

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

3 Boaz Sobrado<sup>1</sup>

4 <sup>1</sup> 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](mailto: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: 8000

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 (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 and Onnela (2016)]. Today millions of individuals carry smartphones with themselves every  
82 day and do not need to be encouraged to do so by researchers. Smartphones are equipped  
83 with a range of sensors that can be used to track the location of the device at all times.  
84 These smartphones can collect and store hundreds of location measurements a day. For  
85 instance, Google Location History contains movement information on millions of users, often

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

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

100

## Background

Smartphone location measures are obtained primarily (but not exclusively) by Global Positioning System (GPS) measurements. A GPS sensor uses the distance between a device and several satellites to determine the location of the device. Although using a GPS sensor is the most accurate way to establish location on a smartphone, the GPS sensor is also the most energy consuming sensor on most smartphones (M. Y. Chen et al., 2006; LaMarca et al., 2005). In order to avoid battery depletion and to overcome computational constraints smartphones also use less-accurate heuristics such as WiFi access points and cellphone tower triangulation. Smartphone location logs contain measurements from all of these sources, usually in the form of time-stamped latitude, longitude and accuracy values. The accuracy *a* of any given measurement is given in meters such that it represents the radius of a 67% 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 (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, the CUTOFF method estimates missing  
214 values by incorporating similar observed temporal information from the value's nearest  
215 spatial neighbors (Feng, Nowak, O'Neill, & Welsh, 2014 ). The authors illustrate their  
216 example using rainfall data from gauging stations across Australia. Similarly, Z. Zhang et al.  
217 (2017) use a variety of machine learning methods to impute missing values. The example  
218 provided relates to underground water data. Generally these models assume fixed

219 measurement stations (such as rainfall gauging stations). For this reason they cannot be  
220 easily applied to missing 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 [ Palmius et al. (2017) ;Barnett and Onnela (2016);  
223 wu\_spatial-temporal-semantic\_2017]. Importantly, none of them worked with secondary  
224 logs. Nonetheless we will detail what they did as informative examples. Palmius et al. (2017)  
225 deal 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 On the other hand, wu\_spatial-temporal-semantic\_2017 use what they call a Spatial  
253 Temporal Semantic Neural Network (STS-NN) to predict future human movement. While  
254 the authors are concerned with prediction and not imputation, they devised a method called  
255 the Spatial Temporal Semantic (STS) algorithm which converts raw measurements to  
256 machine learning friendly discrete bins. Working with high-frequency measurements,  
257 wu\_spatial-temporal-semantic\_2017’s method down-samples the raw data temporally and  
258 map-matches the resulting bins to discrete points along pre-established geographical features  
259 such as roads and highways. This minimises measurement error and paves the way for  
260 applying machine learning methods to human mobility problems.

261 In this section we have argued that there is a lack of established practices to follow  
262 with respect to missing data in human mobility logs. Moreover, extensively used  
263 spatiotemporal methods, such as state space models (SSMs), are not well suited to deal with  
264 human mobility patterns in secondary logs. Finally we discussed in detail three approaches  
265 which deal explicitly with mobility patterns from custom or secondary logs [ Palmius et al.  
266 (2017) ;Barnett and Onnela (2016); wu\_spatial-temporal-semantic\_2017].

## 267 Methodology

### 268 Notation

269 Location measurements, such as those produced by GPS sensors, provide us with  
270 coordinates (latitude and longitude) on the surface of the earth, which is ellipsoid shaped.

271 Projecting three dimensional measurements in  $\mathbb{R}^3$  onto a two dimensional plane in  $\mathbb{R}^2$  results  
272 in distortion. For clarity, when we use the term distance we refer to the geodesic distances  
273 on an ellipsoid using the WGS84 ellipsoid parameters.

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

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

## 286 Personalised Map Matching Imputation

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

289 **Modelling assumptions.** First, following Barnett and Onnela (2016) we categorise  
290 all time-points  $t$  as either belonging to the set  $P$  (pause) or set  $F$  (flight). Conceptually  
291 pauses can be understood as periods of time where an individual spends significant amount  
292 of continuous time without moving. Flights are the times where the individual is moving.

293 Let  $t_a$  be a pause of length  $n$ .

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

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

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

<sup>295</sup> Then it follows that between the two pauses there must be a flight indexed by  $t_x$  of length  
<sup>296</sup>  $j - i + n$ .

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

<sup>297</sup> We define pause locations  $G(t_a), G(t_b) | t_a, t_b \in P$  as locations where an individual spends an  
<sup>298</sup> extended amount of time in the same space (e.g. school, home, work, train station, barber  
<sup>299</sup> shop, bar, gym). Importantly, our model assumes period and cyclic human movement such  
<sup>300</sup> that there are many pauses  $t_{a1}, t_{a2}, \dots, t_{an}$  such that  $G(t_{a1}) = G(t_{a2}) = \dots = G(t_{an})$ .  
<sup>301</sup> Moreover, it is possible for  $G(t_a) = G(t_b)$  such that  $t_a \neq t_b$ . For example, if the individual  
<sup>302</sup> leaves home for a run and returns home without stopping anywhere else.

<sup>303</sup> Let us define as  $Flight_{ab}^x$  the set of all points belonging to a flight between  $G(t_a)$  and  
<sup>304</sup>  $G(t_b)$  at time-point  $t_x$ .

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

<sup>305</sup> Again, there are many flights  $t_{x1}, t_{x2}, \dots, t_{xn}$  such that  $Flight_{ab}^{x1} = Flight_{ab}^{x2} = \dots = Flight_{ab}^{xn}$ .  
<sup>306</sup> Then, we can define as  $Path_{ab}$  the set of all flights between  $G(t_a)$  and  $G(t_b)$  at all  
<sup>307</sup> time-points. For simplicity, we assume that  $Path_{ab} = Path_{ba}$ .

<sup>308</sup> In addition, we consider all measurements  $D(t)$  to be imperfect measurements of  $G(t)$ :

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

### <sup>309</sup> Personalised Map Matching Imputation algorithm

<sup>310</sup> Our algorithm performs the following steps:

- <sup>311</sup> 1. *Map building*: Extract from measurements  $D$  all pause location bins and path location  
<sup>312</sup> bins to create a personalised map.
- <sup>313</sup> 2. *Binning*: Assign each measurement  $D$  to a unique discrete location bin.
- <sup>314</sup> 3. *Imputing*: Use a classification method to predict missing measurements based on all  
<sup>315</sup> the available information.

316 **Map building.** Following wu\_spatial-temporal-semantic\_2017's

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

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

- 331 1. The next measurement  $D_{t+1}$  is within  $t_{\text{Pause lim}}$  amount of seconds (so it is not missing)
- 332 2. The next measurement  $D_{t+1}$  is within  $d_{\text{Pause lim}}$  meters.
- 333 3. The duration of the pause is more than  $\delta_{\text{Pause lim}}$  seconds.
- 334 4. Let the measurements of a possible pause which fit the aforementioned criteria be

335  $D_{t,t+1,\dots,t+n}$ . These points are only a pause if the distance between the mean  
336 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  
337 the mean accuracy  $a$  of  $D_{t,t+1,\dots,t+n}$ .

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

341 For all remaining measurements we assume that they belong to paths. In this  
342 implementation of PMMI we use the following algorithm to estimate paths:

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

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

355         **Binning.** wu\_spatial-temporal-semantic\_2017's spatial-temporal-semantic (STS)  
 356         feature extraction algorithm uses map matching as a ground truth to assign noisy  
 357         measurements into discrete bins along roads. In other words, in addition to the measured  
 358         data they also use a geographic database that contains information about the area in which  
 359         the individual is (e.g. where precisely the roads are), and sort measurements into bins based  
 360         on both the measurement and the geographic data base. For example, if an individual is  
 361         measured as moving closely in parallel to a road A in an area where there is no other parallel  
 362         road, wu\_spatial-temporal-semantic\_2017's method will assume that the individual is on  
 363         the road A.

