

1 Handling missing data in smartphone location logs

2 Boaz Sobrado<sup>1</sup>

3 <sup>1</sup> Utrecht University

4 Author Note

5 Department of Methodology & Statistics

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

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

8 [boaz@boazsobrado.com](mailto:boaz@boazsobrado.com)

9

## Abstract

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

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

30 Word count: 2297

31 Handling missing data in smartphone location logs

32 **Introduction**

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

53 Despite the importance of mobility measures, the majority of these measures are  
54 obtained through the use of questionnaires (such as the IPIP) or specifically through  
55 pen-and-paper travel diaries. Questionnaires and travel diaries are have well-known  
56 methodological flaws: they are burdensome to collect and rely on accurate self-reporting.  
57 The burdensomeness stems from the fact that participants must be explicitly asked to report

58 on their movement patterns at frequent intervals due to forgetfulness. This makes the data  
59 expensive to collect and limits the extent of data collection in practice. In addition, the  
60 frequent reporting duties of the participant may bias the participants behaviour. Participant  
61 forgetfulness also limits the accuracy of self reported measures. There is clear evidence that  
62 participants are systematically biased when self-reporting mobility measures. For instance,  
63 participants under-report the frequency of short trips (Wolf, Oliveira, & Thompson, 2003)  
64 and underestimate the duration of regular commutes (Delclòs-Alió, Marquet, &  
65 Miralles-Guasch, 2017). These obstacles can be overcome by using objective data on human  
66 mobility.

67 Social scientists now have the unprecedented opportunity to easily obtain objective  
68 data on human mobility from smartphones. Before smartphones the only way to collect  
69 objective data on human mobility involved giving participants an expensive  
70 professional-grade location sensor and convincing them to take it with themselves at all  
71 times. Barnett and Onnela (2016) points out that introducing a new device to the  
72 participant's life may bias their behaviour. Moreover, collecting data in such a way is costly,  
73 places a high burden on participants and therefore the logs do not span a long time [Barnett  
74 and Onnela (2016)]. Today millions of individuals carry smartphones with themselves every  
75 day and do not need to be encouraged to do so by researchers. Smartphones are equipped  
76 with a range of sensors that can be used to track the location of the device at all times.  
77 These smartphones can collect and store hundreds of location measurements a day. For  
78 instance, Google Location History contains movement information on millions of users, often  
79 spanning years (Location History, 2017). Moreover, recent changes in EU regulations with  
80 regard to consumer data-portability rights ensure that a willing participant should be able to  
81 easily share this information with researchers at no cost to either the participant or the  
82 researchers (Commission, 2017). Taken together, this means that researchers now have at  
83 their disposal the ability to easily access the objective human mobility data of millions of  
84 individuals spanning several years at little cost and without a significant burden of the

85 participants.

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

93 **Background**

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

106 Researchers often develop custom-made tracking applications which participants are  
107 instructed to download on their phone. Alternatively participants are given a phone to use  
108 for a given period of time with the custom-made tracking app pre-installed. We call location  
109 logs resulting from these custom-made apps *custom logs*. The advantage of custom logs is  
110 that the researchers can adjust tracking parameters, such as the frequency of

111 measurements and the sensor with which they are made. The disadvantage with this  
112 approach is that researchers have to develop or adapt a custom-made tracking application  
113 (which is not easy given hundreds of different types of smartphone models), distribute it  
114 among research participants and enforce participation. Participants may dislike tracking  
115 apps because they view them as more intrusive and these apps regularly drain the battery of  
116 the device (G. M. Harari et al., 2016). Moreover, researchers have distribute this application  
117 among research participants and convince them not to turn it off.

118 We focus on another solution, which is to take advantage of existing smartphone  
119 location logs (*secondary logs*) . The advantages are clear: repositories such as Google's  
120 Location History contains information on millions of users spanning years (Location History,  
121 2017), participants can share the data by the click of a button and there can be no  
122 behavioural changes due to participation in the study as the participant share past data.  
123 The disadvantage is that researchers have no control over the tracking parameters, often  
124 resulting in logs with sparse and inaccurate measurements. Hence, two important challenges  
125 are dealing with missing data and measurement noise.

126 In order to work with secondary logs, researchers need to be able to handle the data  
127 sparsity that leads to missing data. Missing data is a pervasive issue in secondary logs as it  
128 can arise due to several reasons. Technical reasons include signal loss, battery failure and  
129 device failure. Behavioural reasons include leaving the phone at home or switching the  
130 device off. As a result, secondary logs often contain wide temporal gaps with no  
131 measurements. For instance, several research groups studying mental health report missing  
132 data rates between 30% to 50% (Grünerbl et al., 2015; Palmius et al., 2017; Saeb et al.,  
133 2015). Other researchers report similar trends in different fields (e.g. G. M. Harari et al.,  
134 2016; Jankowska, Schipperijn, & Kerr, 2015). In Figure 1 shows that despite the long  
135 duration of the log the sparsity it is also evident.

136 There is no golden standard for dealing with missing data in GPS logs (Barnett &  
137 Onnela, 2016). Importantly, spatiotemporal data measurements are auto-correlated in both

time and space. This means that best practices with other types of data, such as mean imputation, are unsuitable. For example, imagine an individual who splits almost all her time between work and home. Suppose she spends a small amount of time commuting between the two along a circular path. Using mean imputation to estimate her missing coordinates, we impute her to be at the midpoint between home and work, even though she has never been there. Worryingly, there is little transparency on how researchers deal with missing data (Jankowska et al., 2015).

