



ALMA MATER STUDIORUM
UNIVERSITÀ DI BOLOGNA

ARCHIVIO ISTITUZIONALE
DELLA RICERCA

Alma Mater Studiorum Università di Bologna Archivio istituzionale della ricerca

Categorizing three active cyclist typologies by exploring patterns on a multitude of GPS crowdsourced data attributes

This is the final peer-reviewed author's accepted manuscript (postprint) of the following publication:

Published Version:

Poliziani C., Rupi F., Mbuga F., Schweizer J., Tortora C. (2021). Categorizing three active cyclist typologies by exploring patterns on a multitude of GPS crowdsourced data attributes. RESEARCH IN TRANSPORTATION BUSINESS & MANAGEMENT, 40, 1-8 [10.1016/j.rtbm.2020.100572].

Availability:

This version is available at: <https://hdl.handle.net/11585/796086> since: 2023-01-30

Published:

DOI: <http://doi.org/10.1016/j.rtbm.2020.100572>

Terms of use:

Some rights reserved. The terms and conditions for the reuse of this version of the manuscript are specified in the publishing policy. For all terms of use and more information see the publisher's website.

This item was downloaded from IRIS Università di Bologna (<https://cris.unibo.it/>).
When citing, please refer to the published version.

(Article begins on next page)

RESEARCH IN TRANSPORTATION BUSINESS & MANAGEMENT

Categorizing three active cyclist typologies by exploring patterns on a multitude of GPS crowdsourced data attributes.

Paper for peer review. Special issue: NECTAR2020, Venice

Cristian Poliziani¹, Federico Rupi¹, Felix Mbuga², Joerg Schweizer¹, Cristina Tortora²

¹DICAM- University of Bologna, Viale Risorgimento 2, 40136 Bologna, Italy

²San José State University, One Washington square, 95192 San Jose, CA USA

Abstract

The paper tackles the problem of characterizing cyclist typologies using a data set of GPS traces. The data is crowdsourced and consists of 29,431 traces recorded in the city of Bologna, Italy during the morning rush hours from 7am to 10am of work-days and during the months from April through September 2017. Different criteria to group the heterogenous behavior of cyclists into separate categories have already been described in literature, where studies are generally based on stated preference interviews or on observing a small sample of the population. The novelty of this study is clustering a data set based on GPS traces from 2,135 cyclists, which is the equivalent to a revealed survey. Furthermore, refined pre-processing of the GPS traces allows the determination of dynamic attributes, a comparison of the chosen path respect to the shortest path and the evaluation of other specific trip attributes, which are either difficult or impossible to assess by a classical interview. The applied clusterization process leads to 3 main cyclist typologies, where each type is characterized by different trip attributes and behaviors involving safety, riskiness, precaution, inexperience, knowledge, fear and hastiness: *Risky & Hasty, Inexperienced & Inefficient and Smart & Informed*.

Keywords

Cyclist typology, Cluster analysis, GPS trace, Big data

Highlights

- *Crowdsourced big data of GPS traces for the characterization of cyclist typologies*
- *Cluster analysis with more than 50 trip characterizing attributes.*
- *Three cyclist typologies: Risky & Hasty, Inexperienced & Inefficient and Smart & Informed*

1. Introduction

1.1 Motivation and scope

The availability of a big set of GPS traces from cyclists, together with a range of instruments and tools able to deeply analyze GPS data has led to the idea of identifying different types of cyclists using a cluster analysis method, which is the core of the present work. The motivation behind this approach is to replace interview-based surveys with a large quantity of revealed data in order to obtain a more objective classification of cyclists. The knowledge of the type of cyclist is relevant for the link-cost and route-choice modeling of cyclists.

37 1.2 Literature review and State of the art

38 In order to properly plan bicycle infrastructure, it is crucial to know and characterize the current
39 transportation demand. However, data collection for the purpose of demand modeling is a challenge faced
40 by many cities. Most municipalities do not systematically monitor cycling activities.

41 Cyclists are heterogeneous and show different riding behavior. The differentiation between individual user
42 groups by targeting more accurately the needs and requirements of different types of users aims to better
43 conduct bicycle infrastructure planning, model cyclist behaviors, identify critical points and reduce barriers
44 for cycling. Planning a network adapted to different cyclist types could be an effective strategy to increase
45 cycling mode share and frequency among the various groups (Damant-Sirois et al., 2014).

46 Many studies are based on survey data to group them according to different aspects of their riding behavior.
47 The main cyclist typology factors found as trip purpose (e.g., Kroesen and Handy, 2014), comfort/perception
48 of safety (Geller, 2006) as well as motivational (e.g., Gatersleben and Haddad, 2010; Damant-Sirois et al.,
49 2014) and social factors such as identification as a cyclist (Damant-Sirois et al., 2014). Kroesen and Handy
50 (2014) considered non cyclists and cyclists for work and non-work purposes and used a latent transition
51 model for clustering people into four groups: non-cyclists, non-work cyclists, all-around cyclists (for work
52 and non-work purposes) and commuter cyclists. Gatersleben and Haddad (2010), using the results of a
53 survey conducted amongst 244 cyclists and non-cyclists in England, distinguished four different bicyclist
54 types based on behavior, motivation and characteristics of the typical bicyclist: responsible, lifestyle,
55 commuter and day-to-day. These types differed between bicyclists and non-bicyclists.

56 Geller (2006) identified four types of cyclists, subjectively developed on the basis on his expert knowledge:
57 strong and fearless, enthused and confident, interested but concerned, and no way no how. This typology is
58 based on perceived safety (comfort level on different types of bikeways and fear of people driving
59 automobiles) and on people's interest in cycling more. This classification, referred to all adults regardless of
60 their current cycling behavior, was subsequently formalized in a method and validated by Dill and McNeil by
61 a random phone survey in a sample of adults in the 50 largest metro regions in the U.S. (Dill and McNeil,
62 2013; 2016). Damant-Sirois et al. (2014), using data from an online survey aimed only at cyclists, propose a
63 multidimensional cyclist typology based on seven factors including weather conditions and effort, time
64 efficiency, street design, bicycle facilities, personal identity toward cycling and past cycling history. They
65 distinguished four distinct cyclist types: dedicated cyclists, path-using cyclists, fairweather utilitarian and
66 leisure cyclists. More recently, Cabral and Kim (2020) question the classic Four Types of Cyclists proposed by
67 Geller, particularly with respect to perceived comfort. They used an online survey and video clips to classify
68 people into three categories: Uncomfortable or Uninterested, Cautious Majority, and Very Comfortable
69 Cyclists. Their empirical segmentation is based on variables of comfort, cycling intent, and cycling in the
70 previous summer.

