

1 **CLASSIFYING BEHAVIORAL DYNAMICS OF TAXI DRIVERS ROUTE CHOICES**
2 **USING LONGITUDINAL GPS DATA**

3
4

5 **Hamzeh Alizadeh**

6 Department of Civil, Geological and Mining Engineering, Polytechnique Montréal,
7 C.P. 6079, succ. Centre-Ville Montréal (Québec) Canada, H3C 3A7
8 Tel: (514) 340-4711 ext. 4235 Fax: 5143404191 E-mail: hamzeh.alizadeh-aliabadi@polymtl.ca

9

10 **Bilal Farooq**

11 Head of Laboratory of Innovations in Transportation (LITrans),
12 Department of Civil, Geological and Mining Engineering, Polytechnique Montréal,
13 C.P. 6079, succ. Centre-Ville Montréal (Québec) Canada, H3C 3A7
14 Tel: (514) 340-4711 ext. 4802, Fax: (514) 340-3981, E-mail: bilal.farooq@polymtl.ca

15

16 **Catherine Morency**

17 Head of Mobilité Research Chair
18 Department of Civil, Geological and Mining Engineering, Polytechnique Montréal,
19 C.P. 6079, succ. Centre-Ville Montréal (Québec) Canada, H3C 3A7
20 Tel: (514) 340-4711 ext. 4502 Fax: (514) 340-3981 E-mail: cmorency@polymtl.ca

21

22 **Nicolas Saunier**

23 Department of Civil, Geological and Mining Engineering, Polytechnique Montréal,
24 C.P. 6079, succ. Centre-Ville Montréal (Québec) Canada, H3C 3A7
25 Tel: +1 (514) 340-4711 ext. 4962, Fax: (514) 340-3981, E-mail: nicolas.saunier@polymtl.ca

26

27

28 Word count: 6,464 words text + 6 tables/figures x 250 words (each) = 7,946 words

29

30

31

32

33

34 11-November-2016

1 ABSTRACT

2 This study aims to capture the behavioral heterogeneity in route choice by identifying subgroups
3 of drivers based on their actual route choices and factors affecting them. We have studied a highly
4 longitudinal GPS dataset, tracking 1,746 taxi drivers over a period of one year, making more than
5 22,000 trips between the Islands of Montreal and Laval. We opted for a two-step procedure, where
6 in the first step a Principal Component Analysis (PCA) is performed to reduce collinearity among
7 attributes, followed by a Hierarchical Agglomerative Clustering (HAC) to form behavioral clusters
8 in the second step. Results show that four major types of route choice behaviors are observable
9 among taxi drivers. These clusters show significant variations based on the time of day (day/night)
10 and the traveled distance (shorter trips/longer trips) and are labelled: “*Short trips night drivers*”,
11 “*Long trips night drivers*”, “*Short trips day drivers*”, and “*Long trips day drivers*”. Due to the
12 rise of ride-hailing services, the understanding of these patterns are important for city and
13 transportation planners in the context of proposing new laws and policies that safeguard taxi
14 industry as well as encourage sharing economy. The inclusion of similar typologies in route choice
15 models would improve their behavioral aspect as well as their estimation and prediction abilities.

16

17

18 *Keywords:* Route choice behavior, Principal Component Analysis (PCA), Hierarchical
19 Agglomerative Clustering (HAC), Longitudinal GPS data, Taxi

1 INTRODUCTION

2 Individuals have different inclinations toward choosing a route between same Origin-Destination
3 (OD) pairs. These choices form the collective patterns of route choice behaviors, observable in our
4 everyday lives. Understanding these patterns is of crucial importance to the city and transportation
5 planners (1). Route choice models are used to capture the complex process of drivers' route
6 selection behavior. These models are used to estimate and predict the probability of a certain route
7 being chosen between a given OD pair. The complexity of this process is mostly due to factors
8 such as the sophisticated nature of human behavior; ambiguity of decision making process;
9 stochasticity of individuals' preferences, as well as the high density of road networks; large number
10 of possible alternatives between OD pairs; and correlation among these alternatives.

11 Route choices are highly dependent on individuals' characteristics. Therefore, in order to
12 improve the behavioral aspect of route choice models, it is essential to incorporate various sources
13 of heterogeneity influencing drivers' route choice decisions. This heterogeneity mostly comes
14 from having different preferences, experiences, information levels, and attitudes. These sources of
15 heterogeneity are somehow correlated, which makes the understanding of route choice behavior
16 even more complex. Therefore, it might be interesting and necessary to ask "*How factors affecting
17 route choices are correlated?*", and "*How these factors can be grouped into various categories?*".
18 Another relevant question could be "*How do these different groups represent different categories
19 of drivers?*".

20 Unfortunately, a detailed representative classification of drivers according to their route
21 choice decisions is missing from the existing literature. According to (2), dividing drivers into
22 different categories based on their route choice behavior can improve route choice models and can
23 be used in defining different functional forms for traffic assignment models. (3) classified drivers
24 according to their age, gender, ethnicity, education, driving experience, and annual driven miles,
25 and found meaningful differences in their route choice behaviors. They have argued that the
26 incorporation of different groups of drivers can improve the accuracy of route choice models. (4)
27 illustrated that the incorporation of an activity-based segmentation can improve traffic assignment
28 procedures and proposed stratifying traffic assignment models by incorporating different
29 typologies of attitudes toward route choice.

30 A further challenge in route choice modeling is the selection of a proper data collection
31 method. Although travel surveys are one of the most adopted methods (5-8), it is not easy to
32 capture route choice behavior via diary surveys. With the introduction of data collection through
33 Global Positioning Systems (GPS) in transportation studies, researchers are supplied with
34 unprecedented high-resolution spatial and temporal records on drivers' route choices.

35 Earlier route choice studies have mostly focused on traditional cross-sectional analysis of
36 drivers' behavior. Therefore, the decision making behavior for a single day has usually been
37 interpreted to derive important attributes affecting drivers' decision making patterns. Studies
38 focusing on driving patterns over an extended period appear to be minimal, which might be mostly
39 due to the challenges of data collection and analysis. In recent years, with the dissemination of
40 GPS technologies and the progressive utilization of in-vehicle or on-person GPS devices,
41 researchers are supplied with longitudinal choice data over an extended period, which provides
42 better insights on drivers' route choice patterns and heterogeneities than cross sectional data.

43 In this study we aim to extend the understanding of drivers' route choice behaviors by
44 classifying drivers into separate categories, based on their observed route choice decisions, using

1 a longitudinal GPS dataset. Although previous studies have shown that some attributes such as
2 travel time, travel speed, some demographic and socioeconomic attributes, etc., influence drivers'
3 route choices, to the best knowledge of the authors, none of these studies have ever categorized
4 common attributes resulting in the same decision, to identify various classes of route choice
5 behavior. This research relies on a longitudinal GPS dataset, tracking 1,746 taxi drivers making
6 more than 22,000 trips over a period of one year. We focus on trips between the Islands of Montreal
7 and Laval, originating in Montreal.