364             PMMI uses a similar logic, but without using any external geographic database. The  
 365         key modification in PMMI is that whilst wu\_spatial-temporal-semantic\_2017 uses a map  
 366         from outside the persons location logs, we use the total location history of the individual to  
 367         create a personalised map. This map is subsequently used to bin measurements. This is  
 368         feasible for two reasons: humans tend to have repetitive movement habits and secondary logs  
 369         tend to be long. To put it in simpler terms, we consider each measurement at  $D_x$  as a sample

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

371 Thus, all measurements  $D$  were assigned to a discrete bin on the personal  
372 map. This includes previous measurements which were discarded from the map building  
373 exercise due to an accuracy  $a$  value which exceeded  $a_{P\lim}$  or  $a_{F\lim}$ . In this implementation  
374 we used a simple assigning function, whereby the measurements were assigned to the bin  
375 nearest to the measurement.

376 **Classification.** At this point the objective of PMMI is to take all the information  
377 available about the mobility history of an individual and impute the missing value. In this  
378 implementation we trained an artificial neural network (ANN) to do so. For more  
379 information on the precise architecture of the artificial neural network please consult the  
380 appendix. The input variables to the ANN are:

- 381 1. The previous and subsequent observed bin as a binary class matrix.
- 382 2. The distance in time to the next & previous bin.
- 383 3. The time of the day encoded as a cyclical two-dimensional feature.
- 384 4. The day of the month as a binary class matrix.
- 385 5. The month of the year as a binary class matrix.

386 For the encoding of the time of day we took the cosine and the sine transforms of the  
387 amount of seconds that have elapsed after midnight (London, 2016). This is necessary so  
388 that the model can understand that one second past midnight and one second before  
389 midnight are in fact two seconds away from each other. Moreover we scaled the non-binary  
390 values to occupy a range between 0 and 1 in order to ensure convergence.

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

395 **Datasets & Analyses**

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

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

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

413 All analyses were performed using R (Version 3.4.3; R Core Team, 2017) and the  
414 R-packages *bindrcpp* (Version 0.2; Müller, 2017), *dbplyr* (Version 1.2.0; Wickham & Ruiz,  
415 2018), *dplyr* (Version 0.7.4; Wickham, Francois, Henry, & Müller, 2017), *geosphere* (Version  
416 1.5.7; Hijmans, 2017a), *ggplot2* (Version 2.2.1; Wickham, 2009), *ggthemes* (Version 3.4.0;  
417 Arnold, 2017), *kableExtra* (Version 0.7.0; Zhu, 2018), *keras* (Version 2.1.4; Allaire & Chollet,  
418 2018), *knitr* (Version 1.20; Xie, 2015), *leaflet* (Version 2.0.0; Cheng, Karambelkar, & Xie,  
419 n.d.), *padr* (Version 0.4.0; Thoen, 2017), *papaja* (Version 0.1.0.9709; Aust & Barth, 2018),  
420 *raster* (Version 2.6.7; Hijmans, 2017b), *RColorBrewer* (Version 1.1.2; Neuwirth, 2014), *readr*  
421 (Version 1.1.1; Wickham, Hester, & Francois, 2017), *rgdal* (Version 1.2.16; R. Bivand, Keitt,

<sup>422</sup> & Rowlingson, 2017), *scales* (Version 0.5.0; Wickham, 2017), *sp* (Version 1.2.7; Hijmans,  
<sup>423</sup> 2017a; Pebesma & Bivand, 2005), *tibbletime* (Version 0.1.0; Vaughan & Dancho, 2018), and  
<sup>424</sup> *tidyrr* (Version 0.8.0; Wickham & Henry, 2018). All the code is available on a public  
<sup>425</sup> repository (Sobrado, 2018).

## <sup>426</sup> Results & Evaluation Metrics

<sup>427</sup> The results will consist of multiple steps:

- <sup>428</sup> 1. Evaluating the performance of the map building and assigning functions.
- <sup>429</sup> 2. Comparing the performance of PMMI using a) baseline models b) performance with  
<sup>430</sup> randomly removed data in comparison to the aforementioned methods (Barnett &  
<sup>431</sup> Onnela, 2016; Palmius et al., 2017) and c) objective ground-truth data (public  
<sup>432</sup> transportation time-stamps).

<sup>433</sup> **Map building & binning evaluation.** Before we can evaluate the accuracy of the  
<sup>434</sup> imputations, it is essential to evaluate how well noise in the data has been cleaned.  
<sup>435</sup> Otherwise we run the risk of over-fitting the model in the sense that we will measure the  
<sup>436</sup> extent to which an imputation method can correctly impute measurement noise within  $D_t$   
<sup>437</sup> instead of true location  $G(t)$ .

<sup>438</sup> In order to evaluate the map building and binning we will first visually evaluate the  
<sup>439</sup> paths and pause locations. A visual evaluation of paths superimposed on is an established  
<sup>440</sup> way to heuristically check their accuracy (e.g. Brunsdon (2007)). Then, let the average  
<sup>441</sup> distance between the actual measured point and the binned point be the *deviation distance*  
<sup>442</sup>  $\delta_{\text{dev}}$ . With respect to the deviation distance  $\delta_{\text{dev}}$ , we expect:

- <sup>443</sup> 1. A positive relationship between the deviation distance and the accuracy of each  
<sup>444</sup> measurement.
- <sup>445</sup> 2. Roughly 67% of the deviation distances  $\delta_{\text{dev}}$  are within accuracy  $a$  of each  
<sup>446</sup> measurement.

447       **Imputation algorithm performance.** We will compare the performance of PMMI,

448 Palmius et al. (2017) and Barnett and Onnela (2016). In addition, we will also compute a  
449 naive model, which simply imputes as the missing value the previous observed value. The  
450 naive model will serve as a baseline model. To compare the performance of these methods  
451 we will remove 25% of measured time intervals at random within a four week period. We will  
452 make our comparisons for intervals of 5 minutes, 1 hour and 1 day. In other words, we will  
453 remove 25% of time intervals at random, while varying the size of the time intervals removed.

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

457       To compare all methods with each other we will compute a distance measure (how far  
458 was the removed location from the predicted location). For PMMI's imputations we will use  
459 a weighted mean for the distance measures whereby each 5 minute period is weighted equally.  
460 This is necessary because the other two model estimate far fewer values than PMMI. While  
461 Barnett and Onnela (2016)'s and Palmius et al. (2017)'s models estimate 12 measurements  
462 for each missing hour, PMMI estimates as many measurements as there are in the log. When  
463 individuals go to infrequently travelled locations they tend to use their phone's location  
464 services more than 100 more times, which leads to more measurements, and hence more  
465 imputed values. Hence, if the weights are not used PMMI is susceptible to error  
466 overestimates which are related to measurement frequency in a way that the other two  
467 methods are not.

468       In addition, other measures of interest for PMMI are accuracy (in what percentage of  
469 the cases was the appropriate cluster predicted), the *confidence* and the *distance expectation*.  
470 The confidence is the probability with which the model predicts the most likely cluster. For  
471 instance, if the model predicts the missing bin to be bin A with a probability of 0.9 then we  
472 can say it has high confidence in the prediction. Similarly, the distance expectation is the  
473 cross-product of the estimated probabilities that the individual is at any of all given clusters

<sup>474</sup> and the distances between the clusters to the true cluster. For example, if the true location  
<sup>475</sup> of an individual is bin A (bin A is 10 meters away from bin B) and the model assigns a  
<sup>476</sup> probability of 0.9 at bin A and 0.1 at bin B, then the distance expectation would be 1 meter.