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

Noisy data can lead to inaccurate conclusions if it is not accounted for. Suppose we wish to calculate an individual's movement in a day. A simple approach would be to calculate the sum of the distance between each measurement. But if there is noise, the coordinates will vary even though the individual is not moving. If the measurements are frequent and noisy, we will calculate a lot of movement, even if the individual did not move at all! This issue is also visualised in Figure ???. The problem is further complicated because missing data and noisy measurements are related. Methods used by researchers to reduce noise, such as throwing out inaccurate measurements (e.g. Palmius et al., 2017), can exacerbate the severity of the missing data problem.

In this paper we will propose PPMI as a method for dealing with missing data and measurement error in secondary location logs. We will compare PPMI to similar solutions in

the literature by evaluating the distance between points which were simulated as missing and their imputed counterparts. Finally we will calculate time spent at home as a function of the imputation method.

168

## Related Work

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

A lack in methods for missing data imputation for human mobility patterns does not imply there is not a vast literature on modelling movement. The most widespread models are SSMs, therefore we shall detail a few examples and subsequently argue that they are nonetheless unsuited for long term human mobility logs. Ecologists have used SSMs to explain how animals interact with their environment (Patterson, Thomas, Wilcox, Ovaskainen, & Matthiopoulos, 2008). These models can be quite complex. Preisler, Ager, Johnson, and Kie (2004) uses Markovian movement processes to characterise the effect of roads, food patches and streams on cyclical elk movements. The most well studied SSM is the Kalman filter, which is the optimal algorithm for inferring linear Gaussian systems. The extended Kalman filter is the de facto standard for GPS navigation (Z. Chen & Brown,

191 2013). The advantage of state space models is that they are flexible, deal with measurement  
192 inaccuracy, include information from different sources and can be used in real time.

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

204 In the wider realm of spatiotemporal statistics there are numerous missing data  
205 imputation methods. These often come from climate or geological research and rely on  
206 spatiotemporal auto-correlations. For instance, the CUTOFF method estimates missing  
207 values by incorporating similar observed temporal information from the value's nearest  
208 spatial neighbors (Feng, Nowak, O'Neill, & Welsh, 2014 ). The authors illustrate their  
209 example using rainfall data from gauging stations across Australia. Similarly, Z. Zhang et al.  
210 (2017) use a variety of machine learning methods to impute missing values. The example  
211 provided relates to underground water data. Generally these models assume fixed  
212 measurement stations (such as rainfall gauging stations). For this reason they cannot be  
213 easily applied to missing mobility tracks without significant pre-processing.

214 On the other hand, a few researchers have explicitly attempted to impute missing data  
215 from human mobility patterns (???: Barnett & Onnela, 2016; Palmius et al., 2017 ).  
216 Importantly, none of them worked with secondary logs. Nonetheless we will detail what they  
217 did as informative examples. Palmius et al. (2017) deal with the measurement inaccuracy of

<sup>218</sup>  $D$  in custom logs by removing from the data set all unique low-accuracy  $a$  data points that  
<sup>219</sup> had  $\frac{d}{dt}D > 100 \frac{km}{h}$ . Subsequently the researchers down sample the data to a sample rate of  
<sup>220</sup> 12 per hour using a median filter. Moreover, Palmius et al. (2017) explain:

<sup>221</sup> “If the standard deviation of  $[D]$  in both latitude and longitude within a 1 h  
<sup>222</sup> epoch was less than 0.01 km, then all samples within the hour were set to the  
<sup>223</sup> mean value of the recorded data, otherwise a 5 min median filter window was  
<sup>224</sup> applied to the recorded latitude and longitude in the epoch”.

<sup>225</sup> Missing data was imputed using the mean of measurements close in time if the  
<sup>226</sup> participant was recorded within 500m of either end of a missing section and the missing  
<sup>227</sup> section had a length of  $\leq 2h$  or  $\leq 12h$  after 9pm. In cases where the previous conditions are  
<sup>228</sup> not met no values are imputed.

<sup>229</sup> Barnett and Onnela (2016) follow a different approach which is, to the best of our  
<sup>230</sup> knowledge, the only principled approach to dealing with missing data in human mobility  
<sup>231</sup> data. Barnett and Onnela (2016) work with custom logs where location is measured for 2  
<sup>232</sup> minutes and subsequently not measured for 10 minutes. In the words of the authors, Barnett  
<sup>233</sup> and Onnela (2016) handle missing data by first converting data to mobility traces, which are  
<sup>234</sup> defined as a sequence of flights and pauses. Flights are segments of linear movements and  
<sup>235</sup> pauses corresponding to periods of time where a person does not move. Subsequently, the  
<sup>236</sup> authors impute missing data by:

<sup>237</sup> “simulat[ing] flights and pauses over the period of missingness where the direction,  
<sup>238</sup> duration, and spatial length of each flight, the fraction of flights versus the  
<sup>239</sup> fraction of pauses, and the duration of pauses are sampled from observed data.”

<sup>240</sup> This method can be extended to imputing the data based on temporally, spatially or  
<sup>241</sup> periodically close flights and pauses. In other words, for a given missing period, the  
<sup>242</sup> individual’s mobility can be estimated based on measured movements in that area, at that  
<sup>243</sup> point in time or movements in the last 24 hours (*circadian proximity*).

244 On the other hand, (???) use what they call a Spatial Temporal Semantic Neural  
 245 Network (STS-NN) to predict future human movement. While the authors are concerned  
 246 with prediction and not imputation, they devised a method called the Spatial Temporal  
 247 Semantic (STS) algorithm which converts raw measurements to machine learning friendly  
 248 discrete bins. Working with high-frequency measurements, (???)’s method down-samples the  
 249 raw data temporally and map-matches the resulting bins to discrete points along  
 250 pre-established geographical features such as roads and highways. This minimises  
 251 measurement error and paves the way for applying machine learning methods to human  
 252 mobility problems.

