# Short-Term Prediction of the Traffic Status in Urban Places Using Neural Network Models

Georgia Aifadopoulou[1], Charalampos Bratsas[2], Kleanthis Koupidis[2], Aikaterini Chatzopoulou[2], Josep-Maria Salanova[1(✉)], and Panagiotis Tzenos[1]

[1] Centre for Research and Technology Hellas – Hellenic Institute of Transport, 57001 Thessaloniki, Greece
`jose@certh.gr`
[2] Open Knowledge Greece, Semelis 1, 54352 Thessaloniki, Greece

**Abstract.** The last decades the phenomenon of urbanisation has led to crowded and jammed areas, which makes life in cities more stressful. Thus, there is a high interest in the field of Intelligent Transportation Systems in order to prevent the traffic congestion. The most common way to prevent this phenomenon is with the use of short-term forecasting of traffic parameters, such as traffic flow and speed. Nowadays, the accuracy of the estimations has increased significantly due to the use of the latest technological advances, such as probe data in combination with machine learning techniques. Probe data is a type of crowd-sourced data collected from individuals, including vehicles, passengers, travellers or pedestrians. This paper focuses on the data processing component with the use of neural networks, for predicting traffic status in urban areas based on the relation between traffic flows and speed. As a case study is used the traffic status in the city of Thessaloniki, Greece. In this case, data is aggregated after the collection phase, which gives a better representation of the mobility patterns in the city. Two types of test were performed. The first one shows the results of the prediction of eight sequentially quarters of the time, while the second test provides the prediction four steps forward of the date time. The results of both tests provide accurate predictions.

**Keywords:** Neural network · Traffic prediction

## 1 Introduction

Transportation is one of the factors responsible for 26% of Green House Gas (GHG) emissions at European level. The percentage varies depending on the urban area, based on different activities, but mostly on Intelligent Transportation Systems as well as by Big and Open Data within the Smart Cities framework.

Even though, technological advances have been attributed to the increase of quality and quantity of mobility-related data. The challenge of producing the best possible end-products out of these big datasets is still twofold; first there is a need for developing algorithms able to fuse, filter, validate and process big amounts of data (almost) at real-time; and secondly, there is a constant need for developing new applications and

services for providing innovative and advanced traveller information services, traffic management schemes and environmental indicators based on these data and processing capabilities.

This paper focuses on the data processing component by presenting a machine learning algorithm for predicting traffic status in urban areas based on statistical measures of traffic flows and speed. The model is applied to the city of Thessaloniki, Greece. The paper is structured as follows. A review of key contributions in the domain of predicting traffic status is provided in Sect. 2. Sections 3 and 4 deal with the methodological approach, and its application to the traffic status in Thessaloniki, which is monitored by collecting Floating Car Data from a professional fleet. Finally, conclusions are presented in Sect. 5.

## 2   Literature Review

The last decades cities have become more crowded and jammed, which increased the need for accurate traffic and mobility management through the development of solutions based on Intelligent Transport Systems. Therefore, the interest in the short-term forecasting of traffic parameters, such as traffic flow and speed, has been increased.

The accuracy of the estimations and predictions have risen significantly by using more granulate data sources, such as probe data. In this case, data is aggregated after the collection phase, which significantly increases the quality of the collected data and multiplies the capabilities for processing this data and having better representation of the mobility patterns in a city. The main probe data sources are based on detections or Bluetooth-enabled devices [1], mobile cell phones [2] or vehicles telemetric, such as Floating Car Data (FCD) [3]. An overview of data collection technologies is provided in [4]. These sources can be used for measuring traffic characteristics, such as speeds but it may not capture traffic flow correctly [4].

The main issue for the prediction of the traffic flow in road networks comes to the development of an algorithm that will combine computational speed and accuracy for both short and long-term problems. Many ways have been introduced to perform short-term predictions, such as Regression models [5], nearest neighbour [6], ARIMA, discretization modelling approach, as an easier solution to the complicated nonlinear models [7], and neural networks, which are considered to be the best alternative [8, 9].

