

1 Global Positioning System Data to Model Network-Wide Road Segment
2 Level Number of Lanes Using Spatial Analysis and Machine Learning
3

4 Adham Badran ^{a*}, Ahmed El-Geneidy ^b, and Luis Miranda-Moreno ^a

5 ^a Civil Engineering Department, McGill University, 817 Sherbrooke Street West, Montreal, H3A 0C3, Canada

6 ^b School of Urban Planning, McGill University, 815 rue Sherbrooke West, Montreal, H3A0C2, Canada

7 * Corresponding Author
8

9 Contact: adham.badran@mail.mcgill.ca, Department of Civil Engineering, McGill University, 817
10 Sherbrooke Street West, Montreal H3A 0C3, Canada.
11

12
13 Word Count: 5363
14

15 Submitted on July 27th, 2023
16
17
18
19
20
21

22 For Citation Please use: Badran, A., El-Geneidy, A., Miranda-Moreno, L. (2024). *Global positioning*
23 *system data to model network-wide road segment level number of lanes using machine learning*. Paper
24 presented at the 103rd Transportation Research Board Annual Meeting, Washington DC, USA.

1 Abstract

2 One of the main features required in transport network modelling is the number of lanes used to
3 estimate the road capacity and predict vehicular travel times based on traffic flows. Traditionally,
4 the number of lanes information is collected manually or more recently extracted using computer
5 vision techniques, which are two resource intensive methods. This research proposes the use of
6 emerging crowd-sensed Global Positioning System (GPS) data to predict the number of lanes per
7 road segment for large scale transport models through geographic operations and machine
8 learning. The developed method consists of i) a spatial analysis to analyze the GPS trajectory data
9 and estimate predictors and ii) a supervised machine learning model development to build a model
10 able to predict the number of lanes per road segment.

11 It was found that the method predicts the number of lanes at an accuracy of 91% using two
12 predictors: number of GPS points per road segment and a lateral distance variable containing 60%
13 of the GPS data points, centered around the lateral distance distribution median. The best
14 prediction model was obtained using decision trees classifier. It was also found that most of the
15 local roads did not have sufficient data points to obtain a stable lateral distance distribution,
16 therefore, the model was limited to a subset of road segments with sufficient observations. Given
17 the availability of high spatiotemporal coverage GPS data, the method can be adapted and applied
18 to large scale road network models and predict the number of lanes accurately and cost-effectively.

19
20 **Keywords:** Global Positioning System, Transport Model, Road Network, Number of Lanes, Road
21 Capacity, EMME.

23 Introduction

24 Knowledge of the number of lanes on road segments within the transport network is essential for
25 the planning and operation of the transport system. For example, conventional and autonomous
26 vehicle navigation, transport modelling and simulation, road safety applications all require the
27 number of lanes information as an input. In fact, lane-level digital maps are critical for advanced
28 driver assistance systems and continuous research is being performed to improve their
29 development (Guo et al., 2016). Moreover, autonomous vehicle navigation requires prior
30 knowledge of the road network in addition to real-time detection of the road lanes to select the
31 trajectory appropriately (Bounini et al., 2015). In transport modelling, the number of lanes
32 information is essential for all modelling scales. Macroscopic models include the number of lanes
33 information into volume-delay functions to determine the road's vehicular capacity and evaluate
34 road segment level travel time. Meanwhile microscopic transport models consider the number of
35 lanes through lane changing models (Treiber and Kesting, 2013). Another example is the use of
36 the number of lanes when analyzing pedestrian-vehicle interaction at crossings and the relationship
37 with road-user safety (Kadali and Vedagiri, 2020).

38 The challenge in obtaining the number of lanes information is for large-scale road networks and
39 maps. In fact, at large scales, the required resources to develop and maintain detailed digital
40 networks become significant which renders the manual development infeasible. With the advances
41 in technology, new data sources and techniques are emerging and present a potential to extract
42 transport network-related information. Global Positioning System (GPS) trajectory data is being
43 collected by different organizations through GPS-enabled smartphones and stored on servers using
44 cellular internet. For example, the city of Montreal has provided its residents a smartphone
45 application that records their trajectories for a limited period to analyze the trajectory data and
46 improve transport planning and reduce traffic delays (Montréal).

1 Extracting the number of lanes has been tackled in the past using different data sources. The most
2 frequent method to extract the number of lanes for large-scale networks is based on aerial imagery
3 and computer vision techniques. Multiple studies have been looking at extracting road network
4 features automatically using different data sources. First, high-resolution imagery, in combination
5 with computer vision methods have enabled the large-scale detection and extraction of road
6 network-related attributes. One of the research groups has done extensive work using road
7 segmentation to detect different visible features such as the road, sidewalks, vegetation, buildings,
8 and cars to augment OpenStreetMap by adding more features (Mattyus et al., 2015). The main
9 challenges were found to be the presence of trees, shadows, cars, as they increase heterogeneity in
10 the images in addition to misalignment issues with respect to the road centreline file used as a
11 priori of road segments' location. The same research group further expanded the analysis by
12 collecting and incorporating street-level imagery in the number of lanes recognition algorithm
13 which increased its prediction accuracy (Máttyus et al., 2016). Recognizing that collecting street-
14 level imagery presents high collection and processing costs, they proposed a more resource-
15 friendly version that only employs satellite imagery but takes advantage of new methodological
16 advances in deep learning to improve the model accuracy (Máttyus et al., 2017). Another study
17 has also extracted the number of lanes information from satellite imagery using an SVM classifier
18 for lane identification based on brightness levels. Although they predicted the number of lanes at
19 an accuracy of 100%, the experiment was only presented for six road segments (Tang et al., 2014).
20 Although satellite imagery has been used to detect the number of lanes and improved by collecting
21 street level high-resolution imagery, it is not without limitations. Data availability is limited due
22 to the collection costs, moreover, occlusions, illumination variability and unmarked road lines
23 reduce the capacity of such techniques (Kasmi et al., 2018). The best number of lanes prediction
24 accuracy obtained was 83 %.

25 Recent research efforts have been studying the extraction of road networks from GPS data using
26 different spatial analysis algorithms. Three main approaches were used to extract road networks:
27 Clustering, intersection linking, and track alignment. For example, the work by Guo et al. (2021)
28 proposes a clustering method to extract road network centreline and intersections with the accuracy
29 of 92%. Clustering is in fact the most popular method to extract road networks from GPS trajectory
30 data. Another study by Zhang et al. (2019) employs the intersection linking method to detect the
31 road network and intersections at an accuracy greater than 90%. Although not very popular, studies
32 by Leichter and Werner (2019) and Zhongyi et al. (2018) have also used the track alignment
33 method to generate road networks. However, accuracy was either low or not compared to the
34 ground truth. Although these road network inference methods are able in some cases to extract the
35 road network centreline and intersections with high accuracy, they do are not designed to extract
36 more detailed road network features such as the number of lanes.