253 In this section we have argued that there is a lack of established practices to follow  
 254 with respect to missing data in human mobility logs. Moreover, extensively used  
 255 spatiotemporal methods, such as state space models (SSMs), are not well suited to deal with  
 256 human mobility patterns in secondary logs. Finally we discussed in detail three approaches  
 257 which deal explicitly with mobility patterns from custom or secondary logs (???: Barnett &  
 258 Onnela, 2016; Palmius et al., 2017 ).

## 259 Methodology

### 260 Notation

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

266 Subsequently let us simplify by assuming that a persons location is on two-dimensional  
 267 Euclidean plane. Let a person’s true location on this two-dimensional plane be  
 268  $G(t) = [G_x(t)G_y(t)]$  where  $G_x(t)$  and  $G_y(t)$  denote the location of the individual at time  $t$  on  
 269 the x-axis and y-axis respectively. For simplicity, we can assume that the x-axis is the

270 longitude and the y-axis is the latitude. Moreover, let  $D \in \mathbb{R}^2$  be the recorded data  
271 containing the latitude and longitude. In addition, let  $a$  denote the estimated accuracy of  
272 the recorded data. Furthermore,  $G(t)$ ,  $D$  and  $a$  are indexed by time labeled by the countable  
273 set  $t = t_1 < \dots < t_{n+1}$ . The measure of accuracy  $a_t$  is given in meters such that it represents  
274 the radius of a 67% confidence circle. If  $D_t = \emptyset$  it is considered *missing* and it is not missing  
275 otherwise.

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

## 278 Personalised Map Matching Imputation

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

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

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

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

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

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

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

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

289 We define pause locations  $G(t_a), G(t_b) | t_a, t_b \in P$  as locations where an individual spends an  
290 extended amount of time in the same space (e.g. school, home, work, train station, barber  
291 shop, bar, gym). Importantly, our model assumes period and cyclic human movement such

<sup>292</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>293</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>294</sup> leaves home for a run and returns home without stopping anywhere else.

<sup>295</sup> Let us define as  $Flight_{ab}^x$  the set of all points belonging to a flight between  $G(t_a)$  and  
<sup>296</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>297</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>298</sup> Then, we can define as  $Path_{ab}$  the set of all flights between  $G(t_a)$  and  $G(t_b)$  at all  
<sup>299</sup> time-points. For simplicity, we assume that  $Path_{ab} = Path_{ba}$ .

<sup>300</sup> In addition, we consider all measurements  $D(t)$  to be imperfect measurements of  $G(t)$ :

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

### <sup>301</sup> Personalised Map Matching Imputation algorithm

<sup>302</sup> Our algorithm performs the following steps:

- <sup>303</sup> 1. *Map building*: Extract from measurements  $D$  all pause location bins and path location  
<sup>304</sup> bins to create a personalised map.
- <sup>305</sup> 2. *Binning*: Assign each measurement  $D$  to a unique discrete location bin.
- <sup>306</sup> 3. *Imputing*: Use a classification method to predict missing measurements based on all  
<sup>307</sup> the available information.

<sup>308</sup> **Map building.** Following (???)'s spatial-temporal-semantic (STS) feature  
<sup>309</sup> extraction algorithm our aim is to transform pause and path locations into machine learning  
<sup>310</sup> friendly discrete location sequences. There are multiple ways of extracting such measurement  
<sup>311</sup> clusters in the literature, such as Spatio-Temporal Density-Based Spatial Clustering of  
<sup>312</sup> Applications with Noise (ST-DBSCAN) and sequence oriented clustering (SOC)  
<sup>313</sup> [("ST-DBSCAN," 2007)]. We will focus on two methods which explicitly with mobility

<sup>314</sup> patterns from unevenly sampled smartphone logs (Barnett & Onnela, 2016; Palmius et al.,  
<sup>315</sup> 2017). Both of these methods pre-process the data and subsequently use two steps to  
<sup>316</sup> extract pause locations: first they extract pauses and their corresponding locations, then  
<sup>317</sup> they cluster pause locations based on spatial proximity. This implementation of PMMI uses  
<sup>318</sup> a stricter version of Barnett and Onnela (2016)'s approach to extract pauses.

<sup>319</sup> First the measurements  $D$  are filtered such that only measurements with an accuracy  
<sup>320</sup> value lower than  $a_{P\lim}$  remain within the sample. Then, a measurement  $D_t$  belongs to a  
<sup>321</sup> pause if and only if:

- <sup>322</sup> 1. The next measurement  $D_{t+1}$  is within  $t_{P\lim}$  amount of seconds (so it is not missing)
- <sup>323</sup> 2. The next measurement  $D_{t+1}$  is within  $d_{P\lim}$  meters.
- <sup>324</sup> 3. The duration of the pause is more than  $\delta_{P\lim}$  seconds.
- <sup>325</sup> 4. Let the measurements of a possible pause which fit the aforementioned criteria be  
<sup>326</sup>  $D_{t,t+1,\dots,t+n}$ . These points are only a pause if the distance between the mean  
<sup>327</sup> 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  
<sup>328</sup> the mean accuracy  $a$  of  $D_{t,t+1,\dots,t+n}$ .

<sup>329</sup> This set of points were then hierarchically clustered using a distance matrix, such that  
<sup>330</sup> all points within  $d$  meters of each other were clustered into a pause location. Each pause  
<sup>331</sup> location is a bin.

<sup>332</sup> For all remaining measurements we assume that they belong to paths. In this  
<sup>333</sup> implementation of PMMI we use the following algorithm to estimate paths:

- <sup>334</sup> 1. Take all measurements which are not pauses, filter them based on an accuracy  
<sup>335</sup> threshold  $a_{Path\ Lim}$ .
- <sup>336</sup> 2. Create a distance matrix for all remaining measurements  $D_t \in F$  and hierarchically  
<sup>337</sup> cluster it accordingly, such that all points within  $d_{Path\ Lim}$  meters of each other are  
<sup>338</sup> clustered into a single pause point.

