

1 **A Classification Approach Using Machine Learning for Predicting Traffic Flows in Areas**
2 **with Missing Sensors**

3
4 **Martina Fazio**

5 Department of Physics and Astronomy
6 University of Catania, Catania, Italy, 95125
7 Email: martina.fazio@phd.unict.it

8
9 **Piergiorgio Vitello**

10 Department of Engineering
11 University of Luxembourg, Esch-sur-Alzette, Luxembourg, 4365
12 Email: piergiorgio.vitello@uni.lu

13
14 **Juan Pineda-Jaramillo**

15 Department of Engineering
16 University of Luxembourg, Esch-sur-Alzette, Luxembourg, 4365
17 Email: juan.pineda@uni.lu

18
19 **Richard D. Connors**

20 Department of Engineering
21 University of Luxembourg, Esch-sur-Alzette, Luxembourg, 4365
22 Email: richard.connors@uni.lu

23
24 **Francesco Viti**

25 Department of Engineering
26 University of Luxembourg, Esch-sur-Alzette, Luxembourg, 4365
27 Email: francesco.viti@uni.lu

28
29
30 Word Count: 6,851 words + 2 tables = 7,351 words

31
32
33 *Submitted 31th of July*

34

1 **ABSTRACT**

2 The high rate of urbanization and the ever-increasing number of vehicles on roads require multiple efforts
3 in order to correctly identify and evaluate interventions to improve transport systems. This phenomenon
4 pushes cities to invest in new technologies and increasingly advanced support tools. In the era of Big Data,
5 the exploitation of innovative data can allow the resolution of various issues related to the world of transport
6 planning. However, data are not always available, e.g. in many rural areas, and even in the cities of
7 developing countries. In such cases, it is possible to resort to sophisticated models able to fill in the lack of
8 information. In this respect, this paper focuses on the prediction of traffic flows in areas where no data are
9 available. To do this, we develop a classification approach using machine learning (ML) models involving
10 available, but limited, datasets. After the prediction of traffic flows, a validation and similarity analysis
11 have been performed. From the latter it is possible to identify sensor devices and regions that best represent
12 those areas with no data. Moreover, a SHAP analysis has been carried out to identify the impact of features
13 associated to the traffic volume. The main finding of this study is the developing of a support tool able to
14 aid government agencies by providing an exploratory overview of what should be the traffic volume in
15 areas with poor data. We believe that the proposed methodology has the potential of becoming a first-step
16 analysis in the demand forecasting process.

17
18 **Keywords:** Traffic flow prediction, machine learning, data extraction, data analysis

1 INTRODUCTION

2 Transport systems have long been the basis of the economic development of cities. Inefficient management
3 of transport networks reflects negative consequences on the entire community. The rapid increase in traffic
4 volumes generates concrete repercussions in terms of pollution, accidents, and road congestion, constituting
5 a serious risk for urban sustainability and livability. For this reason, the prediction of traffic flows plays a
6 fundamental role for providing an idea to what extent this phenomenon negatively affects system
7 performance (1). Due to the complexity of transport networks, it is crucial to resort to advanced models and
8 tools able to provide support to the already consolidated theoretical models. In this respect, it is worthy to
9 underline the important role of ML models, especially for the prediction of transport systems features (2).
10 Furthermore, the reliability of these models for predicting and simulating the traffic states is strictly
11 depending on the type of data involved in the modeling process. In recent years, new challenges have arisen
12 on how to analyze and address these issues with the support of new technologies. The digitalization can
13 help in providing a large amount of data to monitor transport networks in order to find new solutions and
14 bringing greater efficiency to the transport system. These data can be collected by different sources such as
15 GPS devices, sensors, social networks, smartphones, quickly providing huge datasets constituting of
16 different types of data (structured, unstructured, semi-structured). Precisely with the adjectives velocity,
17 volume, and variety (3Vs model) a first definition of Big Data was given (3).

18 The use of Big Data has totally transformed the world of transport thanks to the spread of Internet
19 of Things (IoT) devices, giving greater importance to Intelligent Transport Systems (ITS) (4). In the
20 transport field, the exploitation of Big Data is now widespread: OD matrix estimation (5, 6), mobility
21 pattern analysis (7, 8), travel mode detection (9, 10), road safety analysis (11, 12). In this respect, Big Data
22 can provide an important aid for both governments and private companies to do in-depth analyses,
23 monitoring, make decisions, and improve transport network' performance (13).

24 The complexity of Big Data requires many efforts to manage the datasets with the aim of obtaining
25 information that can suggest strategic decisions. For this reason, the world of Big Data analysis is now
26 exploiting ML models, which take advantage of this data to identify patterns and predict future events (2).
27 The basic concept of ML models is the capability of learning from data without prior knowledge on
28 disparate application fields. From a transportation point of view, ML models have been involved for route
29 optimization analysis (14–16), travel behavior (17–19), traffic flows prediction (13, 20–22). Regarding this
30 latter, it is worthy to underline the importance of conducting this type of analysis in a more automatic and
31 reliable way, without resorting to manual counting in the fields (23). The purpose of this paper is framed in
32 this context, i.e. the prediction of traffic flows through a classification approach using ML models with free
33 available datasets. The case study is the Grand-Duchy of Luxembourg, in which sensors devices detecting
34 traffic flows are installed and the relative data are freely accessible by the “Portail des Travaux Publics” of
35 Luxembourg government website (24). Sensors are localized almost homogeneously around Luxembourg
36 except for Luxembourg city. For this reason, we decided to develop a methodology capable of predicting
37 traffic flows also in areas without traffic sensors using the information available from another dataset
38 retrieved through TomTom move platform (sample dataset) (25) which provides information also for
39 Luxembourg city but in a limited time period (i.e. January 2020).

40 The main challenge of this study is the prediction of flows in an area without counting devices.
41 Research studies that have dealt with a similar topic usually focused attention on the forecasting of flows
42 in short/long-term (20–22, 26–28). To the best of our knowledge, there are no studies that addressed how
43 to predict traffic flows in absence of other data sources for the study area under consideration. The
44 methodology is capable of providing valuable information for government agencies about the most
45 representative areas and counters useful to obtain a first-step exploratory overview of flows in poor data
46 regions, without resorting to count fields or data purchasing.

47 BACKGROUND

48 Research on the prediction of traffic flows is rooted in several studies starting from the '70s (29–32). One
49 of the first efforts to predict traffic flows has been carried out by Nicholson and Swann (33), who elaborated
50 a short-term prediction of traffic flows across the Mersey tunnel (Liverpool) using spectral analysis and
51

1 collecting data from inductive loop traffic detectors. The methodology provided reasonable results, but it
2 was not able to make predictions for an entire day. Lu (34) proposed an adaptive prediction system which
3 consists of a set of iterative mathematical steps to minimize the mean squared error and then maximize the
4 accuracy of the prediction of traffic flows. Despite the results of these studies are promising, the model
5 seems to have no general applicability, but it is closely linked to the case study considered. An analogous
6 study is the one conducted by Williams et al. (35). They integrated seasonal AutoRegressive Integrated
7 Moving Average (ARIMA) analysis with Winters exponential smoothing models to forecast seasonal traffic
8 flow in urban freeway. The aim was to compare the performance level of the analysis proposed compared
9 to other similar approaches (e.g. k-Nearest Neighbor, Artificial Neural Network, and historical average
10 models).