37 Very few studies have examined the use of sole GPS trajectory data to extract the number of lanes.
38 A study by Arman and Tampere (2020) proposes a method that extracts lane locations on a highway
39 corridor. However, the number of lanes extracted is not explicitly validated by comparing to the
40 ground truth. Therefore, no accuracy was provided. One attempt by Zhang et al. (2010) used GPS
41 traces and a road centreline map from OpenStreetMap to improve the map quality and estimate the
42 number of lanes. The main limitation was the assumption of normal distribution of GPS traces
43 with respect to the road centre, which is not the case and resulted in number of lanes prediction
44 accuracy of less than 60%.

45 Moreover a study by Chen and Krumm (2010) fits Gaussian mixture models to GPS trajectory
46 data to determine the number of lanes. Although the study attempts to preserve the continuous

1 nature of road segments, it is limited by the sample size and the fact that this method requires prior
2 knowledge of the number of Gaussian distributions to fit. Thus, the study resulted in relatively low
3 accuracy predictions.

4 In sum, the main limitations of past studies extracting the number of lanes are the high cost of
5 imagery data collection and the output accuracy. In fact, the high cost reduces the frequency of
6 map updates which can result in maps not representing the continuously evolving nature of the
7 road network. In addition, the output accuracy of past studies can potentially be improved by using
8 large-scale GPS trajectory data.

9 Considering the general availability of road centreline data or algorithms to infer them from
10 different data sources, the objective of this study is to propose a method that uses GPS trajectory
11 data to extract the number of lanes with a relatively high accuracy. This is done through spatial
12 analysis of GPS trajectory points to extract variables that feed into a machine learning
13 classification algorithm that predicts the number of lanes for road segments for use in large-scale
14 transport models.
15

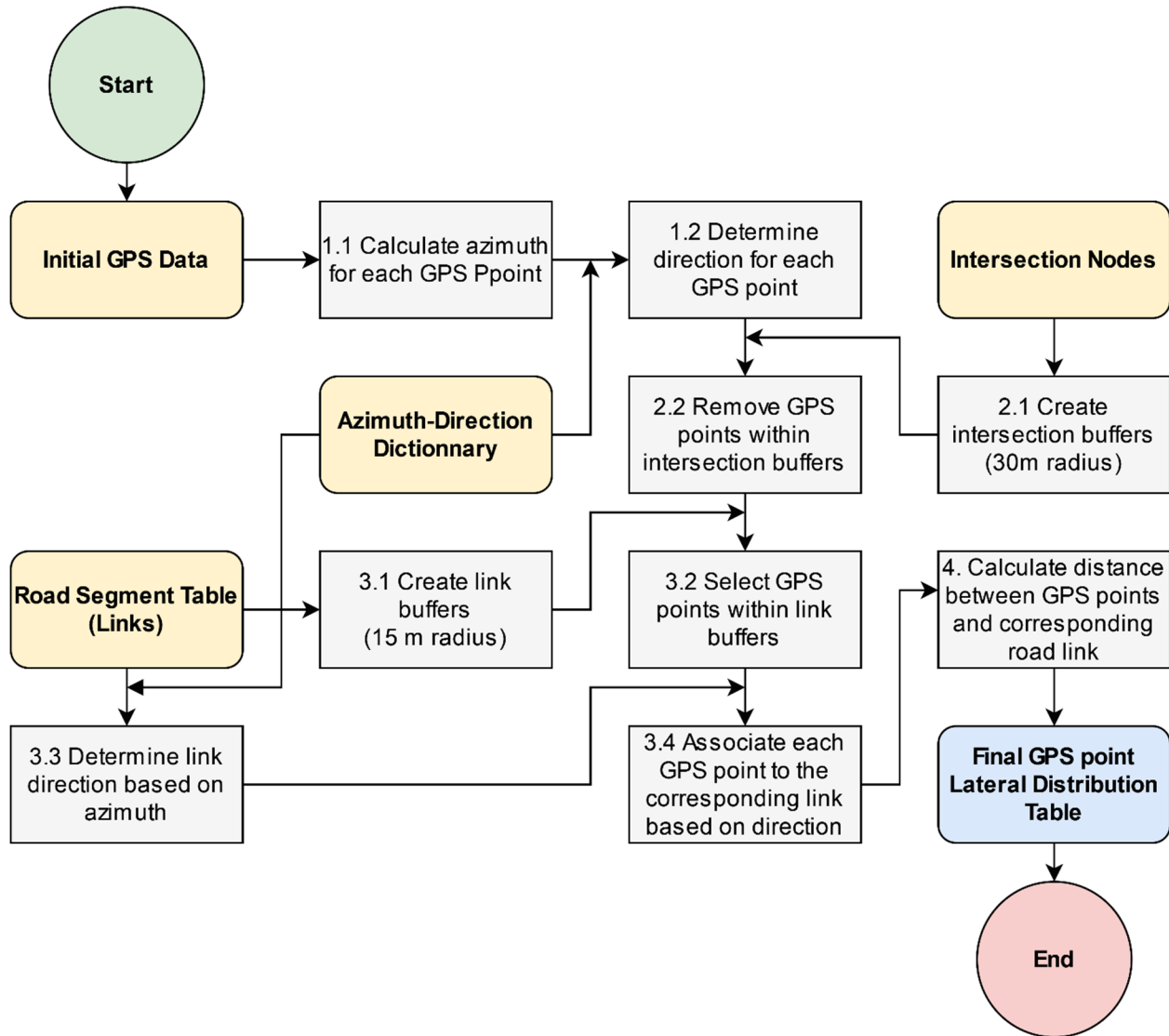
16 Methodology

17 GPS data treatment can be divided into two main parts based on the analysis type. The first part of
18 the analysis was the spatial analysis using Geographic Information System (GIS) software carried
19 out in the FME software. This software was selected since it is a very efficient data integration
20 platform capable of managing, combining, and transforming big data with advanced spatial data
21 analysis capabilities. The second step was the number of lanes prediction model development and
22 visualization carried out in MATLAB. The general assumption of this study is that although GPS
23 accuracy is between 7 and 13 meters (Merry and Bettinger, 2019), GPS trajectory points will be
24 distributed around the middle of traffic lanes when the sample size is large. Therefore, this method
25 proposes to determine the distance distribution of GPS points for each directional link with respect
26 to a reference line and infer the number of lanes based on the distribution properties through a
27 machine learning classification method.