71 Francke et al. (2020) propose a multidimensional typology of cyclists which includes the influence factors of
72 already existing studies complemented by motivational factors. They use an empirical approach through a
73 Germany-wide online survey (10,294 responses) on cycling behavior in order to distinguish four distinct
74 types of cyclists: ambitious, functional, pragmatic, and passionate cyclists.

75 The use of recorded GPS data instead of interviews may be the way forward as GPS data represent objective
76 information about the chosen route and motion of each cyclist. A large number of GPS traces are usually
77 available from bike-tracing campaigns. In recent years, many studies on cycling mobility have made use of
78 GPS data, which is often available at low cost and allows to gain a broad range of information, such as the
79 spatial distribution of cyclists on the city's road network. Such information allows to calibrate the cyclist
80 route choice model, see (Lu et al., 2018; Rupi et al., 2019; Rupi and Schweizer, 2018; Pritchard et al., 2019;
81 Pritchard , 2018; Griffin and Jiao, 2015; Charlton et al., 2011; Dill, 2009; Menghini et al., 2010; Hood et al.,

82 2011; Broach et al., 2012; Zimmermann et al., 2017; Bernardi et al., 2018; Schweizer et al., 2020; Chen et
83 al., 2018). Other studies use GPS traces to obtain information on cyclist speeds (Manum et al., 2018, Flügel
84 et al., 2017, Strauss et al., 2017), speed profiles (Strauss and Miranda-Moreno, 2017; Clarry et al., 2019;
85 Laranjeiro et al., 2019) and waiting times at intersections (Watkins and LeDantec, 2016; Rupi et al. 2020).

86 **1.3 Research contribution**

87 The present study explores a new method to identify and characterize different types of cyclists based on
88 GPS traces and some additional attributes of cyclists. The method is applied to the GPS traces recorded in
89 the city of Bologna, Italy. The novelty is the use of GPS traces and derived quantities (such as trip-length,
90 speeds, waiting times, deviation from shortest path, etc.) instead of interviews for the purpose of
91 identifying types of cyclist, which is equivalent to a revealed preference survey – hence the GPS traces are
92 expected to be more reliable than classical surveys, where interviewees declare their behavior. In addition,
93 GPS traces are often available in large quantities. The described method requires a refined pre-processing of
94 the GPS traces: The traces are matched to a road network extracted from Open Street Map (OSM) and
95 elaborated in order to create a rich database that describes the experience, the performance, choices and
96 behavior of each cyclist as detailed in Rupi et al., 2020.

97 The main focus of the paper is on the successive cluster analysis which uses the results from the pre-
98 processing step with the scope of clustering the data set in groups of cyclists characterized by similar habits.

99 Section 2 describes the applied clustering method and section 3 presents the Bologna GPS data set and
100 briefly explains the data pre-processing. The results in section 4 present the 3 types of cyclists and illustrate
101 their respective characteristics.

102 **2. Cluster analysis**

103 The goal of cluster analysis is to find homogeneous groups of units within the data, i.e. homogeneous
104 groups of cyclists based on their habits. There are many clustering techniques in the statistical literature
105 (see Everitt 2011 for examples), among them model based clustering techniques are popular. Model-based
106 clustering assumes that a population is a convex combination of a finite number of density functions. The
107 multivariate Gaussian distribution is one of the most popular for its simplicity (McLachlan and Basford
108 1988): each cluster has only two parameters, the mean vector that determines the position of the cluster in
109 the space, and the covariance matrix.

110 Figure 1 (a) shows an example of a two-dimensional data set with 3 clusters that follow a bivariate Gaussian
111 distribution. In cluster one (black circles) the variables have unitary variance and no correlation, in cluster
112 two (red triangles) the variables have variance equal to 2 and no correlation, and in cluster three (green
113 plues) the variables have unitary variance and correlation equal to 0.6.

114 Formally, a p-dimensional random vector X follows a finite mixture of distributions if, for all $x \in X$, its

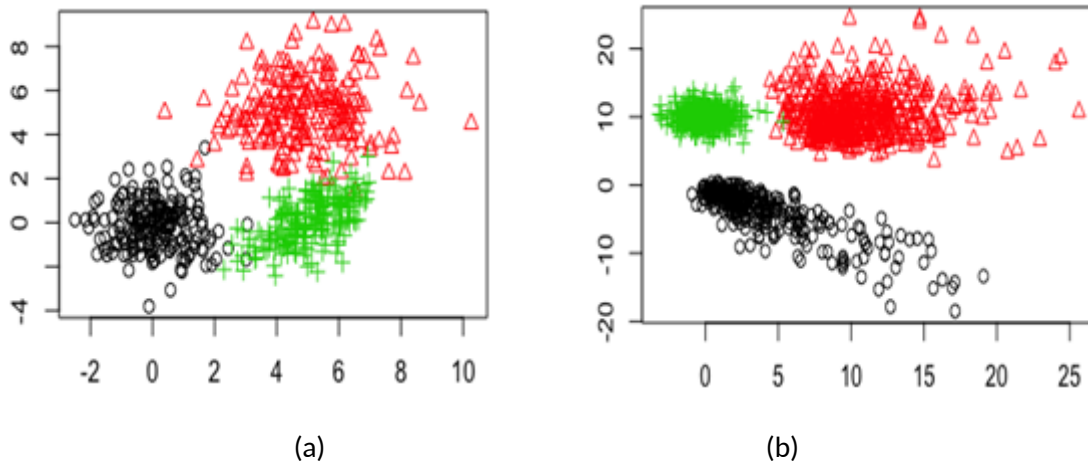
115 density can be written as: $f(x|\vartheta) = \sum_{g=1}^G \pi_g f(x|\theta_g)$, where G is the number of clusters, $\pi_g > 0$, such that

116 $\sum_{g=1}^G \pi_g = 1$, is the g th mixing proportion, $f(x|\theta_g)$ is the g th component density, and

117 $\vartheta = (\pi_1, \dots, \pi_G, \theta_1, \dots, \theta_G)$ is the vector of parameters. One of the advantages of model-based clustering is
118 that, besides the partition in clusters, the method produces an estimate of the parameters θ_g of each
119 cluster that can be used for interpretation.

120 A variety of distributions have been used to model the density functions; McNicholas 2016 contains a good
121 review of them. Among these distributions, the generalized hyperbolic distribution (GHD) (Browne and

122 McNicholas 2015) has the advantage of being extremely flexible. The GHD is characterized by five
123 parameters - mean vector, scale matrix, skewness vector, concentration, and index parameters - and can
124 identify clusters of a different shape. Moreover, many other distributions - like the Gaussian distribution,
125 the multivariate Student t, or the skewed Student t distribution - can be obtained as a special or a limiting
126 case of the GHD, i.e. it can detect clusters that are for example normally distributed. Tortora et al. 2019
127 proposed a flexible extension of the GHD - the coalesced GHD (CGHD)- that adds even more flexibility and
128 can model non-convex clusters. Figure 1 (b) shows an example of a two-dimensional data set with 3 clusters
129 that follow a bivariate CGHD. The Bayesian information criterion (BIC) is recommended to choose the
130 number of clusters for mixtures of CGHD.



