the railway operation relies on cyber-physical systems (CPS) systems. This work presents SPD-Safe—an administration framework for railway CPS, leveraging artificial intelligence for monitoring and managing the system in real-time. The network layer protections integrated provide the core security properties of confidentiality, integrity, and authentication, along with energy-aware secure routing and authorization. The effectiveness in mitigating attacks and the efficiency under normal operation are assessed through simulations with the average delay in real equipment being 0.2–0.6 s. SPD metrics are incorporated together with safety semantics for the application environment. Considering an intelligent transportation scenario, SPD-Safe is deployed on railway critical infrastructure, safeguarding one outdoor setting on the railway's tracks and one in-carriage setting on a freight train that contains dangerous cargo. As demonstrated, SPD-Safe provides higher security and scalability, while enhancing safety response procedures. Nonetheless, emergence response operations require a seamless interoperability of the railway system with emergency authorities' equipment (e.g., drones). Therefore, a secure integration with external systems is considered as future work.

**Keywords:** intelligent transportation; railway; CPS; security; safety; critical infrastructure

**Citation:** Hatzivasilis, G.; Fysarakis, K.; Ioannidis, S.; Hatzakis, I.; Vardakis, G.; Papadakis, N.; Spanoudakis, G. SPD-Safe: Secure Administration of Railway Intelligent Transportation Systems. *Electronics* **2021**, *10*, 92. <https://doi.org/10.3390/electronics10010092>

Received: 18 November 2020

Accepted: 29 December 2020

Published: 5 January 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Railways continue to be one of the main transport systems nowadays [1,2], covering public, private, and military needs over a wide operational area. Thus, railway assets are an attractive target for malicious actors and are exposed to various threats, from natural events to man-made ones, such as terrorism or vandalism (e.g., [3–6]).

The associated risks are exacerbated by the fact that railway infrastructure assets are typically placed along the route, including remote areas where physically protecting them is challenging. Moreover, railway premises have a large attack surface (due to their numerous electronic and electrical parts, such as power supply, switches, scheduling, and other subsystems), but often reside far from the main stations. While auditing for physical threats is quite important [7], the premises are usually inspected remotely through cameras. Sensory equipment is also deployed to monitor environmental parameters. The goal is to prevent potential intruders [8–10], avoid machinery overheating, and detect fires. Since the

interconnection of this monitoring equipment is, at least partly, wireless, it can become a target of several types of attacks.

In this context and considering that a successful attack could damage the railway's operation or even cause severe injuries and deaths, cybersecurity is an important consideration for such interconnected critical systems [11,12]. Attackers can disrupt communications (e.g., through jammers) or even infiltrate the networks and take control of critical equipment [13]. Cyber-attacks on the command-and-control centers (C&C) and the information systems are also feasible [13,14]. Thus, the secure interconnection of all the deployed elements and platforms is important, and the cyber and physical security of the critical infrastructure becomes imperative [9,10].

As sketched above, safety is another design factor and one that is closely related to security. Cargo and passengers are transported in high volumes each day, covering long distances. In the past, railway accidents have caused a number of deaths, along with significant financial losses [15]. While the introduction of electronic controllers reduced the occurrence of such situations [16], safety risks cannot be ignored, considering the wide railway coverage that still includes aspects such as uninspected car-crossings, system malfunctions (like signal loss), and, of course, the human factor [17,18].

Within the ever-changing technological landscape, there is currently a move from automated to intelligent cyber-physical systems (CPS), motivated by the speedy infiltration of the Internet of Things (IoT) and cloud computing and enabled by wireless networking [19–22]. Wireless sensor networks (WSNs) [23] can cover the wide railway operational territory, gathering and processing pieces of ambient knowledge, while gateways can be used to transmit the data to the controlling center or a cloud service. The railway controlling software at the backend can, then, collect and integrate the spatial information and manage the underlying subsystems [24,25]. Therefore, WSNs are an ideal solution for covering the railway operating area, including the railway routes and various scattered shelters.

However, the railway cyber infrastructure and networks currently only adopt rudimentary defenses (e.g., cryptography), which provide protection against the most basic threats, forfeiting effective ways of detecting advanced cyber-attacks [26]. While initially designed as closed systems, current infrastructure networks are vulnerable to various network layer attacks, like blackhole, badmouthing, and jamming attacks [27].

Motivated by the above, this work presents “SPD-Safe”, (security, privacy, and dependability (SPD)), an administration framework for railway CPS, aiming to enhance the security, privacy, dependability, and safety of the intelligent railway infrastructure, while enabling services for monitoring and managing the overall setting. The framework integrates mechanisms for mitigating cyber-attacks attempting to disrupt communications or compromise infrastructure assets, and periodic malfunctioning of assets is also taken into consideration. SPD-Safe can act as an intelligent communications-based train control (CBTC) system for railway CPS, leveraging artificial intelligence (AI) to manage the system at runtime. The system uses standardized solutions, and its building blocks can be easily retrofitted in current deployments.

In addition to the detailed description of the proposed framework, a preliminary implementation is described and evaluated, concentrating on the management of: (a) In-carriage, and (b) on-route sub-systems. WSNs are deployed inside the carriage and by the railway tracks to safeguard carriages that transfer dangerous freight and to help avoid crashes with objects blocking the train's route (like stuck vehicles on rail track crossings), respectively. Furthermore, smart cameras are installed to improve the physical security of the critical infrastructure. In the context of the two use cases (a) and (b) above, through SPD-Safe the railway CPS is configured in real-time to tackle ongoing cyber-attacks and control safety-related incidents. This hands-on validation was developed and demonstrated under the EU-funded project new embedded Systems architecture for multi-Layer Dependable solutions (nSHIELD) [28], with the cooperation of major industrial partners in the railway and defense domains, including Ansaldo STS (<http://www.railway-technology.com/contractors/signal/ansaldo-sts/>), Selex ES (now Leonardo



S.p.A.: <https://www.leonardocompany.com/en/home>), and HAI (<http://www.haicorp.com/en/>). Simulation analysis was also conducted during the design phase, utilizing the security-aware Cyber-Physical Systems (CPS) Simulator Framework (COSSIM) [29], paving the way for the final installation of the proposed system, as presented in the following sections.

The rest of the paper is structured as follows: In Section 2, related work on railway signaling systems is reviewed. In Section 3, the middleware platform and intelligent agent technologies that manage the underlying equipment are presented. In Section 4, the network layer protection mechanisms are detailed. In Section 5, the implementation details of SPD-Safe are provided and the application in the railway setting is demonstrated. The proposed system is also compared with relevant systems in Section 6, while Section 7 features the concluding remarks.

## 2. Materials and Methods—Related Work

Smart transportation ecosystems involve, among others, passenger services as well as critical infrastructure-related applications and the associated safeguards. The fundamental goals in this context include “green” (i.e., environment-friendly) operation, improved performance and efficacy, as well as enhanced security and safety.

Railways, in specific, rely on signaling systems that direct the trains’ traffic. Infrastructure control and management is achieved via various telecommunication means that are installed on carriages and tracks. Communication between track equipment and trains is achieved via CBTC signaling systems [30–32] enabling the railway’s management and infrastructure control. For the European Union (EU), the international wireless communications standard for railways includes the European Train Control System (ETCS) [33]. The communication baseline is implemented by the Global System for Mobile Communications—Railway (GSM-R) [34], which is further enhanced with the General Packet Radio Service (GPRS) [35] and forms the base of an intelligent transportation application. ETCS utilizes trackside equipment that transmits information regarding the route to unified controlling equipment within the train cab. Thus, all lineside data are passed wirelessly to the driver, without requiring the direct observation of lineside visual signals, as was the case in legacy railway settings. The adoption of ETCS results in more and longer running trains, with increased traffic and railway management capabilities.

In addition to the signaling developments, WSNs can now cover a wide railway operational area, gathering ambient data. Embedded systems implement intelligence solutions encompassing the underlying critical assets as well the interlinked smart city ecosystems. Related frameworks for intelligent monitoring of the critical infrastructure have already been proposed in the literature (e.g., [36,37]). The Integrated System for Transport Infrastructure surveillance and Monitoring by Electromagnetic Sensing (ISTIMES) project [36] implements a transport infrastructure surveillance and monitoring system with electromagnetic sensing. Distributed and local sensory equipment (e.g., optic fiber sensors, infrared thermography, low-frequency geographical techniques, etc.) are utilized to perform non-destructive electromagnetic sensing and monitoring of the critical infrastructure. The Cloud to Infrared Thermography (Cloud2IR) [37] deploys an infrared and environmental Structural Health Monitoring (SHM) information system. The software architecture enables multi-sensor connection and the interplay with cloud computing services (e.g., data aggregation, system management, etc.). However, the heterogeneity of the deployed equipment and diverse demands of the various applications make the administration of the underlying infrastructure a challenging task.

In parallel, as Service-oriented Architectures (SoAs) increase in popularity, a continuous effort to deploy SoAs within the Industrial IoT (IIoT) domain and the smart railway CPSs can be observed. Several technologies are proposed that support the required functionality, ranging from agent frameworks and middleware platforms, to communication protocols and data representation standards. Such state-of-the-art solutions are presented in the subsections that follow.

### 2.1. Management Platforms and Reasoning Systems

Agent technologies constitute the typical option for modeling ambient intelligent systems that exchange information with the environment and user [38,39]. Intelligent agents inspect the surrounding setting and react to upcoming events at normal operation. Their AI modules process context-aware data, as collected from the surrounding environment by the attached devices.

Regarding the various agent technologies, 24 frameworks were analyzed in Kravari and Bassiliades [40], including the popular Java Agent DEvelopment framework (JADE), Agent Globe (A-GLOBE), and Jason (the hero's name from Greek mythology). JADE implements the relevant standards for Semantic Web and the Foundation for Intelligent Physical Agents (FIPA: <http://www.fipa.org/>) (e.g., the Agent Communication Language (ACL) [41]). The platform is easy to learn and user-friendly, while offering portability and compatibility with all Java Virtual Machines (JVMs). The open-source and stable developer versions operate with several programming languages, such as Java, Jess, and Prolog. The agent communication is fast, and the overall framework is efficient and scalable. Moreover, JADE supports strong user authentication and cryptographic solutions—i.e., JADE security (JADE-S)—along with Hypertext Transfer Protocol Secure (HTTPS). The framework is widely-used and is deployed in several fields, including reasoning in multiple domains, general purpose applications, mobile computing, and e-commerce. The study of Kravari and Bassiliades [40] also infers that JADE is the most popular framework due to the pure Java design and the co-operation with several web systems. In addition, five respectable organizations (France Telecom, Motorola, Profactor, TILAB, and Whitestein TEchnologies AG) supervise the framework [40].

Regarding middleware systems and messaging protocols, a comparative analysis of relevant IoT solutions (Constrained Application Protocol (CoAP), Message Queuing Telemetry Transport (MQTT), and Devices Profile for Web Services (DPWS)) was conducted in Fysarakis et al. [42]. DPWS [43], by the Organization for the Advancement of Structured Information Standards (OASIS: <https://www.oasis-open.org/>), constitutes the benchmark in terms of ease-of-design. The framework is flexible and robust in terms of service eventing, discovery, and subscription; following initialization, the underlying devices can discover the provided services and communicate in a seamless manner.

Finally, concerning the deductive rule engines that enable the AI reasoning features, National Aeronautics and Space Administration (NASA) examines the capabilities of several related approaches (Jess, Drools, Microsoft Business Rule Engine, Official Production System Java (OPSJ), and Intelligence Logiciel (ILOG)), as described in [44,45]. Jess is efficient and excels in many categories. It works with dynamic facts and dynamism in variables and rules and is appropriate for NASA's mission-critical applications as well as many other research areas [46,47].

### 2.2. Intelligent Railway Systems

To the best of the authors' knowledge, only a few multi-agent systems (MAS) have been developed and tested on actual railway environments [48,49]. In addition, despite strong industrial involvement, not all developed agent capabilities are fully used in practice and thus, the full potential of agent technologies is not exploited fully. Three indicative installations on railway settings are: Train Integrity (TrainIntegrity) [50], Condition Monitoring of a Light Rail Vehicle and Track (CMLRVT) [51], and Sensor Networks for Railways (SENSORAIL) [52].

TrainIntegrity utilizes WSNs to check the integrity of cargo trains [50]. The nodes consist of the RCM 3400 RabbitCore module and sense environmental parameters. The WSN raises an alarm if it infers that an unexpected change has occurred in the train's composition. CMLRVT is built and tested on a tramway operation in Poland. The system consists of a dispersed sensor network that is installed on the vehicles and railway infrastructure, along with the data acquisition component and a data server that maintains the artifacts of the management and analysis procedures. At first, the system collects data during

normal operation, which is stored in the server. Then, the new pieces of knowledge that are sensed by the devices are compared with the nominal values. The detected variations are further analyzed by the system, revealing safety-related incidents (e.g., rail cracks) that are presented to the user through a dedicated application. However, neither of the two systems considers security issues at the network link, nor do they integrate and manage the heterogeneous underlying embedded systems.

SENSORAIL [52] is an early warning system for the railway monitoring infrastructure. WSNs collect and integrate data to enable the detection of structural failures and security threats. Sensor clusters communicate information towards distant controlling centers through GSM-R/GPRS mobile equipment. The integration of heterogeneous sensors is managed by a component referred to as “scalable software architecture for the integration of heterogeneous sensor systems” (SeNsIM) [53], while the detection of events is made by a model-based data correlation component, called “novel framework for the detection of attacks to critical infrastructure” (DETECT) [54]. SENSORAIL specifies the examined threats in the Event Description Language (EDL) [55] and maintains them within a scenario repository. Upcoming events are stored in a history database and model-checking is performed at runtime. Regarding the middleware and agent platform, SeNsIM does not utilize any semantic technologies and does not support related standards, thus lacking in terms of interoperability and ease of integration with existing setups. Moreover, SENSORAIL does not include any protection mechanisms, solely focusing on the detection of threats.

### 2.3. Network Layer Protection

Several schemes are suggested in the literature for protecting communication in WSNs, attempting to address pertinent security concerns (e.g., [56–59]).

The Reputation-based Framework for Sensor Networks (RFSN) [56] authenticates underlying nodes with the Timed Efficient Stream Less Tolerant Authentication protocol ( $\mu$ TESLA) [60], implementing a beta Bayesian formulation for fading and evaluating the reputation of the routing operation and the legitimacy of the reported sensed variables. Ariadne [57] also utilizes a TESLA variant for authentication, collecting feedback regarding the successful delivery of packets to choose optimum communication paths and avoid malicious behavior. The Cooperative Secure routing protocol based on ARAN (CSRAN) [58] integrates digital certificates and asymmetric cryptography for authentication. As in the case of RFSN, it uses a Bayesian distribution for fading and, when a node detects malicious activity, it automatically re-routes communication from that point on. The Secure Resilient Reputation-based Routing (SR3) [59] adopts lightweight cryptography (LWC) and symmetric modules for security and authentication. Fading is accomplished by a First In, First Out (FIFO) finite list. The system combines reputation with a reinforced random walk algorithm, producing enhanced load-balancing at the cost of a high intermediate forwarding node count.

Despite the plethora of proposed solutions, and while most can tackle basic security attacks and malfunctioning cases, there are still various open avenues for attackers, including flooding attacks in congested periods, topology-related attacks, and jamming [61,62].

## 3. Administration of IoT Deployments

Considering the landscape sketched above, this section presents the proposed SPD-Safe solution, and more specifically the deployed platform and the reasoning process of each SPD-Safe agent. The core reasoning engine has been previously presented by the authors in [63]. This version enhances the network layer security and is applied in a mobile setting, forming an intelligent transportation system that complies with the real-time requirements of CBTC for railways.

SPD-Safe comprises a framework that integrates variants of the aforementioned primitives (i.e., agents, middleware, rule engine, and network layer protection) across all system layers to implement an efficient, scalable, practical, and easy-to-deploy and maintain solution, with adequate reasoning and management capabilities. From top-to-bottom, the

system consists of four layers: (i) An overlay with intelligent agents that control distinct subsystems; (ii) a middleware platform that enables communication between the agents and underlying networks; (iii) the network layer that consists of interconnected IoT devices; and (iv) the node layer that represents the devices themselves. The core technological building blocks will be detailed in the subsections that follow.

### 3.1. Agent Technologies & Middleware Solutions

SPD-Safe utilizes JADE [64] as the top-layer multi-agent system. It adopts and implements standardized approaches to agent deployment, such as the ACL [41] by FIPA. The JADE-S add-on [65,66] safeguards communication at the overlay and offers built-in security functionality for confidentiality, integrity, authentication, and authorization.

Then, each agent is ported as a bundle in the middleware platform Open Service Gateway initiative (OSGi) [67]. Through it, an agent can monitor the underlying subsystem, enhancing real-time management. Network gateways also deploy a controlling bundle in the same platform, defining the offered functionality as a service in the DPWS standard [43]. The agent and the related network bundles interchange well-structured semantic data, as defined in the OASIS standardized Common Alerting Protocol (CAP) [68]. The OSGi platform also provides its own built-in security features for the inner-platform communication, limiting bundle functionality to pre-defined capabilities and protecting both the agent and controller bundles.

Here, other than these built-in features that are provided by the deployed platforms at the overlay and middleware layers, SPD-Safe integrates an additional defense mechanism for the network and node layers, namely Secure Route (SecRoute) [61], a security protocol that protects the wireless ad hoc communication of the underlying embedded devices. This protocol counters several types of threats and attacks at the network link, protects the nodes' assets and their resource consumption, and acts as an intrusion detection module for the upper layers. When a security incident is recorded, the network gateway bundle will send related CAP messages to the responsible agent, which may take further action. SecRoute is detailed in the next section.

Metrics that evaluate the various system aspects are now an integral feature of the development cycle. They offer a quantitative indication regarding the compliance with the targeted requirements of the application domain. An evaluation method for the estimation of the security, privacy, and dependability (SPD) properties for configurable embedded systems is presented by the authors in Hatzivasilis et al. [63]. For every configuration option, the metrics derive a triple vector of <Security, Privacy, Dependability>, whereby the vector's factors are assigned a value from 0–100, representing no to full protection respectively. SPD-Safe adopts this methodology to enable a metric-driven SPD- and safety-aware administration, where the reasoning procedure triggers runtime system adaptations to reach specific SPD goals [69].

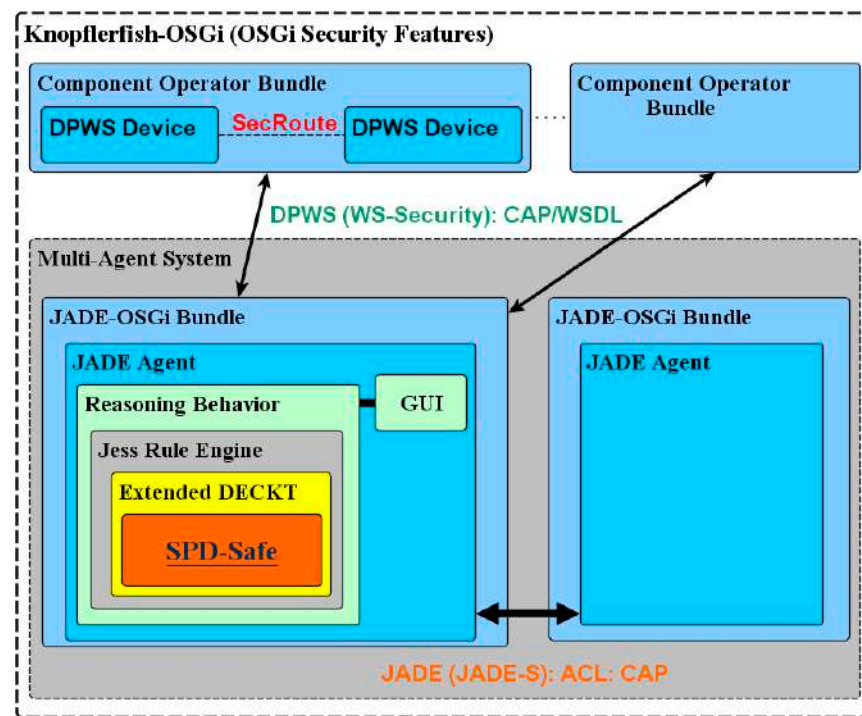
### 3.2. Reasoning Capabilities & Conflict Resolution

The artificial intelligence (AI) behavior of each agent is developed in the rule engine Jess [70,71]. For knowledge representation and reasoning, the Jess-EC [72] is used. The latter is an Event Calculus (EC) [73] implementation in Jess, offering the required semantics modeling. SPD-Safe's software layers are illustrated in Figure 1.

Each agents' AI procedure implements automated temporal, casual, and epistemic reasoning with real-time events, action preconditions, rule priorities, indirect effects, context-sensitive side-effects, as well as the common law of inertia. Moreover, the reasoning capabilities can cope with the requirements of dynamic and partially known or uncertain domains.

However, as agents exchange information, contradictory reasoning results may occur due to the local viewpoint of each entity and the lack of global knowledge. Thus, for resolving conflicts, SPD-Safe introduces the epistemic mechanism of share theories [63]. The participating entities send the involved theory rules to a mediator agent, along with

the recently sensed local events. The mediator combines these elements and performs a reasoning operation that determines the final outcome and the state of the conflicting assets.



**Figure 1.** The software layers of the proposed security, privacy, and dependability (SPD)-Safe framework.

Nevertheless, if an agent utilizes protected data that must be maintained locally and not distributed (e.g., confidential information regarding user policies or system settings), it will not be able to contribute in the share theory with its full knowledge. For this occasion, an alternative relational grading mechanism, called certainty degree [63], resolves the affair quickly and efficiently. The mechanism utilizes subjective criteria as well as the agents' roles and hierarchy, marshaling the problem without constructing the related share theory and retaining the system's coherency. Thus, the certainty degree is applied in affairs where reasoning with locally protected data is involved, otherwise a share theory is constructed.

### 3.3. SPD Measurement

The SPD multi-metric methodology [69] measures the provided protection level of a system and its various configurations. The system's perimeter is identified and the data sources, entry, and exit points are recorded. Then, the mechanisms that protect each of these elements are assessed based on the standardized Criteria Evaluation Methodology (CEM) [69]. This involves the attack potential risk analysis that evaluates the attacker's motive to misuse specific system elements, expertise, and the resources that they are willing to devote for an attack. Henceforth, five parameters are examined for the analysis of a potential threat:

- Required time: The time that it is required to perform a specific attack (e.g., in days or weeks);
- Expertise: The technical skills and knowledge that the attacking group can exhibit (such as copy-cat, advanced, or expert);
- Knowledge of the target: Familiarity with the targeted system and its operation (e.g., public, sensitive, or critical information concerning some subsystems, etc.);
- Window of opportunity: The attacker may require appreciable access to the system in order to exploit a vulnerability and avoid detection;
- Resources: The software, hardware, or other equipment that is necessary to perform an attack (such as specialized or common resources).

The method does not investigate every possible attack but educes a good indication of the defense status in accordance with standard ratings. The protection level for each of the three SPD properties is calculated by integrating the risk analysis with the efficacy of the installed defenses against known attacks and/or other limitations (e.g., based on the latest reports from Computer Emergency Response Teams (CERTs) or Common Vulnerabilities and Exposures (CVE) repositories). The result is a value in the range of 0–100, where 0 represents the absence of defense mechanisms and 100 represents full protection. The final outcome is a vector of <Security, Privacy, Dependability>, which represents the total SPD value of the currently composed setting of the system. The SPDs of different system configurations can be estimated either in advance or at runtime. The first option is leveraged by the AI units of SPD-Safe in order to perform proactive and/or automated changes in the state architecture when a safety or security event occurs. The second option provides indications to the human operator in order to take decisions and make manual interventions.

Therefore, the protection status of all mechanisms and their integration in the demonstration examples are pre-calculated based on this method, as described in Sections 4 and 5. Then, automated administration policies are triggered in response to real-time events, as presented in Section 5.

These features enable the implementation of a relative novel protection strategy, called Moving Target Defenses (MTDs) [74]. When a system is stable, it is seen as a “sitting duck” by the attacker, who has plenty of time to analyze it, detect potential vulnerabilities, and exploit them. With MTD, a system that is aware of the defense level of its various components, their configurations, and the integration of all of them, can alter the setting automatically or semi-automatically in a periodic fashion. The AI modules are always keeping the system in a secure state, while the different configuration and architectural sets increase the system states that have to be analyzed by the attacker. In addition, the time that a specific setting remains active is determined by the time required for an average hacker to analyze it (i.e., based on the “Required Time” factor of the attack potential risk analysis). Performing attacks is becoming quite hard, while the window of opportunity for the malicious entities has significantly decreased.

### 3.4. AI Processing & Performance

The reasoning component of Jess implements the RETE algorithm (Latin word for net, meaning network in this domain) [75]. This is the most widely-used pattern matching technique for rule-based systems and is optimized for speed. Scalability and performance are affected by the three factors of: (i) The rules’ volume (R), (ii) the average number of patterns in the left-hand-side of each rule (P), and (iii) the facts in the working memory (F). Computational complexity is linear to the working memory size and in the order of  $O(RPF)$ . For each SPD agent in the railway mission-critical applications that are examined in the following sections, the theory rules volume (R) is very low (around 30 rules per scenario). In order to reduce the pattern-matching space, unique identifiers are assigned to every modeled entity, and therefore, occurring events affect specifically defined parameters, keeping the pattern-matching ration low (P) and in the order of 1–3. Performance is mostly influenced by the number of facts (F). In the demonstrated cases, it requires 10–20 facts per scenario. The computational overhead for an SPD agent is in the range of a nanosecond with additionally 50 bytes in memory.

For a central agent that collects information from the whole railway system, it requires around 500 facts and 40 rules to model the underlying setting. At boot time, the reasoning engine takes to run around 1.6 s, 87 MB for code, and 45 MB in RAM. Then, a reasoning process for a theory and a few hundreds of facts would require 0.002 s on average, representing the actual delay that affects the applications.

### 3.5. Relevant Methodologies for Secure IoT Modeling

Over time, several solutions have been proposed that try to resolve the open issues of capturing the security posture of an IoT or other system and facilitate its administration [76–80]. Eby et al. [76] integrated the Simple Modeling Language for Embedded Systems (SMoLES) with the Security Model Analysis Language (SMAL) [76]. SMAL provides security extensions to the composition meta-model of the Domain Specific Modeling Language (DSML) [77] and can express access control policies for IoT applications. The resulting framework is called SMoLES Security (SMoLES-SEC). However, its reasoning capabilities are bounded due to the constrained expressiveness of the underlying SMAL. Furthermore, SMoLES-SEC cannot deduce which security characteristics hold after the compositions of two components or the final security status of the composed system.

Service Dependency Trees (SDTs) [78] support the verification of service secure composition in IoT ecosystems. The IoT devices/nodes construct their own SDT. For each provided service, the relevant SDT defines the potential external service nodes that the service is depending on. The nodes are also aware of all recursive SDTs for their composed services. Thus, secure service composition is performed by enabling integration only with SDTs where all paths and involved entities are trusted. On the other hand, creating a SDT for a real IoT application is not trivial, while trustworthiness and consistency in an actual complex and dynamic environment may be challenging.

Albanese et al. [79] utilize attack surface metrics in order to evaluate the security aspects of system and materialize MTDs strategies. This solution calculates the distance of the security surface of the various system states. The goal is to administrate responses against ongoing attacks as well as to deduce a system setting that exhibits specific desirable parameters. Techniques for assessing and reducing the cost for the defender are also included.

Savola and Sihvonen [80] propose a MTD approach based on a multi-metric-driven management framework. The overall solution has been applied in an e-health digital environment for chronic diseases [80], where three metric types are considered. Risk-driven security assurance and engineering metrics are defined at deployment-time to offer an early assessment on the deployed defense mechanisms and their effectiveness. Continuous security monitoring metrics are determined at operational-time, enabling the security correctness assessment, enhanced systematization, and traceability of the various product requirements and involved metrics. Thereupon, automated adaptive decision-making metrics are assigned at operational-time and accomplish a higher quality security effectiveness understanding in operational security auditing and future versioning of the system. The method supports continuous security monitoring and automated metric-driven security-related actions.

Table 1 presents the outcomes of the qualitative analysis. The modeling expressiveness of SPD-Safe is quite general and can also be utilized in complex and dynamic systems. Moreover, it assesses all three security, privacy, and dependability properties and can evaluate their status both before a composition is performed and after the integration of the system. As with the other relevant approaches, the MTD features are driven by metrics and SPD-Safe provides a concrete implementation of this modern defense type. The overall solution fits with the distributed nature of IoT ecosystems and can resolve conflicts that may arise due to knowledge sharing between the various entities.

**Table 1.** AI (artificial intelligence) modeling features.

Feature	SPD-Safe [This Paper]	SMoLES-SEC [76]	SDT [78]	Attack Surface MTD [79]	Multi-Metric-Driven MTD [80]
System composition					
Expressiveness	Y	N	Y	N	N
generality	Y	Y	Y	N	N
Dynamicity	Y	Y	Y	N	N

Table 1. Cont.

Feature	SPD-Safe [This Paper]	SMoLES-SEC [76]	SDT [78]	Attack Surface MTD [79]	Multi-Metric-Driven MTD [80]
Validation					
Pre-composition	Y	Y	Y	N	N
Post-composition	Y	N	N	N	N
Evaluated Properties					
Security	Y	Y	Y	Y	Y
Privacy	Y	N	N	N	N
Dependability	Y	P	P	N	P
Artificial Intelligence					
Distributed reasoning/processing	Y	N	Y	N	Y
Conflict resolution	Y	N	N	N	N
MTD	Y	N	N	Y	Y

Y(es), N(o), P(artial). Service Dependency Trees (SDTs); Moving Target Defenses (MTDs); Simple Modeling Language for Embedded Systems Security (SMoLES-SEC).

#### 4. Network Layer Security

The protection of the network link is essential in order to safeguard the underlying systems of critical railway infrastructure (e.g., WSN, signaling equipment, surveillance, etc.). For this purpose, as mentioned, SecRoute [61] is developed; a novel defence primitive that provides the core security properties for authentication, integrity, and confidentiality, along with energy-aware secure routing and authorization.

The secure routing protocol protects the involved entities from malicious operations while improving performance and offering load-balancing. It consists of three main primitives:

- The cryptographic service with the Timed Efficient Stream Less Tolerant Authentication protocol ( $\mu$ TESLA) [60], which implements message authentication, confidentiality, and integrity;
- The efficient secure routing service with the Self-Channel Observation Trust and Reputation System (SCOTRES) [62] that safeguards the communication link against ad-hoc routing attacks and network layer vulnerabilities;
- The authorization service with the Policy-Based Access Control framework (PBAC) [42], which offers authorization and access control based on policies.

Table 2 summarizes the overall security properties that are provided by the integrated network layer defense mechanism and the relevant threats and attacks that are countered, while a brief analysis is presented in the subsections that follow. More details regarding the three services are presented in the relevant papers for  $\mu$ TESLA [60], SCOTRES [62], and PBAC [42], respectively.

Figure 2 presents the block diagram of the main SPD-Safe modules and their connection.

Table 2. Protection aspects of the SPD-Safe's network layer security.

Primitive	System Property	Countered Threats
$\mu$ TESLA	Authentication	Impersonation, Sybil attacks
	Integrity	Data tampering, modification, interruption
	Forward security	Replay attacks
	Confidentiality (optional)	Disclosure



Table 2. Cont.

Primitive	System Property	Countered Threats
SCOTRES	Topology-awareness	Attacks on topology-significant entities
	Energy-awareness & Load-balancing	Energy dissipation, overloading attacks on congested periods
	Channel health	Jamming attacks
	Reputation	<p>Malicious or selfish activity on the network operations of:</p> <ul style="list-style-type: none"> <li>- Routing (link spoofing, routing table poisoning, HELLO flooding, false link break, loops, nonexistent paths),</li> <li>- Forwarding (blackhole, grayhole, sleep deprivation),</li> <li>- Or making recommendations (badmouth, ballot-stuffing)</li> </ul>
	Trust	Overall misbehavior on the previously mentioned networking perspectives
	Secure routing	General attacks on the pure routing protocol (e.g., Denial of Service (DoS), inject arbitrary packets)
PBAC	Authorization based on policies	Unauthorized access

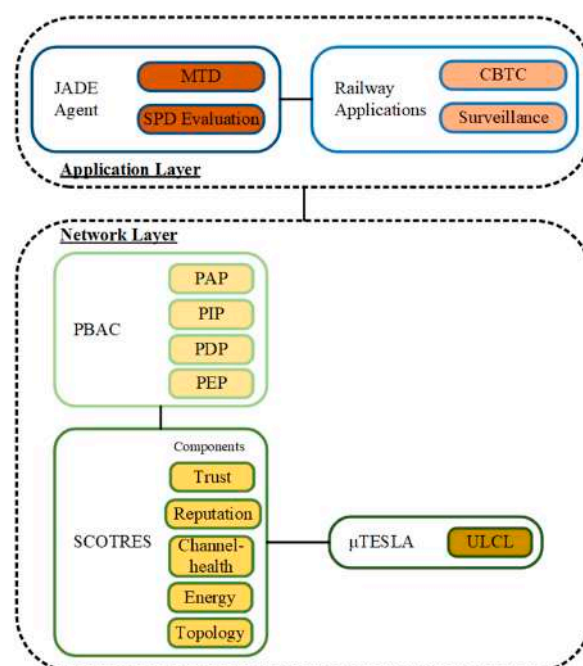


Figure 2. The building blocks of the SPD-Safe framework.

#### 4.1. Cryptographic Service— $\mu$ TESLAs

$\mu$ TESLA is a building-block for the Sensor Protocols for Information via Negotiation (SPIN) [81]. Loose time synchronization is required between the receiver and sender, with  $\mu$ TESLA utilizing broadcast messages and symmetric cryptography to implement the aforementioned core cryptographic properties. The security functionality of asymmetric cryptography is achieved by utilizing keyed Message Authentication Code (MAC) operations. In brief, the sender includes a keyed MAC on every transmitted packet, where this key is initially known only to this entity. Receivers maintain the received packets without authenticating the sender at this point. Shortly after, the key is revealed by the sender and then the receiver authenticates the packet and proceeds to further processing. Otherwise, the receiver discards the unauthenticated packets after a time-slot.

The protocol  $\mu$ TESLA is efficient and exhibits low computational and communicational overheads. It also tolerates packet loss and scales well for large networks. We use the Ultra-Lightweight Cryptographic Library (ULCL) [61] in order to develop the cryptographic functionality of  $\mu$ TESLA, adopting the Secure Hash Algorithm (SHA) with 256-bits message digest (SHA-256) for the MAC computations and the Advanced Encryption Standard (AES) with 256-bits cryptographic keys (AES-256) for the encryption/decryption.

#### 4.2. Secure Routing Service—SCOTRES

After authenticating a package with  $\mu$ TESLA, SCOTRES evaluates the sender's trustworthiness and its contribution to the network [62]. SCOTRES is a secure routing system for wireless ad-hoc systems that is based on trust computing and is designed around the intricacies of CPS solutions. It maximizes the information that is inferred regarding the network state, based on the knowledge that a node already processes. It safeguards communication against Internet-originating attacks or compromised equipment and jammers. The overall setting is utilized for real-time monitoring of IoT and CPS applications and their management through the cloud.

SCOTRES consists of five components that rate different aspects of the networking operation: (i) The topology-aware component improves the traffic load-balancing and defends distant entities from being isolated; (ii) the energy-aware component estimates the remaining energy of each node, defending the network against energy dissipation and other relevant threats; (iii) the channel-health component identifies jamming in the wireless medium, constraining its effects by routing communication through unaffected paths; (iv) the reputation component ranks a node's fair use of the network resources for routing, forwarding, and recommending activities; and (v) finally, the trust component aggregates all these pieces of knowledge and evaluates the trustworthiness and overall cooperativeness of network entities. Performance and security analyses for the five components have been conducted in [62].

#### 4.3. Authorization Service—PBAC

After verifying the message's legitimacy, the receiver node must decide if it will perform the requested action or not. The PBAC framework is used to implement this authorization functionality. The framework manages direct access to a smart device's resources as determined by a pre-defined collection of policies and rules that are modeled on the OASIS standards DPWS [43] and the eXtensible Access control Markup Language (XACML) [82]. PBAC consists of four components that are placed between the backend infrastructure and devices: The Policy Administrator Point (PAP) and Policy Information Point (PIP) that maintain the attribute values for creating and managing policies in a central repository, the Policy Decision Point (PDP) that runs on a trusted gateway node with sufficient computational capabilities, evaluates the request, and renders the authorization decision, and the Policy Enforcement Point (PEP) that enforces authorization at the end device and makes decision requests. These are combined to provide fine-grained, policy-based access control on assets from remote endpoints (like control stations, sensors, or cameras). Therefore, the specification of an active policy set can be used to define the

rights to access to acquired resources (e.g., sensed data and video/audio streams), the rights to update the settings, and even the rights to push notifications of emergency alerts (e.g., blocked routes and train crashes).

#### 4.4. Performance Evaluation

To assess the performance and validate the feasibility of the proposed approach, SecRoute is deployed on an embedded system which features BeagleBone (<http://beagleboard.org/bone>) devices and is integrated with the Distance Source Routing (DSR) protocol (DSR Uppsala University: <http://dsruu.sourceforge.net/>). BeagleBone is a low-cost and credit-card-sized device with ARM architecture, executing compact Linux operating systems (ARM Cortex-A8 processor at 720 MHz, 256 MB RAM, Ubuntu Linux). The devices sense environmental conditions, like humidity and temperature, and exchange data wirelessly with a central processing unit via a USB-WiFi.

We measure the processing overhead for SecRoute under normal operation without attacks taking place. Table 3 details the resource consumption of the proposed network layer defense. As indicated in the results, the calculation of reputation is the most compute-intensive part, as it maintains a history with previous interactions, which also increases the overall resource demands for trust computations. The requirements of authentication, routing, forwarding, as well as policy check are low. For the end-to-end interaction, the network latency is also low, ranging between 0.2–0.6 s on average. In an in-carriage setting where the distance among the nodes is short (a few meters) the transmission overhead is minimal, while the maximum delay is recorded for outdoor deployment where the nodes are placed hundreds of meters away from each other. The two scenarios are detailed in the next section.

**Table 3.** Resource allocation for SecRoute on BeagleBone devices.

Component	ROM (KB)	RAM (KB)	CPU (ms)
Cryptographic service			
Authentication	3.7	2.22	0.0020
Encryption	25.0	10.41	0.0028
Authenticated encryption	28.7	12.63	0.0048
Secure routing service			
Direct trust	5.6	4899.00	677.52
Reputation evaluation	20.0	1621.00	108.97
Indirect trust (recommendations)	30.0	185.00	37.90
Total trust	2.9	45,756.00	9.48
Accept route request	2.0	0.00	104.23
Suitable route selection	15.0	40.00	33.17
Authorization service			
PBAC policy check	24.7	36.00	7.50
Total resource consumption			
Total SecRoute	210.4	46,000.00	1652.50
DSR	310.0	90,000.00	2300.00
SecRoute_DSR	520.4	136,000.00	3952.50

#### 4.5. Comparison with Other Protocols

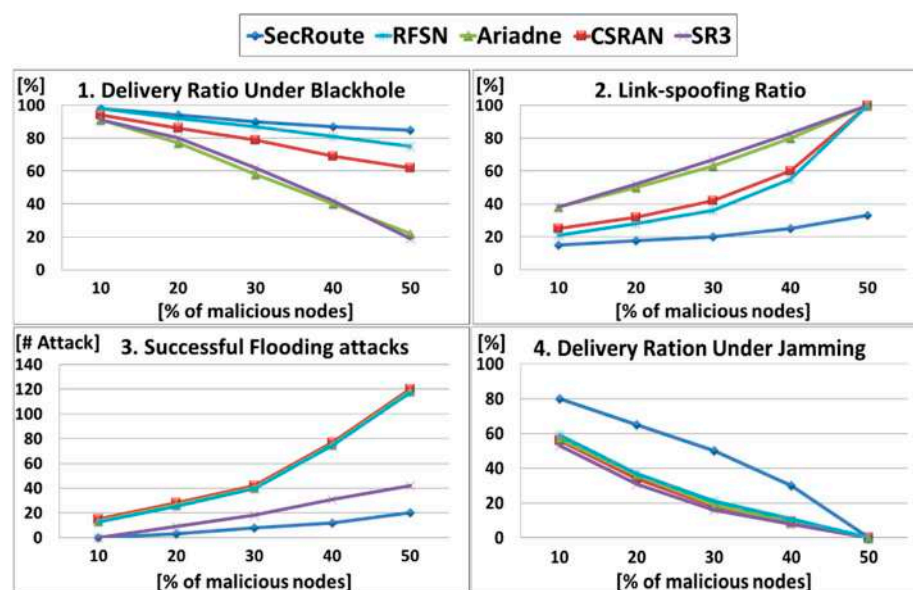
Efficiency and security analysis of the proposed network layer solution and five relevant systems (RFSN, Ariadne, CSRN, and SR3) have been presented by the authors in [62]. Table 4 summarizes the main features of the examined secure routing protocols.

**Table 4.** Secure routing protocols.

Property	SecRoute	RFSN	Ariadne	CSRAN	SR3
Authentication	$\mu$ TESLA	$\mu$ TESLA	TESLA	Certificates	LWC
Routing method	DSR	DSR	DSR	ARAN	Random walk
Reputation	Bayesian	Bayesian	NO	Bayesian	FIFO *
Fading					
Load-balancing	YES	NO	NO	NO	Partially
Energy-aware	YES	NO	NO	NO	NO
Anti-jamming	YES	NO	NO	NO	NO
Authorization	YES	NO	NO	NO	NO

\* FIFO: First In First Out finite list.

The secure network communication link of SPD-Safe is compared with the five most relevant proposals for protecting WSNs. Simulation analysis has been performed in the Network Simulator 2 (NS2: <http://www.isi.edu/nsnam/ns/>), analyzing the performance of each scheme and the provided protection level on a medium-size WSN with 50 nodes [62]. Four attack cases are considered for blackhole, ballot-based attacks for link-spoofing, topology- and energy-aware attacks, and jamming. For each setting, several experiments have been conducted, with the attackers' participation in the network ranging from 10–50%. Figure 3 presents the evaluation of the simulation results. SecRoute counters the attacks and outperforms the relevant schemes, providing the highest level of security and demonstrating the best energy- and load-balancing characteristics.

**Figure 3.** Simulation results for the evaluation of the network layer security solutions against four attack scenarios.

## 5. SPD-Safe Demonstration

### 5.1. Railway CPS Architecture

This section details the demonstration and evaluation of the whole SPD-Safe framework in the context of protecting and managing a railway CPS. In the proof-of-concept setting, our proposal assesses and manages the system and ambient ecosystem with the goal of safeguarding the trains' carriages and railway's routes. The hardware platforms incorporate embedded devices that control smart equipment (e.g., cameras and electronic doors), inspect environmental conditions, and exchange information wirelessly. Furthermore, the PBAC framework is applied for the control of the physical access for personnel, determined by access rights that are specified in XACML policies. Every agent manages a smart sub-system, like a train or a station. Backend agents can also run at the cloud in order to gather

high level information, perform big data analysis, and enable interaction with external systems and actuators. These agents run on virtual machines deployed on the research cloud platform GRNET Virtual Machines (ViMA: <http://vime.grnet.gr/about/info/en/>). Figure 4 illustrates the railway system architecture. The whole setting is administered by a master agent (MA) at the C&C. At the edge, simple and more lightweight agents (SAs) protect the local subsystems (applying access control, lightweight data analysis, incident detection, etc.) and exchange information with the MA (i.e., security/safety events and response strategies). The MA can, optionally, forward data to a cloud SA for storing or in-depth analysis. The cloud SA also presents high-level knowledge to end-users as well as the current SPD status of the railway infrastructure.

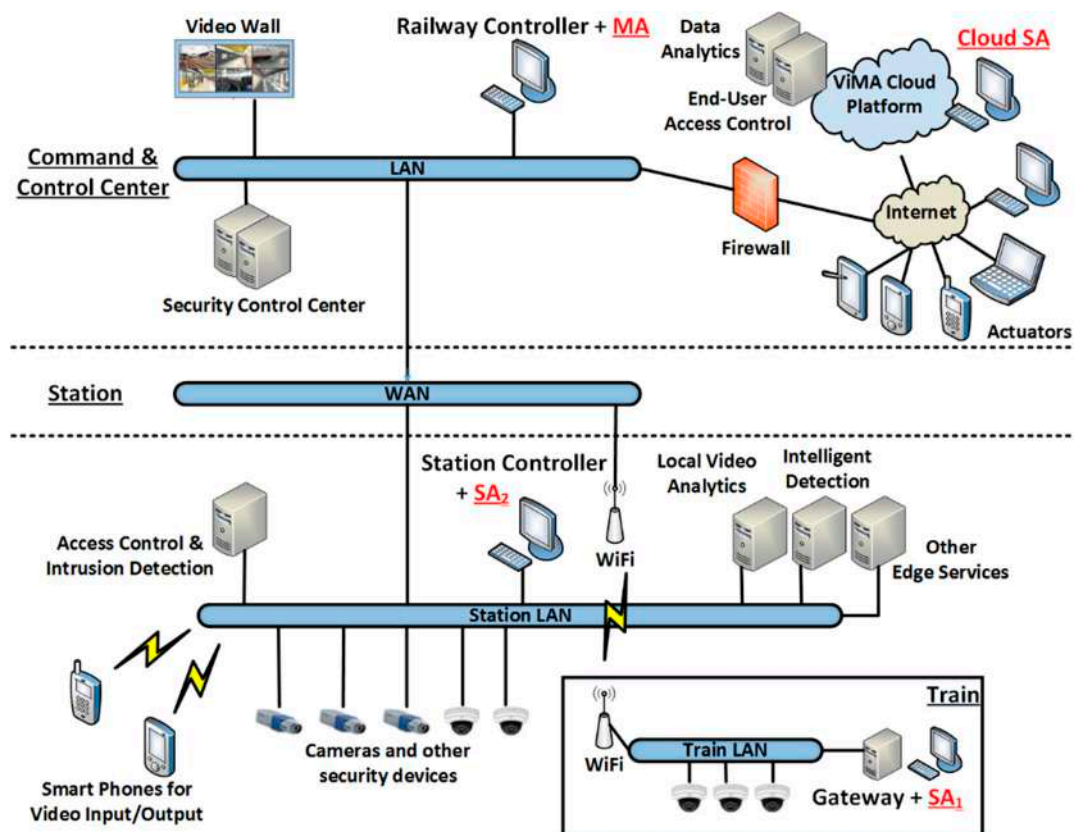


Figure 4. The smart railway use case architecture.

For this demonstration, the MA and the C&C services are deployed on a laptop. Both MA and cloud SA are installed on machines with a 2.1 GHz Intel Core i-7 processor, 8 GB of RAM, and the Ubuntu Linux Operating System (OS). The SAs are deployed on the BeagleBone devices at the edge systems.

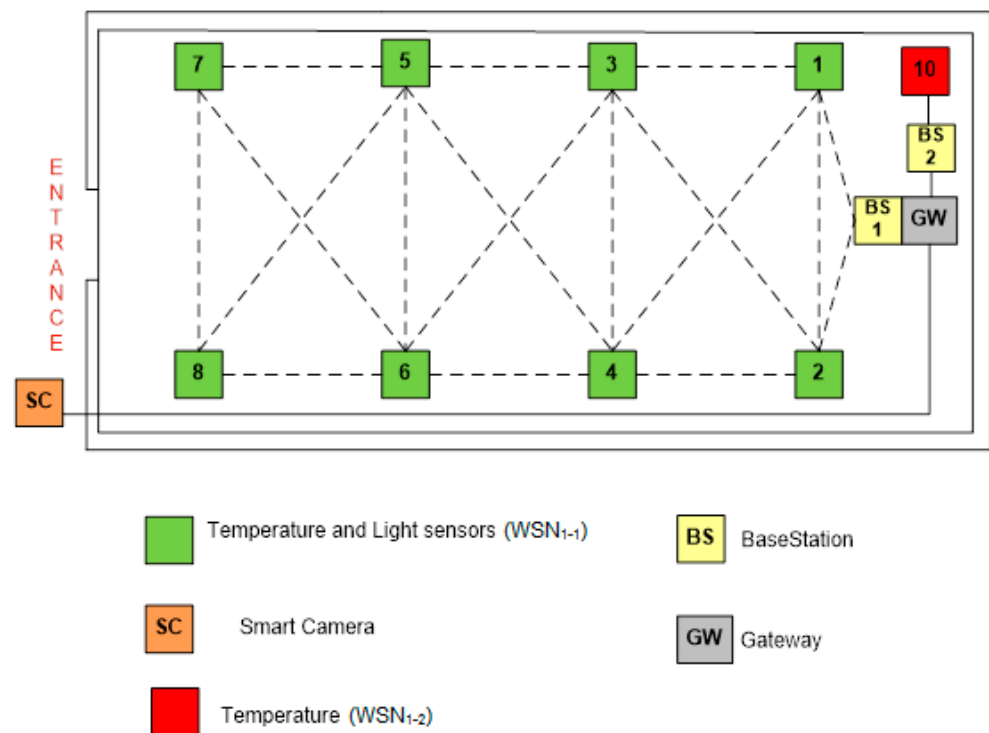
As a case study, two deployments are evaluated. In the first indoor setting, which emulates in-carriage or shelter equipment, we test the system under normal operation and the aforementioned attacks on routing. In the second outdoor scenario, which emulates the on-route equipment, we examine the system's response to safety-related incidents. Both networks run the SecRoute protocol [61] to enable communication, protect the network layer against cyber-attacks, and act as an intrusion detection and incident response system for the upper layers.

### 5.2. Indoor Setting—Cyber-Security

The demonstration setting includes a carriage/shelter inspecting application, which is equipped with a surveillance system and WSNs. Those components are sensitive to

network layer threats, like jamming and blackhole attacks. The deployed network is depicted in Figure 5, where these devices are deployed in a shelter [28]:

- At the entrance, the smart camera inspects for physical intrusion;
- Two WSNs are deployed in the shelter. WSN<sub>1-1</sub> (green color) monitors light and temperature, and WSN<sub>1-2</sub> (red color) senses temperature. WSN<sub>1-1</sub> and WSN<sub>1-2</sub> utilize different hardware to enhance diversity and ensure redundancy for the monitored factors;
- A gateway interconnects the rest of the components with the C&C.