28 Spatial Analysis

29 The first step requires raw GPS trajectory data, a road network model (links and nodes), and an
30 azimuth-direction dictionary table as input. A summary of the spatial analysis steps can be seen in
31 Figure 1. The yellow boxes present the input data sources required to carry out the spatial analysis
32 steps. In this study, each GPS trajectory point had two sets of longitude and latitude points; raw,
33 and map matched coordinates, which were both used at different stages of the analysis.

34 The process can be divided into four main steps: 1. Determine the direction of each GPS trajectory
35 point, 2. Remove GPS trajectory points located at intersections, 3. Associate each GPS trajectory
36 point to a directional road segment, and 4. Calculate the lateral distance between each GPS point
37 and the reference line. The number corresponding to each step is also presented in Figure 1.



1

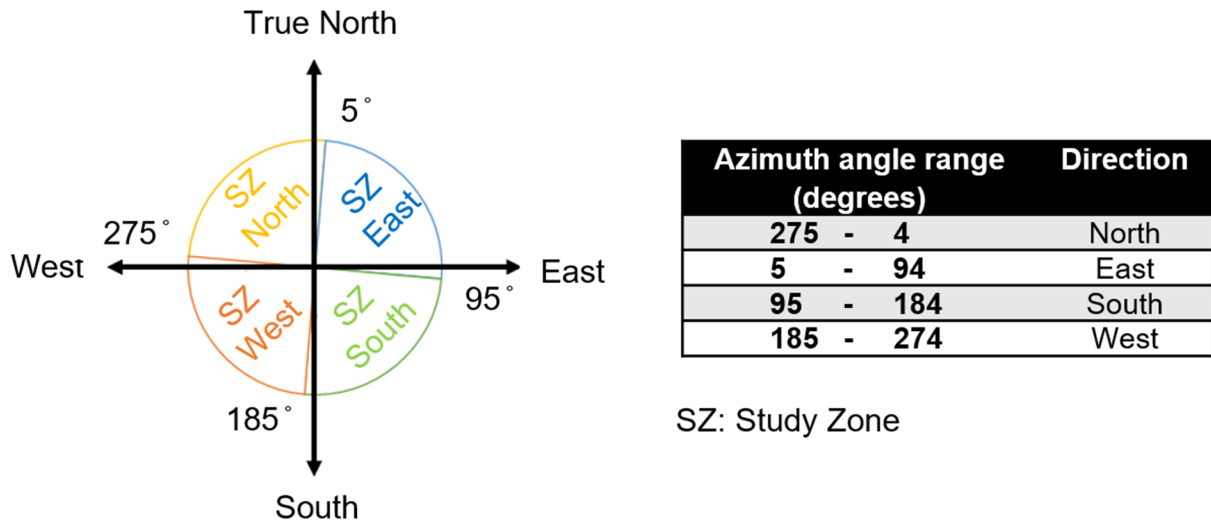
2 *Figure 1 - GPS Trajectory Points GIS Treatment Diagram*

3 Firstly, the azimuth of each GPS trajectory point was calculated based on its location and the
 4 location of the consecutive point within the same trip. The azimuth is defined as the orientation,
 5 in degrees, between two points as the number of degrees clockwise from the north reference. The
 6 azimuth was selected as the measure to define trip segment directions and an azimuth-direction
 7 dictionary was created for that purpose as seen in Figure 2. Map matched coordinates were used
 8 to calculate the azimuth to ensure consistent direction results and remove the fluctuations found in
 9 raw GPS point data. Following the azimuth calculation for each point, the direction was calculated
 10 using the azimuth-direction dictionary.

11 Secondly, intersection buffers were used to remove GPS points that fall within the vicinity of
 12 intersections. Given that this study aims to determine the mid-block road segment number of lanes,
 13 the GPS trajectory points in the vicinity of intersections were removed since the number of lanes
 14 near an intersection is sometimes different to allow for upstream dedicated turning lanes or
 15 downstream insertion lanes. Following visual inspection of the road network, a buffer size of 30-
 16 meter radius with respect to the intersection centres was used to filter GPS trajectory points within

1 intersection areas. This ensured that the remaining GPS trajectory points correspond to travel
 2 within the road segment.

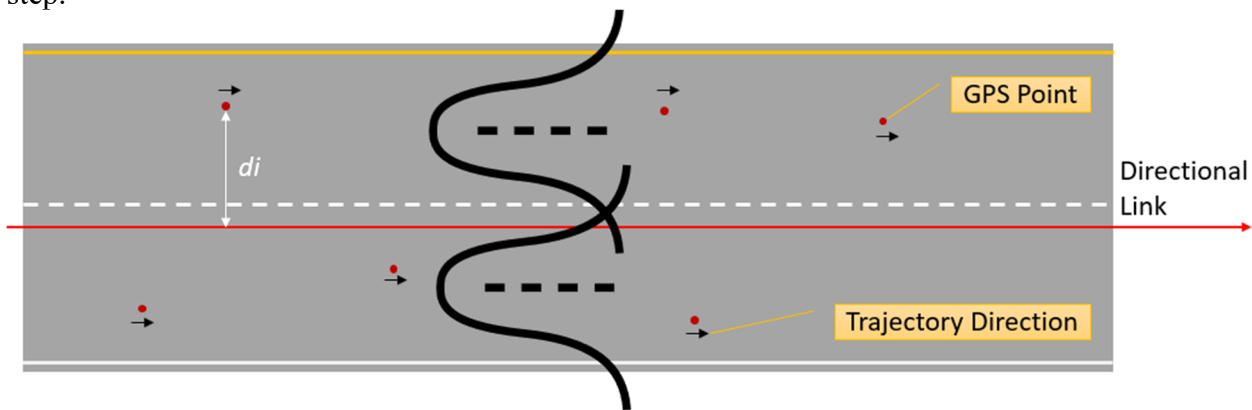
Study Zone Azimuth Definition



3
 4 *Figure 2 - Azimuth-Direction Correspondence*

5 Thirdly, to remove noisy GPS trajectory points, a road segment (link) buffer of 15-meter radius
 6 was created to select GPS points associated to each link through nearest neighbour analysis. This
 7 buffer size was selected to ensure that the GPS points' lateral distribution profile with the respect
 8 to the directional link is captured entirely while minimizing the number of outliers. This was
 9 validated in the following steps of the analysis by examining all lateral distance distribution
 10 histograms and kernel density estimators. The link direction was also obtained based on the
 11 azimuth to add an extra criterion when selecting the nearest neighbour and ensure that every GPS
 12 point is associated to the correct directional link.