8 This work distinguishes itself from the existing literature based on two original
9 contributions. First, to the extent of our knowledge, drivers' route choice behaviors have never
10 been stratified using drivers' actual route choices. Second, route choice preferences have never
11 been looked at through longitudinal GPS data over a one-year period. Previous datasets were not
12 large and/or detailed enough to allow this type of stratification.

13 The rest of this paper is organized as follows. First, we overview factors affecting drivers'
14 route choice decisions. Second, we present the studied regional context and describe our dataset.
15 Then, we discuss the process of deriving indicators and our clustering methodology. Finally, we
16 present descriptive analysis of studied trips, discuss the clustering results and underline the
17 findings of this study.

18

19 **STATE-OF-THE-ART**

20 The aim of this study is to classify drivers' route choice behaviors based on their observed choices
21 over a long duration of time. In a perfect world, drivers make route choices by minimizing a certain
22 perceived cost. This cost function is unique for each driver and situation, resulting in different
23 chosen routes, and can be attributed to several factors, such as different levels of information,
24 different capacities to process them, and different computational and prediction abilities (9). A
25 large body of literature exists on modelling route choice and incorporating factors affecting them.
26 These factors are correlated and form different types of preferences, which in turn can be used to
27 classify drivers into different categories.

28 Classifications of driving behaviors have been mostly cited in traffic safety studies to
29 distinguish between conflicts and collisions, risk-taking versus safe drivers and experienced versus
30 novice drivers (10-13), as well as car-following research classifying drivers according to their
31 driving behavior (14). A classification of driving patterns has been proposed by (15), which divides
32 drivers into three categories, namely "*passionate*", "*every day*", and "*leisure time*" drivers, based
33 on their attitudes toward the environment. Behavioral classification studies have also been
34 conducted on bike users. (16) studied the causality relation between non-work-related trips and
35 commuters, and classified cyclist into four different clusters, namely *Non-cyclists*, *Non-work*
36 *cyclists*, *All-around cyclists*, and *Commuter cyclists*. Another major work on cyclists classification
37 is the work by (17), in which bike users are classified into four categories: *Strong and Fearless*,
38 *Enthusied and Confident*, *Interested but concerned*, and *No Way No How*. The reader is referred to
39 (18) for more reflections on this work. Another recent work by (19) classified bike-sharing users
40 into two groups of *Commuters* and *Recreational users*. Drivers' perceptions and experiences have
41 been studied by (20), using a driving simulator and two initial and final questionnaires. They found
42 observable differences between drivers' route choice behaviors and categorized drivers based on
43 their learning skills into four different categories. In a further study, (21) have demonstrated that
44 the inclusion of those learning clusters would improve disaggregated route choice models. In an

1 interesting study by (22), cab drivers have been classified into *Top Drivers* and *Ordinary Drivers*
2 based on their incomes, and their trip characteristics have been studied.

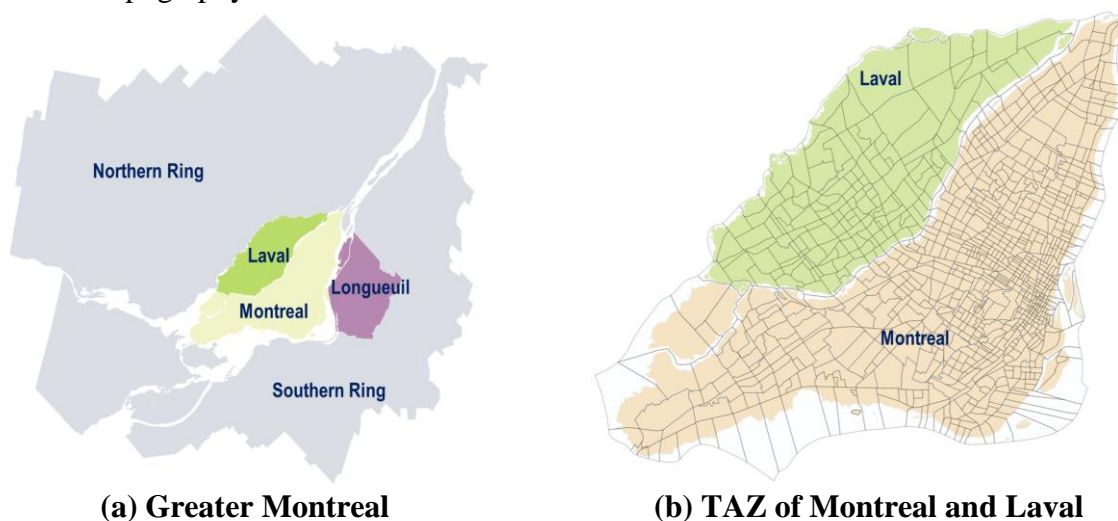
3 Although previous studies showed that different categories of road users are observable,
4 an explicit classification of drivers based on their actual choices is missing from the literature. Due
5 to the availability of very large and rich datasets from mobility services (e.g. taxi companies), that
6 have detailed GPS trajectories of their fleets for a long duration of time, it has become possible to
7 address this critical research need.

9 CONTEXT AND DATA DESCRIPTION

10 Regional Context

11 This study is performed in the context of the metropolitan region of Greater Montreal, depicted in
12 Figure 1(a). This Island city is separated from its suburbs by two large rivers; Prairies River and
13 Saint Lawrence River in the north and south of the city, respectively. In this study we focus on
14 trips originating in Montreal with a destination in Laval, the largest suburb of Montreal located
15 north of the city, across the Prairies River. These two islands contain roughly 2.3 million
16 inhabitants, and cover a total surface of 632.3 km² (23).

17 Traffic Analysis Zones (TAZ), delineated by Quebec's Ministry of Transport, are used in
18 this study to determine the characteristics of taxi trips' origin and destination points (Figure 1(b)).
19 TAZs are geographical areas, which divide the city into smaller similar areas based on various
20 factors such as population, demography, socioeconomic information, road network, transit access,
21 land use and topography.



22

23 **FIGURE 1 Regional context of the study.**

24

25 GPS Dataset

26 We investigated GPS traces collected over a period of one year (2015) by a major Taxi company
27 operating in Montreal. This taxi company constitutes around 25% of Montreal's taxi fleet, and its
28 operation is restricted to trips starting or ending in the central part of the island. Every taxi is
29 equipped by a data logger and GPS data are collected continuously for operational purposes. For
30 our study we extracted a subset of the main dataset, consisting of around 750,000 GPS records
31 collected from 1,746 taxi drivers making a total of 22,394 trips. Drivers are associated with a

1 unique ID so that we can distinguish between trips made by a same car and different drivers. We
2 should mention that personal information on drivers' demographic and socioeconomic
3 characteristics are not available due to privacy issues.

4 In order to explore factors affecting drivers' route choices, we first need to derive the
5 observed route from the recorded GPS points. Therefore, every record has to be associated with a
6 link in the network, a process which is called map-matching. This step is crucial, since it
7 determines the accuracy of reconstructed trajectories, and accordingly their derived attributes.