**Figure 5.** The internal wireless sensor network (WSN) for the carriage setting.

WSN<sub>1-1</sub> consists of eight Memsic Iris sensor nodes (16 MHz Atmel ATmega 1281 processor, 8 KB RAM, Contiki OS). The devices are battery powered and measure light and temperature. Furthermore, the smart camera is controlled by the node at the carriage's entrance. WSN<sub>1-2</sub> is installed for redundancy and is comprised of Zolertia Z1 sensor nodes (16 MHz MSP430 processor, 8 KB RAM, and Contiki OS) that collect temperature data. The two WSNs are monitored by two relevant simple agents (SA<sub>1-1</sub> and SA<sub>1-2</sub> respectively). Every device executes the PEP module of the PBAC framework. The devices also exchange data with the gateway, which runs the access policies for PBAC, and communicates with an MA which administrates the whole network.

The devices gather environmental information and send data to the relevant base station (laptop with WiFi connectivity). This component integrates and processes the received information. It also runs an application with which the user accesses and manages the overall testbed.

The different components are evaluated by the corresponding agents, who also estimate the aggregate SPD value of the whole system. The agents inspect their underlying domain, managing it based on an SPD-aware reasoning operation. Furthermore, the system is re-configurable at runtime according to the SPD protection and performance goals defined in the activated policy. Affected agents configure their subsystem's settings to raise the SPD value when attacks are performed and then return to normal when the attacks are over (to save resources). Regarding the adaptation capabilities integrated within the proof-of-concept, the cryptographic service provides three communication states: Plain-text, authenticated, as well as authenticated encryption. Additionally, the trust scheme

supports two trust evaluation states: Direct trust only, as well as a combination of direct and indirect trust.

The system begins with a moderate SPD configuration to conserve resources (i.e., authenticated communication and direct trust). If SPD-Safe observes malicious activity, it informs the system entities to raise their protection level. The relevant response actions are specified in a security policy (applicable to the specific device type), such as applying authenticated and encrypted communication with combined direct and indirect trust information. The SPD value and status of each system component is then altered as a response to the launched attacks, so as to achieve a sufficient level of protection. The WSNs comply with the current policies, becoming stricter to misbehavior and isolating the compromised nodes. The main protection mechanism against cyber-attacks (i.e., blackhole or link-spoofing) is provided by SecRoute, while the smart camera enhances physical protection. In the same way, the system returns to the previous (initial) state when the triggering conditions are over.

For WSN<sub>1-1</sub>, we emulate scenarios where: (i) A node is malfunctioning due to low battery, and (ii) a compromised node launches a badmouth attack. In (i) the node is protected when a low energy level is observed by not including in traffic forwarding operations. The administrator gets notified accordingly. When the issue is fixed, the trust level is restored and the nodes' operational status returns back to normal. In case (ii), the compromised entity is detected when the attack rate reaches a threshold and it is blocked from routing operations. For WSN<sub>1-2</sub>, we launch blackhole and jamming attacks against congested or topology significant components. The secure routing mechanism successfully identifies both attacks and mitigates them. Table 5 presents in detail the above-mentioned scenario phases. The SPD levels are depicted with: (i) Red for values of 0–50—i.e., a situation where the provided protection is low, the proper functionality may not be available, and the operator must take immediately the related countermeasures; (ii) yellow for values of 51–70—i.e., moderate protection but still safe operation; and (iii) green for values of 71–100—i.e., high levels of protection.

**Table 5.** Scenario steps of the smart transportation use case.




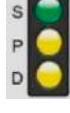






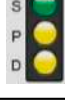
Event	Description	SPD State	Total <S, P, D> Value	SPD Visualization
1	Start of all components and services. Discovery/registration operations.	Initial State	<80, 70, 65>	
2	Bad-mouthing attack to WSN <sub>1-1</sub> . MA (master agent) is alerted for the attack and commands the rest agents to increase security.	Security level decreases	<60, 70, 65>	
3	Security status is enhanced on all SAs (simple agents). MA is notified.	Security level increases	<85, 70, 65>	
4	WSN <sub>1-1</sub> counters bad-mouthing and SA <sub>1-1</sub> informs the MA. The MA requests from the SAs to restore the normal state (to conserve resources).	Security level returns to initial state	<80, 70, 65>	
5	Blackhole attack to WSN <sub>1-2</sub> . MA is alerted for the attack and commands the rest agents to increase security.	Security level decreases	<50, 70, 65>	
6	Security status is enhanced on all SAs. MA is notified.	Security level increases	<85, 70, 65>	

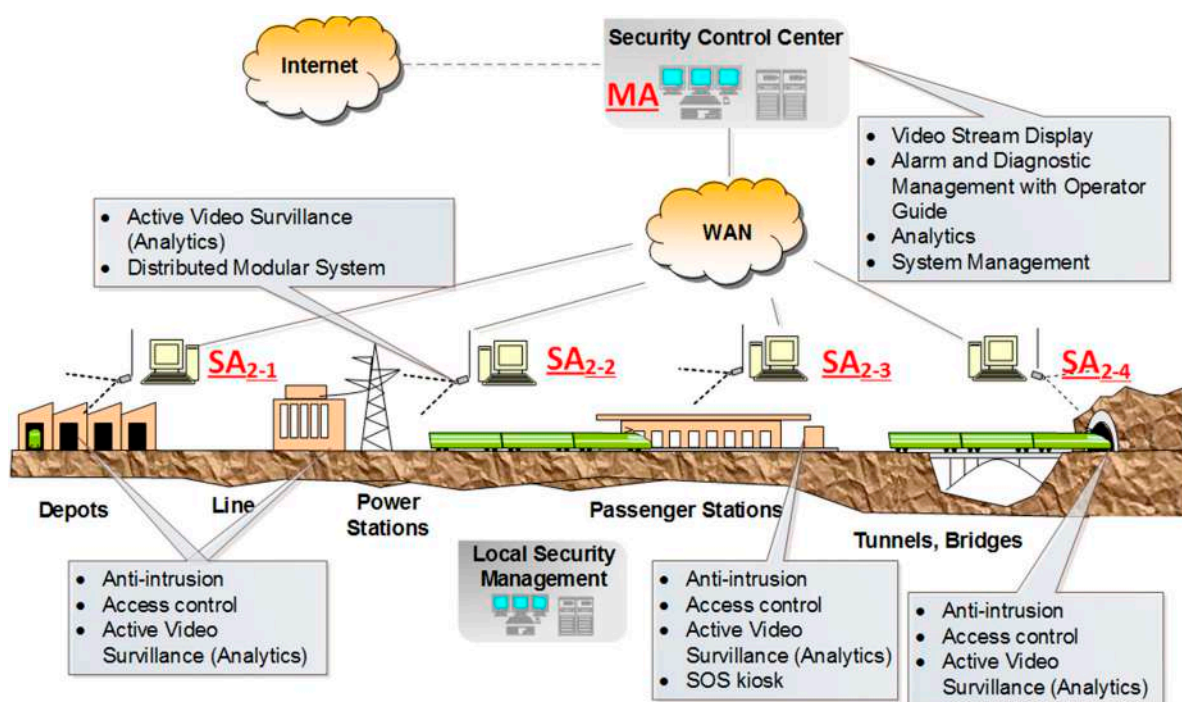


Table 5. Cont.

Event	Description	SPD State	Total <S, P, D> Value	SPD Visualization
7	WSN <sub>1-2</sub> counters the blackhole attack and SA <sub>1-2</sub> informs the MA. The MA requests from the SAs to restore the normal state (to conserve resources).	Security level returns to initial state	<80, 70, 65>	
8	A node has died in WSN <sub>1-2</sub> . MA is informed.	Dependability level decreases	<80, 70, 30>	
9	The dead node is replaced by the personnel. Dependability is restored. SA <sub>1-1</sub> reports the new status to MA.	Dependability level returns to initial state	<80, 70, 65>	
10	Simulated jamming attack against the network layer of WSN <sub>1-2</sub> . MA is informed.	S & D levels decrease	<40, 70, 40>	
11	The trust-based routing component counters the attack. SA <sub>1-2</sub> reports new state to MA.	S & D levels return to initial state	<80, 70, 65>	

### 5.3. Outdoor Setting—Safety Scenario

For outdoor on-route defense, a similar WSN with four BeagleBone nodes is installed. The nodes are connected with a mains power supply and control a smart camera as well as weather sensors. In the emulated use-case, the nodes and related SAs are deployed on: (i) The passenger's station; (ii) the track; (iii) the carriage departure, and; (iv) all bridges and tunnels along the track. Figure 6 illustrates the on-route WSN<sub>2</sub> [28] along with the central MA and underlying SA<sub>2-1</sub>–SA<sub>2-4</sub>.



**Figure 6.** The outdoor WSN for the on-route scenario. The MA is deployed in the security control center and the four SAs are installed in the edge system.



Through the responsible SA, the networking components (e.g., sensors and cameras) send real-time information to a security control center and the related master agent. Figure 7 depicts the graphical user interface and the visualization of the information that is collected by the on-route equipment, as developed by Ansaldo STS.

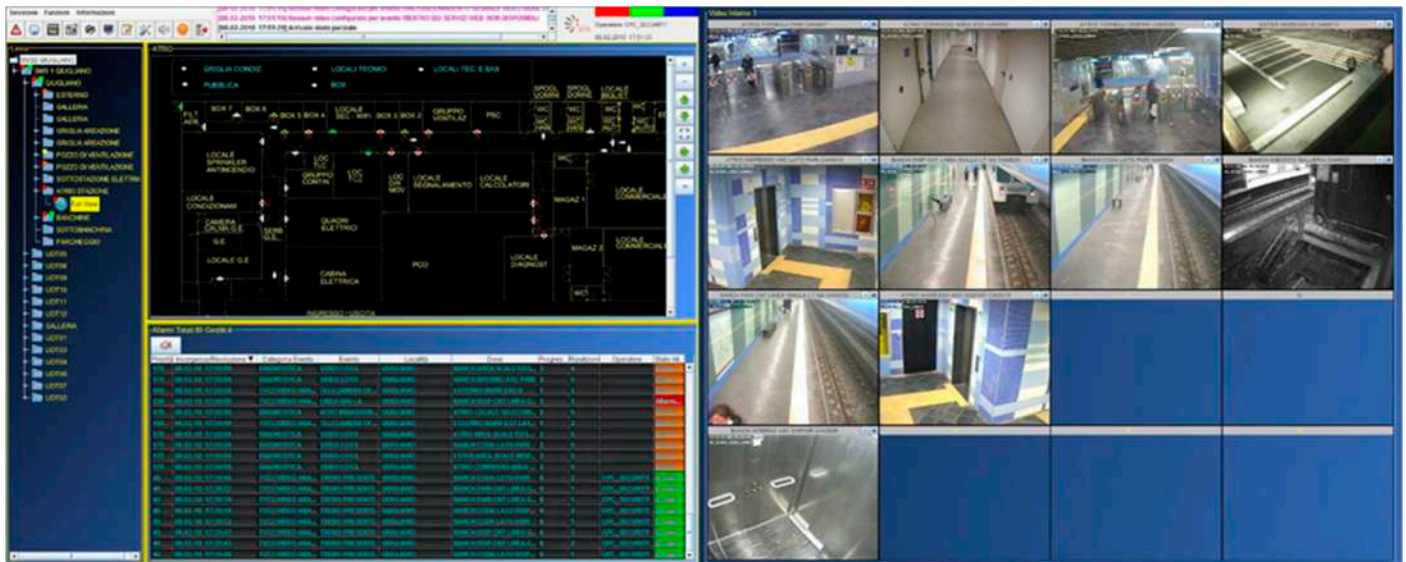


Figure 7. The railway on-route WSN graphical user interface.

In case of an emergency, the agents manage the system components to advise the personnel and assist the passengers. The demonstrated incident emulates the response strategy for a fire alarm, where decisions concerning both safety and security must be taken. In Appendix A, the code sample Figure A1 describes the CAP message that indicates the fire alarm.

Normally, for the indoor setting, the personnel and passengers are allowed to open doors based on their access rights (as determined by safety and security rules). When fire is detected by the sensors, an alarm is triggered, and the associated agent is notified. The agent takes the decision to degrade the security status by unlocking all doors, therefore enabling the unhindered evacuation of the train. Furthermore, via GSM, the agent automatically transmits an SMS to the responsible authorities concerning this incident (including situation's severity, GPS coordinates) and alerts the neighboring entities to be aware (e.g., agents on nearby trains). The train agents that cross the area are also notified to perform related actions (such as stop to the nearest station or change route). Moreover, it is assumed that during normal operation the smart cameras capture frames at a low rate to preserve bandwidth. When the alarm is raised, the setting is reconfigured at runtime, offering a high framerate and continuous monitoring of the affected area. As the fire is extinguished and the damaged components are restored, the normal status is restored. The code shown in Figure A2 summarizes the main processing flow and the emergency response rules that perform the described actions (for more information regarding EC, please refer to Mueller [73]).

## 6. Discussion

### 6.1. Comparison

This subsection compares SPD-Safe with the related works presented in Section 2.2 (i.e., TrainIntegrity, CMLRVT, and SENSORAIL) in terms of features. Table 6 summarizes the comparison results.

**Table 6.** Smart railway systems.

Property	SDP-Safe	TrainIntegrity	CMLRVT	SENSORAIL
AI technologies	JADE/Jess	NO	NO	SeNsIM
Reasoning & processing	Distributed	Centralized	Centralized	Centralized
Conflict resolution	YES	NO	NO	NO
Security management	YES	NO	NO	NO
Safety management	YES	YES	YES	YES
Middleware	OSGi	NO	NO	NO
Network layer protection	SecRoute	NO	NO	NO
Cloud management	ViMA	NO	NO	NO

All the related smart railway systems identified adopt semantic representation and reasoning. The service-oriented approaches conform to the specific application aspects and, therefore, in all relevant systems the agents are uniquely responsible for specific operations. The conflicting patterns are also not examined in most of these designs, limiting their applicability to specific deployments.

Furthermore, the three related systems do not use any management middleware for embedded devices. This approach is quite limiting in the IoT era, where high volumes of heterogeneous equipment have to be deployed and co-function. The systems also neglect the popular agent frameworks which, among others, provide efficient agent-related functionality and implement relevant standards. The reasoning operation is developed with general purpose programming languages, ignoring the advantages offered by the deductive rule-based techniques. Mechanisms for resolving conflicts, when implemented, are based either on epistemic or relational reasoning. More importantly, these related systems do not safeguard security, privacy, and dependability, and do not utilize any built-in protection technologies.

Conversely, SPD-Safe is a solution focusing on the SPD management of IoT and CPS settings. The SPD modeling is based on well-structured metrics that analyze the various configuration options of a multi-layered system. The AI process adjusts the railway CPS and counters attacks at runtime. SPD-Safe integrates state-of-the-art technological building blocks and platforms for the implementation of reasoning, as well as the management of devices and agents. Epistemic and relational reasoning are incorporated for resolving conflicts. Furthermore, the proposed framework adopts standardized technologies, from semantic standards to communication protocols and authorization schemes.

## 6.2. Future Work

SPD-Safe integrates several technologies in a secure manner. It preserves the SPD properties, enables active defenses and countermeasures, and can facilitate emergency response operations.

Active and offensive types of defenses are proposing nowadays, as the next step to enhance protection and mitigate threats, that the mainstream passive mechanisms (e.g., cryptography, network slicing, anti-viruses, etc.) cannot tackle. MTD is such an approach. It is becoming harder to analyze the system and exploit its vulnerabilities. Furthermore, in conjunction with other intrusion detection techniques, it can mitigate or even block some type of ongoing attacks. Nevertheless, more research is needed in order to make guidelines for the implementation of effective MTD policies as well as strategies to mitigate more advance attacks.

Moreover, safety-related events require the participation of relevant authorities. In modern settings, emergency authorities possess their own equipment, which is utilized during safety incidents. The cooperation of the involved systems becomes vital when it comes to rescuing lives. For the effectiveness of the response services, the systems must authenticate and authorize the various participants and exchange information (e.g., sensors' data,

surveillance video, etc.) in real time. The seamless interoperation will be examined in future extensions of SPD-Safe.

## 7. Conclusions

This paper introduced SPD-Safe, an administration framework for IoT settings in ambient secure and safety-critical domains, applied to protect a railway CPS. For secure connectivity, an innovative secure routing protocol was integrated in the network layer. The protocol covers all core security properties (confidentiality, integrity, and authentication) and features policy-based authorization. It was found to be energy efficient and could effectively counter a variety of attacks, providing defense against several threats that are not mitigated by existing solutions. For smart monitoring and automatic adaptation, smart agents were deployed at the edge systems and backend infrastructure, and performed the required AI processes. A multi-agent system was developed in the JADE platform and integrated on the OSGi middleware for the management of DPWS-enabled equipment, also utilizing various built-in protection mechanisms. The core reasoning process was implemented in Event Calculus. The SPD validation and metric-driven administration were modeled as a heuristic framework in a security-related theory. The implementation of MTDs was enabled, providing extra protection against attacks that were not mitigated by passive defenses. Furthermore, the system models a safety-related theory and implemented associated AI ambient strategies and plans. The two features were incorporated to administrate the underlying components, considering both the SPD and safety aspects. To validate the proposed approach, SPD-Safe was deployed to administrate WSNs on a complex railway CPS testbed, where the underlying components were successfully configured at runtime and mitigate security-related attacks, while AI reactive plans preserved the safety of personnel and passengers in emergency situations. The average delay in real equipment was around 0.2–0.6 s.

In terms of future work, advances in MTD solutions and integration with emergency response services were considered. MTDs are coming to the foreground nowadays and are expected to play a significant role in future defense strategies as AI becomes an integral part of new generation systems. Safety critical systems, such as the railway ones, must provide an adequate means to collaborate with emergency authorities and support their operations. Facilitating emergency response and a rapid restoration of service must also be considered by modern smart railway installations.

**Author Contributions:** Conceptualization, G.H. and K.F.; methodology, G.H. and S.I.; software, G.H.; validation, K.F., N.P., G.V., and I.H.; writing—original draft preparation, G.H.; writing—review and editing, K.F. and S.I.; Supervision and Project administration, S.I. and G.S. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work has received funding from the European Union Horizon’s 2020 research and innovation program under the grant agreements No. 786890 (THREAT-ARREST), No. 830927 (CONCORDIA), and No. 269317 (nSHIELD).

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** Data sharing not applicable.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## Appendix A

The code sample Figure A1 describes the CAP message that indicates the fire alarm.

```

<soap:Envelope xmlns:soap="http://schemas.xmlsoap.org/soap/envelope/">
  <soap:Body>
    <ns3:Notify xmlns="http://docs.oasis-open.org/wsrf/bf-2"
      xmlns:ns2="http://www.w3.org/2005/08/addressing"
      xmlns:ns3="http://docs.oasis-open.org/wsn/b-2"
      xmlns:ns4="http://protectrail.eu/model/events/resource"
      xmlns:ns5="urn:oasis:names:tc:emergency:cap:1.2"
      xmlns:ns6="http://docs.oasis-open.org/wsn/t-1">
      <ns3:NotificationMessage>
        <ns3:Topic Dialect="http://docs.oasis-open.org/wsn/t-1/TopicExpression/Full">FireDetection</ns3:Topic>
        <ns3:Message>
          <ns5:alert>
            <ns5:identifier>urn:rixf:com.tuc.SPD-Safe:id/FireAlert_01</ns5:identifier>
            <ns5:sender>WSN2</ns5:sender>
            <ns5:sent>2018-10-19T10:43:09.000+02:00</ns5:sent>
            <ns5:status>Actual</ns5:status>
            <ns5:msgType>Alert</ns5:msgType>
            <ns5:source>urn:rixf:com.tuc.fireprotection/devices/WSN</ns5:source>
            <ns5:scope>Public</ns5:scope>
            <ns5:info>
              <ns5:category>Fire</ns5:category>
              <ns5:event>FireDetection</ns5:event>
              <ns5:responseType>Evacuate</ns5:responseType>
              <ns5:urgency>Immediate</ns5:urgency>
              <ns5:severity>Extreme</ns5:severity>
              <ns5:certainty>Observed</ns5:certainty>
              <ns5:parameter>
                <ns5:valueName>Area</ns5:valueName>
                <ns5:value>51.468928,16.858863 0.01</ns5:value>
              </ns5:parameter>
            </ns5:info>
          </ns5:alert>
        </ns3:Message>
      </ns3:NotificationMessage>
    </ns3:Notify>
  </soap:Body>
</soap:Envelope>

```

**Figure A1.** Simple Object Access Protocol (SOAP) message that contains the CAP alert for the fire alarm.

The code shown in Figure A2 summarizes the main processing flow and the emergency response rules that perform the described actions (for more information regarding EC, please refer to Mueller [73]).

```

Rule 1: Happens(EventDitection(WSN,Fire),t) ∧ HoldsAt(ControllingAgent(SA,WSN),t)
=>
Happens(InformAgent(SA,WSN,Fire),t+1)

Rule 2: Happens(InformAgent(SA,WSN,Fire),t)
=>
Happens(InformTrainOperator(SA,SA_TrainOperator,Fire),t+1) ∧ Happens(UnlockDoors(SA, SA_Train),t+1) ∧
Happens(InformMA(SA,MA,Fire),t+1) ∧
Happens(InformNearbyAgents(SA,Fire),t+1)

Rule 3: Happens(InformNearbyAgents(SAi,Fire),t) ∧ HoldsAt(LocatedIn(SAi,area),t) ∧ HoldsAt(NearbyAgents(SAi, SAi),t)
=>
Happens(AvoidArea(SAi, SAi_TrainOperator, area),t+1)

Rule 4: Happens(InformMA(SA,MA,Fire),t)
=>
Happens(InformSystemOperator(MA,SystemOperator,Fire),t+1) ∧ Happens(InformAuthorities(MA,FireBrigate,Fire),t+1)

```

**Figure A2.** EC rules that perform the main safety reasoning behavior of each intelligent agent.

## References

- Xu, S.; Zhu, G.; Ai, B.; Zhong, Z. *A Survey on High-Speed Railway Communications: A Radio Resource Management Perspective*. *Computer Communications*; Elsevier: Amsterdam, The Netherlands, 2016; Volume 86, pp. 12–28.
- Chamoso, P.; González-Briones, A.; Rodríguez, S.; Corchado, J.M. Tendencies of Technologies and Platforms in Smart Cities: A State-of-the-Art Review. *Wirel. Commun. Mobile Comput.* **2018**, *2018*, 3086854.
- Boudi, Z.; El Koursi, E.M.; Ghazel, M. The New Challenges of Rail Security. *J. Traffic Logist. Eng.* **2016**, *4*, 56–60. [CrossRef]
- Kour, R.; Thaduri, A.; Karim, R. Railway Defender Kill Chain to Predict and Detect Cyber-Attacks. *J. Cyber Secur. Mobil.* **2019**, *9*, 47–90. [CrossRef]
- Luxton, A.; Marinov, M. Terrorist Threat Mitigation Strategies for the Railways. *Sustainability* **2020**, *12*, 3408. [CrossRef]
- Zhang, J.; Hu, F.; Wang, S.; Dai, Y.; Wang, Y. Structural vulnerability and intervention of high speed railway networks. *Phys. A Stat. Mech. Appl.* **2016**, *462*, 743–751. [CrossRef]
- González-Briones, A.; Garcia-Martin, R.; de AlbaJuan, F.L.; Corchado, M. Agent-Based Platform for Monitoring the Pressure Status of Fire Extinguishers in a Building. In *International Conference on Practical Applications of Agents and Multi-Agent Systems (PAAMS)*; Springer: Berlin/Heidelberg, Germany, 2020; Volume 1233, pp. 373–384.
- Catalano, A.; Bruno, F.A.; Galliano, C.; Pisco, M.; Persiano, G.V.; Cutolo, A.; Cusano, A. An optical fiber intrusion detection system for railway security. *Sens. Actuators A Phys.* **2017**, *253*, 91–100. [CrossRef]
- Fraga-Lamas, P.T.; Fernández-Caramés, M.; Castedo, L. Towards the Internet of Smart Trains: A Review on Industrial IoT-Connected Railways. *Sensors* **2017**, *17*, 1457. [CrossRef]
- Wang, Y.; Zhu, L.; Yu, Z.; Guo, B. An adaptive track segmentation algorithm for a railway intrusion detection system. *Sensor* **2019**, *19*, 2594. [CrossRef]
- Gai, K.; Qiu, M.; Hassan, H. Secure Cyber Incident Analytics Framework using Monte Carlo Simulations for Financial Cybersecurity Insurance in Cloud Computing. In *Concurrency and Computation: Practice and Experience*; Wiley: Hoboken, NJ, USA, 2017; Volume 29, issue 7.
- Chang, S.E.; Liu, A.Y.; Lin, S. Exploring privacy and trust for employee monitoring. *Ind. Manag. Data Syst.* **2015**, *115*, 88–106. [CrossRef]
- Paganini, P. Modern Railroad Systems Vulnerable to Cyber Attacks. Security Affairs. 2016. Available online: <http://securityaffairs.co/wordpress/43196/hacking/railroad-systems-vulnerabilities.html> (accessed on 18 November 2020).
- Bababeik, M.; Khademi, N.; Chen, A.; Nasiri, M.M. Vulnerability analysis of railway networks in case of multi-link blockage. *Transp. Res. Procedia* **2017**, *22*, 275–284. [CrossRef]
- Khanmohamadi, M.; Bagheri, M.; Khademi, N.; Ghannadpour, S.F. A security vulnerability analysis model for dangerous goods transportation by rail—Case study: Chlorine transportation in Texas-Illinois. *Saf. Sci.* **2018**, *110*, 230–241. [CrossRef]
- Salmane, H.; Khoudour, L.; Ruichek, Y. A Video-Analysis-Based Railway–Road Safety System for Detecting Hazard Situations at Level Crossings. *IEEE Trans. Intell. Transp. Syst.* **2015**, *16*, 596–609. [CrossRef]
- Chernov, A.V.; Savvas, I.K.; Butakova, M.A. Detection of Point Anomalies in Railway Intelligent Control System Using Fast Clustering Techniques. In *Proceedings of the 3rd International Scientific Conference Intelligent Information Technologies for Industry*, Sochi, Russia, 17–21 September 2018; pp. 267–276.
- Coppola, P.; Silvestri, F. Assessing travelers’ safety and security perception in railway stations. *Case Stud. Transp. Policy* **2020**, *8*, 1127–1136. [CrossRef]
- Mrazovic, P.; Eser, E.; Ferhatosmanoglu, H.; Larriba-Pey, J.L.; Matskin, M. Multi-vehicle Route Planning for Efficient Urban Freight Transport. In *Proceedings of the 2018 International Conference on Intelligent Systems (IS)*, Funchal, Madeira, Portugal, 25–27 September 2018; pp. 744–753.
- Zhu, C.; Shu, L.; Leung, V.C.M.; Guo, S.; Zhang, Y.; Yang, L.T. Secure multimedia Big Data in trust-assisted sensor-cloud for smart city. *IEEE Commun. Mag.* **2017**, *55*, 24–30. [CrossRef]
- Chamoso, P.; De La Prieta, F.; De Paz, F.; Corchado, J.M. Swarm Agent-Based Architecture Suitable for Internet of Things and Smartcities. In *Distributed Computing and Artificial Intelligence, 12th International Conference*; Springer: Berlin/Heidelberg, Germany, 2015; Volume 373, pp. 21–29.
- Zhang, Q.; Chen, Z.; Leng, Y. Distributed fuzzy c-means algorithms for big sensor data based on cloud computing. *Int. J. Sens. Networks* **2015**, *18*, 32–39. [CrossRef]
- Tsaramirsis, G.; Karamitsos, I.; Apostolopoulos, C. Smart Parking: An IoT application for Smart City. In *Proceedings of the 10th INDIACom-2016 International Conference*, New Delhi, India, 16–18 March 2016; pp. 2271–2275.
- Yin, W.; He, S.; Zhang, Y.; Hou, J. A Product-Focused, Cloud-Based Approach to Door-to-Door Railway Freight Design. *IEEE Access* **2018**, *6*, 20822–20836. [CrossRef]
- Dong, Q.; Hayashi, K.; Kaneko, M. An Optimized Link Layer Design for Communication-Based Train Control Systems Using WLAN. *IEEE Access* **2017**, *6*, 6865–6877. [CrossRef]
- Fanian, F.; Rafsanjani, M.K. Cluster-based routing protocols in wireless sensor networks: A survey based on methodology. *J. Netw. Comput. Appl.* **2019**, *142*, 111–142. [CrossRef]
- Khanna, N.; Sachdeva, M. Study of trust-based mechanism and its component model in MANET: Current research state, issues, and future recommendation. *Int. J. Commun. Syst.* **2019**, *32*, 1–23. [CrossRef]
- Cesena, M. SHIELD Technology Demonstrators. In *Measurable and Composability Security, Privacy, and Dependability for Cyberphysical Systems*; CRC Press: Boca Raton, FL, USA, 2017; pp. 381–434.

29. Brokalakis, A.; Tampouratzis, N.; Nikitakis, A.; Andrianakis, S.; Papaefstathiou, I.; Dollas, A. An Open-Source Extendable, Highly-Accurate and Security Aware CPS Simulator. In Proceedings of the 2017 13th International Conference on Distributed Computing in Sensor Systems (DCOSS), Ottawa, ON, Canada, 5–7 June 2017; pp. 81–88.
30. Farooq, J.; Soler, J. Radio Communication for Communications-Based Train Control (CBTC): A Tutorial and Survey. *IEEE Commun. Surv. Tutor.* **2017**, *19*, 1377–1402. [CrossRef]
31. Sun, W.; Yu, F.R.; Tang, T.; Bu, B. Energy-Efficient Communication-Based Train Control Systems with Packet Delay and Loss. *IEEE Trans. Intell. Transp. Syst.* **2016**, *17*, 452–468. [CrossRef]
32. Garcia-Loygorri, J.M.; Val, I.; Arriola, A.; Briso-Rodriguez, C. 2.6 GHz Intra-Consist Channel Model for Train Control and Management Systems. *IEEE Access* **2017**, *5*, 23052–23059. [CrossRef]
33. Fotso, S.J.T.; Frappier, M.; Laleau, R.; Mammari, A. Modeling the Hybrid ERTMS/ETCS Level 3 Standard Using a Formal Requirements Engineering Approach. In *International Conference on Abstract State Machines*; Alloy, B., Ed.; Springer: Berlin/Heidelberg, Germany, 2018; pp. 262–276.
34. Chetty, K.; Chen, Q.; Woodbridge, K. Train monitoring using GSM-R based passive radar. In Proceedings of the 2016 IEEE Radar Conference (RadarConf), Philadelphia, PA, USA, 1–6 May 2016; pp. 1–4.
35. Bates, R.J. GPRS: General Packet Radio Service. In *Book GPRS: General Packet Radio Service*; McGraw-Hill, Professional Telecom: New York, NY, USA, 2001.
36. Proto, M.; Bavusi, M.; Bernini, R.; Bigagli, L.; Bost, M.; Bourquin, F.; Cottineau, L.-M.; Cuomo, V.; Della Vecchia, P.; Dolce, M.; et al. Transport Infrastructure Surveillance and Monitoring by Electromagnetic Sensing: The ISTIMES Project. *Sensors* **2010**, *10*, 10620–10639. [CrossRef] [PubMed]
37. Crinière, A.; Dumoulin, J.; Mevel, L.; Andrade-Barroso, G. Cloud2IR an Infrared and Environmental SHM Information System. In Proceedings of the 13th Quantitative Infrared Thermography Conference (QIRT), Gdansk, Poland, 4–8 July 2016; pp. 226–235.
38. Xie, J.; Liu, C.-C. Multi-agent systems and their applications. *J. Int. Counc. Electr. Eng.* **2017**, *7*, 188–197. [CrossRef]
39. De la Prieta, F.; Rodríguez-González, S.; Chamoso, P.; Corchado, J.M.; Bajo, J. Survey of agent-based cloud computing applications. *Future Gener. Comput. Syst.* **2019**, *100*, 223–236. [CrossRef]
40. Kravari, K.; Bassiliades, N. A Survey of Agent Platforms. *J. Artif. Soc. Soc. Simul.* **2015**, *18*, 11–29. [CrossRef]
41. FIPA, “FIPA ACL Message Structure Specification,” Foundation for Intelligent Physical Agents. 2002. Available online: <http://www.fipa.org/specs/fipa00061/SC00061G.html> (accessed on 18 November 2020).
42. Fysarakis, K.; Askoxylakis, I.; Soultatos, O.; Papaefstathiou, I.; Manifavas, C.; Katos, V. Which IoT Protocol? Comparing Standardized Approaches over a Common M2M Application. In Proceedings of the 2016 IEEE Global Communications Conference (GLOBECOM), Washington, DC, USA, 4–8 December 2016; pp. 1–6.
43. OASIS. “Devices Profile for Web Services Version 1.1,” Organization for the Advancement of Structured Information Standards. 2009. Available online: <http://docs.oasis-open.org/ws-dd/dpws/1.1/os/wsdd-dpws-1.1-spec-os.pdf> (accessed on 18 November 2020).
44. Thirumalainambi, R. Pitfalls of Jess for dynamic systems. In Proceedings of the International Conference on Artificial Intelligence and Pattern Recognition (AIPR), Orlando, FL, USA, 9–12 July 2007; Volume 1, pp. 491–494.
45. Kumar, S.; Prasad, R. Importance of expert system shell in development of expert system. *Int. J. Innov. Res. Dev.* **2015**, *4*, 128–133.
46. Semmel, G.; Davis, S.; Leucht, K.; Rowe, D.; Kelly, A.; Boloni, L. Launch commit criteria monitoring agent. In Proceedings of the 4th International Joint Conference on Autonomous Agents and MultiAgent Systems (AAMAS), Utrecht, The Netherlands, 25–29 July 2005; pp. 3–10.
47. Goseva-Popstojanova, K.; Tyo, J. Experience Report: Security Vulnerability Profiles of Mission Critical Software: Empirical Analysis of Security Related Bug Reports. In Proceedings of the 28th International Symposium on Software Reliability Engineering (ISSRE), Toulouse, France, 23–26 October 2017; pp. 152–163.
48. Leitao, P.; Karnouskos, S. *Industrial Agents: Emerging Applications of Software Agents in Industry*, 1st ed.; Elsevier Science: Amsterdam, The Netherlands, 2015; pp. 1–476.
49. Ghadimi, P.; Wang, C.; Lim, M.K.; Heavey, C. Intelligent sustainable supplier selection using multi-agent technology: Theory and application for Industry 4.0 supply chains. *Comput. Ind. Eng.* **2019**, *127*, 588–600. [CrossRef]
50. Scholten, H.; Westenberg, R.; Schoemaker, M. Sensing Train Integrity. In Proceedings of the IEEE Sensors Conference, Christchurch, New Zealand, 25–28 October 2009.
51. Firlik, B.; Chudzikiewicz, A. Condition monitoring of a light rail vehicle—From concept to implementation. *Key Eng. Mater.* **2012**, *518*, 66–75. [CrossRef]
52. Flammini, F.; Gaglione, A.; Ottello, F.; Pappalardo, A.; Pragliola, C.; Tedesco, A. Towards wireless sensor networks for railway infrastructure monitoring. In Proceedings of the Electrical Systems for Aircraft, Railway and Ship Propulsion (ESARS), Bologna, Italy, 19–21 October 2010.


53. Casola, V.; Gaglione, A.; Mazzeo, A. A reference architecture for sensor networks integration and management. In Proceedings of the 3rd International Conference on Geosensor Networks, Oxford, UK, 13–14 July 2009; pp. 158–168.
54. Flammini, F.; Gaglione, A.; Mazzocca, N.; Pragliola, C. DETECT: A novel framework for the detection of attacks to critical infrastructures. In *Safety, Reliability and Risk Analysis: Theory, Methods and Applications*; Taylor & Francis: Abingdon, UK, 2008; pp. 105–112.
55. Chakravarthy, S.; Mishra, D. Snoop: An expressive event specification language for active databases. *Data Knowl. Eng.* **1994**, *14*, 1–26. [CrossRef]
56. Ganeriwal, S.; Balzano, L.; Srivastava, M. Reputation-based framework for high integrity sensor networks. *ACM Trans. Sen. Netw.* **2008**, *4*. [CrossRef]
57. Hu, Y.-C.; Perrig, A.; Johnson, D.B. Ariadne: A secure on-demand routing protocol for ad hoc networks. *Wirel. Netw.* **2005**, *11*, 21–38. [CrossRef]
58. Zhang, Y.; Xu, L.; Wang, X. A Cooperative Secure Routing Protocol based on Reputation System for Ad Hoc Networks. *J. Commun.* **2008**, *3*, 43–50. [CrossRef]
59. Altisen, K.; Devismes, S.; Jamet, R.; Lafourcade, P. SR3: Secure resilient reputation-based routing. In Proceedings of the 2013 IEEE International Conference on Distributed Computing in Sensor Systems (DCOSS), Cambridge, MA, USA, 20–23 May 2013; pp. 258–265.
60. Dhaheri, A.A.; Yeum, C.Y.; Damiani, E. New Two-Level  $\mu$ TESLA Protocol for IoT Environments. In Proceedings of the 2019 IEEE World Congress on Services (SERVICES), Milan, Italy, 8–13 July 2019; pp. 84–91.
61. Hatzivasilis, G.; Papaefstathiou, I.; Askoxylakis, I.; Fysarakis, K. SecRoute: End-to-end secure communications for wireless ad-hoc networks. In Proceedings of the 22nd IEEE Symposium on Computers and Communications (ISCC), Heraklion, Crete, Greece, 3–6 July 2017; pp. 558–563.
62. Hatzivasilis, G.; Papaefstathiou, I.; Manifavas, C. SCOTRES: Secure Routing for IoT and CPS. *IEEE Internet Things J.* **2017**, *4*, 2129–2141. [CrossRef]
63. Hatzivasilis, G.; Papaefstathiou, I.; Plexousakis, D.; Manifavas, C.; Papadakis, N. AmbISPDm: Managing embedded systems in ambient environment and disaster mitigation planning. In *Applied Intelligence*; Springer: Berlin/Heidelberg, Germany, 2017; pp. 1–21.
64. Java Agent Development (JADE) Framework. Available online: <http://jade.tilab.com/> (accessed on 18 November 2020).
65. Tilab, S.P.A. JADE Security Add-On Guide. 2005. Available online: [http://jade.tilab.com/doc/tutorials/JADE\\_Security.pdf](http://jade.tilab.com/doc/tutorials/JADE_Security.pdf) (accessed on 18 November 2020).
66. Ali, B.; Manzoor, U.; Zafar, B. eJADE-S: Encrypted JADE-S for Securing Multi-Agent Applications. In Proceedings of the International Conference on Artificial Intelligence (ICAI), Athens, Greece, 27–30 July 2015; pp. 548–554.
67. Open Services Gateway Initiative (OSGi). Available online: <http://www.osgi.org/> (accessed on 18 November 2020).
68. OASIS. Common Alerting Protocol Version 1.2, Organization for the Advancement of Structured Information Standards. 2010. Available online: <http://docs.oasis-open.org/emergency/cap/v1.2/CAP-v1.2-os.pdf> (accessed on 18 November 2020).
69. Hatzivasilis, G.; Papadakis, N.; Hatzakis, I.; Ioannidis, S.; Vardakis, G. AI-driven composition and security validation of an IoT ecosystem. *Appl. Sci.* **2020**, *10*, 4862. [CrossRef]
70. Friedman-Hill, E.J. Jess: The Rule Engine for Java Platform. Sandia National Laboratories. 2008. Available online: <http://www.jessrules.com/docs/71/> (accessed on 18 November 2020).
71. Lu, Y.; Li, Q.; Zhou, Z.; Deng, Y. Ontology-based knowledge modeling for automated construction safety checking. *Saf. Sci.* **2015**, *79*, 11–18. [CrossRef]
72. Patkos, T.; Plexousakis, D.; Chibani, A.; Amirat, Y. An event calculus production rule system for reasoning in dynamic and uncertain domains. In *Theory and Practice of Logic Programming*; Cambridge University Press: Cambridge, UK, 2016; Volume 16, pp. 325–352.
73. Mueller, E.T. *Commonsense Reasoning*, 2nd ed.; Kaufmann, M., Ed.; Elsevier: Amsterdam, The Netherlands, 2015; pp. 1–516.
74. Lei, C.; Zhang, H.-Q.; Tan, J.-L.; Zhang, Y.-C.; Liu, X. Moving Target Defense Techniques: A Survey. *Secur. Commun. Netw.* **2018**, *2018*, 1–25. [CrossRef]
75. Berstel, B. Extending the RETE algorithm for event management. In Proceedings of the 9th International Symposium on Temporal Representation and Reasoning, Manchester, UK, 7–9 July 2002; pp. 49–51.
76. Eby, M.; Werner, J.; Karsai, G.; Ledeczi, A. Integrating security modeling into embedded system design. In Proceedings of the 14th Annual IEEE International Conference and Workshops on the Engineering of Computer-Based Systems (ECBS), Tucson, AZ, USA, 26–29 March 2007; pp. 221–228.
77. Kelly, S.; Tolvanen, J.-P. *Domain-Specific Modeling: Enabling Full Code Generation*; Wiley-IEEE Computer Society Pr.: Hoboken, NJ, USA, 2008; pp. 1–444.
78. Ko, H.; Jin, J.; Keoh, S.L. Secure Service Virtualization in IoT by Dynamic Service Dependency Verification. *IEEE Internet Things J.* **2016**, *3*, 1006–1014. [CrossRef]
79. Albanese, M.; Battista, E.; Jajodia, S.; Casola, V. Manipulating the Attacker’s View of a System’s Attack Surface. In Proceedings of the IEEE Conference on Communications and Network Security, San Francisco, CA, USA, 29–31 October 2014; pp. 472–480.

80. Savola, R.M.; Sihvonen, M. Metrics driven security management framework for e-health digital ecosystem focusing on chronic diseases. In Proceedings of the MEDES '12: International Conference on Management of Emergent Digital EcoSystems, Addis Ababa, Ethiopia, 28–31 October 2012; pp. 75–79. [CrossRef]
81. Ayyappan, B.; Kumar, P.M. Security protocols in WSN: A survey. In Proceedings of the 2017 Third International Conference on Science Technology Engineering & Management (ICONSTEM), Chennai, India, 23–24 March 2017; pp. 301–304.
82. Parducci, B.; Lockhart, H. *eXtensible Access Control Markup Language (XACML) Version 3.0*; OASIS Standard: Burlington, MA, USA, 2013; pp. 1–154.



## Article

# Exploratory Data Analysis and Data Envelopment Analysis of Urban Rail Transit

Guillermo L. Taboada <sup>1,\*</sup>  and Liangxiu Han <sup>2</sup><sup>1</sup> Computer Architecture Group, CITIC, Universidade da Coruña, Campus de Elviña, 15071 A Coruña, Spain<sup>2</sup> Department of Computing and Mathematics, Manchester Metropolitan University, Manchester M1 5GD, UK; l.han@mmu.ac.uk

\* Correspondence: guillermo.lopez.taboada@udc.es

Received: 6 July 2020; Accepted: 5 August 2020; Published: 7 August 2020

**Abstract:** This paper deals with the efficiency and sustainability of urban rail transit (URT) using exploratory data analytics (EDA) and data envelopment analysis (DEA). The first stage of the proposed methodology is EDA with already available indicators (e.g., the number of stations and passengers), and suggested indicators (e.g., weekly frequencies, link occupancy rates, and CO<sub>2</sub> footprint per journey) to directly characterize the efficiency and sustainability of this transport mode. The second stage is to assess the efficiency of URT with two original models, based on a thorough selection of input and output variables, which is one of the key contributions of EDA to this methodology. The first model compares URT against other urban transport modes, applicable to route personalization, and the second scores the efficiency of URT lines. The main outcome of this paper is the proposed methodology, which has been experimentally validated using open data from the Transport for London (TfL) URT network and additional sources.

**Keywords:** urban rail transit (URT); exploratory data analysis (EDA); data envelopment analysis (DEA); sustainable transport systems; intelligent transportation systems (ITS); big-data applications

## 1. Introduction

Rail is one of the most energy-efficient transport modes [1], accounting for approx. 8% of global freight and motorized passenger movements but only 2% of transport energy use, being the transport mode with highest percentage of electric penetration. Thus, the continuous decarbonization of power production will allow zero-emission rail transport in the medium term. This is especially relevant for urban environments, where fuel-based transport modes impact the most on people's health. For these reasons, urban rail transit (URT) plays a key role in a context of a significant rise of urban population, particularly in emerging economies, which increases pollution, congestion, and city-center traffic restrictions.

URT is ideally suited for high passenger throughput, and although investment is especially high per kilometer, costs per throughput capacity are lower than for urban road infrastructure [2]. Shifting passengers from private cars to public transport, particularly in large cities, is key to reducing net energy use and emissions to be able to meet the mobility challenges within the sustainable development goals (SDG) [3].

Nearly 200 cities worldwide have metro systems (URT with the highest capacity), whose length exceed 32,000 km, whereas around 400 cities have light rail systems (URT with less investment requirements, less speed, and more modest capacity). Most recent (in the 2010s decade) URT developments have been, in the case of metro, which requires the highest investments, in Asia (34 of 46 new cities with metro). In the case of light rail, 28 new projects have been developed in

Europe and another 37 new projects roughly equally among Asia, North America, the Middle East, and North Africa.

URT, on the one hand, has multiple benefits such as mitigating CO<sub>2</sub> emissions and local air pollution, and wider social and economic benefits. One of these is reducing commuting time, therefore expanding the urban/suburban areas to those directly communicated by URT. This way, labor force can live at a higher distance from the city center in new urban developments, which are more cost-effective (less expensive).

On the other hand, one of the limitations of URT networks, the high infrastructure investment and its capillarity degree, could be partially addressed through multi-modal communication, using URT for the main journey combined with another communication mode, typically walking for distances up to 1 km, cycling, bus, or private car from a park-and-ride facility. However, there is still room for improvement as living more than 1 km away from an URT station requires either a frequent bus service network or private car ownership and confronting parking costs.

The objective of this paper is to characterize the efficiency and sustainability of URT using a proposed methodology based on exploratory data analysis (EDA) [4] and data envelopment analysis (DEA) [5,6]. First, the available data is explored using EDA to identify the main factors influencing URT for, secondly, derive efficiency scores, both across different transport modes and different lines within a rail network.

The reminder of the paper is structured as follows: the current state of the research field is presented in Section 2 and our methodology for assessing the efficiency and sustainability of URT using EDA and DEA is explained in Section 3. Section 4 experimentally validates the approach using open data from Transport for London (TfL) URT, particularly its underground network data. Section 5 discusses how additional big-data sources can improve the efficiency and sustainability of URT. In Section 5, we draw the main conclusions, stressing the main contributions of this work.

## 2. State of the Art

The recent increase in quantity and quality of data in public transport systems has fueled the adoption of data-driven solutions, mainly based on EDA or artificial intelligence (AI)/machine learning, to make public transport systems more intelligent, green, and safe. However, research in this area is, compared to applications in roads and private services/vehicles, more challenging due to the scarcity of available data, and the difficulties in testing research hypotheses in the real world. Thus, it is relatively common to have projects such as TEMA Big-Data Platform [7], which monitored 28,000 fuel vehicles with on-board GPS in Modena and Firenze (Italy) for one month to obtain 4.5 million trips and parking events, whereas most public transport analysis are based on significantly lower data records.

In an earlier study [8] a method and a software has been developed to estimate an URT passenger origin–destination trip matrix using an automatic data collection system. The method was experimentally assessed with automatic fare collection (AFC) origin-only data from Chicago Transport Authority (CTA), inferring the destination to replace the manual and costly origin–destination surveys. TfL, although it has a different fare collection scheme, an origin–destination system vs. the origin-only of CTA, collaborated in this research. In [9] smart-card bus data (metro was not included), travel surveys, and passengers' addresses have been used to measure commuting efficiency in Beijing in 2008–2010 as a function of commuting time, and residence/work location. In [10] public transport users' behaviors have been explored, whereas [11] analyzes individual mobility choices in carpooling. In [12] multi-modal transportation systems are presented as a way to increase efficiency through economies of scale, claiming that a multi-modal system combining a fast and efficient URT with other mobility options can provide more potential gains than optimizing single modal transport systems. These early works are EDA or data mining single cases studies using only a few data sources.

As the available data has been increasing during the last decade, particularly thanks to the integration of sensors in intelligent transportation systems (ITS) [13], especially in roads and connected vehicles, the number of projects started growing exponentially. A recent review paper [14] presents

almost a hundred EDA, data mining, AI, and machine-learning applications, challenges, and limitations, particularly for management, traffic safety, public transportation, and urban mobility. However, when it comes to public transport only tackles route planning, aviation, on-demand bus, and shared mobility, but no references to URT. Another reference paper [15] covers big-data projects and technologies in transportation and mobility, highlighting the scarcity of references in maritime and rail transport systems, with only a few works on predictive maintenance, risk management, and railway accidents.

URT, especially underground, faces important capacity limitations, especially in city centers at peak times. This has been the focus of [16] for forecasting passenger flows using Artificial Neural Networks (ANN) on a single metro line in Naples with a simulated dataset. Short-term forecasting on urban metros has also been studied along with other methods, such as Kalman filter in [17] and ARIMA (autoregressive integrated moving average) models in [18]. Li et al. [19] proposed a Multi-Scale Radial Basis Function (MSRBF) for forecasting short-term metro passenger flows on special occasions, such as sporting events and concerts. In this case, passenger flow is very irregular, and predictions are more difficult to obtain. Ling et al. [20] used smart-card data for predicting passenger flows in the subway of Shenzhen (China); they analyzed four predictive models: a historical average model, ANN, regression model, and a gradient-boosted regression tree model. Liu et al. [21] proposed a deep learning method for short-term forecasting of metro inbound/outbound passenger flows, while Wang et al. [22] proposed a Novel Markov-Grey model for solving the same problem.

A novel model of Multi-scale Mixture Feedback Wavelet Neural Network (MMFWNN) has been proposed in [23] to predict the short-term entrance flow of Shanghai subway stations, distinguishing passengers into commuter (more predictable) and non-commuter (more dependent on the weather). In [24] the factors affecting Seoul Metro boarding have been analyzed using regression analyses against the station environment (density, employment, commercial/office area), external connectivity (through metro and roads) and intermodal (bus and metro). This, and previous models, can predict highly accurately the short-term entrance flow, as it corresponds with regular patterns. However, the lack of historical data limits the behavior in special situations/events.

Relevant research works on other transportation modes are [25], where four classification algorithms have been used to model the relationship in London between weather and short cycling journeys using docked bikes. In addition, [26], proposing the application of deep learning methods to a Bus Rapid Transit (BRT) system (Xiamen, China) to forecast the hourly flow, adopting a three-stage architecture. This paper also analyses the literature, identifying four different approaches: (1) traditional classical algorithms; (2) regressive models; (3) machine-learning-based models, including ANNs; (4) hybrid models. In [27] a novel Context Neural Network framework has been proposed for the prediction of road traffic flow showing better long-term predictions than previous well-established models. All studied cases, however, refer to short-term or long-term time periods, without considering the spatial dimension.