13 Fourthly, the shortest distance between each GPS point and the associated directional link is
 14 calculated and serves in the following step develop a number of lanes prediction model. This
 15 distance corresponds to the length of the perpendicular line, d_i , between the GPS location point
 16 and the directional link as seen in Figure 3. It was the main variable carried to the next modelling
 17 step.



18
 19 *Figure 3 - Distance from Point to Link Calculation*

1 The location of the directional link with respect to the actual road segment is approximate since it
 2 is based on a simple street centreline shapefile. Moreover, for bidirectional road segments, the
 3 links for both directions are superimposed. This was taken into account while determining road
 4 segment buffer size.

5 Road Segment Number of Lanes Prediction Model

6 Following spatial analysis, the second part of the method consisted of creating the number of lanes
 7 prediction model. Assuming that road segments with different numbers of lanes have different GPS
 8 trajectory data characteristics (such as spatial distribution pattern and number of points), roads
 9 with different numbers of lanes are seen as distinct categories and the question is formulated as a
 10 number of lanes classification problem. A classification model was calibrated using input variables
 11 derived from GPS trajectory points data to output the number of lanes for each road segment. For
 12 each road segment, input variables were compiled following the spatial analysis part and were
 13 used to train the model to predict the number of lanes as a categorical variable. The two main GPS
 14 trajectory points descriptors, used to derive input variables to the classification tree model, were
 15 the lateral distance d_i and the number of points per directional road segment. A frequency
 16 histogram and a kernel density estimator were fitted to the distance variable to visualize the
 17 distribution with respect to the reference line (directional link) and determine model parameters.
 18 First, it was observed that for some of the links, sample size was too low and resulted in unstable
 19 and unmeaningful distributions. Following inspection of the kernel density estimators and
 20 frequency histogram, the sample size was limited to a minimum of 500 GPS points per directional
 21 link to produce stable results in terms of distribution shape. Road segments with fewer GPS point
 22 observations were removed.

23 Based on the observed distributions and preliminary tests and aiming to create variables that reflect
 24 the lateral distribution of GPS trajectory points with respect to the link, distance percentiles, d_{ipc} ,
 25 were calculated for different percentiles, i , of 5%, 10%, 15%, 20%, 80%, 85%, 90%, and 95%. To
 26 standardize these values and render them comparable across different links, new variables were
 27 created by calculating the variable D_p defined as the lateral distance containing a proportion, p , of
 28 the GPS points data. D_p is calculated using lateral distance percentiles to ensure that this new
 29 variable is centred around the median distance value. The following are the lateral distance
 30 variables that were calculated:

$$31 \quad D_{90} = d_{95pc} - d_{5pc}$$

$$32 \quad D_{80} = d_{90pc} - d_{10pc}$$

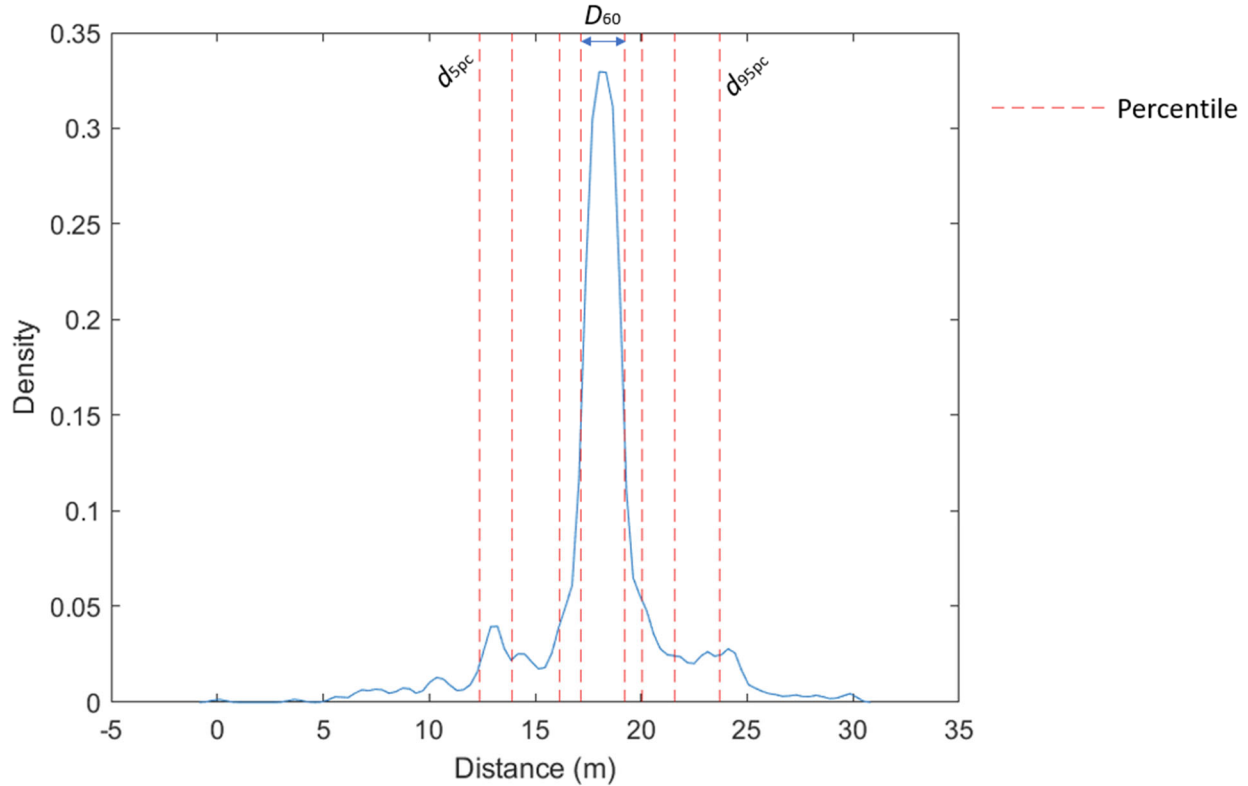
$$33 \quad D_{70} = d_{85pc} - d_{15pc}$$

$$34 \quad D_{60} = d_{80pc} - d_{20pc}$$

35 For example, D_{60} corresponds to the difference between the 80th percentile distance and 20th
 36 percentile distance, therefore it contains 60% of the GPS points data. A visual illustration is
 37 provided in Figure 4. The number of GPS points per link and the standard deviation of lateral
 38 distance per link were also calculated to be tested in the model specification.

39 Following the creation of the variables for each road segment, supervised machine learning
 40 classification methods were tested. In fact, classification tree analysis was carried out to determine
 41 if it can create an accurate model that can be used for prediction. This method is a good option
 42 when ground truth data is available for the learning step. Moreover, it is non-parametric and does
 43 not require prior knowledge of the distribution of each variable. Another advantage of this method
 44 compared to other machine learning techniques such as neural networks classification is its
 45 transparency which makes the model easy to interpret (Ian et al., 2017).