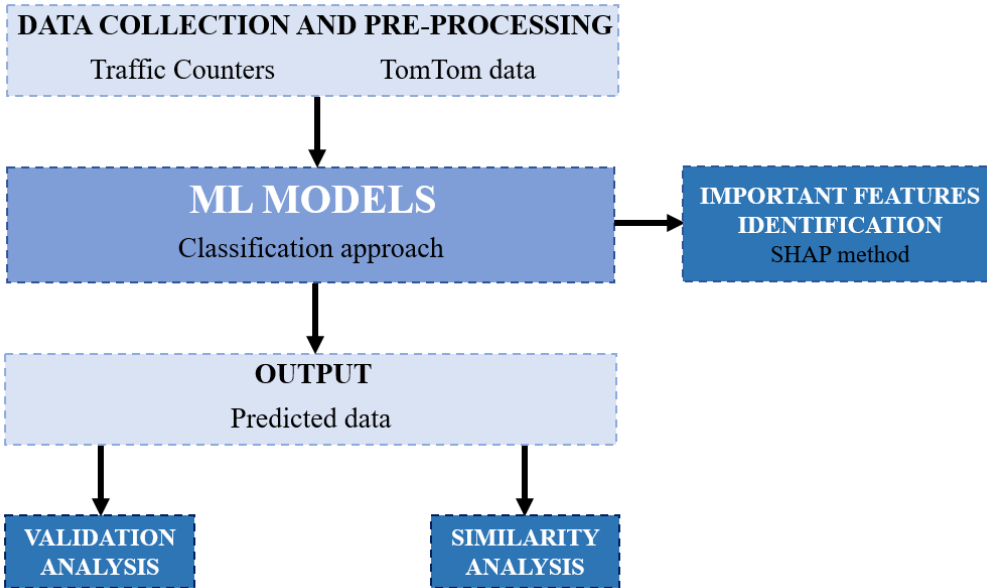
11 The accuracy of the prediction models increases with the amount of data that are engaged in the
12 analysis process. The spread of new technologies has generated many sources of mobility data (e.g. sensors,
13 smartphones) that make available a large amount of real-time information that enhances the quality of the
14 analysis carried out (36–38). Exploiting innovative data for prediction analysis is extensively addressed in
15 the literature. Moghaddam and Hellinga (39), developed a data-driven model for the prediction of travel
16 times in arterial roadways using Bluetooth detectors data, Iqbal et al. (40) elaborated a combined
17 optimization-based approach with a microscopic traffic simulation analysis for OD matrices estimation,
18 collecting data from mobile phones and traffic counts devices. Similarly, the study of Cantelmo and Viti
19 (37) also demonstrated the effectiveness of the combined use of mobile phone data and Floating Car Data
20 (FCD) for the estimation of OD matrices. The importance of exploiting innovative data is also underlined
21 by Di Donna et al. (41), who carried out a user mobility pattern analysis gathering data from Call Detail
22 Records (CDR).

23 Over the years, pushed by the technologies progresses and the advent of Big Data, research moved
24 to develop methodologies based on more sophisticated ML models. There are several studies that have dealt
25 with the prediction of traffic flows adopting a ML approach. Chikaraishi et al. (20) focused on traffic
26 prediction during disastrous events, conducting an exploratory analysis to test the applicability of a set of
27 ML models. Lv et al. (21) used a Stacked Autoencoder model (SAE) for the estimation of traffic flows of
28 roads with different levels of traffic loads. Yi and Bae (22) proposed a Deep Neural Network (DNN) model
29 in order to predict and analyze link-based traffic conditions. In this study the authors focused on
30 distinguishing the different levels of congestion using real-time traffic data. A similar study was conducted
31 by Xu et al. (42) focusing on bikesharing trips dynamic forecasting. Through the use of a deep learning
32 model, authors gained accurate predictions in terms of produced and attracted trips in different time
33 intervals. Other studies moved in this direction using different ML models: Fuzzy-neural model, Deep
34 Belief Networks, Long-Short-Term Memory network (LSTM), Convolutional Neural Network (CNN), k-
35 Nearest Neighbor, among others (26, 43–46).

36 This paper proposes a classification approach using ML models able to predict traffic flows in areas
37 where no data are available, and the only way to obtain it would require field surveys or data purchase. This
38 is an aspect that has not been investigated in the literature, and this study is aiming at bridging this gap.
39 Bauer et al. (47) conducted a similar study dealing with the estimation of O/D matrices based on flow count
40 data from sensors without resorting to a prior O/D matrices. Combining least squares with maximum
41 entropy methods authors proposed a statistical tool to calibrate path flows even in absence of prior
42 knowledge. However, the data that the authors used have been collected by devices installed inside the
43 study area.

44 **METHODOLOGY**

45 In this section we will describe the data and the methodology adopted to predict traffic flows within areas
46 in the absence of sensor devices, through a classification approach using ML models with insufficient
47 datasets. **Figure 1** shows the framework steps of the proposed methodology that will be thoroughly detailed
48 below.
49



1
2 **Figure 1** Flowchart of the proposed methodology

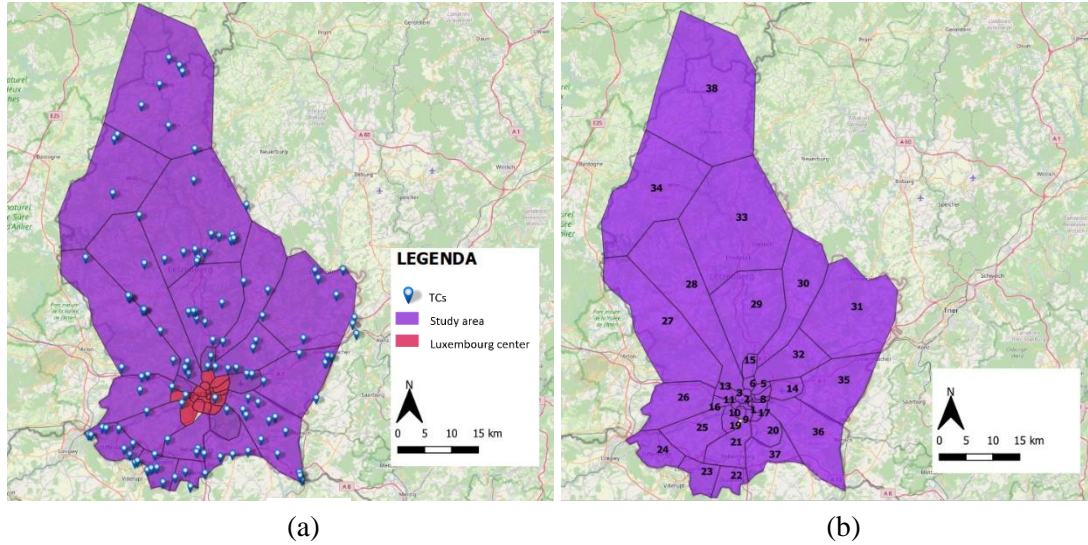
3 **Data collection and pre-processing**

4 Data has been collected from two data sources. The first collected data is provided by the “Portail
5 des Travaux Publics” of the Luxembourg government website. This open data provides car flows detected
6 by 135 traffic counters (TCs) distributed throughout the country of Luxembourg (**Figure 2a**). TCs are
7 classified according to the name and type of road in which they are localized (e.g., Highway, National,
8 Secondary Roads, and mobile counting). Each TC collects daily traffic flows per hour for the two road
9 directions, which we classified as “inflow” (i.e., toward the center) and “outflow” (i.e., away from the
10 center). Nevertheless, TCs have a different time availability depending on the year in which the sensor was
11 created, besides this data can be only accessed with a limit of the past 5 years.

12 The second collected data comprises a sample of the traffic flows provided by the TomTom Move
13 platform. TomTom collects FCD information from different devices: smartphones, car navigators, black
14 boxes. This data includes information on flows into and out from different user-defined regions, with a time
15 range of 3 hours. The sample, freely downloadable from the platform, consists of data for the month of
16 January 2020. For our analysis we divided the study area (i.e. the Grand-Duchy of Luxembourg) into 38
17 regions (**Figure 2b**). We also collected aggregate traffic flows coming from outside, into the study area.

18 As can be seen from **Figure 2a** the TCs are distributed throughout the country except for the center-
19 south (Luxembourg City), where just one TC is located. Notwithstanding, TomTom data provides
20 information also for the center regions (**Figure 2b**), the main issue with using this information is the limited
21 time period in which is possible to collect the data, whereas the datasets from TCs are available for any
22 calendar date. In this respect, the purpose of this study is to develop a ML model able to combine both
23 datasets in order to predict traffic flows in areas in which there are not TCs (i.e., Luxembourg center). **Table**
24 **1** presents the composition of the final datasets after the pre-processing analysis.