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

346 **Binning.** (???)'s spatial-temporal-semantic (STS) feature extraction algorithm uses  
 347 map matching as a ground truth to assign noisy measurements into discrete bins along roads.  
 348 In other words, in addition to the measured data they also use a geographic database that  
 349 contains information about the area in which the individual is (e.g. where precisely the roads  
 350 are), and sort measurements into bins based on both the measurement and the geographic  
 351 data base. For example, if an individual is measured as moving closely in parallel to a road  
 352 A in an area where there is no other parallel road, (???)'s method will assume that the  
 353 individual is on the road A.

354 PMMI uses a similar logic, but without using any external geographic database. The  
 355 key modification in PMMI is that whilst (???) uses a map from outside the persons location  
 356 logs, we use the total location history of the individual to create a personalised map. This  
 357 map is subsequently used to bin measurements. This is feasible for two reasons: humans  
 358 tend to have repetitive movement habits and secondary logs tend to be long. To put it in  
 359 simpler terms, we consider each measurement at  $D_x$  as a sample of  $Path_{ab}$ , and by  
 360 aggregating many measurements we can use them to map out  $Path_{ab}$ .

361 Thus, all measurements  $D$  were are then assigned to a discrete bin on the personal  
 362 map. This includes previous measurements which were discarded from the map building  
 363 exercise due to an accuracy  $a$  value which exceeded  $a_{P\lim}$  or  $a_{F\lim}$ . In this implementation  
 364 we used a simple assigning function, whereby the measurements where assigned to the bin  
 365 nearest to the measurement.

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

- 371 1. The previous and subsequent observed bin as a binary class matrix.
- 372 2. The distance in time to the next & previous bin.
- 373 3. The time of the day encoded as a cyclical two-dimensional feature.
- 374 4. The day of the month as a binary class matrix.
- 375 5. The month of the year as a binary class matrix.

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

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

### 385 Datasets & Analyses

386 The secondary location log used to train the imputation methods was collected  
387 between 2013 and 2017 on different Android devices from a single individual. About 54% of  
388 the data is missing for the entire duration of the log. This may be misleading as there are  
389 several long periods with no measurements whatsoever. For days which were not entirely  
390 missing, approximately 22% of all five minute segments were missing. The structure of  
391 missingness of a day with measurements is shown in Figure 4. As you can see, there are

392 several long periods over the course of the log for which there are no measurements. The  
393 median sampling frequency per day for non-missing days is around 0.006 Hz.

394 For simplicity, we subsequently used a time period when the individual was living in  
395 the Netherlands. This leaves us with 156,000 measurements over a period of less than six  
396 months. Analyses as well as implementations and adaptions Palmius et al. (2017)'s and  
397 Barnett and Onnela (2016)'s model were performed using R and a multitude of other  
398 statistical packages R (Version 3.4.3; R Core Team, 2017) and the R-packages *ggplot2*  
399 (Version 2.2.1; Wickham, 2009), *ggthemes* (Version 3.4.0; Arnold, 2017), *knitr* (Version 1.20;  
400 Xie, 2015), and *papaja* (Version 0.1.0.9709; Aust & Barth, 2018). All the code is available on  
401 a public repository (???).

## 402 Results & Evaluation Metrics

403 The results will consist of multiple steps:

- 404 1. Evaluating the performance of the map building and assigning functions.
- 405 2. Comparing the performance of PMMI using a) baseline models b) performance with  
406 randomly removed data in comparison to the aforementioned methods (Barnett &  
407 Onnela, 2016; Palmius et al., 2017) and c) objective ground-truth data (public  
408 transportation time-stamps).

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

414 In order to evaluate the map building and binning we will first visually evaluate the  
415 paths and pause locations. A visual evaluation of paths superimposed on is an established  
416 way to heuristically check their accuracy (e.g. (??)). Then, let the average distance

<sup>417</sup> between the actual measured point and the binned point be the *deviation distance*  $\delta_{\text{dev}}$ .

<sup>418</sup> With respect to the deviation distance  $\delta_{\text{dev}}$ , we expect:

<sup>419</sup> 1. A positive relationship between the deviation distance and the accuracy of each  
<sup>420</sup> measurement.

<sup>421</sup> 2. Roughly 67% of the deviation distances  $\delta_{\text{dev}}$  are within accuracy  $a$  of each  
<sup>422</sup> measurement.

<sup>423</sup> **Imputation algorithm performance.** We will compare the performance of PMMI,

<sup>424</sup> Palmius et al. (2017) and Barnett and Onnela (2016). In addition, we will also compute a

<sup>425</sup> naive model, which simply imputes as the missing value the previous observed value. The

<sup>426</sup> naive model will serve as a baseline model. To compare the performance of these methods

<sup>427</sup> we will remove 25% of measured time intervals at random within a four week period. We will

<sup>428</sup> make our comparisons for intervals of 5 minutes, 1 hour and 1 day. In other words, we will

<sup>429</sup> remove 25% of time intervals at random, while varying the size of the time intervals removed.

<sup>430</sup> For the Barnett and Onnela (2016) and PMMI models we will use all the available

<sup>431</sup> data to train the models with the exception of the time periods being investigated. Palmius

<sup>432</sup> et al. (2017)'s model does not require training.

<sup>433</sup> To compare all methods with each other we will compute a distance measure (how far

<sup>434</sup> was the removed location from the predicted location). For PMMI's imputations we will use

<sup>435</sup> a weighted mean for the distance measures whereby each 5 minute period is weighted equally.