8 Data has been stored in a PostgreSQL database, and the PostGIS spatial extension has been
9 added to support geographical datatypes and queries. To associate GPS records to the road
10 network, a direction-based nearest link point-to-curve map matching algorithm has been adopted,
11 in which every record is matched onto the closest link in the network with respect to its azimuth.
12 To reconstruct the complete trajectory, consecutive GPS points have been connected using a
13 distance-based shortest-path algorithm.

14 **Network Dataset**

15 The road network has been extracted from the OpenStreetMap project in the format of
16 geographical layers ("shapefiles"). It contains more than 156,000 nodes and 89,000 links. The
17 network has been made routable, through a geospatial extension of PostgreSQL named
18 "pgrouting", in order to enable the calculation of shortest paths. It has also been segmented on an
19 intersection-to-intersection basis; therefore, links are defined to be road segments between two
20 consecutive intersections.

21 **METHODS**

22 **Process of Deriving Indicators**

23 The first step toward classifying drivers' route choice behaviors is to explore a large range of trips'
24 (choices) characteristics and factors affecting these choices. These factors include temporal and
25 environmental attributes, drivers' attitudes and preferences, network familiarity, personal
26 demographic and socio-economic characteristics, and route related attributes among others. Since
27 personal characteristics are not observable through GPS traces, the main focus of this section is to
28 derive explanatory factors from the GPS dataset to thoroughly describe observed trajectories.
29 Factors are classified into five broad categories, namely *Temporal Indicators*, *Route*
30 *Specifications*, *Driver Characteristics*, *Land Use*, and *Route Similarities*. A concise description of
31 derived attributes follows.

32 *Temporal Indicators*

33 The timestamp information of the origin point of a trip, containing information regarding the date
34 and time of the record, has been used to derive the following indicators:

- 35 • Peak/Off-peak: Trips starting between 6:00 to 9:00, and 16:00 to 19:00 are considered as
36 peak hour trips, while all other trips are regarded as off-peak trips.
- 37 • Weekday/Weekends
- 38 • Day of the Week
- 39 • Day/Night: Day trips consist of trips starting between 6:00 to 21:00 and the rest are
40 considered as night trips.
- 41 • Season

42 *Route specifications*

1 These attributes characterize observed routes between given OD pairs, based on their physical
2 specifications and their similarity levels with their respective distance-based shortest paths.

- 3 • Length: The total traveled distance.
- 4 • Travel time: The timestamp information at origins and destinations have been used to
5 extract the observed travel time.
- 6 • % Highway: Specifies the percentage length of trips made on highways.
- 7 • # Turns / km: The total number of turns per kilometer.
- 8 • # Segments / km: The total number of road segments per kilometer.
- 9 • % Shortest path: The percentage of overlap with the distance-based shortest path.

10 *Driver characteristics*

11 The total number of trips per driver and trips made by the same driver between the same TAZ pairs
12 have been inspected and several attributes, pertinent to the regularity of drivers between these TAZ
13 pairs, have been derived:

- 14 • Total trips: The total number of trips made by each driver.
- 15 • Average number of trip per TAZ pair: The total number of trips is divided by the total
16 number of TAZ pairs, between which a driver has traveled. The lower bound of this
17 attribute is one, indicating only one travel per TAZ for a driver. Higher values signify
18 higher propensity of traveling between the same TAZ pairs.
- 19 • TAZ pairs with more than 4 trips: The total number of TAZ pairs between which the driver
20 has made more than four trips, which has been selected based on the distribution of trips
21 between same TAZ pairs, as an indicator of regularity of drivers.
- 22 • Maximum same TAZ trips: The maximum number of trips made between a given TAZ
23 pair by the same driver. This factor is also interpreted as an indicator of regular trips
24 between same TAZ pairs.

25 *Land use*

26 Although GPS data allows to locate the exact pick up and drop off position of travellers, we have
27 no information regarding their exact trip purposes. In this work we associated GPS points to their
28 respective TAZ to obtain information regarding the land use of trips' origin and destination points.
29 Land use data was prepared by local government agencies. Four categories of land uses are
30 considered in this study, namely residential, commercial, work/study, and recreational.
31 Accordingly, the following attributes were derived:

- 32 • Origin land use
- 33 • Destination land use

34 *Route Similarities:*

35 To evaluate the degree of similarity for trips taken by a particular driver between a given pair of
36 TAZ, we have used a measure called path-size (PS). This measure has been proposed by (24) to
37 account for similarities between routes in logit based discrete choice models. The following
38 formulation has been adopted in this study:

$$39 \quad PS_{in} = \sum_{a \in \Gamma_i} \frac{L_a}{L_i} \frac{1}{\sum_{j \in \varphi_n} \delta_{aj}} \quad (1)$$

where PS_{in} denotes the path-size factor for driver n and route i , L_a and L_i represent the length of link a and route i , Γ_i is the set of road segments in route i , φ_n denotes the observed routes for driver n between the same pair of TAZ, $\sum_{j \in \varphi_n} \delta_{aj}$ indicates the total number of alternatives in φ_n sharing link a (δ_{aj} is the link-path incident binary variable which is 1 if link a is on route i , and 0 otherwise). The upper bound value of this formulation is 1, indicating that observed routes are completely independent and do not share any links. However, smaller values of PS indicate longer overlaps and higher dependencies between trips. The following indicator has been derived:

- **PS:** The average PS factor calculated for all trips made by the same driver, between all TAZ pairs.

Cluster Analysis

Clustering is an unsupervised categorization technique that aims to segregate multivariate data sets into more meaningful clusters, according to their main describing attributes. In transportation studies datasets usually have very high dimensions. Multivariate analysis techniques and dimension reduction techniques, such as Principal Component Analysis (PCA), are adopted to better interpret the data and to improve the quality of cluster analysis. PCA is an unsupervised dimension reduction technique, which preserve most of the initial information within a smaller numbers of mutually uncorrelated factors.

We adopt a two-step procedure in order to classify driver behaviors. First, in order to reduce collinearity among attributes, a PCA for mixed data is performed (25). The number of principal components extracted from an analysis depends on the level of correlation between attributes and the amount of variance that can be explained by each principal component. The total number of components is equal to the number of attributes, but only the first few ones are important and are considered in the next step (26). Then, a Hierarchical Agglomerative Clustering (HAC) using Ward's criterion is used to classify driver behaviors. In HAC method, clusters are formed hierarchically by merging the two closest clusters at each step. All PCA and clustering analysis are performed in the TANAGRA statistical package (27).