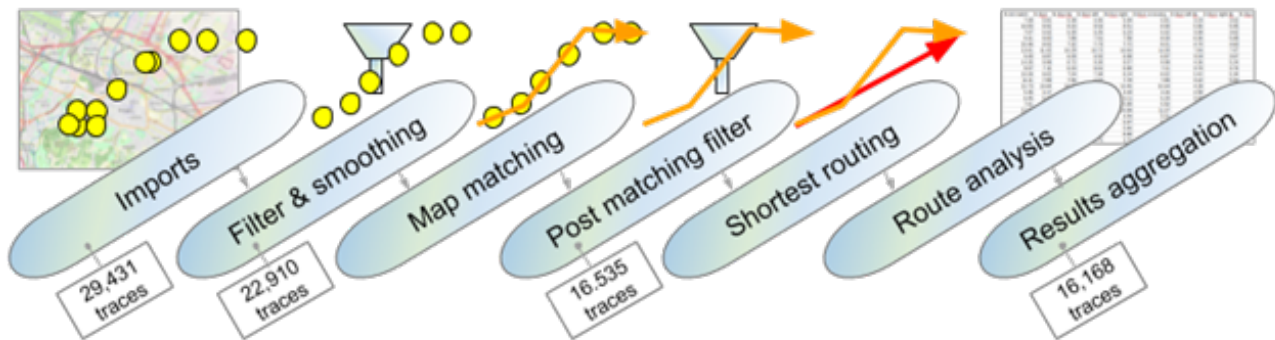
133 *Fig.1 Two dimensional data sets where each cluster was generated from a Gaussian distribution (a) or from*
134 *a CGHD (b).*

135 3. The Bologna data set

136 The city of Bologna, Italy, hosted in 2017 - from April 1st to September 30th - the 'Bella Mossa' initiative
137 (BM), funded by the EU and Bologna's local Government. The initiative had the objective to promote
138 sustainable mobility by rewarding people (with coupons for local shops) for recording their GPS traces of
139 sustainable trips - meaning trips done by transit, bike or walking . The smartphone application 'Betterpoints'
140 (Betterpoints, 2020) has been used to record and collect the data. The full data sample contains
141 approximately 270,000 bike GPS traces, which consist of more than 62 million points: the smartphone
142 application records one GPS point every 5 seconds when the bike is in motion. When the bike stops, for
143 example at intersections, the recording stops also. The present study focuses only on bikes GPS traces
144 recorded during the morning peak of work days - from 7am to 10am - because of computational time
145 problems, but also trying to englobe especially working trips during the morning peak hour and avoid as
146 much as possible trips for other purposes.

147 The following data processing steps have been implemented using the SUMOPy environment (Schweizer,
148 2020), an open source extension of the software SUMO (Eclipse SUMO, 2020). In a first step, the open
149 street map (OSM) network covering the urban area of Bologna (OpenStreetMap, 2020) has been imported
150 into SUMO. This SUMO network is attribute-rich and contains information on road width, road access (e.g.
151 reserved bikeways, shared access, with pedestrians, etc.) and speed-limits. From these basic attributes
152 SUMO derives a road priority (1-10), where low priority roads are taken values from 1 to 7. Successively the
153 network has been manually improved in order to eliminate errors due to an imperfect OSM representation
154 as well as conversion errors. Next, unrealistic GPS traces have been deleted: trips outside the study area and
155 traces which have probably not been recorded while riding a bike. (See Fig.2). Valid traces must also satisfy
156 criteria such as a certain total distance, a minimum number of points, and a minimum and maximum

157 distance between successive points. This trace filtering step does ensure the GPS traces can be successfully
 158 matched to the road network by the map-matching process. During the map-matching the most likely route
 159 (as sequence of network links) can be identified for each GPS trace (Schweizer et al., 2020).



160
 161 *Fig.2 Evaluating cyclist attributes starting from GPS traces registered with Smartphones.*

162 After a further filter that ensures the quality of the map-matching process, the shortest route connecting
 163 same link of origin and destination as the matched route has been identified for each trip. Subsequently,
 164 the matched routes have been analyzed - also from a dynamic point of view (see Rupi et al., 2020) - and
 165 compared with the respective shortest routes. There has been a further filtering of GPS traces that are
 166 suitable for the dynamic analysis. Finally, the attributes related to all the trips carried out from the same
 167 users have been aggregated in 5 different ways: mean, mean weighted by trips distances, median, standard
 168 deviation and mean absolute difference. These averaged attributes describe the experience of each user.
 169 After carrying out all the aforementioned pre-processing step the cyclist population is composed of 2,135
 170 users (see tab.1), recording a total of 16,168 trips. This cyclist population belongs mainly to the firsts two
 171 age groups - 16-44 years - based on the phases of life described by Wittwer et al., 2014.

N users	N trips	Male	16→29	30→44	45→54	55→65	>65
2135	16168	46.50%	43.90%	34.90%	12.20%	8.30%	0.60%

172 *Tab.1 Sample characterization in terms of size, gender and age - based on the phases of life described by*
 173 *Wittwer.*

174 The following aggregate attributes have been obtained from the GPS trace analysis: trip length, trip duration
 175 and average speed; number of road priority changes per km; share of trip inside the center area, in roads
 176 reserved for cyclists, in mixed roads - shared with taxi, bus or pedestrians - and in low-priority roads;
 177 number of passed nodes, left turns, right turns and crossing - both in total and only in presence of a traffic
 178 light system (TLS) - per km; average numbers of maneuvers in the traveled intersections; all the route-
 179 attributes minus the same attributes referred to the shortest route; share of the shortest route length
 180 respect to the matched route; average in-motion speed; waiting time, as well as the registered waiting time
 181 minus the expected waiting time - based on the waiting times registered by the other users - at nodes, left
 182 turns, right turns, crossings - both in total and only in presence of a TLS - and at edges, per km. In addition
 183 to these attributes, the Bellamossa database contains additional information on the users such as: age,
 184 gender and number of recorded traces by each participant.

185 Tab.2 shows a statistical description of all the trip attributes aggregated for each cyclist as the median of the
 186 attributes related to their recorded trips. The table highlights also the dispersion of the various attributes
 187 referred to all cyclists. This suggests that cyclists have different habits and there might be the possibility of
 188 grouping the data set in some cyclist clusters. In particular, the average speed in motion is 4.85 ± 1.04 m/s,
 189 the average speed is 3.77 ± 1.11 m/s and the waiting time represents on average 14.6 % of the whole trip

190 duration; these results are very similar to the dynamic results presented in (Rupi et al. 2020), where a
 191 different set of GPS-traces from Bologna, Italy, has been analyzed.