## 3   Methodological Approach

### 3.1   Overview

The main scope is the development of an algorithm that will predict the traffic status in the city of Thessaloniki. This approach integrates machine learning techniques using the travel times, traffic counts and speeds as well as the skewness, kurtosis and standard deviation of the speed, to train an appropriate NN Model for efficient and robust traffic speed prediction. The considerable amount and the nature of the data and the advantage of multiple learning algorithms led us to use Artificial Neural Networks (ANN).

We want to detect all possible interactions and complex nonlinear or linear relationships, to provide better speed predictions.

For this work, we have created a package (TrafficBDE[1]) in R Software, available on Github. The user selects the road and the date time to predict the wanted variable, either the mean speed or the entries, for this road at this date time and also how many steps forward they want the predicted value to cover.

### 3.2 Dataset and Experimental Set up

Two main datasets are used, a dataset composed of floating car data composed of pulses generated by vehicles and a dataset composed of the road network segments.
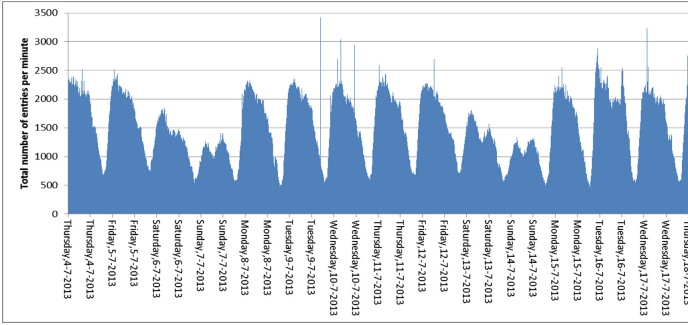
The road network segments (Fig. 1) were extracted from the urban mobility model for Thessaloniki [10], which is composed of 47.807 intersections and 137.854 directed links, both elements bearing geometric (length, location in the network) and traffic related characteristics (number of lanes, free flow speed, capacity, direction, allowed transportation modes, existence of dedicated lanes, parking prohibition).



**Fig. 1.** Overview of Thessaloniki's modeled road network [10].

The floating car data in Thessaloniki is obtained from a fleet of 1200 taxis, which represent the 50% of the taxis in the city and collect location (lat, long), orientation, status (empty or occupied) and speed every 100 meters. The total amount of data varies between 500 and 2.500 datasets per minute. The data is directly collected by the Taxi association and provided for this particular research. The temporal distribution of the data presents a peak early in the morning which is reduced slowly during the day until the afternoon. It also presents a significant reduction during the weekends, as it can be observed in Fig. 2. Saturday still present a higher peak in the morning and a stable period in the afternoon, while two peaks are clearly observed on Sundays, one in the morning and one in the afternoon.

---

[1] TrafficBDE package imports the following packages on RStudio, caret, data.table, dplyr, graphics, grDevices, jsonlite, lubridate, RCurl, readr, reshape, stats, zoo, and it is available for the R version 3.3.1 or later. https://github.com/okgreece/TrafficBDE.

**Fig. 2.** Distribution of taxi pulses during the period 4/7/2013–18/7/2017 [11].

The authors decided to use only internal factors aiming at building a data-driven model, which will be able to predict the impact of rain in the network speed without knowing that it is raining, only form the drop on the speed of the taxis. This will allow for predicting non-recurrent congestion such as the one generated locally in the surroundings of a stadium twice a month, when there is a game, again without being informed about the game taking place.

### 3.3 Pre-processing and Data Analysis

The road, the date time, the wanted variable to be predicted and how many steps forward the prediction will be, are being selected. In Table 1 is a short description of the inputs that must be defined, to be used the TrafficBDE package.