To assess the significance level of attributes in the clustering result, the "Test Value" (TV) criterion has been used. Following formulas are used to calculate TV for continuous (t_c), and discrete (t_d) values for each cluster (28):

$$t_c = \frac{\mu_g - \mu}{\sqrt{\frac{n - n_g}{n - 1} \times \frac{\sigma^2}{n_g}}} \quad (2)$$

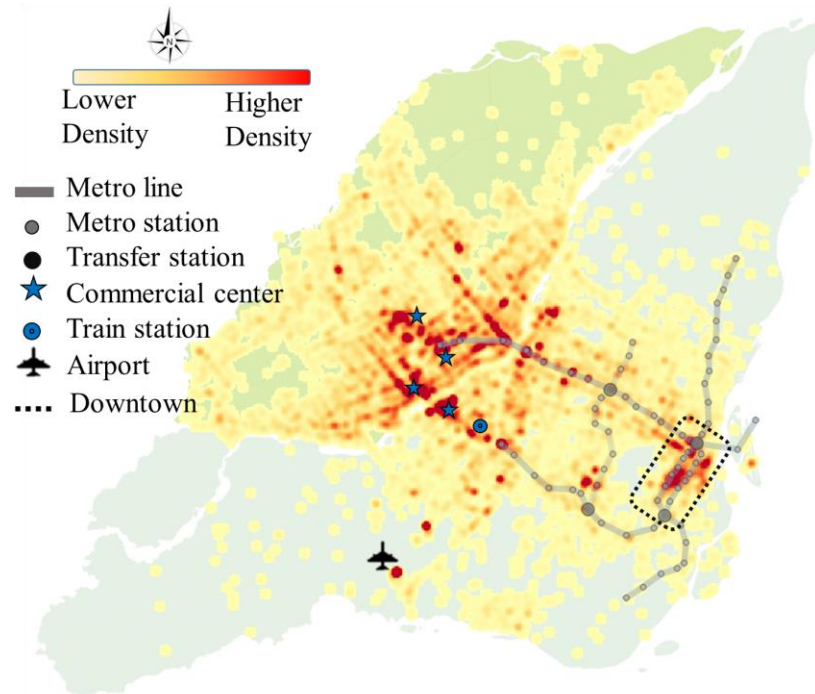
$$t_d = \frac{n_{jg} - \frac{n_g \times n_j}{n}}{\sqrt{\frac{n - n_g}{n - 1} \times \left(1 - \frac{n_j}{n}\right) \times \frac{n_g \times n_j}{n}}} \quad (3)$$

where μ and μ_g are attributes' means in the cluster and group, respectively; n and n_g denote the size of the cluster and the group, respectively; σ^2 represents the attribute variance in the cluster; and n_{jg} is the number of observations corresponding to the discrete attribute j in cluster g .

1 RESULTS

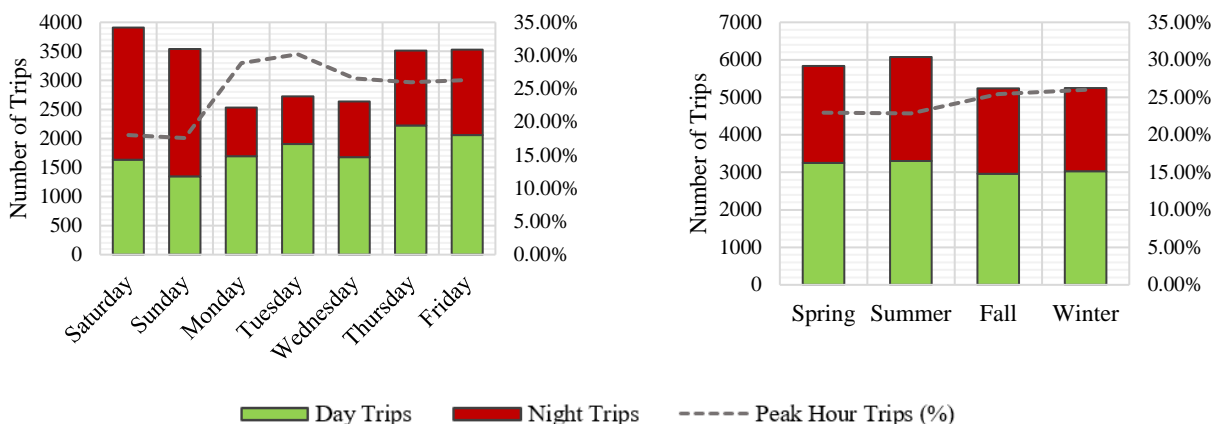
2 Trip Characteristics

3 The spatial distribution of origin and destination points is used to visualize the dispersion of taxi
 4 demand for trips between Montreal and Laval. As depicted in the heat map of Figure 2, taxi trips
 5 mostly originate from downtown Montreal, the airport, major commercial centers, and around train
 6 stations. Destination regions with high travel densities are more dispersed in Laval, which is
 7 probably due to the higher dispersion of commercial centers and overall lower density. Also,
 8 considering that Laval is a suburb of Montreal, the higher segregation between its dense residential
 9 areas contributes to the dispersion of high density destination points.
 10



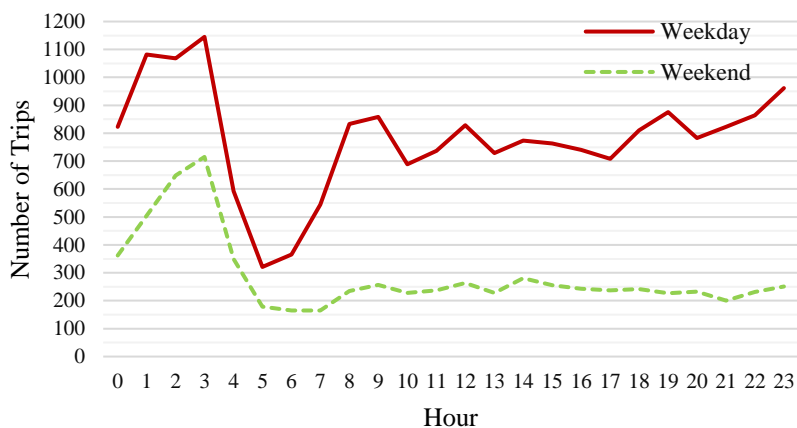
11
 12 **FIGURE 2 Heat map of origin-destination points for trips between Montreal and Laval.**

13
 14 A detailed comparison of trip frequencies across different days, the proportion of day trips
 15 versus night trips, and the percentage of trips made in peak hours are presented in Figure 3(a). A
 16 quick look reveals that taxi trips are more frequent during weekends, night trips are more common
 17 on Fridays and weekends, and peak hour trips are more frequent during weekdays. The same
 18 comparison is illustrated for trips made in different seasons (Figure 3(b)). Slightly higher numbers
 19 of trips are observed during spring and summer. However, the overall percentage of peak hour
 20 trips remains around 24 % throughout the whole year. Overall, around 56 % of trips are day trips
 21 and the remaining 44 % are made during nights. A closer look at the hourly distribution reveals
 22 that the peak demand occurs around 3 AM, drops significantly early in the morning, augments and
 23 stays steady over the day, and starts to increase again by the end of the day, around 9 PM (Figure
 24 3(c)). The same demand pattern holds for both weekdays and weekends, although the overall
 25 weekday demand is much higher.
 26