Attribute	Unit	Aver.	St.Dev.	Min	1stQu.	2ndQu.	3rdQu.	Max
Average length	m	2864	1607	189	1721	2505	3603	13650
Share of trip in the center area	%	42.57	35.51	0.00	0.00	40.92	71.62	100.00
Average speed	m/s	3.58	0.91	1.14	2.98	3.52	4.10	10.14
Number of priority changes	1/km	0.75	0.84	0.00	0.06	0.54	1.06	6.59
Share of reserved cycleway	%	23.79	20.10	0.00	6.98	20.07	36.47	100.00
Share of mixed road	%	30.52	18.43	0.00	16.84	27.91	41.27	100.00
Share of low-priority roads	%	69.91	22.70	0.00	54.37	72.59	89.08	100.00
Number of nodes	1/km	16.60	3.01	3.73	14.78	16.57	18.39	32.15
Number of left turns	1/km	2.11	0.98	0.00	1.45	2.02	2.68	6.97
Number of right turns	1/km	2.37	1.08	0.00	1.68	2.30	2.97	15.91
Number of crossings	1/km	11.36	2.81	0.00	9.45	11.25	13.18	28.37
Number of TLS nodes	1/km	2.82	1.60	0.00	1.70	2.76	3.73	10.92
Number of TLS left turns	1/km	0.39	0.45	0.00	0.00	0.31	0.60	4.14
Number of TLS right turns	1/km	0.38	0.38	0.00	0.00	0.34	0.60	3.07
Number of TLS crossings	1/km	1.89	1.21	0.00	0.98	1.82	2.68	8.32
Average number of maneuvers per node	/	10.30	1.19	6.00	9.64	10.22	10.82	22.25
Share shortest length	%	84.73	7.83	32.58	80.60	85.84	90.30	100.00
Number of priority changes*	1/km	-0.09	0.95	-6.18	-0.33	0.00	0.32	4.64
Share of reserved cycleway*	%	4.03	14.98	-59.80	-1.27	0.00	9.06	85.06
Share of mixed road*	%	-1.86	12.50	-72.05	-6.87	-0.31	4.16	52.28
Share of low-priority roads*	%	4.05	16.64	-63.12	-2.83	1.36	9.98	77.68
Number of nodes*	1/km	-0.73	2.33	-19.60	-1.75	-0.50	0.49	10.77
Number of left turns*	1/km	0.37	0.88	-5.74	-0.09	0.36	0.83	4.27
Number of right turns*	1/km	0.45	0.98	-5.17	-0.08	0.45	0.97	7.01
Number of crossings*	1/km	-1.59	2.33	-14.14	-2.70	-1.39	-0.27	10.11
Number of TLS nodes*	1/km	0.02	0.96	-5.13	-0.40	-0.04	0.43	7.50
Number of TLS left turns*	1/km	0.11	0.39	-2.73	-0.02	0.00	0.24	4.14
Number of TLS right turns*	1/km	0.11	0.36	-1.98	-0.01	0.00	0.28	2.66
Number of TLS crossings*	1/km	-0.22	0.81	-4.15	-0.57	-0.13	0.07	3.54
Average number of maneuvers per node*	/	0.08	0.87	-3.49	-0.31	0.04	0.39	5.64
Average in motion speed	m/s	4.70	0.92	1.54	4.13	4.63	5.19	11.11
Share of waiting time whole trip	%	14.71	9.96	0.00	7.56	13.13	19.88	71.74
Average waiting time per node	s	1.70	1.77	0.00	0.55	1.24	2.27	17.31
Average waiting time per left turn	s	1.94	4.67	0.00	0.00	0.00	1.88	67.00
Average waiting time per right turn	s	1.35	3.44	0.00	0.00	0.00	1.21	54.00
Average waiting time per crossing	s	1.61	2.11	0.00	0.33	1.00	2.15	35.73
Average waiting time per TLS node	s	4.12	6.32	0.00	0.00	2.25	5.70	96.30
Average waiting time per TLS left turn	s	3.22	8.95	0.00	0.00	0.00	1.00	85.00
Average waiting time per TLS right turn	s	2.82	9.29	0.00	0.00	0.00	0.00	122.00
Average waiting time per TLS crossing	s	4.44	8.20	0.00	0.00	1.71	6.00	147.17
Average waiting time per edge	s/km	13.17	25.75	0.00	0.69	5.34	14.34	334.25
Average waiting time per edge**	s/km	-	22.77	-67.66	-7.65	-4.01	0.92	316.39
Average waiting time per node**	s/km	-	30.62	-76.55	-13.16	-3.97	9.44	521.03

Average waiting time per left turn**	s/km	-	15.20	-57.42	-7.80	-2.59	3.49	242.91
Average waiting time per right turn**	s/km	-	8.94	-45.64	-2.73	-0.84	0.00	187.51
Average waiting time per crossing**	s/km	-	12.83	-34.14	-2.07	-0.71	0.00	472.56
Average waiting time per TLS node**	s/km	-	20.96	-65.45	-9.43	-3.28	4.01	239.87
Average waiting time per TLS left turn**	s/km	-	7.47	-25.23	-2.13	-0.33	0.00	187.51
Average waiting time per TLS right turn**	s/km	-	6.90	-32.55	-1.57	-0.44	0.00	115.30
Average waiting time per TLS crossing**	s/km	-	12.17	-49.67	-5.55	-1.94	1.54	213.77
Average waiting time whole trip**	s/km	-	39.48	-97.21	-16.49	-5.07	9.73	428.71
* Value referred to the matched trip minus the value referred to the shortest trip								
** real minus expected waiting time, based on the average waiting times of other cyclists that passed the same elements of the network								

192 *Tab.2 Descriptive statistics of cyclist attributes considered as the median of values referred to their*
193 *performed trips.*

194 From a first graphical analysis the Bologna cyclist data set does not clearly show spherical, symmetric, or
195 convex clusters, therefore, a flexible distribution, like the CGHD, is preferred for the cyclist clusterization.
196 The CGHD can capture differently shaped clusters, as well as, symmetric or spherical clusters.

197 **4. Results**

198 **4.1 Data Analysis**

199 Since the goal of the analysis is to find homogeneous groups of cyclists based on their cycling habits, the
200 demographic variables have only been used for interpretation of the clusters after the analysis.