25



1
2 **Figure 2** TCs' location in the study area (a); study area zoning used for the retrieval of TomTom
3 **data (b)**

4 The two datasets were pre-processed using common processes traditionally performed in data
5 mining (e.g., creating dummy variables for categorical features, normalizing the numerical features by using
6 the z-score standardization method, among others) (48). Moreover, some additional features were created
7 using the features of the dataset (e.g., time_range (ranges of 3 hours) in the Traffic_counters dataset, in
8 order to merge both datasets) and the two datasets were merged using 4 matching features existing in both
9 datasets (i.e., region_of_origin, region_of_destination, day, time_range). This resulted in a unique dataset
10 with 101,872 rows and 19 columns. Some features have many “null values” for different reasons (e.g., no
11 available data for some TCs, no available volume data for some regions or some TCs, etc.). The target
12 feature, volume in TC, had a right-skewed distribution, with a range between 0 and 3732 vehicles, a mean
13 of 230, a standard deviation of 306 and a median of 136 vehicles. Therefore, we decided to adopt a
14 classification approach by creating the “volume_counter_ranges” feature and selecting 12 volume ranges
15 to the initial “volume_counter” feature based on the distribution of their values. Therefore, the target feature
16 (i.e., volume_counter_ranges) has the following ranges (in vehicles/hour): (a) 0, (b) 1-50, (c) 51-100, (d)
17 101-200, (e) 201-300, (f) 301-400, (g) 401-500, (h) 501-600, (i) 601-700, (j) 701-800, (k) 801-900, (l) more
18 than 900 (see **Table 1**). By setting specified input features, the model will predict the range of traffic volume
19 in the counter.

20

21 **TABLE 1** Composition of the final dataset. 101,872 observations.

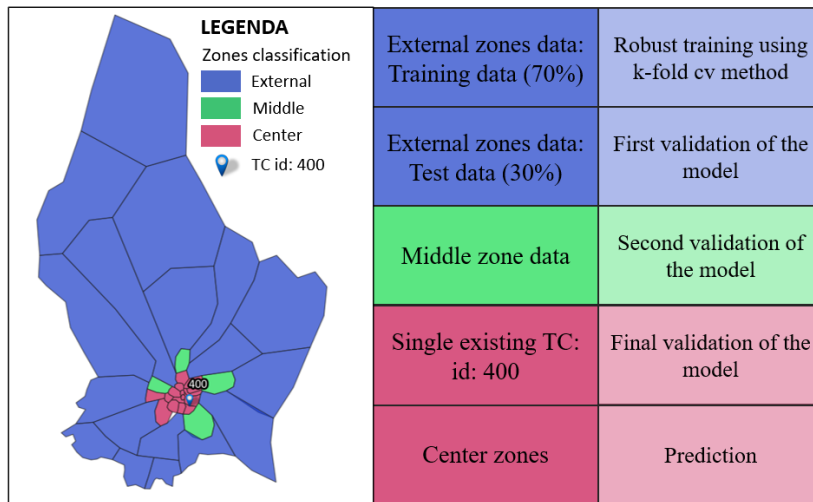
Feature	Description	type	null values
TARGET: volume_ranges in Traffic_counters	Traffic flow recorded in the counter [vehicles] in 12 ranges	categorical	71824
ID_counter	ID of the counter	categorical	52744
time	time of the registered volume in the Traffic_counters dataset	categorical	52744
direction	direction of the traffic flow(Inflow, Outflow)	categorical	52744
region_localization_counter	region where is located the traffic counter	categorical	52744
route_name	name of the road in which the traffic counter is located (e.g. A7, N2, CR184, etc..)	categorical	52744

Feature	Description	type	null values
route_type	type of the road in which the traffic counter is located (Autoroute, Route nationale, Chemin repris, Comptage mobile)	categorical	52744
region_of_origin	region of origin of the traffic flows (based on the zoning shown in Figure 2b)	categorical	0
region_of_destination	region of destination of the traffic flows (based on the zoning shown in Figure 2b)	categorical	0
day	day of the week of the registered volume	categorical	0
time_range (ranges of 3 hours)	time range (range of 3 hours) of the registered volume	categorical	0
volume in TomTom_move	volume registered in TomTom_move dataset in the range of 3 hours (time_range)	numerical	0
peak_time	weight indicating if the time is peak time or not	numerical	0
night	weight indicating if the time is night or not	numerical	0
regions_OD	Feature indicating the OD flow	categorical	0
sum_volume_counter_origin_tomtom	This feature adds all the volume for the same origin region in the same day and the same time_range	numerical	0
sum_volume_counter_destination_tomtom	This feature adds all the volume for the same destination region in the same day and the same time_range	numerical	0

1

2 **ML models**

3 The data was divided in three subsets (external, middle, center) according to the regions presented
 4 in **Figure 3**. This zoning has been performed with the aim of using the external regions data to train the ML
 5 models, and the middle regions data to validate the models. Then the model’s output will be used to perform
 6 a final validation using the single existing TC in Luxembourg city center (i.e., id: 400) and finally the best
 7 model will be applied to predict the traffic volume in the latter regions.
 8



9

10 **Figure 3 Area division for the validation analysis**

11 With the aim of predicting traffic volume ranges in Luxembourg city in absence of TCs’ data, we
 12 selected the subset of external regions to train the models. To achieve this goal, data from external regions
 13 was randomly split into training data and test data (70% and 30%, respectively) where each subdivision
 14 was composed by the traffic volume ranges and the input features in the same proportion for both
 15 subdivisions. The *SMOTE* method was used to handle imbalanced classes within the target feature (49). A
 16 set of 11 ML models were trained to predict the traffic volume ranges: Logistic Regression (LR), K

1 Neighbors Classifier (KNN), Naive Bayes (NB), Decision Tree (DT), Linear Discriminant Analysis (LDA),
 2 Quadratic Discriminant Analysis (QDA), Ada Boost (ADA), Extra Trees (ET), Random Forest (RF),
 3 Gradient Boosting (GB) and Light Gradient Boosting Machine (lightGBM). These models, whose detailed
 4 descriptions can be found in numerous studies (17, 18, 50–55), were selected since they have been widely
 5 and successfully used across a variety of classification problems.

6 The stratified k-fold cross-validation method has been implemented to validate the trained models.
 7 This method is commonly used to assess the performance of classification models performed, thanks to its
 8 capability of reducing any bias produced by the models (56). Moreover, each ML model has a set of
 9 hyperparameters that need to be tuned in order to improve its performance. This process, commonly known
 10 as “hyperparameter tuning”, is carried out by implementing the random search method, which allows
 11 assessing the values of the hyperparameter with larger impact on model performance (56, 57). To assess
 12 the performance of the models, the Area Under the Receiver Operating Characteristics (ROC) Curve,
 13 generally known as *AUC*, the Recall, Precision and F1-score evaluation metrics were employed as loss
 14 functions. Additionally, the learning curve method was implemented to identify possible *underfitting* or
 15 *overfitting* problems in the models by analyzing the model performance as more observations are used in
 16 the training process (58).

18 Validation analysis

19 After analyzing the ML models according to the evaluation metrics, the best model was chosen
 20 based on the confusion matrix, where it is possible to analyze the performance of each model in its
 21 prediction of each range. Then, once we obtained the predicted traffic volumes from the best ML model,
 22 we still need to analyze the accuracy of such results. As shown in **Figure 3** the validation phase involves
 23 the data of the middle regions in order to compare the predicted traffic volumes with the real ones (i.e., the
 24 data collected by TCs) and then the process is also repeated for the single existing TC in Luxembourg
 25 center. Therefore, we establish the goodness of our model in the test data of the external regions, in the
 26 middle regions, and in one TC in Luxembourg center.

28 Identification of important features

29 After identifying the best model based on the evaluation metrics and the confusion matrices, the
 30 impact of the features associated to the traffic volume range was obtained using the SHapley Additive
 31 exPlanation method (SHAP) (59). This method obtains the importance and the direct impact of each feature
 32 in order to better interpret the results of the model (60), and allows to identify the direct effect and the
 33 magnitude of the contributions of each input feature in the traffic volume ranges.