<sup>477</sup> Finally, we will take objective real world data and compare it to predicted values. We  
<sup>478</sup> will use information from the Dutch public transportation card. The Dutch public  
<sup>479</sup> transportation service provides users with time-stamped data of when and where they board,  
<sup>480</sup> change lines or leave public transportation. To be able to make a comparison between  
<sup>481</sup> models we will remove all measurements from within the 5 minute period that a  
<sup>482</sup> time-stamped measurement is available. Then we will use each model to impute the location  
<sup>483</sup> of the individual within that period.

## <sup>484</sup> Results

### <sup>485</sup> Map building & binning evaluation

<sup>486</sup> We used the following parameters to extract pauses: an accuracy limit  $a_{\text{Pause Lim}}$  of 250  
<sup>487</sup> meters, a time limit  $t_{\text{Pause Lim}}$  of 300 seconds, a distance limit of  $d_{\text{Pause Lim}}$  50 meters and a  
<sup>488</sup> minimum pause duration limit  $\delta_{\text{Pause Lim}}$  of 100 seconds. Moreover, to extract path clusters  
<sup>489</sup> we used the parameters:  $a_{\text{Path Lim}} = 150$ meters and  $d_{\text{Path Lim}} = 300$ meters.

<sup>490</sup> The selection of these parameters was more-or-less heuristically driven based on their  
<sup>491</sup> on their ability to extract a meaningful personalised map. When selecting parameters there  
<sup>492</sup> is a trade off is between bias and precision. This is because an increase in precision in the  
<sup>493</sup> form of a higher for high density locations resolution comes at the expense of precision as  
<sup>494</sup> assigning measurements to bins becomes more difficult. For instance, by increasing the  
<sup>495</sup> clustering distance parameter  $d_{\text{Pause Lim}}$  we can extract more valid pause locations at the  
<sup>496</sup> expense of falsely categorising certain measurements to the wrong cluster. This is illustrated  
<sup>497</sup> in Figure 5.

<sup>498</sup> Map building results in a personalised map with pause and path clusters. An excerpt  
<sup>499</sup> can be seen in Figure 6. It is important to remind the reader that PMMI is map agnostic and

uses no information from the map. Therefore, the close overlap with features on the map, such as pause bins at relevant buildings and transportation clusters, as well as the flight bins following roads and railway lines indicate a high degree of precision in personal map building. As expected, PMMI's path extraction yields greater accuracy for frequently occurring paths. For example, the frequently travelled Amsterdam-Utrecht railway line has been extracted almost perfectly, while the less frequently travelled Utrecht-Enschede line is far sparser.

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

For the narrower period of March 2017 approximately 74% of the deviance values are within their corresponding accuracy values. Again, this is not very far from the theoretically expected 67%. The raw unweighted mean and median deviance are 38 and 14 meters respectively. If we weight them such that each 5 minute interval is weighted equally to ease comparison with the other two methods, we find a mean deviance of 35 meters with a median of 12.8.

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

The key difference between temporal and spatial down-sampling is shown in Figure ???. Temporal down-sampling is much more sensitive to noise in sparsely measured periods because it averages out values within five minute periods. Often there are only a few noisy measurements in those periods (see Figure 1 ), which leads to a noise in the down-sampled

527 values. Unsurprisingly, there is a positive relationship between deviance and the amount of  
 528 measurements in each down-sampled interval. The fact that over 90% of deviance values are  
 529 within accuracy (substantially higher than the expected theoretical 67%) confirms that  
 530 temporal down-sampling is not sufficiently filtering out the noise.

531 **Imputation evaluation.** The Palmius et al. (2017) model failed to impute 3% of  
 532 all removed values for both the 5 minute and over 10% for the 1 hour tests. Palmius et al.  
 533 (2017)' model failed to impute a single value for the day tests. Similarly, Barnett and Onnela  
 534 (2016)'s method failed to impute over 11% of the missing values. PPMI made an imputation  
 535 for all missing values. Table 1 shows the results of the distance metrics for each method.

536 In terms of PPMI's accuracy in this period, the prediction accuracy was 88%, 73% and  
 537 47% for the 5 minute, 1 hour and 1 day periods respectively. As a comparison, the naive  
 538 model's accuracy ratings were 87%, 68% and 24% respectively. The distance expectation  
 539 values were very similar to the distance scores, albeit approximately 5% higher. Confidence  
 540 scores ranged from 0.006 to 1.

541 **Comparison with objective data.** Once we removed the measurements within  
 542 the 5 minute period of each time-stamp, we employed the models to predict the location of  
 543 the individual at the time-stamped period. The naive model reaches an accuracy of 16%,  
 544 with a median distance of 555 meters and mean distance of 3303 meters. On the other hand  
 545 PPMI has an accuracy of 24%, with a median distance of 1037 meters and a mean distance  
 546 of 5637. Again the expectation of the distance was quite similar to the distance with a mean  
 547 of 6432 and a median of 1037 meters.

548 Palmius et al. (2017)'s model failed to impute 38 out of 97 periods. For those it did  
 549 impute, it showed a mean distance of 1517 meters and a median distance of 1617. Barnett  
 550 and Onnela (2016)'s model failed to impute 18 out of 97 periods. For those it did impute, it  
 551 had a mean of 5506 and a median of 1342.

552 **Example: effect on aggregate measures.** Social scientists are most interested in  
 553 aggregating spatiotemporal data to more socially relevant metrics, such as the amount of

554 time spent at home. As an example we calculated the time spent at home of the user in the  
555 month of March (Figure 8) without any model and with all three of the investigated methods.

556 Interestingly, all three models suggest that the user spent approximately 60% of their  
557 time at home. However, without using any form of missing data imputation approximately  
558 12% of the time is unaccounted for. Given that the amount of time spent at home has been  
559 found to be a reliable predictor of extroversion (G. M. Harari et al., 2016) and the onset of  
560 depressive episodes in bipolar patients (Palmius et al., 2017) obtaining an accurate value for  
561 mobility metrics like this is highly important to social scientists.

## 562 Discussion & Conclusion

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

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

576 PPMI's increased coverage and flexibility comes from its ability to make complex  
577 non-linear predictions. For instance, in a given missing period it might make sense to predict  
578 that the individual is either at home, or at the office, or at a shop with equal probability.  
579 While PPMI can make such an imputation, none of the alternative methods can do this.

580 Moreover, the ability to take the prediction probability values from the neural network also  
581 helps in dealing with uncertainty. A known-drawback of single imputation is that it takes an  
582 imputed value and treats it as observed. Simple rule-based methods such as Palmius et al.  
583 (2017)' are essentially algorithmic single imputation methods. With PPMI it is possible to  
584 model uncertainty using the predicted probabilities of each estimate. For instance, in the  
585 previous example, we could choose to only take estimates with a high degree of confidence,  
586 thus creating confidence intervals by adding and subtracting the amount of cases where the  
587 location of the individual is ambiguous.

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

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

607 There are multiple methodological limitations in this paper. Most importantly, the  
608 evaluation methods are imperfect. The golden standard would be to use at least one highly  
609 accurate professional grade GPS device with high sampling frequency to compare our data  
610 to. Until that is available, the use of public transport data and cross-validation is just a  
611 substitute.

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

614 In map building, the probability of a pause at a given location is certainly related to  
615 other factors, such as the time of the day as well as the prior history of pauses at that  
616 location. These factors are not taken into account in the pause extraction function.

617 Improved methods would do well to do so. As for paths, a drawback of the current method  
618 is that the density of the clusters is a function of the clustering parameter  $d$ , the distance  
619 between the observed points and their sampling density. It does not take into account the  
620 length of the path as well as the average sampling frequency of the path. This is an issue  
621 because it can lead to bins to which data points are seldom assigned. For example, while the  
622 Amsterdam-Utrecht line has been mapped out almost perfectly, many of the clusters along  
623 the route have only been assigned few measurements. This leads to difficulties in the  
624 classification part of the model, as infrequently observed clusters are hard for the model to  
625 predict.