The relevant number of references of the application of DEA to different public transport modes [28], contrasts with the scarcity of research on the efficiency of rail networks, [29], especially URT. These references in rail transport systems generally compare different public transport agencies at regional or national levels. In [30] 17 European URT networks have been evaluated using a two-stage methodology focusing on the relationship between the operational performance and their socioeconomic contexts. In [31] urban public transport systems of 652 Chinese cities have been analyzed, highlighting the high efficiency of URT. In [32] DEA has been used to assess the efficiency of 31 railway companies across multiple countries. In [33] the efficiency of 20 representative URT systems, among them London, Hong Kong, and New York, have been analyzed concluding that the higher the number of stations, the higher the efficiency. This conclusion is also supported by a recent study [34] on Chinese URT.

In [35] DEA has been used to assess the performance of the bus lines of a single transport authority in a suburban area in California Central Coast. In [36] the efficiency of Seoul Arterial Bus Route has been analyzed using DEA considering a wide variety of factors, including total rides, service satisfaction, and CO<sub>2</sub> emissions. This latter work was expanded in [37] with a network DEA model,

also validated with bus companies in Seoul. In addition, finally, in [38], DEA has been used to compare different transport options and investments on a single route.

The selection of input and output variables in DEA is regarded as an important step that is normally conducted before the DEA model is implemented. Available techniques are, on the one hand, based on expert intervention, using heuristic decision-making, and expert judgement (e.g., using Delphi), and, on the other hand, fully automatic approaches [39] which in turn maximize efficiencies and lose discrimination power without a full understanding of the domain. There is a lack of data-based methodologies and use cases that avoid bias of experts and at the same time provide useful, repeatable, and interpretable results. The proposed methodology in this paper, using EDA for a thorough selection of a limited number of variables, addresses this need by combining both approaches.

A review of the related literature on efficiency analysis in urban public transport [28] shows a quite homogenous selection of input and output variables, guided by experts, with a fairly narrow perspective. Thus, state-of-the-art variables are (in parenthesis the percentage of the papers in the literature that reported each variable):

- Input (physical measure): number of vehicles (61%), number of employees (40%), fuel consumption (36%), worked hours (11%), drivers (3%), non-driving employees (2%), seat capacity (total seats of the fleet) (2%), and number of depots (1%).
- Input (CAPEX): Price of capital (21%), and investment (4%).
- Input (OPEX): Price of labor (39%), price of fuel (29%), OPEX (12%), material costs (11%), fuel costs (9%), operating cost of vehicles (11%), maintenance costs (6%), and operating labor expenses (4%).
- Output (Service supply): overall traveled kilometers by all vehicles (53%), seats offered multiplied by overall traveled kilometers (26%), vehicles multiplied by hours of operation (4%), and revenue-vehicle kilometer (3%).
- Output (Service consumption): passengers multiplied by traveled kilometers (24%), number of passengers (17%), and number of trips (6%).
- Output (revenue): operating revenues (11%).

Furthermore, in the literature there are additional variables, neither considered inputs nor outputs, but sometimes considered external variables, which characterize public transport systems. Representative examples of these variables are (listed together with their presence, in percentage, in the analyzed related literature):

- Quality and characteristics of service: length of the network (24%), average commercial speed (21%), average fleet age (15%), service frequency (7%), and number of stops (6%).
- Socio/demographic: population density (16%), location (13%), car ownership (7%), population served (7%), and area served (5%).
- Managerial: public company (15%), contract type of the operator (11%), size of the company (3%).
- Subsidies: subsidies from public funds (10%), subsidies to operating expenses (8%), local subsidy (2%).
- Externalities: number of accidents (7%), and emissions (4%).

With regards to URT, the variables used in the related literature are:

- In [30] the network length, the number of stations and cars are the inputs (CAPEX), whereas the number of employees is considered the only input (OPEX), due to the scarcity of materials and energy consumption information, two relevant inputs (OPEX). Additional variables considered to be inputs are ratios between these variables (e.g., the network length divided by the number of cars), historical data, as well as socioeconomic variables, such as area, population density of the core city, average household size, unemployment rate, GDP (Gross Domestic Product) per capita, and diesel pump price. In [30] two models are computed: (i) efficiency, using the number of cars-kilometers produced as output, and (ii) effectiveness, considering the number of transported

passengers. The large number of variables and the limited number of analyzed URT networks (17) ends up with most of the evaluated systems considered highly efficient (here most URT networks excel in some, disjoint parameters, increasing its efficiency). The impact (elasticity) of the variables has also been considered, but the work fails in selecting the most representative ones.

- In [32] six inputs have been considered, the annual cost of operation as input (OPEX), and the network length, and the number of employees, traction vehicles, passenger cars, and cargo cars as inputs (CAPEX). Additionally, five outputs have been defined, revenues earned, transported passengers, transported passengers per kilometer, transported cargo tons, and transported cargo tons per kilometer.
- In [33] the number of employees and the labor costs are the selected inputs (OPEX), whereas the number of cars in operation and non-labor costs are used as inputs (both OPEX and CAPEX). The selected outputs are car-kilometers and transported passengers. Historical data has also been considered. Furthermore, additional variables have been used in the Tobit models phase, after DEA, such as population density, the number of stations, distance between stations, geographic location, and the type of URT (light/rapid or heavy).

So far, the use of input and output variables in DEA URT models relies on a wide range of state-of-the-art variables from the related literature, generally with limited selection and statistical analysis. Moreover, the access to these variables incurs relevant collection costs, such as accessing to unstructured reports, limiting the viability of comparing additional URTs.

This paper overcomes these latter limitations through:

- Selecting a limited number of representative variables through EDA, both state-of-the-art and new variables, increasing the discrimination power of DEA by bringing forward the statistical and visual analysis, prior to the variable selection (previous works [40] only suggested EDA after DEA, to understand the impact of variables on the models, so using EDA as first stage is one of the key contributions of this work).
- Automating data collection from public sources (e.g., open-data and online services), thus supporting the direct comparison across different URT systems.
- Comparing, for the first time, to the best of our knowledge, a single URT system at the line level, and also against other transport models from the traveler perspective, focusing on the efficiency and sustainability, and skipping the wide range of sociodemographic variables that require two-step modeling, as for [30,33].

Thus, the combination of EDA and DEA will be able to monitor, understand, and improve URT management.

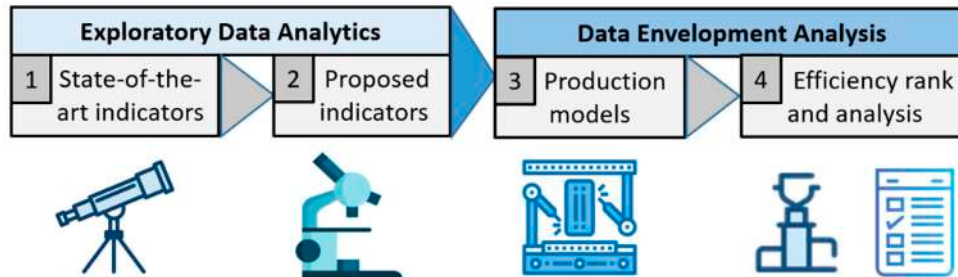
### 3. Methods and Materials

Despite the relevance of URT for the development of sustainable cities, there is a lack of research on the efficiency and sustainability of URT systems and their management. The increasing availability of data, both personal (e.g., GPS location) and Internet-of-Things (IoT) big data, is expected to play a key role in the development of tailor-made mobility solutions, also known as Mobility-as-a-Service (MaaS) [41], based on convenience, sustainability, and resource efficiency to meet passengers' individual needs.

This paper introduces a methodology for assessing the efficiency and sustainability of an URT network based on large-scale data analytics consisting of four stages: (1st) EDA using state-of-the-art indicators; (2nd) EDA using new proposed indicators that deepen the analysis; (3rd) DEA using several original transport models, and (4th) rank transport modes and URT network elements according to efficiency measures, analyzing the results. Figure 1 summarizes graphically the methodology.

As URT systems, particularly in large cities, are a combination of complex interrelations, the proposed methodology aims at better capturing the most relevant efficiency and sustainability

indicators to optimize transport infrastructures, from planning to real-time operation. Furthermore, as an additional outcome of this methodology, open-data repositories could be enriched with new data sources such as occupancy rates, queueing times, URT network elements capacities (e.g., stations) as well as CO<sub>2</sub> footprints.

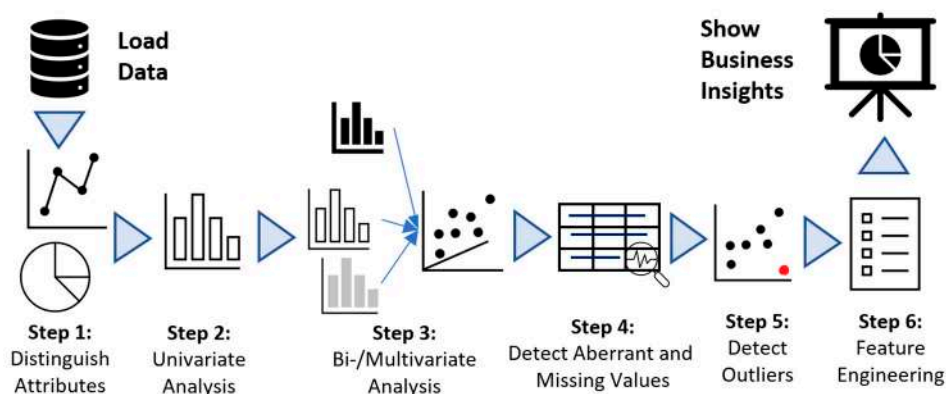


**Figure 1.** Overview of the proposed methodology.

### 3.1. Exploratory Data Analysis (EDA) of URT Data

The first stage of the proposed methodology uses EDA for deriving state-of-the-art quantitative indicators [30]: network length, number of stations, the number of trains, the number of frequencies, the number of employees, the number of operated kilometers, and the number of passengers. This data is usually publicly available at transport operator level, useful for comparing operator's efficiency, but it is more difficult to find at line level, limiting the analysis of the efficiency of rail network elements. However, thanks to big-data technologies (e.g., logging API requests/responses, queueing transport events, and web scraping) these indicators can be potentially estimated using models at a more fine-grained level. In the absence of data from operators (according to [28] only 9% of the research papers in this area has access to official data, generally open data) relying on big data is a much more scalable and cost-effective solution than ad hoc surveys. This approach will contribute to deepening the analysis of transport operators, thus increasing the limited number of research papers with city coverage (only 6% in [28]).

EDA, also known as Visual Analytics, is a heuristic search technique for finding significant relationships between variables in large datasets. Its simplicity and efficiency are key to derive insights from big data, in fact, it is usually the first technique when approaching data, particularly unstructured. According to Tufféry [42] EDA usually consists of six steps (see Figure 2) namely: (i) Distinguish/Identify Attributes; (ii) Univariate Data Analysis to characterize the data of the dataset; (iii) Detect Interactions Among Attributes performing bivariate and multivariate analysis; (iv) Detect and minimize impact of Missing and Aberrant Values; (v) Detect Outliers (further analysis or errors), and finally (vi) Feature Engineering, where features are transformed or combined to generate new features.



**Figure 2.** Exploratory data analysis (EDA) steps.

There is a large number of tools for performing EDA (50 of them are analyzed in [43]) with different functionalities to assist both with the identification of hidden patterns and correlations among attributes, but also with the formulation of hypotheses from the data and their validation. EDA can also be performed using R, python (used in our research work, programming ELTs—Extract, Load, and Transforms—followed by Datawrapper visualization) or any other programming language oriented to data preparation and exploration. Additionally, due to the geographical dimension of transport it is relevant that the tool includes Geographical Information Systems (GIS) support and a strong set of visualization capabilities.

### 3.2. Efficiency and Sustainability Key Performance Indicators (KPIs) for URT

The output of EDA is the estimation of state-of-the-art Key Performance Indicators (KPIs), as well as defining new ones based on large-scale data. For instance, new KPIs that can be defined are the number of trains per line that could be estimated based on the travel time and the rail frequencies. Moreover, another KPI, the number of passengers per line, can be estimated from the number of trains and the entry/exit numbers at the stations of a line. Finally, URT CO<sub>2</sub> footprint can be estimated from the annual supply (in GWh) and the breakdown by source of the consumed electricity and their CO<sub>2</sub> respective footprints.

Additional candidate KPIs that can be modeled after big-data sources are:

- Occupancy ratio: considering 100% occupancy ratio equals to all seated places plus 4 standing people per m<sup>2</sup>.
- Trains Per Hour (TPH): the number of trains that enter and leave a given station per hour, monitoring real-time traffic conditions.
- Travel time: from origin to destination, involving an URT journey.
- Excess Journey Time: is the additional time on top of scheduled time.

The definition, measurement, and analysis of the evolution of KPIs is key to improve the efficiency, security, convenience, and sustainability of existing URT. In fact, the public availability of these KPIs might support that personalized preferences for route selection can be expanded to, for instance, the occupancy rate, for ensuring the availability of seating space, CO<sub>2</sub> footprint, or risk of an excess journey time higher than 10 min. Currently the preferences for route selection are quite rigid, the faster route, or manifest preference for a transport mode, although eventually passengers are considering additional factors, as seen when analyzing, anonymously, their routes using Wi-Fi data.

### 3.3. Data Envelopment Analysis (DEA) for Assessing Efficiency and Sustainability of Public Transport

DEA is a non-parametric method to measure the performance of entities, called Decision-Making Units (DMUs). A DMU can be a factory, a bank branch, a hospital and, as in our paper, a transport mode, an URT line, or an URT station. The initial DEA models consider Constant Return to Scale (CRS or CCR for Charnes, Cooper, and Rodhes), which ignores the fact that different DMUs could be operating at different scales. In our scenario, it would not make any distinction between two URT lines, one with 6 stations and another with 60 stations. To overcome the drawback the Variable Returns to Scale (VRS or BCC for Banker, Charnes, and Cooper) mode [44] was introduced, ensuring that DMUs are only benchmarked against DMUs of similar size. Figure 3 presents an example of four DMUs and both CRS and VRS efficiency frontiers. DMU 1 is the only one in CRS efficiency frontier (the only efficient in CRS), maximizing the output/input ratio, whereas DMUs 1, 2 and 3 are in VRS efficiency frontier (the three are efficient in VRS, DMU 2 in low input values and DMU 3 in high input values). Further to VRS, a wide range of DEA models have been designed for measuring efficiency and capacity specializing the original models into different types of problems.

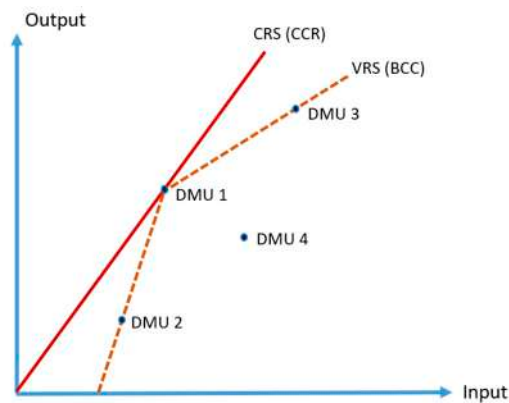


Figure 3. DEA CRS and VRS efficiency frontiers and four DMUs.

DEA models can be classified in either input-oriented or output-oriented models. Figure 4 shows an inefficient DMU (DMU 4 or C) to exemplify both approaches. Input-oriented efficiency is  $BA/CA$ . Output-oriented efficiency is  $CD/ED$ . With input-oriented DEA, a DMU computes the potential savings of inputs in case of operating efficiently (in Figure 4 reducing the inputs from C to B while providing the same output). In contrast, with output-oriented DEA, a DMU measures its potential output increase given its inputs do not vary (in Figure 4 increasing the outputs from C to E while using the same amount of input, D). If C were in the frontier, so  $C = B = E$ , the efficiency would be 1.

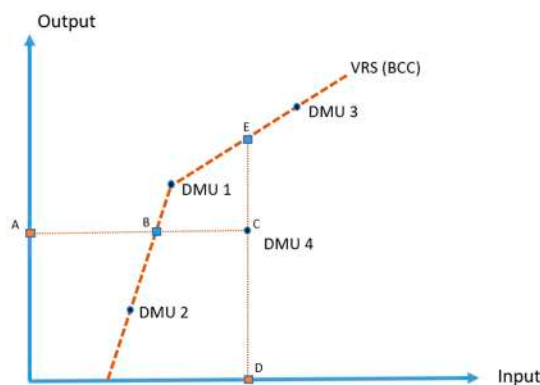


Figure 4. DEA VRS efficiency frontier and DMU 4 efficiencies.

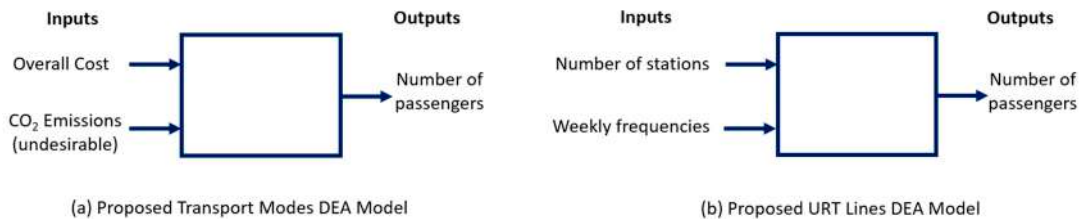
The bad/undesirable outputs, in our case  $\text{CO}_2$  emissions, have been treated as inputs reversing traditional DEA models [45,46]. This technique is based on the fact that undesirable outputs can be treated as inputs when there is a combination of undesirable and desirable outputs. The objective is to minimize the undesirable output, so considering it as input the function looks for its minimization.

A DEA model is a particular selection of inputs and outputs to analyze the efficiency of DMUs. In previous DEA assessments of transit lines, labor, capital, and energy have been used as inputs and vehicle-kms and passenger-kms have been used as outputs. In the absence of actual costs of labor, fuel/energy, and other operational expenses for individual transport lines, it is reasonable to assume that the cost of operating a line is related to its travel time, round-trip distance, and the number of stations/bus stops [35]. Additional, when alternative transport options are being considered, the cost is usually the single input whereas travel time savings, patronage (people for each transport mode), and car trips removed are outputs, as shown in [38], a study that implemented a constant returns of scale–output-oriented (CRS–O) model.

Figure 5 presents our candidate DEA models for: (a) assessing different transport modes from the traveler's viewpoint (for route planning), and (b) analyzing URT lines from the operator/local authority perspective. The analyzed DMUs are the available transport modes (e.g., URT, bus, car, taxi, walking, and cycling) for the first model, and the available URT lines, usually in the range of 1 to 24 lines (e.g.,



New York has the highest number of metro lines, 24, followed by Beijing 23, Seoul 23, Shanghai 17, Paris 16, Moscow 14, and Tokyo 13). CRSs are considered for both models, in the transport modes model because route planning is generally used for one traveler (or a small group) and DMUs operate in the same scale, whereas URT lines, for a given URT network, are usually directly comparable.



**Figure 5.** (a) Transport modes efficiency and (b) URT lines efficiency DEA models.

The selection of inputs and outputs is especially relevant in this scenario due to the large number of available variables and the modest sample size. Following the cardinality constraints introduced in [39], the recommended number of variables for these two CRS models is 3 (in case of considering VRS it would be 2). The selected variables depend eventually on EDA on the available data; however, a tentative output is the number of passengers, and the models can be considered input-oriented, designed for minimizing inputs when moving a given number of people.

In the first model CO<sub>2</sub> emissions, an undesirable output, have been treated as input, as already mentioned, while considering the overall cost the only true input [38]. Here, as the model is from the passenger viewpoint, the overall cost is the transport fare (or the direct costs incurred) plus the monetary value of the passenger time. As most of the mobility is associated with commuting to work, the passenger time value can be estimated at the cost of unskilled working time, although this can be configured on a per-passenger basis for personalized route planning. The selection of these two inputs, which combined with the output reaches the recommended number of variables (3), is original, selected after using EDA on the available data, which contrasts with state-of-the-art indicators for route planning such as travel time and fare cost.

With regards to the URT lines model, two tentative inputs, subject to change due to EDA conclusions on the available data for a given URT network, are considered: (1) the number of stations per line as estimate of the capital costs (CAPEX); and (2) weekly frequencies as operating costs (OPEX).

The related literature in public transport generally uses the actual investment as CAPEX; however, when considering URT lines it is neither directly disaggregated per lines, nor comparable across the time (e.g., 20th century vs. 21st century URT lines). The number of stations per line has been selected as input due to the wide availability of this KPI, although generally indirectly, derived from the longest URT route obtained from online route planning/maps services and applications. The line length, although it is a more popular metric and it is also widely available, has not been selected after EDA on the available data (data sourced from [47]) as long lines usually have generally lower investment due to a higher ratio of above-ground to underground construction, especially in suburban areas where distance between stations tend to be higher. In fact, in [47], a reference paper in CAPEX in Urban Rail only considers costs per kilometer, a state-of-the-art KPI, which shows a higher variability than the cost per station (e.g., in 16 European URT projects, after discarding 3 outliers, the cost per kilometer ranges from 26.7 to 88.3 M USD\$, whereas the cost per station ranges, for the same projects, from 39.4 M to 83.1 USD\$), with lower standard deviation. Additionally, stations have a share of 25–30% of the infrastructure costs, which favors the selection of the number of stations versus line length as CAPEX.

Regarding OPEX inputs, the related literature in public transport generally uses the price of labor and the price of fuel. However, they are not particularly useful for comparing different lines within the same URT system, as they are set at the operator level. Labor and energy consumption can vary per line, although this level of detailed data is generally not available. Nevertheless, a variable directly related to OPEX that is generally available per line is the number of weekly frequencies. EDA on route

planning data shows different patterns for weekdays and for weekends, so the week is the selected period. This input, in combination with the number of stations (the other input), are the selected variables for this DEA model after EDA on publicly available data from URT systems.

The shortlisted inputs (e.g., number of stations and weekly frequencies) have a relevant positive correlation with most of the state-of-the-art inputs, such as the line length, labor force, and number of URT cars, as shown from EDA on [47] and validated in Section 4 (e.g., using LU lines key parameters), thus making it a highly representative selection, with higher discriminatory power and simplicity thanks to minimizing redundancy. Furthermore, the shortlisted inputs are directly obtained from route planning services and applications (e.g., Apple Maps, Bing Maps, Google Maps, and services such as Rome2rio.com that, as of today, includes worldwide 176,885 rail lines from 4151 operators), significantly easier than collecting data from other sources, some of them not available publicly. Finally, in case of availability of data, our candidate inputs for a more representative model would be car capacities, consider line branches, and breakdown passengers into time bands (a.m./p.m. peak versus off-peak). The selection of these additional candidate inputs, which add relevant information about URT efficiency, is one of the outcomes of the previous step, defining new indicators from EDA.

These DEA models have been computed using the solver software that comes with the reference DEA book by Cooper [6]. To illustrate DEA concepts this subsection concludes with an example of DEA analysis, computing the efficiency of London Underground (LU) lines, the use case to validate the proposed models, using for clarity purposes a simplified URT lines DEA model, with a single input, the number of stations, and a single output, the number of passengers. Table 1 summarizes the input and output data, as well as the results provided by the solver. As there is a single input/output the resolution is direct.

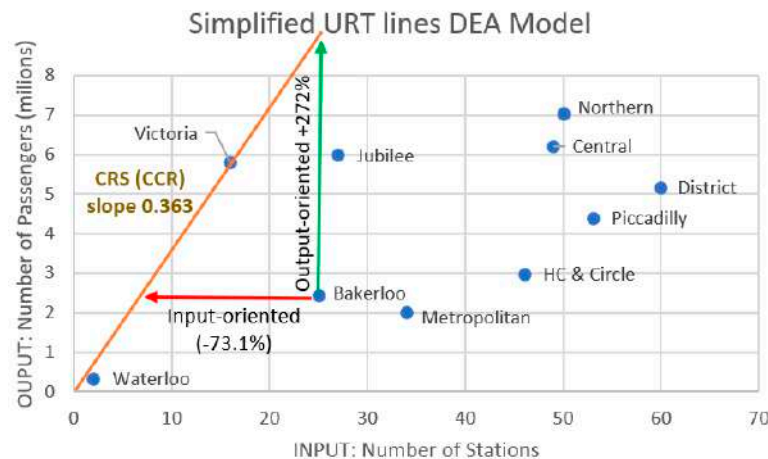
The DMU Victoria maximizes the production function (weekly passengers per station), 363,000, so it scores 1. Compared to the first DMU of the list, Bakerloo, with 98,000 passengers per station, 26.9% of 363,000, thus scoring 0.269. This is a CRS model, similar to the two proposed models, so the production function is the same for all DMUs, not varying at scale (as for VRS). Since there is a fixed number of stations, the key parameter is the number of passengers that maximizes the efficiency for each line, so the model has been computed as output-oriented. In fact, the highest ratio, 363,000 passengers per station, has been used to compute the projection of passengers, presented in Table 1, as well as the difference between the projection and the actual line passengers. Thus, for Bakerloo, ranking 6th in Efficiency, the projection is 9.08 million passengers, +272% over the actual number of passengers, 2.44 million passengers. Alternatively, models can be computed following an input-oriented approach, thus minimizing the required number of stations to achieve the maximum ratio. Thus, for Bakerloo line, it would need to carry 2.44 million passengers with 6.7 stations ( $2.44/0.353$ ), which is 73.1% less stations (1 minus its efficiency score, 0.269).

**Table 1.** Results of simplified URT lines DEA model of LU lines.

LU Line (DMU)	Num. of Stations Longest Route (I)	Weekly Passengers (millions) (O)	Ratio Passeng./Station (thousands)	Efficiency Score (Solver Output)	Rank	Passengers Projection (millions)	Difference (%)
Bakerloo	25	2.44	98	0.269	6th	9.08	272%
Central	49	6.22	127	0.349	5th	17.8	186%
District	60	5.17	86	0.237	7th	21.8	322%
H&C and Circle	46	2.99	65	0.179	9th	16.71	459%
Jubilee	27	5.99	222	0.61	2nd	9.81	64%
Metropolitan	34	2	59	0.162	10th	12.35	517%
Northern	50	7.03	141	0.387	4th	18.17	158%
Piccadilly	53	4.4	83	0.229	8th	19.26	337%
Victoria	16	5.81	363	1	1st	5.81	0
Waterloo City	2	0.33	166	0.456	3rd	0.73	119%

Figure 6 represents graphically the 10 DMUs (URT LU lines) using their coordinates (number of passengers as y axis and number of stations as x axis). The production function, CRS, achieves its maximum value for Victoria, thus scoring 1 in efficiency. Please note that the CRS function starts at the

origin (0,0). The remaining DMUs score below 1, depending on its ratio passengers/station compared to the optimal. The least efficient is Metropolitan, graphically it can be seen that it has the minimum slope to the origin. The figure also helps to understand how to measure inefficiency. Using Bakerloo as a sample, on the one hand, for input-oriented, the CRS optimal function requires 73.1% less stations (6.7 stations) for moving 2.44 million passengers. On the other hand, for output-oriented, CRS optimal function can move 9.08 million passengers, +272%, with 25 stations.



**Figure 6.** Simplified URT lines DEA model computed with 10 LU DMUs, CRS function, and efficiency measures for Bakerloo DMU.

### 3.4. Ranking DEA Models URT Lines According to Efficiency Indicators

The fourth and later stage of our methodology is to rank both transport modes and URT lines using the results of the DEA models. The efficiency of the transport models, from the traveler's viewpoint, can be used for personalized route planning, suggesting different transport modes depending on the time band, the travel distance and the user preferences (e.g., their own estimate of its value of time, and the usage of new mobility solutions such as private electric scooter, or bike/moto/car-sharing).

With regards to URT lines, ranking them according to their efficiency scores instead of less sustainable metrics, such as the number of car-kilometers or the increase in the number of passengers, contributes to align the public transport operation with sustainability goals. In fact, the most efficient URT lines will be those with a reduced number of stations and weekly frequencies that are able to transport more passengers. This model/rank can be complemented with the personalized route planning, as the frequency between URT services could be modified (increased/decreased) up to a point where URT is still the preferred transport choice.

### 3.5. Big Data and Sustainable URT

Public transport services, particularly URT systems, due to the economies of scale, are among the most efficient activities. However, they confront huge initial capital investments, and variables such as the number of stations, length, speed, are determined by this capital investment. Therefore, it is key to characterize their efficiency and sustainability, key to monitor its management.

Big data can gather, store, and process large amounts of heterogeneous, large-scale data to assist regulators, cities, transport operators, and travelers to improve the efficiency, regulation enforcement, and sustainability of their mobility solutions. So far route planning (e.g., Masivo model [48]) and public transport timetable optimization [49] are based on simulation models which can greatly benefit from the incorporation of big-data analysis into their models. Additional big-data applications are personalized route planning and smart taxation (based in the polluters-pay principle) such as dynamic tolling depending on the specific CO<sub>2</sub> footprint of cars and their usage (kilometers) in city centers, where air quality has one of the highest impacts on people's health.

#### 4. Case Study: Efficiency and Sustainability of London Underground (LU)

This section presents the validation of the proposed methodology by analyzing the efficiency and sustainability of a reference URT network, the LU, selected because of the complexity of its network (3 million daily journeys, served by 540 trains across 10 lines covering 402 Km and 263 stations. Figure 7 presents the core of the LU network), and its open-data NUMBAT database (see Appendix A), one of the few publicly available and successful [50] datasets on URT.

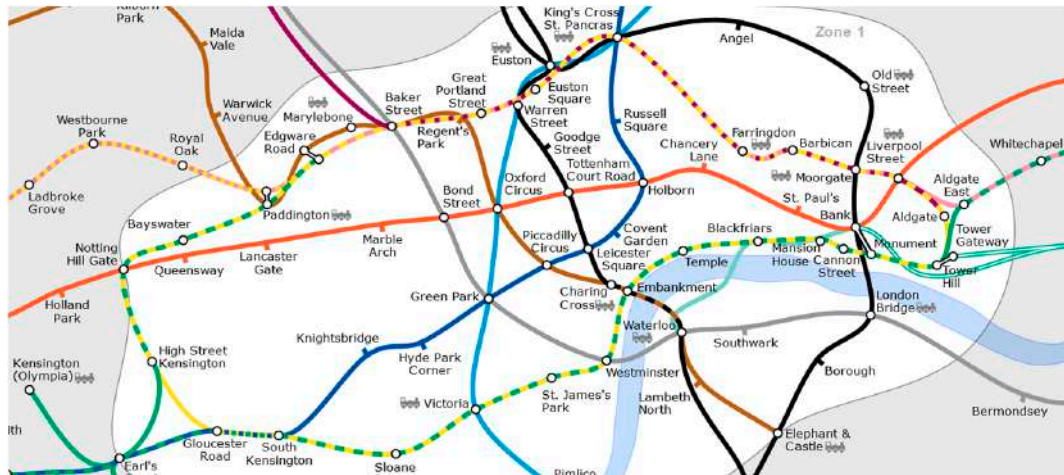


Figure 7. Map of LU lines (colored using the official palette) in Central London.

NUMBAT provides entry/exit/interchange passenger count for 263 stations and the number of trains per station every quarter hour. Additionally, it provides a  $263 \times 263$  origin station–destination station matrix, covering all journeys and the annualized number of passengers for each line. However, NUMBA data is based on real data, but it is not real data. It is the output of a synthetic model used to research LU usage and travel patterns. Moreover, it assumes a perfect train schedule being operated and that all passengers board on the first train arriving at the station. This synthetic model is based on sampling real data from smartcards and gateline entry/exit totals for each station. Data is provided in quarter hours, grouped also by time bands (Early 3–7, AM Peak 7–10, Midday 10–16, PM peak 16–19, Evening 19–22, Late 22–3). Finally, data has been provided in a differentiated way for Fridays, Saturdays, Sundays and for the average of the remaining days (from Monday to Thursday).

As NUMBAT is quite limited (e.g., there is no information about schedules and LU lines, neither descriptive, nor the stations that belong to a line nor the capacity of the trains), we have extended this database with four major data incorporations: (i) train schedules; (ii) a table that relates lines with all their stations; (iii) a table that relates lines with their capacity (seated plus standing at 4 passengers per  $m^2$ ), with data collected from TfL website (TfL open data does not include this data); and (iv) include GPS location for all the stations, obtained from Open StreetMap [51]. See Appendix A for further details. Figure 8 presents some key descriptive metrics of LU which are not originally available in its open-data repository.

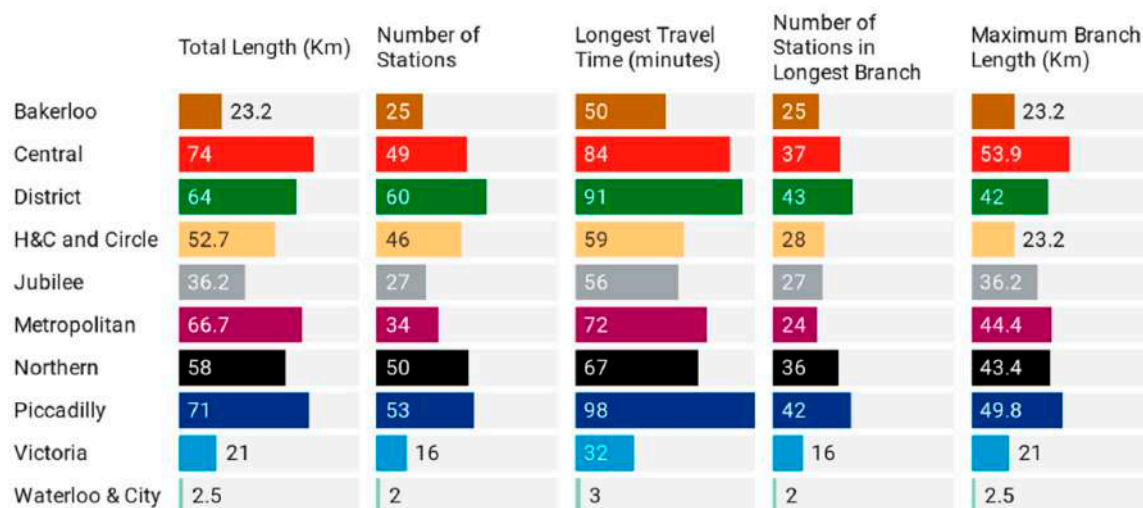


Figure 8. LU Key Descriptive Indicators.

#### 4.1. Assessing the Efficiency and Sustainability of LU Using EDA

The first step of EDA is to distinguish attributes. Table 2 gathers LU key attributes: 3-letter LU line code (in the same order as Figure 8); the longest travel time in the line, it is the average scheduled time of the longest service, usually from the first until the last station of the line; and the length, in kilometers and stations, of the longest route. Additionally, the table contains the scheduled weekly LU frequencies at the station with the highest number of frequencies (usually stations at the middle part of the line), and the weekly passengers per line. A passenger counts as one passenger for each of the lines traveled. On average, a LU passenger uses 1.6 lines per journey (42.4 Weekly passengers in lines and 26 million weekly LU journeys).

The next parameters in Table 2 are metrics/KPI derived from the previous data. Figure 9 presents the scatter plot graphs of the number of passengers versus the number of stations (left), two variables that correlate positively with  $R^2 = 0.55$  (the higher the number of stations, the more travelers it captures). Figure 9 also shows the number of passengers versus the line length (right), with  $R^2 = 0.33$  (a long LU line might be reaching areas with less population density, so this correlation is weaker than the previous one). Additional parameters are the average number of passengers per service and station (included as it contributes to explain the variability with  $R^2 > 0.5$ , discarding the line length). Finally, Speed, in terms of km per hour and minutes per station is presented to illustrate key metrics of LU operation. Based on these analyses, two parameters, the number of stations of the longest route and the weekly frequencies, have been selected to be used in the second phase of the proposed methodology, efficiency scoring using DEA.

So far, the analyzed metrics are average numbers, not considering a relevant source of variability, the day of the week and especially the time band. Figure 10 presents the number of passengers per line and day of the week. The dataset provides an average number from Monday to Thursday. Fridays, except for the Metropolitan and Waterloo & City lines, is the busiest day, whereas Sundays is the day with the lowest number of passengers.

Table 2. LU lines key parameters.

LU Line	Longest Travel Time (minutes)	Longest Length (km)	Num. of Stations Longest Route	Weekly LU Frequencies (scheduled)	Weekly Passengers	Avg. Pass. per Service	Avg. Pass. per Service and Station	Speed (km/h)	Speed (min. per Station)
BAK	50	23.2	25	4836	2,444,910	506	21	28	2
CEN	84	53.9	37	6202	6,218,138	1003	27	39	2.3
DIS	91	42	43	4828	5,166,660	1070	25	28	2.1
HAM	59	23.2	28	3184	2,988,540	939	34	24	2.1
JUB	56	36.2	27	6624	5,985,450	904	33	39	2.1
MET	72	44.4	24	4718	2,003,527	425	18	37	3
NOR	67	43.4	36	10,610	7,028,737	662	18	39	1.9
PIC	98	49.8	42	5710	4,404,640	771	18	30	2.3
VIC	32	21	16	7480	5,813,439	777	49	39	2
WAC	3	2.5	2	3402	331,156	97	49	50	1.5

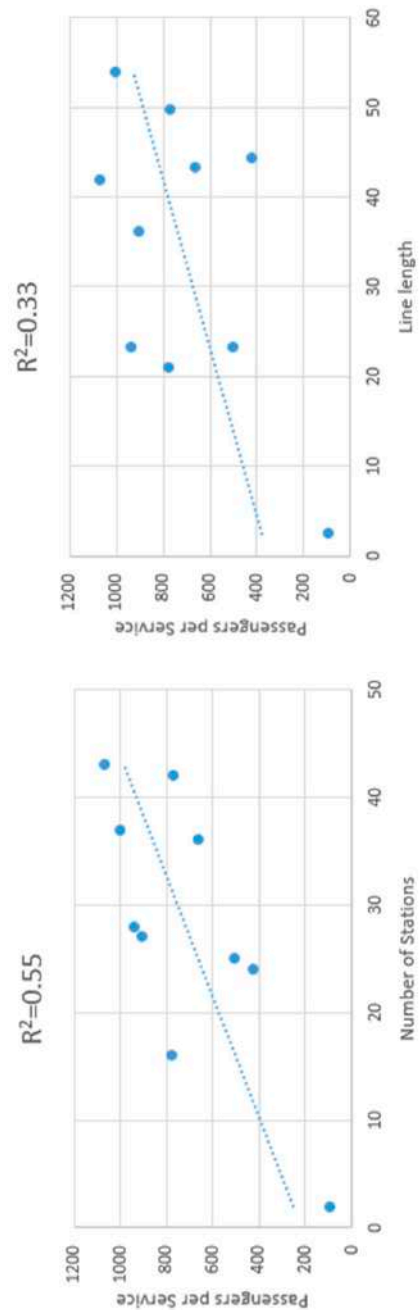


Figure 9. Linear regression of passengers per service versus line stations (left) and line length (right).



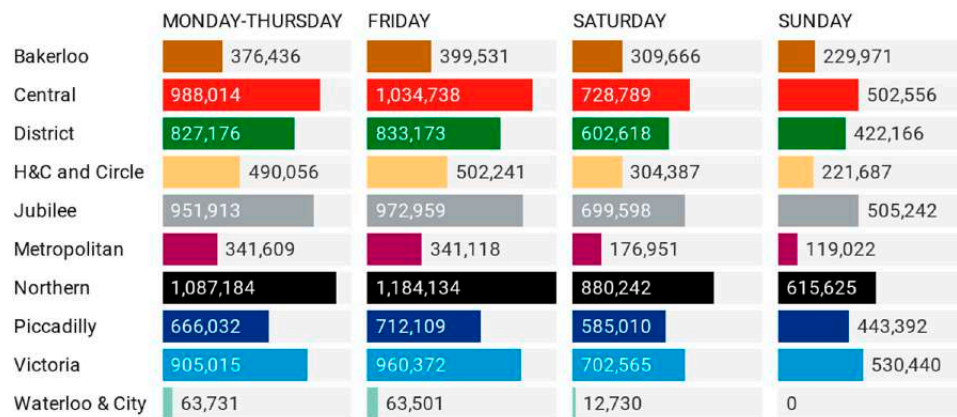


Figure 10. LU Daily passengers per Line.

Figure 11 presents the distribution of passengers per day of week and time bands (Early 3–7, AM Peak 7–10, Midday 10–16, PM peak 16–19, Evening 19–22, Late 22–3). AM and PM peak hours (3 h each) concentrate most of the use from Monday to Friday, whereas Midday (6 h) is the preferred time band for weekend passengers. WAC (Waterloo & City) only operates from Monday to Saturday and has the highest use during Monday to Friday peak hours. Late traffic is higher on Saturday and also Friday, which motivates the different traffic pattern of Friday versus Monday to Thursday, and the higher number of passengers of Friday than Monday to Thursday (except MET and WAC lines, see Figure 10).

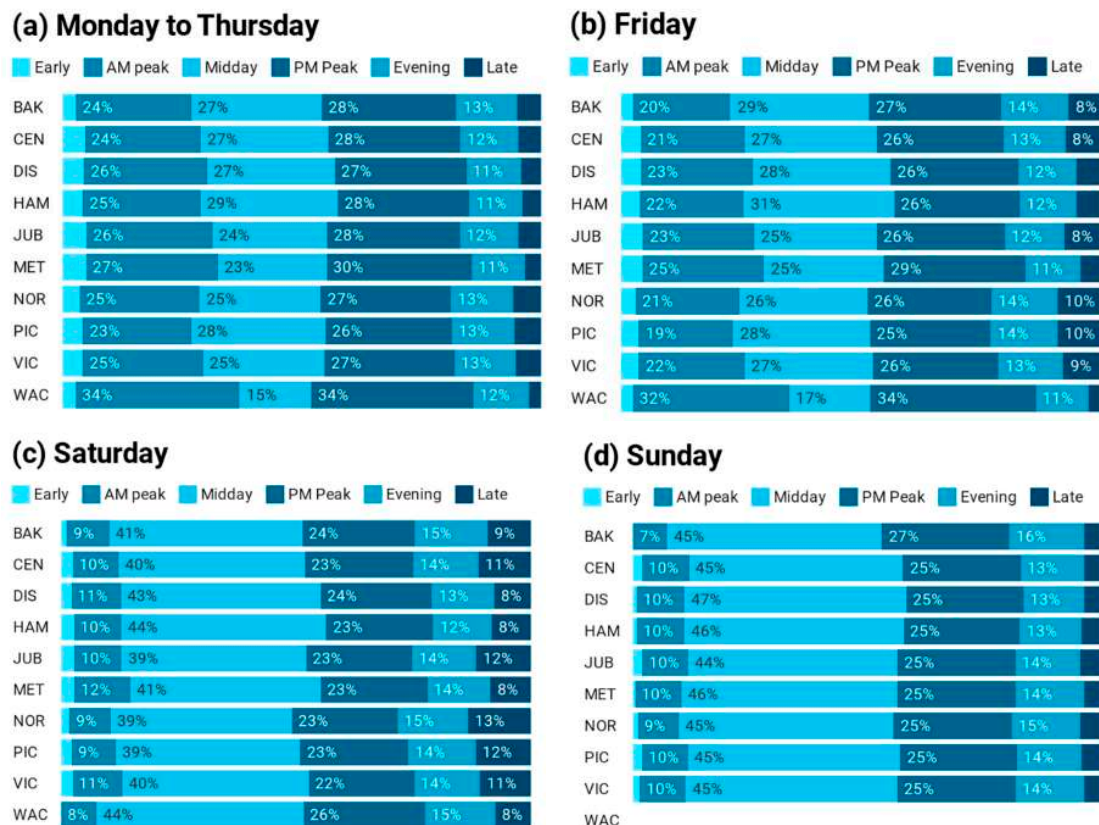


Figure 11. LU passengers per day of week, line, and time band.

A new metric, occupancy rate (usually not reported by URT operators), has been computed dividing the number of passengers by the capacity of the line by time band. To compute this KPI the underground capacity has been considered (seated spaces plus 4 standing passengers per m<sup>2</sup>,

see Appendix A). Figure 12 presents the occupancy rate, sometimes higher than 1 (e.g., Central and District lines). This means that a train, when going from the beginning to the end of the line, can move more passengers than its theoretical capacity. This is possible because these lines, Central and District, have branches and multiple exchanges with other lines, so each seat/standing space can be occupied by more than one passenger per service. A model that estimates the maximum capacity of a line based on an origin–destination trip matrix has been already suggested [52]. However, in our work we will capture these differences in the DEA efficiency model, without providing specific weights to the behavior of line travelers. However, the availability of actual origin–destination data (not the model-based NUMBAT dataset) would increase the interest of this research.

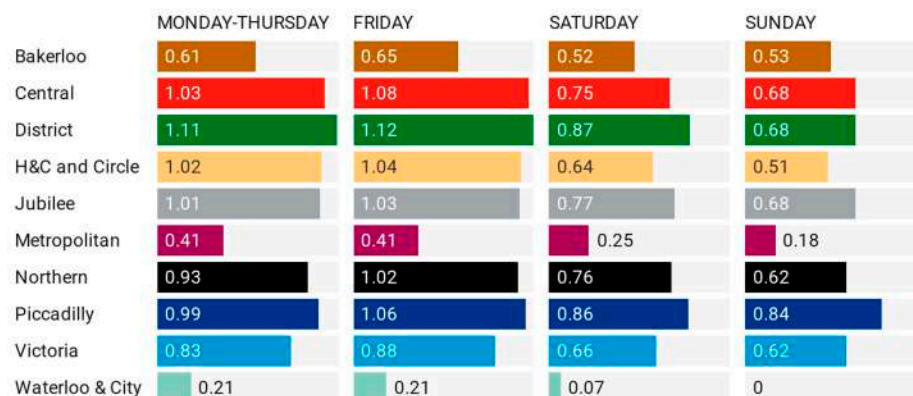


Figure 12. LU occupancy rate by line.

Figure 13 presents the occupancy rate by line, day of the week, and time band. On the one hand, the highest occupancy rates are in PM peak band (4–7 p.m.) from Monday to Thursday, particularly in Central, District and H&C and Circle lines, with rates over 2. As mentioned, on average a LU travel involves 1.6 lines, and these three lines cross Central London, so they might be capturing a relevant number of travels from/to an exchange to another line. In fact, the most crowded line, H&C at PM peak time, has lower traffic at Early time band (before 7 a.m.), which means that is a line close to weekday main destinations (Central London). On the other hand, the occupancy rate of Metropolitan and WAC is the lowest.

Next step is to explore the occupancy rate between two contiguous stations, to characterize the real occupancy rate experienced by travelers. The number of station links is the number of stations minus one for each line, thus 352 station links. The most relevant information analyzing occupancy rates are those extreme values, the lowest and highest, particularly the latter. Figure 14 shows the most crowded station links at the quarter hours with the highest occupancy rates during AM peak (left), 8:30–8:45 a.m., and PM peak (right), 5:30–5:45 p.m. These numbers have been derived from our dataset, combining passengers, line schedules, and line capacities. However, these are estimates as the real flow of passengers and train delays are not publicly available. As the objective of this paper is to characterize the efficiency and sustainability of LU, EDA finishes with the analysis of occupancy rates of stations links, relevant for assessing that the LU carriages theoretical capacity (with 4 standing people per m<sup>2</sup>) can be considered its maximum capacity.



