

1 Personalised Map Matched Imputation: imputing missing data from smartphone location
2 logs

3 Boaz Sobrado¹

4 ¹ Utrecht University

5 Author Note

6 Department of Methodology & Statistics

7 Submitted as a research report conforming to APA manuscript guidelines (6th edition).

8 Correspondence concerning this article should be addressed to Boaz Sobrado, . E-mail:

9 boaz@boazsobrado.com

10

Abstract

11 Personal mobility, or how people move in their environment, is associated with a vast range
12 of behavioural traits and outcomes, such as socioeconomic status, personality and mental
13 health. The widespread adoption of location-sensor equipped smartphones has generated a
14 wealth of objective personal mobility data. Nonetheless, smartphone collected personal
15 mobility data has remained underutilised in behavioural research, partly due to the practical
16 difficulties associated with obtaining the data and partly because of the methodological
17 complexity associated with analysing it. Recent changes in European regulation have made
18 it easier for researchers to obtain this data, but the methodological difficulties remain. The
19 difficulty lies in that smartphone location data is irregularly sampled, sparse and often
20 inaccurate. This results in a high proportion of missing data and significant noise. In this
21 paper we present a method called Personal Map Matched Imputation (PPMI) to deal with
22 missing data and noise in smartphone location logs. The main innovation of PPMI is that it
23 creates a personalised spatial map for each individual based on all the available data. In
24 doing so PPMI leverages the regularity of human mobility in order to smoothen noisy
25 measurements and impute missing data values. By simulating missing periods in real data we
26 find that a simple implementation of PPMI performs as well as existing methods for short (5
27 minute) missing intervals and substantially better for longer (1 day) missing intervals. When
28 imputing a subset of real missing data where travel logs are available as a reference points,
29 we find that PPMI performs substantially better than existing models.

30

Keywords: Missing Data, Measurement Bias, GPS, Human Mobility

31

Word count: 8319

32 Personalised Map Matched Imputation: imputing missing data from smartphone location
33 logs

34 **Introduction**

35 Why is human mobility important? Human mobility measures are quantified metrics of
36 how people move about in their environment. Human mobility affects a wide range of
37 outcomes, such as health, income and social capital (Goodchild & Janelle, 2010). The most
38 widely administered personality questionnaires ask individuals to what extent they agree
39 with statements such as “I love large parties”, “I prefer going to the movies to watching
40 videos at home” and “I love to travel to places that I have never been before” (Goldberg et
41 al., 2006). At their root these questions are mobility measures. In economic research, the
42 postal code of an individual’s home address is often used as a proxy for socioeconomic status
43 (e.g. Villanueva & Aggarwal, 2013). Economists are interested not only in where people live,
44 but also where they go to work: geographic labour mobility is related to the income of
45 individuals, the well-being of a country’s economy and even informs the policy of bodies such
46 as the European Comission (Tatsiramos, 2009). The extensive use of measures like these
47 across different domains in social science strongly suggests that mobility metrics are linked
48 to important real world outcomes. Perhaps it is unsurprising that behavioural researchers
49 have found that mobility measures can be used to predict academic performance (R. Wang,
50 Harari, Hao, Zhou, & Campbell, 2015), the incidence of obesity (Zenk, Schulz, &
51 Odoms-Young, 2009) or even the onset of a depressive episode in bipolar depression patients
52 (Palmius et al., 2017). Indeed, there is an argument to be made that perhaps psychologists
53 have not been studying mobility enough. When studying behavioural differences within
54 individuals behavioural scientists have often neglected the fact that individuals vary not only
55 over time but also over space. To fully understand behaviour we must understand how
56 behaviour can vary across environments.

57 Despite the importance of mobility measures, the majority of these metrics are
58 obtained through the use of questionnaires (such as the aforementioned personality

59 questionnaire) or specifically through pen-and-paper travel diaries. As a thought experiment
60 we encourage the reader to try to remember where exactly he or she was five weeks ago on
61 Saturday at midday. Most people find this sort of task very difficult (if not impossible),
62 which illustrates that questionnaires and travel diaries are have well-known methodological
63 flaws. This method of collecting data is difficult and relies on accurate self-reporting. Travel
64 diaries are burdensome to collect because participants must be explicitly asked to write down
65 on their movement patterns at frequent intervals due to human forgetfulness. This makes the
66 data collection expensive and limits the duration of data collection in practice. In addition,
67 the frequent reporting duties of the participant may bias the participants behaviour. The
68 limits of human cognition also limit the accuracy of self-reported measures. There is
69 evidence that participants are systematically biased when self-reporting mobility measures.
70 For instance, participants under-report the frequency of short trips (Wolf, Oliveira, &
71 Thompson, 2003) and underestimate the duration of regular commutes (Delclòs-Alió,
72 Marquet, & Miralles-Guasch, 2017). These obstacles can only be overcome by using
73 objective data on human mobility.

74 Social scientists now have the unprecedented opportunity to easily obtain objective
75 data on human mobility from smartphones. Before smartphones the only way to collect
76 objective data on human mobility involved giving participants an expensive
77 professional-grade location sensor and convincing them to take it with themselves at all
78 times. Barnett and Onnela (2016) points out that introducing a new device to the
79 participant's life may bias their behaviour. Moreover, collecting data in such a way is costly,
80 places a high burden on participants and therefore the logs do not span a long time (Barnett
81 & Onnela, 2016). Today millions of individuals carry smartphones with themselves every day
82 and do not need to be encouraged to do so by researchers. Smartphones are equipped with a
83 range of sensors that can be used to track the location of the device at all times. These
84 smartphones can collect and store hundreds of location measurements a day. For instance,
85 Google Location History contains movement information on millions of users, often spanning

86 years (Location History, 2017). Moreover, recent changes in EU regulations with regard to
87 consumer data-portability rights ensure that a willing participant should be able to easily
88 share this information with researchers at no cost to either the participant or the
89 researchers (Commission, 2017). Taken together, this means that researchers now have at
90 their disposal the ability to easily access the objective human mobility data of millions of
91 individuals spanning several years at little cost and without a significant burden of the
92 participants.

93 This paper wishes to achieve four objectives: First, we have argued that understanding
94 human mobility is important. Secondly, we argued that social scientists should leverage data
95 logs from smartphones to study human mobility, instead of relying on out-dated
96 pen-and-paper questionnaires. Now we will explore the practical difficulties in using
97 smartphone location logs. Finally we will introduce Personal Map Matched Imputation
98 (PPMI), a method for surmounting these difficulties. We will compare PPMI to existing
99 methods in the literature.

100

Background

101 Smartphone location measures are obtained primarily (but not exclusively) by Global
102 Positioning System (GPS) measurements. A GPS sensor uses the distance between a device
103 and several satellites to determine the location of the device. Although using a GPS sensor is
104 the most accurate way to establish location on a smartphone, the GPS sensor is also the
105 most energy consuming sensor on most smartphones (M. Y. Chen et al., 2006; LaMarca et
106 al., 2005). In order to avoid battery depletion and to overcome computational constraints
107 smartphones also use less-accurate heuristics such as WiFi access points and cellphone tower
108 triangulation. Smartphone location logs contain measurements from all of these sources,
109 usually in the form of time-stamped latitude, longitude and accuracy values. The accuracy *a*
110 of any given measurement is given in meters such that it represents the radius of a 67%
111 confidence circle (Location History, 2017). In other words, the true location of a device

112 should be within the radius a of the measurement 67% of the time.

113 Researchers often develop custom-made tracking applications which participants are
114 instructed to download on their phone. Alternatively participants are given a phone to use
115 for a given period of time with the custom-made tracking app pre-installed. We call location
116 logs resulting from these custom-made apps *custom logs*. The advantage of custom logs is
117 that the researchers can adjust tracking parameters, such as the frequency of
118 measurements and the sensor with which they are made. The disadvantage with this
119 approach is that researchers have to develop or adapt a custom-made tracking application
120 (which is not easy given hundreds of different types of smartphone models), distribute it
121 among research participants and enforce participation. Participants may dislike tracking
122 apps because they view them as more intrusive and these apps regularly drain the battery of
123 the device (G. M. Harari et al., 2016). Moreover, researchers have distribute this application
124 among research participants and convince them not to turn it off.

125 We focus on another solution, which is to take advantage of existing smartphone
126 location logs (*secondary logs*) . The advantages are clear: repositories such as Google's
127 Location History contains information on millions of users spanning years (Location History,
128 2017), participants can share the data by the click of a button and there can be no
129 behavioural changes due to participation in the study as the participant share past data.
130 The disadvantage is that researchers have no control over the tracking parameters, often
131 resulting in logs with sparse and inaccurate measurements. Hence, two important challenges
132 are dealing with missing data and measurement noise.

133 In order to work with secondary logs, researchers need to be able to handle the data
134 sparsity that leads to missing data. Missing data is a pervasive issue in secondary logs as it
135 can arise due to several reasons. Technical reasons include signal loss, battery failure and
136 device failure. Behavioural reasons include leaving the phone at home or switching the
137 device off. As a result, secondary logs often contain wide temporal gaps with no
138 measurements. For instance, several research groups studying mental health report missing

139 data rates between 30% to 50% (Grünerbl et al., 2015; Palmius et al., 2017; Saeb et al.,
140 2015). Other researchers report similar trends in different fields (e.g. G. M. Harari et al.,
141 2016; Jankowska, Schipperijn, & Kerr, 2015). In Figure 1 shows that despite the long
142 duration of the log the sparsity it is also evident.

143 There is no golden standard for dealing with missing data in GPS logs (Barnett &
144 Onnela, 2016). Importantly, spatiotemporal data measurements are auto-correlated in both
145 time and space. This means that best practices with other types of data, such as mean
146 imputation, are unsuitable. For example, imagine an individual who splits almost all her
147 time between work and home. Suppose she spends a small amount of time commuting
148 between the two along a circular path. Using mean imputation to estimate her missing
149 coordinates, we impute her to be at the midpoint between home and work, even though she
150 has never been there. Worryingly, there is little transparency on how researchers deal with
151 missing data (Jankowska et al., 2015).

152 Another methodological problem is related to the noise in the measurements that are
153 collected. The accuracy of smartphone location measurements is substantially lower than
154 that of professional GPS location trackers because smartphones often use less accurate
155 sensors. In professional GPS trackers less than 80% of measurements fall within 10 meters of
156 the true location. GPS measures are most inaccurate in dense urban locations and indoors
157 (S. Duncan et al., 2013; Schipperijn et al., 2014). Unfortunately for researchers, this is where
158 people in the developed world spend most of their time. Figure 2 shows how accuracy can
159 vary as a function of user behaviour, time and location. Most notably, low accuracy is often
160 (but not always) associated with movement (see Figure 3).

161 Noisy data can lead to inaccurate conclusions if it is not accounted for. Suppose we
162 wish to calculate an individual's movement in a day. A simple approach would be to
163 calculate the sum of the distance between each measurement. But if there is noise, the
164 coordinates will vary even though the individual is not moving. If the measurements are
165 frequent and noisy, we will calculate a lot of movement, even if the individual did not move

166 at all! This issue is also visualised in Figure 7. The problem is further complicated because
167 missing data and noisy measurements are related. Methods used by researchers to reduce
168 noise, such as throwing out inaccurate measurements (e.g. Palmius et al., 2017), can
169 exacerbate the severity of the missing data problem.

170 In this paper we will propose PPMI as a method for dealing with missing data and
171 measurement error in secondary location logs. We will compare PPMI to similar solutions in
172 the literature by evaluating the distance between points which were simulated as missing and
173 their imputed counterparts. Finally we will calculate time spent at home as a function of the
174 imputation method.