(a) Daily variations

(b) Seasonal variations



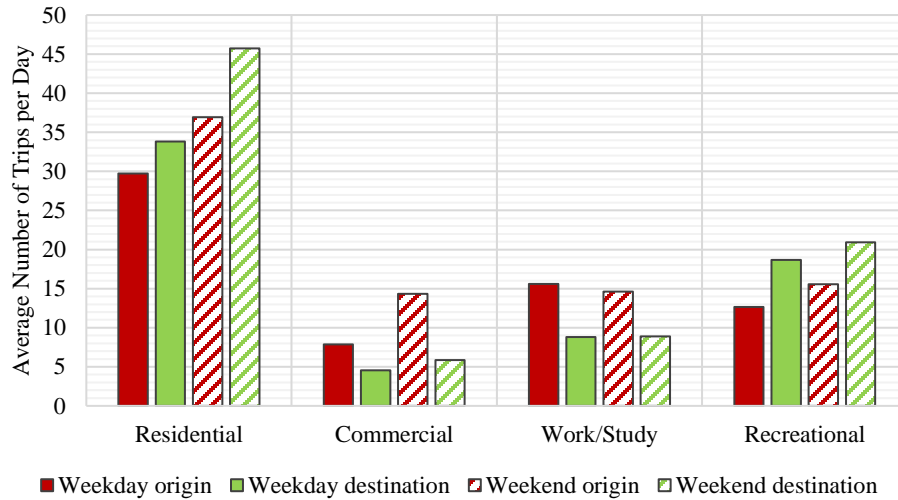
(c) Hourly variations

FIGURE 3 Daily, seasonal and hourly variations of taxi trips.

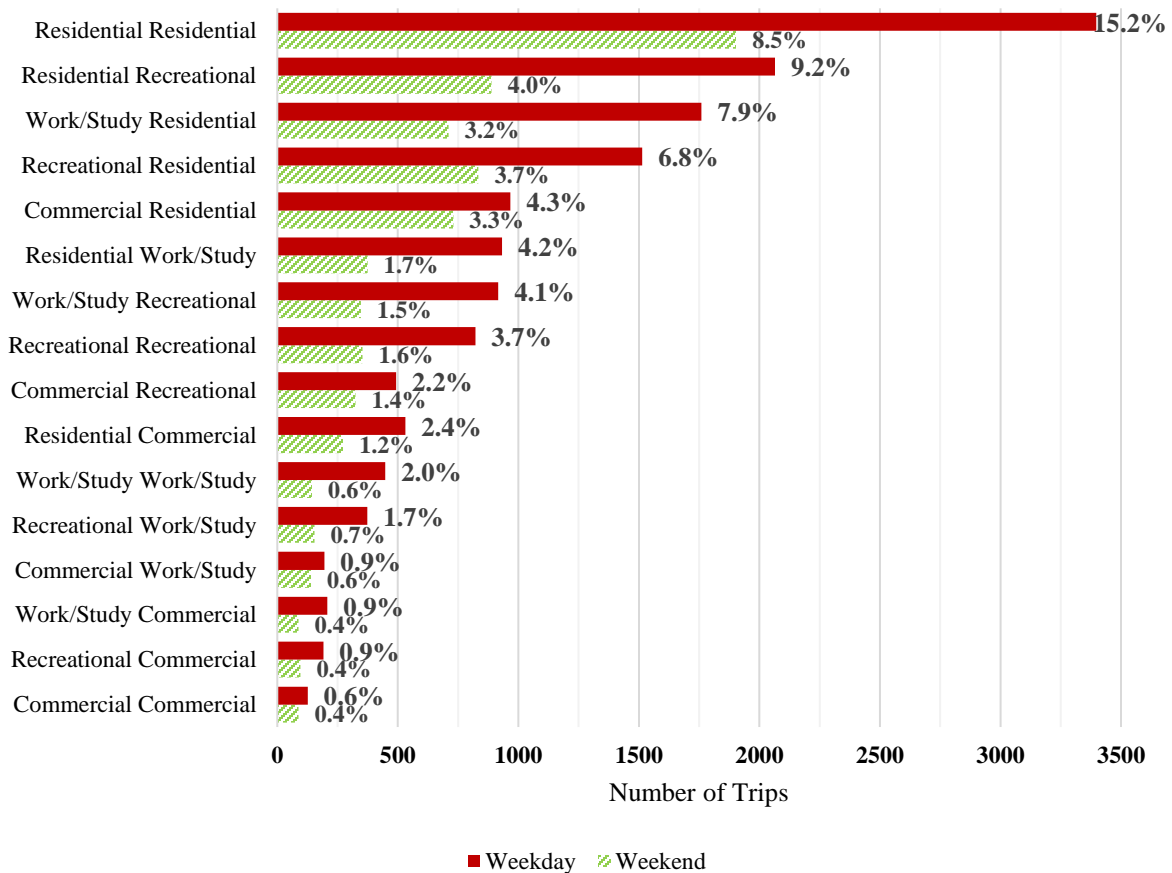
The total traveled distance is around 321,240 km. Around 20 % of taxi trips are shorter than 5 km; another 20 % have a length of 5 to 10 km; 35 % consist of trips from 10 to 20 km long; and the remaining 25 % are longer than 20 km. The mean, median and standard deviation of trip lengths are 14.3, 11.8, and 10.2 km, respectively. Although these values do not show any significant differences between trips made on weekdays compared to weekends, night trips and off-peak trips are slightly longer than day trips and peak hour trips, respectively. Another notable trend is that highway usage is considerably higher for night trips and off-peak trips, which is expected given lower congestion during these periods.

Figure 4(a) illustrates the average number of trips per day, classified based on the land use specification of the TAZ where the passenger was picked-up and dropped-off. Most of the trips originate from and/or have a destination in residential areas. It also highlights the high number of residential-end taxi trips for trips toward Laval. Trips are grouped based on OD land uses and the percentage of trips in each group is shown in Figure 4(b). Number of trips are represented by the length of the bar and labels represent the total share in percentages. A great majority of trips (around 46 %) were residential-based, meaning that trips start from a residential location, while around 20 % start from work/study regions.

1



(a) Average number of trips per weekday/weekend per land use



(b) Number and percentage of trips based on OD land uses

2

3

FIGURE 4 Statistics on trips' land uses.

4

1 **Driver Characteristics:**

2 The dataset used in this study includes a subset of 1,746 drivers comprising 204 night drivers and
3 498 day drivers, who only operate during nights and days, respectively. A total of 68 drivers only
4 operate during off peak hours, while 551 drivers work exclusively during peak hours. Around 38
5 % of peak hour trips take place during morning peak hours and the remaining 62 % occur during
6 evening peak hours. An average of 13 trips and a median of 7 trips per year (between the islands
7 of Montreal and Laval) are recorded for each driver, where the maximum number is 286.