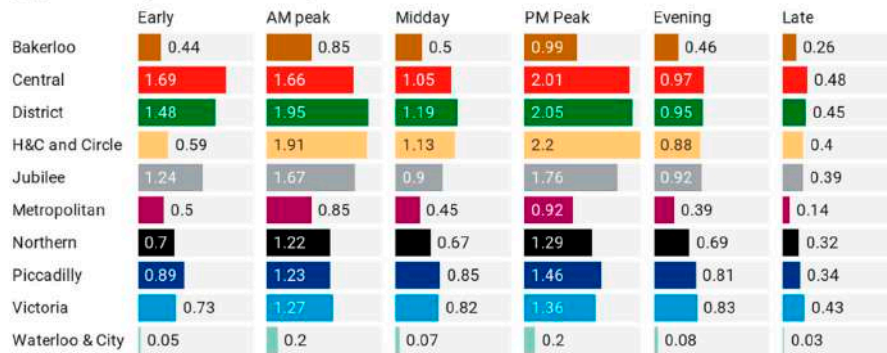
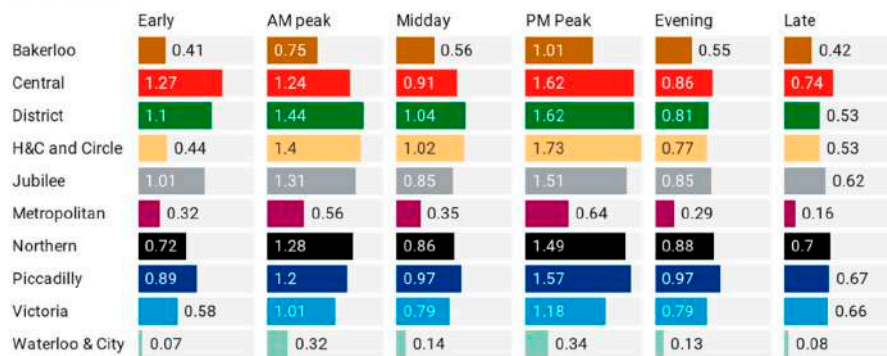
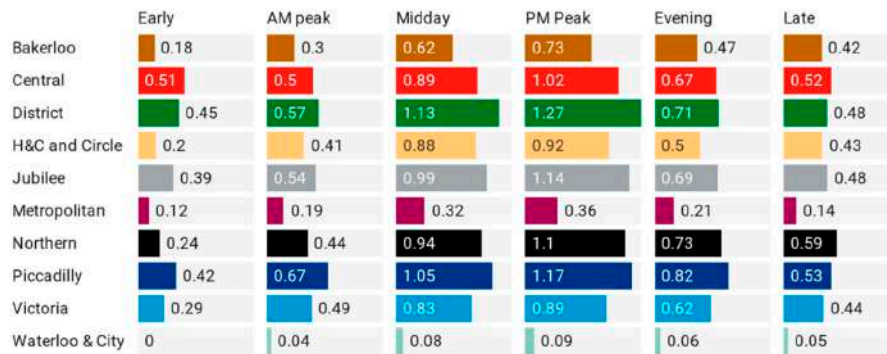
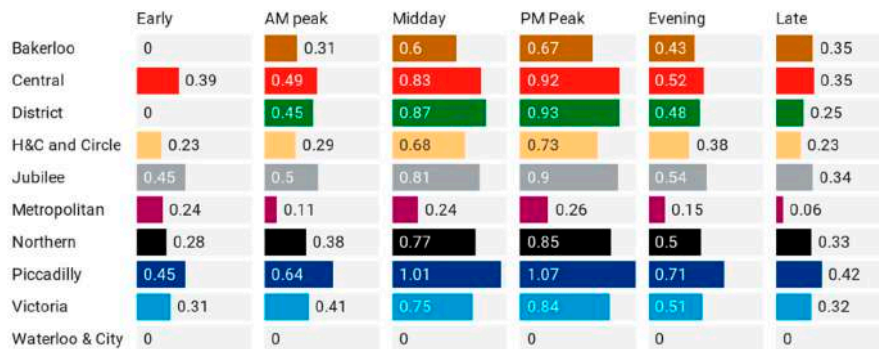
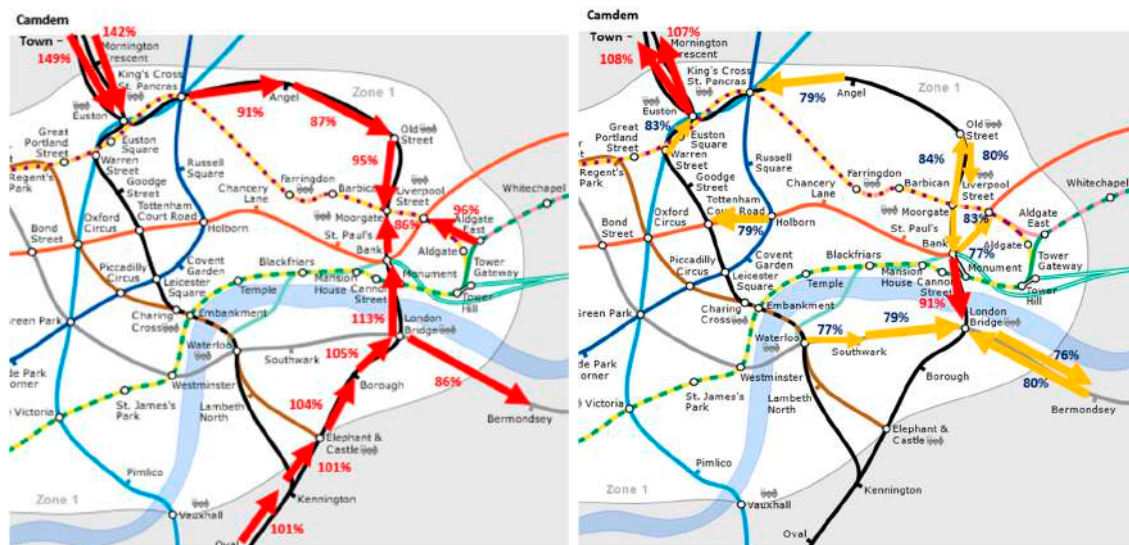
**(a) Monday to Thursday****(b) Friday****(c) Saturday****(d) Sunday**

Figure 13. LU occupancy rate by line and time band.



**Figure 14.** Map of Central London showing the highest occupancy rate of LU stations links at 8:30–8:45 a.m. (**left**) and 5:30–5:45 p.m. (**right**). Those links with 85% or higher occupancy rates are in red.

#### 4.2. LU Additional KPIs

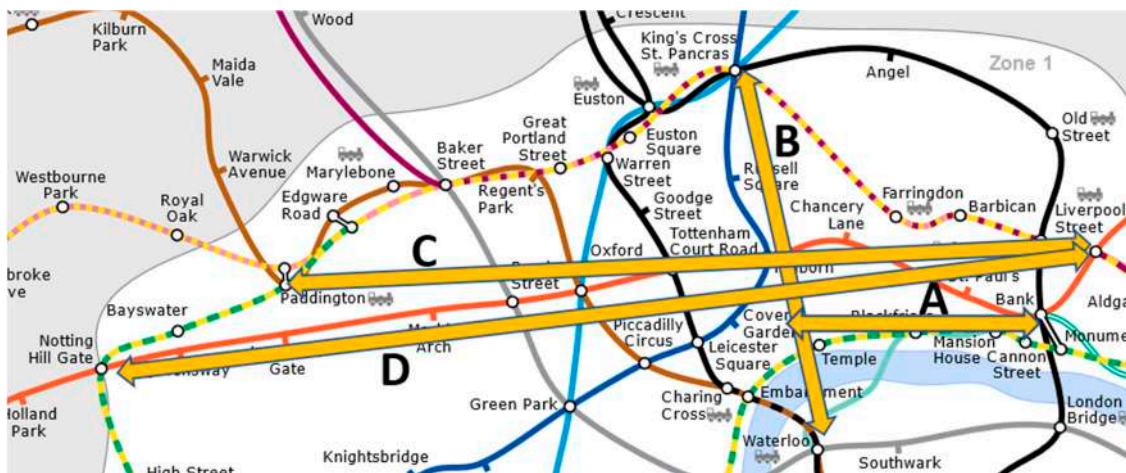
TfL considers additional LU KPIs in its reports [53], focused on service provision, reliability, and journey times, such as the percentage of scheduled kilometers operated (95.8% of the 88.7 million kilometers scheduled), and the excess journey time, and the average delay or (4.6 min, 11% of the average journey time which is 41.6 min). The average delay is formally defined as excess journey time, the additional time on top of scheduled time for access/egress/interchange, platform wait time and on train (the latest figure is 4.6 min for LU for 2018/2019 (table 12.5 in [53])). TfL has reduced the excess journey time since 2008/09, from 6.6 min to 4.6 min by increasing the frequency of the services around 20% higher.

Finally, from the attributable CO<sub>2</sub>-equivalent emissions of operating LU (372,000 tons) and 12 billion annually passenger-km [53], a footprint of 31 g of CO<sub>2</sub>-equivalent has been estimated by us. Previously, TfL released, outside of the open-data repository, its CO<sub>2</sub> footprints with out-of-date higher estimates [54]. Additional non-official estimates exist [55,56], although also out-of-date. Although the number of operated kilometers raised a 20% over the last 10 years, the CO<sub>2</sub> footprint has decreased far more than 20% (LU operates with power and UK National Grid has been reducing more than 20% its CO<sub>2</sub> footprint during this decade). Thus, LU is more sustainable than a decade ago, and more sustainable than buses (97% fuel-based), which have and 90 g CO<sub>2</sub> footprint per passenger per km (480 million vehicle-km, 4.45 billion passenger-km, and an average CO<sub>2</sub> emission of 822 g/km per vehicle, accounting for around 400,000 CO<sub>2</sub> tons).

#### 4.3. Assessing the Efficiency and Sustainability of London Transport Modes Using DEA

This subsection presents the efficiency of the proposed DEA models, first transport modes, and second URT lines.

Figure 15 presents four routes to evaluate five transport modes (LU, bus, car/taxi, walking, and cycling), and potential combinations of these five transport modes, in Central London, from the shortest to the longest: (A) Bank–Covent Garden, (B) King's Cross St. Pancras–Waterloo, (C) Paddington–Liverpool Street, and (D) Notting Hill Gate–Liverpool Street. These are quite popular routes, connecting national rail stations, and commercial, leisure, and residential areas. However, apart from D, they are not directly connected via LU. Here the optimal route (minimizing travel time) for each transport mode, has been suggested by online services for multi-modal route planning (e.g., Rome2Rio, selected for reporting LU and bus distances and fares).



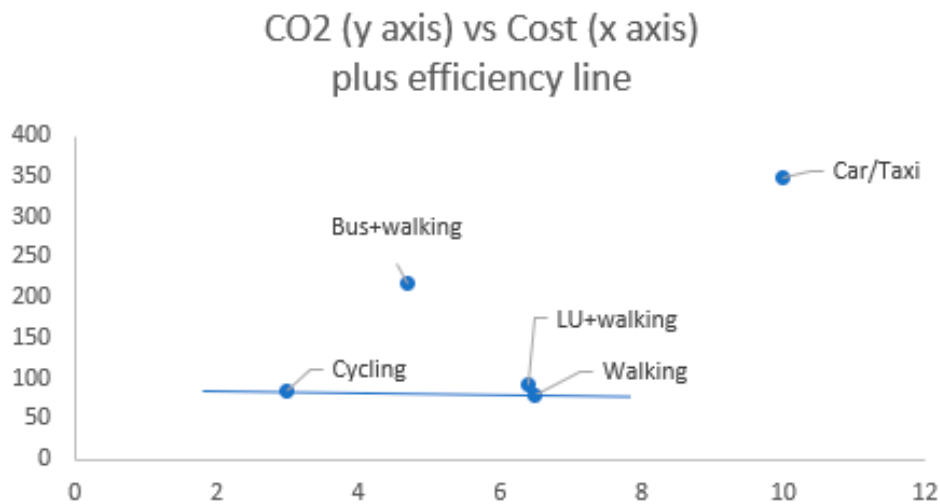
**Figure 15.** Routes under consideration for analyzing transport modes in Central London.

Table 3 presents the key parameters of the five analyzed transport modes for the Route A, and a sixth mode, the combination LU+bus. To be able to run DEA no missing values (or 0) are allowed, so it has been assigned a transport cost for cycling (0.20 GBP, the daily cost of an annual London cycle hiring subscription), and for walking (0.10 GBP per 2.4 km, an estimate of the cost of shoe wear). The estimated value of time is 12.00 GBP per hour, an estimate of unskilled pay rate in London, to consider the time factor. LU fare in Central London (Zone 1) is 2.40 GBP, and TfL Bus fare is 1.50 GBP. Costs are provided in the local currency. Moreover, for cycling and walking the additional physical activity has been also considered, estimating 1 g of CO<sub>2</sub>-equivalent emission per additional Kcal of energy. This number varies with the diet and weight of the traveler, although it is usually in the range 0.5–2 g CO<sub>2</sub>-equiv. per Kcal [57]. The additional cost of walking for a 70 Kg person at 5 km/h in a flat route has been estimated in 150 Kcal/h, and cycling at 15 km/h results in an additional consumption of 360 Kcal/h (these values are average of online calculators). Bus CO<sub>2</sub> emissions are 90 g per passenger per km and LU footprint 31 g per passenger per km. Private car/taxi estimates are 120 g per km, the maximum for driving within the Ultra-Low Emission Zone (ULEZ) of Central London. The number of passengers has been set to 1. These and other values are being used only for illustrative purposes, they can be adapted for personalized route planning and personalized efficiency analysis. Nevertheless, to the best of our knowledge they could be valid estimates.

Computed DEA efficiencies for Route A are 100% efficiency for cycling and walking, particularly for its lowest CO<sub>2</sub> footprint, followed by the combination LU+walking (there is no direct LU link for Route A). Although bus+walking has the second lowest overall cost, its emissions are more than double the most efficient and it scores 64% efficiency. Car/taxi is the least efficient. DEA shows that the limiting factor for improving the efficiency of bus+walking and car/taxi is CO<sub>2</sub> footprint, which can be seen graphically in Figure 16. Shifting from fossil fuel to electric transport can reduce emissions by 75% (according to CO<sub>2</sub> footprint of electricity mix in the UK). Thus, bus+walking would reach the efficiency line whereas car/taxi would increase its efficiency significantly.

**Table 3.** Route A. Results in descending order of efficiency of the transport modes DEA model.

Transport Mode	Distance (km)	Transport Cost (GBP)	Time (min)	CO <sub>2</sub> -Equ. Emissions (l) (g)	Overall Cost (l) (GBP)	Passengers (O)	Efficiency (%)	Rank
Cycling	3.4	0.20	14	84	3.00	1	100%	1st
Walking	2.7	0.10	32	80	6.50	1	100%	1st
LU+Walking	1.7 + 1.3	2.40	20	92	6.40	1	88%	3rd
Bus+Walking	2.3 + 0.4	1.50	16	218	4.70	1	64%	4th
Car/Taxi	2.9	8.00	10	348	10.00	1	30%	5th



**Figure 16.** Scatter plot of the two inputs (output is always 1) of the transport modes DEA for Route A.

Table 4 presents the key parameters for Route B, in descending order of efficiency, from cycling, 100%, down to car/taxi, 21%. However, if 4 passengers go by car/taxi, CO<sub>2</sub> footprint is the same (for clarity purposes we will consider the same), and the overall cost rises from 16 to 30 GBP, whereas for the other transport modes both CO<sub>2</sub> footprint and costs are four times higher than the cost of one passenger. In this scenario, see Table 5, car/taxi jumps to the third efficiency position, rivaling with LU.

**Table 4.** Route B. Results in descending order of efficiency of the transport modes DEA model.

Transport Mode	Distance (km)	Transport Cost (GBP)	Time (min)	CO <sub>2</sub> -Equ. Emissions (l) (g)	Overall Cost (l) (GBP)	Passengers (O)	Efficiency (%)	Rank
Cycling	4.5	0.20	18	108	3.80	1	100%	1st
Walking	4.5	0.20	46	120	9.40	1	90%	2nd
LU	4.3	2.40	23	133	6.80	1	81%	3rd
LU+Bus	1.6 + 2.0	2.40 + 1.50	18	230	7.50	1	51%	4th
Bus	3.4	1.50	35	306	8.50	1	48%	5th
Car/Taxi	4.4	14.00	20	528	18.0	1	21%	6th

**Table 5.** Route B with 4 passengers. Results of efficiency of the transport modes DEA model.

Transport Mode	Distance (km)	Transport Cost (GBP)	Time (min)	CO <sub>2</sub> Emissions (l) (g)	Overall Cost (l) (GBP)	Passengers (O)	Efficiency (%)	Rank
Cycling	4.5	0.20	18	432	15.20	4	100%	1st
Walking	4.5	0.20	46	480	37.60	4	90%	2nd
Car/Taxi	4.4	14.00	20	528	30.00	4	82%	3rd
LU	4.3	9.60	23	532	27.20	4	81%	4th
LU+Bus	1.6 + 2.0	9.60 + 6.00	18	920	30.00	4	51%	5th
Bus	3.4	6.0	35	1224	34.00	4	45%	6th

Table 6 presents the key parameters for Route C, in descending order of efficiency, from cycling, 100%, closely followed by LU, 94%. In this scenario the limiting factor is the cost. Considering 18 GBP per hour as time value then the fastest transport modes increase their efficiencies (LU rises to 100%, Car/Taxi to 40%), whereas slower transport modes reduce their efficiencies (Bus goes down to 52%).

Table 7 presents the key parameters for Route D, in descending order of efficiency, where LU and cycling are both 100% efficient. LU has the lowest overall cost (Cycling has 27% higher cost, 7.60 versus 6.00) and the second lowest CO<sub>2</sub> footprint (254 g., 14% higher than Cycling, the option with the lowest CO<sub>2</sub> emissions with 222 g.). In this scenario both CO<sub>2</sub> footprints and costs are the limiting factors. Thus, electrification of vehicles will have a limited impact if the overall cost remains unaltered.



The main cost reduction would come from reducing even more travel times in buses and car/taxi. This might be feasible reducing traffic in Central London, for instance imposing higher restrictions to polluting vehicles in ULEZ.

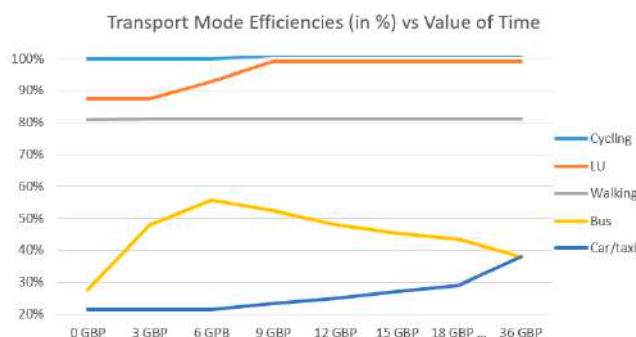
**Table 6.** Route C. Results in descending order of efficiency of the transport modes DEA model.

Transport Mode	Distance (km)	Transport Cost (GBP)	Time (min)	CO <sub>2</sub> -Equ. Emissions (l) (g)	Overall Cost (l) (GBP)	Passengers (O)	Efficiency (%)	Rank
Cycling	8.3	0.20	33	198	6.80	1	100%	1st
LU	7.4	2.40	24	229	7.20	1	94%	2nd
Walking	7.2	0.30	86	215	17.50	1	90%	3rd
Bus	8.3	1.50	55	747	12.50	1	54%	4th
Car/Taxi	8.3	17.00	23	996	21.6	1	31%	5th

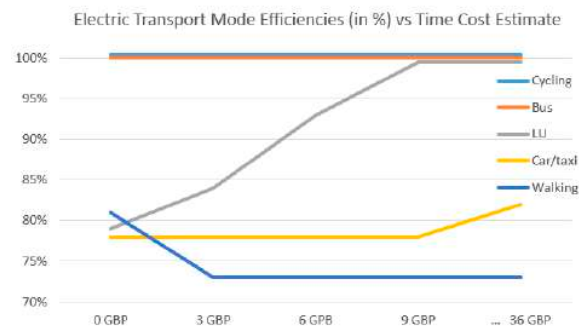
**Table 7.** Route D. Results in descending order of efficiency of the transport modes DEA model.

Transport Mode	Distance (km)	Transport Cost (GBP)	Time (min)	CO <sub>2</sub> -Equ. Emissions (l) (g)	Overall Cost (l) (GBP)	Passengers (O)	Efficiency (%)	Rank
LU	8.2	2.40	18	254	6.00	1	100%	1st
Cycling	9.4	0.20	37	222	7.60	1	100%	1st
Walking	8.6	0.40	109	273	22.20	1	81%	3rd
Bus	8.9	1.50	55	801	12.50	1	48%	4th
Car/Taxi	8.6	19.00	26	1032	24.20	1	25%	5th

Figure 17 presents, considering the latter Route D, an analysis of the sensitivity of the value of time, the main factor impacting the transport mode efficiency. The range considered, from 0 GBP to 36 GBP per hour, shows that the fastest transport modes, Car/taxi and LU, gain efficiency as the value of time increases, slower for the Car/Taxi due to the higher transport cost of this mode. It is remarkable that walking, the slowest transport mode, keeps its efficiency due to its low CO<sub>2</sub> footprint. Thus, a shift from fuel to electric vehicles, reducing the CO<sub>2</sub> footprint by 75%, according to the energy generation mix in the UK, has been considered in Figure 18, also for Route D, together with a varying value of time. Now electric Bus is always 100% efficient, due to its low transport costs and low emissions, very similar to those of cycling, whereas LU, faster than LU and cycling but with a more expensive fare, is also efficient for passengers who value their time from 9 GBP/h on. In a scenario with electric cars/taxis this transport mode (Car/Taxi) is more efficient than walking from 1 GBP/h of value of time. As DEA is a relative (non-absolute) efficiency measure, improvements in some DMUs might impact the efficiency of other DMUs.



**Figure 17.** Route D: Impact of the value of time in the efficiency of the transport modes efficiency.



**Figure 18.** Route D: impact of the value of time and electricity in the transport modes efficiency.

#### 4.4. Assessing the Efficiency and Sustainability of LU Lines Using DEA

This subsection presents the DEA efficiencies of the URT lines model, a CRS model computed with the same DEA software solver as in the previous subsection.

Table 8 presents the key parameters of the ten analyzed LU lines, the two input parameters, the number of stations of the longest route and weekly frequencies, and the output, weekly passengers. Then the efficiency, ranging from 44% for the Metropolitan line to four 100% efficient lines (Central, District, Jubilee, and Victoria line). In addition, finally, four KPIs considered in EDA to characterize and compare LU lines. Although DEA is a non-parametric technique, so efficiency is not a linear combination of the inputs, it looks as if the best performers are those lines with the highest average passengers per service, the highest passengers per service and station, and the highest speed. WAC efficiency (46%, 9th) is limited by the weekly frequencies, with just 426 weekly frequencies (87% lower than the current number), retaining the number of passengers, it would be 100% efficient. The rest of the inefficient lines, those scoring below 100%, are limited both by the number of stations and the weekly frequencies. Table 9 shows the optimal projections of the inputs of the LU lines.

An efficient line would maximize the number of passengers with the lowest number of stations (proxy variable of the capital expenses, CAPEX), and the lowest weekly frequencies (proxy variable of the operating expenses, OPEX), an analysis in tune with previous works [30]. However, to increase the efficiency, closing stations is not an option. URT management can only influence operating expenses, reducing/increasing the weekly frequencies. Thus, the efficiency of a LU line will increase if reducing a given percentage the number of weekly frequencies (e.g., 10%) the number of passengers reduces significantly less than the reduction of the frequencies. Further analysis of real transport data, actual number of passengers and actual schedule of LU trains, will help to understand the relationship between frequencies and the number of passengers of a line, particularly in such a complex network as LU, with multiple exchanges and different lines sharing the same rail section/station links.

Table 10 presents an alternative DEA model with two additional input variables, the longest travel time, and the longest length in km. The new ranking that comes out of this extended model only interchanges positions 8th and 9th, as the new variables are highly correlated with the previous input variables. Thus, now WAC ranks 8th and BAK ranks 9th as the new model favors the short length and travel time of WAC, although BAK also increases its efficiency.

Finally, Tables 11–13 present the efficiency of the proposed URT lines DEA model (the original with 2 input variables) using the data, frequencies, and passengers, for AM, Midday, and PM peak time bands, from Monday to Thursday, respectively. Efficiency results for the mentioned time bands (bands of 3, 6, and 3 h, respectively) are in tune with overall line efficiencies presented in Table 8. However, some differences arise, such as WAC is the 7th in efficiency during peak times, but the 10th during Midday. WAC connects a national rail station and transport hub, Waterloo Station, with Bank tube station, in the heart of the financial area in the City of London. Therefore, its traffic pattern shows more activity during AM and PM peak hours. Moreover, VIC is the only line 100% efficient in the three time bands. Finally, except for WAC, LU lines score similarly across the analyzed time bands.

**Table 8.** Results of the URT lines DEA model of LU lines in descending order of efficiency.

LU Line	Num. of Stations Longest Route (l)	Weekly LU Frequencies (Scheduled) (l)	Weekly Passengers (O)	Efficiency (%)	Rank	Avg. Pass. per Service	Avg. Pass./Service and Station	Speed (km/h)	Speed (min. per Station)
CEN	37	6202	6,218,138	100%	1st	1,003	27	39	2.3
DIS	43	4828	5,166,660	100%	1st	1,070	25	28	2.1
JUB	27	6624	5,985,450	100%	1st	904	33	39	2.1
VIC	16	7480	5,813,439	100%	1st	777	49	39	2
HAM	28	3184	2,988,540	88%	5th	939	34	24	2.1
NOR	36	10,610	7,028,737	77%	6th	662	18	39	1.9
PIC	42	5710	4,404,640	76%	7th	771	18	30	2.3
BAK	25	4836	2,444,910	53%	8th	506	21	28	2
WAC	2	3402	331,156	46%	9th	97	49	50	1.5
MET	24	4718	2,003,527	44%	10th	425	18	37	3

**Table 9.** Optimal LU input projections of the URT lines DEA model.

LU Line	Num. of Stations Longest Route (l)	Weekly LU Frequencies (Scheduled) (l)	Weekly Passengers (O)	Efficiency (%)	Rank	Optimal Num. of Stations	Diff. (%)	Optimal Weekly LU Freq.	Diff. (%)
CEN	37	6202	6,218,138	100%	1st	37	0%	6202	0%
DIS	43	4828	5,166,660	100%	1st	43	0%	4828	0%
JUB	27	6624	5,985,450	100%	1st	27	0%	6624	0%
VIC	16	7480	5,813,439	100%	1st	16	0%	7480	0%
HAM	28	3184	2,988,540	88%	5th	24.6	-12%	2799	-12%
NOR	36	10,610	7,028,737	77%	6th	27.8	-23%	8182	-22%
PIC	42	5710	4,404,640	76%	7th	28.8	-24%	4325	-24%
BAK	25	4836	2,444,910	53%	8th	13.2	-47%	2544	-47%
WAC	2	3402	331,156	46%	9th	0.9	-54%	426	-87%
MET	24	4718	2,003,527	44%	10th	10.6	-56%	2095	-56%

**Table 10.** Results of the URT lines DEA model of LU with two additional input parameters.

LU Line	Num. of Stations Longest Route (I)	Weekly LU Frequencies (Scheduled) (I)	Longest Travel Time (Minutes) (I)	Longest Length (km) (I)	Weekly Passengers (O)	Efficiency (%)	Rank
CEN	37	6202	84	53.9	6,218,138	100%	1st
DIS	43	4828	91	42	5,166,660	100%	1st
JUB	27	6624	56	36.2	5,985,450	100%	1st
VIC	16	7480	32	32	5,813,439	100%	1st
HAM	28	3184	59	23.2	2,988,540	94%	5th
NOR	36	10,610	67	43.4	7,028,737	79%	6th
PIC	42	5710	98	49.8	4,404,640	76%	7th
WAC	2	3402	3	2.5	331,156	60%	8th
BAK	25	4836	50	23.2	2,444,910	58%	9th
MET	24	4718	72	44.4	2,003,527	44%	10th

**Table 11.** URT lines DEA model scores of LU Monday to Thursday (MTT) AM peak traffic.

LU Line	Num. of Stations Longest Route (I)	MTT LU Frequencies (Scheduled) (I)	MTT Daily Passengers (O)	Efficiency (%)	Rank
JUB	27	176	250,037	100%	1st
VIC	16	210	227,188	100%	1st
DIS	43	128	212,294	100%	1st
HAM	28	74	120,077	98%	4th
CEN	37	170	239,754	94%	5th
NOR	36	262	273,060	77%	6th
WAC	2	128	21,875	77%	7th
PIC	42	144	150,140	65%	8th
MET	24	130	93,937	49%	9th
BAK	25	126	90,913	49%	10th



**Table 12.** URT lines DEA model scores of LU Monday to Thursday (MTT) Midday traffic.

LU Line	Num. of Stations Longest Route (I)	MTT LU Frequencies (Scheduled) (I)	MTT Daily Passengers (O)	Efficiency (O) (%)	Rank
CEN	37	294	262,604	100%	1st
VIC	16	328	229,735	100%	1st
DIS	43	218	220,130	100%	1st
HAM	28	146	140,918	96%	4th
JUB	27	298	227,774	95%	5th
PIC	42	258	187,219	76%	6th
NOR	36	482	274,669	74%	7th
BAK	25	240	102,527	51%	8th
MET	24	202	77,589	44%	9th
WAC	2	150	9537	33%	10th

**Table 13.** URT lines DEA model scores of LU Monday to Thursday (MTT) PM peak traffic.

LU Line	Num. of Stations Longest Route (I)	MTT LU Frequencies (Scheduled) (I)	MTT Daily Passengers (O)	Efficiency (%)	Rank
CEN	37	160	273,109	100%	1st
JUB	27	176	263,784	100%	1st
VIC	16	212	245,455	100%	1st
HAM	28	72	134,825	100%	1st
DIS	43	130	227,226	96%	5th
NOR	36	270	295,882	78%	6th
WAC	2	126	21,460	70%	7th
PIC	42	142	175,828	70%	8th
BAK	25	126	105,896	52%	9th
MET	24	132	103,654	50%	10th

## 5. Conclusions

This paper has analyzed the efficiency and sustainability of URT using EDA and DEA. The main contributions of this work are: (1) propose and compute new indicators for EDA of URT sustainability and efficiency (e.g., occupancy rate by URT line, station links, and time band, and CO<sub>2</sub> footprint per journey); (2) design and propose a methodology for DEA performance assessment based on the selection of input and output variables using EDA on publicly available data; (3) develop two original DEA production models, the first one for characterizing the sustainability of different transport modes, and the second one for measuring the efficiency of URT lines; (4) validating the methodology with open data from TfL and online services; and (5) ranking URT against other transport modes and analyzing DEA efficiency scores of URT lines.

The main conclusions of the paper are: (1) EDA plays a key role analyzing URT efficiency and sustainability indicators, as well as defining new indicators; (2) DEA variable selection can be done in a semi-automated and repeatable way relying on EDA; and (3) DEA is a simple and straightforward non-parametric technique to score multiple transport modes and URT lines efficiency to monitor, understand, and improve its management, even focusing on time bands and URT line sections for the latter scenario.

To sum up, the introduced big-data-based methodology supports the advance of efficiency and sustainability in public transport, particularly in URT, through disseminating data, KPIs, and assessments based on them. Thus, both operators and travelers alike are encouraged to improve their decision-making, from transport network management to route planning, to meet the Sustainable Development Goal target of having a more sustainable transport by 2030.

**Author Contributions:** Conceptualization, G.L.T. and L.H.; Investigation, G.L.T. and L.H.; Writing—original draft, G.L.T.; Writing—review & editing, G.L.T. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by the Ministry of Economy, Industry and Competitiveness of Spain, Project TIN2016-75845-P (AEI/FEDER/EU) and SNEO-20161147 (CDTI) and by Xunta de Galicia and FEDER funds of the EU (Centro de Investigación de Galicia accreditation 2019–2022, ref. ED431G2019/01, and Consolidation Programme of Competitive Reference Groups, ref. ED431C 2017/04).

**Acknowledgments:** The authors would like to thank CITIC and Project ref. ED431G2019/01 for supporting the Postdoc visit of Guillermo L. Taboada to Manchester Metropolitan University to develop part of the collaborative work presented in this paper.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## Appendix A

The sources of data used in the study are openly available online on TfL open data site: <https://data.tfl.gov.uk> (accessed on 25 June 2020).

The main source of information is Urban Rail Passengers Count and Travel Flow Dataset (codenamed project NUMBAT), <https://crowding.data.tfl.gov.uk>, which is based on data, from smartcards (Oyster and contactless bank/NFC cards), gatelines, and automatic passenger counters and services from timetables, combined through a model to assign journeys to routes using generalized journey time. Data has been obtained during the autumn of each year (at the time of writing, autumn 2018 is the last one) and provides an average for weekdays from Monday to Thursday, and additionally data for Friday, for Saturday and for Sunday, each of them independently. Includes only aggregated data (~100 MB data annually). Additional information from TfL origin–destination dataset is described here: <https://data.london.gov.uk/dataset/tfl-rolling-origin-and-destination-survey>.

This dataset, published as open data using an open TfL license, provides:

- Journeys per day by all URT modes in London, LU, London Overground (LO), DLR, Crossrail EZL, and London Trams.

- -- Values for each 15-minute period of the day.
- -- Number of passengers' entries/exits at all stations.
- -- Number of interchanging passengers at all stations.

This data has been enriched including GPS location for all the stations, obtained from Open StreetMap [51], (ii) creating a table that relates lines with all their stations (TfL open data does not include this table), and (iii) a table that relates lines with their capacity (seated plus standing at 4 passengers per m<sup>2</sup>), with data collected from TfL website (once again, TfL open data does not include this table). Table A1 presents LU lines and its associated train capacities.

**Table A1.** LU train capacity per line as of 2018.

Line Name	Capacity
Bakerloo	851
Central	1047
District	1045
H&C and Circle	1045
Jubilee	964
Metropolitan	1176
Northern	752
Piccadilly	798
Victoria	986
Waterloo & City	506

## References

1. International Association of Public Transport UITP. Energy Efficiency: Contribution of Urban Rail Systems. Available online: <https://www.uitp.org/energy-efficiency-contribution-urban-rail-systems> (accessed on 25 June 2020).
2. International Energy Agency. "Tracking Transport". Available online: <https://www.iea.org/reports/tracking-transport-2019/rail> (accessed on 25 June 2020).
3. Sustainable Mobility for All (Sum4all). Global Mobility Report 2017. Available online: [https://sustainabledevelopment.un.org/content/documents/2643Global\\_Mobility\\_Report\\_2017.pdf](https://sustainabledevelopment.un.org/content/documents/2643Global_Mobility_Report_2017.pdf) (accessed on 25 June 2020).
4. Tukey, J.W. *Exploratory Data Analysis*; Addison Wesley: Reading, MA, USA, 1977.
5. Charnes, A.; Cooper, W.W.; Rhodes, E. Measuring the efficiency of decision making units. *Eur. J. Oper. Res.* **1978**, *2*, 429–444. [CrossRef]
6. Cooper, W.W.; Seiford, L.M.; Tone, K. *Data Envelopment Analysis*; Springer: New York, NY, USA, 1999.
7. Gennaro, M.D.; Paffumi, E.; Martini, G. Big Data for Supporting Low-Carbon Road Transport Policies in Europe: Applications, Challenges and Opportunities. *Big Data Res.* **2016**, *6*, 11–25. [CrossRef]
8. Zhao, J.; Rahbee, A.; Wilson, N.H.M. Estimating a Rail Passenger Trip Origin-destination Matrix Using Automatic Data Collection Systems. *Comput. Aided Civ. Infrastruct. Eng.* **2007**, *22*, 376–387. [CrossRef]
9. Zhou, J.; Murphy, E.; Long, Y. Commuting efficiency in the Beijing metropolitan area: An exploration combining smartcard and travel survey data. *J. Transp. Geogr.* **2014**, *41*, 175–183. [CrossRef]
10. Tao, S.; Corcoran, J.; Mateo-Babiano, I.; Rohde, D. Exploring bus rapid transit passenger travel behaviour using big data. *Appl. Geogr.* **2014**, *53*, 90–104. [CrossRef]
11. Galland, S.; Knapen, L.; Yasar, A.; Gaud, N.; Janssens, D.; Lamotte, O.; Koukam, A.; Wets, G. Multi-agent simulation of individual mobility behavior in carpooling. *Transp. Res. Part C Emerg. Technol.* **2014**, *45*, 83–98. [CrossRef]
12. Zhang, M.; Wiegman, B.; Tavasszy, L. Optimization of multimodal networks including environmental costs: A model and findings for transport policy. *Comput. Ind.* **2013**, *64*, 136–145. [CrossRef]
13. Guerrero-Ibáñez, J.; Zeadally, S.; Contreras-Castillo, J. Sensor Technologies for Intelligent Transportation Systems. *Sensors* **2018**, *18*, 1212. [CrossRef]
14. Abduljabbar, R.; Dia, H.; Liyanage, S.; Bagloee, S.A. Applications of artificial intelligence in transport: An overview. *Sustainability* **2019**, *11*, 189. [CrossRef]

15. Torre-Bastida, A.I.; Del Ser, J.; Laña, I.; Ilardia, M.; Bilbao, M.N.; Campos-Cordobés, S. Big data for transportation and mobility: Recent advances, trends and challenges. *IET Intell. Transp. Syst.* **2018**, *12*, 742–755. [CrossRef]
16. Gallo, M.; De Luca, G.; D’Acierno, L.; Botte, M. Artificial Neural Networks for Forecasting Passenger Flows on Metro Lines. *Sensors* **2019**, *19*, 3424. [CrossRef] [PubMed]
17. Jiao, P.; Li, R.; Sun, T.; Hou, Z.; Ibrahim, A. Three revised kalman filtering models for short-term rail transit passenger flow prediction. *Math. Probl. Eng.* **2016**, 9717582. [CrossRef]
18. Cai, C.; Yao, E.; Wang, M.; Zhang, Y. Prediction of urban railway station’s entrance and exit passenger flow based on multiply ARIMA model. *J. Beijing Jiaotong Univ.* **2014**, *38*, 135–140.
19. Li, Y.; Wang, X.; Sun, S.; Ma, X.; Lu, G. Forecasting short-term subway passenger flow under special events scenarios using multiscale radial basis function networks. *Transp. Res. Part C Emerg. Technol.* **2017**, *77*, 306–328. [CrossRef]
20. Ling, X.; Huang, Z.; Wang, C.; Zhang, F.; Wang, P. Predicting subway passenger flows under different traffic conditions. *PLoS ONE* **2018**, *13*, e0202707. [CrossRef] [PubMed]
21. Liu, Y.; Liu, Z.; Jia, R. DeepPF: A deep learning based architecture for metro passenger flow prediction. *Transp. Res. Part C Emerg. Technol.* **2019**, *101*, 18–34. [CrossRef]
22. Wang, Y.; Ma, J.; Zhang, J. Metro Passenger Flow Forecast with a Novel Markov-Grey Model. *Period. Polytech. Transp. Eng.* **2019**, *48*, 70–75. [CrossRef]
23. Zhang, B.; Li, S.; Huang, L.; Yang, Y. An improved feedback wavelet neural network for short-term passenger entrance flow prediction in Shanghai subway system. In *Lecture Notes in Computer Science (LNCS), Proceedings of the International Conference on Neural Information Processing (ICONIP), Guangzhou, China, 14–18 November 2017*; Springer: Berlin/Heidelberg, Germany, 2017; Volume 10638, pp. 35–45.
24. Sohn, K.; Shim, H. Factors generating boardings at Metro stations in the Seoul metropolitan area. *Cities* **2010**, *27*, 358–368. [CrossRef]
25. Chin, J.; Callaghan, V.; Lam, I. Understanding and personalising smart city services using machine learning, The Internet-of-Things and Big Data. In *Proceedings of the IEEE 26th International Symposium on Industrial Electronics (ISIE)*, Edinburgh, UK, 18–21 June 2017; pp. 2050–2055.
26. Liu, L.; Chen, R.-C. A novel passenger flow prediction model using deep learning methods. *Transp. Res. Part C Emerg. Technol.* **2017**, *84*, 74–91. [CrossRef]
27. Bartlett, Z.; Han, L.; Nguyen, T.T.; Johnson, P. A Novel Online Dynamic Temporal Context Neural Network Framework for the Prediction of Road Traffic Flow. *IEEE Access* **2019**, *7*, 153533–153541. [CrossRef]
28. Daraio, C.; Diana, M.; Di Costa, F.; Leporelli, C.; Matteucci, G.; Nastasi, A. Efficiency and effectiveness in the urban public transport sector: A critical review with directions for future research. *Eur. J. Oper. Res.* **2016**, *248*, 1–20. [CrossRef]
29. Wanke, P.; Azad, A.K. Efficiency in Asian railways: A comparison between data envelopment analysis approaches. *Transp. Plan. Technol.* **2018**, *41*, 573–599. [CrossRef]
30. Lobo, A.; Couto, A. Technical Efficiency of European Metro Systems: The Effects of Operational Management and Socioeconomic Environment. *Netw. Spat. Econ.* **2016**, *16*, 723–742. [CrossRef]
31. Han, J.; Hayashi, Y. A Data Envelopment Analysis for Evaluating the Performance of China’s Urban Public Transport Systems. *Int. J. Urban Sci.* **2008**, *12*, 173–183. [CrossRef]
32. Kutlar, A.; Kabasakal, A.; Sarikaya, M. Determination of the efficiency of the world railway companies by method of DEA and comparison of their efficiency by Tobit analysis. *Qual. Quant.* **2013**, *47*, 3575–3602. [CrossRef]
33. Tsai, C.H.P.; Mulley, C.; Merkert, R. Measuring the cost efficiency of urban rail systems: An international comparison using DEA and Tobit models. *J. Transp. Econ. Policy* **2015**, *49*, 17–34.
34. Zhang, H.; You, J. An Empirical Study of Transport Efficiency of Urban Rail Transit Based on Data Envelopment Analysis and Tobit Model. *J. Tongji Univ.* **2019**, *46*, 1306–1311.
35. Lao, Y.; Liu, L. Performance evaluation of bus lines with data envelopment analysis and geographic information systems. *Comput. Environ. Urban* **2009**, *33*, 247–255. [CrossRef]
36. Hahn, J.-S.; Kim, H.-R.; Kho, S.-Y. Analysis of the efficiency of Seoul arterial bus routes and its determinant factors. *KSCE J. Civ. Eng.* **2011**, *15*, 1115–1123. [CrossRef]
37. Hahn, J.-S.; Kim, D.-K.; Kim, H.-C.; Lee, C. Efficiency analysis on bus companies in Seoul city using a network DEA model. *KSCE J. Civ. Eng.* **2013**, *17*, 1480–1488. [CrossRef]

38. Caulfield, B.; Bailey, D.; Mullarkey, S. Using Data Envelopment Analysis as a public transport project appraisal tool. *Transp. Policy* **2013**, *29*, 74–85. [CrossRef]
39. Peyrache, A.; Rose, C.; Sicilia, G. Variable selection in Data Envelopment Analysis. *Eur. J. Oper. Res.* **2020**, *282*, 644–659. [CrossRef]
40. Charnes, A.; Cooper, W.W.; Lewin, A.Y.; Seiford, L.M. The DEA Process, Usages, and Interpretations. In *Data Envelopment Analysis: Theory, Methodology, and Applications*; Springer: Dordrecht, The Netherlands, 1994.
41. Goodall, W.; Fishman, T.D.; Bornstein, J.; Bonthron, B. The rise of Mobility as a Service. *Deloitte Rev.* **2017**, *20*, 112–129.
42. Tufféry, S. *Data Mining and Statistics for Decision Making*; Wiley: Chichester, UK, 2011; Volume 2.
43. Ghosh, A.; Nashaat, M.; Miller, J.; Quader, S.; Marston, C. A comprehensive review of tools for exploratory analysis of tabular industrial datasets. *Vis. Inform.* **2018**, *2*, 235–253. [CrossRef]
44. Banker, R.D.; Charnes, A.; Cooper, W.W. Some models for estimating technical and scale inefficiencies in data envelopment analysis. *Manag. Sci.* **1984**, *30*, 1078–1092. [CrossRef]
45. Fare, R.; Grosskopf, S. Modelling undesirable factors in efficiency evaluation: Comment. *Eur. J. Oper. Res.* **2004**, *157*, 242–245. [CrossRef]
46. Seiford, L.M.; Zhu, J. Modeling undesirable factors in efficiency evaluation. *Eur. J. Oper. Res.* **2002**, *142*, 16–20. [CrossRef]
47. Flyvbjerg, B.; Bruzelius, N.; Wee, B.V. Comparison of Capital Costs per Route-Kilometre in Urban Rail. *Eur. J. Trans. Infrastruct. Res.* **2008**, *8*, 17–30.
48. Ruiz-Rosero, J.; Ramirez-Gonzalez, G.; Khanna, R. Masivo: Parallel Simulation Model Based on OpenCL for Massive Public Transportation Systems' Routes. *Electronics* **2019**, *8*, 1501. [CrossRef]
49. Hartmann Tolić, I.; Nyarko, E.K.; Ceder, A.A. Optimization of Public Transport Services to Minimize Passengers' Waiting Times and Maximize Vehicles' Occupancy Ratios. *Electronics* **2020**, *9*, 360. [CrossRef]
50. Stone, M.; Aravopoulou, E. Improving journeys by opening data: The case of Transport for London (TfL). *Bottom Line* **2018**, *31*, 2–15. [CrossRef]
51. Open Street Map. List of London Underground Station. Available online: [https://wiki.openstreetmap.org/wiki/List\\_of\\_London\\_Underground\\_stations](https://wiki.openstreetmap.org/wiki/List_of_London_Underground_stations) (accessed on 25 June 2020).
52. Shuai, H.; Haiying, L. A urban rail transport network carrying capacity calculation method based on the logit model. In Proceedings of the International Conference on Transportation, Mechanical, and Electrical Engineering (TMEE), Changchun, China, 16–18 December 2011; pp. 191–194.
53. TfL Travel in London Edition 12. Available online: <http://content.tfl.gov.uk/travel-in-london-report-12.pdf> (accessed on 25 June 2020).
54. TfL FOI Request Detail. Available online: <https://tfl.gov.uk/corporate/transparency/freedom-of-information/foi-request-detail?referenceId=FOI-0880-1819> (accessed on 25 June 2020).
55. Carbon Independent Calculator Based on AEA Report. Available online: <https://www.carbonindependent.org/20.html> (accessed on 25 June 2020).
56. EC. AEA Report. Handbook of External Costs of Transport. Available online: [https://ec.europa.eu/transport/sites/transport/files/handbook\\_on\\_external\\_costs\\_of\\_transport\\_2014\\_0.pdf](https://ec.europa.eu/transport/sites/transport/files/handbook_on_external_costs_of_transport_2014_0.pdf) (accessed on 25 June 2020).
57. Tom, M.S.; Fischbeck, P.S.; Hendrickson, C.T. Energy Use, Blue Water Footprint, and Greenhouse Gas Emissions for Current Food Consumption Patterns and Dietary Recommendations in the US. *Environ. Syst. Decis.* **2016**, *36*, 92–103. [CrossRef]



© 2020 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).



## Article

# A Data Mining and Analysis Platform for Investment Recommendations

Elena Hernández-Nieves <sup>1,\*</sup> , Javier Parra-Domínguez <sup>1</sup> , Pablo Chamoso <sup>1</sup> , Sara Rodríguez-González <sup>1</sup>   
and Juan M. Corchado <sup>1,2,3,4</sup> 

- <sup>1</sup> BISITE Research Group, University of Salamanca, Edificio Multiusos I+D+i, 37007 Salamanca, Spain; javierparra@usal.es (J.P.-D.); chamoso@usal.es (P.C.); srg@usal.es (S.R.-G.); corchado@usal.es (J.M.C.)
- <sup>2</sup> Air Institute, IoT Digital Innovation Hub, Carbajosa de la Sagrada, 37188 Salamanca, Spain
- <sup>3</sup> Department of Electronics, Information and Communication, Faculty of Engineering, Osaka Institute of Technology, Osaka 535-8585, Japan
- <sup>4</sup> Pusat Komputeran dan Informatik, Universiti Malaysia Kelantan, Bachok 16300, Kelantan, Malaysia
- \* Correspondence: elenahn@usal.es

**Abstract:** This article describes the development of a recommender system to obtain buy/sell signals from the results of technical analyses and of forecasts performed for companies operating in the Spanish continuous market. It has a modular design to facilitate the scalability of the model and the improvement of functionalities. The modules are: analysis and data mining, the forecasting system, the technical analysis module, the recommender system, and the visualization platform. The specification of each module is presented, as well as the dependencies and communication between them. Moreover, the proposal includes a visualization platform for high-level interaction between the user and the recommender system. This platform presents the conclusions that were abstracted from the resulting values.

**Keywords:** artificial intelligence; Big Data analytics; forecasting systems; recommender system; Fintech

**Citation:** Hernández-Nieves, E.; Parra-Domínguez, J.; Chamoso, P.; Rodríguez-González, S.; Corchado, J.M. A Data Mining and Analysis Platform for Investment Recommendations. *Electronics* **2021**, *10*, 859. <https://doi.org/10.3390/electronics10070859>

Academic Editor: Amir Mosavi

Received: 16 March 2021

Accepted: 31 March 2021

Published: 4 April 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Data analysis is a process of inspecting, cleaning, transforming, sorting, and modelling data for the purpose of finding useful information, reaching conclusions, and making appropriate decisions. In statistics, data analysis is divided into descriptive analytics, exploratory analytics, and predictive analytics.

Predictive analytics is defined as the branch of analytics that is used to make predictions regarding future events facing, for example, an organization. To do so, it will use various methods, such as data mining, text mining, artificial intelligence, statistics, or data modelling, among others. In addition, predictive analytics manages information technologies, analysis methods, and business process modelling with the purpose of anticipating future events that may happen to the organization in question.

In this research the focus is on predictive analytics with a specific approach to stock market analysis. It is assumed that a stock market prediction is considered successful if it achieves the best results using the minimum data input and the least complex stock market model [1]. Within the field of Artificial Intelligence, the emergence of Machine Learning and the increasing computing performance have allowed developing new services on the basis of traditional financial products, providing financial-economic instruments that provide higher versatility and greater speed [2]. As Jigar Patel et al. point out in [3], forecasting a stock's value is difficult because of the uncertainty of prediction due to the large number of potential determinants. The authors suggest a method that includes both fundamental and technical analysis combined with Machine Learning algorithms, which is an approach to prediction that tries to improve its efficiency.

In contrast to other research that focuses on a single model for investment recommendation such as Artificial Neural Networks (ANN) or optimised decision trees, in this research, a series of algorithms (Random Forest Regressor, Gradient Boosting Regressor, SVM-LinearSVR, MLP Regressor and kNNNeighbors Regressor) are applied. In addition, technical analysis is used, combining Momentum Indicators and Moving Averages. The proposed recommendation system will remove subjectivity from the process after evaluating and validating the algorithms and will provide the user with the algorithm with the best accuracy. However, the main advantage of the investigated system consists in the possibility of consulting the whole process that the system has carried out (analysis, prediction and investment recommendations).

The research conducted in this study has led to the development of a platform that integrates different modules. The modular approach favours not only the overall research, but it is also good for achieving scalability, flexibility, and usability. The modules that make up the system are:

1. Analysis and data mining. The initial objective was to draw up a document that breaks down the functioning of the Spanish continuous market. The goal of this analysis was to determine the needs to be met by the prediction and recommendation model. Therefore, the market analysis has served as a starting point for the development of the platform. Regarding data extraction, given that the operation of the prediction and recommendation platform is based on a dataset containing the historical data of the companies in the continuous market, it has been necessary to create a system that is in charge of extracting the data in real time. Thus, this implies the need to find a reliable data source that contains the information that is required by the platform. It is possible to either make use of an Application Programming Interface (API) to allow for data retrieval, or to develop a system based on Web Scraping for the extraction and formatting of data. The analysis of the data greatly facilitates the subsequent development of a forecasting and recommendation system and the calculation of the technical analysis factors.
2. The forecasting system. The objective of this system is to predict the closing value of a share in the Spanish continuous market from its opening value on the same day. This minimizes the error of the prediction model as much as possible, which will be presumably based on Machine Learning regression algorithms. The forecasting system will be developed on the basis of the extracted historical data of the shares of the Spanish continuous market companies.
3. The technical analysis. On the basis of the premise that the forecasting system relies on a series of historical market opening values to predict its closing values, the addition of a system that is based on the calculation of technical analysis factors (widely used in economics, specifically in the field of investment) is proposed, in order to combine Artificial Intelligence with the human calculation of technical factors. This brings a distinctive value to the prediction system, which is based on the combination of a series of techniques to determine the recommendation that is to be made to the user.
4. The recommendation system. It is proposed to create a recommender system, which, based on the values that result from the aforementioned objectives, is capable of recommending the decision to buy or sell a share in the Spanish continuous market to the user. Therefore, the recommender system is based on calculating the outputs of the rest of the modules and combining them in order to abstract a decision that benefits the user. Thus, the recommendation system is the most crucial and delicate phase of all the modules that make up the platform.
5. The visualization platform. In addition, we propose the creation of a platform that allows for the visualization of the information, recommendations, and predictions for each company in the Spanish continuous market that the end user wishes to consult. The visualization platform graphs the previously made calculations and predictions, so that the end user can consult how the platform operates.