201 The total number of variables used for the analysis is 41 and the missing values have been imputed using
202 multivariate imputation by chained equations (MICE) (Buuren and Groothuis-Oudshoorn 2010), i.e. each
203 missing value has been imputed 5 times, obtaining 5 different complete data sets; the same analysis has
204 been therefore performed on the 5 data sets, and the adjusted Rand Index (ARI) (Huber and Arabie 1985)
205 has been used as a measure of pairwise cluster agreement in order to assess the robustness of the models,
206 i.e. if all the partitions are similar, then the analysis is robust. The ARI is equal to 1 when there is perfect
207 agreement between two partitions and the expected value is 0 for random classification: in the case study,
208 the ARI was between 0.71 and 0.76, indicating a good level of robustness. The reported results refer to one
209 of the data sets.

210 The analysis has been run using the software R (R Core team, 2020) by varying the number of clusters
211 between 2 and 5, the BIC selected three clusters. The algorithm was initialized using a robust clustering
212 technique called k-medoids (Kaufman and Rousseeuw 1990), and the package MixGHD (Tortora et al. 2019)
213 has been used for the cluster analysis.

214 **4.2 Cluster analysis results**

215 The three clusters of cyclists present similar ages, but different trip attributes and cyclist choices and
216 behaviors. The cyclist groups are composed of 806, 749 and 580 cyclists, with respectively 51%, 63% and
217 46% of women. According to the cluster's attributes, the three cyclist categories have been named as
218 follow:

- 219 **1. Risky & Hasty cyclists (RHC):** they prefer going straight along the shortest path, traveling on
220 unsafe roads and also large roads. This type of cyclist accepts also roads without reserved
221 cycleway and a high density of intersections. The RHC is in average fast when in motion,
222 but he/she loses time encountering many traffic lights, even if the average waiting time at
223 intersections is fairly low, indicating that the RHC does often ignore the red signal.

224 • *Keywords: no deviations from the shortest path - low waiting times - fast*

225 2. **Inexperienced & Inefficient** cyclists (IIC): they travel at a low speed and they showed
226 considerably higher waiting times respect to the other cyclists in all infrastructure elements
227 (roads, intersections, traffic lights), probably because they are precautionary, but also
228 afraid and not completely convinced to use a bike. They make deviations from the shortest
229 path like the SIC group, but not to travel on safer roads with more exclusive cycleways, less
230 traffic and fewer intersections.

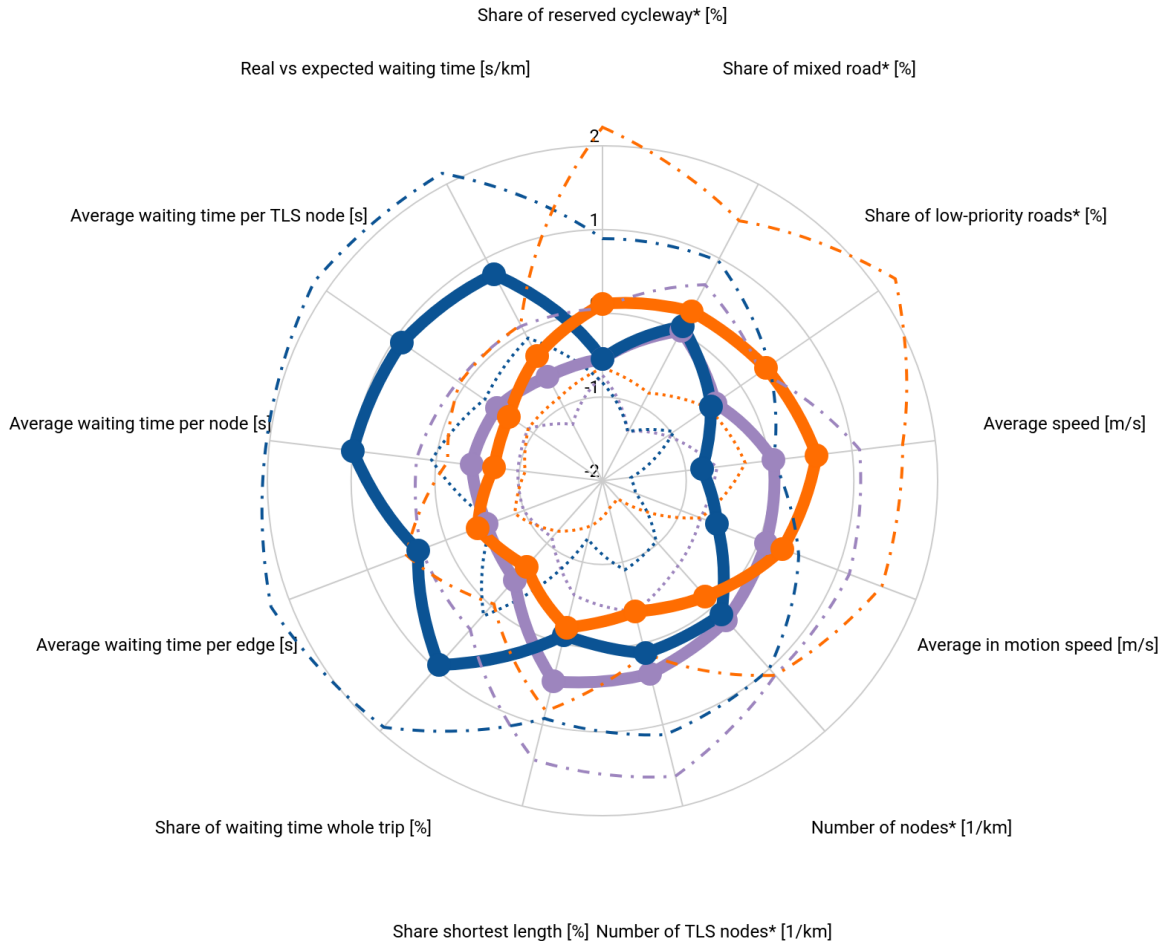
231 • *Keywords: considerable waiting times - low speeds*

232 3. **Smart & Informed** cyclists (SIC): they know how to make deviations from the shortest path
233 in an efficient way, searching frequently for roads with more reserved cycleways, fewer
234 traffic and fewer intersections, thus gaining on both, safety and speed. They showed fairly
235 low waiting times at intersections - like the RHC group - because they encounter less traffic
236 lights and they probably pass also with the red signal (supposedly when it is safe).

237 • *Keywords: smart deviations from the shortest path for safer or faster roads - low waiting
238 times - fast*

239 The radar graphs show the differences in terms of behavior (Fig.3) and trip characteristics (Fig.4) for the
240 three different types of cyclists. The graphs contain the normal standardization of the first, second and third
241 quartiles of the attributes referred to the cyclists of three clusters, - of which cyclist attributes are
242 considered as the median value of the attributes referred to their recorded traces - in order to better
243 visualize the differences in the same scale, while Tab.3 shows a descriptive analysis of the same attributes.