626 The current assignment function simply assigns each measurement to the nearest bin.  
627 It does not take into account any contextual information that can be gleaned from the entire  
628 movement history of the information, such as what path they are on. For instance, assume  
629 that it is known that an individual is travelling from point A to point B along path AB, and  
630 there is an inaccurate measurement closest to a cluster which belongs to path AC. By only  
631 taking distance into account, the measurement can get assigned to the wrong cluster on path  
632 AC. An improvement would be use a Bayesian method, whereby assignment is a function of  
633 both the measurement and a model of the individuals movement history. In terms of

634 state-space models the state equation would represent a probabilistic representation of where  
635 the individual could be at that given time based on the individuals entire movement history.  
636 The space side of the model would be a measurement equation representing the measurement  
637 and the uncertainty surrounding it in the form of  $a$ .

638 As for the simplicity of the classification method, the neural network which was used to  
639 generate predictions used no information on sequence patterns longer than the previous and  
640 next bin. A more sophisticated recurrent neural network (RNN), or a long short-term  
641 memory recurrent neural network (LSTM) would likely perform significantly better.

642 That there is room for improvement is to be expected given that we are in the early  
643 days of using smartphone location measurements in social science. Nonetheless, the  
644 methodological advantages are clear: millions of individuals have years long location logs  
645 containing objective measurements which can be easily obtained. This is an unprecedented  
646 opportunity given that objective location logs are vastly superior to alternatives such  
647 as questionnaires which rely on accurate self-reporting. Social science researchers must take  
648 advantage of regulatory changes in with regard to data portability and put the vast wealth of  
649 data collected by commercial entities to scientific use.

**References**

- 650
- 651 Allaire, J., & Chollet, F. (2018). *Keras: R interface to 'keras'*. Retrieved from  
<https://CRAN.R-project.org/package=keras>
- 652
- 653 Arnold, J. B. (2017). *Ggthemes: Extra themes, scales and geoms for 'ggplot2'*. Retrieved  
from <https://CRAN.R-project.org/package=ggthemes>
- 654
- 655 Aust, F., & Barth, M. (2018). *papaja: Create APA manuscripts with R Markdown*.  
Retrieved from <https://github.com/crsh/papaja>
- 656
- 657 Barnett, I., & Onnela, J.-P. (2016). Inferring mobility measures from GPS traces with  
missing data. *arXiv:1606.06328 [Stat]*. Retrieved from  
<http://arxiv.org/abs/1606.06328>
- 658
- 659
- 660 Bivand, R., Keitt, T., & Rowlingson, B. (2017). *Rgdal: Bindings for the 'geospatial' data  
abstraction library*. Retrieved from <https://CRAN.R-project.org/package=rgdal>
- 661
- 662 Brunsdon, C. (2007). Path estimation from GPS tracks. *Proceedings of the 9th International  
Conference on GeoComputation*. Retrieved from  
<http://eprints.maynoothuniversity.ie/6148/>
- 663
- 664
- 665 Chen, M. Y., Sohn, T., Chmelev, D., Haehnel, D., Hightower, J., Hughes, J., ... Varshavsky,  
A. (2006). Practical metropolitan-scale positioning for GSM phones. In *UbiComp  
2006: Ubiquitous computing* (pp. 225–242). Springer, Berlin, Heidelberg.  
doi:10.1007/11853565\_14
- 666
- 667
- 668
- 669 Chen, Z., & Brown, E. N. (2013). State space model. *Scholarpedia*, 8(3), 30868.  
doi:10.4249/scholarpedia.30868
- 670
- 671 Cheng, J., Karambelkar, B., & Xie, Y. (n.d.). *Leaflet: Create interactive web maps with the  
javascript 'leaflet' library*. Retrieved from <http://rstudio.github.io/leaflet/>
- 672
- 673 Commission, E. (2017). *Protecting your data: Your rights - european commission*. Retrieved  
from [http://ec.europa.eu/justice/data-protection/individuals/rights/index\\_en.htm](http://ec.europa.eu/justice/data-protection/individuals/rights/index_en.htm)
- 674
- 675 Delclòs-Alió, X., Marquet, O., & Miralles-Guasch, C. (2017). Keeping track of time: A  
smartphone-based analysis of travel time perception in a suburban environment.
- 676

- 677 *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)

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

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

685 Goldberg, L. R., Johnson, J. A., Eber, H. W., Hogan, R., Ashton, M. C., Cloninger, C. R., &  
686 Gough, H. G. (2006). The international personality item pool and the future of  
687 public-domain personality measures. *Journal of Research in Personality*, 40(1),  
688 84–96.

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

691 Grünerbl, A., Muaremi, A., Osmani, V., Bahle, G., Ohler, S., Tröster, G., ... Lukowicz, P.  
692 (2015). Smartphone-based recognition of states and state changes in bipolar disorder  
693 patients. *IEEE Journal of Biomedical and Health Informatics*, 19(1), 140–148.  
694 doi:[10.1109/JBHI.2014.2343154](https://doi.org/10.1109/JBHI.2014.2343154)

695 Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D.  
696 (2016). Using smartphones to collect behavioral data in psychological science:  
697 Opportunities, practical considerations, and challenges. *Perspectives on Psychological  
698 Science*, 11(6), 838–854. doi:[10.1177/1745691616650285](https://doi.org/10.1177/1745691616650285)

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

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

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

- 704 physical activity and sedentary behavior studies. *Exercise and Sport Sciences*  
705 *Reviews*, 43(1), 48–56. doi:[10.1249/JES.00000000000000035](https://doi.org/10.1249/JES.00000000000000035)
- 706 LaMarca, A., Chawathe, Y., Consolvo, S., Hightower, J., Smith, I., Scott, J., ... Schilit, B.  
707 (2005). Place lab: Device positioning using radio beacons in the wild. In *Pervasive  
708 computing* (pp. 116–133). Springer, Berlin, Heidelberg. doi:[10.1007/11428572\\_8](https://doi.org/10.1007/11428572_8)
- 709 Location History, G. (2017). *Timeline*. Retrieved from  
710 <https://www.google.com/maps/timeline?pb>
- 711 London, I. (2016). *Encoding cyclical continuous features - 24-hour time*. *Ian london's blog*.  
712 Retrieved April 27, 2018, from  
713 [//ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/](https://ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/)
- 714 Müller, K. (2017). *Bindrcpp: An 'rcpp' interface to active bindings*. Retrieved from  
715 <https://CRAN.R-project.org/package=bindrcpp>
- 716 Neuwirth, E. (2014). *RColorBrewer: ColorBrewer palettes*. Retrieved from  
717 <https://CRAN.R-project.org/package=RColorBrewer>
- 718 Palmius, N., Tsanas, A., Saunders, K. E. A., Bilderbeck, A. C., Geddes, J. R., Goodwin, G.  
719 M., & Vos, M. D. (2017). Detecting bipolar depression from geographic location data.  
720 *IEEE Transactions on Biomedical Engineering*, 64(8), 1761–1771.  
721 doi:[10.1109/TBME.2016.2611862](https://doi.org/10.1109/TBME.2016.2611862)
- 722 Patterson, T. A., Thomas, L., Wilcox, C., Ovaskainen, O., & Matthiopoulos, J. (2008).  
723 State-space models of individual animal movement. *Trends in Ecology & Evolution*,  
724 23(2), 87–94. doi:[10.1016/j.tree.2007.10.009](https://doi.org/10.1016/j.tree.2007.10.009)
- 725 Pebesma, E. J., & Bivand, R. S. (2005). Classes and methods for spatial data in R. *R News*,  
726 5(2), 9–13. Retrieved from <https://CRAN.R-project.org/doc/Rnews/>
- 727 Preisler, H. K., Ager, A. A., Johnson, B. K., & Kie, J. G. (2004). Modeling animal  
728 movements using stochastic differential equations. *Environmetrics* 15: P. 643-657.  
729 Retrieved from <https://www.fs.usda.gov/treesearch/pubs/33038>
- 730 R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna,