**Table 1.** Short description of the inputs in the algorithm.

| Input | Description |
|---|---|
| Path | The path with the historical data available |
| Link_id | The Link_id for the road to be predicted |
| Direction | The direction of the road to be predicted |
| Datetime | The time and date for the prediction |
| Steps | How many steps forward will the prediction be |
| Predict | The variable to be predicted. Either "Mean_speed" or "Entries" |

Each road is deviated by its geographical position and length so to each part has been given a unique link id. Namely, each link id represents a specific part of a road. The roads are either one-way or two-way based on that the direction value is 1 or 2.

The path mentioned above contains the data that will be used for the prediction. The input data consists of the min, max and mean speed of the roads, the date time these speeds were observed, the entries of each road at this date time, and the unique entries, i.e. how many were the different entries. Finally, there are also provided some statistical measures of the mean speed of each Link id in a particular quarter,

these measures are standard deviation, skewness and kurtosis, which are going to be used as features.

Firstly, the algorithm filters the historical data of the roads based on the selected Link id and direction, and then the algorithm keeps the data of the previous two weeks from the date time wanted and calculates the features, mentioned below, of the speed. Afterwards, the algorithm checks if all the quarters exist, this means 1344 quarters for two weeks. If there are missing quarters, they are created, and linear interpolation fills the rest data values. When the data are completed, they are split into the train set and test set, and they are processed and normalised between 0 and 1.

As mentioned above the statistical measures are going to be used as features to train a more accurate model. Namely, these features refer to standard deviation, skewness and kurtosis of each Link id in a particular quarter. The features are, also, processed and normalised between 0 and 1 (Fig. 3).

|  | Min_speed | Max_speed | Stdev_speed | Skewness_speed | Kurtosis_speed | Entries | UniqueEntries | Mean_speed |
|---|---|---|---|---|---|---|---|---|
| 2017-01-16 22:15:00 | 0.04838710 | 0.4915254 | 0.34838094 | 0.5029376 | 0.4147874 | 0.154411765 | 0.42857143 | 0.34693878 |
| 2017-01-16 22:30:00 | 0.32258065 | 0.3898305 | 0.26030382 | 0.5844592 | 0.3162458 | 0.088235294 | 0.24489796 | 0.32653061 |
| 2017-01-16 22:45:00 | 0.11290323 | 0.2542373 | 0.23749948 | 0.2927972 | 0.6167640 | 0.091911765 | 0.28571429 | 0.32653061 |

**Fig. 3.** Sample of the data after being normalized.

## 3.4 Multilayer Perceptron Model

After the preprocessing, the data are divided into the train and the test set. In the created algorithm, the input layer consists of 7 nodes-features; min and max speed of the road skewness, kurtosis and standard deviation of the speed and entries and unique entries. The algorithm used to train the NN is Resilient back propagation (Rprop).

Rprop focuses on eliminated the influence of the size of the partial derivative on the weight step. Therefore, for the indication of the direction of the weight update, only the sigh of the derivative is considered. The size of the weight change is determined by the update-value.

The model used is Multilayer perceptron (MLP). MLP is a feedforward ANN model that is known for the simplicity and the performance of nonlinear patterns. Due to those characteristics, it is used in similar applications. MLP represents a directed graph of multi-layers of nodes, there are three types of layers the input, hidden and an output; each layer is connected to the next one.

In general, there is no restriction on the number of nodes and hidden layers to be used. We are testing in each step which combination provides the best results (cv result), and that is chosen as our final model.

## 3.5 Cross-Validation

Different combinations of the number of neurons in each hidden layer are checked with 10-fold cross-validation, and the model with the minimum error is used as the train model. The train set is separated into 10 datasets, which are consists of 134, original sample is divided into 10 equal subsamples (1344/10), observations each. Nine of them

will be used to train the NN and will predict the other observations. This process will be repeated 10 times and according to the smallest error the NN, which provides the best predictions, will be chosen.