175 **Related Work**

176 How have researchers dealt with missing data in human mobility logs thus far?
177 Unfortunately there is no golden standard in how to deal with this type of missing data.
178 Researchers are generally vague about what practices they follow (Jankowska et al., 2015).
179 This vagueness is worrisome as it invites solutions which contain significant researcher
180 degrees of freedom (Simmons, Nelson, & Simonsohn, 2011). The vagueness is possibly also
181 due to the fact that most researchers are unfamiliar with possible solutions. Most researchers
182 simply down-sample temporally and remove missing observations or use some sort of
183 rule-based common sense imputations (e.g. Palmius et al. (2017)). The only principled
184 approach that we know of that aims to solve the issue of missing data in location logs as
185 they relate to human mobility is that of Barnett and Onnela (2016). We will explore the
186 methods of Barnett and Onnela (2016) and Palmius et al. (2017) in detail subsequently,
187 after introducing exploring other spatiotemporal methods.

188 A lack in methods for missing data imputation for human mobility patterns does not
189 imply there is not a vast literature on modelling movement. The most widespread models
190 are SSMs, therefore we shall detail a few examples and subsequently argue that they are
191 nonetheless unsuited for long term human mobility logs. Ecologists have used SSMs to

192 explain how animals interact with their environment (T. A. Patterson, Thomas, Wilcox,
193 Ovaskainen, & Matthiopoulos, 2008). These models can be quite complex. Preisler, Ager,
194 Johnson, and Kie (2004) uses Markovian movement processes to characterise the effect of
195 roads, food patches and streams on cyclical elk movements. The most well studied SSM is
196 the Kalman filter, which is the optimal algorithm for inferring linear Gaussian systems. The
197 extended Kalman filter is the de facto standard for GPS navigation (Z. Chen & Brown,
198 2013). The advantage of state space models is that they are flexible, deal with measurement
199 inaccuracy, include information from different sources and can be used in real time.

200 For secondary logs the main limitation of SSM implementations is that they ignore
201 movement routines. For instance, humans tend to go to work on weekdays and sleep at night.
202 Because SSMs are based on the Markov property, they cannot incorporate this information.
203 In other words, the estimated location $G(t)$ at time-point t is often based only upon
204 measurements D_t , D_{t-1} and ignores all $D_{t-i}|i \geq 2$. Hierarchical structuring and conditioning
205 on a larger context have been suggested as ways to add periodicity to Markovian models.
206 These solutions are often computationally intractable or unfeasible (Sadilek & Krumm,
207 2016). Moreover, these models often assume time and space invariance (location is not a
208 direct function of time or space). These mathematical assumptions are violated in the case
209 of human movement patterns. For this reason we do not consider existing SSMs to be useful
210 for imputing missing data in this case.

211 In the wider realm of spatiotemporal statistics there are numerous missing data
212 imputation methods. These often come from climate or geological research and rely on
213 spatiotemporal auto-correlations. For instance, Feng, Nowak, O'Neill, and Welsh (2014)
214 estimate missing values by incorporating similar observed temporal information from the
215 value's nearest spatial neighbors. The authors illustrate their example using rainfall data
216 from gauging stations across Australia. Similarly, Z. Zhang et al. (2017) use a variety of
217 machine learning methods to impute missing values. The example provided relates to
218 underground water data. Generally these models assume fixed measurement stations (such

219 as rainfall gauging stations). For this reason they cannot be easily applied to missing
220 mobility tracks without significant pre-processing.

221 On the other hand, a few researchers have explicitly attempted to impute missing data
222 from human mobility patterns (Barnett & Onnela, 2016; Palmius et al., 2017 ; @ Wu, Fu,
223 Wang, Xiao, & Fu, 2017). Importantly, none of them worked with secondary logs.
224 Nonetheless we will detail what they did as informative examples. Palmius et al. (2017) deal
225 with the measurement inaccuracy of D in custom logs by removing from the data set all
226 unique low-accuracy a data points that had $\frac{d}{dt}D > 100 \frac{km}{h}$. Subsequently the researchers
227 down sample the data to a sample rate of 12 per hour using a median filter. Moreover,
228 Palmius et al. (2017) explain:

229 “If the standard deviation of $[D]$ in both latitude and longitude within a 1 h
230 epoch was less than 0.01 km, then all samples within the hour were set to the
231 mean value of the recorded data, otherwise a 5 min median filter window was
232 applied to the recorded latitude and longitude in the epoch”.

233 Missing data was imputed using the mean of measurements close in time if the
234 participant was recorded within 500m of either end of a missing section and the missing
235 section had a length of $\leq 2h$ or $\leq 12h$ after 9pm. In cases where the previous conditions are
236 not met no values are imputed.

237 Barnett and Onnela (2016) follow a different approach which is, to the best of our
238 knowledge, the only principled approach to dealing with missing data in human mobility
239 data. Barnett and Onnela (2016) work with custom logs where location is measured for 2
240 minutes and subsequently not measured for 10 minutes. In the words of the authors, Barnett
241 and Onnela (2016) handle missing data by first converting data to mobility traces, which are
242 defined as a sequence of flights and pauses. Flights are segments of linear movements and
243 pauses corresponding to periods of time where a person does not move. Subsequently, the
244 authors impute missing data by:

245 “simulat[ing] flights and pauses over the period of missingness where the direction,
246 duration, and spatial length of each flight, the fraction of flights versus the
247 fraction of pauses, and the duration of pauses are sampled from observed data.”

248 This method can be extended to imputing the data based on temporally, spatially or
249 periodically close flights and pauses. In other words, for a given missing period, the
250 individual’s mobility can be estimated based on measured movements in that area, at that
251 point in time or movements in the last 24 hours.

252 It is also worth mentioning that there is a body of work that incorporates road features
253 and other background information to generate predictions or imputations. For instance, Liao,
254 Patterson, Fox, and Kautz (2007) use hierarchical Markov models in combination with
255 knowledge about the transportation system to bridge the gap between raw GPS sensor
256 measurements and high level information such as a user’s destination and mode of
257 transportation. This can be then used to impute or predict missing steps. Along similar lines
258 Wu et al. (2017) use what they call a Spatial Temporal Semantic Neural Network (STS-NN)
259 to predict future human movement. While the authors are concerned with prediction and
260 not imputation, they devised a method called the Spatial Temporal Semantic (STS)
261 algorithm which converts raw measurements to machine learning friendly discrete bins.

262 Working with high-frequency measurements, Wu et al. (2017)’s method down-samples the
263 raw data temporally and map-matches the resulting bins to discrete points along
264 pre-established geographical features such as roads and highways. This minimises
265 measurement error and paves the way for applying machine learning methods to human
266 mobility problems. Subsequently we will focus on methods which do not incorporate external
267 information as they are more generally applicable.

268 In this section we have argued that there is a lack of established practices to follow
269 with respect to missing data in human mobility logs. Moreover, extensively used
270 spatiotemporal methods, such as state space models (SSMs), are not well suited to deal with
271 human mobility patterns in secondary logs. Finally we discussed in detail three approaches

272 which deal explicitly with mobility patterns from custom or secondary logs (Barnett &
273 Onnela, 2016; Palmius et al., 2017 ; Wu et al., 2017).

274 **Methodology**

275 **Notation**

276 Location measurements, such as those produced by GPS sensors, provide us with
277 coordinates (latitude and longitude) on the surface of the earth, which is ellipsoid shaped.
278 Projecting three dimensional measurements in \mathbb{R}^3 onto a two dimensional plane in \mathbb{R}^2 results
279 in distortion. For clarity, when we use the term distance we refer to the geodesic distances
280 on an ellipsoid using the WGS84 ellipsoid parameters.

281 Subsequently let us simplify by assuming that a persons location is on two-dimensional
282 Euclidean plane. Let a person's true location on this two-dimensional plane be
283 $G(t) = [G_x(t)G_y(t)]$ where $G_x(t)$ and $G_y(t)$ denote the location of the individual at time t on
284 the x-axis and y-axis respectively. For simplicity, we can assume that the x-axis is the
285 longitude and the y-axis is the latitude. Moreover, let $D \in \mathbb{R}^2$ be the recorded data
286 containing the latitude and longitude. In addition, let a denote the estimated accuracy of
287 the recorded data. Furthermore, $G(t)$, D and a are indexed by time labeled by the countable
288 set $t = t_1 < \dots < t_{n+1}$. The measure of accuracy a_t is given in meters such that it represents
289 the radius of a 67% confidence circle. If $D_t = \emptyset$ it is considered *missing* and it is not missing
290 otherwise.

291 When several data sets are available from individuals living in overlapping areas we
292 can construct a $t \times i$ matrix M where the entry $M(t, i)$ contains $G(t)$ for the individual i .

293 **Personalised Map Matching Imputation**

294 Our algorithm is designed to leverage the periodic nature of human movement along
295 with the long span of secondary to deal with measurement sparsity and inaccuracy.

296 **Modelling assumptions.** First, following Barnett and Onnela (2016) we categorise

297 all time-points t as either belonging to the set P (pause) or set F (flight). Conceptually
298 pauses can be understood as periods of time where an individual spends significant amount
299 of continuous time without moving. Flights are the times where the individual is moving.

300 Let t_a be a pause of length n .

$$t_a = t_i < \dots < t_{i+n}$$

301 Let t_b be a pause of length m such that there is no temporal overlap between t_a and t_b :

$$t_b = t_j = < \dots < t_{j+m} | t_{i+n} < t_j$$

302 Then it follows that between the two pauses there must be a flight indexed by t_x of length
303 $j - i + n$.

$$t_x = t_{i+n} < \dots < t_j | t_x \in F$$

304 We define pause locations $G(t_a), G(t_b) | t_a, t_b \in P$ as locations where an individual spends an
305 extended amount of time in the same space (e.g. school, home, work, train station, barber
306 shop, bar, gym). Importantly, our model assumes period and cyclic human movement such
307 that there are many pauses $t_{a1}, t_{a2}, \dots, t_{an}$ such that $G(t_{a1}) = G(t_{a2}) = \dots = G(t_{an})$.

308 Moreover, it is possible for $G(t_a) = G(t_b)$ such that $t_a \neq t_b$. For example, if the individual
309 leaves home for a run and returns home without stopping anywhere else.

310 Let us define as $Flight_{ab}^x$ the set of all points belonging to a flight between $G(t_a)$ and
311 $G(t_b)$ at time-point t_x .

$$Flight_{ab}^x = G(t_x) | t_x \in F = \{G(t_{i+n}), \dots, G(t_j)\}$$

312 Again, there are many flights $t_{x1}, t_{x2}, \dots, t_{xn}$ such that $Flight_{ab}^{x1} = Flight_{ab}^{x2} = \dots = Flight_{ab}^{xn}$.

313 Then, we can define as $Path_{ab}$ the set of all flights between $G(t_a)$ and $G(t_b)$ at all
314 time-points. For simplicity, we assume that $Path_{ab} = Path_{ba}$.

315 In addition, we consider all measurements $D(t)$ to be imperfect measurements of $G(t)$:

$$G(t) = D(t) + \text{Measurement Error}$$

316 **Personalised Map Matching Imputation algorithm**

317 Our algorithm performs the following steps:

- 318 1. *Map building*: Extract from measurements D all pause location bins and path location
- 319 bins to create a personalised map.
- 320 2. *Binning*: Assign each measurement D to a unique discrete location bin.
- 321 3. *Imputing*: Use a classification method to predict missing measurements based on all
- 322 the available information.