8 Taking a closer look at trips between same TAZ pairs reveals that certain drivers are more
9 frequent among certain OD pairs, while others have a higher diversity of traveled TAZ pairs. More
10 than 77 % of drivers make a single trip, and only around 5 % of them make four or more trips
11 between a given TAZ pair. Among these more frequent drivers, 52 drivers make more than four
12 trips between more than one pair of TAZ. The maximum number of trips between the same pair
13 of TAZ by the same driver is found to be 40 trips.

14 **Cluster Characteristics**

15 A major challenge in clustering analysis is the selection of a good number of clusters. Although
16 visualization may be an effective way to verify results, it is highly difficult to visualize data with
17 more than three dimensions. In this study, a series of two to eight clusters have been experimented
18 and the optimal number of clusters has been defined by the maximum value of Between-group
19 Sum of Squares (BSS ratio) and GAP-statistic (29), as well as the behavioral interpretation that
20 can be associated with the clusters. Furthermore, the dendrogram representing the HAC process,
21 has been inspected to verify the plausibility of results. From this exercise a set of four clusters has
22 been found to provide the best results. These clusters are based on the first two principal
23 components found in the first step (PCA analysis), which explained more than 60 % of the
24 variation between attributes. Clustering has also been performed on more principal components,
25 comprising more than 70 % and 80 % of variations; however, results were not found to be stable,
26 meaningful, and reliable.

27 Two extra steps were undertaken to measure the results' stability: 1) the dataset was
28 randomly divided in two halves and the clustering process was separately performed on each half
29 for the same parameter settings; 2) observations were randomly permuted in our dataset. Since
30 results were not significantly different in both cases, it was concluded that the four-cluster solution
31 has a high degree of stability and reliability (30).

32 Results showed significant variations of drivers' behavior toward shorter versus longer
33 routes, and routes taken during day versus night. However, the exact hour of the trip, variation of
34 trip days, seasonal variations, and different OD land uses showed no significant impacts on drivers'
35 route choice behaviors in these clusters. Some descriptive statistics of factors, based on which the
36 final clustering has been made, are illustrated in Table 1(a). Correlations between these attributes
37 are presented in Table 1(b). Clusters are presented in Table 1(c) and characterized through their
38 significant attributes. To assess the level of significance, the Test Value criterion is evaluated.

39 First, two types of drivers are distinguished, representing drivers working during nights
40 and during days. Each of these types are further divided into two separate clusters based on trip
41 lengths. The resulting four clusters are labelled: *Short trips night drivers*, *Long trips night drivers*,
42 *Short trips day drivers*, and *Long trips day drivers*, and represent 20.4 %, 34.9 %, 29.7 %, and 15
43 % our sample population, respectively.

44

1 **TABLE 1 Descriptive Analysis, Correlations Matrix, and Clusters' Specifications.**
 2

Attribute	Mean (SD)	Median	SD ^a	Min.	Max.
Day	57.3	60.0	36.5	0.0	100.0
Weekday	64.4	66.7	27.9	0.0	100.0
Peak	24.3	20.0	25.5	0.0	100.0
Length (km)	18.5	18.3	7.3	1.2	65.6
Travel time (min)	21.4	20.9	7.3	1.2	62.4
% Highway	54.4	56.7	20.5	0.0	100.0
Path-Size	0.97	1.0	0.1	0.4	1.0
# Links/km	5.69	5.6	1.3	2.5	12.3
# Turns/km	0.64	0.6	0.3	0.0	2.4
% Shortest path	43.9	43.1	18.2	0.0	99.4
# TAZ	12.7	8.0	15.3	1.0	154.0

^a Standard Deviation

3 **(a) Descriptive analysis of significant factors**
 4

Attribute	1	2	3	4	5	6	7	8	9	10	11
1 Day	-										
2 Weekday	0.31	-									
3 Peak	0.61	0.17	-								
4 Length	-0.11	-0.03	-0.07	-							
5 Travel time	0.08	0.06	0.08	0.69	-						
6 % Highway	-0.08	-0.01	-0.07	0.65	0.19	-					
7 Path-Size	-0.06	-0.09	-0.05	0.20	0.25	0.13	-				
8 # Links/km	0.02	-0.05	0.05	-0.58	-0.14	-0.92	-0.08	-			
9 # Turns/km	0.18	0.07	0.15	-0.35	-0.12	-0.61	-0.14	0.66	-		
10 % Shortest path	0.15	0.08	0.08	-0.39	-0.38	-0.17	-0.12	0.10	-0.05	-	
11 # TAZ	-0.03	0.06	-0.01	-0.35	-0.37	-0.20	-0.59	0.15	0.22	0.08	-

5 **(b) Correlation matrix of significant factors**
 6

	Short trip night drivers		Long trip night drivers		Short trip day drivers		Long trip day drivers		All
	TV ^a	Mean (SD) ^b	TV	Mean (SD)	TV	Mean (SD)	TV	Mean (SD)	Mean (SD)
Day	-16.7	28.5 (28.2)	-20.0	33.5 (27.7)	22.4	87.5 (16.2)	16.9	92.6 (13.6)	57.3 (36.5)
Weekday	-10.5	50.5 (28.1)	-10.0	55.3 (25.9)	10.3	75.0 (22.9)	12.2	83.7 (21.5)	64.4 (27.9)
Peak	-13.4	8.2 (10.6)	-16.9	10.3 (12.0)	17.5	40.8 (24.4)	15.2	46.4 (24.3)	24.3 (25.5)
Length (km)	-15.1	13.3 (4.6)	18.2	22.8 (5.5)	-17.4	13.8 (5.1)	15.2	24.9 (6.5)	18.6 (7.3)
Travel time	-11.3	17.5 (6.3)	5.4	22.7 (5.3)	-5.9	19.8 (7.2)	13.1	26.8 (8.4)	21.4 (7.3)
% Highway	-17.4	37.5 (15.8)	19.8	67.6 (11.4)	-17.5	41.1 (16.9)	15.7	72.6 (13.6)	54.4 (20.5)
Path-Size	-4.0	0.9 (0.1)	6.6	1.0 (0.5)	-3.8	0.9 (0.1)	6.2	1.0 (0.1)	1.0 (0.1)
# Links/km	17.6	6.7 (1.1)	-17.8	4.9 (0.8)	15.5	6.5 (1.2)	-16.0	4.5 (0.7)	5.7 (1.3)
# Turns/km	7.9	0.8 (0.3)	-15.1	0.5 (0.2)	15.8	0.8 (0.3)	-9.0	0.5 (0.2)	0.7 (0.3)
%SP ^c	4.5	47.8 (18.4)	-9.0	38.6 (16.3)	6.9	48.6 (17.4)	-1.8	42.0 (19.9)	44.0 (18.2)
# TAZ	6.8	17.6 (22.5)	-5.5	10.0 (7.4)	4.6	15.3 (17.3)	-6.2	7.3 (8.5)	12.7 (15.3)
# Obs. (%)	357 (20.4 %)		610 (34.9 %)		517 (29.7 %)		262 (15 %)		1746 (100)