35 Similarity analysis

36 Similarity analysis is focused on identifying similarity trends between predicted values and the data
 37 collected by TCs. This analysis allows identification of the most representative existing TCs compared to
 38 predicted traffic flows in the center regions. The aim is to depict the inflow and the outflow traffic volumes
 39 from Luxembourg City center only using the TCs in the surrounding areas. The first step of this phase is to
 40 divide the real TCs data into two categories, depending on which type of road they are installed: highways
 41 or national roads. For both categories, we compute the similarity between the predicted TCs in center
 42 regions and the real TCs. To measure similarity between two temporal profiles of TCs, we exploit the
 43 symmetric index of Jensen-Shannon divergence (JSD) that outperforms the asymmetric Kullback-Leibler
 44 divergence (KLD) that always returns a finite value. The similarity of two TCs *i* and *j* is calculated through

45 **Equation 1:**

$$46 \quad J(D_i, D_j) = H\left(\frac{D_i + D_j}{2}\right) - \left(\frac{H(D_i) + H(D_j)}{2}\right) \quad (1)$$

47 where *J* represents the divergence of two temporal profiles, *H* is the Shannon entropy (61) and *D* is the
 48 temporal profile of a TC. The similarity can assume values in the range [0-1]. The value 0 represents the
 49
 50

1 maximum similarity (e.g., the temporal pattern of a TC with itself) and 1 represents the maximum
 2 divergence. In our dataset we have 8 TCs from highways, for each of them we compute an average value
 3 of similarity with all the TCs in center regions. The similarity for the k -th TC located in a highway (S_h) is
 4 computed as follows (**Equation 2**):

$$5 \quad S_h = \frac{\sum_{k=1}^N J(D_h, D_k)}{N} \quad (2)$$

7 where J is the JSD, and N is the number of TCs predicted in the center. From the similarity values, we were
 8 able to detect the highway TCs that have the highest similarity with the predicted traffic flows, and
 9 consequently we can identify which TCs are the most representative for the traffic flows of center regions.

10 Regarding the TCs installed on National roads, since these are the majority in our dataset (89), we
 11 choose a different approach by combining TCs for a region and hence defining the similarity at a regional
 12 level. The similarity for a region r (S_r) is defined as follow (**Equation 3**):

$$14 \quad S_r = \frac{\sum_{i=1}^{K_r} S_i}{K_r} \quad (3)$$

15 where S_r represents the similarity value between a TC and the center TCs described in **Equation 2**, and K_r
 16 is the number of TCs belonging to region r . Looking at the similarity values at the regional level we can
 17 identify which areas have traffic trends similar to center regions, this tell us where we should focus more
 18 in order to detect the general traffic flows of center regions.

21 **RESULTS AND DISCUSSION**

22 This section provides the results of the models to predict the traffic flows within a region in the absence of
 23 devices' data, through a classification approach using ML models. The process involved four main stages:
 24 (a) training a set of 11 ML models to identify the best in terms of the aforementioned evaluation metrics;
 25 (b) validation of the models using middle regions' data and the single TC in Luxembourg center; (c)
 26 identification of the direct effect and importance of input features in the predicted traffic volume ranges;
 27 (d) application of the best model to Luxembourg center in order to perform the similarity analysis.

28 **ML models results**

29 The ML models were trained with the training data and implementation of the k-fold cross-
 30 validation method (with $k=10$). Evaluation metrics were then obtained using the test data described. **Table**
 31 **2** summarizes the results of the initial models: lightGBM, RF and ET are the best models in terms of the
 32 evaluation metrics. Therefore, the hyperparameter tuning method was applied using the random search
 33 method in order to tune the performance of these models, analyzing the behavior of their learning curves to
 34 ensure they are correct (i.e., the necessary clear convergence trend between training and cross-validation
 35 scores, allowing to foresee if adding more observations to the training of these models will probably
 36 improve their performance, decreasing the risks of *overfitting* and *underfitting*). Those results are also
 37 presented in **Table 2**.

38 **TABLE 2 Results of the ML models assessed using the test data.**

Initial results				
Model	AUC	Recall	Precision	F1-score
Light Gradient Boosting (lightGBM)	0.9478	0.5933	0.7166	0.7140

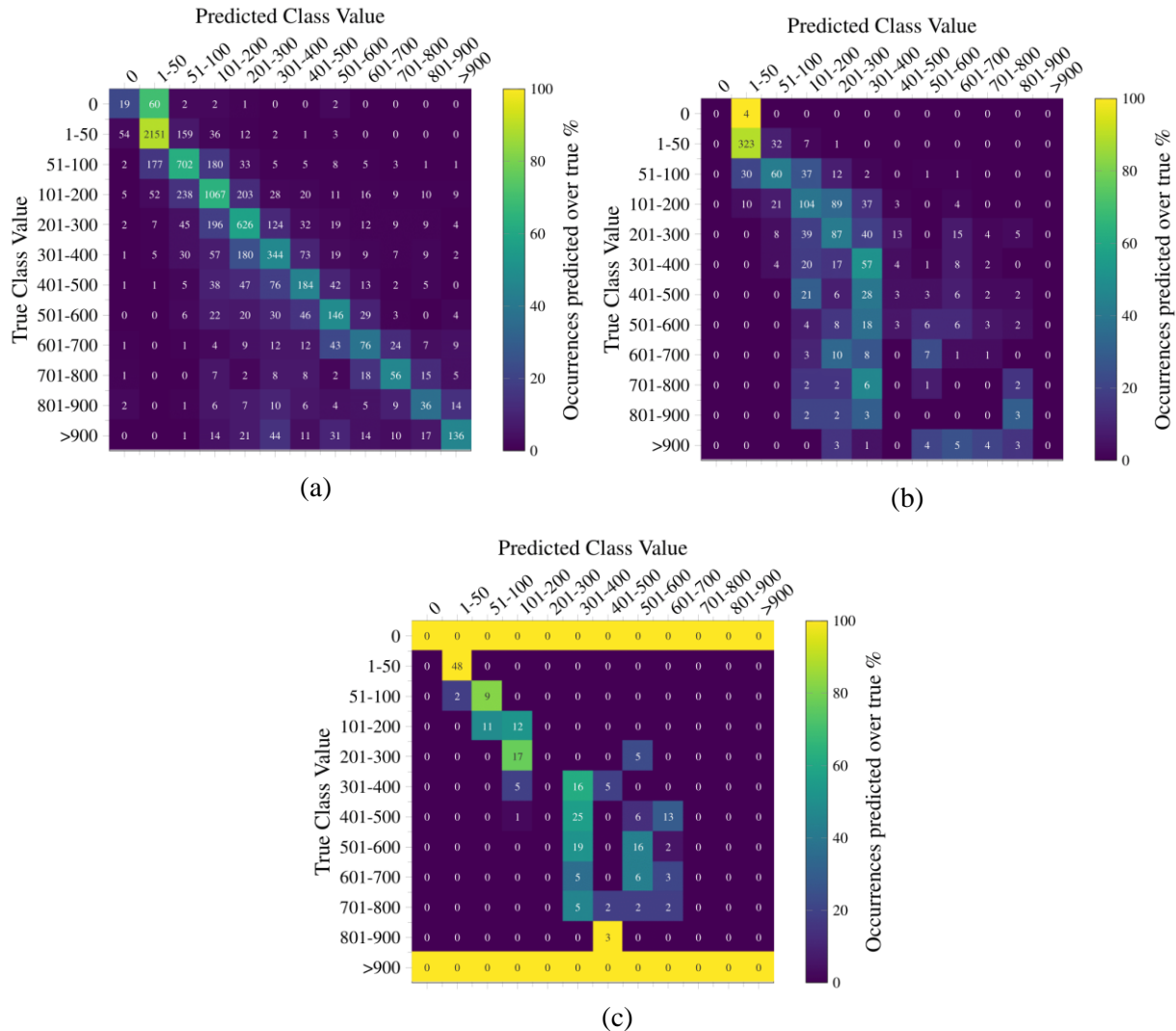
Random Forest (RF)	0.9446	0.6020	0.7446	0.7425	
Extra Trees (ET)	0.9166	0.6075	0.7473	0.7458	
Logistic Regression (LR)	0.9012	0.4965	0.5988	0.5584	
K Neighbors Classifier (KNN)	0.9005	0.5881	0.7068	0.6933	
Gradient Boosting (GB)	0.8948	0.4843	0.5888	0.5753	
Decision Tree (DT)	0.8873	0.5999	0.7412	0.7391	
Linear Discriminant Analysis (LDA)	0.8849	0.4511	0.5711	0.5126	
Naive Bayes (NB)	0.7064	0.2683	0.4755	0.2091	
Ada Boost (ADA)	0.5936	0.1791	0.3530	0.2025	
Quadratic Discriminant Analysis (QDA)	0.5023	0.0881	0.2861	0.0685	
Tuned versions after hyperparameter tuning					
Model	Behavior of the Learning Curve	AUC	Recall	Precision	F1-score
<i>Initial lightGBM</i>	<i>Good</i>	<i>0.9478</i>	<i>0.5933</i>	<i>0.7166</i>	<i>0.7140</i>
Tuned lightGBM	Good	0.9569	0.6066	0.7456	0.7441
Initial RF	Good	0.9446	0.6020	0.7446	0.7425
Tuned RF	Good	0.8909	0.5126	0.6132	0.5809
Initial ET	Good	0.9166	0.6075	0.7473	0.7458
Tuned ET	Bad	0.8688	0.4740	0.5690	0.5031