323 **Map building.** Following Wu et al. (2017)'s spatial-temporal-semantic (STS) feature extraction algorithm our aim is to transform pause and path locations into machine learning friendly discrete location sequences. There are multiple ways of extracting such measurement clusters in the literature, such as Spatio-Temporal Density-Based Spatial Clustering of Applications with Noise (ST-DBSCAN) and sequence oriented clustering (SOC) ("ST-DBSCAN," 2007). We will focus on two methods which explicitly work with mobility patterns from unevenly sampled smartphone logs (Barnett & Onnela, 2016; Palmius et al., 2017). Both of these methods pre-process the data and subsequently use two steps to extract pause locations: first they extract pauses and their corresponding locations, then they cluster pause locations based on spatial proximity. This implementation of PMMI uses a stricter version of Barnett and Onnela (2016)'s approach to extract pauses.

334 First the measurements D are filtered such that only measurements with an accuracy value lower than $a_{P \lim}$ remain within the sample. Then, a measurement D_t belongs to a pause if and only if:

- 337 1. The next measurement D_{t+1} is within $t_{\text{Pause lim}}$ amount of seconds (so it is not missing)
- 338 2. The next measurement D_{t+1} is within $d_{\text{Pause lim}}$ meters.
- 339 3. The duration of the pause is more than $\delta_{\text{Pause lim}}$ seconds.
- 340 4. Let the measurements of a possible pause which fit the aforementioned criteria be $D_{t,t+1,\dots,t+n}$. These points are only a pause if the distance between the mean

342 coordinates of $D_{t,t+1,\dots,t+n}$ and the furthest away points of $D_{t,t+1,\dots,t+n}$ is within 2 times
 343 the mean accuracy a of $D_{t,t+1,\dots,t+n}$.

344 This set of points were then hierarchically clustered using a distance matrix, such that
 345 all points within d meters of each other were clustered into a pause location. Each pause
 346 location is a bin.

347 For all remaining measurements we assume that they belong to paths. In this

348 implementation of PMMI we use the following algorithm to estimate paths:

- 349 1. Take all measurements which are not pauses, filter them based on an accuracy
 350 threshold $a_{\text{Path Lim}}$.
- 351 2. Create a distance matrix for all remaining measurements $D_t \in F$ and hierarchically
 352 cluster it accordingly, such that all points within $d_{\text{Path Lim}}$ meters of each other are
 353 clustered into a single pause point.

354 At this point all empirically observed path bins and pause bins are extracted. However,
 355 there may be some overlap between pause bins and path bins. Thus, the bins are clustered
 356 again, such that the pause bins retain priority. This means that if a pause bin and a path
 357 bin are within less than d meters of each other, the path bin is removed. The reasoning for
 358 this is that the threshold for not being in a pause cluster should be higher, as individuals
 359 spend the majority of time at a pause cluster. The end result is a discretised map which
 360 contains pause and flight bins based on the entire log history of the individual.

361 **Binning.** Wu et al. (2017)'s spatial-temporal-semantic (STS) feature extraction
 362 algorithm uses map matching as a ground truth to assign noisy measurements into discrete
 363 bins along roads. In other words, in addition to the measured data they also use a
 364 geographic database that contains information about the area in which the individual is
 365 (e.g. where precisely the roads are), and sort measurements into bins based on both the
 366 measurement and the geographic data base. For example, if an individual is measured as
 367 moving closely in parallel to a road A in an area where there is no other parallel road, Wu et
 368 al. (2017)'s method will assume that the individual is on the road A.

369 PMMI uses a similar logic, but without using any external geographic database. The

370 key modification in PMMI is that whilst Wu et al. (2017) uses a map from outside the

371 persons location logs, we use the total location history of the individual to create a

372 personalised map. This map is subsequently used to bin measurements. This is feasible for

373 two reasons: humans tend to have repetitive movement habits and secondary logs tend to be

374 long. To put it in simpler terms, we consider each measurement at D_x as a sample of $Path_{ab}$,

375 and by aggregating many measurements we can use them to map out $Path_{ab}$.

376 Thus, all measurements D were assigned to a discrete bin on the personal

377 map. This includes previous measurements which were discarded from the map building

378 exercise due to an accuracy a value which exceeded $a_{P\lim}$ or $a_{F\lim}$. In this implementation

379 we used a simple assigning function, whereby the measurements were assigned to the bin

380 nearest to the measurement.

381 **Classification.** At this point the objective of PMMI is to take all the information

382 available about the mobility history of an individual and impute the missing value. In this

383 implementation we trained an artificial neural network (ANN) to do so. For more

384 information on the precise architecture of the artificial neural network please consult the

385 appendix. The input variables to the ANN are:

386 1. The previous and subsequent observed bin as a binary class matrix.

387 2. The distance in time to the next & previous bin.

388 3. The time of the day encoded as a cyclical two-dimensional feature.

389 4. The day of the month as a binary class matrix.

390 5. The month of the year as a binary class matrix.

391 For the encoding of the time of day we took the cosine and the sine transforms of the

392 amount of seconds that have elapsed after midnight (London, 2016). This is necessary so

393 that the model can understand that one second past midnight and one second before

394 midnight are in fact two seconds away from each other. Moreover we scaled the non-binary

395 values to occupy a range between 0 and 1 in order to ensure convergence.

396 For a missing time-point at $D_t \in \emptyset$, the output of the model is a set of probability
 397 estimates associated with every location cluster. That is, for each missing time-point the
 398 model returns a vector of probability estimates (with one estimate per bin) associated with
 399 where the individual is.

400 Datasets & Analyses

401 The secondary location log used to train the imputation methods was collected
 402 between 2013 and 2017 on different Android devices from a single individual. About 54% of
 403 the data is missing for the entire duration of the log. This may be misleading as there are
 404 several long periods with no measurements whatsoever. For days which were not entirely
 405 missing, approximately 22% of all five minute segments were missing. The structure of
 406 missingness of a day with measurements is shown in Figure 4. As you can see, there are
 407 several long periods over the course of the log for which there are no measurements. The
 408 median sampling frequency per day for non-missing days is around 0.006 Hz.

409 For simplicity, we subsequently used a time period when the individual was living in
 410 the Netherlands. This subset contains 156,000 measurements over a period of less than six
 411 months.

412 Palmius et al. (2017)'s algorithm was implemented in R based on the original
 413 MATLAB code and pseudocode kindly provided by the author. Barnett and Onnela (2016)'s
 414 method was slightly adapted in R to fit the the data structure from the original R code
 415 provided by the author. When executing their models we used the same parameters as the
 416 authors did. To represent Barnett and Onnela (2016)'s model we used the variant where
 417 movements were sampled based on spatial proximity.

418 All analyses were performed using R (Version 3.4.3; R Core Team, 2017)¹.

¹We, furthermore, used the R-packages *bindrcpp* (Version 0.2; Müller, 2017), *dbplyr* (Version 1.2.0; Wickham & Ruiz, 2018), *dplyr* (Version 0.7.4; Wickham, Francois, Henry, & Müller, 2017), *geosphere* (Version 1.5.7; Hijmans, 2017a), *ggplot2* (Version 2.2.1; Wickham, 2009), *ggthemes* (Version 3.4.0; Arnold, 2017), *kableExtra* (Version 0.7.0; Zhu, 2018), *keras* (Version 2.1.4; Allaire & Chollet, 2018), *knitr* (Version 1.20; Xie,

419 All the code is available on a public repository (Sobrado, 2018).

420 **Results & Evaluation Metrics**

421 The results will consist of multiple steps:

- 422 1. Evaluating the performance of the map building and assigning functions.
- 423 2. Comparing the performance of PMMI compared to alternative methods [(Barnett &
- 424 Onnela, 2016; Palmius et al., 2017)] and a baseline model using cross-validation.
- 425 3. Comparing the performance of PPMI compared to the aformentioned methods using
- 426 objective ground-truth data (public transportation time-stamps locations).

427 **Map building & binning evaluation.** Before we can evaluate the accuracy of the
 428 imputations, it is essential to evaluate how well noise in the data has been smoothed. All
 429 methods smoothen noise differently in an attempt to extract true location $G(t)$ from
 430 measurements D_t . In the subsequent cross-validation step we will remove random
 431 smoothened blocks and compare the imputed value to the smoothened value. For this reason
 432 we want to ensure that the smoothened method actually does reflect the true location $G(t)$.
 433 Otherwise we run the risk of over-estimating the accuracy of a model in the sense that we
 434 will compute the extent to which an imputation method can correctly impute a smoothened
 435 value which has little to do with the true location $G(t)$.

436 In order to evaluate the map building and binning we will first visually evaluate the
 437 paths and pause locations. A visual evaluation of paths superimposed on is an established
 438 way to heuristically check their accuracy (e.g. Brunsdon, 2007). Then, let the average

2015), *leaflet* (Version 2.0.0; Cheng, Karambelkar, & Xie, n.d.), *padr* (Version 0.4.0; Thoen, 2017), *papaja* (Version 0.1.0.9709; Aust & Barth, 2018), *raster* (Version 2.6.7; Hijmans, 2017b), *RColorBrewer* (Version 1.1.2; Neuwirth, 2014), *readr* (Version 1.1.1; Wickham, Hester, & Francois, 2017), *rgdal* (Version 1.2.16; R. Bivand, Keitt, & Rowlingson, 2017), *scales* (Version 0.5.0; Wickham, 2017), *sp* (Version 1.2.7; Hijmans, 2017a; Pebesma & Bivand, 2005), *tibbletime* (Version 0.1.0; Vaughan & Dancho, 2018), and *tidyverse* (Version 0.8.0; Wickham & Henry, 2018).

439 distance between the actual measured point and the binned point be the *deviation distance*
440 δ_{dev} . With respect to the deviation distance δ_{dev} , we expect:

- 441 1. A positive relationship between the deviation distance and the accuracy of each
442 measurement.
- 443 2. Roughly 67% of the deviation distances δ_{dev} are within accuracy a of each
444 measurement.

445 **Imputation algorithm performance in cross-validation.** We will compare the
446 performance of PMMI, Palmius et al. (2017) and Barnett and Onnela (2016). In addition,
447 we will also compute two baseline models. The first one is a naive model, which simply
448 imputes the previous observed value. The second one we will call the “home” model, which
449 imputes that the individual is at home at all missing timepoints. Using these two models as
450 baselines makes sense given the high degree of autocorrelation in the data and that
451 individuals tend to spend most of their time at home. To compare the performance of the
452 aforementioned methods we will remove 25% of measured time intervals at random within a
453 four week period. We will make our comparisons for intervals of 5 minutes, 1 hour and 1 day.
454 In other words, we will remove 25% of time intervals at random, while varying the duration
455 of the time intervals removed.

456 For the Barnett and Onnela (2016) and PMMI models we will use all the available
457 data to train the models with the exception of the time periods being investigated. Palmius
458 et al. (2017)’s model does not require training.

459 To compare all methods with each other we will compute a distance measure (how far
460 was the removed location from the predicted location) and the coverage. The coverage refers
461 to the percentage of missing cases which were imputed. For PMMI’s imputations we will use
462 a weighted mean for the distance measures whereby each 5 minute period is weighted equally.
463 This is necessary because the other two models downsample temporally to 5 minutes and
464 because the frequency of measurements in a 5 minute period is correlated with imprecise

465 measurements and visiting infrequently travelled locations. Hence weights must be used to
466 avoid inflating error rates for PMMI in comparison to the other models.