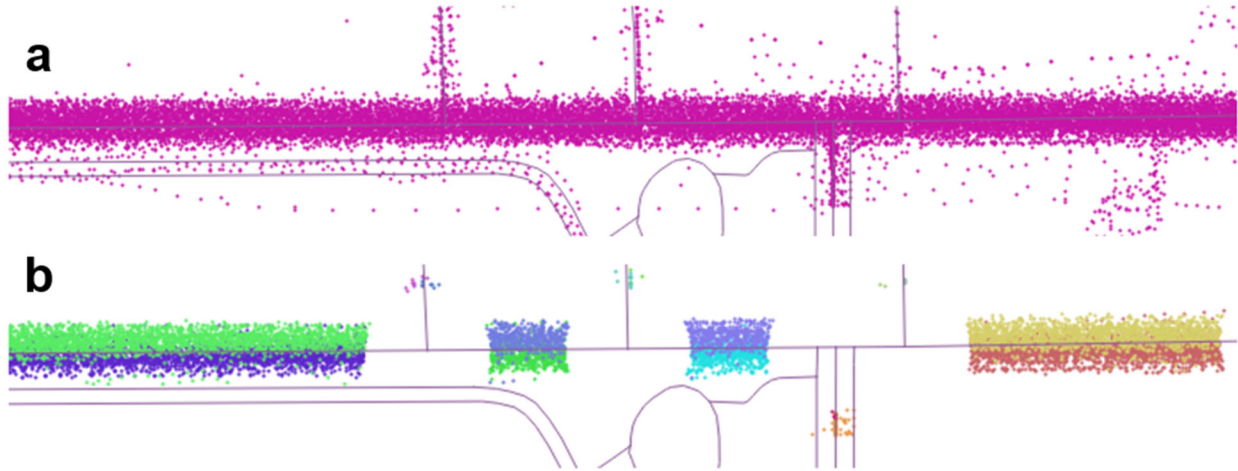
1 To ensure protection against overfitting, model validation was carried out using a 5-fold cross-
2 validation. This validation method divides the dataset randomly into five groups. At each step, one
3 of the five groups is held out to be used for validation while the other four groups are used to train
4 the model. Once the model is specified, it is used to make predictions on the group that was held
5 out. For a 5-fold cross validation, this process is repeated five times.



6
7 *Figure 4 - Example of Percentile Visualization*

8
9 **Data**

10 Three main input datasets are used: 1) GPS trajectory points, 2) Modelled directional road network
11 (links and nodes), and 3) Google maps and Street View. GPS data was collected during the spring
12 of 2014 in Quebec City, Canada. It was collected for 21 days by 2000 voluntary users through the
13 Mon Trajet smartphone app, made available by the Municipality. Each point is described by the
14 following attributes: map matched X and Y coordinates, trip ID, speed, and timestamp (Year-
15 Month-Day-Hour-Minute-Second). Following the preprocessing steps, 245 links were selected as
16 the experimental data to model number of lanes, which included 120 000 GPS points (excluding
17 GPS points within the intersection buffers. This study area was selected based on its urban setting
18 since it is in the city centre where more GPS trajectories were available. Figure 5 presents a sample
19 of the study area where part a shows the raw GPS trajectory points, and part b shows the processed
20 GPS trajectory points for the same road corridor following spatial analysis steps 1 to 3. In part b
21 of the figure, GPS trajectory points are colored differently depending on the link to which they
22 were associated.



1
2 *Figure 5 - Sample of GPS Points Data in Study Area - Before and After Spatial Processing*

3 The directional road network was created using an initial road centreline shapefile which was
4 converted in a network model compatible the EMME transport modelling software to obtain
5 directional links and augmented using the same GPS trajectory data to ensure that road topology
6 and connectivity are valid. Each link is defined by an origin and a destination node. The possible
7 number of lanes per directional link was one, two, or three lanes, for which the ground truth was
8 manually extracted using Google Maps and Street View.

9 For a given road segment, it is important to note that the number of lanes available for traffic can
10 vary spatially and temporally. The presence of lanes dedicated to transit vehicles or high-
11 occupancy vehicles at a fixed schedule on concerned road segments reduces temporally the number
12 of lanes available to general traffic. This is also the case for lanes that are used for parking at fixed
13 schedules. Throughout a road segment, the number of lanes can also change spatially. For example,
14 it is common to see a higher number of lanes at the two extremities of a road segment to allow for
15 traffic insertion and for dedicated turning lanes. The complex nature of traffic lanes can be seen in
16 Figure 6 where a reserved bus lane (highlighted in green) is present at a fixed schedule and the
17 number of lanes at the intersection level is different (usually greater) than the mid-block number
18 of lanes to accommodate turning movement flows. This paper examines the mid-block number of
19 lanes and does not consider reserved lanes.