The article is structured, as follows: Section 2 reviews the existing solutions for forecasting stock ratings, Section 3 considers the proposed system, including the data mining and analysis modules, the prediction system, the technical analysis, and the recommendation system. Section 4 outlines the results of the whole research process. Section 5 covers the discussion and the obtained results, as well as future research.

## 2. State of the Art

Throughout this section, the main contributions made in the field of stock prediction will be reviewed. The review begins with the study by Atsalakis, G.S. et al. in [1] who focused their study on stock forecasting through soft computing techniques. After classifying and processing the sample and applying the type of technique to the fuzzy set, the authors concluded that ANNs (Artificial Neural Networks) and neuro fuzzy models were valid for predicting stock market values. It should be noted that, despite being an exhaustive analysis, their research may be outdated as it was published in the period from 1992 to 2006. Another research that establishes ANNs as the best performing machine learning technique for stock market prediction is the research by Soni, S., in [4]. In the research it is compiled various studies applying machine learning and artificial intelligence techniques.

Beyond the research that proposes ANN as a method for stock prediction, during the review of the state of the art it was observed research that highlighted the need for historical stock market data after reviewing various machine learning techniques for stock prediction [5]. In [3] it was also highlighted that predicting stock market values is challenging due to the lack of certainty, perhaps in relation to the conclusions already drawn in the article discussed above. The authors attribute the lack of certainty to the unpredictability of a changing environment and provided a mixed approach that uses both machine learning algorithms and fundamental and technical analysis.

The research discussed above was the first mixed approach that was identified during the review of the state of the art. Once in this line and seeing that it was perhaps the starting point for the research proposed here, the research of [6] was found. On this occasion, the use of various techniques was focused on integrating collaborative and content-based filtering techniques, where the optimal investment recommendation was given by the investor's preferences, trends, macroeconomic factors, etc. To conclude the review, the research conducted in [7] where a mixed approach is also presented, is analysed. The research concerns the use of a decision tree of technical indicators optimised by GA-SVM. The result is a recommender system capable of detecting stock price fluctuations and suggesting a decision to the investor.

Table 1 summarises the contributions considered in this research.

**Table 1.** State of the art of AI applied to Stock investment recommendations.

References	Approaches
[1]	Stock forecasting through soft computing techniques. The authors concluded that ANNs and neuro fuzzy models were valid for predicting stock market values
[4]	compiled various studies applying machine learning and artificial intelligence techniques
[5]	Highlighted the need for historical stock market data after reviewing various machine learning techniques for stock prediction.
[3]	The author provided a mixed approach that uses both machine learning algorithms and fundamental and technical analysis.
[6]	The authors proposal mixed collaborative and content-based filtering techniques.
[7]	The research concerns the use of a decision tree of technical indicators optimised by GA-SVM.

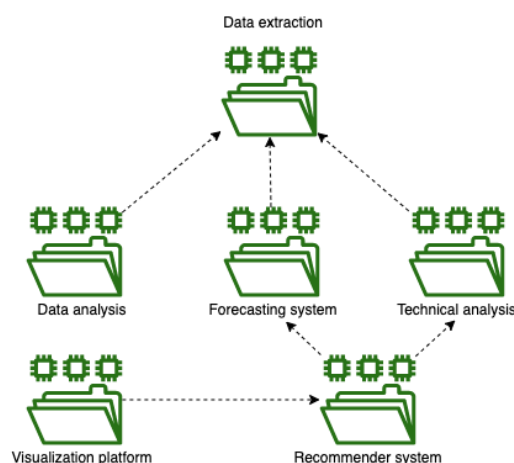
## 3. Proposed Model

Once the state of the art has been reviewed, throughout this section the proposal is presented, more specifically the software architecture that results in the forecasting system. Specifically, the following modules are described and analyzed: data extraction

package, data analysis package, forecasting system module, technical analysis module, recommender system package and the visualization platform. are described.

### 3.1. Software Architecture

This subsection presents the design specification of all the packages that form the software system and how they communicate and interrelate with each other. Figure 1 shows the dependencies between the packages that make up the proposed model, thus showing the different modules that make up the system and, consequently, the interrelationships and dependencies between them.



**Figure 1.** Software system package diagram.

Figure 2 shows the relationships between the use cases and the relationship between the different actors (user and system) and the resulting system. This is intended to provide a clearer and more exemplified understanding of the design specification of each and every one of the modules.

#### 3.1.1. Data Extraction Package

The data extraction module works as follows: first, a user request is received through an API (Application Programming Interface) endpoint or through a request to a method in the package developed in Python for versions 3.x. Secondly, if the system has received that request, it will include in the header the name of the company and the date range (if the historical data has been requested) or it will only extract the name of the company (if the historical data has not been requested).

Figure 3, shows the Python package that has been created for data extraction from Investing.com (after prior authorisation from the company on 28 January 2019). It supports different versions and has been loaded into PyPI (Python Package Indexer). The second significant aspect is continuous integration. The package is monitored, also, unit tests and code coverage are checked (this functionality determines the number of lines of code, identifying unusable lines of code) through Travis CI. In addition, the developed Python package supports more banking products such as funds or ETFs (Exchange Traded Funds), enabling the future implementation of additional functionalities in the platform.

Once the HTML DOM Tree structure has been analyzed to determine which elements of the HTML are to be retrieved and how they can be identified, the development of the Web Scraper begins. Two main steps have been taken:

1. Web request: The HTML of the web have been retrieved and also requests of GET or POST type have been made; principally, urllib3<sup>14</sup> and requests<sup>15</sup>.
2. HTML parsing: consisted in recovering and formatting the data from the previously retrieved HTML. To parse and obtain the information from the HTML, two Python utilities have been employed, BeautifulSoup<sup>16</sup> and lxml<sup>17</sup>.

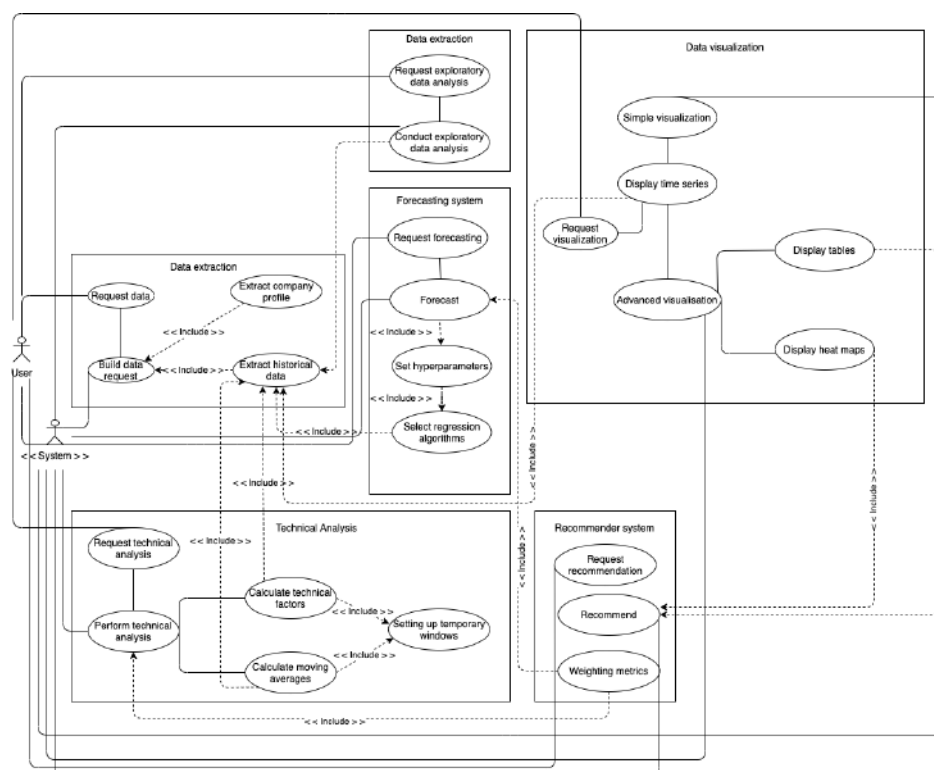


Figure 2. Related Use Case Packages.

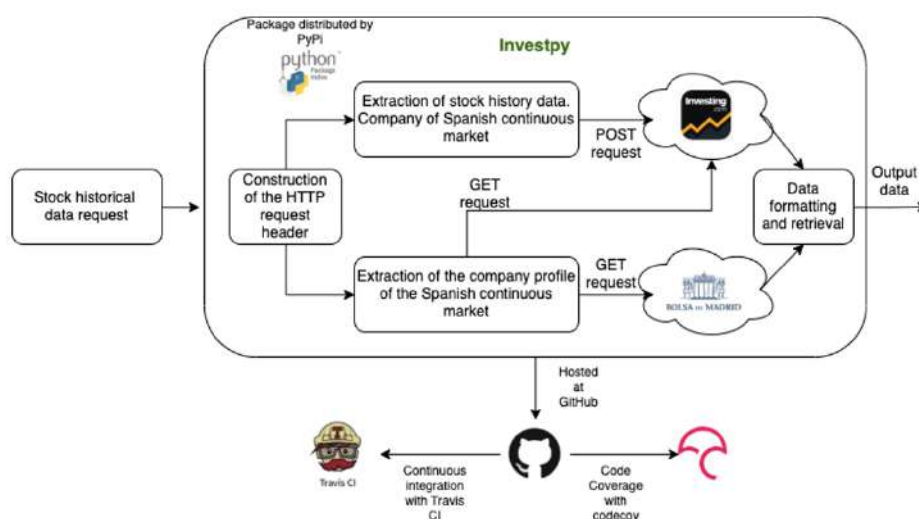
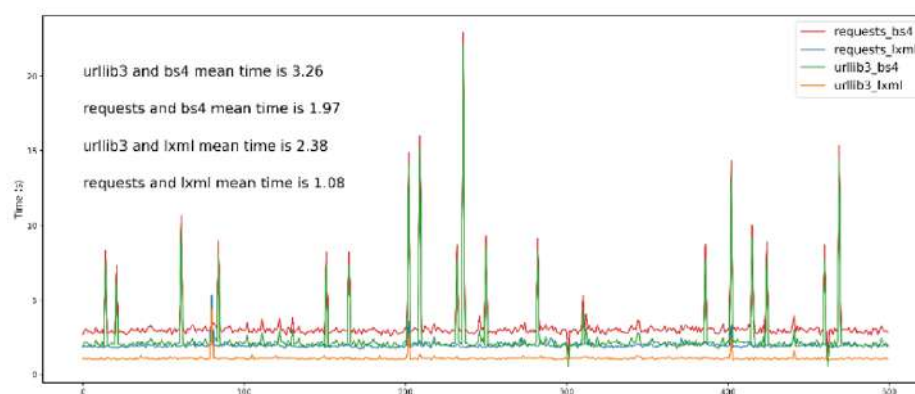


Figure 3. Data extraction Package Design Specification.

Figure 4 shows the graphic representation of the combination of possible Python package times for each of the different phases that are involved in Web Scraping. The combinations are shown in best to worst scaling, as follows: request-lxml, request-bs4, urllib3-lxml and urllib3-bs4. Therefore, to send the request to Investing.com and extract the HTML, either GET or POST type requests are optimal, while lxml is optimal for historical data extraction and parsing.

Finally, the resulting scripts give form to an extensible and open Python package, called investpy [8], intended for data extraction from investing.com. The package facilitates the extraction of data from various financial products, such as: stocks, funds, government bonds, ETFs, certificates, commodities, etc.

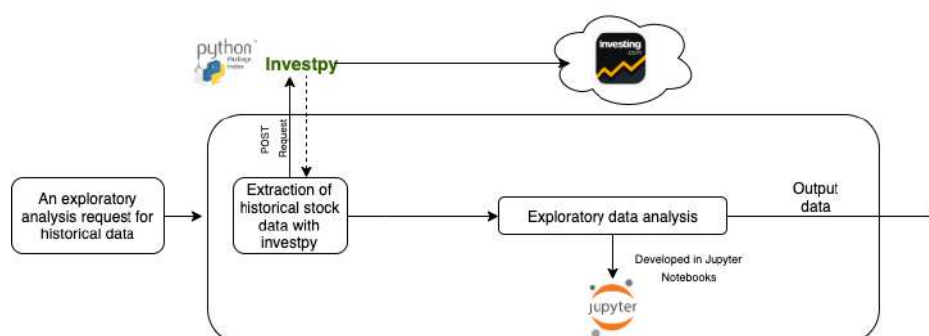


**Figure 4.** Best web scraping combination.

### 3.1.2. Data Analysis Package

Once the historical data for a stock has been extracted, the analysis of the data can be undertaken. All of the packages depend directly or indirectly on the data extraction package, as shown in Figure 1.

Exploratory data analysis is the set of graphical and descriptive tools used for the discovery of data behavior patterns and the establishment of hypotheses with as little structure as possible. Throughout this subsection, the design for the study of the structure of the data and the relationship between them is shown. A representation of how this module operates can be seen in Figure 5.



**Figure 5.** Data Analysis Package Design Specification.

### 3.1.3. Forecasting System

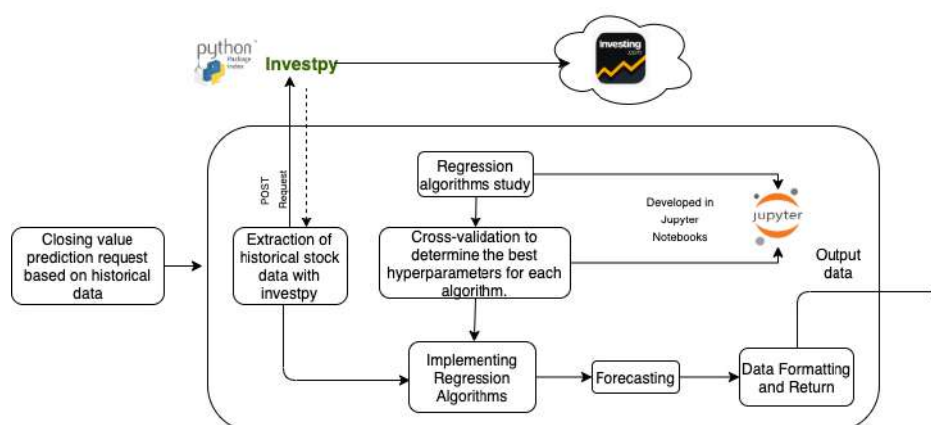
After obtaining the historical data from the last five years of a Spanish continuous market company share through the previously created Python package [8], the Prediction System's design specification is made, as shown in Figure 6.

To predict the future behavior of a stock, Machine Learning regression algorithms [4,9–11] are applied. The objective is to determine the closing price of the stock market, for this the set of opening values has been defined as the input variables and the set of closing values as the output variables, i.e. the closing values are the objective variable of the algorithm. Given the nature of the problem, regression algorithms must be applied. This is because when working with continuous data, regression algorithms can indicate the pattern in a given set of data. Consequently, these algorithms are applied in cases where the relationship between a scalar dependent variable or objective variable  $Y$  and one or more input variables  $X$  is to be modelled. The following section describes the algorithms that were used by the system to predict the last (unknown) closing value based on historical market data, from the last (known) opening value:

1. **Random Forest Regressor:** these algorithms are an automated learning method for classification, regression, and other tasks. A random forest is a meta-stimulus that fits

a series of classification decision trees into various sub-samples of the data set and uses the means to improve productive accuracy and control over fit.

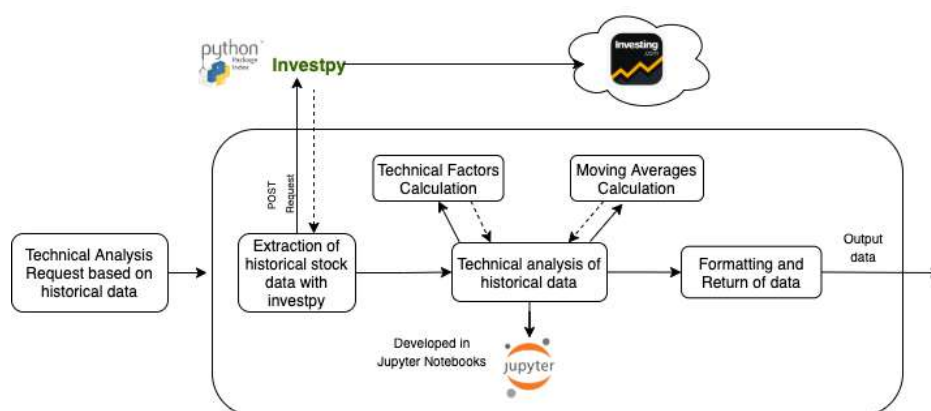
2. Gradient Boosting Regressor: it is an automated learning technique that builds the model in a scenic way, just like methods that are based on reinforcement. It generalizes models allowing for the optimization of an arbitrary and differentiable loss function.
3. SVM-LinearSVR: learning models that analyze data for classification and regression analysis. An SVM training algorithm builds a model that assigns new examples to one or another category, which makes it a non-probabilistic binary linear classifier. In SVR we try to adjust the error within a certain threshold. In this case, it is similar to SVR with the kernel = linear parameter.
4. MLP Regressor: a kind of artificial feedback neural network. MLP uses a supervised learning technique, called backpropagation, for the construction of the network. In addition, its multiple layers and non-linear activation distinguish MLP from a linear perceptron. It also allows for distinguishing data that are not linearly separable.
5. KNNNeighbors Regressor: non-parametric method used for classification and regression



**Figure 6.** Forecasting System Package Design Specification.

### 3.1.4. Technical Analysis

Based on the Spanish continuous market companies' historical stock data, a technical analysis of the market is carried out, in this case combining Momentum Indicators and Moving Averages. This is done for several previously defined time windows for each of the different factors to be calculated based on the standard of the size of the time windows; Figure 7 describes its design specification.



**Figure 7.** Technical Analysis Package Design Specification.

To calculate the factors for the technical analysis, the TA-Lib library has been used through the wrapper written in Python with the same name. Pandas' utilities have been used to calculate the moving averages. Technical Analysis is an analysis that is used

to weigh and evaluate investments. It identifies opportunities to acquire or sell stocks based on market trends. Unlike fundamental analysis, which attempts to determine the exact price of a stock, technical analysis focuses on the detection of trends or patterns in market behavior for the identification of signals to buy or sell assets, along with various graphical representations that help to evaluate the safety or the risk of a stock [12]. This type of analysis can be used in any financial product as long as historical data are available. It is required to include both share prices and volume. Technical analysis is very often employed when a short-term analysis is required, thus, it can help to adequately address the problem presented in this research, where the closing value of a share in a day is predicted. The following indicators are considered in the analysis [13]:

1. Relative Strength Index (RSI): it is a Momentum Indicator (these indicators reflect the difference between the current closing price and the closing price of the previous N days), which measures the impact of frequent changes in the price of a stock, identifying the signs of overbuying or overselling. The representation of the RSI is shown on an oscillator, which is, a line whose value oscillates between two extremes, which, in this case, is between 0 and 100.

$$RSI_{\text{step one}} = 100 - \left[ \frac{100}{1 + \frac{\text{Average gain}}{\text{Average loss}}} \right] \quad (1)$$

2. Stochastic Oscillator (STOCH): it is a Momentum Indicator that compares the closing price of a stock on a given day with the range of closing values of that stock over a certain period of time, defined by the time window. It also allows to adjust the sensitivity of the oscillator either by adjusting the time window or by calculating the moving average of the STOCH result. Like RSI, it identifies the signals of over-bought or oversold stock within a range of 0 to 100 possible values.

$$\%K = 100 - \left( \frac{C - L14}{H14 - L14} \right) \times 100 \quad (2)$$

where C is the most recent closing price, L14 is the lowest price traded of the 14 previous trading sessions, H14 is the highest price traded during the same 14-day period, and %K is the current value of the stochastic indicator.

3. Ultimate Oscillator (ULTOSC): it is a Momentum Indicator used to measure the evolution of a stock over a series of time frames using a weighted average of three different windows or time frameworks. Therefore, it acquires a lower volatility and identifies fewer buy-sell signals than other oscillators that only depend on a single time frame. When the lines generated by ULTOSC diverge from the closing values of a stock, buy and sell signals are identified for it.

$$UO = \left[ \frac{(A_7 \times 4) + (A_{14} \times 2) + A_{28}}{4 + 2 + 1} \right] \times 100 \quad (3)$$

where UO is the Ultimate Oscillator and A is the average. The average calculation follows the next formulas.

$$A_7 = \left[ \frac{\sum_{p=1}^7 BP}{\sum_{p=1}^7 TR} \right] \quad (4)$$

$$A_{14} = \left[ \frac{\sum_{p=1}^{14} BP}{\sum_{p=1}^{14} TR} \right] \quad (5)$$

$$A_{28} = \left[ \frac{\sum_{p=1}^{28} BP}{\sum_{p=1}^{28} TR} \right] \quad (6)$$

where BP is the Buying Pressure and PC is the Prior Close

$$BP = \text{Close} - \text{Min}(\text{Low}, \text{PC}) \quad (7)$$

where TR is the True Range

$$TR = \text{Max}(\text{High}, \text{Prior Close}) - \text{Min}(\text{Low}, \text{Prior Close}) \quad (8)$$

where TR is the True Range

4. Williams %R (WILLR): also known as the Williams Percent Range, is a Momentum Indicator that fluctuates between  $-100$  and  $0$  and measures and identifies levels of stock overbuying or overselling. WILLR is very similar to the STOCH in its use, as it is used for the same purpose. This indicator compares the closing value of a stock with the range between the maximum and minimum values within a given time frame.

$$\text{Williams\%K} = \frac{\text{Highest High} - \text{Close}}{\text{Highest High} - \text{Lowest Low}} \quad (9)$$

where the Highest High is the highest price in the look-back period, typically 14 days, Close is the most recent closing price, and Lowest Low is the lowest price in the look-back period, typically 14 days.

Moving averages are also used in Technical Analysis, as it also represents the Momentum or value change in a timeframe  $N$ . Hence, moving averages help to understand the market trend and, like Momentum Indicators, allow to identify buy and sell signals from the historical data of a stock in a previously mentioned timeframe  $N$ . In this research, we have applied the simple moving average (SMA) and the exponential moving average (EMA) for timeframes of 5, 10, 20, 50, 100, and 200 days, so there will be indicators in different periods.

1. Simple Moving Average (SMA): it is an arithmetic moving average. It is calculated by adding the recent closing values of an action for a window of size  $N$  and dividing that sum by the size of the window. Thus, when the size of the timeframe  $N$  is low, it responds quickly to changes in the value of the stock; if the size of the window  $N$  is high, it responds more slowly.

$$\text{SMA} = \frac{A_1 + A_2 + \dots + A_n}{n} \quad (10)$$

where  $A_n$  the price of an asset at period  $n$  and  $n$  is the number of total periods.

2. Exponential Moving Average (EMA): also called Exponentially Weighted Moving Average, since it weights recent observations, i.e., closing prices of a stock closer to the current one. It can be said that EMAs respond better than SMAs to recent changes in a share's price.

$$\text{EMA}_{\text{Today}} = \left( \text{Value}_{\text{Today}} \times \left( \frac{\text{Smoothing}}{1 + \text{Days}} \right) \right) + \left( \text{EMA}_{\text{Yesterday}} \times \left( \frac{\text{Smoothing}}{1 + \text{Days}} \right) \right) \quad (11)$$

where EMA is the exponential moving average. The smoothing factor is calculated, as follows:

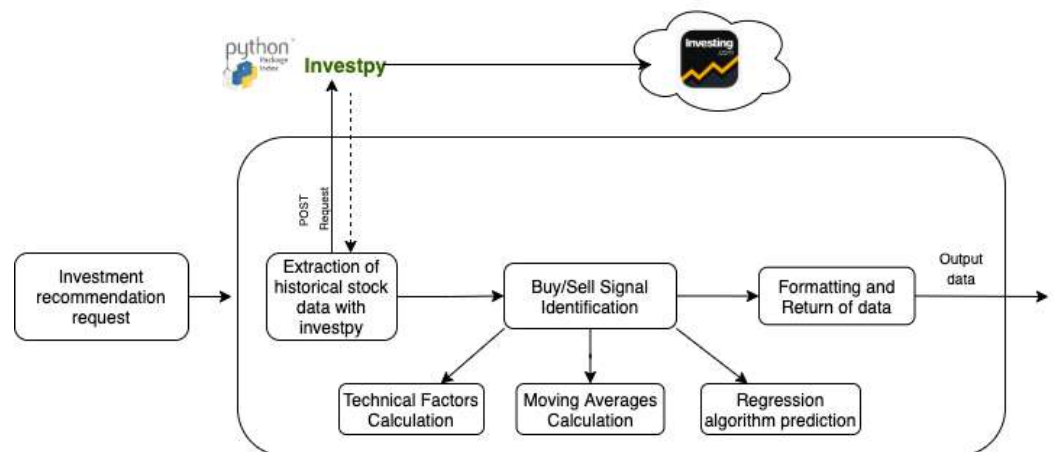
$$\text{Smoothing} = \frac{2}{n + 1} \quad (12)$$

where  $n$  represents the number of periods the EMA uses.

Because both the algorithmic predictions and the results of the technical factor and moving average calculations result in the next closing value of a stock, the recommendation is based on identifying buy and sell signals based on the comparison of the predicted value with the value that the stock has at the current time.

### 3.1.5. Recommender System

Based on the results that are obtained from the forecasting and technical analysis systems, the Recommendation System design specification proceeds, in which the obtained results are weighted to identify buy/sell signals in order to be able to make a recommendation. Figure 8 shows the functionality of the Recommender System and, consequently, the process of signal extraction and the dependencies/relationships between it and the rest of the modules on which it depends.



**Figure 8.** Recommender System Package Design Specification.

The package design proposes the creation of a neutral system, which, based on the analysis of buy/sell signals, determines the action to be taken for/with a stock. This is intended to eliminate the burden of subjectivity.

In addition to the calculation of moving averages and technical analysis ratios, an analysis using regression algorithms is also included, as can be seen in Figure 8. Regression algorithms are used when a prediction is to be made on a continuous dataset. This is the case with the historical time series data of a stock. The output of the algorithm is a quantity that can be measured in a flexible way, depending on the inputs that are passed to the algorithm. Sorting algorithms would be limited to a set of labels.

Linear regression can be defined as an approach to modelling the relationship between a dependent scalar variable  $y$ , and one or more explanatory variables, named  $x$ . Mathematically, it is expressed in the form that is presented in Equation (13).

$$y_i = \beta_0 + \beta_1 x_1 \quad (13)$$

where the variable to be predicted  $y_i$  is distinguished, as well as the constant  $\beta_0$ , the slope  $\beta_1$ , and the input variable  $x_i$ . In the current scenario, given that a stock's historical data set is available, the explanatory variable  $x$  gives the market opening values, and the target variable  $y$  gives the market closing values. Thus, the model input is  $x$  and the expected output  $y$ , where  $y$  is the dependent variable and  $x$  the independent variable, so that the market opening value conditions the market closing value.

The machine learning algorithms that are used by the system to predict the last (unknown) closing value, based on historical market data, from the last (known) opening value are:

1. Random Forest Regressor
2. Gradient Boosting Regressor
3. SVM-LinearSVR
4. MLP Regressor
5. KNeighbors Regressor

These algorithms are applied using sklearn library (Python library that compiles machine learning algorithms) by means of cross-validation (a technique that is used to



evaluate the results of a statistical analysis and ensure that they are independent of the partition between train and test data). In this way, the best hyperparameters of the algorithms can be determined and, thus, the best combination can be identified.

The results have generated a series of heat maps have been generated. The accuracy of the algorithm is represented by the hue. Lighter shading corresponds to a worse result, so darker areas indicate that the resulting hyperparameter combination is better. However, sometimes it is not known which is the best combination because the shades are very similar. The results of these combinations for each action are stored in a JSON file that will be used later by the platform when applying the models of the action prediction system. In this way, the result of applying the cross-validation of the hyperparameters to all of the stocks in the Spanish continuous market is a data file with the best hyperparameters and they are shown in the respective heat maps, to justify the decision. This allows for a more accurate decision to be made, as the user can compare the effectiveness of some hyperparameters against others.

Figures 9–13 present the heat maps resulting from the cross-validation of BBVA's historical data from the last five years. These figures represent the accuracy of the model with each combination of hyperparameters for the algorithms: MLPRegressor, SVMLinearSVR, RandomForestRegressor, GradientBoostingRegressor, and KNNNeighborsRegressor.

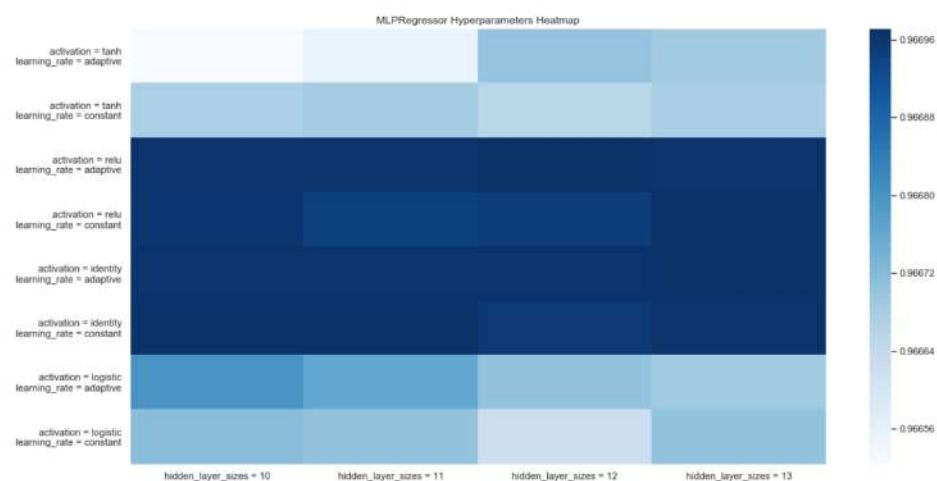


Figure 9. MLPRegressor Hyperparameter Heat Map.

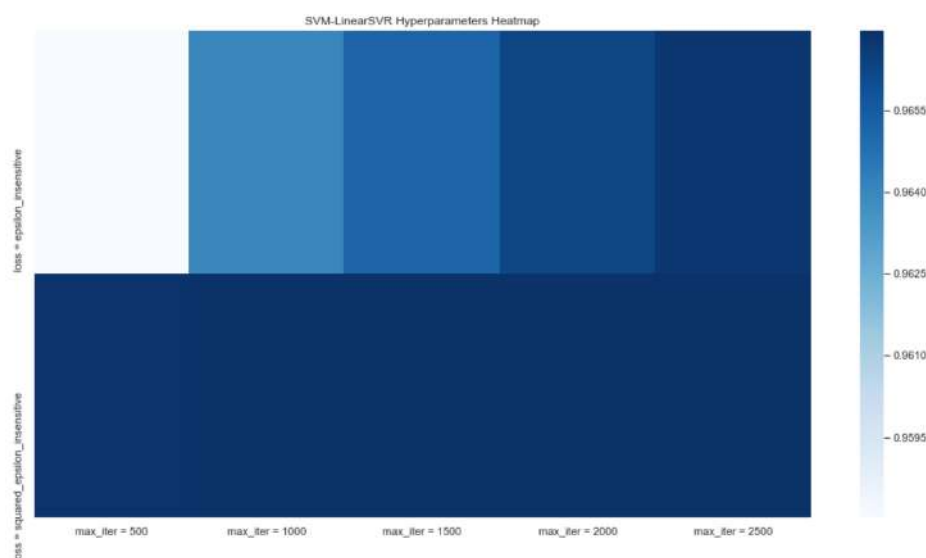


Figure 10. SVM-LinearSVR Hyperparameter Heat Map.



Figure 11. RandomForestRegressor Hyperparameter Heat Map.

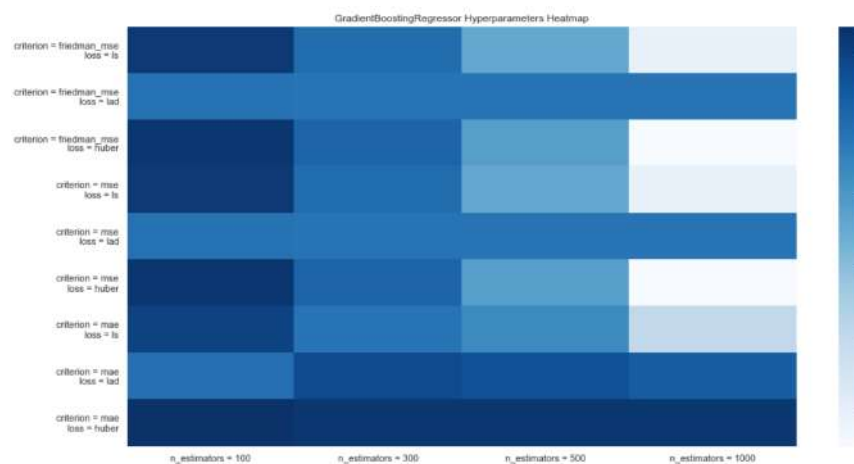


Figure 12. GradientBoostingRegressor Hyperparameter Heat Map.

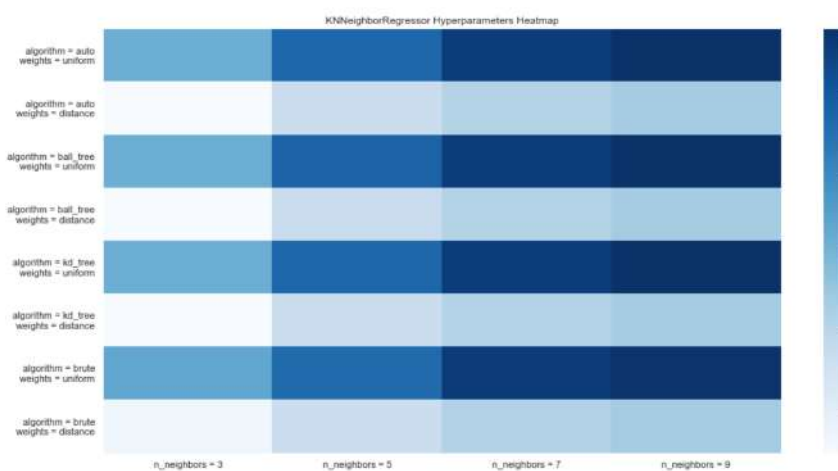


Figure 13. KNNNeighborsRegressor Hyperparameter Heat Map.

Once the evaluation process and the cross-validation of the algorithms have been completed, a graph showing the best algorithm has been drawn up. The top-25 equities of the Spanish Continuous Stock Market have been chosen. Figure 14 shows that the algorithms that best fit the proposed problem are SVM-LinearSVR and MLP Regressor.

Thus, these algorithms are the ones that have the highest accuracy after being trained and tested with an 80/20 split of the dataset.

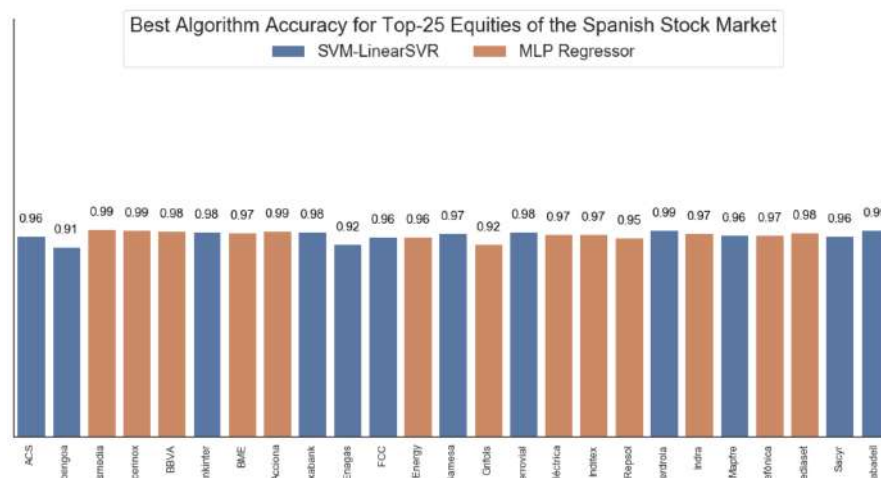


Figure 14. Best algorithm accuracy for top-25 equities of the Spanish Continuous Stock Market.

#### 4. Platform Visualization

Finally, after detailing each of the software modules created, the description of the visualization platform follows as a deduction of the integration of the most relevant aspects of the rest of the modules, so that the end user can interact with the platform.

Therefore, the visualization platform takes up the conclusions of the research carried out in the rest of the modules, so that only two options are given at user level for visualization, either an overview as a result of the exploratory analysis of the data, or the result of the underlying recommender system.

In this way, the different phases or tools that are used for the development of the platform architecture are detailed, based on the results of the study of the rest of the modules.

It is worth mentioning that the development of a visualization platform only aims to bring the results closer to the user, without being the central part of the proposed system.

The result is a platform that provides a user interface for both data visualization, analysis, prediction, and investment recommendation. It has been determined that the platform will be developed using Django, as shown in Figure 15. Django is a Python framework for creating web services, in this case it has been used to communicate the back-end with the front-end. The web application created combines the use of Python for data management and communication, and HTML, CSS, and JavaScript for the visualization of both the platform and the data.

The design pattern used, called MVC (Model-View-Controller), focuses on the division of the web project according to the functionalities of each of its parts. However, Django does not use the MVC pattern, but rather the MVT (Model-View-Template), which is an abstraction of the MVC model. It is worth mentioning that Django works with templates, not with views, being oriented to the development of web applications, as explained in [14], where the author not only teaches aspects of using of Django, but it also lists the different design patterns that can be followed in order to structure to follow the web application to be developed.

The platform's objective is not only to be usable and intuitive, but also to enable any user, whether an expert or not in the stock market, to abstract their own conclusions from the data and evaluate the information analyzed by the system. The created platform completely depends on the Python package developed for data extraction: investpy. The web platform initially shows a screen where the overview option is given on one side and the overview and recommendation option on the other (Figure 16). The overview functionality covers the extraction and basic visualization of the data. The system retrieves

the company profile and the historical data for the last five years of the stock. On the basis of those data, it produces a series of representations:

1. Time series: offers a graphic representation of the retrieved historical data, where the X and Y axes represent the value of the stock in euros, and the date on which the stock reached that value, respectively.
2. Candlestick chart: this representation shows the opening and closing values for each date and the difference between the maximum and minimum values for the same date.
3. Data table: represents the available values. They are called OHLC (Open-High-Low-Close).

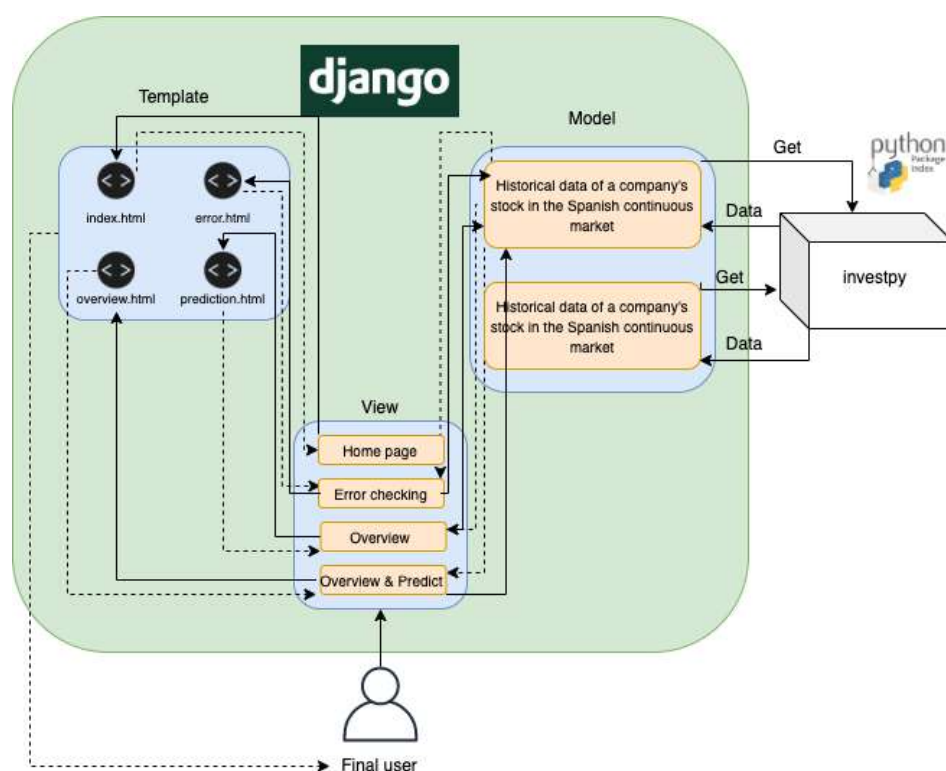


Figure 15. Django Design Architecture.



Figure 16. Main view of the web platform.

The Overview & Recommendation functionality is the same as the user input check, in that it also extracts the company profile and historical data. However, this functionality also includes technical factors and moving averages with the consequent buy/sell recommendation. The generated graphs are visualized on the platform, among them are graphs that compare the different algorithms that the system has applied to make the prediction. This

enables the user to identify those that have had a better precision. The platform presents the conclusions abstracted from the resulting values. It shows the buy/sell recommendation that is based on those values. The process of prediction and recommendation made by the system is transparent to the user.

The novelties that are presented by the module are the graphs generated, in which a comparison between the different algorithms applied by the system to make the prediction can be observed, thus being able to contrast which one has had the best accuracy (Figure 17).

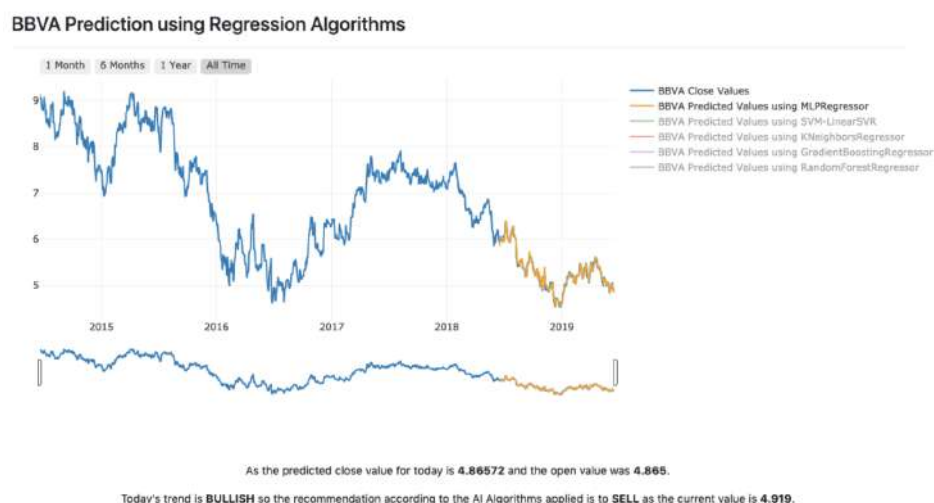


Figure 17. Prediction using regression algorithms.

Additionally, there is an option that allows the user to observe which algorithms have been applied, what they consist of, and which hyperparameters have been used based on the results in the form of a heat map of the cross validation carried out by the system.

Once the justification of the regression algorithms used in the platform by the system has been shown, the results of the different algorithms applied are displayed, where the “best” algorithm (the one with better precision than the rest) is the one that shows its results by default (Figure 18). Even so, the platform gives the option of displaying different time windows and visualising the results of all the algorithms. Finally, the platform displays a paragraph, in which it indicates the conclusions drawn from the study of the values resulting from the prediction and, therefore, shows the buy/sell recommendation based on these values.

Therefore, the platform displays the recommendations based on the results of the prediction, which it will combine with the results of the financial technical analysis, which includes the calculation of moving averages and technical factors. Finally, the system calculates the technical factors, called Momentum Indicators, which indicate the market trend based on calculations taking different time windows (Figure 19).

In this way, the system not only makes the recommendation, but it also supports this recommendation and each of the predictions and calculations that give rise to it, with the data used throughout the process. Therefore, the prediction and recommendation process carried out by the system is transparent to the user at a technical level, so that the user is aware of what has happened in each of the stages of the process, being able to trust that the prediction has not been altered for the benefit of third parties, for example.

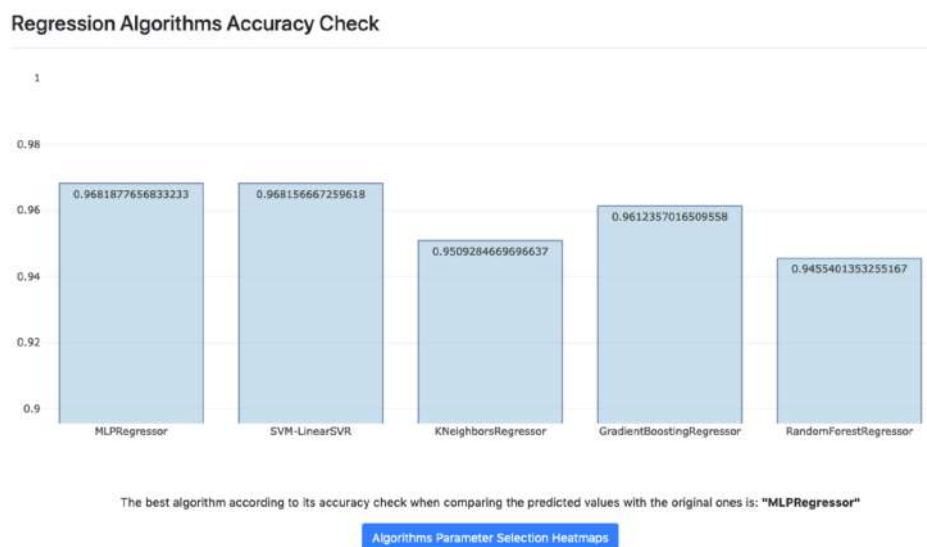


Figure 18. Regression Algorithms Accuracy Check.

#### • Technical Factors

Tech. Factor	Result	Signal
RSI - 14 days	43.63047	SELL
STOCH - 9,6 days	38.29643	SELL
ULTOSC - 7,14,28 days	53.80424	BUY
WILLR - 14 days	-55.2381	SELL

[Back](#)

Figure 19. Buy/sell recommendations.

## 5. Discussion and Results

The conducted research provides an initial approach to data analysis and the combined use of Machine Learning algorithms and techniques, with traditional market analysis. Their use enables the proposed platform to arrive at conclusions regarding future market behavior. Thus, it can be concluded that, when Machine Learning algorithms are trained with a sufficiently large amount of data, it is possible to successfully predict the closing value on the basis of the current opening value of the market. Thus, after identifying buy and sell signals, it has been possible to create a system that recommends the user to buy, hold, or sell a stock at a certain time of day, according to the prediction obtained by the regression algorithms.

Although the recommender system operates well and meets the initial objectives of this study, system extensions will be considered in future research. The breadth of the platform in terms of functionalities was the most significant complication that arose during the research, therefore it was decided to approach it with a modular architecture. Thanks to the modular, highly scalable design it is possible to provide the system with more functionalities; the combination of Natural Language Processing (NLP) techniques could be used in an opinion mining process, the recommender system will be able to abstract the future market trend based on the sentiment analysis. In addition, the use of NLP techniques is also proposed for the classification of companies into sectors based on their company profiles, thus being able to group companies into sectors based on the description that each company in the Spanish continuous market proposes. Therefore, additional functionality must be added to the Python package that was created for the extraction of the Investing

data, called investpy. The enhancement will consist of retrieving all the data provided freely by Investing.com. Additionally, a study of the algorithms applied to other markets should be carried out, as the proposed system is oriented towards a very specific market; the Spanish continuous market. It will be necessary to carry out a study to determine the best algorithms for the stock markets of each of the countries to be incorporated. It is considered to be viable given that all historical stock data previously go through a GridSearchCV, which consists in cross validating the optimal hyperparameters to be used by an algorithm from a specific dataset. In addition, further research is considered on event identifications that can be used to better choose the operation performed (buy/sell) and the social characteristics of the different communities [15,16].