244

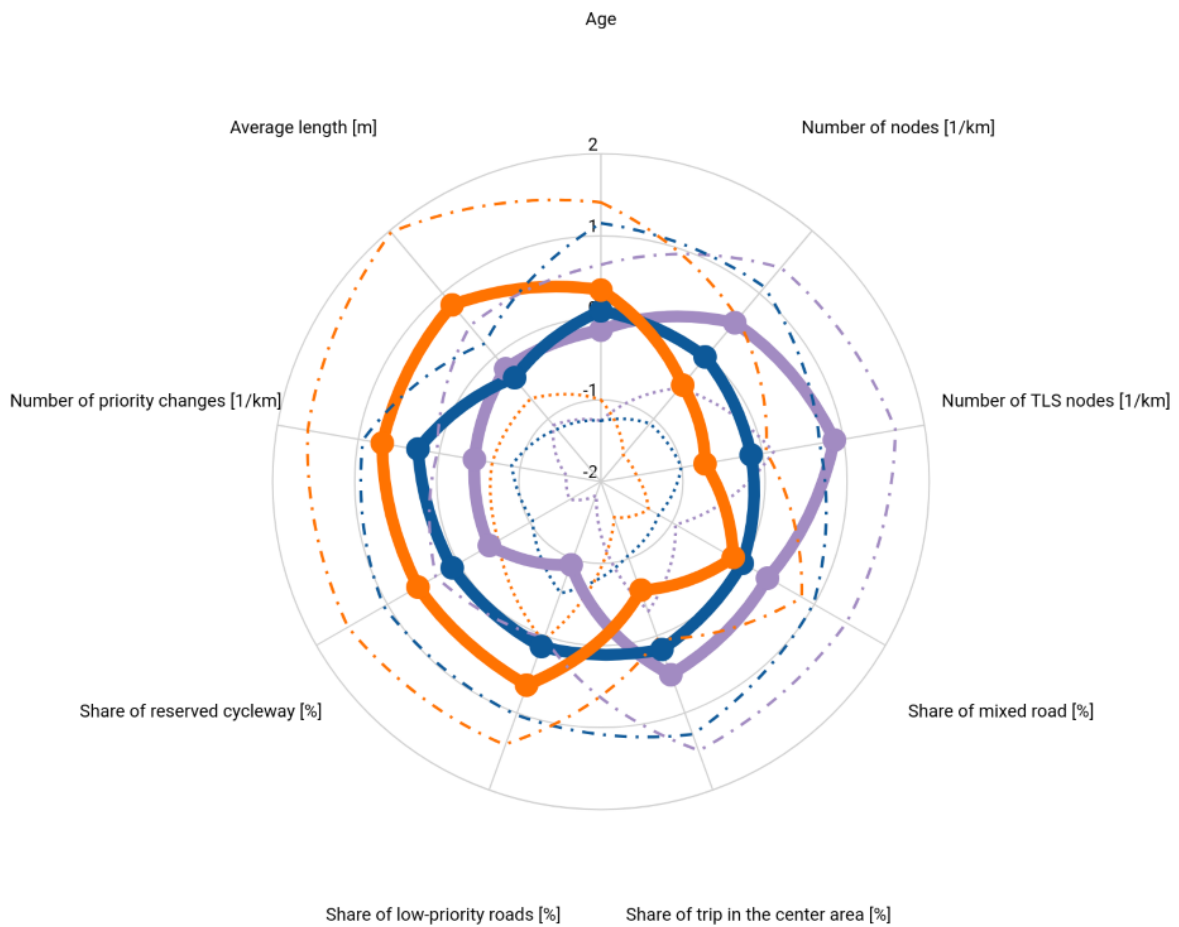


245

246 Fig.3 Behavioral differences of the three cyclist types - RHC in purple, IIC in blue and SIC in orange. The
 247 dotted lines represent the first and third quartiles, while the solid line represent the median of cyclist
 248 attributes considered as the median of values referred to their performed trips. * Value referred to the
 249 matched trip minus the value referred to the shortest trip.

250 The RHC type - who prefers to go straight on the shortest route - travel mainly on unsafe roads while
 251 encountering many intersections - of which most are with traffic lights - in particular in central areas of
 252 Bologna. The IIC type characterized by considerably high waiting times and low speed, indicating an
 253 insecure and cautious behavior. The other trips characteristics correspond to average values with respect
 254 the other cyclist types. The SIC type smartly deviates from the shortest route, often travels longer distances
 255 and changes frequently the priority of the road: his deviations contain safer roads with reserved cycleways,
 256 low priority roads and fewer intersections. The SIC also prefers to travel outside the center.

257



258

259 Fig.4 Differences in terms of trip characteristics of the three cyclist types - RHC in purple, IIC in blue and SIC
 260 in orange. The dotted lines represent the first and third quartiles, while the solid line represent the median
 261 of cyclist attributes considered as the median of values referred to their performed trips.

262

Attributes		1st Quartile			2nd Quartile			3rd Quartile		
Name	Unit	RCU	IIC	SIC	RCU	IIC	SIC	RCU	IIC	SIC
Share of reserved cycleway*	%	-1.2	-1.8	-0.7	0	0	4.1	3.8	9	17.3
Share of mixed road*	%	-7.9	-7.9	-5.2	-0.8	-0.5	0.6	2.5	4.2	7
Share of low-priority roads*	%	-4.9	-4.2	0	0.8	0.2	6.7	6	7.2	22
Average speed	m/s	3.2	2.6	3.4	3.6	3.1	3.9	4.2	3.6	4.5

Average in motion speed	m/s	4.3	3.9	4.3	4.7	4.4	4.8	5.2	4.9	5.4
Number of nodes*	1/km	-1.3	-1.7	-2.4	-0.4	-0.5	-0.8	0.5	0.4	0.5
Number of TLS nodes*	1/km	-0.2	-0.4	-0.7	0.1	0	-0.2	0.6	0.4	0
Share shortest length	%	82.8	79.7	79	87.4	84.9	84.5	91.7	89.4	89
Share of waiting time whole trip	%	5.8	14.6	5.1	10.6	20.2	9.1	16.2	27.2	13.3
Average waiting time per edge	s	0	2.9	0.8	3.4	10	4.3	9.6	23.9	11.1
Average waiting time per node	s	0.4	1.5	0.3	1	2.5	0.7	1.7	3.6	1.3
Average waiting time per TLS node	s	0	2.1	0	1.5	6.1	0.9	3.4	10.3	3.3
Real vs expected waiting time	s/km	-22.9	-3.7	-17	-12.1	10.8	-7.7	-0.8	33.3	-0.6
Age	-	13	15	16	30	31	35	39	43	45
Number of nodes	1/km	15.5	14.7	13.8	17.2	16.5	15.7	18.9	18.4	17.5
Number of TLS nodes	1/km	2.8	1.6	0.9	3.5	2.5	1.8	4.5	3.4	2.8
Share of mixed road	%	19.2	16.5	15	31.9	26.1	26	45.6	39.4	38.5
Share of trip in the center area	%	31	0	0	56.3	42.6	13.3	82.6	77.8	41.5
Share of low-priority roads	%	41.3	61.3	72.5	56.7	77.4	85.2	72.7	91	95.4
Share of reserved cycleway	%	2.9	8.1	15.5	12.3	21	28.9	28.5	37.3	43.6
Number of priority changes	1/km	0	0.2	0.4	0.3	0.6	0.8	0.7	1.1	1.4
Average length	m	1680	1612	2075	2346	2286	3325	3318	3038	4990
* Value referred to the matched trip minus the value referred to the shortest trip										