467 In addition, other measures of interest for PMMI are *accuracy* (in what percentage of
468 the cases was the appropriate cluster predicted), the *confidence* and the *distance expectation*.
469 The confidence is the probability with which the model predicts the most likely cluster. For
470 instance, if the model predicts the missing bin to be bin A with a probability of 0.9 then we
471 can say this is a high confidence prediction. Similarly, the distance expectation is the
472 cross-product of the estimated probabilities that the individual is at any of all given clusters
473 and the distances between the clusters to the true cluster. For example, if the true location
474 of an individual is bin A (bin A is 1000 meters away from bin B) and the model assigns a
475 probability of 0.9 at bin A and 0.1 at bin B, then the distance expectation would be 100
476 meters.

477 **Imputation algorithm performance using external benchmarks.** We will
478 compare the performance of PMMI, Palmius et al. (2017), Barnett and Onnela (2016) and
479 the naive baseline model with regards to imputing real missing measurements. To do will use
480 information from the Dutch public transportation card. The Dutch public transportation
481 service provides users with time-stamped location data of when and where they board,
482 change lines or leave public transportation. In this paper we used 97 such events in a single
483 month.

484 To be able to make a comparison between models we will remove all measurements
485 from within the 5 minute period that a time-stamped measurement is available. Then we
486 will use each model to impute the location of the individual within that period and compare
487 it to the true location.

488

Results

489 **Map building & binning evaluation**

490 We used the following parameters to extract pauses: an accuracy limit $a_{\text{Pause lim}}$ of 250
491 meters, a time limit $t_{\text{Pause lim}}$ of 300 seconds, a distance limit of $d_{\text{Pause lim}}$ 50 meters and a
492 minimum pause duration limit $\delta_{\text{Pause lim}}$ of 100 seconds. Moreover, to extract path clusters
493 we used the parameters: $a_{\text{Path Lim}} = 150$ meters and $d_{\text{Path Lim}} = 300$ meters.

494 The model was tuned based on a visual inspection of the resulting personal map
495 meeting the following criteria:

- 496 1. Plausible extraction of pause locations. For example, a pause in a field with no path
497 leading to it is less plausible than a pause at a train station, airport, or office building.
- 498 2. A suitable extraction of paths. For instance, are ensuring that individual paths are
499 complete and not merged unnecessarily in the case of closely parallel paths.
- 500 3. Less than a fifth of extracted clusters are assigned only one measurement. Unique
501 clusters are unavoidable as people sometimes visit a location briefly and only once,
502 nonetheless they pose problems for the classification algorithms.

503 When selecting parameters there is a trade off is between bias and precision. This is
504 because an increase in precision (in the form of a higher resolution of locations location)
505 comes at the expense of precision (assigning measurements to bins becomes more difficult).
506 For instance, by increasing the clustering distance parameter $d_{\text{Pause lim}}$ we can extract more
507 valid pause locations at the risk of falsely categorising certain measurements to the wrong
508 bin. This is illustrated in Figure 5.

509 Map building results in a personalised map with pause and path clusters. An excerpt
510 can be seen in Figure 6. It is important to remind the reader that PMMI is map agnostic and
511 uses no information from the map. Therefore, the close overlap with features on the map,
512 such as pause bins at relevant buildings and transportation clusters, as well as the flight bins
513 following roads and railway lines indicate a high degree of precision in personal map building.

514 As expected, PMMI’s path extraction yields greater accuracy for frequently occurring paths.
 515 For example, the frequently travelled Amsterdam-Utrecht railway line has been extracted
 516 almost perfectly, while the less frequently travelled Utrecht-Enschede line is far sparser.

517 For the entire period examined period the we find a deviance of 40 meters and a
 518 median deviance of 15 meters. Around 69% of the deviance values are within their
 519 corresponding accuracy value, which is close to the theoretical 67% value that is expected.
 520 Approximately 9% of values were not taken into account when creating the bins,given that
 521 their the accuracy a exceeded $a_{\text{Path Lim}} = 300$.

522 In comparison, Palmius et al. (2017)“s method has a median deviance δ_{dev} of 3 meters,
 523 with a mean of 115 due to high deviance outliers. On the other hand, Barnett and Onnela
 524 (2016)’s method has mean deviance δ_{dev} of 343 meters and a median deviance of 8 meters.
 525 Barnett and Onnela (2016)’s deviance is necessarily higher than Palmius et al. (2017)”, as
 526 they down-sample temporally (like Palmius et al. (2017)) and subsequently aggregate into
 527 pauses and linear flights.

528 The key difference between temporal and spatial down-sampling is shown in Figure ??.
 529 Temporal down-sampling is much more sensitive to noise in sparsely measured periods
 530 because it averages out values within five minute periods. Often there are only a few noisy
 531 measurements in those periods (see Figure 1), which leads to a noise in the down-sampled
 532 values. Unsurprisingly, there is a positive relationship between deviance and the amount of
 533 measurements in each down-sampled interval. The fact that over 90% of deviance values are
 534 within accuracy (substantially higher than the expected theoretical 67%) confirms that
 535 temporal down-sampling is not sufficiently filtering out the noise.

536 **Imputation algorithm performance evaluation in cross-validation.** The first
 537 factor we must take into account is the coverage of each imputation model. The results show
 538 us that while PPMI and the naive baseline model impute all missing values, the other two
 539 models have difficulties doing so (Table 1). In particular, the models of Palmius et al. (2017)
 540 and Barnett and Onnela (2016)’s fail to impute an increasing amount of missing values with

541 an increase in the duration of the missing period. Palmius et al. (2017)'s method in
542 particular failed to impute a single missing time period. The longer a log, the higher the
543 probability that long time periods will be missing. Secondary logs are typically long. Hence
544 missing data imputation methods designed to work with secondary logs must be capable of
545 imputing longer time periods.

546 With respect to the distance metrics, Table 1 shows the results of the distance metrics
547 for each method. We can see that PPMI predicted the location perfectly for over half of the
548 removed values in all three time categories. To be precise PPMI's prediction accuracy was
549 88%, 73% and 47% for the 5 minute, 1 hour and 1 day periods respectively. As a comparison,
550 the baseline naive model's prediction accuracy was 87%, 68% and 22%, and the home
551 baseline model's were 46%, 52% and 40%. Although the home baseline model and PPMI
552 achieved similar distance metrics in the 1 day condition, this does not mean that PPMI
553 defaulted to simply predicting the individual is home at all times (only 40% of predictions
554 were home). PPMI outperformed the alternative models when taking into account the
555 median scores, but not the mean. The explanation is that the alternative models failed to
556 impute more difficult time periods.

557 On the other hand, the confidence score for each time period were heavily skewed left,
558 withh a median of 1. Accordingly, the distance expectation (the crossproduct of the model
559 probabilities of an estimate with the distances of the true value to the estimates) values were
560 quite similar to the distance scores, albeit approximately 5% higher. Based on this we can
561 say that this implementation of PPMI was overconfident in its predictions.

562 **Imputation algorithm performance evaluation using external benchmarks.**

563 Once we removed the measurements within the 5 minute period of each time-stamp in the
564 public transportation log, we employed each of the aforementioned methods to impute the
565 location of the individual at the time-stamped period. Then we calculated the distance
566 metrics between the true location (based on the public transport log) and the imputed
567 location. We did not compute the home baseline model as it would have always been wrong

568 for this subset of data (travelling and being at home are mutually exclusive).

569 First, the median distance between the imputed value and the true value for the
570 baseline naive model was 555 meters (with a mean of 3303). For this method, the prediction
571 accuracy was 16%. This means that the binning method places the individual at the location
572 reported by the public transportation travel log in the five minutes prior to the measurement
573 16% of the time. Needless to say the coverage was 100% as all missing values were imputed.

574 For the PPMI the median distance metric was 1037 meters (mean 5637). The
575 prediction accuracy was higher than that of the baseline naive model at 24%. Although
576 PPMI imputed more points correctly than the naive baseline model, those it imputed
577 incorrectly were significantly more incorrect. All missing values were imputed, hence
578 resulting in a coverage of 100%.

579 The same could not be said about Palmius et al. (2017)'s method, as it failed to
580 impute almost 40% of the time periods. This results in a coverage of 60%. For the time
581 periods the method did impute, there was a median distance of 1617 meters between the
582 true location and the imputed value (with a mean of 1517 meters).

583 Barnett and Onnela (2016)'s method performed better in terms of coverage at 81%.
584 The median distance metric was 1342 meters and the mean value was 5506 meters.

585 **Example: effect on aggregate measures.** Social scientists are most interested in
586 aggregating spatiotemporal data to more socially relevant metrics, such as the amount of
587 time spent at home, frequency of travels to new locations, the amount of trips made or
588 distance covered. As an example we calculated the time spent at home of the user without
589 imputing any missing data and with all three of the investigated methods. The results can
590 be seen in Figure 8. Without any form of missing imputation the individual's time spent at
591 home is unknown for 12% of the time, which adds up to almost four days over the course of
592 a 30 day month.

593 The amount of time spent at home has been found to be a reliable predictor of
594 extroversion (G. M. Harari et al., 2016) and the onset of depressive episodes in bipolar

595 patients (Palmius et al., 2017). The 12% of time that remains unaccounted for could mean
596 the difference between an introvert and an extrovert, or a healthy individual and an
597 individual suffering of depression. Filling in these gaps is exactly the sort of problem
598 that can be solved using missing data imputation. Interestingly, despite their radically
599 different imputation methods, all three examined models suggest that the user spent
600 approximately 60% of their time at home.

601 Discussion & Conclusion

602 Overall the PMMI performed better than the alternative models, particularly during
603 longer missing periods and with objective data gathered from the Dutch public transport
604 service. However, PPMI did not perform substantially better than the baseline models for all
605 time periods and with the objective Dutch public transport data, suggesting that the
606 classification method can be improved. In addition, the comparison to the performance of
607 the Barnett and Onnela (2016) and Palmius et al. (2017) models is somewhat unfair, as they
608 were created for custom logs, not secondary logs. Nonetheless, the comparison remains valid
609 as they are the closest we found to a missing data imputation methods in smartphone GPS
610 logs.

611 In addition to higher accuracy under the conditions typical of secondary logs, the
612 advantages of PPMI are increased coverage and flexibility for missing data imputation,
613 robustness to irregular sampling, the ability to model complex non-linear interactions in its
614 imputations, and the ability to use historical records to smooth movement noise.

615 PPMI's increased coverage and flexibility comes from its ability to make complex
616 non-linear predictions. For instance, in a given missing period it might make sense to predict
617 that the individual is either at home, or at the office, or at a shop with equal probability.
618 While PPMI can make such an imputation, none of the alternative methods can do this.
619 Moreover, the ability to take the prediction probability values from the neural network also
620 helps in dealing with uncertainty. A known-drawback of single imputation is that it takes an

621 imputed value and treats it as observed. Simple rule-based methods such as Palmius et al.
622 (2017)' are essentially algorithmic single imputation methods. With PPMI it is possible to
623 model uncertainty using the predicted probabilities of each estimate. For instance, in the
624 previous example, we could choose to only take estimates with a high degree of confidence,
625 thus creating confidence intervals by adding and subtracting the amount of cases where the
626 location of the individual is ambiguous.