**Author Contributions:** Conceptualization, E.H.-N., J.P.-D. and P.C.; Funding acquisition, E.H.-N.; Investigation, E.H.-N., J.P.-D. and P.C.; Methodology, S.R.G. and J.M.C.; Software, P.C.; Supervision, S.R.-G. and J.M.C.; Writing—original draft, E.H.-N. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was partially Supported by the project “Computación cuántica, virtualización de red, edge computing y registro distribuido para la inteligencia artificial del futuro”, Reference: CCTT3/20/SA/0001, financed by Institute for Business Competitiveness of Castilla y León, and the European Regional Development Fund (FEDER). The research of Elena Hernández-Nieves is funded by Ministry of Education of the Junta de Castilla y León and the European Social Fund grant number EDU/556/2019.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Atsalakis, G.S.; Valavanis, K.P. Surveying stock market forecasting techniques—Part II: Soft computing methods. *Expert Syst. Appl.* **2009**, *36*, 5932–5941. [CrossRef]
2. Nicoletti, B.; Nicoletti, W.; Weis. *Future of FinTech*; Palgrave Studies in Financial Services Technology: Rome, Italy, 2017.
3. Patel, J.; Shah, S.; Thakkar, P.; Kotecha, K. Predicting stock and stock price index movement using trend deterministic data preparation and machine learning techniques. *Expert Syst. Appl.* **2015**, *42*, 259–268. [CrossRef]
4. Soni, S. Applications of ANNs in stock market prediction: A survey. *Int. J. Comput. Sci. Eng. Technol.* **2011**, *2*, 71–83.
5. Yoo, P.D.; Kim, M.H.; Jan, T. Machine learning techniques and use of event information for stock market prediction: A survey and evaluation. In Proceedings of the International Conference on Computational Intelligence for Modelling, Control and Automation and International Conference on Intelligent Agents, Web Technologies and Internet Commerce (CIMCA-IAWTIC’06), Vienna, Austria, 28–30 November 2005; Volume 2, pp. 835–841.
6. Taghavi, M.; Bakhtiyari, K.; Scavino, E. Agent-based computational investing recommender system. In Proceedings of the 7th ACM Conference on Recommender Systems, Hong Kong, China, 12–16 October 2013; pp. 455–458.
7. Nair, B.B.; Mohandas, V. An intelligent recommender system for stock trading. *Intell. Decis. Technol.* **2015**, *9*, 243–269. [CrossRef]
8. Del Canto, A. Investpy—Financial Data Extraction from Investing.com with Python. Available online: <https://github.com/alvarobartt/investpy> (accessed on 4 April 2020).
9. Arora, N. Financial analysis: Stock market prediction using deep learning algorithms. In Proceedings of the International Conference on Sustainable Computing in Science, Technology and Management (SUSCOM), Amity University Rajasthan, Jaipur, India, 26–28 February 2019.
10. Khaidem, L.; Saha, S.; Dey, S.R. Predicting the direction of stock market prices using random forest. *arXiv* **2016**, arXiv:1605.00003.
11. Pimprikar, R.; Ramachandran, S.; Senthilkumar, K. Use of machine learning algorithms and twitter sentiment analysis for stock market prediction. *Int. J. Pure. Appl. Math.* **2017**, *115*, 521–526.
12. Edwards, R.D.; Magee, J.; Bassetti, W.C. *Technical Analysis of Stock Trends*; Routledge, Taylor & Francis Group: New York, NY, USA, 2018.
13. Dash, R.; Dash, P.K. A hybrid stock trading framework integrating technical analysis with machine learning techniques. *J. Financ. Data Sci.* **2016**, *2*, 42–57. [CrossRef]
14. Ravindran, A. *Django Design Patterns and Best Practices*; Packt Publishing Ltd.: Birmingham, UK, 2015.
15. Di Girolamo, R.; Esposito, C.; Moscato, V.; Sperli, G. Evolutionary game theoretical on-line event detection over tweet streams. *Knowl. Based Syst.* **2021**, *211*, 106563. [CrossRef]
16. Mercorio, F.; Mezzanzanica, M.; Moscato, V.; Picariello, A.; Sperli, G. DICO: A graph-db framework for community detection on big scholarly data. *IEEE Trans. Emerg. Top. Comput.* **2019**. [CrossRef]





## Article

# A Mathematical Study of Barcelona Metro Network

Irene Mariñas-Collado <sup>1</sup>, Elisa Frutos Bernal <sup>2</sup>, Maria Teresa Santos Martin <sup>3</sup>, Angel Martín del Rey <sup>4,\*</sup>, Roberto Casado Vara <sup>5</sup> and Ana Belen Gil-González <sup>5</sup>

- <sup>1</sup> Department of Statistics and Operations Research and Mathematics Didactics, University of Oviedo, 33007 Oviedo, Spain; marinasirene@uniovi.es  
<sup>2</sup> Department of Statistics, University of Salamanca, 37007 Salamanca, Spain; efb@usal.es  
<sup>3</sup> Department of Statistics, Institute of Fundamental Physics and Mathematics, University of Salamanca, 37007 Salamanca, Spain; maysam@usal.es  
<sup>4</sup> Department of Applied Mathematics, Institute of Fundamental Physics and Mathematics, University of Salamanca, 37007 Salamanca, Spain  
<sup>5</sup> BISITE Research Group, University of Salamanca, 37007 Salamanca, Spain; rober@usal.es (R.C.V.); abg@usal.es (A.B.G.-G.)  
\* Correspondence: delrey@usal.es

**Abstract:** The knowledge of the topological structure and the automatic fare collection systems in urban public transport produce many data that need to be adequately analyzed, processed and presented. These data provide a powerful tool to improve the quality of transport services and plan ahead. This paper aims at studying, from a mathematical and statistical point of view, the Barcelona metro network; specifically: (1) the structural and robustness characteristics of the transportation network are computed and analyzed considering the complex network analysis; and (2) the common characteristics of the different subway stations of Barcelona, based on the passenger hourly entries, are identified through hierarchical clustering analysis. These results will be of great help in planning and restructuring transport to cope with the new social conditions, after the pandemic.

**Keywords:** complex network analysis; centrality measures; network robustness; ridership patterns; clustering analysis; passenger flow; Barcelona underground

**Citation:** Mariñas-Collado, I.; Frutos Bernal, E.; Santos Martín, M.T.; Martín del Rey, A.; Casado Vara, R.; Gil-González, A.B. A Mathematical Study of Barcelona Metro Network. *Electronics* **2021**, *10*, 557. <https://doi.org/10.3390/electronics10050557>

Academic Editor: Myung-Sup Kim

Received: 29 January 2021

Accepted: 22 February 2021

Published: 27 February 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Sustainable urban mobility is one of the most distinct characteristics of Smart Cities. Specifically, intelligent public urban transport planning plays an important role in the design of the future cities and in the sustainable development of the environment (in this sense, it has become one of the most powerful tools in the fight against air pollution in cities); moreover, it is well known that efficient mass transit systems have a highly beneficial impact on economic development and social integration. Particularly, the subway is the best choice in big cities since it exhibits many advantages including reducing traffic congestion, saving energy and non-renewable resources, reducing the number of traffic accidents and therefore deaths, large capacity, time reliability, etc. [1].

Hundreds of millions of passengers commute in public transport daily in large cities, hence failures in the network can cause major problems to commuters and business activities with significant economic and social losses. In addition, the COVID-19 pandemic has changed the security measures on the transport network in order to maintain the sanitary requirements. Proper social distancing between passengers is hard to ensure in public transport if it is not well planned (taking into account the different characteristics of the different stations and lines). To avoid overcrowded stations and trains, it is crucial to know transit trip patterns. This will also allow better network planning, demand forecasting and, ultimately, a more effective use of the available resources in general.

Two main goals are addressed in this work: (1) study the structural and robustness characteristics of Barcelona subway network; and (2) identify ridership patterns at its

stations. In the first case, the basic techniques of Complex Network Analysis are used (centrality measures, structural indices, robustness coefficients, etc.), whereas, in the second case, a hierarchical cluster analysis is performed to group stations according to their boarding patterns. Barcelona's metro is Spain's second largest city subway system: there are a total of 13 lines and 151 stations in the network. Its length is 119 km, and during 2018 more than 400 million people used it.

In recent years, the complex network approach has been used to analyze the subway rail networks of several cities around the world. Since 2002, when Latora and Marchiori studied the topological properties of the Boston subway [2], many other works have appeared. Lu and Shi found that the public transportation network in China had scale-free and small world characteristics [3]. Zhang et al. studied the topological characteristics of some subway networks around the world and investigated network failures to discuss the vulnerability of these subway networks [4]. Liu and Song [5] studied the topology of Guangzhou subway network using L-space method, and the value and distribution of the network's degree, clustering coefficient and average shortest path length were computed and analyzed. Cats [6] conducted a longitudinal analysis of the topological evolution of a multimodal rail network by investigating the dynamics of its topology for the case of Stockholm during 1950–2025.

The robustness of subway networks has also been discussed by many other researches. For example, Derrible and Kennedy studied the complexity and robustness of 33 metro networks [7]. Using network science and graph theory, ten theoretical and four numerical robustness metrics and their performance in quantifying the robustness of metro networks under random failures or targeted attacks were investigated by Wang et al. [8]. Zhang et al. [9] investigated the connectivity, robustness and reliability of the Shanghai subway network of China. Forero-Ortiz et al. [10] gave insights for stakeholders and policymakers to enhance urban flood risk management, as a reasonable approach to tackle this issue for Metro systems worldwide. De Bona et al. [11] proposed a novel methodology called Reduced Model as a simple method of network reduction that preserves the network skeleton (backbone structure) by properly removing 2-degree nodes of weighted and unweighted network representations. In [12], a new perspective for understanding vulnerability of metro networks is shown with the aims of improving operation reliability and stability of the network, designing emergency strategies to protect the network, etc.

In this work, the topological characteristics of the metro network are investigated considering the complex network approach. Specifically a brief analysis of the Barcelona subway network is provided from the computation of the most important centrality measures: (i) degree centrality  $C_D$ ; (ii) average degree  $E[D]$ ; (iii) degree distribution  $p(k)$ ; (iv) average path length  $L$ ; (v) closeness centrality  $C_{CL}$ ; and (vi) betweenness centrality  $C_B$ . In addition, to assess the robustness of the subway network, eight theoretical robustness metrics are investigated: (i) normalized robustness indicator  $\bar{r}^T$ ; (ii) effective graph conductance  $C_G$ ; (iii) average efficiency  $E[\frac{1}{H}]$ ; (iv) clustering coefficient  $C_{CG}$ ; (v) normalized algebraic connectivity  $\bar{\mu}_{N-1}$ ; (vi) average degree  $E[D]$ , (vii) normalized natural connectivity  $\bar{\lambda}$ ; and (viii) degree diversity  $\kappa$ .

Most public transit networks use automated fare collection (AFC) systems. The interest in this kind of technology is because it is perceived as a secure method of user validation and fare payment. Moreover, it improves the quality of the data, gives transit a more modern look and provides new opportunities for innovative and flexible fare structuring [13]. While the main purpose of AFC systems is to collect revenue, they also produce very large quantities of very detailed data of on-board transactions. These data are very useful to transit planners, from the day-to-day operation of the transit system to the strategic long-term planning of the network [14].

AFC systems are classified into two types according to the fare charge mode of transit: flat-rate fare systems and distance-based fare systems. In flat-rate fare systems, only entry swipes are registered, while, in distance-based fare systems, entry and exit swipes are registered. Barcelona metro uses a flat-rate fare system, therefore only metro boarding is

available in this study. This has the inconvenience of not knowing where the passenger's journey ends, e.g., the trip's purpose. The destination of the trip helps understand peak hours. For instance, most of the work and education trips start in the morning peak from home and return back to home in the evening peak. While not within the scope of this paper, the destination estimation of public transport is one of the major concerns for the implementation of smart card data and there exist several approaches (see, e.g., [15–18]).

Every day, depending on the size of the network, millions of transactions are registered by the AFC systems, which can be used to analyze human mobility. It has been determined that human trajectories and trips generated with human mobility show a high degree of temporal and spatial regularity [19]. Passenger flow of the urban subway varies according to time and space, including working days, holidays, seasons, residential areas, business centers, workplaces and other factors such as weather, as well as other forms of transportation that connect to the subway network. In this regard, several methods have been developed in the literature for this type of analysis, most using clustering approaches [20].

Two viewpoints can be considered when a cluster analysis using smart card data is performed. The first one clusters stations based on the temporal-spatial distribution characteristics of subway ridership. The second one identifies groups of passengers that have similar boarding times aggregated into weekly profiles [21].

From the first point of view, Chen et al. [22] studied the diurnal pattern of subway ridership in New York City using the k-means algorithm. Wang et al. [23] analyzed eight metro stations in the central area of Hong Kong using the hierarchical cluster analysis. The k-means algorithm was also employed by Kim et al. [24] to identify the daily travel patterns at subway stations of Seoul Capital Area. Ding et al. [25] applied gradient boosting decision trees to investigate the non-linear effects of built environment variables on station boarding in the Washington metropolitan area. Langlois et al. [26] proposed a longitudinal representation of user's multi-week activity and identified 11 travel patterns from London's public transport network.

The study and analysis of different characteristics of subway networks have been tackled by means of other different paradigms. For example, risk analysis has been addressed in some recent works (see, e.g., [10,27–29]), the GIS-based technologies improves the analysis performed using mathematical methods [30], modern statistical and mathematical techniques can be also applied [31–34], the study of bus–metro transfers is considered in [35,36], etc. Moreover, techniques based on the Artificial Intelligence paradigm have also been used to study different aspects of subway networks (see, e.g., [37–39]).

The rest of the paper is organized as follows. Section 2 describes the data used in the study. Section 3 is devoted to presenting the methodology used for the analysis of travel patterns. Finally, the results obtained and the discussion are presented in Section 4 and the conclusions in Section 5.

## 2. Structural and Transit Data of Barcelona Subway Network

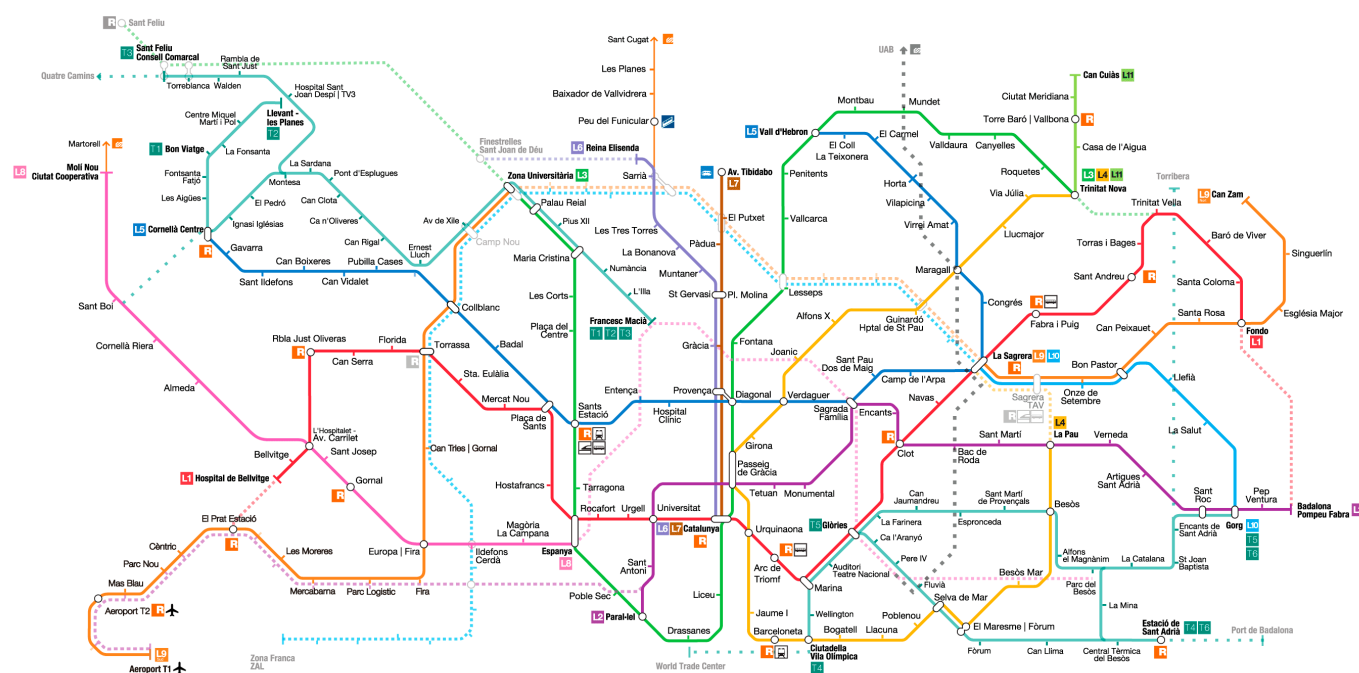
### 2.1. Study Area

Barcelona is considered a significant success in urban development across Europe. As the second largest city of Spain, it has been growing and transforming itself to be a knowledge-intensive city and, more importantly, a pioneer in being a smart city [40]. In addition, it has been one of the Spanish cities with the most confirmed cases of coronavirus. This is why it is an excellent case to explore.

Barcelona has an area of 102 km<sup>2</sup> and a resident population of more than 1.62 million. The city has a diverse public transport system composed of metro, urban and intercity buses, commuter trains, tramway, funicular cable tramway and taxis.

The Barcelona Metro is a metropolitan railway network that gives service to Barcelona and the municipalities of its metropolitan area: Badalona, Cornellà de Llobregat, L'Hospitalet de Llobregat, Montcada i Reixac, El Prat de Llobregat, Sant Adrià de Besòs, Sant Boi de Llobregat and Santa Coloma de Gramanet. It comprises 13 lines with a length of 119 km (see Figure 1):

- L1: Hospital de Bellvitge–Fondo
- L2: Paral·lel–Badalona Centre
- L3: Zona Universitària–Trinitat Nova
- L4: La Pau–Trinitat Nova
- L5: Cornellà Centre–Vall d’Hebron
- L6: Plaça Catalunya–Reina Elisenda
- L7: Plaça Catalunya–Avinguda Tibidabo
- L8: Plaça Espanya -Molí Nou Ciutat Cooperativa
- L9 Nord: La Sagrera–Can Zam
- L9 Sud: Aeroport T1–Zona Universitària
- L10 Nord: La Sagrera–Gorg
- L10 Sud: Foc–Collblanc
- L11: Trinitat Nova–Can Cuiàs



**Figure 1.** The 2019 Barcelona subway (Available online: <https://www.metrobarcelona.es/mapas.html> (accessed on 15 February 2021)).

## 2.2. Transit Data

The data used in this research correspond to the ridership (number of entries) in each station from 5 March 2018 to 11 March 2018. The reason this week was selected is because it is a week without public holidays or summer or winter holidays, and, therefore, it can reflect the general station ridership characteristics under normal circumstances. There was no extreme weather associated with that week either (e.g., heavy storms or very hot temperatures).

A statistical analysis of daily transit data was performed to analyze hourly inbound ridership of the 151 stations of Barcelona subway. The Barcelona metro operates from Sunday to Thursday from 5:00 to 24:00. On Fridays, the metro schedule is extended until 2:00, while on Saturdays it offers continuous service for 24 h. Thus, there are 140 variables for each station.

There are some aspects that need to be taken into account when addressing the analysis. First, it is important to notice there are two time-related patterns: the inbound ridership patterns on weekdays and at weekends. While they are both highly correlated on their own, the correlation between the ridership on weekdays and on the weekend is

relatively low (see Figure 2). Second, from the analysis of the inbound ridership, it can be deduced that the highest peak hour during weekday mornings is between 7:00 and 8:00. During the evening rush hour, the highest peak hours are between 14:00 and 15:00 and between 18:00 and 19:00. Meanwhile, the rush hours during the weekend are from 13:00 to 14:00 and from 18:00 to 19:00 (see Figure 3). Figure 4, where the total number of entries at each hour is added up for all the days in the selected week for 35 randomly selected stations, illustrates how the different rush hours change depending on the station, and that both the time and the number of validations that represent a peak for a station vary. In addition, the total number of passengers significantly differs from one station to another. For instance, taking the daily ridership of 5 March, Diagonal station has a total of 54,636 passengers, while, at Casa de l'agua, there were only 207 boardings that day. These are the stations with the maximum and minimum total number of boardings and illustrate the huge difference there can be. Finally, as shown in Figure 5, the distribution of passenger flow decreases significantly on Saturdays and Sundays, which is why it was decided to focus on the data from Monday to Friday.



**Figure 2.** Pearson's correlation coefficients of daily ridership.



Figure 3. Time-varying diagram of passenger flow (total counts of boarding).

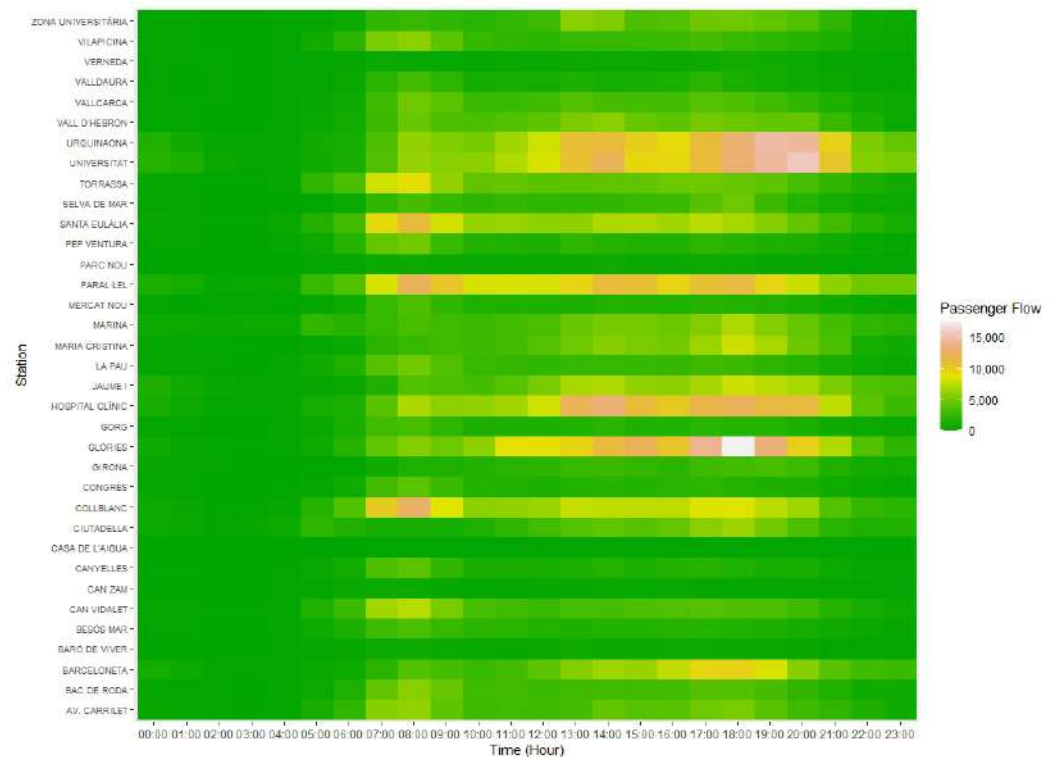
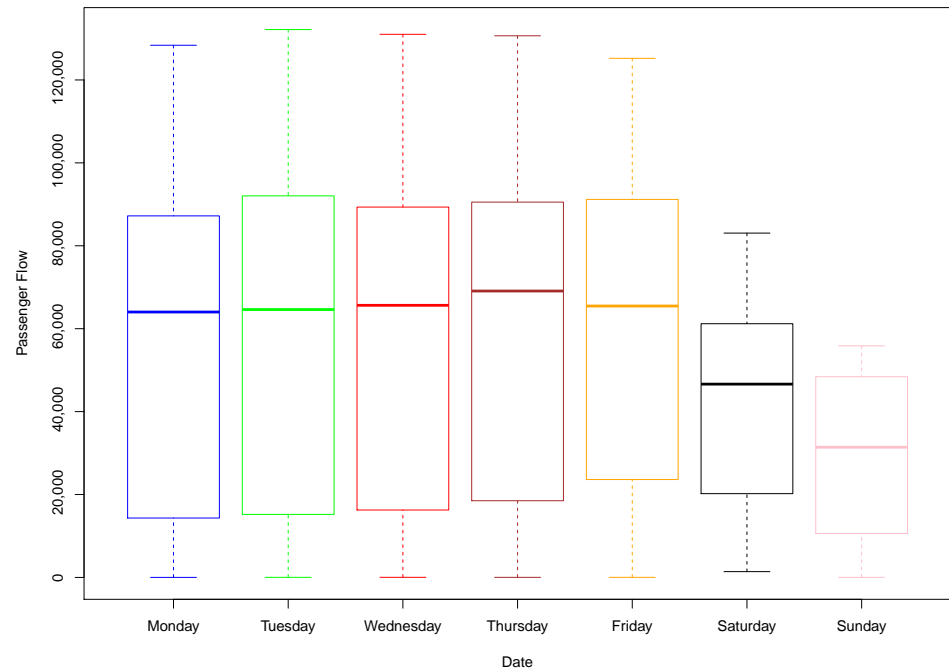


Figure 4. Heatmap with the total number of validations per hour for 35 randomly selected stations.



**Figure 5.** Passenger flow boxplots.

### 3. Methodology

#### 3.1. Complex Network Analysis

In this study, the L-space representation of the network is considered. Hence, the stations of the subway network are represented by nodes of a graph and the tracks connecting two stations are represented by edges of the graph. Therefore, the subway network is represented by a undirected graph  $G = (V, E)$  where  $V = \{v_1, v_2, \dots, v_N\}$  is the set of nodes, and  $E = \{e_{ij} = (v_i, v_j), v_i, v_j \in V\}$  is the set of edges, where  $|E| = M$ .

The adjacency matrix of  $G$ ,  $A_G = (a_{ij})_{1 \leq i, j \leq N}$ , is a  $N \times N$  symmetric matrix such that the coefficient  $a_{ij}$  takes the value 1 or 0 depending on whether or not there is a link between nodes  $v_i$  and  $v_j$ . The degree of a node  $v_i$  is the number of adjacent nodes to  $v_i$  and can be computed as follows:  $d_i = \sum_{j=1}^N a_{ij}$ .

The Laplacian matrix  $Q_G = \Delta - A_G$  is an  $N \times N$  matrix, where  $\Delta = \text{diag}(d_1, \dots, d_N)$  is the  $N \times N$  diagonal degree matrix. The eigenvalues of  $Q_G$  play a very important role in robustness analysis; they are non-negative and can be ordered as  $0 = \mu_N \leq \mu_{N-1} \leq \dots \leq \mu_1$ .

##### 3.1.1. Centrality Measures

The analysis of a complex network is performed through the computation and analysis of several structural coefficients of the network topology. Specifically, the most important are the following [41]: degree centrality, average degree, degree distribution, average path length, closeness centrality and betweenness centrality.

The degree centrality of  $v_i$  is the average number of incident edges to  $v_i$ :

$$C_D(v_i) = \frac{d_i}{N}, \quad (1)$$

and the normalized average degree of the network  $G$  is given by:

$$\bar{E}[D] = \frac{\sum_{i=1}^N d_i}{N(N-1)}. \quad (2)$$

Moreover, the degree distribution of the network,  $P(k)$ , is the probability distribution of degrees over the whole network.

The shortest path length or distance between two nodes  $v_i, v_j \in V$  is denoted by  $d(v_i, v_j)$  and is defined as the minimum number of links necessary to go from node  $v_i$  to node  $v_j$ . The average path length of the network is defined as the average distance between two nodes:

$$L = \frac{2}{N(N-1)} \sum_{1 \leq i < j \leq N} d(v_i, v_j). \quad (3)$$

The diameter  $D$  of  $G$  is the greatest distance between any pair of nodes:

$$D = \max\{d(v_i, v_j), v_i, v_j \in V\}. \quad (4)$$

The closeness centrality of the node  $v_i$  measures the mean distance from  $v_i$  to the rest of the nodes of the network:

$$C_{CL}(v_i) = \frac{1}{\sum_{i \neq j} d(v_i, v_j)}. \quad (5)$$

The greater is the value of closeness centrality, the smaller is the length of the shortest paths to all other nodes.

Finally, the *betweenness centrality* of the node  $v_i \in V$  measures the number of shortest paths between two nodes that run through node  $v_i$ . Mathematically it is defined as follows:

$$C_B(v_i) = \frac{2}{(N-1)(N-2)} \sum_{r \neq s \neq i} \frac{\ell_{rs}(v_i)}{\ell_{rs}}, \quad (6)$$

where  $\ell_{rs}$  is the total number of shortest paths from  $v_r$  to  $v_s$  and  $\ell_{rs}(v_i)$  is the the number of shortest paths between  $v_r$  and  $v_s$  that pass through  $v_i$ . In networks, the greater is the number of paths that pass through a node, the greater is the importance of this node and more central it is.

### 3.1.2. Theoretical Robustness Metrics

Robustness can be defined as the network's ability to survive random failures or deliberate attacks consisting of the elimination of nodes and/or edges [42]. In this sense, several robustness measures have been proposed to quantitatively determine this characteristic. The most important ones are described in what follows:

The *normalized robustness indicator*  $\bar{r}^T$  measures the ratio between the number of alternative paths in the network topology and the total number of stations [8]:

$$\bar{r}^T = \frac{\ln(M - N + 2)}{\ln\left(\frac{N(N-1)}{2} - N + 2\right)}. \quad (7)$$

Note that  $\bar{r}^T$  is higher in the case there are alternative routes to reach a destination and it is smaller in large systems.

The *effective graph resistance*  $R_G$  estimates the robustness of a network from the number of parallel paths (i.e., redundancy) and the length of each path between each pair of nodes. The effective graph resistance is calculated in terms of the eigenvalues of the Laplacian matrix as follows:

$$R_G = N \sum_{i=1}^{N-1} \frac{1}{\mu_i}. \quad (8)$$

In this work, the normalized version of the the effective graph resistance, called *effective graph conductance* [43], is used:

$$C_G = \frac{N-1}{R_G}. \quad (9)$$



Note that  $0 \leq C_G \leq 1$  and a larger  $C_G$  indicates a higher level of robustness. The average efficiency  $E[\frac{1}{H}]$  is defined as follows [44]:

$$E[\frac{1}{H}] = \frac{2}{N(N-1)} \sum_{i,j=1, i \neq j}^N \frac{1}{d(v_i, v_j)}. \quad (10)$$

Note that the greater is the value of the average efficiency, the greater is the robustness of the network (recall that the global efficiency of the complete network is 1).

The *clustering coefficient* is used to assess how the neighbors of a node are connected with one another [41]. For node  $v_i$ , it is mathematically defined as follows:

$$C_C(v_i) = \frac{2E_i}{d_i(d_i - 1)}, \quad (11)$$

where  $E_i$  is the number of edges linked to the neighbors of node  $v_i$ . The clustering coefficient shows the fault tolerance characteristic: in a subway network, when one station is out of function, the traffic will not be affected if the neighboring stations are connected. Thus, a larger value of  $C_C$  implies a better tolerance to fault in a local scale. The *average clustering coefficient* is the average of all the individual clustering coefficients:

$$C_{CG} = \frac{1}{N} \sum_{i=1}^N C_C(v_i). \quad (12)$$

The *algebraic connectivity*  $\mu_{N-1}$  is the second smallest eigenvalue of the Laplacian matrix  $A_G$ . It has been shown that the larger  $\mu_{N-1}$  is, the higher the robustness of a network is [43]. The *normalized algebraic connectivity* is obtained dividing by the total number of nodes:  $\bar{\mu}_{N-1} = \frac{\mu_{N-1}}{N}$ .

The *normalized natural connectivity*  $\bar{\lambda}$  is defined as:

$$\bar{\lambda} = \frac{\ln[\frac{1}{N} \sum_{i=1}^N e^{\lambda_i}]}{N - \ln N}, \quad (13)$$

where  $\lambda_i$  is the  $i$ th eigenvalue of the adjacency matrix  $A_G$ . It measures the redundancy in terms of alternative paths and is considered as a measure of structural robustness [45].

Finally, the *degree diversity*  $\kappa$  is defined as:

$$\kappa = \frac{\sum_{i=1}^N d_i^2}{\sum_{i=1}^N d_i}. \quad (14)$$

The greater  $\kappa$  is, the more nodes must be removed from the network to disintegrate it [46]. In this work, we take the inverse of the degree diversity  $\bar{\kappa} = \frac{1}{\kappa}$  in order to scale the value in the interval  $[0, 1]$ .

### 3.2. Normalization and Dimensionality Reduction

Given the large differences in the number of passengers from station to station, the entries are normalized. The normalization consists in using the ratio of hourly passengers to the total number of passengers that day at each station, instead of the total amount of passengers per hour [24].

On the other hand, the number of variables used to classify the stations is large and they are also highly correlated; therefore, it was decided to perform a Principal Component Analysis (PCA). PCA is a technique for reducing the dimensionality of large datasets, increasing interpretability and minimizing information loss [47]. PCA is defined

as an orthogonal linear transformation which transforms the data into a new system of coordinates such that the first coordinate (called the first principal component) represents the largest variance, the second coordinate the second greatest, etc. PCA can be thought of as fitting an  $n$ -dimensional ellipsoid to the data, where each axis of it represents a principal component. If an axis of the ellipse is small, then the variance along that axis is also small. To find the axes of the ellipse, first the mean of each variable from the dataset must be subtracted to center the data around the origin. Then, the covariance matrix of the data is computed. The covariance between two data is calculated as:

$$\sigma_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \quad (15)$$

The principal components are calculated from the eigen-vectors and eigenvalues of this matrix. The eigenvectors represent the directions, whereas the eigenvalues are the numbers representing how much variance there is in the data in each particular direction. The eigenvector with the highest eigenvalue is taken as the first principal component. More details can be found in the work of Dunteman [48].

### 3.3. Clustering Analysis

Cluster analysis is an exploratory technique which is used to classify objects into groups, known as clusters, in such way that observations belonging to a cluster are more similar to each other than observations assigned to different clusters. Nevertheless, clustering is rather a subjective statistical analysis and there are several possible algorithms that may be used. The decision of which technique to apply should be made depending on the kind of data or the type of problem to be solved. The  $k$ -means algorithm is known to be computationally fast and has the ability to handle large datasets. However, one needs to know the number of clusters in advance, it is sensitive to outliers and different initial centroids produce different results [49]. Hierarchical clustering is one of the most popular clustering techniques. Although it may be computationally slower when the dataset size increases and clusters depend on the distance metric used, the authors consider that the result of a hierarchical clustering is a structure that is more informative and interpretable than the unstructured set of flat clusters returned by  $k$ -means. Hence, it is easier to determine the optimal number of clusters by looking at the dendrogram of a hierarchical clustering than trying to predict this optimal number in advance in case of  $k$ -means. For these reasons, the agglomerative hierarchical clustering technique is used [50]. The basic algorithm consists of the following steps:

1. Initially, each observation is considered as a single-element cluster.
2. An iterative process is then initiated in which the two clusters that are the most similar are combined into a new bigger cluster. This is done by computing the dissimilarities between every pair of observations. This procedure is iterated until all points are members of one single big cluster.
3. Finally, one needs to determine where to cut the hierarchical tree into clusters. This creates a partition of the data.

The distance between clusters can be calculated using different methods [51,52]. In this study, the Ward method was used, which has been very widely used since its first description by Ward Jr [53], it and has outperformed other methods in several comparison studies [54,55]. The Ward method is the only one among the agglomerative clustering methods that is based on a classical sum-of-squares criterion, producing groups that minimize within-group dispersion at each binary [56]. In the Ward method, the distance between two clusters,  $A$  and  $B$ , is how much the sum of squares will increase once they are merged:

$$\begin{aligned}
\Delta(A, B) &= \sum_{i \in A \cup B} \|\vec{x}_i - \vec{m}_{A \cup B}\|^2 - \sum_{i \in A} \|\vec{x}_i - \vec{m}_A\|^2 - \sum_{i \in B} \|\vec{x}_i - \vec{m}_B\|^2 \\
&= \frac{n_A n_B}{n_A + n_B} \|\vec{m}_A - \vec{m}_B\|^2,
\end{aligned} \tag{16}$$

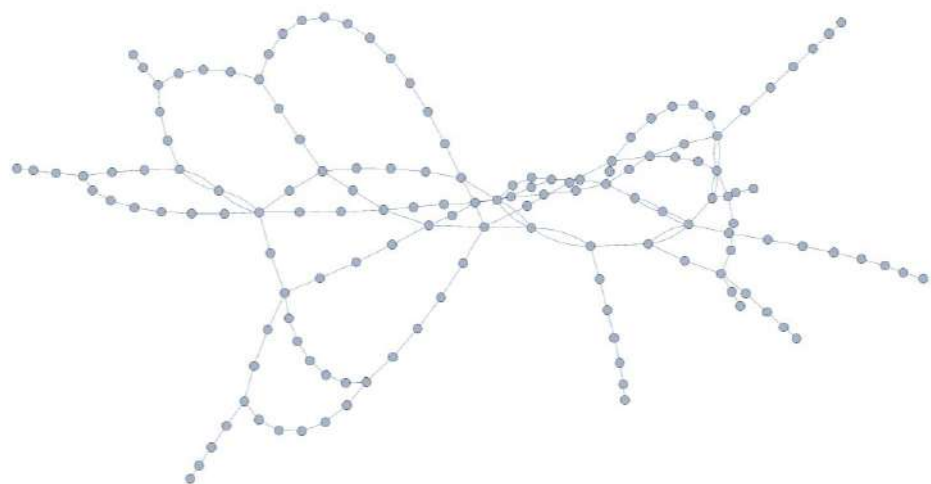
where  $\vec{m}_j$  is the center of cluster  $j$  and  $n_j$  is the number of points in it.  $\Delta$  is called the merging cost of combining the clusters  $A$  and  $B$ . In this method, in each step, the variability within clusters is minimized.

In addition, the agglomerative coefficient (AC), measuring the clustering structure of the dataset, is calculated [57]. For each observation  $i$ , let  $m(i)$  represent its dissimilarity to the first cluster it is merged with, divided by the dissimilarity of the merger in the final step of the algorithm. The AC is the average of all  $1 - m(i)$ . Generally speaking, the AC describes the strength of the clustering structure that has been obtained by group average linkage. However, the AC tends to become larger when  $n$  increases, so it should not be used to compare datasets of very different sizes. The coefficient takes values from 0 to 1, and it is actually the mean of the normalized lengths at which the clusters are formed. A coefficient close to 1 points to a pretty reasonable cluster structure in the data.

## 4. Mathematical and Statistical Analysis

### 4.1. Structural Network Analysis

As previously mentioned, the topology of Barcelona subway network is established using the L-space method, where each station stands for a node of the graph and the edges are defined by means of the direct connections by rail ways between the stations. The number of nodes is  $N = 151$  and the number of edges is  $M = 177$  and therefore the density of the subway network is  $d \approx 0.0157$ . In Figure 6, the graph corresponding to Barcelona subway network using Mathematica is shown (note that the exact placing and positioning of the stations is not taken into account).



**Figure 6.** The graph representing the Barcelona subway network using Mathematica.

#### 4.1.1. Basic Structural Characteristics

In this subsection, the most usual coefficients and centrality measures, introduced in Section 3.1.1, are computed and associated to the Barcelona subway network.

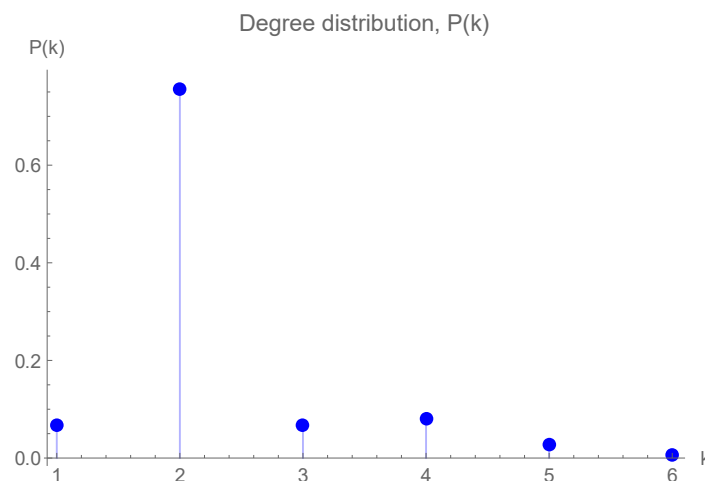
As shown in Table 1 the five stations with the highest degree are “Passeig de Gràcia” with degree 6 and “Diagonal”, “Espanya”, “Catalunya” and “La Sagrera” with degree 5.

Note that the first four stations belong to Line 3; in addition, three of the top five are on Line 1.

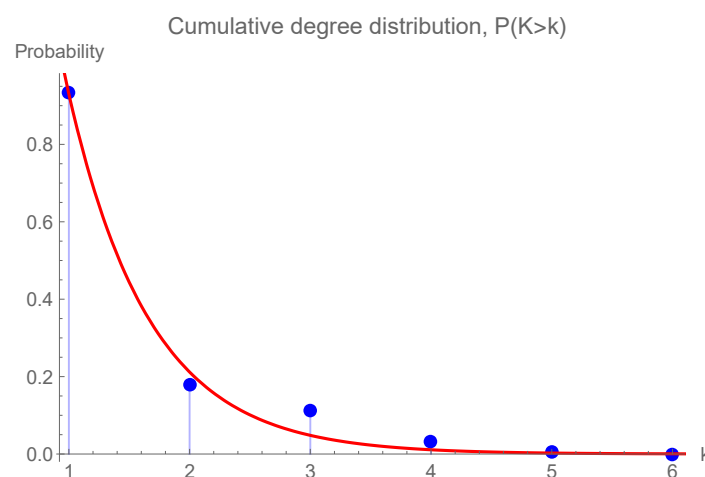
**Table 1.** The five stations with the highest degree.

Station	Subway Lines	Degree
Passeig de Gràcia	2, 3, 4	6
Diagonal	3, 5	5
Espanya	1, 3	5
Catalunya	1, 3	5
La Sagrera	1, 5, 9N, 10N,	5

The average degree of the network is  $E[D] \approx 2.2649$  and the degree distribution  $p(k)$  is shown in Figure 7, while the cumulative degree distribution is illustrated in Figure 8. A simple calculus shows that the fitting function of the cumulative degree distribution is  $h(x) = 4.0834e^{-1.4796x}$ .



**Figure 7.** Degree distribution of Barcelona subway network.



**Figure 8.** Cumulative degree distribution of Barcelona subway network.

The maximum travel distance of the network is no more than 31 stops (diameter), while the average shortest path is 11.0032 stops.

Table 2 shows the results obtained from the computation of the closeness centrality. The station with the highest closeness centrality is “Diagonal” with  $C_{CL} \approx 0.1424$ , and the next four stations (“Verdaguer”, “Hospital Clínic”, “Passeig de Gràcia” and “Provença”)

have similar closeness centrality. In this case, the most centrality subway line is Line 5 and, to a lesser extent, Line 3.

**Table 2.** The five stations with the highest closeness centrality.

Station	Subway Lines	Closeness Centrality
Diagonal	3, 5	0.1424
Verdaguer	4, 5	0.1372
Hospital Clínic	5	0.1362
Passeig de Gràcia	2, 3, 4	0.1358
Provença	3, 5, 6, 7	0.1323

Finally, the results obtained when the betweenness centrality was computed are displayed in Table 3. It is important to note that all the stations with the highest coefficient belong to Line 5.

**Table 3.** The five stations with the highest betweenness centrality.

Station	Subway Lines	Betweenness Centrality
Diagonal	3, 5	0.4298
Verdaguer	4, 5	0.3333
Sants Estació	3, 5	0.2610
Hospital Clínic	5	0.2593
Entença	5	0.2553

From these results, it can be seen that some specific stations play a central role in the structural definition of the network. For example, “Diagonal” and “Verdaguer” are very important structural pieces of the subway network since they have the highest values of closeness and betweenness centralities. In addition the most central lines are Lines 5, 3 and 1.

#### 4.1.2. Network Robustness

Failures of subway networks can have enormous impact on our society, so the analysis of the robustness is very important when studying subway networks. The robustness of networks reflects the extent to which the networks can solve possible (intentional or unintentional) failures by offering alternative routes that overcome the attacked edges or nodes.

In this section, eight robustness metrics (introduced in Section 3.1.2 are computed for the Barcelona subway network and compared with those obtained for the Madrid subway network.

In Table 4, the stations with the highest clustering centrality are illustrated. The most central are “Catalunya” ( $C_C = 0.2$ ), “Universitat” and “Urquinaona” with  $C_C \approx 0.1666$  and “Passeig de Gràcia” with  $C_C \approx 0.1333$ . As a consequence, they have better tolerance to fault in a local scale. The first three stations belong to Line 1, and Lines 2–4 have a couple of stations on this list. Moreover, the mean clustering coefficient is 0.0044, which is significantly lower than that of other metro networks such as London ( $C_C = 0.0409$ ), Tokyo ( $C_C = 0.0285$ ) or Paris ( $C_C = 0.0163$ ) [58].

Table 5 shows the values of the eight robustness metrics computed using Equations (7)–(14) for the Barcelona subway network and the Madrid subway network [59].

According to the reduced robustness indicator  $\bar{r}^T$ , the Barcelona metro network is slightly more robust than the Madrid metro network, probably because there are more alternative paths between any pair of nodes.

**Table 4.** Stations with non-zero clustering centrality.

Station	Subway Lines	Clustering Centrality
Catalunya	1, 3	0.2
Universitat	1, 2	0.1666
Urquinaona	1, 4	0.1666
Passeig de Gràcia	2, 3, 4	0.1333

According to the effective graph conductance  $C_G$ , the Barcelona subway network also has a slightly higher value than that of Madrid. Note that the effective graph conductance takes into account not only the number of alternative paths but also the length of each alternative path, hence effective graph conductance favors networks with the smallest length of the shortest paths.

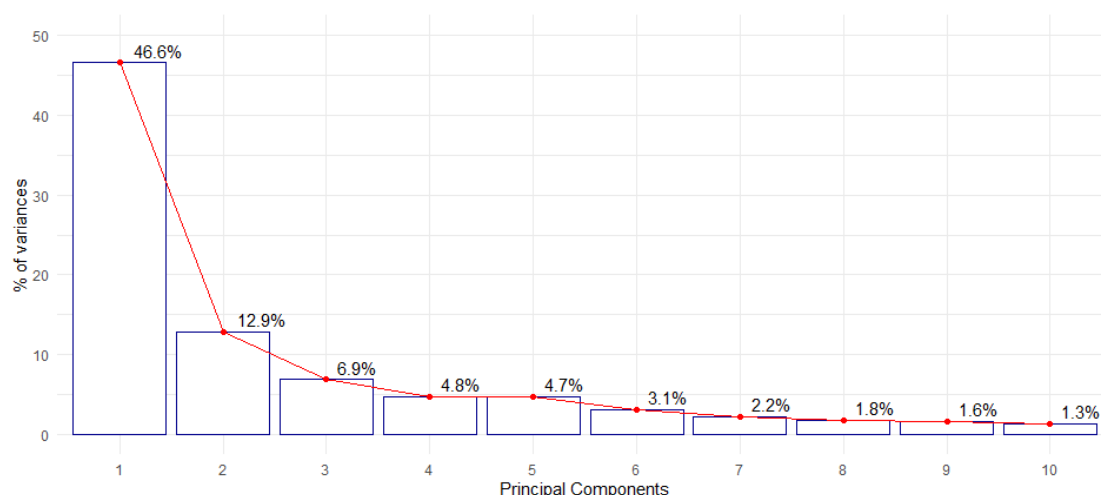
In general, according to all the metrics except the clustering coefficient  $C_{CG}$  and the normalized degree diversity  $\bar{\kappa}$ , Barcelona has a higher robustness level than Madrid.

**Table 5.** Robustness metrics in Barcelona and Madrid subway networks.

Coefficients	Barcelona Subway	Madrid Subway
Nodes, $N$	151	243
Edges, $M$	177	280
Normalized robust indicator, $\bar{r}^T$	0.35747	0.35635
Effective graph conductance, $C_G$	0.00221	0.00086
Average efficiency, $E[\frac{1}{H}]$	0.13524	0.10533
Average clustering coefficient, $C_{CG}$	0.00441	0.00774
Normalized algebraic connectivity, $\bar{\mu}_{N-1}$	0.00006	0.00001
Normalized average degree, $\bar{E}[D]$	0.01562	0.00952
Normalized natural connectivity, $\bar{\lambda}$	0.00770	0.00441
Normalized degree diversity, $\bar{\kappa}$	0.35975	0.37135

#### 4.2. Data Analysis Results

Principal component analysis was performed to study the data from the working days (Monday to Friday) of the selected week. The first three principal components are able to explain 66.32% of the variability in the data (PC1 = 46.56%, PC2 = 12.88% and PC3 = 6.89%). Figure 9 shows the total variability explained by each principal component.

**Figure 9.** Total variance explained by each principal component (weekdays).

In Figure 10, the top plot shows the contributions of 18 variables to the first three components. The six variables which most contribute to each component are chosen. In the bottom plot, the correlations of these 18 variables to each component are shown. The contribution is represented both by the color scale and the circle size, while, for the correlation, the direction of the correlation is represented by color and the circle size represents the strength of the relationship. The variables which contribute the most to the first component are those corresponding to 7 a.m., and they are strongly negatively correlated with it. Regarding the second component, the variables which contribute the most are the ones corresponding to 11 a.m. and noon. Finally, the variables contributing to the third component are the ones from 1 and 11 p.m. The second and third components have a positive correlation with the variables that contribute the most to them.



**Figure 10.** Contributions (**top**) and correlation (**bottom**) for the first three components.

A hierarchical cluster analysis was performed over the coordinates from the first three principal components. The resulting AC is 0.9811, which indicates a pretty reasonable cluster structure in the data. The dendrogram in Figure 11 shows that two clustering solutions are possible. The four-cluster solution is chosen as it provides a more detailed segmentation of the stations.

Statistical properties of the four clusters are summarized in Table 6. The diameters represent the maximum within cluster distances. The average and median distances are the within cluster average and median distances. Separation is the minimum distance of a point in the cluster to a point of another cluster and average to other is the average distance of a point in the cluster to the points of other clusters.

**Table 6.** Statistical properties of the four clusters (weekdays).

Cluster	1	2	3	4
Size	49	41	35	4
Diameter	17.49	9.39	9.32	8.46
Average distance	7.55	4.15	3.58	5.42
Median distance	7.30	4.05	3.45	5.21
Separation	1.64	1.08	1.08	7.70
Average to other	12.76	9.32	11.79	19.93

In Table 7, the stations belonging to each cluster are listed. For a better understanding of the clusters, the different stations of each cluster are located in the Barcelona map, making use of a Voronoi diagram (based on Euclidean distance) to partition the city map. In Figure 12, each Voronoi cell representing a station is colored by cluster. It may be noted

that stations from the same cluster are not necessarily close in space, but their behavior pattern is similar. This may be due to, e.g., the business activities taking place in the area or being residential neighborhoods.