263 *Tab.3 Descriptive analysis of the clusters with cyclist attributes considered as the median of values referred*
264 *to their performed trips.*

265 5. Conclusions

266 The present study identified three types of cyclists by applying a cluster analysis to attributes calculated
267 from GPS traces. The data set is composed of 16,168 GPS traces recorded between 7-10 am on work days
268 from April to September 2017. Based on the results of the cluster analysis, the types have been named
269 **Risky & Hasty** cyclists (RHC), **Inexperienced & Inefficient** cyclists (IIC) and **Smart & Informed** cyclists (SIC).
270 The novelty of the present study, is the use of crowdsourced GPS data that allows to analyze a large data
271 sample as objective, revealed survey data. The pre-processing of the GPS traces has produced a large and
272 variegated set of attributes which characterize the cyclist experience, leading to a meticulous description of
273 the different types of cyclists. The most significant attributes in the cyclist clusterization are: the amount of
274 trip deviation made for improving the cyclist experience in terms of both safety and speed as well as the
275 waiting times detected in various parts of the road network.

276 Rather than providing threshold values to split other cyclist database into the three groups, a model-based
277 approach has been used. Each cluster is modeled using a CGHD distribution with different parameters. The
278 model can be used to find the belonging probability to each cluster of cyclists in the same area not included
279 in the original sample.

280 The proposed paper presents on what these groups differ and how much, thus presenting the differences in
281 terms of experience, performances, choices and behavior of each group of cyclists for design both new
282 surveys and infrastructures in the city - adapted to each group of cyclists - as well as initiatives that can
283 allow to improve the trip experience of certain cyclists, making it safer and more efficient as well as
284 reducing barriers for cycling.

285 One current limitation of the data set is that it may not be representative. In future research it is possible to
286 attribute more weight to the data of underrepresented parts of the population.

287 Acknowledgements

288 We are grateful to SRM Bologna srl for providing the GPS data of the Bella Mossa Campaign 2017.

289 **References**

- 290 Betterpoints. (2020). <https://www.betterpoints.ltd/behaviour-change/> Accessed 1 May 2020.
- 291 Bernardi, S., La Paix-Puello, L., Geurs, K., (2018). Modelling route choice of Dutch cyclists using smartphone
292 data. *J. Trans. Land Use* 11 (1), 883–900.
- 293 Broach, J., Dill, J., Gliebe, J., (2012). Where do cyclists ride? A path choice model developed with revealed
294 preference GPS data. *Transp. Res. Part A* 46, 1730–1740. 10.1016/j.tra.2012.07.005
- 295 Browne, R. P., McNicholas, P. D. (2015). A mixture of generalized hyperbolic distributions. *Canadian Journal*
296 *of Statistics*, 43(2), 176-198. 10.1002/cjs.11246
- 297 Buuren, S. V., Groothuis-Oudshoorn, K. (2010). mice: Multivariate imputation by chained equations in R.
298 *Journal of statistical software*, 1-68. 10.18637/jss.v045.i03
- 299 Cabral, L, Kim, A. M., (2020). An empirical reappraisal of the four types of cyclists. *Transportation Research*
300 *Part A: Policy and Practice*, Volume 137, 206-221. 10.1016/j.tra.2020.05.006
- 301 Charlton B., Sall E., Schwartz M., Hood J. (2011). Bicycle Route Choice Data Collection using GPS-Enabled
302 Smartphones, *TRB 2011 Annu. Meet.* 1– 10.
- 303 Chen P., Shen Q., Childress S., (2018). A GPS data-based analysis of built environment influences on bicyclist
304 route preferences, *Int. J. Sustain. Transp.* 12(3), 218–231. 10.1080/15568318.2017.1349222
- 305 Clarry, A., Faghih Imani, A., Miller, E.J., (2019). Where we ride faster? Examining cycling speed using
306 smartphone GPS data. *Sustainable Cities and Society* 49, 101594. 10.1016/j.scs.2019.101594.
- 307 Damant-Sirois G., Grimsrud M., El-Geneidy A. M. (2014). What's your type: a multidimensional cyclist
308 typology. *Transportation*. 41(6), 1153–1169. 10.1007/s11116-014-9523-8
- 309 Dill, J., (2009). Bicycling for Transportation and Health: The Role of Infrastructure. *J. Public Health Policy* 30,
310 95–110. 10.1057/jphp.2008.56
- 311 Dill J., McNeil N. (2013). Four Types of Cyclists? Examination of Typology for Better Understanding of
312 Bicycling Behavior and Potential. *Transportation Research Record: Journal of the Transportation Research*
313 *Board*. 2387, 129-138. 10.3141/2387-15
- 314 Dill J., McNeil N. (2016). Revisiting the Four Types of Cyclists: Findings from a National Survey.
315 *Transportation Research Record: Journal of the Transportation Research Board*. 2587, 90-99. 10.3141/2587-
316 11
- 317 Eclipse SUMO. (2020). <http://sumo.sourceforge.net/userdoc/> Accessed 1 May 2020.
- 318 Everitt, Brian S., et al. (2011). Cluster analysis. *John Wiley & Sons*.
- 319 Flügel S., Hulleberg N., Fyhri A., Weber C., Evarsson G. (2017). Empirical speed models for cycling in the
320 Oslo road network. *Transportation (Amst)*. 1.
- 321 Francke A., Anke J., Lißner S., Schaefer L.-M., Becker T., Petzoldt T. (2020). Are you an ambitious cyclist?
322 Results of the cyclist profile questionnaire in Germany. *Traffic Inj. Prev.* 0(0), 1–6.
- 323 Gatersleben B., Haddad H. (2010). Who is the typical bicyclist? *Transp Res Part F Traffic Psychol Behav*.
324 13(1), 41–48. 10.1016/j.trf.2009. 10.003

325 Geller, R. (2006) Four Types of Cyclists. Portland Bureau of Transportation, Portland, Ore.
326 <http://www.portlandoregon.gov/transportation/article/264746> Accessed 31 May 2020.