<sup>436</sup> This is necessary because the other two model estimate far fewer values than PMMI. While

<sup>437</sup> Barnett and Onnela (2016)'s and Palmius et al. (2017)'s models estimate 12 measurements

<sup>438</sup> for each missing hour, PMMI estimates as many measurements as there are in the log. When

<sup>439</sup> individuals go to infrequently travelled locations they tend to use their phone's location

<sup>440</sup> services more than 100 more times, which leads to more measurements, and hence more

<sup>441</sup> imputed values. Hence, if the weights are not used PMMI is susceptible to error

<sup>442</sup> overestimates which are related to measurement frequency in a way that the other two

<sup>443</sup> methods are not.

444 In addition, other measures of interest for PMMI are accuracy (in what percentage of

445 the cases was the appropriate cluster predicted), the *confidence* and the *distance expectation*.

446 The confidence is the probability with which the model predicts the most likely cluster. For

447 instance, if the model predicts the missing bin to be bin A with a probability of 0.9 then we

448 can say it has high confidence in the prediction. Similarly, the distance expectation is the

449 cross-product of the estimated probabilities that the individual is at any of all given clusters

450 and the distances between the clusters to the true cluster. For example, if the true location

451 of an individual is bin A (bin A is 10 meters away from bin B) and the model assigns a

452 probability of 0.9 at bin A and 0.1 at bin B, then the distance expectation would be 1 meter.

453 Finally, we will take objective real world data and compare it to predicted values. We

454 will use information from the Dutch public transportation card. The Dutch public

455 transportation service provides users with time-stamped data of when and where they board,

456 change lines or leave public transportation. To be able to make a comparison between

457 models we will remove all measurements from within the 5 minute period that a

458 time-stamped measurement is available. Then we will use each model to impute the location

459 of the individual within that period.

460

## Results

### 461 Map building & binning evaluation

462 We used the following parameters to extract pauses: an accuracy limit  $a_{\text{Pause Lim}}$  of 250

463 meters, a time limit  $t_{\text{Pause Lim}}$  of 300 seconds, a distance limit of  $d_{\text{Pause Lim}}$  50 meters and a

464 minimum pause duration limit  $\delta_{\text{Pause Lim}}$  of 100 seconds. Moreover, to extract path clusters

465 we used the parameters:  $a_{\text{Path Lim}} = 150$ meters and  $d_{\text{Path Lim}} = 300$ meters.

466 The selection of these parameters was more-or-less heuristically driven based on their

467 on their ability to extract a meaningful personalised map. When selecting parameters there

468 is a trade off is between bias and precision. This is because an increase in precision in the

469 form of a higher for high density locations resolution comes at the expense of precision as

470 assigning measurements to bins becomes more difficult. For instance, by increasing the  
 471 clustering distance parameter  $d_{\text{Pause Lim}}$  we can extract more valid pause locations at the  
 472 expense of falsely categorising certain measurements to the wrong cluster. This is illustrated  
 473 in Figure 5.

474 Map building results in a personalised map with pause and path clusters. An excerpt  
 475 can be seen in Figure 6. It is important to remind the reader that PMMI is map agnostic and  
 476 uses no information from the map. Therefore, the close overlap with features on the map,  
 477 such as pause bins at relevant buildings and transportation clusters, as well as the flight bins  
 478 following roads and railway lines indicate a high degree of precision in personal map building.  
 479 As expected, PMMI's path extraction yields greater accuracy for frequently occurring paths.  
 480 For example, the frequently travelled Amsterdam-Utrecht railway line has been extracted  
 481 almost perfectly, while the less frequently travelled Utrecht-Enschede line is far sparser.

482 For the entire period examined period the we find a deviance of 40 meters and a  
 483 median deviance of 15 meters. Around 69% of the deviance values are within their  
 484 corresponding accuracy value, which is close to the theoretical  $\frac{2}{3}$  value that is expected.  
 485 Approximately 9% of values were not taken into account when creating the bins,given that  
 486 their the accuracy  $a$  exceeded  $a_{\text{Path Lim}}$ .

487 For the narrower period of March 2017 approximately 74% of the deviance values are  
 488 within their corresponding accuracy values, which is close to the theoretical  $\frac{2}{3}$  value that is  
 489 expected. The raw unweighted mean and median deviance are 38 and 14 meters respectively.  
 490 If we weight them such that each 5 minute interval is weighted equally to ease comparison  
 491 with the other two methods, we find a mean deviance of 35 meters with a median of 12.8.

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

497        The key difference between temporal and spatial down-sampling is shown in Figure ??.

498        Temporal down-sampling is much more sensitive to noise in sparsely measured periods  
499        because it averages out values within five minute periods. Often there are only a few noisy  
500        measurements in those periods (see Figure 1 ), which leads to a noise in the down-sampled  
501        values. Unsurprisingly, there is a positive relationship between deviance and the amount of  
502        measurements in each down-sampled interval. The fact that over 90% of deviance values are  
503        within accuracy (much higher than the expected theoretical 67%) confirms that temporal  
504        down-sampling is not sufficiently filtering out the noise.

505        **Imputation evaluation.** The Palmius et al. (2017) model failed to impute 3% of

506        all removed values for both the 5 minute and over 10% for the 1 hour tests. Palmius et al.  
507        (2017)' model failed to impute a single value for the day tests. Similarly, Barnett and Onnela  
508        (2016)'s method failed to impute over 11% of the missing values. PPMI made an imputation  
509        for all missing values. Table 1 shows the results of the distance metrics for each method.

510        In terms of PPMI's accuracy in this period, the prediction accuracy was 88%, 73% and  
511        47% for the 5 minute, 1 hour and 1 day periods respectively. As a comparison, the naive  
512        model's accuracy ratings were 87%, 68% and 24% respectively. The distance expectation  
513        values were very similar to the distance scores, albeit approximately 5% higher. Confidence  
514        scores ranged from 0.006 to 1.