731 Austria: R Foundation for Statistical Computing. Retrieved from

732 <https://www.R-project.org/>

733 Sadilek, A., & Krumm, J. (2016). Far out: Predicting long-term human mobility. *Microsoft*  
734 *Research*. Retrieved from <https://www.microsoft.com/en-us/research/publication/far-predicting-long-term-human-mobility/>

735  
736 Saeb, S., Zhang, M., Karr, C. J., Schueller, S. M., Corden, M. E., Kording, K. P., & Mohr, D.  
737 C. (2015). Mobile phone sensor correlates of depressive symptom severity in daily-life  
738 behavior: An exploratory study. *Journal of Medical Internet Research*, 17(7), e175.  
739 doi:[10.2196/jmir.4273](https://doi.org/10.2196/jmir.4273)

740 Schipperijn, J., Kerr, J., Duncan, S., Madsen, T., Klinker, C. D., & Troelsen, J. (2014).  
741 Dynamic accuracy of GPS receivers for use in health research: A novel method to  
742 assess GPS accuracy in real-world settings. *Frontiers in Public Health*, 2, 21.  
743 doi:[10.3389/fpubh.2014.00021](https://doi.org/10.3389/fpubh.2014.00021)

744 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology:  
745 Undisclosed flexibility in data collection and analysis allows presenting anything as  
746 significant. *Psychological Science*, 22(11), 1359–1366. doi:[10.1177/0956797611417632](https://doi.org/10.1177/0956797611417632)

747 Sobrado, B. (2018). *My thesis on missing data in GPS measurements*. Retrieved from  
748 <https://github.com/sobradob/thesis>

749 ST-DBSCAN: An algorithm for clustering spatial-temporal data. (2007). *Data & Knowledge*  
750 *Engineering*, 60(1), 208–221. doi:[10.1016/j.datwk.2006.01.013](https://doi.org/10.1016/j.datwk.2006.01.013)

751 Tatsiramos, K. (2009). Geographic labour mobility and unemployment insurance in europe.  
752 *Journal of Population Economics*, 22(2), 267–283. doi:[10.1007/s00148-008-0194-7](https://doi.org/10.1007/s00148-008-0194-7)

753 Thoen, E. (2017). *Padr: Quickly get datetime data ready for analysis*. Retrieved from  
754 <https://CRAN.R-project.org/package=padr>

755 Vaughan, D., & Dancho, M. (2018). *Tibbletime: Time aware tibbles*. Retrieved from  
756 <https://CRAN.R-project.org/package=tibbletime>

757 Villanueva, C., & Aggarwal, B. (2013). The association between neighborhood

- 758 socioeconomic status and clinical outcomes among patients 1 year after  
759 hospitalization for cardiovascular disease. *Journal of Community Health*, 38(4),  
760 690–697. doi:[10.1007/s10900-013-9666-0](https://doi.org/10.1007/s10900-013-9666-0)
- 761 Wang, R., Harari, G., Hao, P., Zhou, X., & Campbell, A. T. (2015). SmartGPA: How  
762 smartphones can assess and predict academic performance of college students. In  
763 *Proceedings of the 2015 ACM international joint conference on pervasive and*  
764 *ubiquitous computing* (pp. 295–306). New York, NY, USA: ACM.  
765 doi:[10.1145/2750858.2804251](https://doi.org/10.1145/2750858.2804251)
- 766 Wickham, H. (2009). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.  
767 Retrieved from <http://ggplot2.org>
- 768 Wickham, H. (2017). *Scales: Scale functions for visualization*. Retrieved from  
769 <https://CRAN.R-project.org/package=scales>
- 770 Wickham, H., & Henry, L. (2018). *Tidyr: Easily tidy data with 'spread()' and 'gather()'*  
771 *functions*. Retrieved from <https://CRAN.R-project.org/package=tidyr>
- 772 Wickham, H., & Ruiz, E. (2018). *Dbplyr: A 'dplyr' back end for databases*. Retrieved from  
773 <https://CRAN.R-project.org/package=dbplyr>
- 774 Wickham, H., Francois, R., Henry, L., & Müller, K. (2017). *Dplyr: A grammar of data*  
775 *manipulation*. Retrieved from <https://CRAN.R-project.org/package=dplyr>
- 776 Wickham, H., Hester, J., & Francois, R. (2017). *Readr: Read rectangular text data*.  
777 Retrieved from <https://CRAN.R-project.org/package=readr>
- 778 Wolf, J., Oliveira, M., & Thompson, M. (2003). Impact of underreporting on mileage and  
779 travel time estimates: Results from global positioning system-enhanced household  
780 travel survey. *Transportation Research Record: Journal of the Transportation*  
781 *Research Board*, 1854, 189–198. doi:[10.3141/1854-21](https://doi.org/10.3141/1854-21)
- 782 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida:  
783 Chapman; Hall/CRC. Retrieved from <https://yihui.name/knitr/>
- 784 Zenk, S. N., Schulz, A. J., & Odoms-Young, A. (2009). How neighborhood environments

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

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

787 Zhang, Z., Yang, X., Li, H., Li, W., Yan, H., & Shi, F. (2017). Application of a novel hybrid  
788 method for spatiotemporal data imputation: A case study of the minqin county  
789 groundwater level. *Journal of Hydrology*, 553(Supplement C), 384–397.

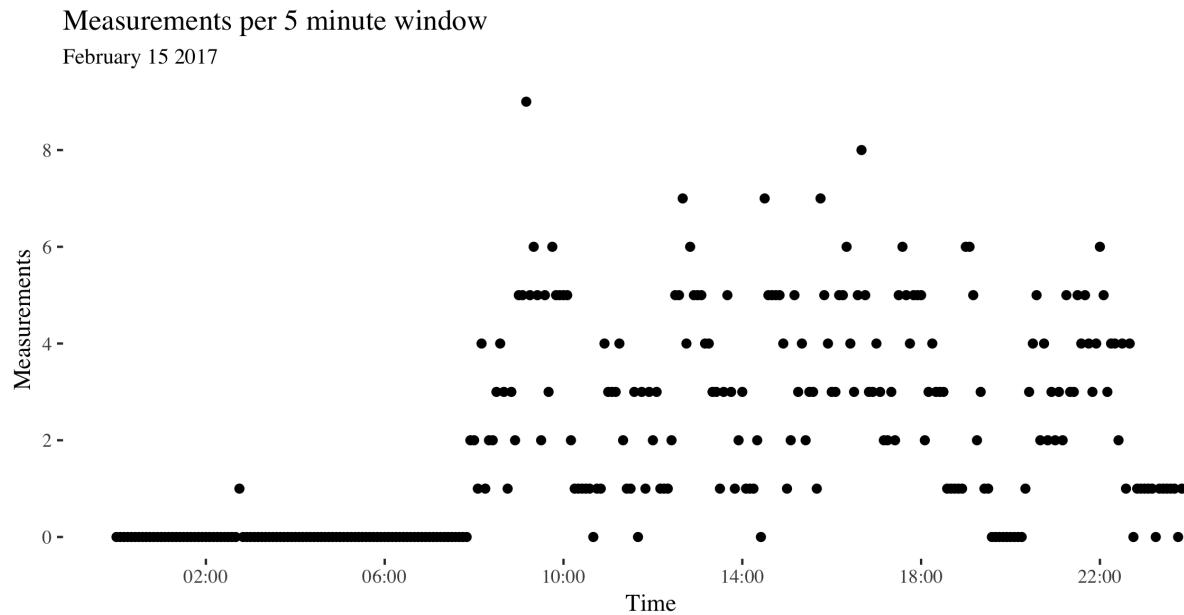
790 doi:10.1016/j.jhydrol.2017.07.053

791 Zhu, H. (2018). *KableExtra: Construct complex table with 'kable' and pipe syntax*. Retrieved  
792 from <https://CRAN.R-project.org/package=kableExtra>