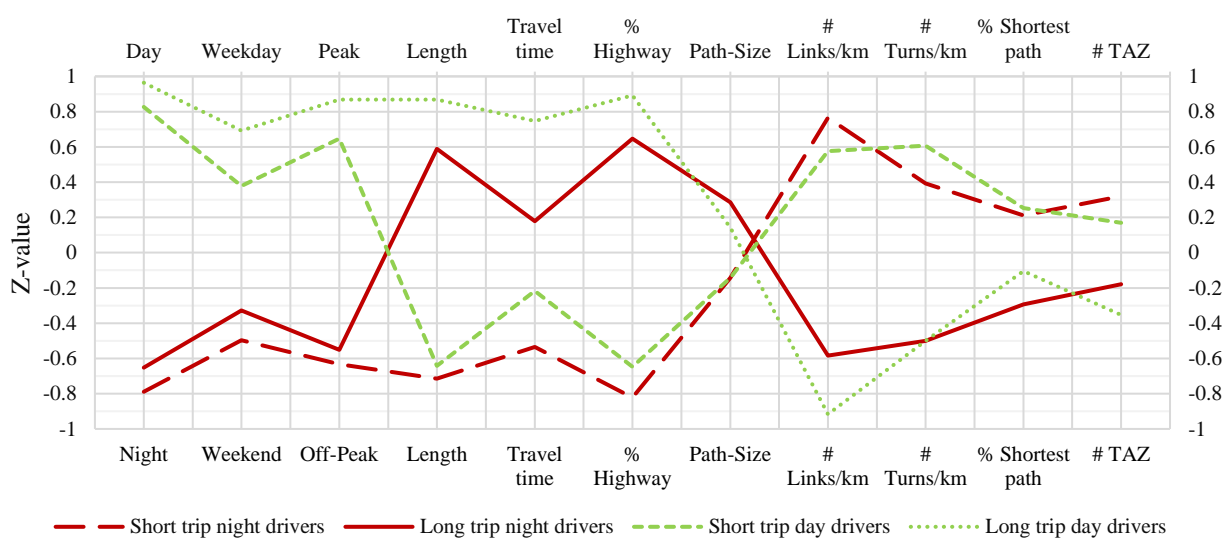
^a Test Value

^b Standard Deviation

^c Shortest Path

(c) Clusters' specifications

1 In general, two types of behavior are observed across these clusters. Drivers who most
 2 frequently take shorter trips, usually do not use highways and prefer to take local routes with higher
 3 number of intersections and turns, which accentuate their familiarity with the road network and
 4 their awareness of traffic conditions. These drivers tend to choose similar routes between same
 5 OD pairs, which have higher proportions of identical segments with the distance-based shortest
 6 path. On the other hand, longer trip drivers showed higher propensity toward using highways to
 7 avoid intersections and turns, which is intuitive. These drivers revealed more willingness in taking
 8 more diverse routes between the same OD pairs, which are less similar to the distance-based
 9 shortest path than routes taken in shorter trips. This might be partly due to the higher number of
 10 feasible alternatives between distanced OD pairs compared to closer ones. In order to visualize
 11 these clusters, the Z-score (cluster mean minus population mean, divided by population standard
 12 deviation), has been calculated (see Figure 5).
 13



14
 15

16 **FIGURE 5 Profile plot of the four clusters.**

17

18 CONCLUSIONS

19 Understanding heterogeneous route choice behaviors is very important for city and transportation
 20 planners. Although previous studies showed that different categories of road users are observable,
 21 there is a lack of representative classifications of drivers based on their actual route choices over
 22 a long duration of time. New possibilities have opened up due to the availability of very large and
 23 rich datasets from mobility services (e.g. taxi companies) that have detailed GPS trajectories of
 24 their fleets for a long duration of time.

25 The main objective of this research is to improve the understanding of drivers' route choice
 26 behaviors, by classifying their behaviors based on their observed route choices, using these
 27 longitudinal datasets. A GPS dataset comprising more than 22,000 trips, made by 1,746 taxi drivers
 28 over a period of one year, for trips originating in Montreal with a destination in Laval, has been
 29 studied. It is worth mentioning that the destination choice is not the focus of this research and this
 30 specific regional context has been chosen to ensure a wide range of traveled distances.

1 We first present important statistical properties of taxi trips. Since various degrees of
2 correlation were observed among attributes, a PCA analysis has been performed to improve the
3 efficiency of the clustering algorithm and to reduce the probability of the algorithm getting stuck
4 in a local optima (26). Then, a HAC algorithm has been performed to classify drivers' behaviors.

5 Due to the significant behavioral variations found in trips made during days and nights,
6 and between short trips and long trips, these four clusters were labelled "*Short trips night drivers*",
7 "*Long trips night drivers*", "*Short trips day drivers*", and "*Long trips day drivers*". Intuitively,
8 drivers prefer highways for longer trips and take local routes with higher number of turns and
9 intersections for shorter trips.

10 Although it is not possible to encompass all variations of route choice behaviors based on
11 GPS traces alone (due to the lack of some explanatory variables, such as demographics and
12 preferences), this study shed some light on the variation of route choice behaviors and the
13 possibility of classifying them using clustering algorithms.

14 Findings of this research pave the route for several future directions. A major limitation of
15 this work was the lack of personal information such as demographic and socio-economic variables.
16 Similar longitudinal datasets including these variables can be used to enrich the findings with more
17 personal information. Population synthesis can also be used as a complementary tool to
18 compensate for the lack of these information (31).

19 In this study, we limited our sampled population to taxi drivers. Taxi drivers are considered
20 as a well-informed group of drivers who have acquired higher knowledge of the road network and
21 its travel time variations. From a policy and planning perspective, an improved understanding of
22 taxi drivers' route choice behavior is important, since taxi, as an important mode of transportation
23 in big cities, provides further insights on human mobility patterns and urban structures. A future
24 direction could be the inclusion of regular car drivers and to examine the validity of the four
25 behavioral clusters proposed in this study. Since the experience level of regular car drivers varies
26 widely from taxi drivers, we might expect different driving patterns. Another interesting
27 elaboration of this study could be the incorporation of these clusters into route choice models and
28 to examine the enhancement of their estimation and prediction abilities. Since our findings suggest
29 that discrete route choice behaviors exist, the application of latent class discrete choice models
30 would be appropriate to capture this behavioral heterogeneity. Lastly, the recent rise of ride-hailing
31 services has pushed the cities and transportation planners to develop new laws and policies for
32 urban mobility. This study can be very helpful in giving insights to decision makers and can help
33 them propose new laws and policies which safeguard taxi industry as well as encourage sharing
34 economy.