515        **Comparison with objective data.** Once we removed the measurements within

516        the 5 minute period of each time-stamp, we employed the models to predict the location of  
517        the individual at the time-stamped period. The naive model reaches an accuracy of 16%,  
518        with a median distance of 555 meters and mean distance of 3303 meters. On the other hand  
519        MYMETHOD has an accuracy of 24%, with a median distance of 1037 meters and a mean  
520        distance of 5637. Again the expectation of the distance was quite similar to the distance  
521        with a mean of 6432 and a median of 1037 meters.

522        Palmius's model failed to impute 38 out of 97 periods. For those it did impute, it

523        showed a mean distance of 1517 meters and a median distance of 1617. Barnett's model

524 failed to impute 18 out of 97 periods. For those it did impute, it had a mean of 5506 and a  
525 median of 1342.

526       **Example: effect on aggregate measures.** Social scientists are most interested in  
527 aggregating spatiotemporal data to more socially relevant information, such as the amount of  
528 time spent at home. As an example we calculated the time spent at home of the user in the  
529 month of March (Figure 8) without any model and with all three of the investigated methods.

530       Interestingly, all three models suggest that the user spent approximately 60% of their  
531 time at home. However, without using any form of missing data imputation approximately  
532 12% of the time is unaccounted for. Given that the amount of time spent at home has been  
533 found to be a reliable predictor of extroversion (G. M. Harari et al., 2016) and the onset of  
534 depressive episodes in bipolar patients (Palmius et al., 2017) obtaining an accurate value for  
535 mobility metrics like this is highly important to social scientists.

536       

## Discussion & Conclusion

537       Overall the PMMI performed better than the alternative models, particularly during  
538 longer missing periods and with objective data. However, it did not perform substantially  
539 better than the naive baseline model. In addition, the comparison to the performance of the  
540 Barnett&Onella and Palmius models is somewhat unfair, as they were created for custom  
541 logs, not secondary logs. Nonetheless, the comparison remains valid as they are the closest  
542 we found to a missing data imputation methods in smartphone GPS logs.

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

547       PPMI's increased coverage and flexibility comes from its ability to make complex  
548 non-linear predictions. For instance, in a given missing period it might make sense to predict  
549 that the individual is either at home, or at the office, or at a shop with equal probability.

550 While PPMI can make such an imputation, none of the alternative methods can do this.  
551 Moreover, the ability to take the prediction probability values from the neural network also  
552 helps in dealing with uncertainty. A known-drawback of single imputation is that it takes an  
553 imputed value and treats it as observed. Simple rule-based methods such as Palmius' are  
554 essentially algorithmic single imputation methods. With PPMI it is possible to model  
555 uncertainty using the predicted probabilities of each estimate. For instance, in the previous  
556 example, we could choose to only take estimates with a high degree of confidence, thus  
557 creating confidence intervals by adding and subtracting the amount of cases where the  
558 location of the individual is ambiguous.

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

568 While both Barnett and Onnela (2016) and Palmius et al. (2017) use historical data to  
569 smoothen pause locations by clustering pause locations with a close degree of spatial  
570 proximity, neither of them do the same for non-pause locations. This may be feasible with  
571 high frequency, regularly sampled short duration logs but creates noise with secondary logs.  
572 Moreover, with secondary logs it is feasible to spatially “average out” multiple samples of the  
573 same path in order to recreate it in its entirety. For instance, although the mean sampling  
574 frequency during train travels on the Amsterdam-Utrecht line is low (about 0.01 Hz) the  
575 personalised map manages to recreate the train line almost perfectly, despite being  
576 completely map agnostic.

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

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

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

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

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

604 state-space models like the Kalman filter the state equation would represent a probabilistic  
605 representation of where the individual could be based on the individuals entire movement  
606 history. The space side of the model would be a measurement equation representing the  
607 measurement and the uncertainty surrounding it in the form of  $a$ .

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

612 That there is room for improvement is unsurprising given that it is still early days for  
613 using smartphone location measurements in social science. Nonetheless, the methodological  
614 advantages are clear: millions of individuals have a location log containing objective  
615 measurements spanning years which can be accessed for free. This is vastly superior to  
616 alternatives such questionnaires which rely on accurate self-reporting. Social science  
617 researchers must take advantage of regulatory changes in data portability regulations and  
618 put the vast wealth of data collected by commercial entities to public use.

619

## References

```
r_refs(file = "r-references.bib")
```

- 620 Arnold, J. B. (2017). *Ggthemes: Extra themes, scales and geoms for 'ggplot2'*. Retrieved  
 621 from <https://CRAN.R-project.org/package=ggthemes>
- 622 Aust, F., & Barth, M. (2018). *papaja: Create APA manuscripts with R Markdown*.  
 623 Retrieved from <https://github.com/crsh/papaja>
- 624 Barnett, I., & Onnela, J.-P. (2016). Inferring mobility measures from GPS traces with  
 625 missing data. *arXiv:1606.06328 [Stat]*. Retrieved from  
 626 <http://arxiv.org/abs/1606.06328>
- 627 Chen, M. Y., Sohn, T., Chmelev, D., Haehnel, D., Hightower, J., Hughes, J., ... Varshavsky,  
 628 A. (2006). Practical metropolitan-scale positioning for GSM phones. In *UbiComp  
 629 2006: Ubiquitous computing* (pp. 225–242). Springer, Berlin, Heidelberg.  
 630 doi:[10.1007/11853565\\_14](https://doi.org/10.1007/11853565_14)
- 631 Chen, Z., & Brown, E. N. (2013). State space model. *Scholarpedia*, 8(3), 30868.  
 632 doi:[10.4249/scholarpedia.30868](https://doi.org/10.4249/scholarpedia.30868)
- 633 Commission, E. (2017). *Protecting your data: Your rights - european commission*. Retrieved  
 634 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)
- 635 Delclòs-Alió, X., Marquet, O., & Miralles-Guasch, C. (2017). Keeping track of time: A  
 636 smartphone-based analysis of travel time perception in a suburban environment.  
 637 *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)
- 638 Duncan, S., Stewart, T. I., Oliver, M., Mavoa, S., MacRae, D., Badland, H. M., & Duncan,  
 639 M. J. (2013). Portable global positioning system receivers: Static validity and  
 640 environmental conditions. *American Journal of Preventive Medicine*, 44(2), e19–29.  
 641 doi:[10.1016/j.amepre.2012.10.013](https://doi.org/10.1016/j.amepre.2012.10.013)
- 642 Feng, L., Nowak, G., O'Neill, T., & Welsh, A. (2014). CUTOFF: A spatio-temporal  
 643 imputation method. *Journal of Hydrology*, 519, 3591–3605.