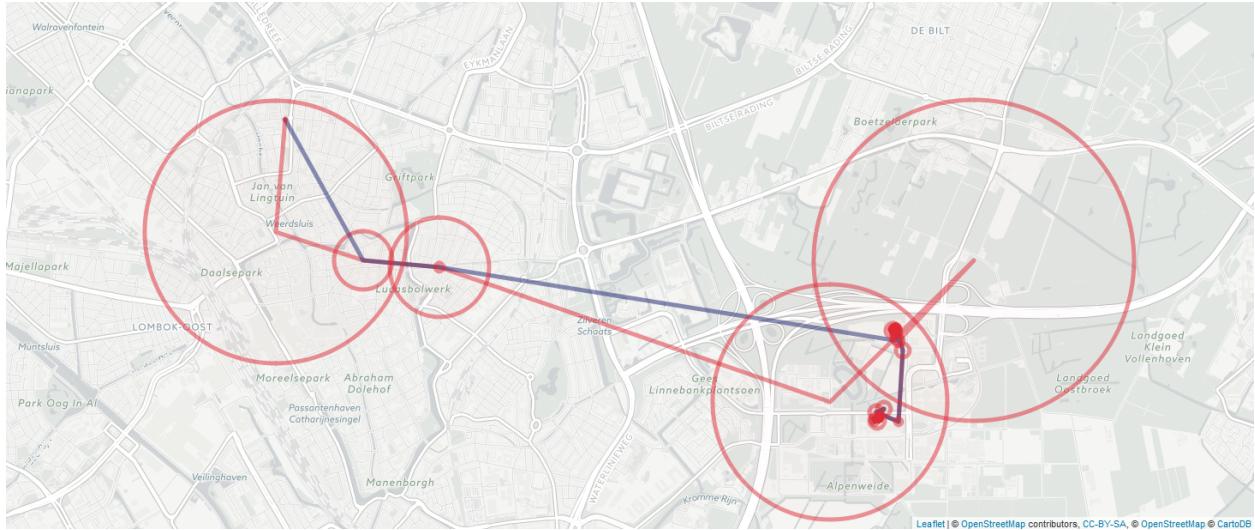
Table 1

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

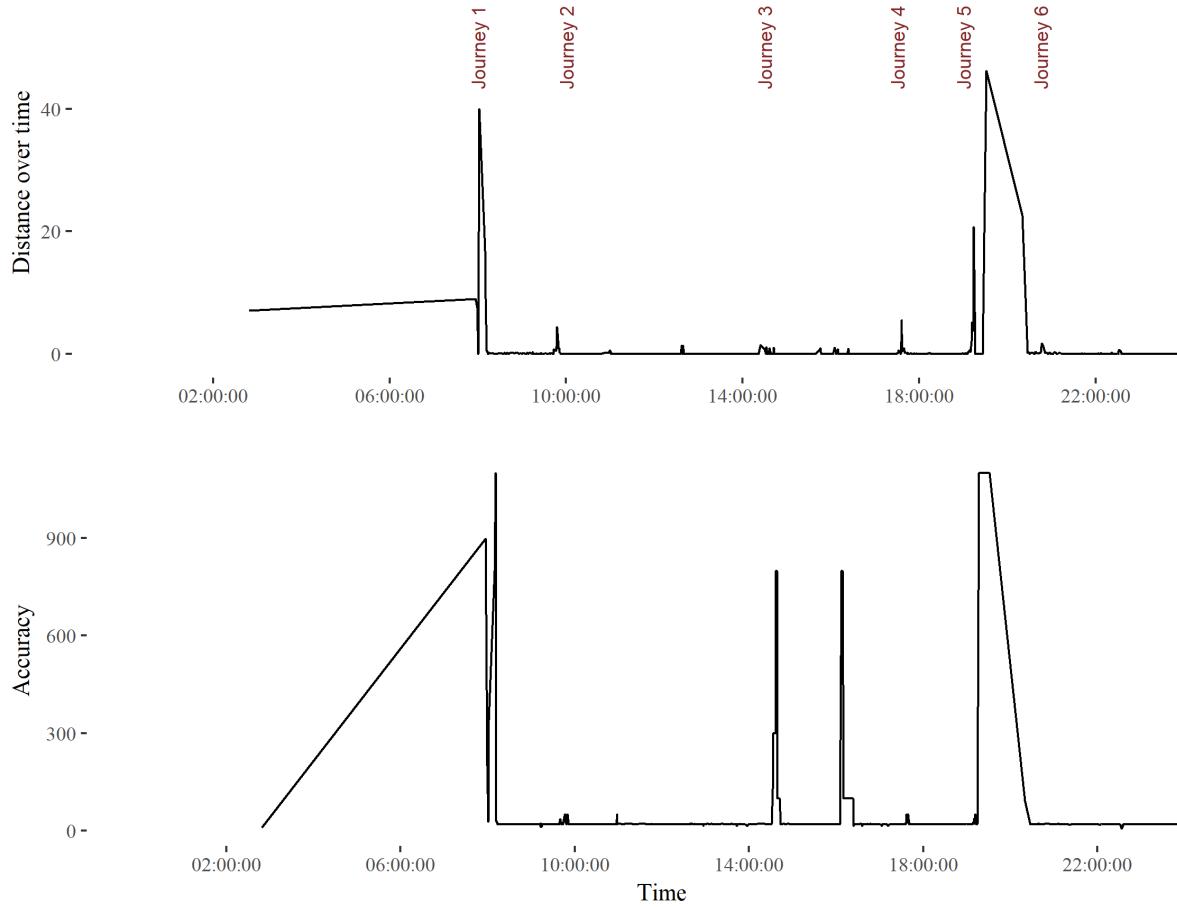
	Five minutes		One Hour		One Day	
	Mean	Median	Mean	Median	Mean	Median
Barnett & Onella	82	4	345	6	9,273	12
Palmius	43	0	497	4	NA	NA
PPMI	269	0	908	0	5,757	0
Naive Baseline	426	0	1,502	0	14,266	1,288



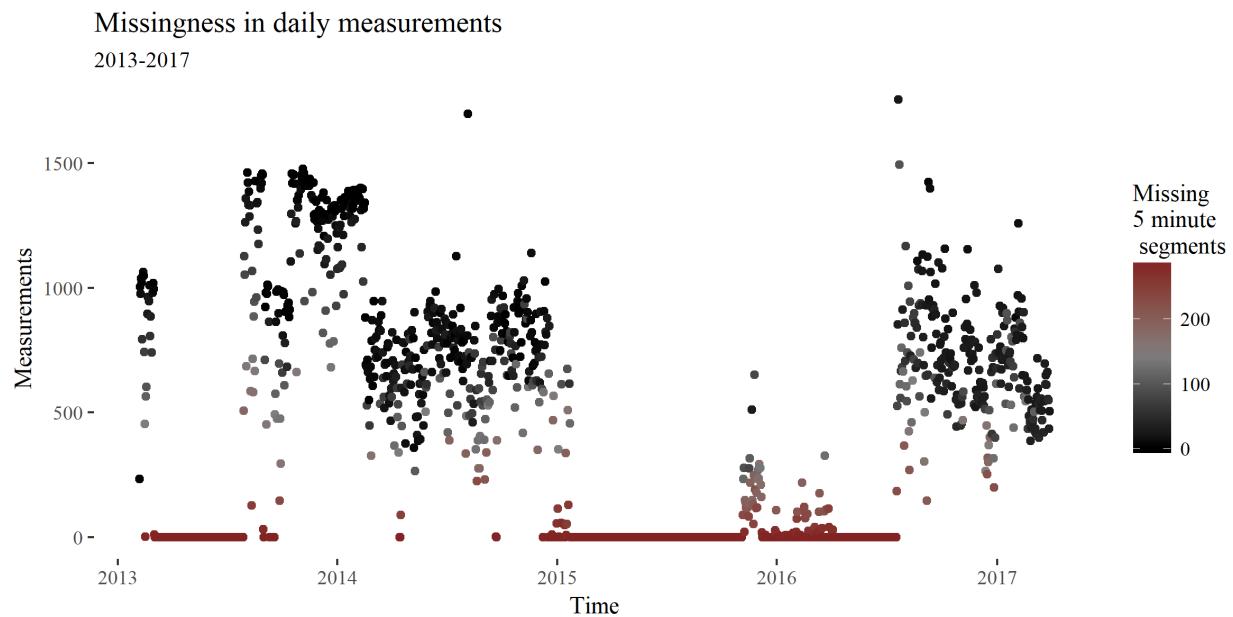
*Figure 1.* Example of missing data over the entire duration of a secondary log. The x-axis denotes time, the y-axis shows how many measurements are made and each point is a five minute window. For this day there were several periods with no information. These points lie on the x-axis.