35 36 REFERENCES

- 37 [1] Manley, E., J. Addison, and T. Cheng. Shortest path or anchor-based route choice: a large-scale
38 empirical analysis of minicab routing in London. *Journal of Transport Geography*, Vol. 43, 2015,
39 pp. 123-139.
- 40 [2] Parkany, E., J. Du, L. Aultman-Hall, and R. Gallagher. Modeling stated and revealed route
41 choice: consideration of consistency, diversion, and attitudinal variables. *Transportation Research*
42 *Record: Journal of the Transportation Research Board*, No. 1985, 2006, pp. 29-39.
- 43 [3] Tawfik, A. M., H. A. Rakha, and S. D. Miller. An experimental exploration of route choice:
44 identifying drivers choices and choice patterns, and capturing network evolution. In *Intelligent*

- 1 *Transportation Systems (ITSC), 2010 13th International IEEE Conference on*, IEEE, 2010. pp.
- 2 1005-1012.
- 3 [4] Ramaekers, K., S. Reumers, G. Wets, and M. Cools. Modelling route choice decisions of car
- 4 travellers using combined GPS and diary data. *Networks and Spatial Economics*, Vol. 13, No. 3,
- 5 2013, pp. 351-372.
- 6 [5] Koller-Matschke, I., H. Belzner, and R. Glas. Traffic Information and Individual Driver
- 7 Behaviour. In *Proceedings of the FISITA 2012 World Automotive Congress*, Springer, 2013. pp.
- 8 179-185.
- 9 [6] Gan, H., and S. Chen. Why Do Drivers Change Routes? Impact of Graphical Route Information
- 10 Panels. *Institute of Transportation Engineers. ITE Journal*, Vol. 83, No. 8, 2013, p. 38.
- 11 [7] Prato, C. G., S. Bekhor, and C. Pronello. Latent variables and route choice behavior.
- 12 *Transportation*, Vol. 39, No. 2, 2012, pp. 299-319.
- 13 [8] Cools, M., E. Moons, L. Creemers, and G. Wets. Changes in travel behavior in response to
- 14 weather conditions: do type of weather and trip purpose matter? *Transportation Research Record:*
- 15 *Journal of the Transportation Research Board*, No. 2157, 2010, pp. 22-28.
- 16 [9] Ben-Akiva, M., A. De Palma, and K. Isam. Dynamic network models and driver information
- 17 systems. *Transportation Research Part A: General*, Vol. 25, No. 5, 1991, pp. 251-266.
- 18 [10] Ulleberg, P. Personality subtypes of young drivers. Relationship to risk-taking preferences,
- 19 accident involvement, and response to a traffic safety campaign. *Transportation Research Part F:*
- 20 *Traffic Psychology and Behaviour*, Vol. 4, No. 4, 2001, pp. 279-297.
- 21 [11] Lucidi, F., A. M. Giannini, R. Sgalla, L. Mallia, A. Devoto, and S. Reichmann. Young novice
- 22 driver subtypes: relationship to driving violations, errors and lapses. *Accident Analysis &*
- 23 *Prevention*, Vol. 42, No. 6, 2010, pp. 1689-1696.
- 24 [12] Marengo, D., M. Settanni, and G. Vidotto. Drivers' subtypes in a sample of Italian
- 25 adolescents: relationship between personality measures and driving behaviors. *Transportation*
- 26 *Research Part F: Traffic Psychology and Behaviour*, Vol. 15, No. 5, 2012, pp. 480-490.
- 27 [13] Saunier, N., N. Mourji, and B. Agard. Mining microscopic data of vehicle conflicts and
- 28 collisions to investigate collision factors. *Transportation Research Record: Journal of the*
- 29 *Transportation Research Board*, No. 2237, 2011, pp. 41-50.
- 30 [14] Higgs, B., and M. Abbas. A two-step segmentation algorithm for behavioral clustering of
- 31 naturalistic driving styles. In *16th International IEEE Conference on Intelligent Transportation*
- 32 *Systems (ITSC 2013)*, IEEE, 2013. pp. 857-862.
- 33 [15] Jensen, M. Passion and heart in transport—a sociological analysis on transport behaviour.
- 34 *Transport Policy*, Vol. 6, No. 1, 1999, pp. 19-33.
- 35 [16] Kroesen, M., and S. Handy. The relation between bicycle commuting and non-work cycling:
- 36 results from a mobility panel. *Transportation*, Vol. 41, No. 3, 2014, pp. 507-527.
- 37 [17] Geller, R. *Four types of cyclists*. PortlandOnline, 2009.
- 38 [18] Dill, J., and N. McNeil. Four types of cyclists? Examining a topology to better understand
- 39 bicycling behavior and potential. In *92nd Annual Meeting of the Transportation Research Board,*
- 40 *Washington, DC*, 2013.
- 41 [19] Reinoso, G., and B. Farooq. URBAN PULSE ANALYSIS USING BIG DATA.
- 42 [20] Tawfik, A. M., H. A. Rakha, and S. D. Miller. Driver route choice behavior: Experiences,
- 43 perceptions, and choices. In *Intelligent Vehicles Symposium (IV), 2010 IEEE*, 2010. pp. 1195-1200.
- 44 [21] Tawfik, A. M., J. Szarka, L. House, and H. A. Rakha. Disaggregate route choice models based
- 45 on driver learning patterns and network experience. In *2011 14th International IEEE Conference*
- 46 *on Intelligent Transportation Systems (ITSC)*, 2011. pp. 445-450.
- 47 [22] Liu, L., C. Andris, and C. Ratti. Uncovering cabdrivers' behavior patterns from their digital
- 48 traces. *Computers, Environment and Urban Systems*, Vol. 34, No. 6, 2010, pp. 541-548.

- 1 [23] Communauté métropolitaine de montréal, C. *An Attractive, Competitive and Sustainable*
2 *Greater Montreal*, 2012.
- 3 [24] Ben-Akiva, M., and M. Bierlaire. Discrete choice methods and their applications to short term
4 travel decisions. In *Handbook of transportation science*, Springer, 1999. pp. 5-33.
- 5 [25] Pagès, J. Analyse factorielle de données mixtes. *Revue de statistique appliquée*, Vol. 52, No.
6 4, 2004, pp. 93-111.
- 7 [26] Poucin, G., B. Farooq, and Z. Patterson. Pedestrian Activity Pattern Mining in WiFi-Network
8 Connection Data. In *Transportation Research Board 95th Annual Meeting*, 2016.
- 9 [27] Rakotomalala, R. TANAGRA: un logiciel gratuit pour l'enseignement et la recherche. In *EGC*,
10 2005. pp. 697-702.
- 11 [28] Lebart, L., A. PIRON, M. Lebart, A. Morineau, and M. Piron. *Statistique exploratoire*
12 *multidimensionnelle*. Dunod, 2000.
- 13 [29] Everitt, B. S., D. Stahl, M. Leese, and S. Landau. Cluster analysis. 2011.
- 14 [30] Mooi, E., and M. Sarstedt. *Cluster analysis*. Springer, 2010.
- 15 [31] Farooq, B., M. Bierlaire, R. Hurtubia, and G. Flötteröd. Simulation based population
16 synthesis. *Transportation Research Part B: Methodological*, Vol. 58, 2013, pp. 243-263.
- 17