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28

Validation of the models in all regions

After the ML analysis, we created confusion matrices for the six best performing models in order to select the single best model, based on its predictive performance for each range. **Figure 4** shows the confusion matrices of the best model, the Initial lightGBM model. LightGBM is an open-source gradient boosting framework initially developed by Microsoft, which has gained popularity due to its advantages in training speed, high efficiency and accuracy, and low memory usage, among others. The gradient boosting model is a tree-based learning algorithm that combines numerous weak-learner models to build a strong-learner model, where every single weak-learner model meaningfully improves the performance of the entire model (62, 63). Moving to the description of the matrices, in **Figure 4** each cell contains the number of occurrences of a predicted class when testing true inputs. The legend bars colors represent the percentage of correct predicted occurrences over the total of true class values. The columns show predicted class values, while the sum of all values in each row indicates the total occurrences for each such class. The diagonal presents the occurrences of correct predictions for each class, and the accuracy is the sum of all elements on the diagonal divide by all elements of the matrix. In more detail, **Figure 4a** shows the performances of the classification algorithm for external regions; as explained earlier the model was trained on the data of these regions. In this case, the matrix depicts a satisfactory accuracy for this area (64.8%), and most of the predicted occurrences are on the diagonal, except for class 0, in which occurrences are often predicted as class 1-50. The reason is clearly connected with the length of this class that includes only one value. Analyzing the middle regions, **Figure 4b** shows that the total accuracy for this area (49.7%) is lower than for the external regions. Despite the accuracy, this result is relevant for this study since mispredictions of our algorithm are often between classes close to each other. Since the number of classes is rather high, and the focus of the analysis is to estimate the general trend of the TCs, errors between neighboring classes is acceptable. Note that the highest number of prediction errors occur for higher classes (more than 500), this is likely due to having fewer occurrences of high traffic volumes, so the model is less trained for those values. The last confusion matrix in **Figure 4c** focuses on the center regions. In this area we have only one real TC, therefore the comparison of true and predicted traffic volumes is based on very few values. The accuracy for this area (43.3%) is similar to the one obtained from middle regions with the difference that

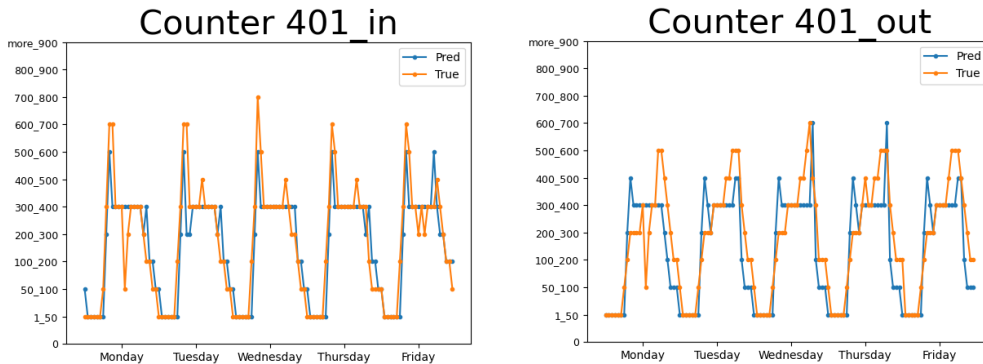
1 in the matrix we have some classes with 0 true occurrences, such as the class 0 and the class >900. As we
 2 noticed for the middle regions's predictions, also for center areas it is clear how the model finds more
 3 difficulties with higher volumes. The matrix shows us how accuracy decreases as traffic volumes increase.
 4



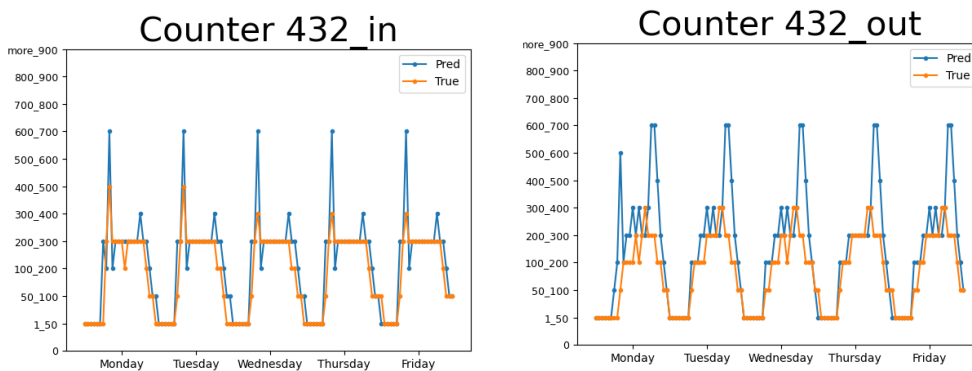
5
 6 **Figure 4 Confusion matrices for external (a), middle (b) and center (c) regions created using the**
 7 **Initial lightGBM model**

8 To show the results in **Figure 4** more intuitively, below we present the predicted temporal profile at TC
 9 level. First, we order the prediction output by TCs, time, and direction of TCs, then we compare the
 10 predicted temporal profiles with the real ones. **Figure 5** shows the results for two TCs (ids: 401 and 432)
 11 that are located in middle region. The plots are divided by direction of the TC (inflow and outflow). From
 12 figure **Figure 5a**, we can observe how the predicted traffic volumes of TC 401 resemble the true values for
 13 both directions, especially the inflow plot shows high accuracy throughout the week. It is interesting to
 14 notice that the prediction are accurate also on different day profiles, e.g. we can see how the model is able
 15 to detect the differences between the Monday and the Friday profile. **Figure 5b** presents the temporal profile
 16 of TC 432 in which the prediction results are characterized by a lower accuracy compared to TC 401,
 17 especially for the outflow direction. Despite this, we can notice how for both directions the general trends
 18 are similar, the peaks occur in the same time intervals, while what our model is not able to estimate correctly

1 is the intensity of the peaks. For instance, the 432 Inflow plot shows that the maximum peaks reach the
 2 400-500 class, while the predicted values indicate higher volumes (class 600-700).
 3



(a)

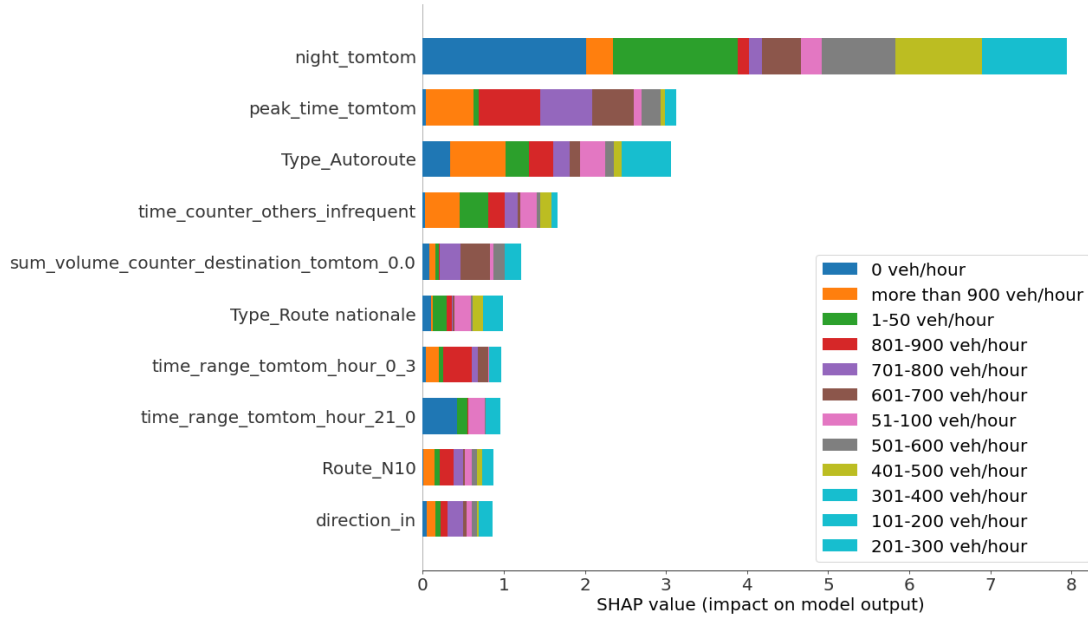


(b)