627 With respect to irregular sampling, alternative methods use temporally based
628 down-sampling in order to reduce noise. This leads to deterioration in resolution not only
629 over space, but also over time. A combination of irregular sampling with the fluctuating
630 accuracy values can lead to nonsensical results. For instance, consider a case where there are
631 two inaccurate measurements in movement at 12:00:01 and 12:04:59. Down-sampling over 5
632 minute periods will lead to a value that will be the mean of the two inaccurate samples,
633 which is likely to be a location the individual is certainly not. PPMI instead down-samples
634 spatially, which ensures that the binned location is one which is composed of the mean of
635 hundreds of observations, not just the few that happen within a single period.

636 While both Barnett and Onnela (2016) and Palmius et al. (2017) use historical data to
637 smoothen pause locations by clustering pause locations with a close degree of spatial
638 proximity, neither of them do the same for non-pause locations. This may be feasible with
639 high frequency, regularly sampled short duration logs but creates noise with secondary logs.
640 Moreover, with secondary logs it is feasible to spatially “average out” multiple samples of the
641 same path in order to recreate it in its entirety. For instance, although the mean sampling
642 frequency during train travels on the Amsterdam-Utrecht line is low (about 0.01 Hz) the
643 personalised map manages to recreate the train line almost perfectly, despite being
644 completely map agnostic. Map agnisticism makes the model more flexible as it can be used
645 for areas where no good geographical databases are available.

646 There are multiple methodological limitations in this paper. Most importantly, the
647 evaluation methods are imperfect. The golden standard would be to use at least one highly

648 accurate professional grade GPS device with high sampling frequency to compare our data
649 to. Until that is available, the use of public transport data and cross-validation is just a
650 substitute.

651 Furthermore, PPMI can be further developed. The map building function, the
652 assignment function and the classification model remain simplistic and could be improved.

653 In map building, the probability of a pause at a given location is certainly related to
654 other factors, such as the time of the day as well as the prior history of pauses at that
655 location. These factors are not taken into account in the pause extraction function.
656 Improved methods would do well to do so. As for paths, a drawback of the current method
657 is that the density of the clusters is a function of the clustering parameter d , the distance
658 between the observed points and their sampling density. It does not take into account the
659 length of the path as well as the average sampling frequency of the path. This is an issue
660 because it can lead to bins to which data points are seldom assigned. For example, while the
661 Amsterdam-Utrecht line has been mapped out almost perfectly, many of the clusters along
662 the route have only been assigned few measurements. This leads to difficulties in the
663 classification part of the model, as infrequently observed clusters are hard for the model to
664 predict.

665 The current assignment function simply assigns each measurement to the nearest bin.
666 It does not take into account any contextual information that can be gleaned from the entire
667 movement history of the information, such as what path they are on. For instance, assume
668 that it is known that an individual is travelling from point A to point B along path AB, and
669 there is an inaccurate measurement closest to a cluster which belongs to path AC. By only
670 taking distance into account, the measurement can get assigned to the wrong cluster on path
671 AC. An improvement would be use a Bayesian method, whereby assignment is a function of
672 both the measurement and a model of the individuals movement history. In terms of
673 state-space models the state equation would represent a probabilistic representation of where
674 the individual could be at that given time based on the individuals entire movement history.

675 The space side of the model would be a measurement equation representing the measurement
676 and the uncertainty surrounding it in the form of a . This would be a similar implementation
677 to that used in Liao et al. (2007), but rather taking the personalised map instead of
678 information about road networks.

679 As for the simplicity of the classification method, the neural network which was used to
680 generate predictions used no information on sequence patterns longer than the previous and
681 next bin. It also has no contextual information about how far away the bins are from each
682 other. A model with more input variables and a more sophisticated design, such as a
683 recurrent neural network (RNN), or a long short-term memory recurrent neural network
684 (LSTM) would likely perform significantly better. Moreover, more inputs can be added from
685 alternative sensors in smartphones such as the accelerometer, barometer, measures of phone
686 activity and so on. These can be used to supplement or improve the raw location
687 measurements, or indeed inform the predictive model.

688 In conclusion, with this paper we aim to start filling the noticeable gap in the social
689 science methodology literature in using smartphone location logs to study human movement.
690 Clearly there is room for improvement, but this is to be expected with new methods. We are
691 in the early days of using smartphone location measurements in social science. Nonetheless,
692 the methodological advantages are clear: millions of individuals have years long location logs
693 containing objective measurements. In addition, these measurements can be easily obtained.
694 This is an unprecedented opportunity. However, with great opportunity comes a great
695 responsibility, as this data raises difficult questions about privacy that researchers must be
696 prepared to answer. Privacy concerns should be addressed, but should not dissuade
697 researchers from developing new methods. Objective location logs are vastly superior to
698 alternatives, such as questionnaires and travel diaries, which rely on accurate self-reporting.
699 Social science researchers must take advantage of regulatory changes in with regard to data
700 portability and put the vast wealth of data collected by commercial entities to scientific use.

701 **References**

- 702 Allaire, J., & Chollet, F. (2018). *Keras: R interface to 'keras'*. Retrieved from
703 <https://CRAN.R-project.org/package=keras>
- 704 Arnold, J. B. (2017). *Ggthemes: Extra themes, scales and geoms for 'ggplot2'*. Retrieved
705 from <https://CRAN.R-project.org/package=ggthemes>
- 706 Aust, F., & Barth, M. (2018). *papaja: Create APA manuscripts with R Markdown*.
707 Retrieved from <https://github.com/crsh/papaja>
- 708 Barnett, I., & Onnela, J.-P. (2016). Inferring mobility measures from GPS traces with
709 missing data. *arXiv:1606.06328 [Stat]*. Retrieved from
710 <http://arxiv.org/abs/1606.06328>
- 711 Bivand, R., Keitt, T., & Rowlingson, B. (2017). *Rgdal: Bindings for the 'geospatial' data
712 abstraction library*. Retrieved from <https://CRAN.R-project.org/package=rgdal>
- 713 Brunsdon, C. (2007). Path estimation from GPS tracks. *Proceedings of the 9th International
714 Conference on GeoComputation*. Retrieved from
715 <http://eprints.maynoothuniversity.ie/6148/>
- 716 Chen, M. Y., Sohn, T., Chmelev, D., Haehnel, D., Hightower, J., Hughes, J., ... Varshavsky,
717 A. (2006). Practical metropolitan-scale positioning for GSM phones. In *UbiComp
718 2006: Ubiquitous computing* (pp. 225–242). Springer, Berlin, Heidelberg.
719 doi:10.1007/11853565_14
- 720 Chen, Z., & Brown, E. N. (2013). State space model. *Scholarpedia*, 8(3), 30868.
721 doi:10.4249/scholarpedia.30868
- 722 Cheng, J., Karambelkar, B., & Xie, Y. (n.d.). *Leaflet: Create interactive web maps with the
723 javascript 'leaflet' library*. Retrieved from <http://rstudio.github.io/leaflet/>
- 724 Commission, E. (2017). *Protecting your data: Your rights - european commission*. Retrieved
725 from http://ec.europa.eu/justice/data-protection/individuals/rights/index_en.htm
- 726 Delclòs-Alió, X., Marquet, O., & Miralles-Guasch, C. (2017). Keeping track of time: A
727 smartphone-based analysis of travel time perception in a suburban environment.

- 728 *Travel Behaviour and Society*, 9(Supplement C), 1–9. doi:[10.1016/j.tbs.2017.07.001](https://doi.org/10.1016/j.tbs.2017.07.001)

729 Duncan, S., Stewart, T. I., Oliver, M., Mavoa, S., MacRae, D., Badland, H. M., & Duncan,
730 M. J. (2013). Portable global positioning system receivers: Static validity and
731 environmental conditions. *American Journal of Preventive Medicine*, 44(2), e19–29.
732 doi:[10.1016/j.amepre.2012.10.013](https://doi.org/10.1016/j.amepre.2012.10.013)

733 Feng, L., Nowak, G., O'Neill, T., & Welsh, A. (2014). CUTOFF: A spatio-temporal
734 imputation method. *Journal of Hydrology*, 519, 3591–3605.
735 doi:[10.1016/j.jhydrol.2014.11.012](https://doi.org/10.1016/j.jhydrol.2014.11.012)

736 Goldberg, L. R., Johnson, J. A., Eber, H. W., Hogan, R., Ashton, M. C., Cloninger, C. R., &
737 Gough, H. G. (2006). The international personality item pool and the future of
738 public-domain personality measures. *Journal of Research in Personality*, 40(1),
739 84–96.

740 Goodchild, M. F., & Janelle, D. G. (2010). Toward critical spatial thinking in the social
741 sciences and humanities. *GeoJournal*, 75(1), 3–13. doi:[10.1007/s10708-010-9340-3](https://doi.org/10.1007/s10708-010-9340-3)

742 Grünerbl, A., Muaremi, A., Osmani, V., Bahle, G., Ohler, S., Tröster, G., ... Lukowicz, P.
743 (2015). Smartphone-based recognition of states and state changes in bipolar disorder
744 patients. *IEEE Journal of Biomedical and Health Informatics*, 19(1), 140–148.
745 doi:[10.1109/JBHI.2014.2343154](https://doi.org/10.1109/JBHI.2014.2343154)

746 Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D.
747 (2016). Using smartphones to collect behavioral data in psychological science:
748 Opportunities, practical considerations, and challenges. *Perspectives on Psychological
749 Science*, 11(6), 838–854. doi:[10.1177/1745691616650285](https://doi.org/10.1177/1745691616650285)

750 Hijmans, R. J. (2017a). *Geosphere: Spherical trigonometry*. Retrieved from
751 <https://CRAN.R-project.org/package=geosphere>

752 Hijmans, R. J. (2017b). *Raster: Geographic data analysis and modeling*. Retrieved from
753 <https://CRAN.R-project.org/package=raster>

754 Jankowska, M. M., Schipperijn, J., & Kerr, J. (2015). A framework for using GPS data in

- 755 physical activity and sedentary behavior studies. *Exercise and Sport Sciences*
756 *Reviews*, 43(1), 48–56. doi:[10.1249/JES.00000000000000035](https://doi.org/10.1249/JES.00000000000000035)
- 757 LaMarca, A., Chawathe, Y., Consolvo, S., Hightower, J., Smith, I., Scott, J., ... Schilit, B.
758 (2005). Place lab: Device positioning using radio beacons in the wild. In *Pervasive
759 computing* (pp. 116–133). Springer, Berlin, Heidelberg. doi:[10.1007/11428572_8](https://doi.org/10.1007/11428572_8)
- 760 Liao, L., Patterson, D. J., Fox, D., & Kautz, H. (2007). Learning and inferring transportation
761 routines. *Artificial Intelligence*, 171(5), 311–331. doi:[10.1016/j.artint.2007.01.006](https://doi.org/10.1016/j.artint.2007.01.006)
- 762 Location History, G. (2017). *Timeline*. Retrieved from
763 <https://www.google.com/maps/timeline?pb>
- 764 London, I. (2016). *Encoding cyclical continuous features - 24-hour time. Ian london's blog*.
765 Retrieved April 27, 2018, from
766 [//ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/](https://ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/)
- 767 Müller, K. (2017). *Bindrcpp: An 'rcpp' interface to active bindings*. Retrieved from
768 <https://CRAN.R-project.org/package=bindrcpp>
- 769 Neuwirth, E. (2014). *RColorBrewer: ColorBrewer palettes*. Retrieved from
770 <https://CRAN.R-project.org/package=RColorBrewer>
- 771 Palmius, N., Tsanas, A., Saunders, K. E. A., Bilderbeck, A. C., Geddes, J. R., Goodwin, G.
772 M., & Vos, M. D. (2017). Detecting bipolar depression from geographic location data.
773 *IEEE Transactions on Biomedical Engineering*, 64(8), 1761–1771.
774 doi:[10.1109/TBME.2016.2611862](https://doi.org/10.1109/TBME.2016.2611862)
- 775 Patterson, T. A., Thomas, L., Wilcox, C., Ovaskainen, O., & Matthiopoulos, J. (2008).
776 State-space models of individual animal movement. *Trends in Ecology & Evolution*,
777 23(2), 87–94. doi:[10.1016/j.tree.2007.10.009](https://doi.org/10.1016/j.tree.2007.10.009)
- 778 Pebesma, E. J., & Bivand, R. S. (2005). Classes and methods for spatial data in R. *R News*,
779 5(2), 9–13. Retrieved from <https://CRAN.R-project.org/doc/Rnews/>
- 780 Preisler, H. K., Ager, A. A., Johnson, B. K., & Kie, J. G. (2004). Modeling animal
781 movements using stochastic differential equations. *Environmetrics* 15: P. 643–657.