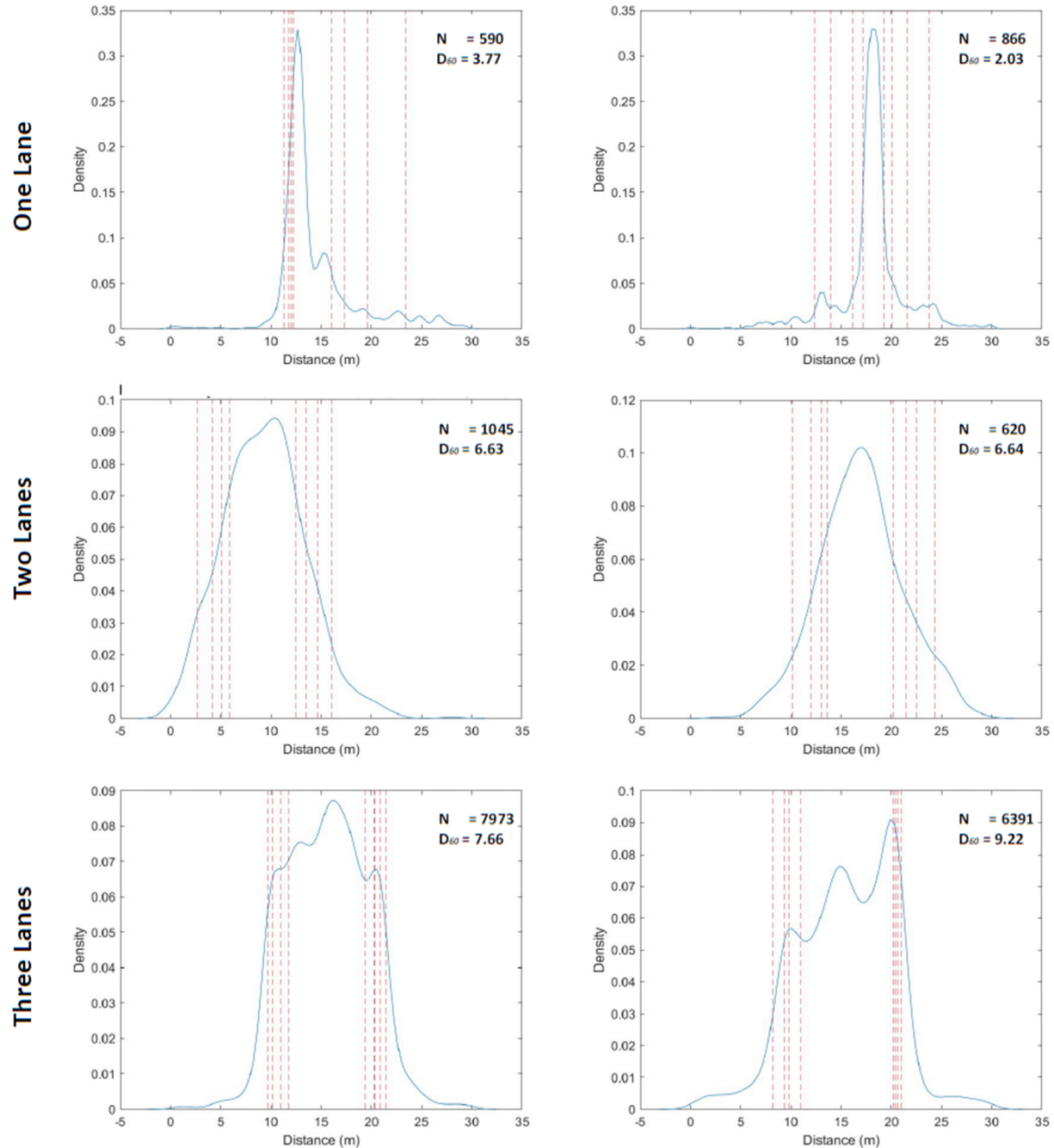


20
21 *Figure 6 - Example of a Complex Road Geometry*

23 Results

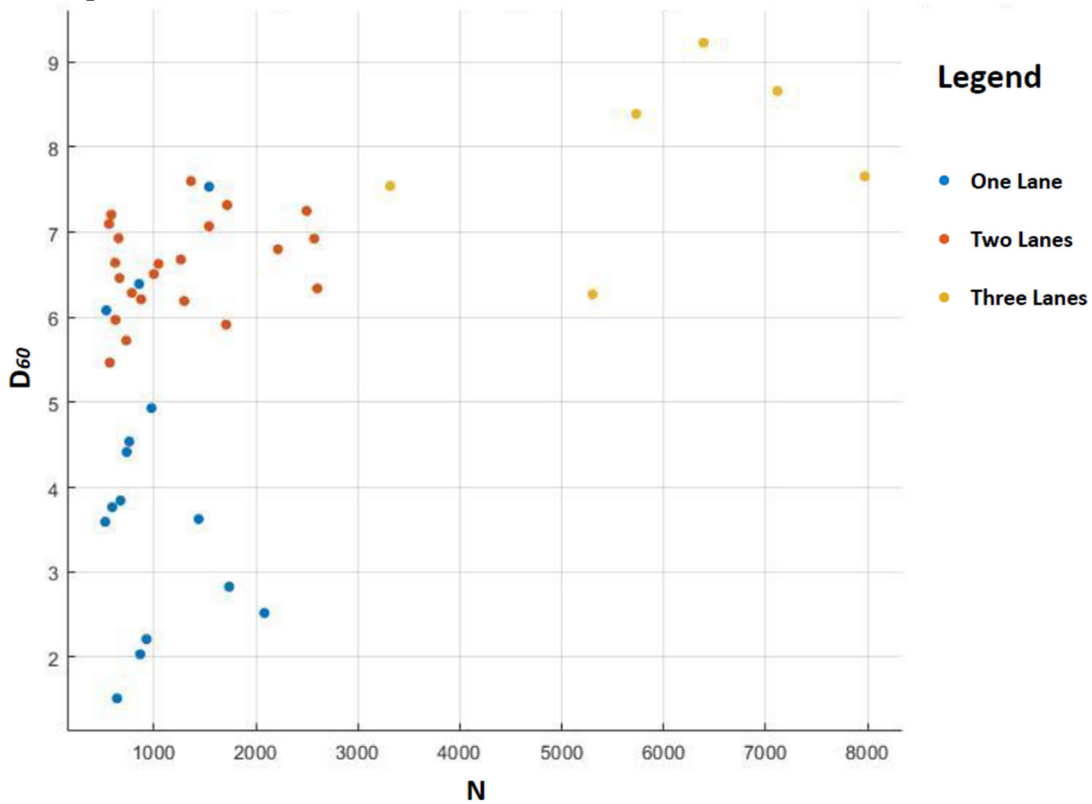
24 With the proposed steps and parameters, it was possible to extract GPS points for road segments
25 and associate each point to the correct directional link based on the trajectory direction. The sample
26 size filter limited the number of analyzed directional links included in the analysis to 43 links. The
27 buffer sizes were also validated based on the frequency distributions of GPS points' lateral distance
28 with respect to the link since the entire distribution profile is captured. This can also be noted in

1 Figure 7 which presents the kernel density estimator fitted to the lateral distance variable
 2 distribution for six different links of varying number of lanes. The sample size, N , and the D_{60}
 3 values are also presented for each link.



4
 5 *Figure 7 - Sample Kernel Density Estimator of Lateral Distance for One, Two, and Three Lanes*
 6 In addition to the distribution profile of lateral distance, the figure also demonstrates the significant
 7 difference in the distribution profile between links having one, two, and three lanes. Through
 8 observation, it was possible to identify that road segments with fewer lanes have lower values of
 9 N and smaller D_{60} values. This can be explained by the fact that roads with fewer users are designed
 10 to have fewer lanes, and GPS points are concentrated in a narrower area.

1 Model specification was carried out to determine the best model and best predictors for the number
2 of lanes. Given the relatively low number of road segments, a 5-fold cross-validation method was
3 performed to avoid overfitting the data. The highest classification prediction accuracy was found
4 using a decision tree classifier at 91% using two predictors, the sample size N and D_{60} . The
5 optimizable decision tree classifier tested iteratively different numbers of splits and different split
6 criteria to reach the minimum classification parameters and error.
7 Figure 8 Figure 9 presents a plot of the two selected predictors, showing a clear delimitation
8 between the predictor values for roads with one, two, or three lanes.