*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.

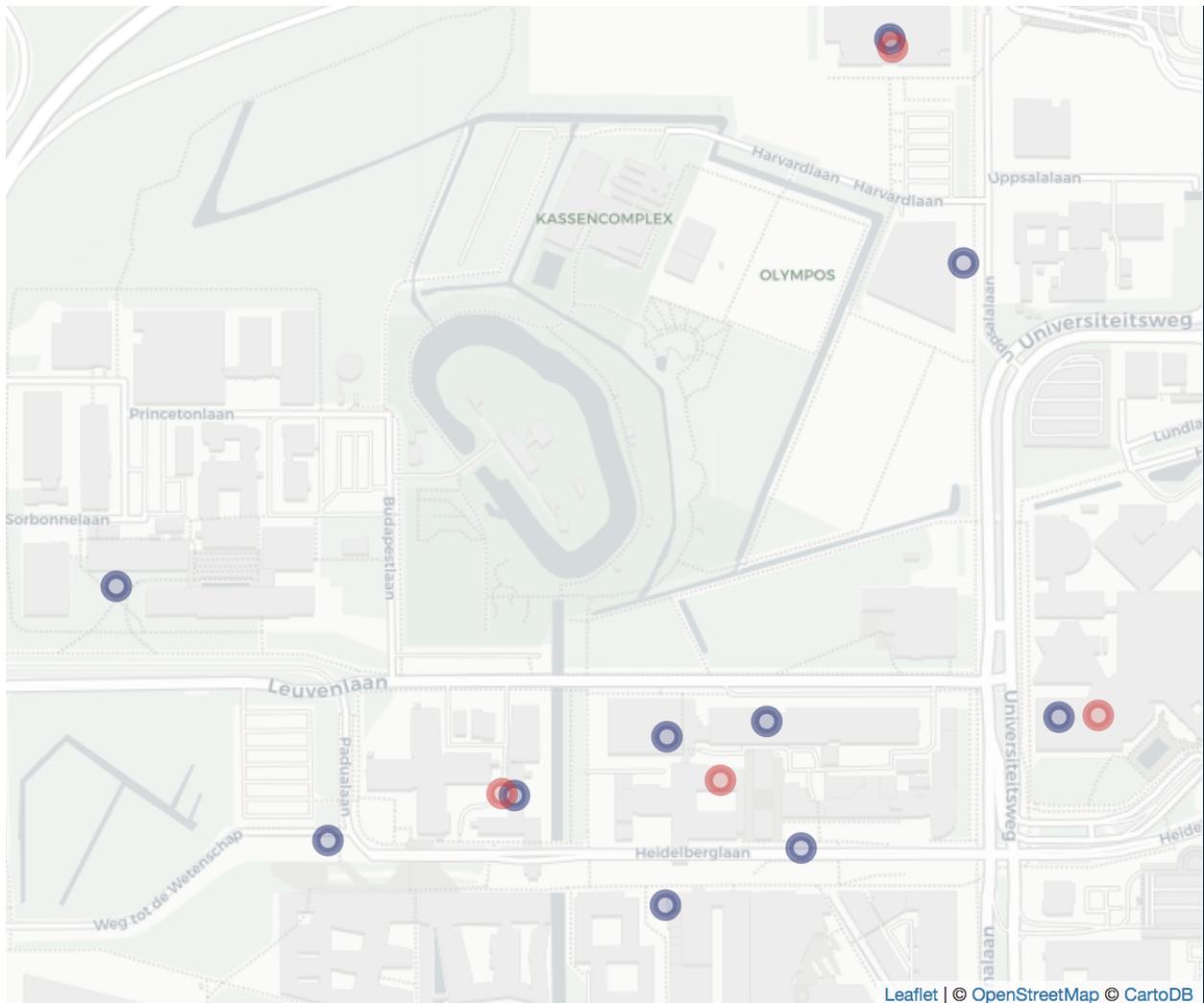


Figure 5. Example of pause locations in De Uithof university campus using 150 meters (blue) and 400 meters (red) as clustering parameters.