### 3.6    Neural Network Output

The output of the algorithm is the predicted value of the wanted variable, either the Mean speed or the Entries. For example, the structure of the model presented below (Fig. 4) is characterised by the 7 inputs, 4 neurons in each hidden layer and one output, the mean speed.
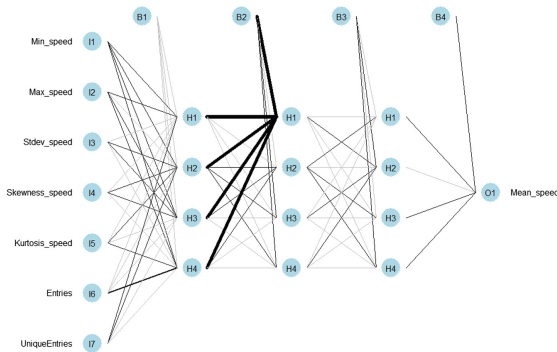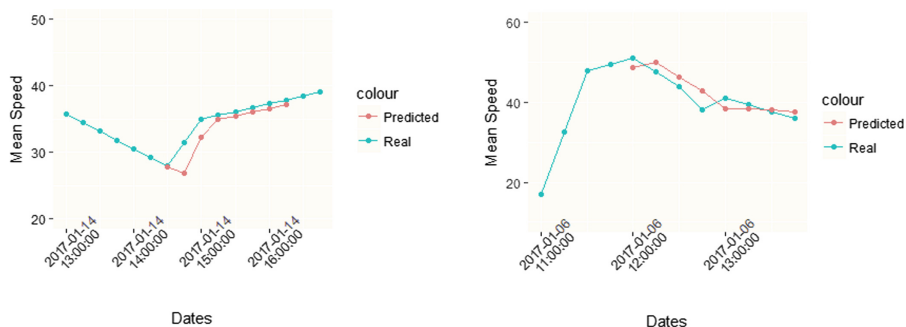


**Fig. 4.**  Example of NN model.

## 4    Results and Discussion

The above methodology is applied to the FCD collected in Thessaloniki [10]. Thessaloniki is the second largest city in Greece, with a total of more than 1 million citizens in its greater area, covering a total of 1,500 km$^2$ with an average density of 665 inhabitants per km$^2$. The total number of vehicles in the city exceeds 750,000, including private cars, heavy vehicles and motorcycles.

The available data are the historical data of two random roads with links for January in 2017. The features mentioned were calculated and algorithm chose the NN model based on 10-fold cross-validation.

We perform two types of test. The first one shows the results of the prediction of eight sequentially date times. As the second type provides the prediction 4 steps forward of the date time, i.e. in the first step, the algorithm uses the historical data to predict the speed value of the first quarter. The predicted value will be used in the second step to predict the speed value of the next quarter. The same process continues until we reach the last quarter.

The following figures are the results of the first test. In both figures, the predicted values follow the pattern. It is worth noting that in cases where the pattern changed abruptly, giving a slightly increased Root Mean Square (RMSE) but still not higher than 10.5 km/h and less than 6.5 km/h in the two selected links, the algorithm recognised the changed pattern and followed it in the next predictions (Fig. 5).

**Fig. 5.** Plot with the Real (blue) and the Predicted (red) values.

The following Table 2 shows the results of the 4-step forward case. Taking a look at the results, we can see the prediction value does not significantly differ from the real, and the algorithm captures the changes in mean speed.

**Table 2.** Results of the prediction in the two selected links.

| Link | Date time | Predicted speed | Real speed | RMSE |
|---|---|---|---|---|
| 1 | 2017-01-12 19:30:00 | 17.07 | 16.71 | 0.35 |
| 1 | 2017-01-12 19:45:00 | 16.88 | 16.14 | 0.74 |
| 1 | 2017-01-12 20:00:00 | 16.02 | 15.57 | 0.45 |
| 1 | 2017-01-12 20:15:00 | 15.69 | 15.00 | 0.69 |
| 2 | 2017-01-14 16:00:00 | 36.75 | 37.28 | 0.53 |
| 2 | 2017-01-14 16:15:00 | 37.26 | 37.85 | 0.59 |
| 2 | 2017-01-14 16:30:00 | 37.77 | 38.42 | 0.65 |
| 2 | 2017-01-14 16:45:00 | 38.31 | 39.00 | 0.68 |