- 644 doi:[10.1016/j.jhydrol.2014.11.012](https://doi.org/10.1016/j.jhydrol.2014.11.012)
- 645 Goodchild, M. F., & Janelle, D. G. (2010). Toward critical spatial thinking in the social  
646 sciences and humanities. *GeoJournal*, 75(1), 3–13. doi:[10.1007/s10708-010-9340-3](https://doi.org/10.1007/s10708-010-9340-3)
- 647 Grünerbl, A., Muaremi, A., Osmani, V., Bahle, G., Ohler, S., Tröster, G., . . . Lukowicz, P.  
648 (2015). Smartphone-based recognition of states and state changes in bipolar disorder  
649 patients. *IEEE Journal of Biomedical and Health Informatics*, 19(1), 140–148.  
650 doi:[10.1109/JBHI.2014.2343154](https://doi.org/10.1109/JBHI.2014.2343154)
- 651 Harari, G. M., Lane, N. D., Wang, R., Crosier, B. S., Campbell, A. T., & Gosling, S. D.  
652 (2016). Using smartphones to collect behavioral data in psychological science:  
653 Opportunities, practical considerations, and challenges. *Perspectives on Psychological  
654 Science*, 11(6), 838–854. doi:[10.1177/1745691616650285](https://doi.org/10.1177/1745691616650285)
- 655 Jankowska, M. M., Schipperijn, J., & Kerr, J. (2015). A framework for using GPS data in  
656 physical activity and sedentary behavior studies. *Exercise and Sport Sciences  
657 Reviews*, 43(1), 48–56. doi:[10.1249/JES.0000000000000035](https://doi.org/10.1249/JES.0000000000000035)
- 658 LaMarca, A., Chawathe, Y., Consolvo, S., Hightower, J., Smith, I., Scott, J., . . . Schilit, B.  
659 (2005). Place lab: Device positioning using radio beacons in the wild. In *Pervasive  
660 computing* (pp. 116–133). Springer, Berlin, Heidelberg. doi:[10.1007/11428572\\_8](https://doi.org/10.1007/11428572_8)
- 661 Location History, G. (2017). *Timeline*. Retrieved from  
662 <https://www.google.com/maps/timeline?pb>
- 663 London, I. (2016). *Encoding cyclical continuous features - 24-hour time*. *Ian london's blog*.  
664 Retrieved April 27, 2018, from  
665 [//ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/](https://ianlondon.github.io/blog/encoding-cyclical-features-24hour-time/)
- 666 Palmius, N., Tsanas, A., Saunders, K. E. A., Bilderbeck, A. C., Geddes, J. R., Goodwin, G.  
667 M., & Vos, M. D. (2017). Detecting bipolar depression from geographic location data.  
668 *IEEE Transactions on Biomedical Engineering*, 64(8), 1761–1771.  
669 doi:[10.1109/TBME.2016.2611862](https://doi.org/10.1109/TBME.2016.2611862)
- 670 Patterson, T. A., Thomas, L., Wilcox, C., Ovaskainen, O., & Matthiopoulos, J. (2008).

671 State-space models of individual animal movement. *Trends in Ecology & Evolution*,  
672 23(2), 87–94. doi:[10.1016/j.tree.2007.10.009](https://doi.org/10.1016/j.tree.2007.10.009)

673 Preisler, H. K., Ager, A. A., Johnson, B. K., & Kie, J. G. (2004). Modeling animal  
674 movements using stochastic differential equations. *Environmetrics* 15: P. 643-657.  
675 Retrieved from <https://www.fs.usda.gov/treesearch/pubs/33038>

676 R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna,  
677 Austria: R Foundation for Statistical Computing. Retrieved from  
678 <https://www.R-project.org/>