4 **Figure 5 Comparison between predicted and true traffic flows for TCs 401 (a), and 432 (b) in**
 5 **middle regions**
 6

7
 8 **Identification of the importance of the input features in the predicted traffic volume ranges**

9 **Figure 6** presents the importance of the most relevant input features shown in **Table 2** for the traffic
 10 volume ranges in Luxembourg and the direct impact of each one on the output magnitude for the initial
 11 lightGBM model. All features are presented in descending order according to the impact of the features on
 12 the output of the model, and it is also possible to see the magnitude of the prediction power of the feature
 13 for each traffic volume range. According to this, the results suggest that indicators for night, peak time, and
 14 road type are the three features that have a greatest impact on predicting traffic volume in Luxembourg. It
 15 is also possible to visibly identify the different impact of each feature for each range (e.g., night has a
 16 greater impact for traffic volume ranges below 50 vehicles/hour, than for the other ranges; while peak time
 17 has a greater impact on the traffic volume ranges greater than 700 vehicles/hour). Moreover, it is interesting
 18 to note that among road types (Highway, National road, Secondary road, and mobile counting), Highways
 19 (“Autoroute”) and National roads (“Route nationale”) seem to be the most important to identify the traffic
 20 volume range.



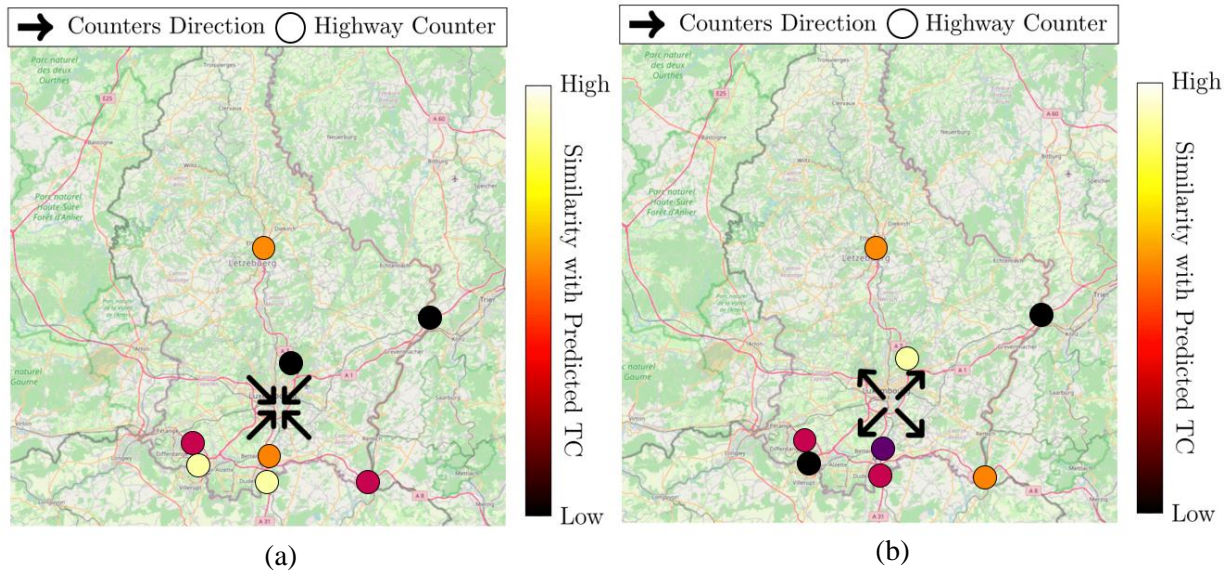
1

2 **Figure 6 SHAP plot with the global feature importance in the lightGBM model**

3 **Similarity analysis**

4 Based on the findings shown in **Figure 6** we decided to focus our attention on TCs installed on
 5 highways and national roads. In particular, a similarity analysis has been performed in order to identify TCs
 6 whose data are the most similar to those estimated for the central regions. **Figure 7** shows the results of the
 7 similarity analysis for highways TCs. Following the color bar, the lighter TCs are the higher is the similarity
 8 level. From this type of analysis, it is possible to understand which TC provides the most similar information
 9 compared with the predicted inflows (**Figure 7a**) and outflows (**Figure 7b**) in the central areas. For national
 10 road TCs, we considered the sum of flows for each TC in each region, to identify the most similar regions
 11 (external and middle) compared to center regions **Figure 8**. Similarly to the case of highway TCs, in this
 12 way it is possible to infer which are the most representative regions of the center in terms of inflow (**Figure**
 13 **8a**) and outflow (**Figure 8b**). More precisely, in order to obtain a first overview of the amount of inflows
 14 and outflows in center regions, this analysis is able to select which TCs, or regions contain the most similar
 15 traffic flows value to those predicted for central areas. Such analysis can be framed as a support tool for
 16 government agencies offering two services: (i) providing information about which TCs can be considered
 17 the most representative for a fist overview of the possible traffic flows of the target area; (ii) the regional
 18 similarity values can be useful to suggest the best area where to locate new TCs that can better describe the
 19 traffic flows of the target area.

1



2

3

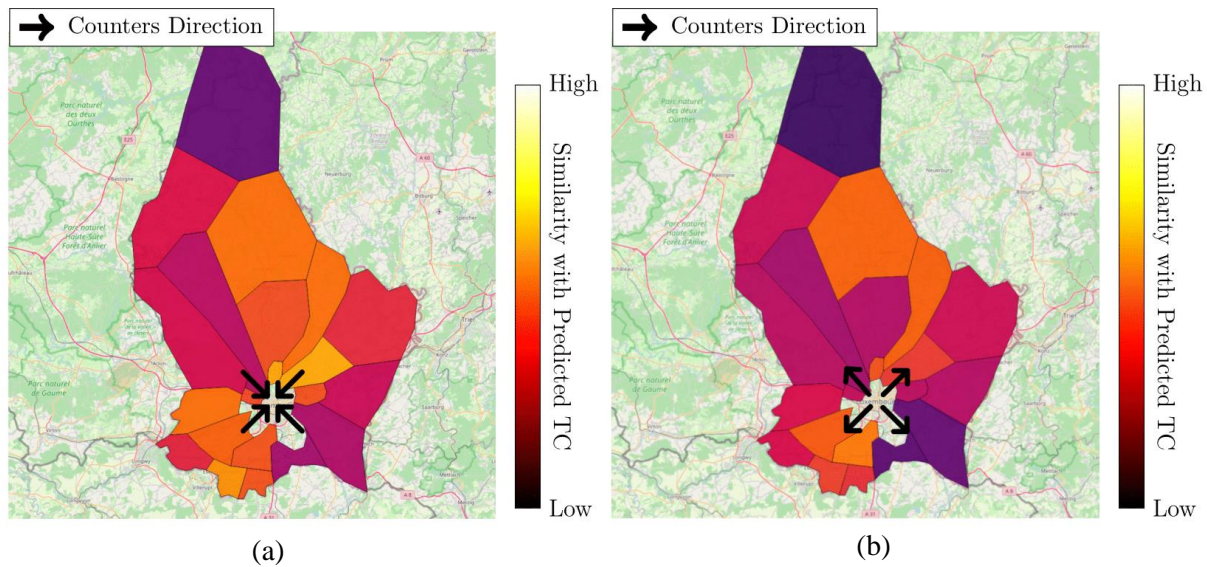
4

Figure 7 Similarity analysis between highway TCs and the predicted flows in center regions considering inflow (a) and outflow (b) direction

5

6

7



8

9

CONCLUSIONS

10

11

12

13

14

15

16

17

The prediction of traffic flows plays a fundamental role in the management of cities and in transport planning; the increasing availability of Open Data information (64) is empowering rapid development of new approaches. In this respect, this paper aims at developing a methodology capable of predicting traffic flows in areas with missing sensors devices using free but limited open datasets whose extracted information has been used to train ML models for prediction, and for validation. Subsequently, SHAP and similarity analysis have been used to identify the most important features and the most representative TCs and regions, to predict traffic flows in center regions. From the analysis carried out we demonstrate that the model is able to predict traffic flows with an acceptable accuracy level. The final purpose is to provide to

1 government agencies a first-step exploratory approach to gain an overview of the traffic flows in areas with
2 little or no other data sources, suggesting also the most appropriate location (i.e., road type and regions)
3 where to install new TCs. However, the methodology has some limitations that can be investigated in further
4 research, e.g. finding data from other data sources in order to have more heterogeneous features at a
5 microscopic level that can improve the accuracy level of the predicted flows (e.g. travel times, path-choice).