Cluster Dendrogram

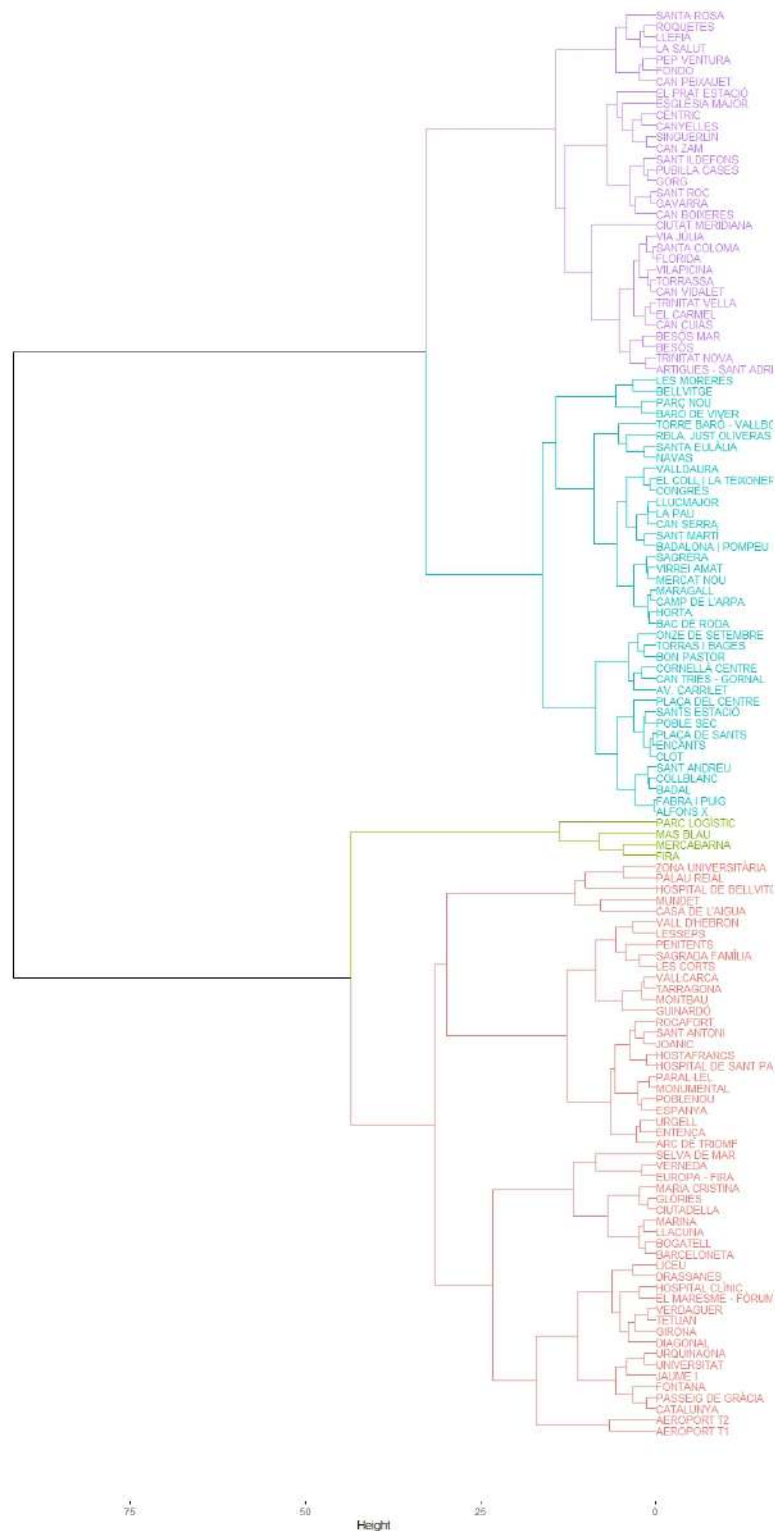
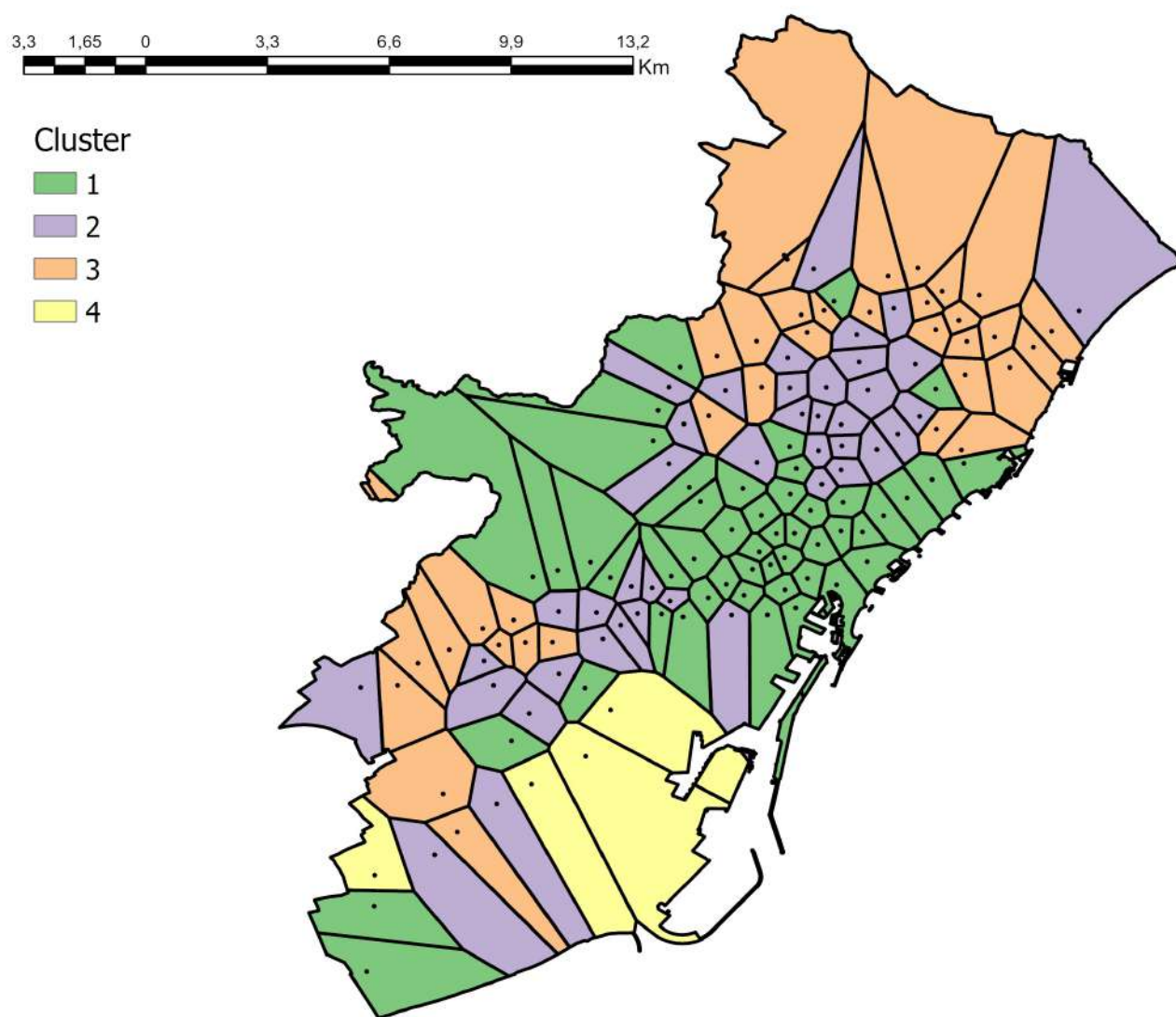


Figure 11. Clusters: Hierarchical clustering.





**Figure 12.** Map of the different stations colored by cluster (weekdays).

In the case of Cluster 1, most of the stations are located in the district of L'Eixample, Ciutat Vella and Sant Martí, where some of the most popular beaches of Barcelona are located, as well as important monuments such as Casa Milà, popularly known as *La Pedrera*, the Cathedral, Park Güell and Casa Batlló. Moreover, this cluster includes the zoo, the Maritime Museum of Barcelona and the museum *Poble Espanyol*. The hospital stations Vall d'Hebron, Hospital Clínic, Sant Pau and Hospital de Bellvitge are included in this cluster too, as well as those belonging to university campuses, such as Mundet, Palau Reial, Universitat and Zona Universitària. There are also two stations from the airport and some stations from the districts Les Corts, Sants, Montjuic and Gracia, all of them located in the city center. In Figure 13, passenger flow per hour is shown for some of the stations in Cluster 1. All of them have peak hours at 8 a.m., 2 p.m. and 7 p.m.

Table 7. Results of station classification.

Cluster	Stations	Number
Cluster 1	Aeroport T1, Aeroport T2, Arc de Triomf, Barceloneta, Bogatell, Casa de l'aigua, Catalunya, Ciutadella, Diagonal, Drassanes, El Maresme-Fórum, Entença, Espanya, Europa-Fira, Fontana, Girona, Gloriès, Guinardó, Hospital Clínic, Hospital de Bellvitge, Hospital de Sant Pau, Hostafrancs, Jaume I, Joaníc, Les Corts, Lesseps, Liceu, Llacuna, Maria Cristina, Marina, Monumental, Mundet, Palau Reial, Parallel, Passeig de Gràcia, Penitents, Poble-nou, Rocafort, Sagrada Família, Sant Antoni, Selva de Mar, Tetuan, Universitat, Urgell, Urquinaona, Vall d'Hebron, Verdaguer, Verneda, Zona Universitària	49
Cluster 2	Alfons X, Av. Carrilet, Bac de roda, Badal, Badalona - Pompeu Fabra, Baró de viver, Bellvitge, Bon pastor, Camp de l'arpa, Can tries - Gornal, Clot, Collblanc, Congrés, Cornellà Centre, El Coll - La Teixonera, Encants, Fabra i Puig, Horta, La Pau, Les Moreres, Lluçmajor, Maragall, Mercat Nou, Montbau, Navas, Onze De Setembre, Parc Nou, Plaça De Sants, Plaça Del Centre, Poble Sec, Rbla. Just Oliveras, Sagrera, Sant Andreu, Sant Martí, Santa Eulàlia, Sants Estació, Tarragona, Torras i Bages, Torre Baró - Vallbona, Vallcarca, Virrei Amat	41
Cluster 3	Artigues - Sant Adrià, Besòs, Besòs Mar, Can Boixeres, Can Cuiós, Can Peixauet, Can Serra, Can Vidalet, Can Zam, Canyelles, Cèntric, Ciutat Meridiana, El Carmel, El Prat Estació, Església Major, Florida, Fondo, Gavarra, Gorg, La Salut, Llefia, Pep Ventura, Pubilla Cases, Roquetes, Sant Ildefons, Sant Roc, Santa Coloma, Santa Rosa, Singuerlín, Torrassa, Trinitat Nova, Trinitat Vella, Valldaura, Via Júlia, Vilapicina	35
Cluster 4	Fira, Mas Blau, Mercabarna, Parc Logístic	4

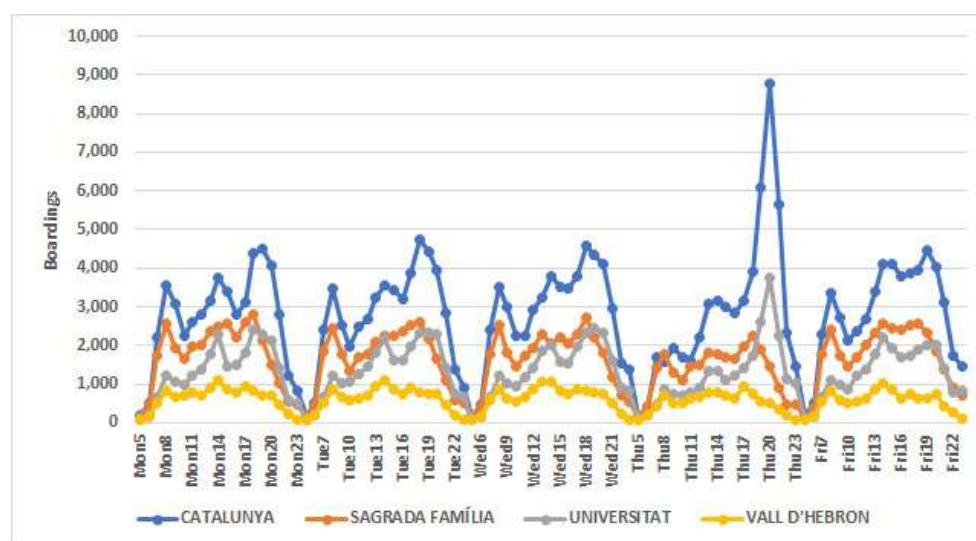


Figure 13. Pattern of boardings in stations of Cluster 1 (weekdays).

The stations in Cluster 2 are mainly around the central area of Barcelona, with some in the north and some in the south. These are traditional, residential, well-communicated neighborhoods, with many markets and shops. The stations in the north are from the districts Sant Andreu, Horta-Ginardó and Nou Barris. The stations in the south belong to L'Eixample and are the furthest from the city center together with the stations from Sant Andreu, one of the entrances to Barcelona with a large cultural and sports offer. The hours with the largest number of passengers in this cluster are 7 a.m., 8 a.m. and 6 p.m. The pattern of boarding per hour for some stations in this cluster is shown in Figure 14.

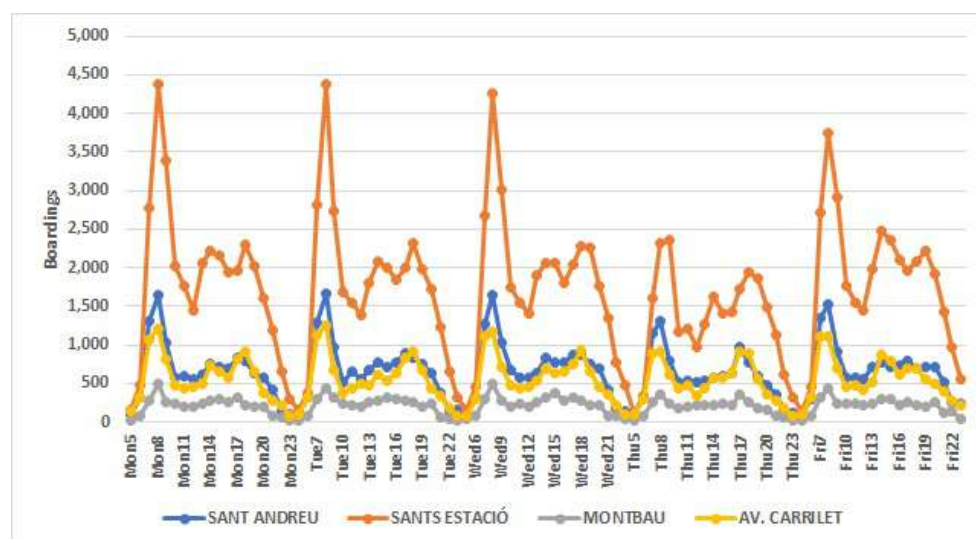


Figure 14. Pattern of boarding in stations of Cluster 2 (weekdays).

Cluster 3 contains mostly stations located outside of the city. There are two stations in El Prat de Llobregat and eight bordering the north side of L'Hospitalet de Llobregat. The rest are gathered in the north urban periphery of the city, linking to different small municipalities or towns, such as Badalona, Santa Coloma de Gramenet or Sant Adrià de Besòs. These belong to what is known as the metropolitan area of Barcelona, which is a geographical area that goes beyond the administrative area. Given the growth of the city of Barcelona, some of these municipalities are now essentially suburbs of Barcelona. Badalona is, however, the third largest city in Catalonia. Moreover, there are also stations in Ciutat Meridiana, which is the poorest neighborhood of the city. In Figure 15, the peak hours of the stations of this cluster can be seen. The hours with the highest number of boardings are 8 a.m., 2 p.m. and 6 p.m.

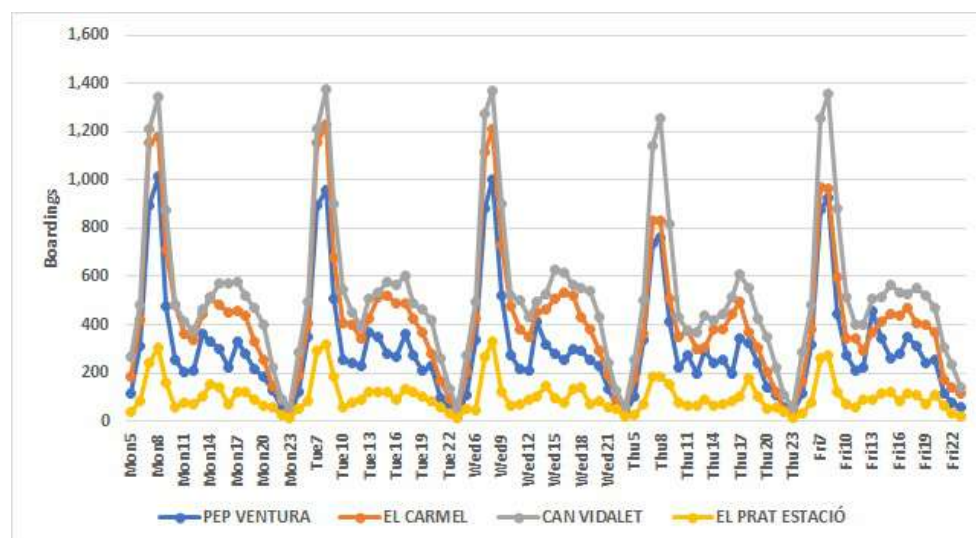


Figure 15. Pattern of boarding in stations of Cluster 3 (weekdays).

The stations that form Cluster 4 have the particular characteristics of the area they give access to: Fira is the entry to one of the largest and most modern fairgrounds of Europe; Mas Blau corresponds with the industrial park closest to Barcelona's airport; Mercabarna is considered the most important central market in Europe, as it is a reference center in the Mediterranean Sea for the distribution of fresh products at the international level; and Parc Logístic serves the logistics park of the city's Free economic zone. Overall, 2 p.m., 5 p.m.

and 6 p.m. have the highest number of boardings. The peak hours of these stations are shown in Figure 16.

All the analysis here presented were performed with RStudio Team [60].

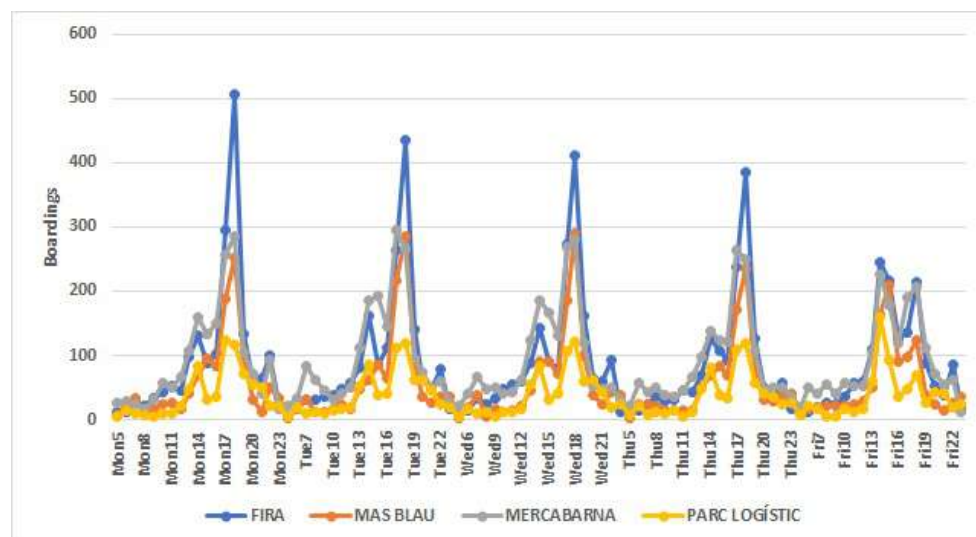


Figure 16. Pattern of boarding stations of Cluster 4 (weekdays).

## 5. Conclusions

In Barcelona, as in any major urban area, many people use the public transport network, which is why it is necessary to have as much information as possible to forecast and plan the subway trip.

Moreover, in the bibliography studied, there are no previous studies that analyze not only the structural and robustness characteristics but also travel patterns of the Barcelona metro network.

In this study, a detailed analysis of Barcelona subway network was done using Complex Network Analysis. To achieve this goal, the most important centrality measures and coefficients were computed. In this sense, the important role of stations such as “Diagonal” and “Verdaguer” to control the flow of passengers was shown. It was also shown that the stations “Catalunya”, “Universitat”, “Urquinaona” and “Passeig de Gràcia” have high fault tolerance in a local scale. Moreover, L5 and L3 are the most central subway lines.

In addition, the robustness of the Barcelona subway network was investigated by analyzing several robustness metrics and compared with the robustness of the Madrid subway network. The results indicate that the Barcelona subway network is slightly more robust than the Madrid subway network according to most of the robustness metrics. A previous study [8] analyzed Barcelona subway robustness using ten theoretical robustness metrics, but only taking into account terminals and transfer stations. The results in the former study cannot be compared with ours since in our study all Barcelona subway stations are used.

The data collected at the entry of the metro stations in Barcelona provide a vast quantity of data with very valuable information about the ridership patterns in them. The set of real data was provided by the Barcelona Metropolitan Network, providing information on the number of entries per hour in each of the 151 stations. There are no data related to the passenger’s journey or personal data (age, sex, fare, etc.).

The statistical techniques used in this study allowed observing the following: in the first place, there are differences in behavior between working days, which are highly correlated with each other, and over the weekend, with which the correlation decreases. The hours with the highest number of passengers correspond mainly to the hours of entry and exit of work and school hours. However, these rush hours are not the same at all stations, nor are the number of passengers each have, reaching a difference of more than

54,000 daily entries between some stations. It is because of this reason that the data were normalized, using the proportion of passengers per hour with respect to the total number of entries in that particular day at each particular station.

The principal component analysis performed reduced the dimensionality of the dataset. The first three principal components explain most of the variability in the data. Moreover, it was observed which hours have a higher effect in each of them.

The cluster analysis carried out revealed, for working days, the existence of four groups with similar characteristics. The first conglomerate gathers the stations of the downtown area, the most touristic and monumental. In the second cluster, the stations that surround the center of Barcelona are grouped. They are, mainly, traditional and residential neighborhoods. The periphery stations, which link the center with the nearest municipalities, are those found in the third cluster. In the fourth cluster, the stations of the fairgrounds, large markets and logistics parks appear. Within each cluster, one can see the same pattern of behavior that reflects the similarities of the stations that form it, as can be seen at peak times, which differ between clusters.

The patterns observed reflect the daily activities of the urban area of Barcelona, which are related to the spatial structuring of the city and its characteristics, and are highly correlated with general daily routines.

The results of this work provide relevant information for the “Transports Metropolitan of Barcelona” company for public transport planning. These studies allow us to discover patterns of behavior needed to make decisions to improve the metro service. Nowadays, in the new post-pandemic normality, it is imperative to travel safely so as to stop the coronavirus spreading. It is important to avoid rush hours travels; people may choose to get on and off at subway stations with fewer travelers and do part of their journey by foot. Moreover, it is the task of public transport companies to increase the number of subway cars at a certain time if it gets too crowded, improve the infrastructure of stations with high passenger flow and reduce the time in-between metro services, among other security measures. For instance, the station “Sant Andreu”, from Cluster 2, has the highest number of passengers between 7:00 and 8:00 a.m., and, therefore, it is one of the stations where increasing the number of subway cars or the frequency of the service would be imperative. On the other hand, the station “Fira”, from Cluster 4, has peak hours at 14:00, 17:00 and 18:00 (p.m.), although with a much smaller number of passengers than “Sant Andreu”, and, thus, depending on the capacity of the station, the measures may not be as crucial as in the first one.

Future work involves relating these results to population, climate and economic variables that reflect other social circumstances that may influence the characteristics of the metro network stations. Moreover, annual data shall be analyzed to detect seasonality in behavior patterns. Further lines of investigations will also include a structural and robustness analysis of the network, using complex network analysis to determine critical nodes using different centrality measures. In addition, a detailed analysis of the structural characteristics of this subway network considering other different topological representations such as reduced L-space, P-space, C-space, etc. must be tackled. In addition, a theoretical framework must be proposed in which the notion of “subway line” is used as the basis to define new structural and robustness coefficients. Furthermore, additional transport lines (light rail network, bus network, etc.), can be considered in the analysis to obtain more realistic results. It would also be interesting to analyze the data post-COVID-19 and compare how the use of the public transport has changed, once the data become available.

**Author Contributions:** Conceptualization, E.F.B., M.T.S.M. and A.M.d.R.; methodology, E.F.B., I.M.-C. and R.C.V.; software, I.M.-C.; writing—original draft preparation, I.M.-C., E.F.B., A.M.d.R. and A.B.G.-G.; and writing—review and editing, M.T.S.M., A.M.d.R., R.C.V. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was funded by Ministerio de Ciencia, Innovación y Universidades (MCIU, Spain), Agencia Estatal de Investigación (AEI, Spain), and Fondo Europeo de Desarrollo Regional



(FEDER, UE) under project NOTREDAMME and by Scientific Research Grant of the “Fundación Memoria D. Samuel Solórzano Barruso”, University of Salamanca.

**Data Availability Statement:** Not Applicable.

**Acknowledgments:** The authors extend their gratitude to the Transport Metropolitans of Barcelona.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## Abbreviations

The following abbreviations are used in this manuscript:

AFC	Automated Fare Collection
AC	Agglomerative Coefficient
PCA	Principal Component Analysis

## References

1. Pternea, M.; Kepaptsoglou, K.; Karlaftis, M.G. Sustainable urban transit network design. *Transp. Res. Part A Policy Pract.* **2015**, *77*, 276–291. [CrossRef]
2. Latora, V.; M., M. Is the Boston subway a small-world network? *Phys. A* **2002**, *314*, 109–113. [CrossRef]
3. Lu, H.P.; Shi, Y. Complexity of public transport network. *Tsinghua Sci. Technol.* **2007**, *12*, 204–213. [CrossRef]
4. Zhang, J.H.; Zhao, M.W.; Liu, H.K.; Xu, X.M. Networked characteristics of the urban rail transit networks. *Phys. A* **2013**, *392*, 1538–1546. [CrossRef]
5. Liu, Z.; Song, R. Reliability analysis of Guangzhou subway with complex network theory. *J. Transp. Syst. Eng. Inf. Technol.* **2010**, *10*, 194–200.
6. Cats, O. Topological evolution of a metropolitan rail transport network: The case of Stockholm. *J. Transp. Geogr.* **2017**, *62*, 172–183. [CrossRef]
7. Derrible, S.; Kennedy, C. The complexity and robustness of metro networks. *Phys. A* **2010**, *389*, 3678–3691. [CrossRef]
8. Wang, X.; Koc, Y.; Derrible, S.; Ahmad, S.N.; Kooij, R.E. Multi-criteria robustness analysis of metro networks. *Phys. A* **2017**, *474*, 19–31. [CrossRef]
9. Zhang, J.H.; Xu, X.M.; Hong, L.; Wang, S.; Fei, Q. Networked analysis of the Shanghai subway network, in China. *Phys. A* **2011**, *390*, 4562–4570. [CrossRef]
10. Forero-Ortiz, E.; Martinez-Gomariz, E.; Canas Porcuna, M.; Locatelli, L.; Russo, B. Flood Risk Assessment in an Underground Railway System under the Impact of Climate Change-A Case Study of the Barcelona Metro. *Sustainability* **2020**, *12*, 5291.10.3390/su12135291. [CrossRef]
11. De Bona, A.; de Oliveira Rosa, M.; Ono Fonseca, K.; Lüders, R. A reduced model for complex network analysis of public transportation systems. *Phys. A Stat. Mech. Its Appl.* **2021**, *567*, 125715. [CrossRef]
12. Wang, Y.; Tian, C. Measure Vulnerability of Metro Network under Cascading Failure. *IEEE Access* **2021**, *9*, 683–692. [CrossRef]
13. Dempsey, P.S. *Privacy Issues with the Use of Smart Cards*; The National Academies Press: Washington, DC, USA, 2007.
14. Pelletier, M.P.; Trépanier, M.; Morency, C. Smart card data use in public transit: A literature review. *Transp. Res. Part C Emerg. Technol.* **2011**, *19*, 557–568. [CrossRef]
15. Li, T.; Sun, D.; Jing, P.; Yang, K. Smart card data mining of public transport destination: A literature review. *Information* **2018**, *9*, 18. [CrossRef]
16. Alsger, A.; Tavassoli, A.; Mesbah, M.; Ferreira, L.; Hickman, M. Public transport trip purpose inference using smart card fare data. *Transp. Res. Part C Emerg. Technol.* **2018**, *87*, 123–137. [CrossRef]
17. Alexander, L.; Jiang, S.; Murga, M.; González, M.C. Origin-destination trips by purpose and time of day inferred from mobile phone data. *Transp. Res. Part C Emerg. Technol.* **2015**, *58*, 240–250. [CrossRef]
18. Jun, C.; Dongyuan, Y. Estimating smart card commuters origin-destination distribution based on APTS data. *J. Transp. Syst. Eng. Inf. Technol.* **2013**, *13*, 47–53. [CrossRef]
19. Gonzalez, M.C.; Hidalgo, C.A.; Barabasi, A.L. Understanding individual human mobility patterns. *Nature* **2008**, *453*, 779–782. [CrossRef]
20. Briand, A.S.; Côme, E.; Trépanier, M.; Oukhellou, L. Analyzing year-to-year changes in public transport passenger behaviour using smart card data. *Transp. Res. Part C Emerg. Technol.* **2017**, *79*, 274–289. [CrossRef]
21. El Mahrsi, M.K.; Come, E.; Oukhellou, L.; Verleysen, M. Clustering Smart Card Data for Urban Mobility Analysis. *IEEE Trans. Intell. Transp. Syst.* **2017**, *18*, 712–728. [CrossRef]
22. Chen, C.; Chen, J.; Barry, J. Diurnal pattern of transit ridership: A case study of the New York City subway system. *J. Transp. Geogr.* **2009**, *17*, 176–186. [CrossRef]

23. Wang, W.; Lo, S.; Liu, S. Aggregated metro trip patterns in urban areas of Hong Kong: Evidence from automatic fare collection records. *J. Urban Plan. Dev.* **2015**, *141*, 05014018. [CrossRef]
24. Kim, M.K.; Kim, S.P.; Heo, J.; Sohn, H.G. Ridership patterns at subway stations of Seoul capital area and characteristics of station influence area. *KSCE J. Civ. Eng.* **2017**, *21*, 964–975. [CrossRef]
25. Ding, C.; Cao, X.; Liu, C. How does the station-area built environment influence Metrorail ridership? Using gradient boosting decision trees to identify non-linear thresholds. *J. Transp. Geogr.* **2019**, *77*, 70–78. [CrossRef]
26. Langlois, G.G.; Koutsopoulos, H.N.; Zhao, J. Inferring patterns in the multi-week activity sequences of public transport users. *Transp. Res. Part C Emerg. Technol.* **2016**, *64*, 1–16. [CrossRef]
27. Lu, Y.; Zhang, Y. Toward a Stakeholder Perspective on Safety Risk Factors of Metro Construction: A Social Network Analysis. *Complexity* **2020**, *2020*, 8884304. [CrossRef]
28. Zhou, C.; Kong, T.; Jiang, S.; Chen, S.; Zhou, Y.; Ding, L. Quantifying the evolution of settlement risk for surrounding environments in underground construction via complex network analysis. *Tunn. Undergr. Space Technol.* **2020**, *103*, 103490. [CrossRef]
29. Niu, K.; Fang, W.; Song, Q.; Guo, B.; Du, Y.; Chen, Y. An Evaluation Method for Emergency Procedures in Automatic Metro Based on Complexity. *IEEE Trans. Intell. Transp. Syst.* **2021**, *22*, 370–383. [CrossRef]
30. Chen, S.; Zhuang, D. Evolution and evaluation of the Guangzhou metro network topology based on an integration of complex network analysis and GIS. *Sustainability* **2020**, *12*, 538. [CrossRef]
31. Bernal, E.; del Rey, A.; Villardón, P. Analysis of madrid metro network: From structural to HJ-biplot perspective. *Appl. Sci.* **2020**, *10*, 5689. [CrossRef]
32. Moreno-Pulido, S.; Pavón-Domínguez, P.; Burgos-Pintos, P. Temporal evolution of multifractality in the Madrid Metro subway network. *Chaos Solitons Fractals* **2021**, *142*, 110370. [CrossRef]
33. Meng, Y.; Tian, X.; Li, Z.; Zhou, W.; Zhou, Z.; Zhong, M. Exploring node importance evolution of weighted complex networks in urban rail transit. *Phys. A: Stat. Mech. Its Appl.* **2020**, *558*, 124925. [CrossRef]
34. Yu, W.; Ye, X.; Chen, J.; Yan, X.; Wang, T. Evaluation indexes and correlation analysis of origination-destination travel time of Nanjing metro based on complex network method. *Sustainability* **2020**, *12*, 1113. [CrossRef]
35. Wang, J.; Ren, J.; Fu, X. Research on Bus and Metro Transfer from Perspective of Hypernetwork- A Case Study of Xi'an, China (December 2020). *IEEE Access* **2020**, *8*, 227048–227063. [CrossRef]
36. Wang, W.; Wang, Y.; Correia, G.; Chen, Y. A Network-Based Model of Passenger Transfer Flow between Bus and Metro: An Application to the Public Transport System of Beijing. *J. Adv. Transp.* **2020**, *2020*, 6659931. [CrossRef]
37. Han, Y.; Wang, S.; Ren, Y.; Wang, C.; Gao, P.; Chen, G. Predicting Station-Level Short-Term Passenger Flow in a Citywide Metro Network Using Spatiotemporal Graph Convolutional Neural Networks. *ISPRS Int. J. Geo-Inf.* **2019**, *8*, 243. [CrossRef]
38. Huh, J.; Seo, Y. Understanding Edge Computing: Engineering Evolution With Artificial Intelligence. *IEEE Access* **2019**, *7*, 164229–164245. [CrossRef]
39. Li, W.; Luo, Q.; Cai, Q. A Smart Path Recommendation Method for Metro Systems with Passenger Preferences. *IEEE Access* **2020**, *8*, 20646–20657. [CrossRef]
40. Bakıcı, T.; Almirall, E.; Wareham, J. A smart city initiative: The case of Barcelona. *J. Knowl. Econ.* **2013**, *4*, 135–148. [CrossRef]
41. Kolaczyk, E.D. *Statistical Analysis of Network Data*; Springer Science+Business Media: Berlin/Heidelberg, Germany, 2009.
42. Barabási, A.L. *Network Science*; Cambridge University Press: Cambridge, UK, 2016.
43. Van Mieghem, P. *Graph Spectra for Complex Networks*; Cambridge University Press: Cambridge, UK, 2011.
44. Newman, M.E.J. *Networks: An Introduction*; Oxford University Press: Oxford, UK, 2010.
45. Wu, J.; Barahona, M.; Tan, Y.J.; Deng, H.Z. Spectral measure of structural robustness in complex networks. *IEEE Trans. Syst. Man Cybern. Part A* **2011**, *41*, 1244–1252. [CrossRef]
46. Li, C.; Wang, H.; de Haan, W.; J, S.C.; Van Mieghem, P. The correlation of metrics in complex networks with applications in functional brain networks. *J. Stat. Mech. Theory Exp.* **2011**, *2011*, P11018. [CrossRef]
47. Jolliffe, I.T.; Cadima, J. Principal component analysis: A review and recent developments. *Philos. Trans. R. Soc. A Math. Phys. Eng. Sci.* **2016**, *374*, 20150202. [CrossRef] [PubMed]
48. Duntelman, G.H. *Principal Components Analysis*; Number 69 in Quantitative Applications in the Social Sciences; Sage Publications: Thousand Oaks, CA, USA, 1989.
49. Govender, P.; Sivakumar, V. Application of k-means and hierarchical clustering techniques for analysis of air pollution: A review (1980–2019). *Atmos. Pollut. Res.* **2020**, *11*, 40–56. [CrossRef]
50. Johnson, S.C. Hierarchical clustering schemes. *Psychometrika* **1967**, *32*, 241–254. [CrossRef]
51. Bouguettaya, A.; Yu, Q.; Liu, X.; Zhou, X.; Song, A. Efficient agglomerative hierarchical clustering. *Expert Syst. Appl.* **2015**, *42*, 2785–2797. [CrossRef]
52. Day, W.H.; Edelsbrunner, H. Efficient algorithms for agglomerative hierarchical clustering methods. *J. Classif.* **1984**, *1*, 7–24. [CrossRef]
53. Ward, J.H., Jr. Hierarchical grouping to optimize an objective function. *J. Am. Stat. Assoc.* **1963**, *58*, 236–244. [CrossRef]
54. Mojena, R. Hierarchical grouping methods and stopping rules: An evaluation. *Comput. J.* **1977**, *20*, 359–363. [CrossRef]
55. Blashfield, R.K. Mixture model tests of cluster analysis: Accuracy of four agglomerative hierarchical methods. *Psychol. Bull.* **1976**, *83*, 377. [CrossRef]

56. Murtagh, F.; Legendre, P. Ward's hierarchical agglomerative clustering method: Which algorithms implement Ward's criterion? *J. Classif.* **2014**, *31*, 274–295. [CrossRef]
57. Kaufman, L.; Rousseeuw, P.J. *Finding Groups in Data: An Introduction to Cluster Analysis*; John Wiley & Sons: Hoboken, NJ, USA, 2009; Volume 344.
58. Wu, X.; Tse, C.; Dong, H.; Ho, I.; Lau, F. A Network Analysis of World's Metro Systems. In Proceedings of the 2016 International Symposium on Nonlinear Theory and Its Applications (NOLTA2016), Yugawara, Japan, 27–30 November 2016; The Institute of Electronics, Information and Communication Engineers: Tokyo, Japan, 2016; pp. 606–609.
59. Frutos Bernal, E.; Martín del Rey, A. Study of the Structural and Robustness Characteristics of Madrid Metro Network. *Sustainability* **2019**, *11*, 3486. [CrossRef]
60. RStudio Team. *RStudio: Integrated Development Environment for R*; RStudio, PBC: Boston, MA, USA, 2020.



## Article

# Two Advanced Models of the Function of MRT Public Transportation in Taipei

You-Shyang Chen <sup>1,\*</sup>, Chien-Ku Lin <sup>2,3</sup> , Su-Fen Chen <sup>4,\*</sup> and Shang-Hung Chen <sup>1</sup>

<sup>1</sup> Department of Information Management, Hwa Hsia University of Technology, New Taipei City 235, Taiwan; trtc042@gmail.com

<sup>2</sup> Department of Business Management, Hsiuping University of Science and Technology, Taichung City 412, Taiwan; cklin@hust.edu.tw

<sup>3</sup> Department of Computer Science and Information Management, Hungkuang University, Taichung City 433304, Taiwan

<sup>4</sup> National Museum of Marine Science & Technology, Keelung City 202010, Taiwan

\* Correspondence: ys\_chen@go.hwh.edu.tw (Y.-S.C.); 10201a@gapps.nou.edu.tw (S.-F.C.)

**Abstract:** Tour traffic prediction is very important in determining the capacity of public transportation and planning new transportation devices, allowing them to be built in accordance with people's basic needs. From a review of a limited number of studies, the common methods for forecasting tour traffic demand appear to be regression analysis, econometric modeling, time-series modeling, artificial neural networks, and gray theory. In this study, a two-step procedure is used to build a predictive model for public transport. In the first step of this study, regression analysis is used to find the correlations between two or more variables and their associated directions and strength, and the regression function is used to predict future changes. In the second step, the regression analysis and artificial neural network methods are assessed and the results are compared. The artificial neural network is more accurate in prediction than regression analysis. The study results can provide useful references for transportation organizations in the development of business operation strategies for managing sustainable smart cities.

**Keywords:** passenger traffic; artificial neural network; regression analysis

**Citation:** Chen, Y.-S.; Lin, C.-K.; Chen, S.-F.; Chen, S.-H. Two Advanced Models of the Function of MRT Public Transportation in Taipei. *Electronics* **2021**, *10*, 1048. <https://doi.org/10.3390/electronics10091048>

Academic Editors: Juan M. Corchado, Josep L. Larriba-Pey, Pablo Chamoso and Fernando De la Prieta

Received: 2 April 2021

Accepted: 26 April 2021

Published: 29 April 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Taipei has the highest population density and traffic capacity in Taiwan. The construction of the Mass Rapid Transit (MRT) [1] system has relieved long-standing traffic problems in Taipei's urban area. In addition to its safety and convenience, the public MRT may also control the growth of private vehicles, reduce carbon emissions, and save energy, allowing for the creation of a low-carbon and green energy-based city. The preliminary road network of the public MRT system in the urban area of Taipei was approved by the Executive Yuan in 1986. The Taipei City Government established the Department of Rapid Transit Systems in 1987 and launched the construction of a preliminary road network or extension, following a revision. Taipei Rapid Transit Corporation was incorporated in 1994. The Muzha Line, the first driverless medium-capacity rapid transit line in Taiwan, was opened in March 1996, turning over a new leaf for public transportation in Taiwan. In March 1997, the first high-passenger-capacity system, the Tamsui Line, was opened, with a service scope that extended from Taipei City to New Taipei City. Following the continuous opening of road networks, 21 administrative regions across Taipei (with 12) and New Taipei City (with 9) all came to be included in the MRT routes after the opening of the Songshan Line in November 2014. According to the statistics of the department of the account of the Ministry of Transportation and Communications (MOTC), for traffic indicators from January to June 2015, the daily passenger capacity was 1.943 million on average, indicating that the MRT system is frequently used now and poses an interesting/important positive

issue relating to traffic volume forecasting. However, autonomous vehicles (or self-driving cars) pose another challenging public transport-related issue, and they have been the subject of much research in recent years. According to the study of Wiseman [2], although it is noteworthy that the overall consequences of autonomous vehicles are still unknown, their impacts on other means of transportation and on the transportation infrastructure remain unclear. In particular, the point at which autonomous vehicles penetrated the transportation market needs to be determined, and more information about them needs to be available, as this information will enable more profound studies about the future. More seriously, Wiseman [3] indicated that public transportation systems will soon be obsoleted due to the benefits of autonomous vehicles. While this serious topic is valuable for sustainability studies, it is not the focus of this study. This study is concerned with traffic forecasting.

It is problematic that the rising number of passengers affects the passenger traffic of each MRT station, resulting in the spread or concentration of some passengers in certain transit stations. Therefore, passenger traffic forecasting and analysis are required. In addition to providing a train control and transportation plan, the data on MRT passenger traffic, with the e-ticket or passenger traffic information about the surrounding High-Speed Rail (HSR) [4], Taiwan Railway, and city bus, may also be used as the basis of the public transportation system plan and operational management to enhance the efficiency and service quality of all public transportation tools and derive the maximum benefits from passenger traffic forecasting. Since the volume of passenger traffic has a significant impact on the transportation industry, past passenger traffic trends may be used as the basis for future operational plans and decisions in order to derive maximum benefit. The forecasting of MRT passenger traffic is correlated with the positive or negative economic development of the greater Taipei area in the future, while a precise forecast depends on a full understanding of the environment of urban areas, factors affecting demand, and the proper forecasting tool selection. The primary purpose is to collect past research data and sort out the best variables affecting the passenger traffic demand forecast from a review of the literature. The aims of this study are as follows: (1) to analyze important variables that affect the selection and forecasting of passenger traffic and develop a forecast model using regression analysis [5]; (2) to select important variables affecting MRT passenger traffic using an artificial neural network [6] and create a forecast model; and (3) to compare the results of the regression analysis and artificial neural network research, evaluate the forecast performance of the two projection methods, and then select the best forecasting model.

In this study, a two-step procedure is used to build a predictive model for public transport. In the first step of this study, regression analysis is used to find the correlations between two or more variables and their associated directions and strength, and the regression function is used to predict future changes. In the second step, the regression analysis and artificial neural network methods are assessed, and the results are compared.

## 2. Literature Review

This section introduces the passenger traffic forecast and application, selection of variables, artificial neural network, regression analysis, and mean absolute percentage error (MAPE) evaluation indicators.

### 2.1. Passenger Traffic Forecasts and Application

Passenger traffic forecasts are an important basis for developing a traffic capacity and transportation plan for a transportation organization, as well as the foundation for constructing and modifying various forms of transportation equipment. The transportation demand and qualitative and quantitative features determine the transportation supply plan. Some common approaches to forecasting passenger traffic demand include regression analysis, econometric modeling, time-series modeling, artificial neural networks, and gray theory. This study adopts an artificial neural network to create the model. The artificial neural network is more flexible in terms of its parameter design without limitations relating

to statistical assumptions, it is highly capable of learning, and it is available to precisely forecast passenger traffic under the conditions of a sufficient number of training samples and appropriate parameter design. The other research model used is regression analysis. The approach adopts mathematical statistics to create independent and dependent variables based on massive observation data in order to understand whether two or more variables are correlated and their correlative direction and strength before establishing a regression function to predict future movement.

The literature related to passenger traffic discussed in this study includes the following: Mazanec [7] suggested that the artificial neural network model is better than discriminate analysis. Furthermore, the artificial neural network can also classify a group of input variables (such as the social, economic, demographic, or behavioral attributes) of train passengers, including output data. Law and Au [8] predicted the passenger traffic of the air route between Hong Kong and Japan, using the GDP, foreign exchange, population, marketing expenses, service price, and hotel rate as input variables. The authors randomly picked 20 trained datasets among the data from the period of 1967–1997 and another 10 datasets for comparison. In the meantime, the MAPE, normalized correlation coefficient, and acceptable output percentage were used as the standards to compare the actual value with the forecast value derived from an artificial neural network, exponential smoothing, the moving average, and regression analysis. As a result, the MAPE value of an artificial neural network is only 10%, meaning that it has the best forecast result.

Kulendran and Witt [9] used the leading indicator transfer function (TF) to predict the needs for travel from England to six other countries (including the U.S., Germany, the Netherlands, Spain, Portugal, and Greece) in the period of 1993–1995. The authors first selected the six most common economic indicators for travel testing—relative price, exchange rate, relative exchange, domestic real income, and GDP—to find the correlation between the need to travel to these six countries and all variables. Next, they found the formula, using data from 1978 to 1992 as the input variable and the passenger number as the output variable to forecast the demand in the period of 1993–1995. Meanwhile, the MAPE was used as the criteria to compare the results of the autoregressive integrated moving average (ARIMA) and error correlation models (ECMs). The results show that ARIMA is more precise than TF when the quarter unit is one, four, and eight, while TF is more precise under other circumstances. On the other hand, the ECMs are more precise than TF in terms of long-term conditions (by a unit of four quarters and eight quarters). Grosche et al. [10] used two gravity models to predict the passenger traffic between two cities. The input variables of the first gravity model were population, consumption ability, GDP, distance, and average travel time. The second one was the extended gravity model, which added the number of competitive airports and average distances of competitive airports as input variables. They were verified by the data on Germany and 29 other European countries, and the results show that both gravity models were able to precisely predict the passenger traffic volume.

## 2.2. Variable Selection

The input variables (independent variables) adopted in this study are summarized in Table 1 based on a literature review and an analysis of the demand for public transportation. X1–X20 are independent variables based on the total passenger number of the Taipei MRT, i.e., the output variables. Among them, each datum shall be based on the “year”.

**Table 1.** Variable selection.

Variable	Name of Variable	Variable Code
Y1	MRT passenger traffic	Passengers
X1	Gross National Product	GNP
X2	Gross Domestic Product	GDP
X3	Economic growth rate	Economic
X4	Number and density of registered financial businesses	InFinanden
X5	Number of major tourists	Visitors
X6	Number of listed companies	Companies
X7	Employment population	Employed
X8	Consumer Price Index	CPI
X9	Population density	InPopden
X10	Personal income	PCI
X11	Import Price Index	IPI
X12	Total cargo input amount	CargoInput
X13	MRT train trips	Trains
X14	Number of administrative regions	Administrative
X15	MRT train kilometers	TrainsKM
X16	MRT passenger traffic income	Income
X17	Number of MRT stations	Stations
X18	MRT operation kilometers	OperationsKM
X19	MRT passenger kilometers	PassengersKM
X20	Round trip transfer discount volume	Transfers

The input variables of the study are shown in Table 1, and their definitions are as follows:

1. Y1: The training samples are from the passenger traffic data of the Muzha Line from 1996 to 2015, and the artificial neural network and regression analysis are adopted to forecast the Taipei MRT passenger traffic.
2. X1: The “nominal market value” of “final products and services” produced by “all people in a specific region” within a “certain period”. A prosperous economy will lead the passenger traffic.
3. X2: The total market value of all the “services” and “final products” produced by a country within a certain period. In general, a higher number means higher productivity in the country, which means it is more prosperous and helpful in relation to passenger traffic.
4. X3: The most important economic indicator to determine the macroeconomic situation, which may reflect the change in the scale of the total economic output, and the most physical symbol of economic prosperity, which may affect the capacity of public transportation tools.
5. X4: The higher the density of a financial industry, the more prosperous the business development, which may increase the demand for MRT.
6. X5: The tourism business is an important service industry that contributes much to the national income, employment, and foreign exchange earnings. Meanwhile, the number of tourists is also helpful in evaluating the increase in passenger traffic.
7. X6: The good operation of listed companies with sufficient working capital may promote business activities and increase the demand for transportation.
8. X7: A local employee population may reflect the level of local economic prosperity. The people of a country with greater employment opportunities may have a certain consumption ability and will affect the public transportation volume.
9. X8: The price variation indicator, calculated based on the prices of products and services relating to living in a residence, is one of the major indicators for measuring inflation. The increase in price will affect the public transportation volume.
10. X9: The total population density of Taipei and New Taipei City, obtained using the land area of Taipei and New Taipei City divided by the population, is adopted as the variable. The population will affect the consumption ability of the public.

- The consumption ability will automatically increase with an increasing population. Therefore, the population may affect the demand for transportation in a certain area.
11. X10: Personal income affects the consumption ability of the public. The higher the income, the greater the consumption ability. Therefore, personal income may affect the transportation demand.
  12. X11: This factor measures the change in prices of imported and exported products. The higher the consumption amount of the public, the stronger the consumption ability of the nationals in the country, which may affect the passenger traffic.
  13. X12: The consumption of the public in the country could be observed from the total amount imported to Taiwan. A higher amount means a stronger consumption ability of the people in the country, which may affect the passenger traffic.
  14. X13: The density of the total MRT train trips and time spent waiting for mass transportation are important factors affecting the willingness of commuters to take public transportation. The shorter the waiting time, the higher the willingness of commuters to take public transportation.
  15. X14: The accessibility of the transportation network to administrative regions may likely affect the willingness of the public to take public transportation. If different mass transportation tools can be effectively integrated, commuters will be more willing to take public transportation.
  16. X15: The total travel mileage of all MRT trains within a specific time and period. The higher the passenger traffic is, the more train trips and travel mileage of trains there will be.
  17. X16: Operation income status is critical for the sustainable operation of a transportation organization. If the organization has operation losses, its transportation supply efficiency must be affected.
  18. X17: The actual number of stops in all the routes of the MRT operation and the expansion of the transportation network will increase the accessibility of the expected destinations. If different mass transportation tools can be effectively integrated, commuters will be more willing to take public transportation.
  19. X18: Shorter distances of transportation routes will increase the accessibility of the planned destinations. If different mass transportation tools can be effectively integrated, commuters will be more willing to take public transportation.
  20. X19: The total number of passengers multiplied by the mileage traveled. The mileage for the calculation of passenger kilometers will be based on the mileage for the freight calculation.
  21. X20: The changes in MRT, bus transportation, and transfer volume are used to analyze the effects of decreases in transfers. The time spent waiting for public transportation tools is an important factor affecting the willingness of commuters to take public transportation. The shorter the waiting time, the higher the willingness of the commuters to take public transportation.