9
10 *Figure 8 - D60 vs. Sample Size (N)*

11 Moreover, the optimized decision tree is presented with the three split levels and values in Figure
12 9. Ensemble classifiers, such as boosted trees, bagged trees, and subspace discriminant were also
13 tested to improve prediction accuracy and the best accuracy was using the subspace discriminant
14 ensemble classifier at 91%. Given that the optimized decision tree was able to predict at the same
15 accuracy level it was selected as the best model in this case since it is simpler to visualize and
16 interpret.

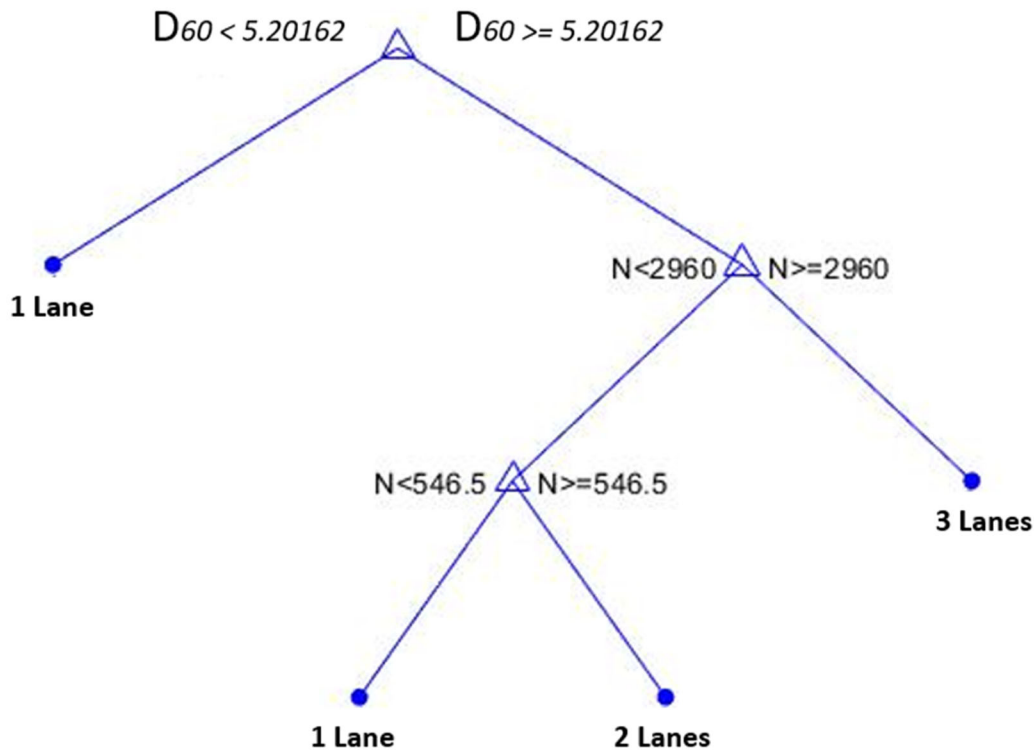


Figure 9 - Selected Classification Decision Tree

Discussion

The proposed methodology predicts the number of lanes per road segment based on the number of GPS points associated to the link and the difference between the 80th and 20th percentile distance, representing a lateral distance measure centered around the median lateral distance.

Given that the best prediction model was obtained using only two variables, an optimized decision tree classifier was sufficient to reach a good model accuracy (91%). However, adding new variables will require retesting ensemble classifier methods to verify if they are able to improve prediction accuracy. Moreover, to use this model, the sample size would need to be translated into relative terms or to be specified with respect to the sample size corresponding to a new dataset. The main hypothesis behind using the sample size as a variable is that for a given period of data collection where we assume a representative sample, it is expected to have a larger number of observations for road segments with a larger number of lanes since they generally have a higher traffic flow.

The spatial analysis and model specification steps were limited by the experimental data available. During the study, it was found that some of the GPS points were map matched in their raw form which signifies that they were snapped to a road centreline at a step prior to accessing the data. Given that this study examines the lateral distribution of raw GPS points with respect to the road link, map matching has a negative impact on data quality. It was also noted that some links had a low number of GPS points, which resulted in unstable lateral distance distribution profiles.

Ideally, larger datasets of uniquely raw GPS points need to be used to have a larger coverage to obtain more realistic distributions and potentially create more predictor variables. The objective is to have more GPS points per link, not necessarily more links as it will also become more complex

1 to obtain the ground truth information. An increase in the number of points per link will also
2 increase the probability of having better coverage for different times of the day, which enables
3 model specification for different time periods to detect the change in the number of traffic lanes
4 temporally.

5 Although some studies have proposed the extraction of the number of lanes using satellite and
6 street-level imagery with a relative high accuracy, they are not without limitations (Nieroda et al.,
7 2022, Mátyus et al., 2016, Kasmí et al., 2018). In fact, the high cost of imagery data collection is
8 an important limitation that is overcome in this study since GPS data is currently being crowd
9 sensed by location-based applications through smartphones. Furthermore this study considers 43
10 road links for model specification which is a larger sample than the work by Tang et al. (2014)
11 which only considers 6 road segments for the analysis.

12 Comparing this study to some studies using GPS data to extract the number of lanes, the prediction
13 accuracy significantly exceeds the 60% accuracy obtained in the study by Zhang et al. (2010). In
14 addition, the method proposed in the current study provides more accurate results and a simpler
15 procedure than the studies by Chen and Krumm (2010) and Arman and Tampere (2020) to obtain
16 the number of lanes for integration in large-scale transport models.

17 18 Conclusion

19 This study proposes a method to predict the number of lanes per road segment using crowd sensed
20 GPS trajectory data as an input in addition to a simple geographic representation of the road
21 network. The proposed framework is composed of two main steps: to predict the number of lanes
22 of road segments using GPS trajectory data while aiming to keep the cost low and to obtain high
23 prediction accuracy.

24 The first step is a spatial analysis process to filter and prepare the GPS trajectory data for variable
25 creation. Due to the noise inherent to GPS trajectory, it was crucial to ensure that raw GPS data
26 points were filtered using buffers. This is also necessary to account for the specificities in road
27 design and for the discrepancies in the road network geographic representation. This step also
28 served to produce variables necessary to derive the predictors for the following step. The two main
29 variables were the number of GPS points per road segment and the lateral distance between each
30 point and the reference line representing the road segment. The second step is the training and
31 validation of a machine learning method using classification tree analysis and ensemble learning.
32 Standardized predictors were derived from the lateral distance variables to ensure that the values
33 are comparable across different road segments.