6

7 **ACKNOWLEDGMENTS**

8 Mrs. Fazio is partially supported by EU4EU fund. Mr. Vitello is supported by the Luxembourg National
9 Research Fund (PRIDE17/12252781/DRIVEN).

10

11 **AUTHOR CONTRIBUTIONS**

12 The authors confirm contribution to the paper as follows: study conception and design: M. Fazio, P. Vitello,
13 J. Pineda-Jaramillo, R. D. Connors, F. Viti; data collection: M. Fazio, P. Vitello, J. Pineda-Jaramillo;
14 analysis and interpretation of results: M. Fazio, P. Vitello, J. Pineda-Jaramillo; draft manuscript preparation:
15 M. Fazio, P. Vitello, J. Pineda-Jaramillo. All authors reviewed the results and approved the final version of
16 the manuscript.

REFERENCES

1. Chen, D. Research on Traffic Flow Prediction in the Big Data Environment Based on the Improved RBF Neural Network. *IEEE Transactions on Industrial Informatics*, 2017. <https://doi.org/10.1109/TII.2017.2682855>.
2. Bhavsar, P., I. Safro, N. Bouaynaya, R. Polikar, and D. Dera. Machine Learning in Transportation Data Analytics. In *Data Analytics for Intelligent Transportation Systems*.
3. Laney, D. 3D Data Management: Controlling Data Volume, Velocity, and Variety. *Application Delivery Strategies*, 2001.
4. Zantalis, F., G. Koulouras, S. Karabetsos, and D. Kandris. A Review of Machine Learning and IoT in Smart Transportation. *Future Internet*.
5. Alexander, L., S. Jiang, M. Murga, and M. C. González. Origin-Destination Trips by Purpose and Time of Day Inferred from Mobile Phone Data. *Transportation Research Part C: Emerging Technologies*, 2015. <https://doi.org/10.1016/j.trc.2015.02.018>.
6. Cantelmo, G., and F. Viti. A Big Data Demand Estimation Framework for Multimodal Modelling of Urban Congested Networks. 2019.
7. González, M. C., C. A. Hidalgo, and A. L. Barabási. Understanding Individual Human Mobility Patterns. *Nature*, 2008. <https://doi.org/10.1038/nature06958>.
8. Wang, F., J. Wang, J. Cao, C. Chen, and X. (Jeff) Ban. Extracting Trips from Multi-Sourced Data for Mobility Pattern Analysis: An App-Based Data Example. *Transportation Research Part C: Emerging Technologies*, 2019. <https://doi.org/10.1016/j.trc.2019.05.028>.
9. Gong, H., C. Chen, E. Bialostozky, and C. T. Lawson. A GPS/GIS Method for Travel Mode Detection in New York City. *Computers, Environment and Urban Systems*, 2012. <https://doi.org/10.1016/j.compenvurbsys.2011.05.003>.
10. Shafique, M. A., and E. Hato. Travel Mode Detection with Varying Smartphone Data Collection Frequencies. *Sensors (Switzerland)*, 2016. <https://doi.org/10.3390/s16050716>.
11. Strauss, J., L. F. Miranda-Moreno, and P. Morency. Mapping Cyclist Activity and Injury Risk in a Network Combining Smartphone GPS Data and Bicycle Counts. *Accident Analysis and Prevention*, 2015. <https://doi.org/10.1016/j.aap.2015.07.014>.
12. Fazio, M., N. Giuffrida, G. Inturri, and M. Ignaccolo. Combining GPS-Tracks and Accident Data to Improve Safety of Cycling Paths. 2020.
13. Mohammed, O., and J. Kianfar. A Machine Learning Approach to Short-Term Traffic Flow Prediction: A Case Study of Interstate 64 in Missouri. 2019.
14. Yang, J., Y. Han, Y. Wang, B. Jiang, Z. Lv, and H. Song. Optimization of Real-Time Traffic Network Assignment Based on IoT Data Using DBN and Clustering Model in Smart City. *Future Generation Computer Systems*, 2020. <https://doi.org/10.1016/j.future.2017.12.012>.
15. Fusco, G., C. Colombaroni, L. Comelli, and N. Isaenko. Short-Term Traffic Predictions on Large Urban Traffic Networks: Applications of Network-Based Machine Learning Models and Dynamic Traffic Assignment Models. 2015.
16. Yu, J., G. L. Chang, H. W. Ho, and Y. Liu. Variation Based Online Travel Time Prediction Using Clustered Neural Networks. 2008.
17. Hagenauer, J., and M. Helbich. A Comparative Study of Machine Learning Classifiers for Modeling Travel Mode Choice. *Expert Systems with Applications*, 2017. <https://doi.org/10.1016/j.eswa.2017.01.057>.
18. Zhao, X., X. Yan, A. Yu, and P. Van Hentenryck. Prediction and Behavioral Analysis of Travel Mode Choice: A Comparison of Machine Learning and Logit Models. *Travel Behaviour and Society*, 2020. <https://doi.org/10.1016/j.tbs.2020.02.003>.
19. Pineda-Jaramillo, J. Travel Time, Trip Frequency and Motorised-Vehicle Ownership: A Case Study of Travel Behaviour of People with Reduced Mobility in Medellín. *Journal of Transport and Health*, 2021. <https://doi.org/10.1016/j.jth.2021.101110>.
20. Chikaraishi, M., P. Garg, V. Varghese, K. Yoshizoe, J. Urata, Y. Shiomi, and R. Watanabe. On the