### 2.3. Artificial Neural Network

An artificial neural network [11] is a parallel calculation model, which is similar to the human neural structure. It is an information-processing technology inspired by brain and nervous system research, and it is also usually referred to as a parallel distributed processing model or a connectionist model. An artificial neural network uses a lot of simply connected artificial neurons [12] to simulate the ability of the biological neural network. An artificial neuron is a simple simulation of a biological neuron, which acquires information from the external environment or other artificial neurons. It performs very simple calculations and then exports the results to the external environment or other artificial neurons for further action. The basic structure of a general neural network is divided into the neuron, layer, and network parts. The layer consists of basic neurons, and the network is constituted by layers. The neuron is also called the processing unit and the basic unit of the artificial neural network. The operation model of an artificial neural

network is mainly divided into the training phase and the recall phase: the training phase refers to learning from training to adjust the weight of the network so that the network can become stable. The recall phase involves determining the output value induced by the network to test whether the output is close to the target output value. The artificial neural network adopted by this study is the back-propagation network [13] in a supervised learning network [14], which is applicable to classification, forecasting, system control, noise filtering, sample identification, and data compression. The input layer units are different in each step. The number of processing units in the hidden layer is determined by the number of input and output layers, and the final output result will be either one or two. The back-propagation network is the most representative and popular model among the currently available artificial neural networks. The basic theory of artificial neural networks involves minimizing the error function, using the gradient steepest descent method (GSDM) [15] to achieve the learning purpose.

#### 2.4. Regression Analysis Method

Regression analysis [16] is a method involving the analysis of data in statistics. It is mainly used to analyze whether there is a specific relationship between one or more independent variables and dependent variables. Regression analysis is a model for establishing the relationship between a response variable  $Y$  and independent variables  $X$ . The purpose is to understand whether two or more variables are related and their correlative direction and strength, as well as to establish a mathematical model that allows for the observation and prediction of specific variables. Since the purpose of prediction regression is not to clarify but rather to establish the best formula, the primary consideration in variable selection is whether there is a maximum practical value, as opposed to the theoretical appropriateness. The theory primarily explains the value of the regression model in practical applications and its mechanism for solving problems in prediction regression. It is expected to achieve the maximum practical value with the lowest cost. The first job of explanation regression is to carefully review the features and relationships among all variables, that is, to examine the correlation among the variables.

#### 2.5. Mean Absolute Percentage Error

MAPE [17–19] refers to the average absolute percentage error, which is the evaluation indicator for whether a prediction model is good or bad. Since MAPE is a relative value that is not affected by the measurement value and estimated value, it can observe the difference between the estimated and evaluated values objectively. The estimation effect is better if the MAPE value is closer to 0. The standards to evaluate the precision of a forecast are shown in Table 2.

**Table 2.** Standards for the precision of MAPE evaluation.

MAPE Value	Standard
MAPE < 10%	Excellent prediction power (the closer to 0, the better)
10% < MAPE < 20%	Good prediction power
20% < MAPE < 50%	Reasonable prediction power
50% < MAPE	Poor prediction power

### 3. Materials and Methods

This section introduces the research structure and research design and steps.

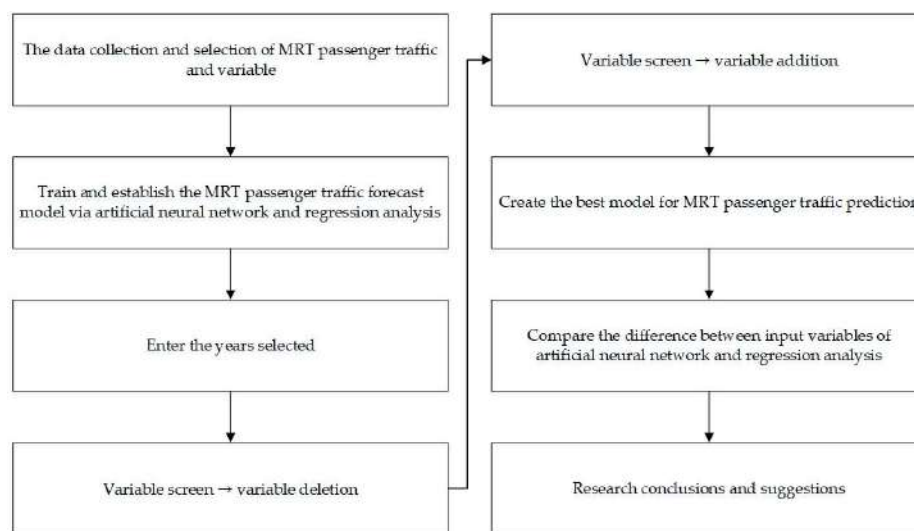
#### 3.1. Research Structure

The flow of this study is shown in Figure 1, below, and the relevant steps are as follows:

1. The literature and selected variables that may affect the MRT passenger traffic are collected via a literature review.
2. The artificial neural network and regression analysis are adopted as research methods for training and establishing a prediction model of the MRT passenger traffic: (1) the

years of training are selected; (2) the input variables of the MRT passenger traffic are deleted one by one, and the best variable is found; (3) the possible input variables that are likely to affect the MRT passenger traffic in the training are added, and the best variable is found; and (4) the MRT passenger traffic forecast model is established via an artificial neural network and regression analysis.

3. The MRT passenger traffic forecast model is established.
4. The advantages and disadvantages of the input variables selected are compared and analyzed, considering the results of the artificial neural network and regression analysis.
5. The research conclusions and suggestions based on the research results are provided.



**Figure 1.** Research structure.

### 3.2. Research Design and Steps

We first collected the relevant literature and selected variables that could potentially affect the MRT passenger traffic via a literature review, as shown in Table 1, and then selected the years of training, as shown in Table 3.

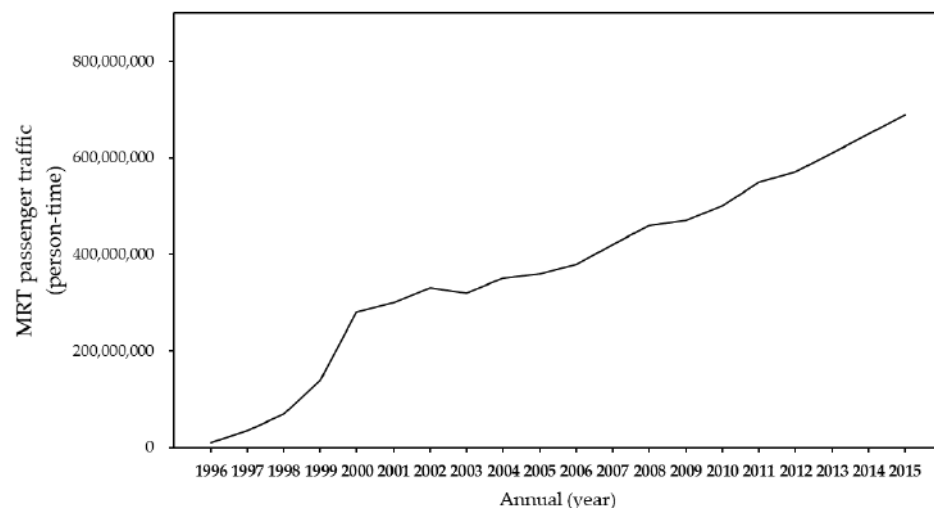
The main structure of this study is established by the idea of an artificial neural network structure, including the input layer, hidden layer, and output layer. X1 to X10 are variables that may affect the Taipei MRT passenger traffic, while Y1 is the Taipei MRT passenger traffic. We first tested independent variables X1 to X10 and then carried out the deletion test, before adding independent variables X11 to X20 one by one. The purpose of this step was to screen the input variables with a lower correlation to the Taipei MRT passenger traffic to avoid interfering with the prediction results. The regression analysis is primarily used to discuss the causal result relationship among the variables and conduct the prediction via a line chart [20,21], a scatter chart [22], correlation analysis [23–25], the enter method [26,27], and stepwise regression [28–30].

**Table 3.** Annual information of the variables in relation to passenger traffic.

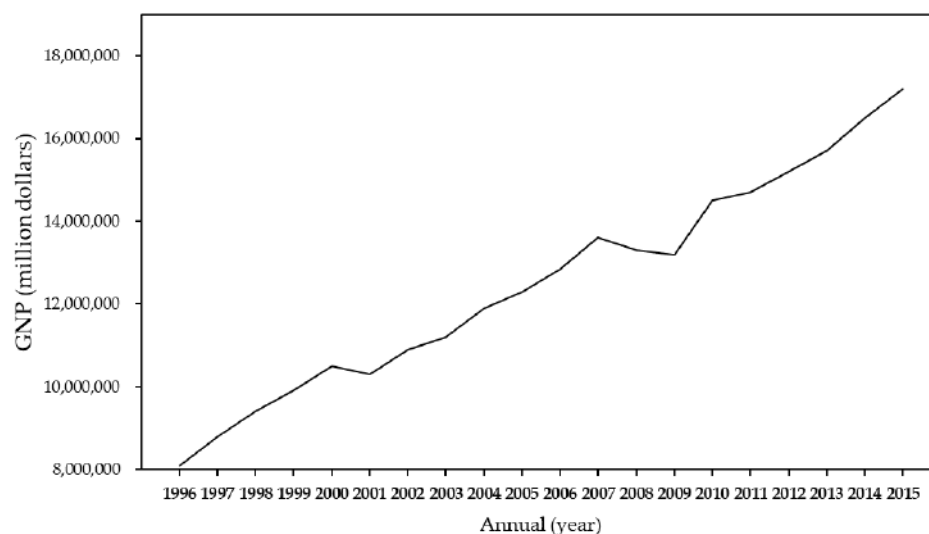
Variable	X1	X2	X3	X4	...	X17	X18	X19	X20	Y1
Unit/Year	Million	Million	%	Km <sup>2</sup> (Train)	...	Station	Km	Persons/Km	Persons (Thousand)	Total Persons
1996	8,146,092	8,036,590	6.18	500	...	12	10.5	57,226,810	102	11,174,359
1997	8,806,852	8,717,241	6.11	509	...	32	32.4	243,676,517	3282	31,081,395
1998	9,449,692	9,381,141	4.21	501	...	39	40.3	512,282,678	12,229	60,737,782
1999	9,906,113	9,815,595	6.72	495	...	56	56.4	1,031,342,472	21,203	126,952,122
2000	10,490,818	10,351,260	6.42	491	...	62	65.1	2,042,303,171	38,138	268,716,740
2001	10,350,233	10,158,209	−1.26	448	...	62	65.1	2,223,486,596	44,368	289,642,714
2002	10,923,385	10,680,883	5.57	437	...	62	65.1	2,469,037,312	53,093	324,433,557
2003	11,294,739	10,965,866	4.12	433	...	62	65.1	2,440,488,934	72,399	316,189,128
2004	12,021,744	11,649,645	6.51	428	...	63	67	2,680,355,529	125,350	350,141,956
2005	12,383,120	12,092,254	5.42	422	...	63	67	2,742,372,258	127,424	360,729,803
2006	12,952,502	12,640,803	5.62	415	...	69	74.4	3,002,988,957	130,916	384,003,220
2007	13,739,828	13,407,062	6.52	417	...	69	74.4	3,298,870,463	140,044	416,229,685
2008	13,465,596	13,150,950	0.70	418	...	70	75.8	3,513,969,060	152,643	450,024,415
2009	13,375,650	12,961,656	−1.57	428	...	82	90.6	3,720,991,244	153,679	462,472,351
2010	14,548,852	14,119,213	10.63	425	...	93	100.8	4,123,189,518	162,098	505,466,450
2011	14,700,572	14,312,200	3.80	425	...	94	101.9	4,607,801,411	170,963	566,404,489
2012	15,141,108	14,686,917	2.06	427	...	102	112.8	4,973,666,983	177,918	594,864,715
2013	15,646,211	15,221,201	2.23	428	...	109	121.3	5,234,160,050	182,193	634,961,083
2014	16,621,378	16,081,798	3.74	428	...	116	129.2	5,589,414,250	180,133	679,506,401
2015	17,485,375	16,881,614	3.78	429	...	117	131.1	5,880,980,257	173,877	717,511,809

### 3.2.1. Line Chart and Scatter Chart

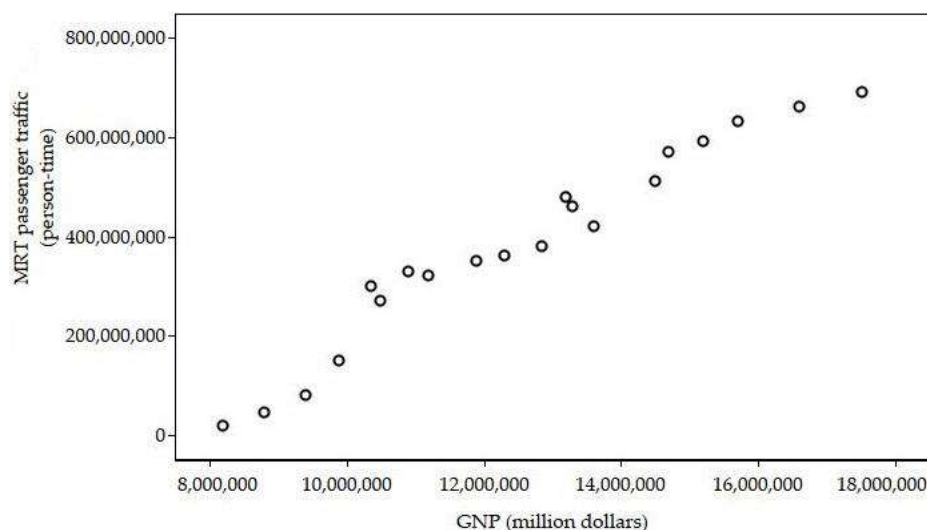
The line chart demonstrates that the growth trend of the Taipei MRT passenger traffic (Figure 2) is relatively similar to the variable X1 GNP (Figure 3), and its effect is therefore assumed to be more important. From the scatter chart (Figure 4), if the variable X1 GNP and the Taipei MRT passenger traffic have a linear distribution, they should be more correlated, and the effect of variable X1 GNP on the Taipei MRT passenger traffic will be obvious.

**Figure 2.** The line chart of the total Taipei MRT passenger traffic.





**Figure 3.** The line chart of the annual GNP.



**Figure 4.** The scatter chart of the MRT passenger traffic and GNP.

### 3.2.2. Correlation Analysis

In correlation analysis, if two variables are significantly correlated, this only means that the strength and direction between the variables are significant. When the coefficient is significant, it only explains that the two variables are correlated at a certain level, including the strength and direction, instead of indicating the existence of a causal relationship. They may both be the cause and result at the same time, or a causal relationship may actually exist.

### 3.2.3. Enter Method

The enter method means that all the prediction variables must be entered into the regression formula, regardless of the significance of the individual variable.

### 3.2.4. Stepwise Regression

When variables are entered into the regression equation, the backward selection method is used to eliminate unimportant variables.

## 4. Empirical Analysis

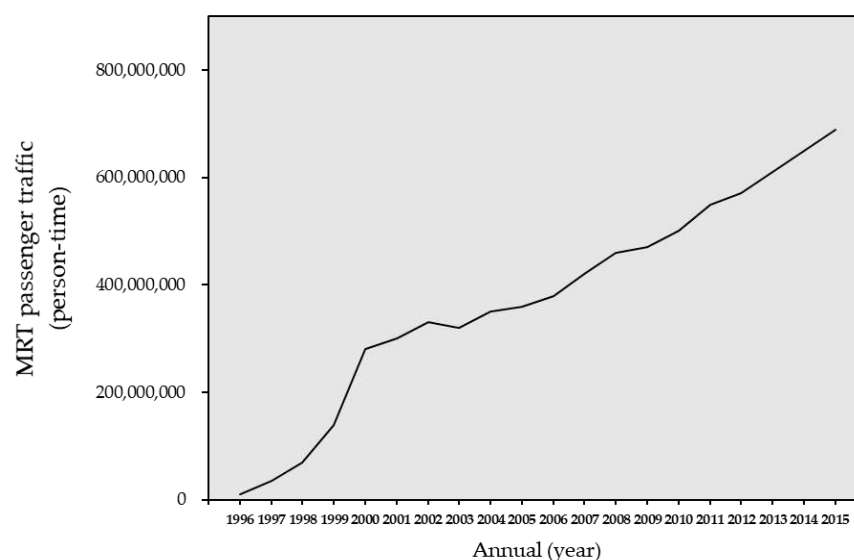
This study uses regression analysis and class neural network methods to select and predict the variables affecting passenger traffic using the Taipei Metro passenger traffic trends from March 1996 to December 2015. Finally, we compare the results of the neural network and regression analysis prediction methods, evaluate and analyze the prediction performance of the two methods, and then select the best prediction model.

### 4.1. Regression Analysis

We adopted statistical software to analyze 20 parametric variables from March 1996 to December 2015 to explore the causal relationship between the variables as well as to make predictions through line graphs, scatter plots, correlation analysis, forced entry, and stepwise regression analysis to explore the causal relationship between the variables.

#### 4.1.1. Line Chart

If the line graphs of the variables are found to be more similar to the growth trend line graphs of Taipei Metro's passenger volume, it is speculated that the effect should be more important (Figure 5).



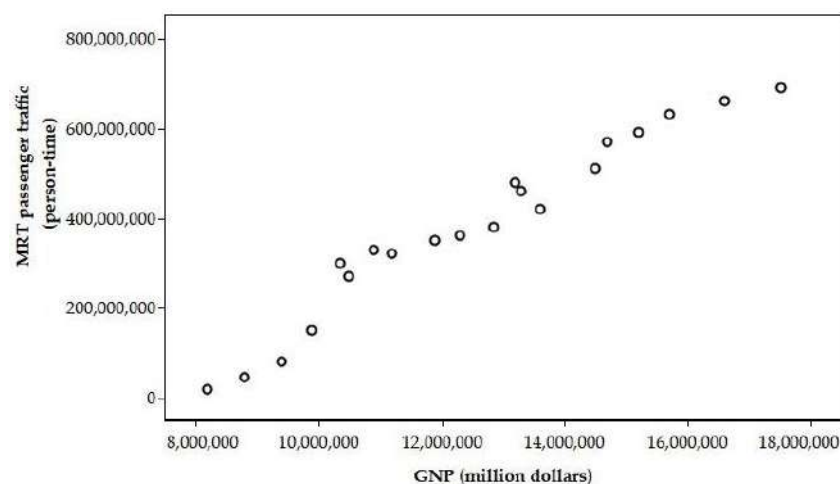
**Figure 5.** Line chart of the MRT passenger traffic.

#### 4.1.2. Scatter Chart

From the scatter plot, we can see that the correlation between the variables and Taipei Metro's passenger volume is greater if they are clearly distributed in a straight line (Figure 6).

#### 4.1.3. Correlation Analysis

The correlation coefficient primarily refers to the degree of correlation between the variables, while it does not verify the influence of the “independent variable” on the “dependent variable”. Therefore, the obtained correlation coefficient (R value) can only indicate that the two variables are positively correlated, negatively correlated, or independent. It cannot be interpreted as the effect of the independent variable on the dependent variable. In the interpretation of the correlation coefficient, positive and negative indicate the direction of the correlation, not the degree of the correlation. If the degree of the correlation is between the R values, i.e., between plus or minus 0.3 (i.e., between 0.3 and  $-0.3$ ), it is considered a low correlation; it is considered a moderate correlation if the value is between plus or minus 0.3 to 0.6 (that is, between 0.3 and 0.6, or between  $-0.3$  and  $-0.6$ ); and if it is between plus or minus 0.6 to 0.9 (i.e., 0.6 to 0.9, or  $-0.6$  to  $-0.9$ ), it is considered a high correlation. If the R value is plus or minus 1, this indicates a complete correlation.



**Figure 6.** Scatter chart of GNP and MRT passenger traffic.

#### 4.1.4. Entry Method

All the predictive variables are incorporated into the regression equation at the same time, and the least squares method and statistical software are used to calculate the complex regression model. The analysis of the explanatory variables is as follows. Table 4, showing the entry variables inputted/removed, provides a list of variables that are subjected to regression analysis, including a total of 16 independent variables.

**Table 4.** Forced entry method variables inputted/removed.

Model	Variable Inputted	Variable Removed	Method
1	X3, X4, X5, X6, X7, X8, X9, X10, X11, X12, X13, X14, X15, X16, X17, X20		Entry

The variance analysis result of the forced entry method shows the overall verification of the model significance to verify the significance of the overall regression model. The quadratic sum of the regression model is 830,852,340,051,751,420, the total sum of squares is 830,902,477,461,923,070, the F value = 3107.157, and the  $p$  value = 0.000 < 0.05, which reaches significance. Furthermore, F is used to verify the regression model,  $F(\alpha, K, N-K-1) = F(0.05, 16, 3) \approx 8.7$ , and the F value of this regression equation = 3107.157 > critical value F value of 8.7, thus rejecting the null hypothesis and indicating that the explanatory power of the overall regression model reaches a significant level.

In Table 5, the coefficients of the forced entry method are the verification results of the multicollinearity of the output independent variables using statistical software. This model has a total of 16 independent variables and 17 eigenvalues, and the minimum tolerance is 0.000. If the tolerance is lower than 0.2, it is determined that there is collinearity between this independent variable and other independent variables. The highest VIF is 4167.205. In general, if the VIF is greater than 4, this means that the independent variable is collinear with other independent variables.

The variables excluded by the forced entry method demonstrate that there are four variables, including the total GNP, GDP, MRT mileage, and MRT extended passenger mileage, that did not enter the regression model. To reduce the collinearity problem of the multiple regression model, we adopted the stepwise method. The advantage of this method is that the variables are selected one by one according to the influence of the independent variable on the dependent variable, thus eliminating the collinearity problem.

**Table 5.** Coefficients of the entry method.

Model		T	Significance	Collinearity Statistical Material	
				Permissible Deviation	VIF
1	(Constant)	−2.374	0.098		
	X1	0.685	0.543	0.094	10.601
	X3	2.787	0.069	0.004	263.957
	X4	1.500	0.231	0.002	457.889
	X5	1.429	0.248	0.001	1047.998
	X6	3.584	0.037	0.000	4167.205
	X8	−1.309	0.282	0.001	835.575
	X9	−3.509	0.039	0.001	1728.353
	X10	−1.115	0.346	0.007	143.916
	X11	2.919	0.062	0.005	206.473
	X12	1.151	0.333	0.007	142.958
	X13	−3.006	0.057	0.003	293.111
	X14	−1.913	0.152	0.003	313.741
	X15	0.030	0.978	0.002	617.120
	X16	11.990	0.001	0.003	374.019
	X17	−0.098	0.928	0.002	656.926
	X20	−2.325	0.103	0.003	348.165

#### 4.1.5. Stepwise Method

To reduce the collinearity problem of multiple regression models, we adopted the stepwise method and performed calculations using statistical software. The following Table 6 was obtained.

**Table 6.** Stepwise regression variables inputted/removed.

Model	Variable Inputted	Method
1	X16	Step by step (criteria: F-to-enter probability $\leq 0.050$ , F-to-remove probability $\geq 0.100$ )
2	X20	Step by step (criteria: F-to-enter probability $\leq 0.050$ , F-to-remove probability $\geq 0.100$ )

Table 6 shows the stepwise regression variables inputted/removed and a list of stepwise method variables. The criteria for selection are: (F-to-enter probability  $\leq 0.050$ , F-to-remove probability  $\geq 0.100$ ). In total, two variables are selected for the regression equation in two steps (models). In Model 1, MRT passenger revenue is selected. In Model 2, two-way transfer preferential volume is added. Therefore, the two variables of the MRT passenger revenue and two-way interchange preferential volume are selected for the model.

The analysis result of stepwise regression variance using Model 1 of the MRT passenger revenue variables shows that the F value of the overall regression model is 15,598.145,  $p = 0.000 < 0.05$ , reaching significance and indicating a significant correlation between the independent variable (the MRT passenger revenue) and the dependent variable (the Taipei MRT passenger traffic). Taking Model 2 as an example, the F value of the regression model of Model 2 is 12,741.962,  $p = 0.000 < 0.05$ , which reaches significance and indicates a significant correlation between the independent variable (the preferential volume of two-way interchange) and the dependent variable (the Taipei MRT passenger traffic). The two variables can effectively predict the passenger traffic of the Taipei MRT.

Table 7 demonstrates a positive relationship between the MRT passenger revenue and the preferential volume of the two-way passenger transfer—two independent variables in the table of stepwise regression coefficients, both of which were consistent with the prediction and significant.

Table 7. Stepwise regression coefficients.

	Model	Unstandardized Coefficients		Standardized Coefficients	T	Significance
		B	Standard Error	Beta		
1	(Constant)	−12,082,965.667	3,520,486		−3.432	0.003
	X16	0.047	0	0.999	124.893	0
2	(Constant)	−9,989,402.575	2,818,390		−3.544	0.002
	X16	0.044	0.001	0.942	53.482	0
	X20	193.099	54.85	0.062	3.52	0.003

Table 8 demonstrates the stepwise mode analysis results of the variable table excluded by stepwise regression, with two variables remaining: the MRT passenger revenue and two-way interchange preferential volume. From Table 8, it can also be understood that the variables may enter the regression model because the correlation value of the dependent variables is the highest in Model 1, where the preferential volume of two-way interchange is 0.649. In Model 2, the correlation value of the total amount of goods imported into Taiwan is −0.310, which is the highest.

Table 8. Variables excluded by the stepwise method.

	Model	Beta	T	Significance	Partial Correlation
1	X1	0.031	0.799	0.436	0.190
	X2	0.024	0.641	0.530	0.154
	X3	0.007	0.780	0.446	0.186
	X4	−0.031	−2.705	0.015	−0.549
	X5	−0.020	−1.234	0.234	−0.287
	X6	0.100	2.949	0.009	0.582
	X7	0.085	2.668	0.016	0.543
	X8	0.009	0.371	0.715	0.090
	X9	0.080	1.411	0.176	0.324
	X10	0.002	0.071	0.944	0.017
	X11	0.025	1.981	0.064	0.433
	X12	0.016	0.782	0.445	0.186
	X13	−0.008	−0.204	0.841	−0.049
	X14	0.012	0.481	0.637	0.116
	X15	−0.039	−0.831	0.418	−0.197
	X17	−0.058	−1.788	0.092	−0.398
	X18	−0.063	−1.873	0.078	−0.414
	X19	0.220	1.473	0.159	0.336
	X20	0.062	3.520	0.003	0.649
2	X1	−0.026	−0.750	0.464	−0.184
	X2	−0.030	−0.900	0.382	−0.219
	X3	0.004	0.562	0.582	0.139
	X4	−0.015	−1.232	0.236	−0.294
	X5	−0.007	−0.519	0.611	−0.129
	X6	0.050	1.239	0.233	0.296
	X7	0.021	0.471	0.644	0.117
	X8	−0.017	−0.804	0.433	−0.197
	X9	−0.076	−1.158	0.264	−0.278
	X10	−0.018	−0.708	0.489	−0.174
	X11	−0.007	−0.447	0.661	−0.111
	X12	−0.026	−1.306	0.210	−0.310
	X13	−0.032	−1.015	0.325	−0.246
	X14	−0.011	−0.519	0.611	−0.129
	X15	−0.024	−0.637	0.533	−0.157
	X17	−0.027	−0.945	0.359	−0.230
	X18	−0.033	−1.093	0.290	−0.264
	X19	0.090	0.698	0.495	0.172

We used statistical software to analyze the inputting of 20 parameter variables from 1996 to 2013 and outputting of the predictive value of the passenger traffic of the Taipei MRT (Table 9). The predicted value was calculated according to the predicted value table of the progressive regression analysis using the MAPE formula. The result obtained using the regression analysis prediction method was MAPE = 6.47%.

**Table 9.** Predictive value of regression analysis.

Year of Observation	Actual Number of MRT Passengers	Predictive Value of Regression Analysis	Residual Value
1996	11,174,359	3,305,117	7,869,242
1997	31,081,395	36,987,970	−5,906,575
1998	60,737,782	68,054,760	−7,316,978
1999	126,952,122	129,264,711	−2,312,589
2000	268,716,740	270,052,789	−1,336,049
2001	289,642,714	287,079,395	2,563,319
2002	324,433,557	319,060,611	5,372,946
2003	316,189,128	312,874,254	3,314,874
2004	350,141,956	351,632,592	−1,490,636
2005	360,729,803	359,697,272	1,032,531
2006	384,003,220	385,978,346	−1,975,126
2007	416,229,685	421,166,927	−4,937,242
2008	450,024,415	449,687,066	337,349
2009	462,472,351	457,533,201	4,939,150
2010	505,466,450	495,603,126	9,863,324
2011	566,404,489	560,822,898	5,581,591
2012	594,864,715	607,214,292	−12,349,577
2013	634,961,083	638,210,638	−3,249,555

#### 4.2. Artificial Neural Network

Using the training data from 1996 to 2013 as the artificial neural network training materials and those from 2014 to 2015 as the sample data, we were able to use a total of 20 parameters as the input variables of the artificial neural network; the output variables were the passenger traffic prediction values of the Taipei MRT. Statistical software was used to analyze 20 parametric variables from 1996 to 2015, before the input variables were converted into an artificial neural network to investigate the causal relationship between the variables and make predictions based on a line chart, scatter chart, correlation analysis, forced entry method, and stepwise method. Therefore, the artificial neural network does not temporarily delete any variables and inputs all of them into the artificial neural network. Then, the optimal combination of input variables was analyzed.

##### 4.2.1. Year of Selection

Basic information on the year of selection for variables X1 to X20 and the output variable Y1 is shown in Table 3, where Y1 stands for the Taipei MRT passenger traffic. Using the backward propagation artificial neural network algorithm model, the average absolute percentile error was used to analyze 7 years/2 years of training, 6 years/2 years of training, and 5 years/2 years of training from 1996 to 2015, and the training results are shown in Table 10. This study uses the training model of 5 years/2 years of training prediction for empirical research.

**Table 10.** Year selection of the artificial neural network.

Evaluation Indicators	7 Years/2 Years of Training (1996–2002/2003–2004)	6 Years/2 Years of Training (2003–2008/2009–2010)	5 Years/2 Years of Training (2009–2013/2014–2015)
MAPE	8.26%	4.35%	2.45%

#### 4.2.2. Deleted Variables

First, we used the 5-year and 2-year training modes (2009 to 2013 and 2014 to 2015) and then deleted the X1 to X10 variables in order, selecting the variables with the highest MAPE value for deletion. The purpose of this was to screen out variables that have little influence on the passenger traffic of the Taipei MRT to avoid interference with the predicted results. After the deletion was complete, we then started the training process to determine the best variables and improve the accuracy of the prediction results.

#### 4.2.3. The Output Results of the Artificial Neural Network

Using the training data from 1996 to 2013 as the artificial neural network training materials and the data from 2014 to 2015 as the sample data, we used a total of 20 parameters as the input variables of the artificial neural network; the output variables were the passenger traffic prediction values of the Taipei MRT (Table 11).

**Table 11.** Predictive values of the artificial neural network.

Year of Observation	Actual Number of MRT Passengers	Predictive Values of the Artificial Neural Network	Residual Value
1996	11,174,359	15,320,813	4,146,454
1997	31,081,395	30,226,981	−854,414
1998	60,737,782	60,188,236	−549,546
1999	126,952,122	126,305,437	−646,685
2000	268,716,740	267,922,813	−793,927
2001	289,642,714	289,012,656	−630,058
2002	324,433,557	324,181,085	−252,472
2003	316,189,128	327,612,737	11,423,609
2004	350,141,956	349,893,508	−248,448
2005	360,729,803	360,534,282	−195,521
2006	384,003,220	383,858,990	−144,230
2007	416,229,685	415,952,914	−276,771
2008	450,024,415	285,870,926	−164,153,489
2009	462,472,351	462,125,697	−346,654
2010	505,466,450	505,270,408	−196,042
2011	566,404,489	565,977,673	−426,816
2012	594,864,715	616,442,418	21,577,703
2013	634,961,083	630,207,775	−4,753,308

According to the prediction values, shown in Table 10, of the artificial neural network, calculated by the MAPE formula, the result obtained by the neural network prediction method is  $\text{MAPE} = 4.82\%$ . To compare the predicted values of regression analysis with those of the artificial neural network, the results showed that the MAPE value of the regression analysis was  $6.47\%$ . The MAPE value of the artificial neural network was  $4.82\%$ , so the artificial neural network has the best prediction outcomes.

#### 4.3. Comparison of Regression Analysis and Artificial Neural Networks

##### 4.3.1. Deletion of the Forced Entry Method of Regression Analysis to Exclude Four Variables

Using the 5-year training for 2 years (2009–2013/2014–2015) model, we deleted the four eliminated variables of the regression analysis, which were the GNP, GDP, MRT mileage, and MRT extended passenger mileage. Using a total of 16 parameters as the input variables for the regression analysis and artificial neural network comparison, it was found that the MAPE value predicted by regression analysis was  $0.94\%$ , and the MAPE value of the artificial neural network was  $0.54\%$ . Both methods showed the best prediction results.

#### 4.3.2. Deletion of Four Variables with High MAPE Values in Artificial Neural Networks

Using the model of 5-year training for 2 years (2009–2013/2014–2015), we deleted the four variables with the highest MAPE values in the artificial neural network, which included the economic growth rate, personal income, total amount of goods imported into Taiwan, and number of MRT stations. In total, 16 parameters were used as the input variables for the regression analysis and artificial neural network comparison, and it was found that the MAPE value predicted by regression analysis was 0.62%. The MAPE value of the artificial neural network was 0.31%. Both methods showed the best prediction results.

#### 4.3.3. Collinearity Verification

In order to verify whether the independent variable of the MRT passenger revenue is collinear with the dependent variable of the MRT passenger traffic, statistical software was used to analyze the inputting of 19 parameter variables from 1996 to 2013 and outputting of the passenger traffic prediction value of the Taipei MRT. Again, the training data from 1996 to 2013 were used as the artificial neural network training materials, and the data from 2014 to 2015 were used as the sample data. In total, 19 parameters were used as the input variables of the artificial neural network, while the output variable was the passenger traffic prediction value of the Taipei MRT. Twenty variables and 19 parameter variables were used to compare the predictive values of the regression analysis and artificial neural network after removing the independent variable, the MRT revenue. The results showed that the MAPE residual value of the regression analysis was 1.28%, and the MAPE residual value of the artificial neural network was 0.19%. We observed no significant difference between the two, which proves that there was no collinearity problem after deleting the independent variable, the MRT passenger revenue.

#### 4.3.4. Monthly Passenger Traffic Experiment

In this study, we collected data on passenger traffic (month) and related variables (month); however, because some data were recorded at different times, 167 regression analyses were conducted from January 2000 to December 2015 (excluding the incomplete data on Cyclone Nari in September 2001). The neural network used the mode of 2000–2013 training 2014–2015 and the following seven parameter variables: the MRT operation mileage, MRT station number, MRT train number, MRT extended vehicle mileage, MRT extended passenger mileage, MRT passenger revenue, and two-way transit preferential volume. The output variable is the passenger traffic prediction value of the Taipei MRT. The results showed that the MAPE value of the regression analysis was 1.45%, and the MAPE value of the artificial neural network was 0.42%. Both methods showed the best prediction results.

#### 4.4. Summary of the Empirical Results

The results of this study are shown in Table 12, which summarizes the empirical results of the regression analysis and artificial neural network.



**Table 12.** Summary of the empirical results of the regression analysis and artificial neural network.

Experiment Module	Regression Analysis MAPE Value	Artificial Neural Network MAPE Value	Results
20 variables	6.47%	4.82%	The artificial neural network is better.
The four variables excluded in the regression analysis were deleted: the GNP, GDP, MRT mileage, and MRT extended passenger mileage.	0.94%	0.54%	Both had the best prediction results.
The four variables with the highest MAPE values were deleted: the economic growth rate, personal income, total amount of goods imported into Taiwan, and number of MRT stations.	0.62%	0.31%	Both had the best prediction results.
19-variable collinearity verification (deletion of the MRT passenger revenue).	5.19%	4.63%	There is no collinearity problem.
Monthly passenger traffic experiment.	1.45%	0.42%	Both had the best prediction results.

## 5. Analysis of Empirical Results

The results and research contributions of this paper in relation to passenger traffic prediction provide a valuable reference for both academics and practitioners. They are summarized in this section.

### 5.1. Analysis of Results

In this study, we used an artificial neural network and regression analysis to construct the Taipei MRT passenger traffic prediction model. Seven findings of this study are worth summarizing:

1. A total of 20 parameters from 1996 to 2013 were used as the input variables of the regression analysis, and the output value was the passenger traffic prediction value of the Taipei MRT. The stepwise regression of the regression analysis, predicting the MAPE value of the Taipei MRT passenger traffic, is 6.47%, which demonstrates an excellent predictive performance.
2. Using the training data from 1996 to 2015 as artificial neural network training materials and the data from 2013 to 2015 as the sample data, a total of 20 parameters were used as the input variables of the artificial neural network; the output value was the passenger traffic prediction value of the Taipei MRT. Using an artificial neural network to predict the MAPE value of the Taipei MRT passenger traffic, the result was 4.82%, which demonstrates an excellent predictive performance.
3. We used the 5-year training for 2 years (2009–2013/2014–2015) model and deleted the already eliminated four variables of regression analysis, which were the GNP, GDP, MRT mileage, and MRT extended passenger mileage. In total, 16 parameters were used as the input variables for regression analysis and artificial neural network comparison. The MAPE value predicted by regression analysis was found to be 0.94%, and the MAPE value of the artificial neural network was 0.54%. Both methods had the best prediction results.
4. We used the model of 5-year training for 2 years (2009–2013/2014–2015) and deleted the four variables with the highest MAPE values in the artificial neural network, including the economic growth rate, personal income, total amount of goods imported into Taiwan, and number of MRT stations. In total, 16 parameters were used as the input variables for regression analysis and artificial neural network comparison. The

- MAPE value predicted by regression analysis was 0.62%, and the MAPE value of artificial neural networks was 0.31%. Both methods had the best prediction results.
5. In order to verify whether the independent variable MRT passenger revenue is collinear with the dependent variable MRT passenger traffic, we applied statistical software to analyze the inputting of 19 parameter variables from 1996 to 2013 and outputting of the passenger traffic prediction value of the Taipei MRT. Again, we used the training data from 1996 to 2015 as the artificial neural network training materials and the data from 2013 to 2015 as the sample data. We used a total of 19 parameters as the input variables of the artificial neural network and the passenger traffic prediction value of the Taipei MRT as the output value in order to compare the predictive values of the regression analysis and artificial neural network. The result showed that the MAPE value of the regression analysis was 5.19%, and the MAPE value of the artificial neural network was 4.63%. Twenty variables and 19 parameter variables were used to compare the predictive values of the regression analysis and artificial neural network after removing the independent variable, the MRT revenue. The results showed that the MAPE residual value of regression analysis was 1.28%, and the MAPE residual value of the artificial neural network was 0.19%. No significant difference was observed between them, which proves that there is no collinearity problem after deleting the independent variable, the MRT passenger revenue.
  6. We collected data on the passenger traffic (month) and related variables (month); however, because some data were recorded at different times, 167 regression analyses were conducted from January 2000 to December 2015 (excluding the incomplete data on Cyclone Nari in September 2001). The neural network used the mode of 2000–2013 training 2014–2015 and the following seven parameter variables: the MRT operation mileage, MRT station number, MRT train number, MRT extended vehicle mileage, MRT extended passenger mileage, MRT passenger revenue, and two-way transit preferential volume. The output value is the passenger traffic prediction value of the Taipei MRT. The results show that the MAPE value of regression analysis was 1.45%, and the MAPE value of the artificial neural network was 0.42%. Both methods had the best prediction results.
  7. This study investigated the passenger traffic prediction of the Taipei MRT and analyzed and constructed the prediction model based on two prediction methods: regression analysis and artificial neural networks. The results are presented to allow transportation organizations to maximize their profits by making plans and decisions relating to future operations based on past passenger traffic trends.

## 5.2. Research Contributions

To improve the utilization rate of the metro area's transportation system and reduce environmental pollution by combining it with the public transportation system of the MRT station, the following research contributions are provided:

1. Choosing the right forecasting tools for business planning and decision making.  
The amount of passenger traffic significantly impacts the transportation industry, and accurate prediction depends on a full understanding of the metropolitan environment and analysis of the factors that influence demand. The commonly used methods for forecasting passenger traffic demand include statistical regression analysis, econometric modeling, time-series modeling, neural-network modeling, gray theory, etc.; selecting the appropriate forecasting tools can allow for the maximization of profits by aiding in planning and decision making relating to future operations based on past passenger traffic trends.
2. Construction and planning of transport systems.  
Passenger traffic prediction is a very important factor in the construction and planning of transportation systems and serves as the essential basis for putting forward requirements for the construction and expansion of transportation equipment. Therefore, the results of

this study can be used as a reference for transportation organizations, such as operation management, manpower allocation, shift distance, transportation demand, etc.

### 3. Curbing the growth of private vehicles.

The excessively increased use of private vehicles will cause such problems as air pollution, noise, road congestion, traffic accidents, etc., and the social costs shall be shared by the public. Therefore, strengthening the control on the increase of private vehicles is necessary, and passenger traffic prediction cooperates perfectly with public transportation transfer system construction, encouraging people to transfer from private vehicles to public transport systems and then improving ridership on public transportation systems.

### 4. Public transportation leads to urban development.

With the public transport system as the backbone of urban development, passenger traffic prediction can establish a different planning method and procedure from traditional urban development, implement the priority concept of public transportation, look toward the future, and actively promote subsequent MRT-related construction based on the existing construction. In addition to being safe and convenient, public transportation systems can also inhibit the growth of private transportation, reduce carbon while promoting energy saving, and build low-carbon and green energy cities.

### 5. Improvement of the efficiency of all public transport vehicles.

In this study, the data on MRT passenger traffic prediction not only provide traffic dispatching and transport strategic planning data but also combine them with the electronic ticket or passenger traffic data of public transport vehicles, such as the high-speed railways, Taiwan's railways, urban buses, and Uikes around the MRT stations, to conduct the overall passenger traffic analysis. This is used as the basis for the planning and operation management of the public transport system to improve the efficiency and service quality of each public transport vehicle and thus maximize the benefit of passenger traffic prediction.

## 6. Conclusions

This study investigated the passenger traffic prediction of the Taipei MRT and analyzed and constructed the prediction model based on two prediction methods: regression analysis and artificial neural networks. The results are presented as a reference for transportation organizations to allow them to maximize profits by making plans and decisions relating to future operations based on past passenger traffic trends.

During this research, we discussed with senior executives of MRT companies and found that accurate passenger traffic prediction can effectively reduce costs. The management implications of this study are as follows:

1. The forecasting of passenger traffic has a great influence on the transportation industry, and transportation organizations can make plans and decisions for future operations based on past passenger traffic trends. Therefore, it is very important for transportation organizations to have a clear and objective forecasting method.
2. Based on past research data and a literature review, we determined the variable that affects the forecasting of passenger traffic demand the most. The forecasting of the MRT passenger traffic affects the future economic development of Taipei City and New Taipei City, so it is necessary to select the appropriate passenger traffic forecasting tool.
3. The data on passenger traffic forecasting in this study not only allow for the planning of traffic dispatching, operation management, manpower allocation, shift distance, transportation demand, etc., but also serve as an essential basis for transportation organizations, enabling them to put forward requirements for the construction and expansion of transportation equipment.

Despite the limitations of this study, we believe that the findings and management implications of our study are intriguing enough to invite future research on the topic of passenger traffic forecasting, as well as future research on other traffic-related topics.

**Author Contributions:** Conceptualization, Y.-S.C. and S.-H.C.; methodology, Y.-S.C.; software, S.-H.C.; validation, S.-H.C.; writing—original draft preparation, S.-H.C.; writing—review and editing, Y.-S.C., S.-F.C. and C.-K.L. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was partially supported by the Ministry of Science and Technology of Taiwan, grant number MOST 109-2221-E-146-003.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

- Verbavatz, V.; Barthelemy, M. Access to mass rapid transit in OECD urban areas. *Sci. Data* **2020**, *7*, 1–6. [CrossRef]
- Wiseman, Y. Autonomous vehicles. In *Encyclopedia of Information Science and Technology*, 5th ed.; IGI Global: Hershey, PA, USA, 2021; pp. 1–11.
- Wiseman, Y. Driverless cars will make passenger rail obsolete [Opinion]. *IEEE Technol. Soc. Mag.* **2019**, *38*, 22–27. [CrossRef]
- Yang, Z.; Li, C.; Jiao, J.; Liu, W.; Zhang, F. On the joint impact of high-speed rail and megalopolis policy on regional economic growth in China. *Transp. Policy* **2020**, *99*, 20–30. [CrossRef]
- Lio, W.; Liu, B. Uncertain maximum likelihood estimation with application to uncertain regression analysis. *Soft Comput.* **2020**, *24*, 9351–9360. [CrossRef]
- Toraman, S.; Alakus, T.B.; Turkoglu, I. Convolutional capsnet: A novel artificial neural network approach to detect COVID-19 disease from X-ray images using capsule networks. *Chaos Solitons Fractals* **2020**, *140*, 110122. [CrossRef]
- Mazanec, J.A. Classifying tourists into market segments: A neural network approach. *J. Travel Tour. Mark.* **1992**, *1*, 39–60. [CrossRef]
- Law, R.; Au, N. A neural network model to forecast Japanese demand for travel to Hong Kong. *Tour. Manag.* **1999**, *20*, 89–97. [CrossRef]
- Kulendran, N.; Witt, S.F. Leading indicator tourism forecasts. *Tour. Manag.* **2003**, *24*, 503–510. [CrossRef]
- Grosche, T.; Rothlauf, F.; Heinzl, A. Gravity models for airline passenger volume estimation. *J. Air Transp. Manag.* **2007**, *13*, 175–183. [CrossRef]
- Shams, S.R.; Jahani, A.; Moeinaddini, M.; Khorasani, N. Air carbon monoxide forecasting using an artificial neural network in comparison with multiple regression. *Model. Earth Syst. Environ.* **2020**, *6*, 1467–1475. [CrossRef]
- Liu, S.; Mocanu, D.C.; Matavalam, A.R.R.; Pei, Y.; Pechenizkiy, M. Sparse evolutionary deep learning with over one million artificial neurons on commodity hardware. *Neural Comput. Appl.* **2020**, 1–16. [CrossRef]
- Sun, W.; Huang, C. A carbon price prediction model based on secondary decomposition algorithm and optimized back propagation neural network. *J. Clean. Prod.* **2020**, *243*, 118671. [CrossRef]
- Liu, R.; Mancuso, C.A.; Yannakopoulos, A.; Johnson, K.A.; Krishnan, A. Supervised learning is an accurate method for network-based gene classification. *Bioinformatics* **2020**, *36*, 3457–3465. [CrossRef] [PubMed]
- Shirokanev, A.S.; Kirsh, D.V.; Kupriyanov, A.V. Research of an algorithm for crystal lattice parameter identification based on the gradient steepest descent method. *Comput. Opt.* **2017**, *41*, 453–460. [CrossRef]
- Hemilä, H.; Chalker, E. Vitamin C may reduce the duration of mechanical ventilation in critically ill patients: A meta-regression analysis. *J. Intensive Care* **2020**, *8*, 15. [CrossRef]
- Şahin, U.; Şahin, T. Forecasting the cumulative number of confirmed cases of COVID-19 in Italy, UK and USA using fractional nonlinear grey Bernoulli model. *Chaos Solitons Fractals* **2020**, *138*, 109948. [CrossRef]
- Prado, F.; Minutolo, M.C.; Kristjanpoller, W. Forecasting based on an ensemble autoregressive moving average-adaptive neuro-fuzzy inference system–neural network-genetic algorithm framework. *Energy* **2020**, *197*, 117159. [CrossRef]
- Longo, G.A.; Righetti, G.; Zilio, C.; Ortombina, L.; Zigliotto, M.; Brown, J.S. Application of an Artificial Neural Network (ANN) for predicting low-GWP refrigerant condensation heat transfer inside herringbone-type Brazed Plate Heat Exchangers (BPHE). *Int. J. Heat Mass Transf.* **2020**, *156*, 119824. [CrossRef]
- Feng, Y.; Yang, T.; Niu, Y. Subpixel computer vision detection based on wavelet transform. *IEEE Access* **2020**, *8*, 88273–88281. [CrossRef]
- Sohn, C.; Choi, H.; Kim, K.; Park, J.; Noh, J. Line Chart understanding with convolutional neural network. *Electronics* **2021**, *10*, 749. [CrossRef]
- Shibzukhov, Z.M. Robust method for finding the center and the scatter matrix of the cluster. *J. Phys. Conf. Ser.* **2020**, *1479*, 012045. [CrossRef]
- Peng, Q.; Wen, F.; Gong, X. Time-dependent intrinsic correlation analysis of crude oil and the US dollar based on CEEMDAN. *Int. J. Financ. Econ.* **2021**, *26*, 834–848. [CrossRef]
- Shi, J.; Liu, B.; Qin, J.; Jiang, J.; Wu, X.; Tan, J. Experimental study of performance of repair mortar: Evaluation of in-situ tests and correlation analysis. *J. Build. Eng.* **2020**, *31*, 101325. [CrossRef]
- Ibrahim, M.S.; Sidiropoulos, N.D. Reliable detection of unknown cell-edge users via canonical correlation analysis. *IEEE Trans. Wirel. Commun.* **2020**, *19*, 4170–4182. [CrossRef]
- Di Tella, M.; Romeo, A.; Benfante, A.; Castelli, L. Mental health of healthcare workers during the COVID-19 pandemic in Italy. *J. Eval. Clin. Pract.* **2020**, *26*, 1583–1587. [CrossRef]

27. Moskowitz, S.; Dewaele, J.M. Is teacher happiness contagious? A study of the link between perceptions of language teacher happiness and student attitudes. *Innov. Lang. Learn. Teach.* **2021**, *15*, 117–130. [CrossRef]
28. Liu, B.; Zhao, Q.; Jin, Y.; Shen, J.; Li, C. Application of combined model of stepwise regression analysis and artificial neural network in data calibration of miniature air quality detector. *Sci. Rep.* **2021**, *11*, 1–12.
29. Hornby, T.G.; Henderson, C.E.; Holleran, C.L.; Lovell, L.; Roth, E.J.; Jang, J.H. Stepwise regression and latent profile analyses of locomotor outcomes poststroke. *Stroke* **2020**, *51*, 3074–3082. [CrossRef]
30. Żogała-Siudem, B.; Jaroszewicz, S. Fast stepwise regression based on multidimensional indexes. *Inf. Sci.* **2021**, *549*, 288–309. [CrossRef]



MDPI  
St. Alban-Anlage 66  
4052 Basel  
Switzerland  
Tel. +41 61 683 77 34  
Fax +41 61 302 89 18  
[www.mdpi.com](http://www.mdpi.com)

*Electronics* Editorial Office  
E-mail: [electronics@mdpi.com](mailto:electronics@mdpi.com)  
[www.mdpi.com/journal/electronics](http://www.mdpi.com/journal/electronics)







MDPI  
St. Alban-Anlage 66  
4052 Basel  
Switzerland

Tel: +41 61 683 77 34  
Fax: +41 61 302 89 18

[www.mdpi.com](http://www.mdpi.com)



ISBN 978-3-0365-3979-9