- 782 Retrieved from <https://www.fs.usda.gov/treesearch/pubs/33038>
- 783 R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna,
784 Austria: R Foundation for Statistical Computing. Retrieved from
785 <https://www.R-project.org/>
- 786 Sadilek, A., & Krumm, J. (2016). Far out: Predicting long-term human mobility. *Microsoft*
787 *Research*. Retrieved from <https://www.microsoft.com/en-us/research/publication/far-predicting-long-term-human-mobility/>
- 788 Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., & Mohr, D.
789 C. (2015). Mobile phone sensor correlates of depressive symptom severity in daily-life
790 behavior: An exploratory study. *Journal of Medical Internet Research*, 17(7), e175.
791 doi:[10.2196/jmir.4273](https://doi.org/10.2196/jmir.4273)
- 792 Schipperijn, J., Kerr, J., Duncan, S., Madsen, T., Klinker, C. D., & Troelsen, J. (2014).
793 Dynamic accuracy of GPS receivers for use in health research: A novel method to
794 assess GPS accuracy in real-world settings. *Frontiers in Public Health*, 2, 21.
795 doi:[10.3389/fpubh.2014.00021](https://doi.org/10.3389/fpubh.2014.00021)
- 796 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology:
797 Undisclosed flexibility in data collection and analysis allows presenting anything as
798 significant. *Psychological Science*, 22(11), 1359–1366. doi:[10.1177/0956797611417632](https://doi.org/10.1177/0956797611417632)
- 799 Sobrado, B. (2018). *My thesis on missing data in GPS measurements*. Retrieved from
800 <https://github.com/sobradob/thesis>
- 801 ST-DBSCAN: An algorithm for clustering spatial-temporal data. (2007). *Data & Knowledge
802 Engineering*, 60(1), 208–221. doi:[10.1016/j.datwk.2006.01.013](https://doi.org/10.1016/j.datwk.2006.01.013)
- 803 Tatsiramos, K. (2009). Geographic labour mobility and unemployment insurance in europe.
804 *Journal of Population Economics*, 22(2), 267–283. doi:[10.1007/s00148-008-0194-7](https://doi.org/10.1007/s00148-008-0194-7)
- 805 Thoen, E. (2017). *Padr: Quickly get datetime data ready for analysis*. Retrieved from
806 <https://CRAN.R-project.org/package=padr>
- 807 Vaughan, D., & Dancho, M. (2018). *Tibbletime: Time aware tibbles*. Retrieved from
808 <https://CRAN.R-project.org/package=tibbletime>

- 809 <https://CRAN.R-project.org/package=tibbletime>
- 810 Villanueva, C., & Aggarwal, B. (2013). The association between neighborhood
811 socioeconomic status and clinical outcomes among patients 1 year after
812 hospitalization for cardiovascular disease. *Journal of Community Health*, 38(4),
813 690–697. doi:[10.1007/s10900-013-9666-0](https://doi.org/10.1007/s10900-013-9666-0)
- 814 Wang, R., Harari, G., Hao, P., Zhou, X., & Campbell, A. T. (2015). SmartGPA: How
815 smartphones can assess and predict academic performance of college students. In
816 *Proceedings of the 2015 ACM international joint conference on pervasive and*
817 *ubiquitous computing* (pp. 295–306). New York, NY, USA: ACM.
818 doi:[10.1145/2750858.2804251](https://doi.org/10.1145/2750858.2804251)
- 819 Wickham, H. (2009). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
820 Retrieved from <http://ggplot2.org>
- 821 Wickham, H. (2017). *Scales: Scale functions for visualization*. Retrieved from
822 <https://CRAN.R-project.org/package=scales>
- 823 Wickham, H., & Henry, L. (2018). *Tidyr: Easily tidy data with 'spread()' and 'gather()'*
824 *functions*. Retrieved from <https://CRAN.R-project.org/package=tidyr>
- 825 Wickham, H., & Ruiz, E. (2018). *Dbplyr: A 'dplyr' back end for databases*. Retrieved from
826 <https://CRAN.R-project.org/package=dbplyr>
- 827 Wickham, H., Francois, R., Henry, L., & Müller, K. (2017). *Dplyr: A grammar of data*
828 *manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- 829 Wickham, H., Hester, J., & Francois, R. (2017). *Readr: Read rectangular text data*.
830 Retrieved from <https://CRAN.R-project.org/package=readr>
- 831 Wolf, J., Oliveira, M., & Thompson, M. (2003). Impact of underreporting on mileage and
832 travel time estimates: Results from global positioning system-enhanced household
833 travel survey. *Transportation Research Record: Journal of the Transportation*
834 *Research Board*, 1854, 189–198. doi:[10.3141/1854-21](https://doi.org/10.3141/1854-21)
- 835 Wu, F., Fu, K., Wang, Y., Xiao, Z., & Fu, X. (2017). A spatial-temporal-semantic neural

836 network algorithm for location prediction on moving objects. *Algorithms*, 10(2), 37.

837 doi:10.3390/a10020037

838 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida:

839 Chapman; Hall/CRC. Retrieved from <https://yihui.name/knitr/>

840 Zenk, S. N., Schulz, A. J., & Odoms-Young, A. (2009). How neighborhood environments

841 contribute to obesity. *The American Journal of Nursing*, 109(7), 61–64.

842 doi:10.1097/01.NAJ.0000357175.86507.c8

843 Zhang, Z., Yang, X., Li, H., Li, W., Yan, H., & Shi, F. (2017). Application of a novel hybrid

844 method for spatiotemporal data imputation: A case study of the minqin county

845 groundwater level. *Journal of Hydrology*, 553(Supplement C), 384–397.

846 doi:10.1016/j.jhydrol.2017.07.053

847 Zhu, H. (2018). *KableExtra: Construct complex table with 'kable' and pipe syntax*. Retrieved

848 from <https://CRAN.R-project.org/package=kableExtra>

Table 1

Distance in meters between the removed time period and the imputed value.

| | Five minutes | | | One Hour | | | One Day | | |
|------------------|--------------|--------|----------|----------|--------|----------|---------|--------|----------|
| | Mean | Median | Coverage | Mean | Median | Coverage | Mean | Median | Coverage |
| Barnett & Onella | 240 | 5 | 97% | 33 | 6 | 88% | 62 | 6 | 85% |
| Palnius | 43 | 9 | 97% | 497 | 4 | 89% | NA | NA | 0% |
| PPMI | 269 | 0 | 100% | 908 | 0 | 100% | 5,757 | 0 | 100% |
| Naive Baseline | 426 | 0 | 100% | 1,502 | 0 | 100% | 14,266 | 1,288 | 100% |
| Home Baseline | 5,599 | 0 | 100% | 5,667 | 0 | 100% | 5,757 | 0 | 100% |

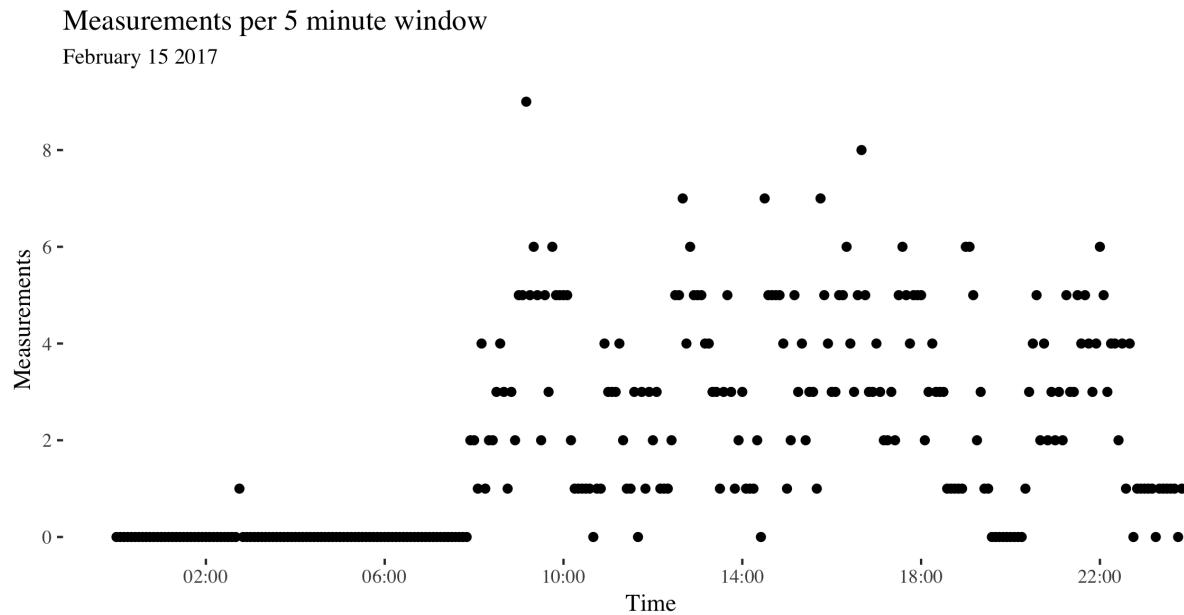


Figure 1. Example of missing data over the entire duration of a secondary log. The x-axis denotes time, the y-axis shows how many measurements are made and each point is a five minute window. For this day there were several periods with no information. These points lie on the x-axis.

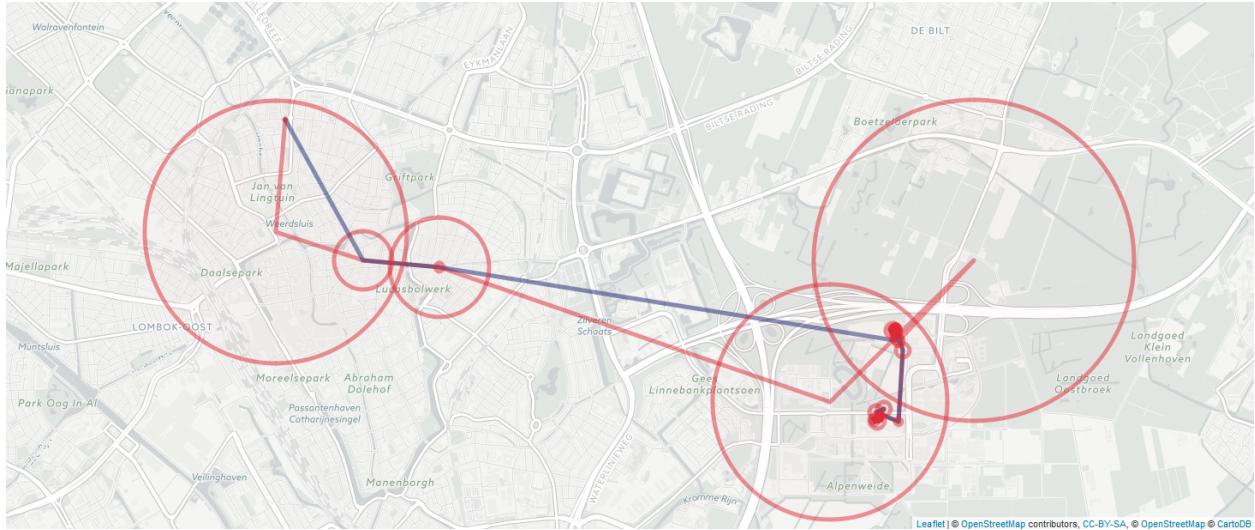


Figure 2. Measurement accuracy of each logged measurement of a morning journey on February 15th 2017. This includes all measurements from midnight to midday. The red circles denote the accuracy of all logged measurement points (the raw data). The points connected in time are connected by a line. The blue line shows the path without the most inaccurate (accuracy > 400 meters) points filtered out. The red line shows the path with all measurements included. In smartphone logs inaccurate location values are interspersed between more accurate location values at higher sample rates per hour. Inaccurate measures are often followed by more accurate measures. There are several recurring low-accuracy points, such as the one in the northwest corner, possibly the result of cellphone tower triangulation.