327 Griffin, G.P., Jiao, J., (2015). Where does bicycling for health happen? Analyzing volunteered geographic
328 information through place and plexus. *J. Transp. Health* 2, 238–247.
329 <https://doi.org/10.1016/j.jth.2014.12.001>

330 Hood, J., Sall, E., Charlton, B., (2011). A GPS-based bicycle route choice model for San Francisco. California.
331 *Transp. Lett. Int. J. Transp. Res.* 3, 63–75. <https://doi.org/10.3328/TL.2011.03.01.63-75>

332 Hubert, L., Arabie, P. (1985). Comparing partitions. *Journal of classification*, 2(1), 193-218.
333 10.1007/BF01908075

334 Joo S. Oh C. (2013). A novel method to monitor bicycling environments. *Transp. Res. Part A Policy Pract.*, 54,
335 1–13.

336 Kaufman, L., Rousseeuw, P. J. (1990). Partitioning around medoids (program pam). *Finding groups in data:
337 an introduction to cluster analysis*, 344, 68-125.

338 Kroesen M., Handy S. (2014). The relation between bicycle commuting and non-work cycling: results from a
339 mobility panel. *Transportation*. 41(3), 507–527. 10.1007/s11116-013-9491-4

340 Laranjeiro, P.F., Merchán, D., Godoy, L.A., Giannotti, M, Yoshizaki, H.T.Y., Winkenbach, M., Cunha, C.B.,
341 (2019). Using GPS data to explore speed patterns and temporal fluctuations in urban logistics: The case of
342 São Paulo, Brazil. *Journal of Transport Geography*, 76, 114-129, 10.1016/j.jtrangeo.2019.03.003

343 Lu, W., Scott D., M., Dalumpines, R. (2018). Understanding bike share cyclist route choice using GPS data:
344 Comparing dominant routes and shortest paths. *Journal of Transport Geography*.71, 172-181.
345 10.1016/j.jtrangeo.2018.07.012

346 Manum, B., Arnensen, P., Nordström, T., Gil, J. (2018). Improving GIS-based models for bicycling speed
347 estimations. *Eur. Transp. Conf.*

348 McLachlan G. J., Basford K. E. (1988). Mixture models: inference and applications to clustering. *Statistics,
349 textbooks and monographs*. 84.

350 McNicholas, P.D. (2016). Model-Based. *Journal of Classification*. 33, 331-373. 10.1007/s00357-016-9211-9

351 Menghini, G., Carrasco, N., Schüssler, N., Axhausen, K.W., (2010). Route choice of cyclists in Zurich. *Transp.
352 Res. Part A* 44, 754–765.

353 OpenStreetMap. (2020). <https://www.openstreetmap.org> Accessed 1 May 2020.

354 Pritchard R. (2018) 'Revealed preference methods for studying bicycle route choice—a systematic review',
355 *Int. J. Environ. Res. Public Health*, 15, 3.

356 Pritchard, R., Frøyen, Y.K., Bernhard, S., (2019). Bicycle Level of Service for Route Choice—A GIS Evaluation
357 of Four Existing Indicators with Empirical Data. *Int. J. Geo-Inf.* 8, 214. doi:10.3390/ijgi8050214

358 R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical
359 Computing, Vienna, Austria. [http://10.3390/ijerph1503047010.3390/ijerph15030470p://www.R-
360 project.org/](http://10.3390/ijerph1503047010.3390/ijerph15030470p://www.R-project.org/). Accessed 31 May 2020.

361 Rupi, F., Poliziani, C. Schweizer, J. (2019). Data-driven bicycle network analysis based on traditional counting
362 methods and GPS traces from smartphone. *ISPRS International Journal of Geo-Information*, 8,
363 10.3390/ijgi8080322

364 Rupi, F.; Schweizer, J. (2018), Evaluating cyclist patterns using GPS data from smartphones. *IET Intell. Trans.*
365 *Syst.*, 12, 279–285. 10.1049/iet-its.2017.0285

366 Rupi, F., Schweizer, J., Poliziani, C., (2020), Analysing the dynamic performances of a bicycle network with a
367 temporal analysis of GPS traces. *Case Studies on Transport Policy*. 10.1016/j.cstp.2020.05.007

368 Schweizer J., Sumopy. (2020). <https://sumo.dlr.de/docs/Contributed/SUMOPy.html> Accessed 1 May 2020.

369 Schweizer, J., Bernardi, S., Rupi, F., (2016), Map-matching algorithm applied to bicycle global positioning
370 system traces in Bologna, *IET Intell. Transp. Syst.*, 10, 244–250. 10.1049/iet-its.2015.0135

371 Schweizer, J, Rupi, F, Poliziani, C. (2020). Estimation of link-cost function for cyclists based on stochastic
372 optimization and GPS traces. *ITE Intelligent Transport Systems*, in press.

373 Strauss, J., Miranda-Moreno, L.F., (2017). Speed, travel time and delay for intersections and road segments
374 in the Montreal network using cyclist Smartphone GPS data. *Transp. Res. Part D: Trans. Environ.* 57, 155–
375 171. doi:10.1016/j.trd.2017.09.001.

376 Strauss J., Zangenehpour S., Miranda-Moreno L. F., Saunier N. (2017). Cyclist deceleration rate as surrogate
377 safety measure in Montreal using smartphone GPS data, *Accid. Anal. Prev.*, 99, 287–296.
378 10.1016/j.aap.2016.11.019

379 Tortora C., ElSherbiny A., Browne R.P., Franczak B.C., McNicholas P.D. (2019). MixGHD: Model Based
380 Clustering, *Classification and Discriminant Analysis Using the Mixture of Generalized Hyperbolic*
381 *Distributions*. <https://CRAN.R-project.org/package=MixGHD>, R package version 2.3.1. Accessed 31 May
382 2020.

383 Tortora, C., Franczak, B. C., Browne, R. P., McNicholas, P. D. (2019). A mixture of coalesced generalized
384 hyperbolic distributions. *Journal of Classification*, 36(1), 26-57. 10.1007/s00357-019-09319-3

385 Watkins K. E. and LeDantec C. (2016). Using Crowdsourcing to Prioritize Bicycle Route Network
386 Improvements. *Report FHWA-GA-16-1439'*

387 Wittwer, R., Zwangsmobilität, Verkehrsmittel. (2014). Orientierung junger Erwachsener: Eine
388 Typologisierung, Schriftenr. Dresden: Technische Universität Dresden Fakultät Verkehrswissenschaften
389 „Friedrich List“ Institut für Verkehrsplanung und Straßenverkehr Geschäftsführender Institutsdirektor: Univ.-
390 Prof. Dr.-Ing. Christian Lippold.

391 Zimmermann, M., Mai, T., Frejinger, E., (2017). Bike route choice modelling using GPS data without choice
392 sets of paths. *Transp. Res. Part C* 75, 183–196.

393