34 This study was able to develop a road segment number of lanes prediction model using GPS
35 trajectory point data with an accuracy of 91% using a decision tree classifier and two predictors.
36 This prediction accuracy is higher than prediction results obtained by previous research. This
37 finding demonstrates that it is possible to extract the number of lanes available for general traffic
38 by using crowd-sensed GPS trajectory data. This will facilitate road transport network model
39 development and update. The proposed method was demonstrated using a case study in Quebec
40 City, Canada.

41 However, the work is not without limitations and can be further developed by having a larger
42 temporal sample coverage to enable the prediction of the number of lanes for different periods
43 allowing the detection of dynamic reserved lanes or parking lanes. This study used manually
44 collected ground truth data which limited the size of the study area, network coverage for model
45 development and validation will be increased in future works by collecting more ground truth data
46 or obtaining this information from another source. Moreover, it is possible to explore adding land

1 use variables that might be correlated with the number of lanes and help in improving the
2 prediction model's accuracy. The potential of this method can also be maximized by automating a
3 procedure that can use GPS trajectory points and other basic input files to create a road network
4 containing the number of lanes per road segment.
5 Eventually, with the arrival of autonomous vehicles, new data sources may also be available in
6 terms of geotagged imagery data that can be automatically collected and treated by these vehicles
7 during their operation. These processed images may in the future be used to mine road network
8 features at a low cost and high accuracy.

1 **Acknowledgement**

2 The authors would like to acknowledge the generous support of McGill University's Faculty of
3 Engineering and the Vadasz Scholars Program.

4

5 **Author Contribution**

6 The authors confirm contribution to the paper as follows:

7 Study conceptualization and design: All authors.

8 Formal Analysis, Investigation, Writing - Original Draft: Adham Badran

9 Supervision: Ahmed El-Geneidy and Luis Miranda-Moreno

10 Interpretation of results and manuscript review and editing: All Authors

References

- 1
2 ARMAN, M. A. & TAMPERE, C. M. J. 2020. Road centreline and lane reconstruction from pervasive GPS
3 tracking on motorways. *Procedia Computer Science*, 170, 8.
- 4 BOUNINI, F., GINGRAS, D., LAPOINTE, V. & POLLART, H. Autonomous Vehicle and Real Time Road
5 Lanes Detection and Tracking. 2015 IEEE Vehicle Power and Propulsion Conference (VPPC), 19-
6 22 Oct. 2015 2015. 1-6.
- 7 CHEN, Y. & KRUMM, J. Probabilistic modeling of traffic lanes from GPS traces. Proceedings of the 18th
8 SIGSPATIAL international conference on advances in geographic information systems, 2010. 81-
9 88.
- 10 GUO, C., KIDONO, K., MEGURO, J., KOJIMA, Y., OGAWA, M. & NAITO, T. 2016. A Low-Cost
11 Solution for Automatic Lane-Level Map Generation Using Conventional In-Car Sensors. *IEEE*
12 *Transactions on Intelligent Transportation Systems*, 17, 2355-2366.
- 13 GUO, Y., LI, B., LU, Z. & ZHOU, J. 2021. A novel method for road network mining from floating car data.
14 *Geo-spatial Information Science*, 16.
- 15 IAN, H., FRANK, E., HALL, M. & CHRISTOPHER, J. 2017. Data mining: Practical machine learning
16 tools and techniques—Part II: More advanced machine learning schemes. Morgan Kaufmann,
17 Burlington, MA.
- 18 KADALI, B. & VEDAGIRI, P. 2020. Role of number of traffic lanes on pedestrian gap acceptance and risk
19 taking behaviour at uncontrolled crosswalk locations. *Journal of Transport & Health*, 19, 100950.
- 20 KASMI, A., DENIS, D., AUFRERE, R. & CHAPUIS, R. Map Matching and Lanes Number Estimation
21 with Openstreetmap. 2018 21st International Conference on Intelligent Transportation Systems
22 (ITSC), 4-7 Nov. 2018 2018. 2659-2664.
- 23 LEICHTER, A. & WERNER, M. 2019. Estimating road segments using natural point correspondences of
24 GPS trajectories. *Applied Sciences-Basel*, 9, 11.
- 25 MÁTTYUS, G., LUO, W. & URTASUN, R. Deeproadmapper: Extracting road topology from aerial
26 images. Proceedings of the IEEE international conference on computer vision, 2017. 3438-3446.
- 27 MÁTTYUS, G., WANG, S., FIDLER, S. & URTASUN, R. Enhancing road maps by parsing aerial images
28 around the world. Proceedings of the IEEE international conference on computer vision, 2015.
29 1689-1697.
- 30 MÁTTYUS, G., WANG, S., FIDLER, S. & URTASUN, R. Hd maps: Fine-grained road segmentation by
31 parsing ground and aerial images. Proceedings of the IEEE Conference on Computer Vision and
32 Pattern Recognition, 2016. 3611-3619.
- 33 MERRY, K. & BETTINGER, P. 2019. Smartphone GPS accuracy study in an urban environment. *PLOS*
34 *ONE*, 14, e0219890.
- 35 MONTRÉAL, V. D. MTL Trajet Study.
- 36 NIERODA, B., WOJAKOWSKI, T., SKRUCH, P. & SZELEST, M. A Heatmap-Based Approach for
37 Analyzing Traffic Sign Recognition and Lane Detection Algorithms. 2022 26th International
38 Conference on Methods and Models in Automation and Robotics (MMAR), 22-25 Aug. 2022 2022.
39 217-221.
- 40 TANG, L., GAN, A. & ALLURI, P. 2014. Automatic Extraction of Number of Lanes from Georectified
41 Aerial Images. *Transportation Research Record*, 2460, 86-96.
- 42 TREIBER, M. & KESTING, A. 2013. Traffic flow dynamics. *Traffic Flow Dynamics: Data, Models and*
43 *Simulation*, Springer-Verlag Berlin Heidelberg, 983-1000.
- 44 ZHANG, C., XIANG, L., LI, S. & WANG, D. 2019. An intersection-first approach for road network
45 generation from crowd-sourced vehicle trajectories. *ISPRS International Journal of Geo-*
46 *Information*, 8, 26.
- 47 ZHANG, L., THIEMANN, F. & SESTER, M. Integration of GPS traces with road map. Proceedings of the
48 Third International Workshop on Computational Transportation Science, 2010. 17-22.
- 49 ZHONGYI, N., LIJUN, X., TIAN, X., BINHUA, S. & YAO, Z. 2018. Incremental road network generation
50 based on vehicle trajectories. *ISPRS International Journal of Geo-Information*, 7, 19.