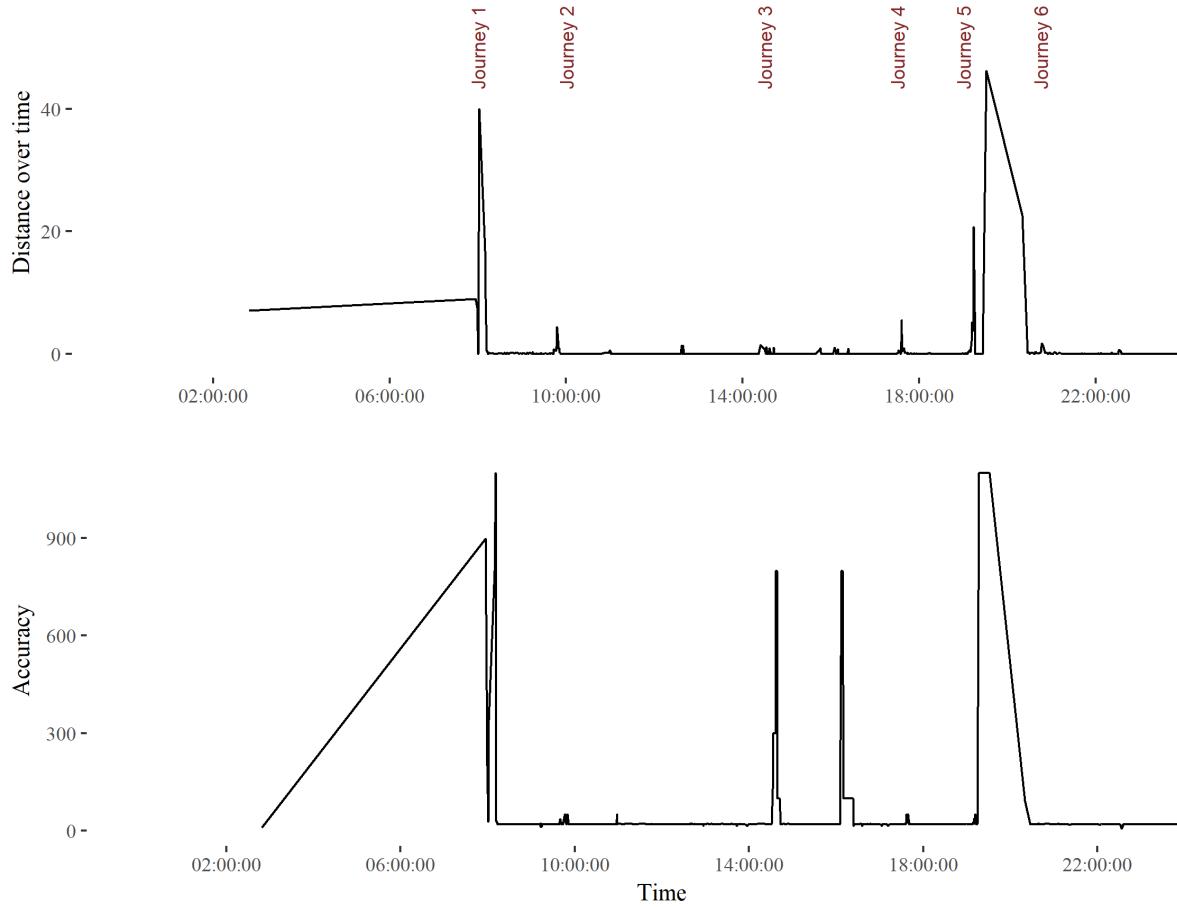


Figure 3. Measures of user activity and measurement accuracy on February 15th 2017. The upper chart shows the distance from the next measured point in meters over the course of the day. The first peak corresponds to the first journey from the user's home to a gym around 8am. The second, smaller peak before 10 reflects a journey from the gym to the nearby lecture theatre. Both journeys can be seen in Figure 2. The large jump between journey 5 and 6 is measurement error. The lower chart shows the accuracy over the course of the day. The figure shows that measurement inaccuracy is sometimes related to the movement of the individual. Stationary accuracy varies depending on phone battery level, wifi connection and user phone use.

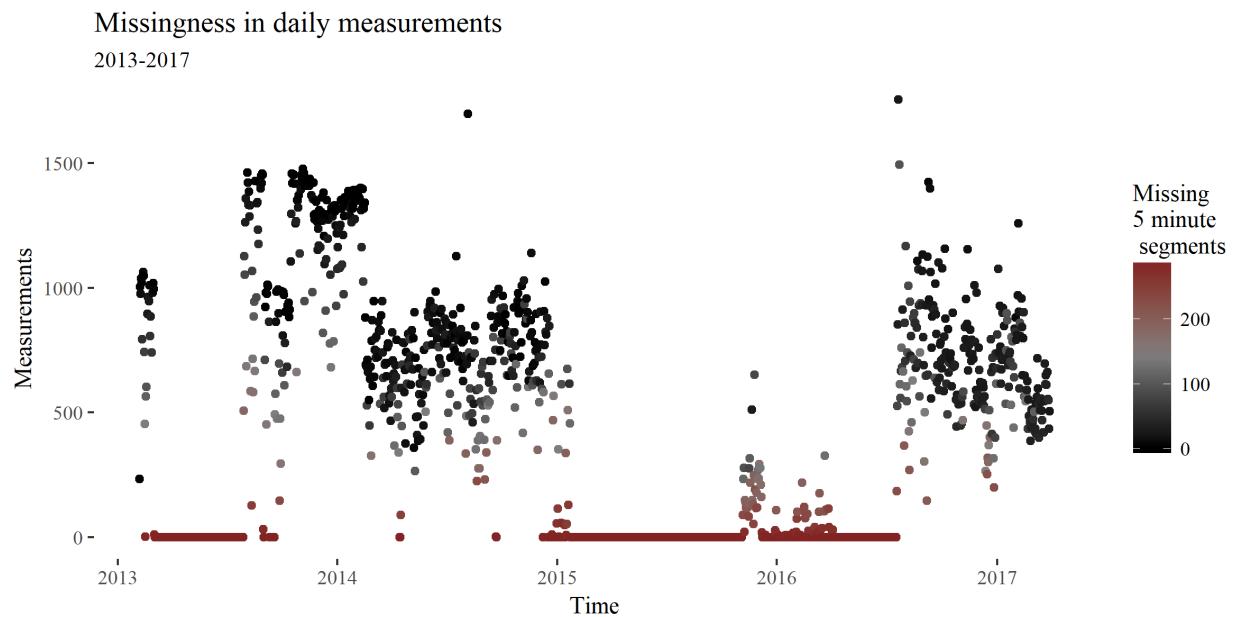


Figure 4. Missing data for the entire duration of the log. The x-axis denotes time, the y-axis shows how many measurements are made and each point is a five minute window. The entire log contains several long periods with no information. These points are filled with red and lie on the x-axis.

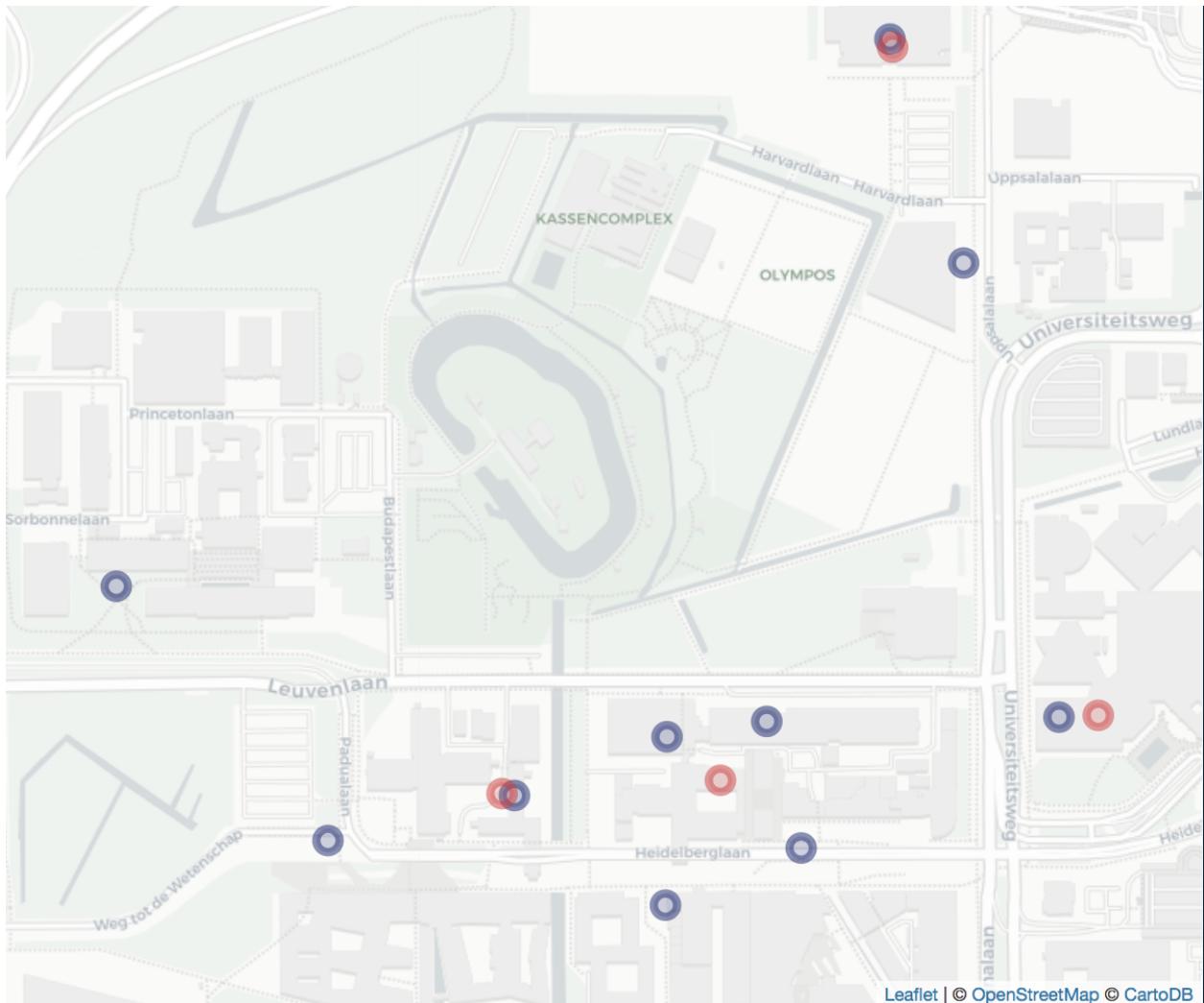


Figure 5. Example of pause locations in De Uithof university campus using 150 meters (blue) and 400 meters (red) as clustering parameters.