## 5 Conclusions

A machine learning algorithm for predicting traffic status (speed) has been presented and applied to the city of Thessaloniki, achieving an accuracy of a few km/h. The raw data is collected by a fleet or 1.200 vehicles with a frequency of 6–10 s, generating large data sets and covering both spatially and temporally the whole city with high granularity. This allows, on the one hand having better accuracy but on the other hand increases significantly the data filtering and processing needs, especially when the predictions are made in real time to fuel mobility services.

The predictions are based on spatial relations of traffic flow in addition to the time series generated for each link, which enriches the dataset significantly due to the propagation properties of traffic flow. The increase in the reliability of the predictions will allow better traffic management in Thessaloniki as well as enhanced information to drivers.

Future research of the authors will deal with including other datasets in the analysis, such as travel time along the main routes in Thessaloniki obtained from the network of Bluetooth detectors. The aforementioned will allow combining travel times along various set of links with instantaneous speed in some of that links, and may end up with more accurate predictions. Finally, the addition of external factors, such as weather, will be evaluated by the authors.

# References

1. Mitsakis, E., Salanova, J.M., Chrysohoou, E., Aifadopoulou, G.: A robust method for real time estimation of travel times for dense urban road networks using point-to-point detectors. Transport **30**(3), 264–272 (2015). Special Issue on Smart and Sustainable Transport

2. Herrera, J.C., Work, D.B., Herring, R., Ban, X., Jacobson, Q., Bayen, A.M.: Evaluation of traffic data obtained via GPS-enabled mobile phones: the Mobile Century field experiment. Transp. Res. Part C **18**, 568–583 (2010)

3. Salanova, J.M., Maciejewski, M., Bischoff, J., Estrada, M., Tzenos, P., Stamos, I.: Use of probe data generated by taxis. In: Schintler, L.A., Chen, Z. (eds.) Big Data for Regional Science. Routledge Advances in Regional Economics, Science and Policy. Taylor & Francis Group, Abingdon (2017)

4. Antoniou, C., Balakrishna, R., Koutsopoulos, H.N.: A Synthesis of emerging data collection technologies and their impact on traffic management applications. Eur. Transp. Res. Rev. **3**, 139–148 (2011)

5. Liebig, T., Piatkowski, N., Bockermann, C., Morik, K.: Dynamic route planning with real-time traffic predictions. Inf. Syst. **64**, 258–265 (2017)

6. Smith, B.L., Billy, M.W., Oswald, R.K.: Comparison of parametric and nonparametric models for traffic flow forecasting. Transp. Res. Part C Emerg. Technol. **10**(4), 303–321 (2002)

7. Li, J.-Q.: Discretization modeling, integer programming formulations and dynamic programming algorithms for robust traffic signal timing. Transp. Res. Part C Emerg. Technol. **19**(4), 708–719 (2011)

8. Vlahogianni, E.I., Karlaftis, M.G., Golias, J.C.: Optimized and meta-optimized neural networks for short-term traffic flow prediction: a genetic approach. Transp. Res. Part C: Emerg. Technol. **13**(3), 211–234 (2005)

9. Ishak, S., Ciprian, A.: Optimizing traffic prediction performance of neural networks under various topological, input, and traffic condition settings. J. Transp. Eng. **130**(4), 452–465 (2004)

10. Mitsakis, E., Stamos, I., Salanova Grau, J.M., Chrysohoou, E., Aifadopoulou, G.: Urban mobility indicators for Thessaloniki. J. Traffic Logist. Eng. (JTLE) **1**(2), 148–152 (2013). ISSN: 2301-3680

11. Salanova, J.M., Toumbalidis, J., Chaniotakis, E., Karanikolaos, N., Aifadopoulou, G.: Correlation between digital and physical world, case study in Thessaloniki. J. Locat. Based Serv. **13**, 1–15 (2018)