- Possibility of Short-Term Traffic Prediction during Disaster with Machine Learning Approaches: An Exploratory Analysis. *Transport Policy*, 2020. <https://doi.org/10.1016/j.tranpol.2020.05.023>.
21. Lv, Y., Y. Duan, W. Kang, Z. Li, and F. Y. Wang. Traffic Flow Prediction with Big Data: A Deep Learning Approach. *IEEE Transactions on Intelligent Transportation Systems*, 2015. <https://doi.org/10.1109/TITS.2014.2345663>.
 22. Yi, H., J. Heejin, and S. Bae. Deep Neural Networks for Traffic Flow Prediction. 2017.
 23. Wang, S. J., and P. Moriarty. Can New Communication Technology Promote Sustainable Transport? 2017.
 24. Portail Des Travaux Publics. <https://travaux.public.lu/fr/infos-traffic/comptage.html>.
 25. TOMTOMmove. <https://move.tomtom.com/login>.
 26. Soua, R., A. Koesdwiady, and F. Karray. Big-Data-Generated Traffic Flow Prediction Using Deep Learning and Dempster-Shafer Theory. 2016.
 27. Polson, N. G., and V. O. Sokolov. Deep Learning for Short-Term Traffic Flow Prediction. *Transportation Research Part C: Emerging Technologies*, 2017. <https://doi.org/10.1016/j.trc.2017.02.024>.
 28. Wu, Y., H. Tan, L. Qin, B. Ran, and Z. Jiang. A Hybrid Deep Learning Based Traffic Flow Prediction Method and Its Understanding. *Transportation Research Part C: Emerging Technologies*, 2018. <https://doi.org/10.1016/j.trc.2018.03.001>.
 29. Baras, J. S., and W. S. Levine. ESTIMATION OF TRAFFIC FLOW PARAMETERS IN URBAN TRAFFIC NETWORKS. 1977.
 30. Chang, M. F., and D. C. Gazis. TRAFFIC DENSITY ESTIMATION WITH CONSIDERATION OF LANE-CHANGING. *Transportation Science*, 1975. <https://doi.org/10.1287/trsc.9.4.308>.
 31. Okutani, I., and Y. J. Stephanedes. Dynamic Prediction of Traffic Volume through Kalman Filtering Theory. *Transportation Research Part B*, 1984. [https://doi.org/10.1016/0191-2615\(84\)90002-X](https://doi.org/10.1016/0191-2615(84)90002-X).
 32. Kawashima, H. Long Term Prediction of Traffic Flow. *IFAC Proceedings Volumes*, 1987. [https://doi.org/10.1016/s1474-6670\(17\)55876-0](https://doi.org/10.1016/s1474-6670(17)55876-0).
 33. Nicholson, H., and C. D. Swann. The Prediction of Traffic Flow Volumes Based on Spectral Analysis. *Transportation Research*, 1974. [https://doi.org/10.1016/0041-1647\(74\)90030-6](https://doi.org/10.1016/0041-1647(74)90030-6).
 34. Lu, J. Prediction of Traffic Flow by an Adaptive Prediction System. *Transportation Research Record*, No. 1287, 1990.
 35. Williams, B. M., P. K. Durvasula, and D. E. Brown. Urban Freeway Traffic Flow Prediction: Application of Seasonal Autoregressive Integrated Moving Average and Exponential Smoothing Models. *Transportation Research Record*, 1998. <https://doi.org/10.3141/1644-14>.
 36. Fazio, M., M. Le Pira, G. Inturri, and M. Ignaccolo. Bus Rapid Transit vs. Metro. Monitoring on-Board Comfort of Competing Transit Services via Sensors. 2020.
 37. Cantelmo, G., and F. Viti. A Big Data Demand Estimation Model for Urban Congested Networks. 2020.
 38. Cantelmo, G., P. Vitello, B. Toader, Antoniou, and F. Viti. Inferring Urban Mobility and Habits from User Location History. 2020.
 39. Moghaddam, S. S., and B. Hellinga. Real-Time Prediction of Arterial Roadway Travel Times Using Data Collected by Bluetooth Detectors. *Transportation Research Record*.
 40. Iqbal, M. S., C. F. Choudhury, P. Wang, and M. C. González. Development of Origin-Destination Matrices Using Mobile Phone Call Data. *Transportation Research Part C: Emerging Technologies*, 2014. <https://doi.org/10.1016/j.trc.2014.01.002>.
 41. Di Donna, S. A., G. Cantelmo, and F. Viti. A Markov Chain Dynamic Model for Trip Generation and Distribution Based on CDR. 2015.
 42. Xu, C., J. Ji, and P. Liu. The Station-Free Sharing Bike Demand Forecasting with a Deep Learning Approach and Large-Scale Datasets. *Transportation Research Part C: Emerging Technologies*, 2018. <https://doi.org/10.1016/j.trc.2018.07.013>.
 43. Yin, H., S. C. Wong, J. Xu, and C. K. Wong. Urban Traffic Flow Prediction Using a Fuzzy-Neural

- Approach. *Transportation Research Part C: Emerging Technologies*, 2002. [https://doi.org/10.1016/S0968-090X\(01\)00004-3](https://doi.org/10.1016/S0968-090X(01)00004-3).
44. Hochreiter, S. Long Short-Term Memory. Vol. 1780, 1997, pp. 1735–1780.
 45. Jo, D., B. Yu, H. Jeon, and K. Sohn. Image-to-Image Learning to Predict Traffic Speeds by Considering Area-Wide Spatio-Temporal Dependencies. *IEEE Transactions on Vehicular Technology*, 2019. <https://doi.org/10.1109/TVT.2018.2885366>.
 46. Munoz-Organero, M., R. Ruiz-Blaquez, and L. Sánchez-Fernández. Automatic Detection of Traffic Lights, Street Crossings and Urban Roundabouts Combining Outlier Detection and Deep Learning Classification Techniques Based on GPS Traces While Driving. *Computers, Environment and Urban Systems*, 2018. <https://doi.org/10.1016/j.compenvurbsys.2017.09.005>.
 47. Bauer, D., G. Richter, J. Asamer, B. Heilmann, G. Lenz, and R. Kölbl. Quasi-Dynamic Estimation of OD Flows from Traffic Counts Without Prior OD Matrix. *IEEE Transactions on Intelligent Transportation Systems*, 2018. <https://doi.org/10.1109/TITS.2017.2741528>.
 48. Pineda-Jaramillo, J. D. A Shallow Neural Network Approach for Identifying the Leading Causes Associated to Pedestrian Deaths in Medellín. *Journal of Transport and Health*, 2020. <https://doi.org/10.1016/j.jth.2020.100912>.
 49. Xue, W., and J. Zhang. Dealing with Imbalanced Dataset: A Re-Sampling Method Based on the Improved SMOTE Algorithm. *Communications in Statistics - Simulation and Computation*, Vol. 45, No. 4, 2016, pp. 1160–1172. <https://doi.org/10.1080/03610918.2012.728274>.
 50. Shafiq, M., Z. Tian, A. K. Bashir, A. Jolfaei, and X. Yu. Data Mining and Machine Learning Methods for Sustainable Smart Cities Traffic Classification: A Survey. *Sustainable Cities and Society*, Vol. 60, 2020, p. 102177. <https://doi.org/10.1016/j.scs.2020.102177>.
 51. Gao, K., Y. Yang, T. Zhang, A. Li, and X. Qu. An Extrapolation-Enhanced Approach for Modeling Travel Decision Making: Integrating Ensemble Machine Learning with Knowledge-Based Decision-Making Theory. *Knowledge-Based Systems*, 2021, p. 106882. <https://doi.org/10.1016/j.knosys.2021.106882>.
 52. Ali, M. PyCaret: An Open Source, Low-Code Machine Learning Library in Python.
 53. Martínez Fernández, P., I. Villalba Sanchís, V. Yepes, and R. Insa Franco. A Review of Modelling and Optimisation Methods Applied to Railways Energy Consumption. *Journal of Cleaner Production*, Vol. 222, 2019, pp. 153–162. <https://doi.org/10.1016/j.jclepro.2019.03.037>.
 54. Wang, F., and C. L. Ross. Machine Learning Travel Mode Choices: Comparing the Performance of an Extreme Gradient Boosting Model with a Multinomial Logit Model. *Transportation Research Record: Journal of the Transportation Research Board*, Vol. 2672, No. 47, 2018, pp. 35–45. <https://doi.org/10.1177/0361198118773556>.
 55. Hillel, T., M. Bierlaire, and Y. Jin. *A Systematic Review of Machine Learning Methodologies for Modelling Passenger Mode Choice*. Lausanne, 2019.
 56. Witten, I. H., E. Frank, M. A. Hall, and C. J. Pal. *Data Mining: Practical Machine Learning Tools and Techniques*. Elsevier, 2016.
 57. Bergstra, J., and Y. Bengio. Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research*, Vol. 13, No. 10, 2012, pp. 281–305.
 58. Meek, C., B. Thiesson, and D. Heckerman. The Learning-Curve Sampling Method Applied to Model-Based Clustering. *Journal of Machine Learning Research*, 2002. <https://doi.org/10.1162/153244302760200678>.
 59. Xu, J., A. Wang, N. Schmidt, M. Adams, and M. Hatzopoulou. A Gradient Boost Approach for Predicting Near-Road Ultrafine Particle Concentrations Using Detailed Traffic Characterization. *Environmental Pollution*, Vol. 265, 2020, p. 114777. <https://doi.org/10.1016/j.envpol.2020.114777>.
 60. Lundberg, S. M., and S. I. Lee. A Unified Approach to Interpreting Model Predictions. 2017.
 61. Lin, J. Divergence Measures Based on the Shannon Entropy. *IEEE Transactions on Information Theory*, 1991. <https://doi.org/10.1109/18.61115>.
 62. Wang, J., B. Liu, T. Fu, S. Liu, and J. Stipanovic. Modeling When and Where a Secondary

- Accident Occurs. *Accident Analysis and Prevention*, Vol. 130, 2019, pp. 160–166.
<https://doi.org/10.1016/j.aap.2018.01.024>.
63. Nti, I. K., A. F. Adekoya, and B. A. Weyori. A Comprehensive Evaluation of Ensemble Learning for Stock-Market Prediction. *Journal of Big Data*, 2020. <https://doi.org/10.1186/s40537-020-00299-5>.
64. Yadav, P., S. Hasan, A. Ojo, and E. Curry. The Role of Open Data in Driving Sustainable Mobility in Nine Smart Cities. 2017.