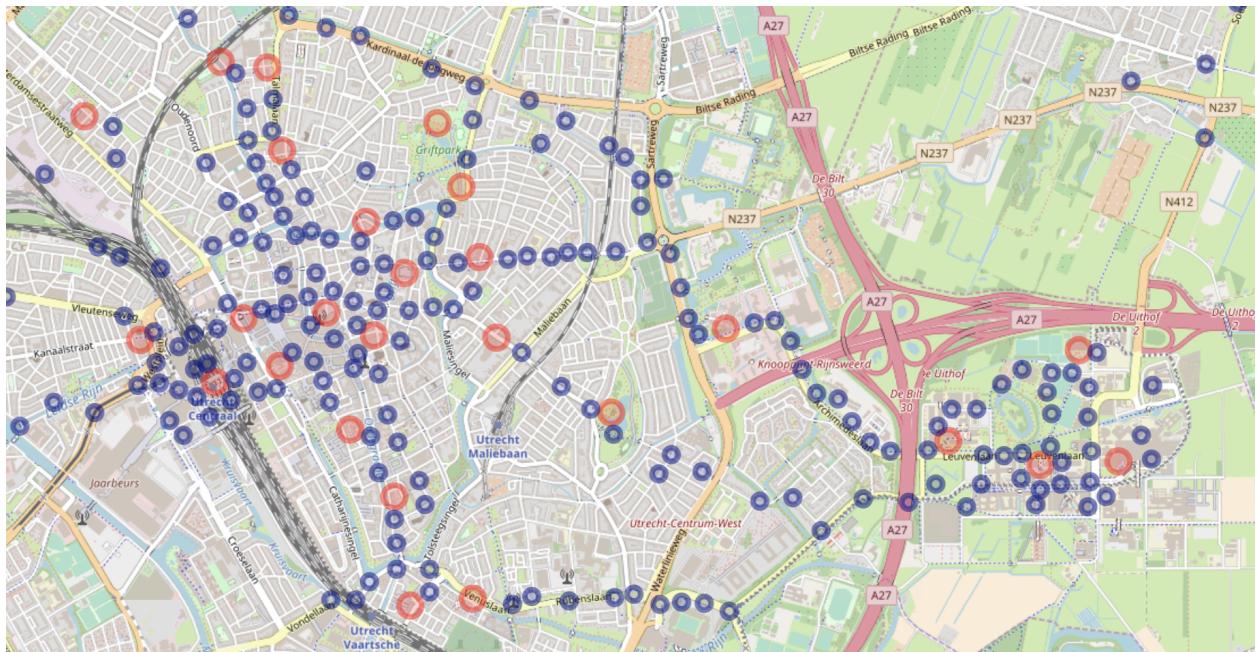


Figure 6. Excerpt of the cluster map of an individual. Red points are pause locations, blue points are path locations.

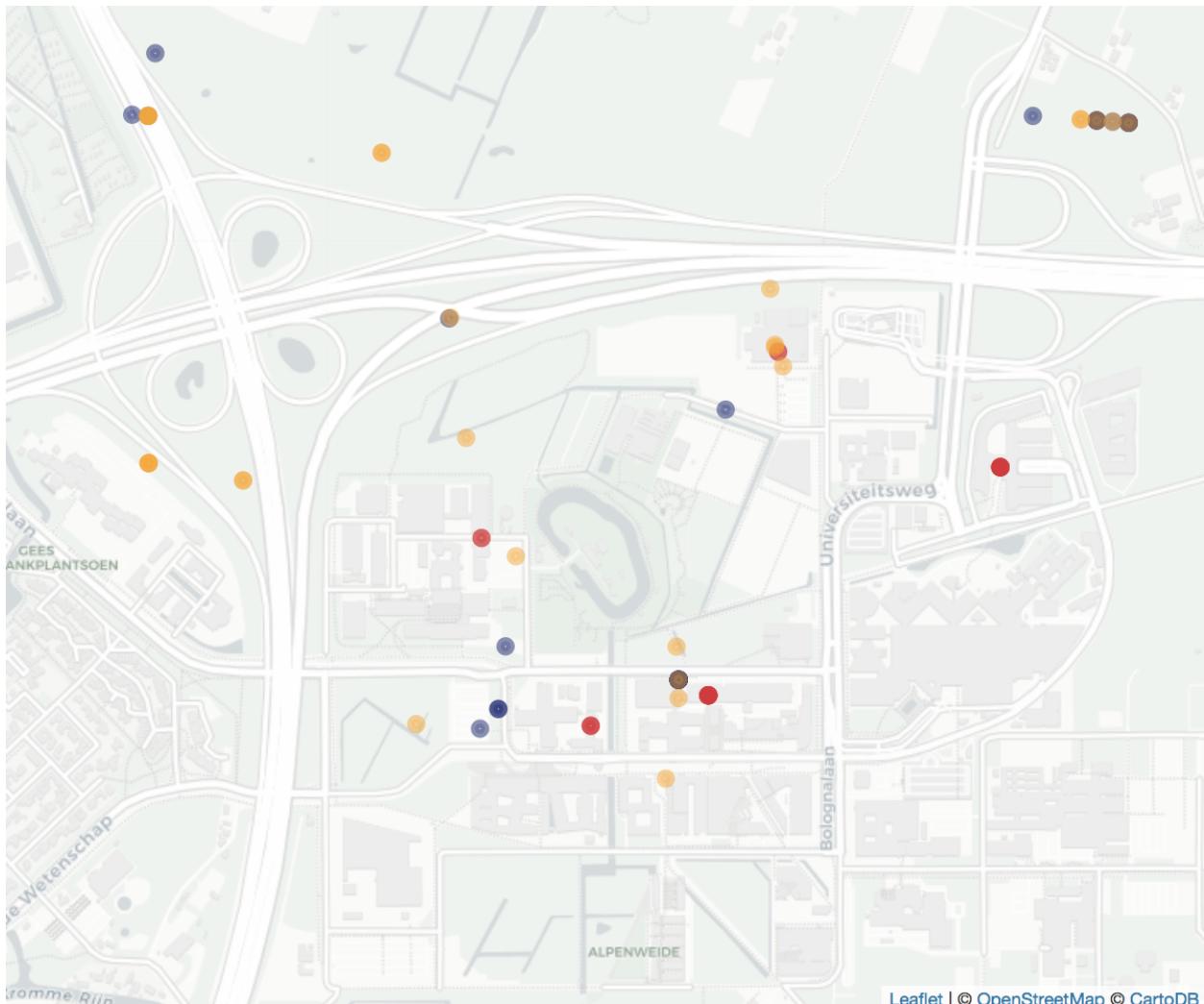


Figure 7. The difference between temporal and spatial downsampling. The blue circles are raw measurements, the yellow circles are temporally downsampled locations. Spatially downsampled locations are in red. Due to measurement sparsity and inaccuracy many of the temporally downsampled locations are in unfeasable locations.

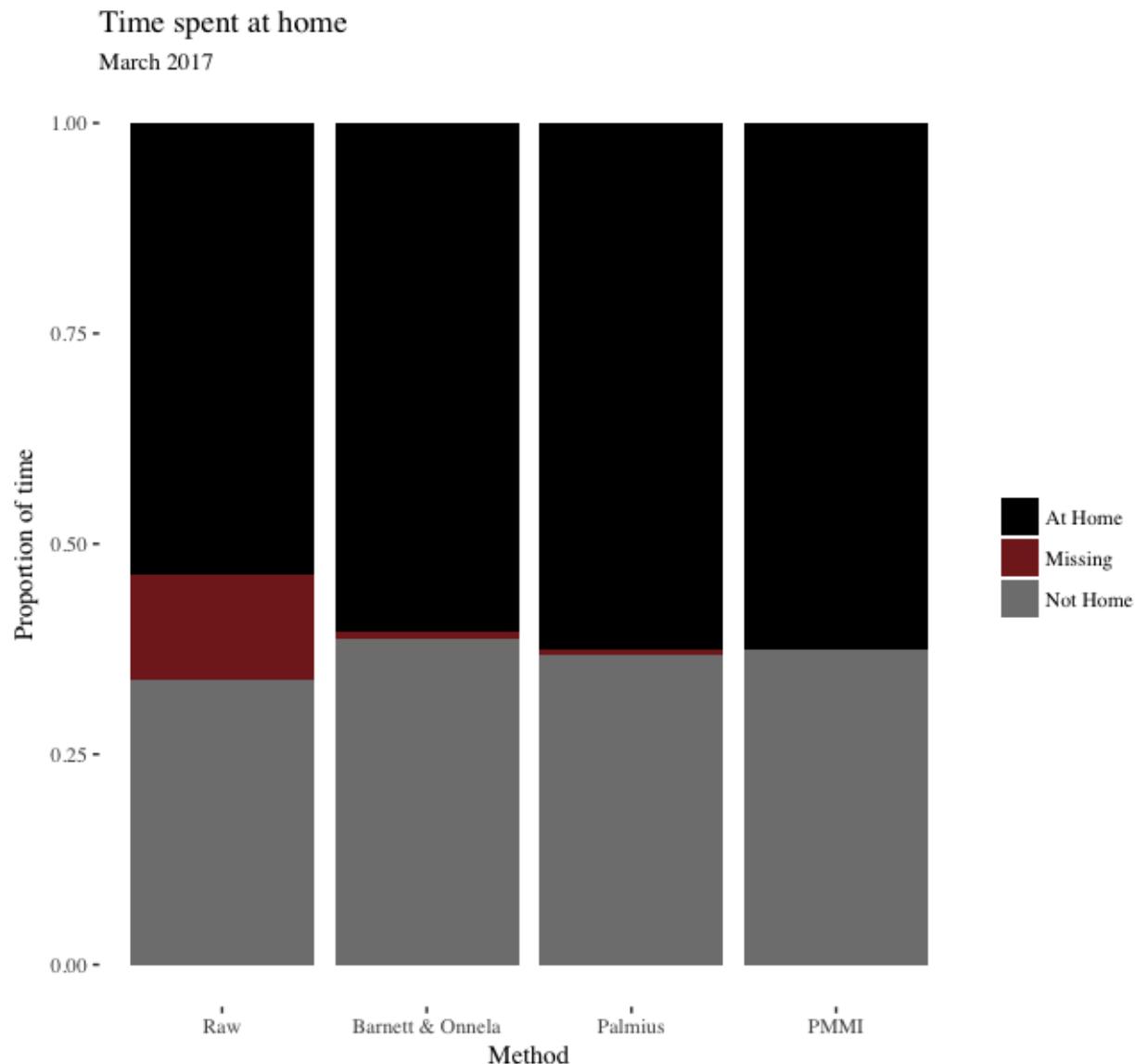


Figure 8. Proportion of time spent at home in March 2017. The raw values are estimated by downsampling temporally the latitude and longitude for every 5 minute time period in the month. We used each method's own binning method and classified as at home if the downsampled measurement was within 250 meters from home.