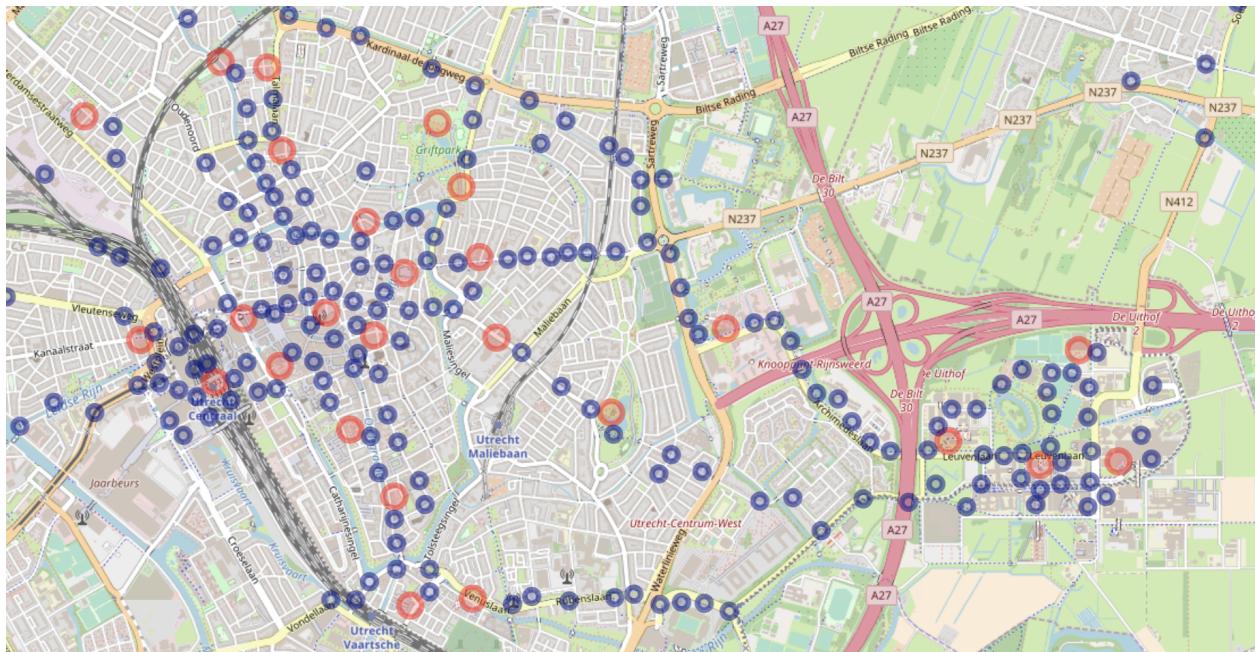
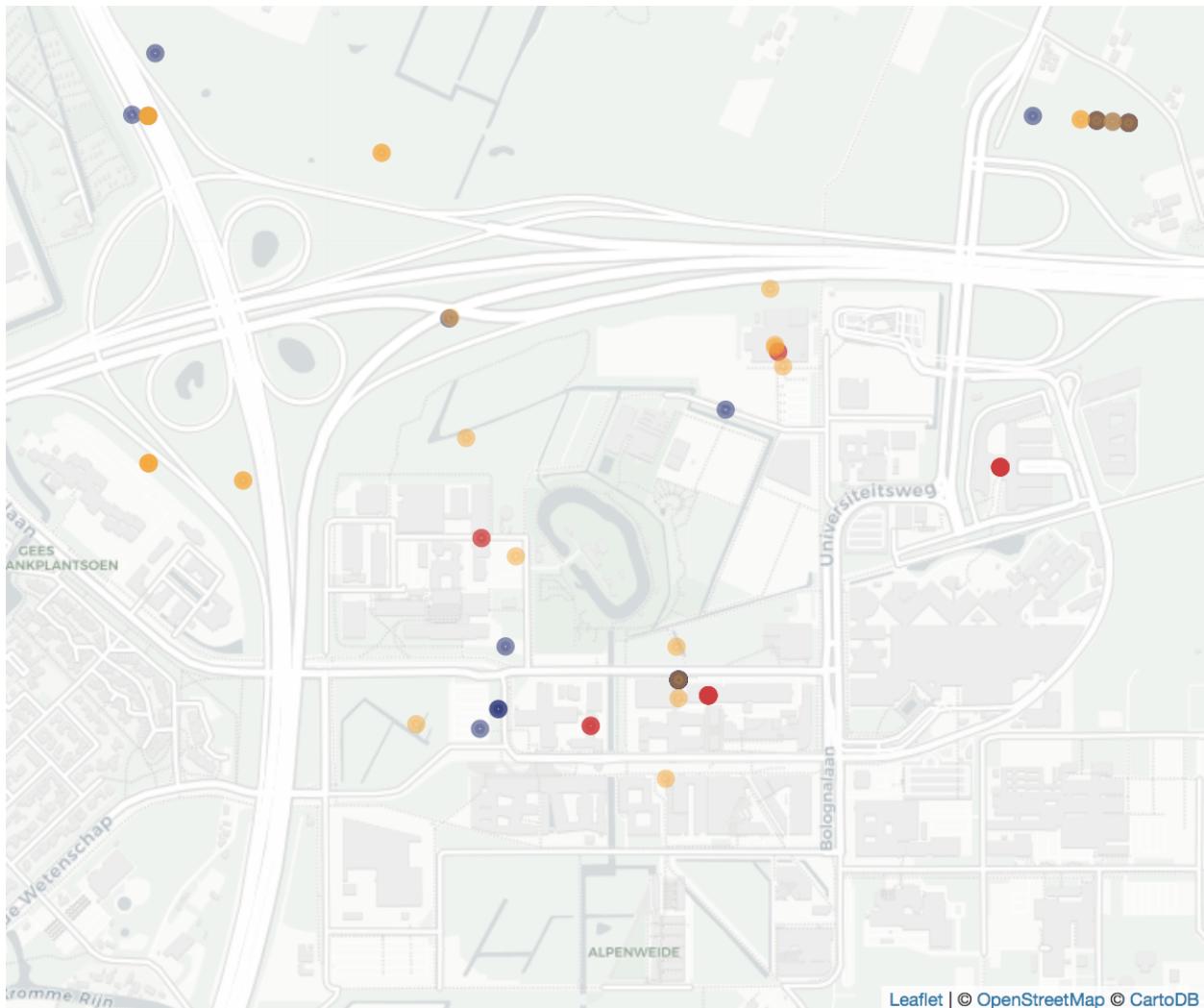
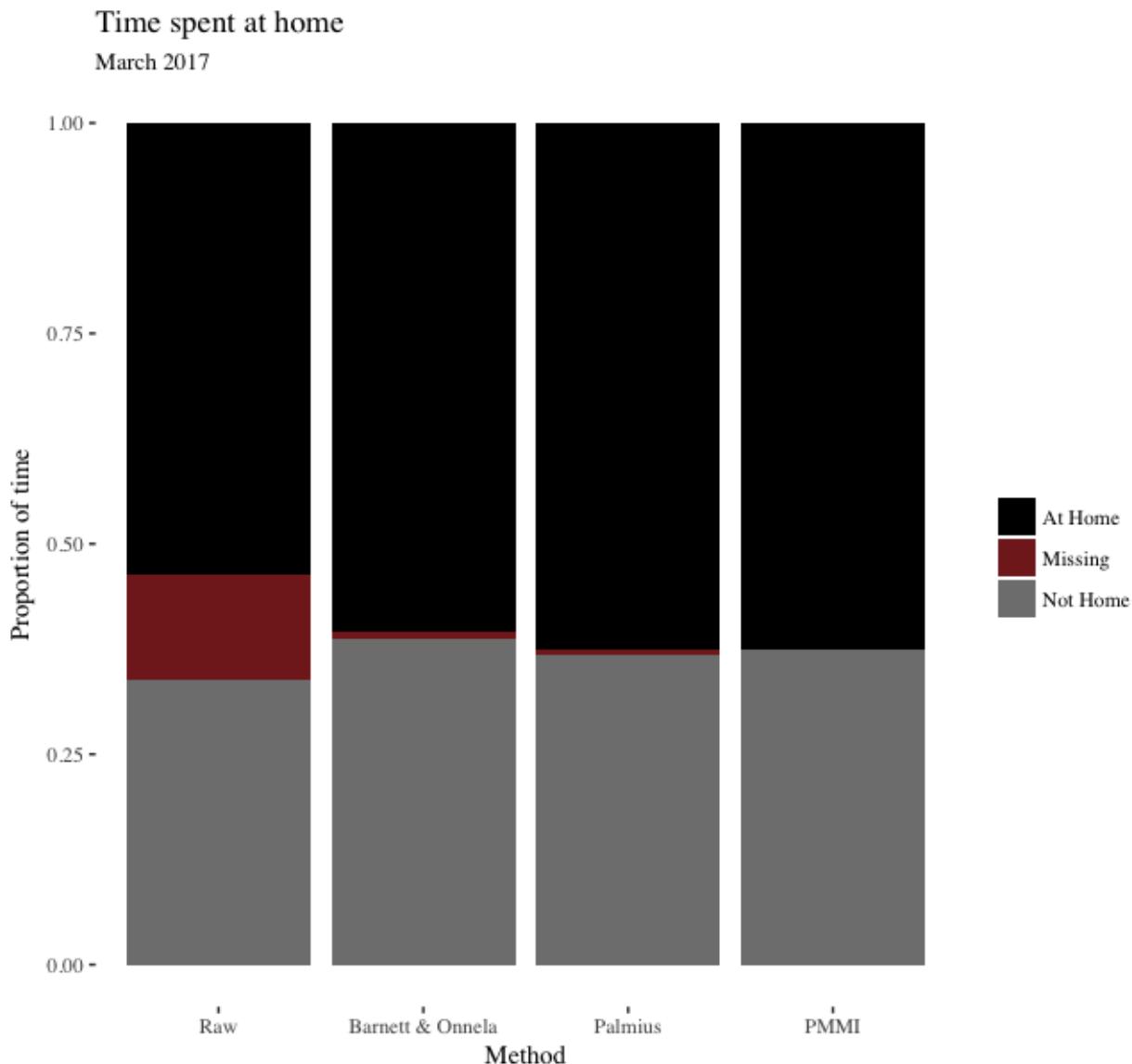


Figure 6. Excerpt of the cluster map of an individual. Red points are pause locations, blue points are path locations.



*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.