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

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

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

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

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

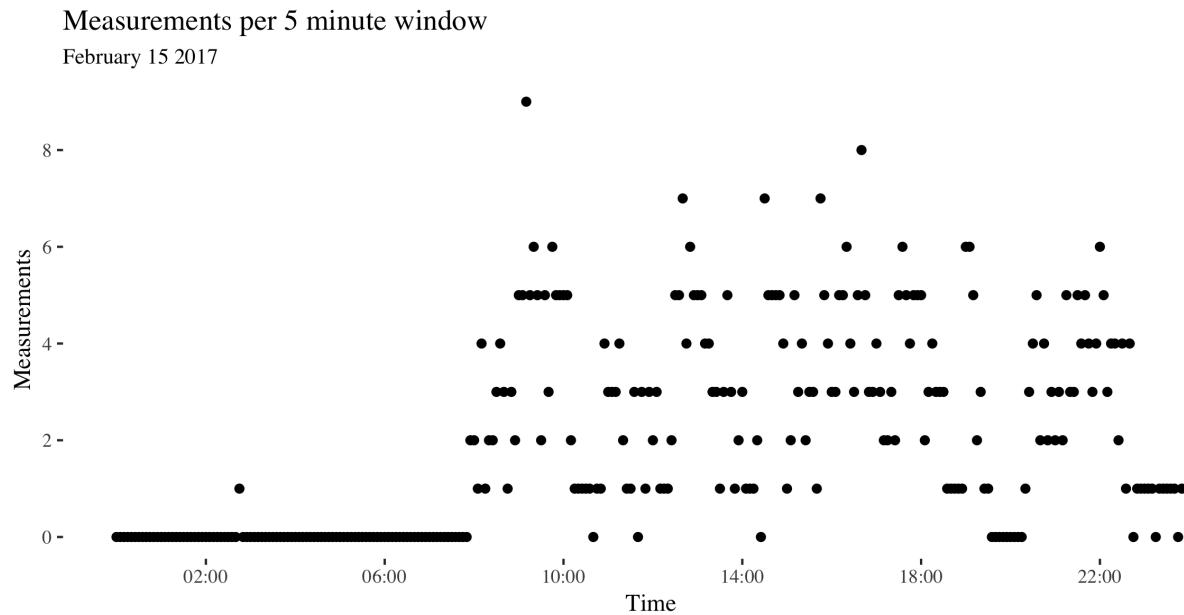
695 Villanueva, C., & Aggarwal, B. (2013). The association between neighborhood  
696 socioeconomic status and clinical outcomes among patients 1 year after  
697 hospitalization for cardiovascular disease. *Journal of Community Health*, 38(4),

- 698 690–697. doi:[10.1007/s10900-013-9666-0](https://doi.org/10.1007/s10900-013-9666-0)
- 699 700 701 702 703 Wang, R., Harari, G., Hao, P., Zhou, X., & Campbell, A. T. (2015). SmartGPA: How smartphones can assess and predict academic performance of college students. In *Proceedings of the 2015 ACM international joint conference on pervasive and ubiquitous computing* (pp. 295–306). New York, NY, USA: ACM.  
doi:[10.1145/2750858.2804251](https://doi.org/10.1145/2750858.2804251)
- 704 705 Wickham, H. (2009). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.  
Retrieved from <http://ggplot2.org>
- 706 707 708 709 Wolf, J., Oliveira, M., & Thompson, M. (2003). Impact of underreporting on mileage and travel time estimates: Results from global positioning system-enhanced household travel survey. *Transportation Research Record: Journal of the Transportation Research Board*, 1854, 189–198. doi:[10.3141/1854-21](https://doi.org/10.3141/1854-21)
- 710 711 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida:  
Chapman; Hall/CRC. Retrieved from <https://yihui.name/knitr/>
- 712 713 Zenk, S. N., Schulz, A. J., & Odoms-Young, A. (2009). How neighborhood environments contribute to obesity. *The American Journal of Nursing*, 109(7), 61–64.  
doi:[10.1097/01.NAJ.0000357175.86507.c8](https://doi.org/10.1097/01.NAJ.0000357175.86507.c8)
- 715 716 717 718 Zhang, Z., Yang, X., Li, H., Li, W., Yan, H., & Shi, F. (2017). Application of a novel hybrid method for spatiotemporal data imputation: A case study of the minqin county groundwater level. *Journal of Hydrology*, 553(Supplement C), 384–397.  
doi:[10.1016/j.jhydrol.2017.07.053](https://doi.org/10.1016/j.jhydrol.2017.07.053)

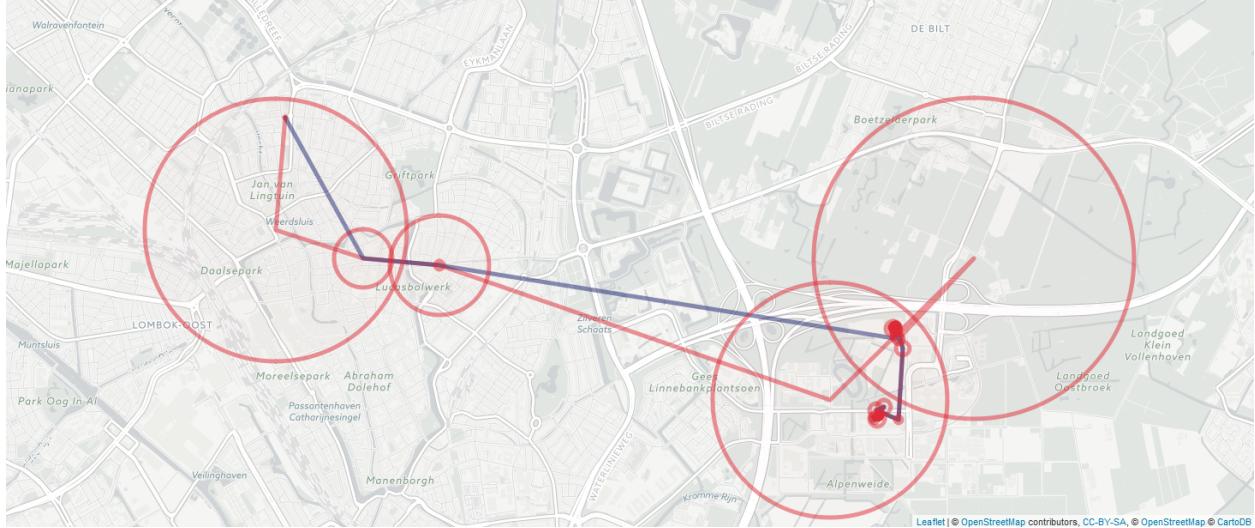
Table 1

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

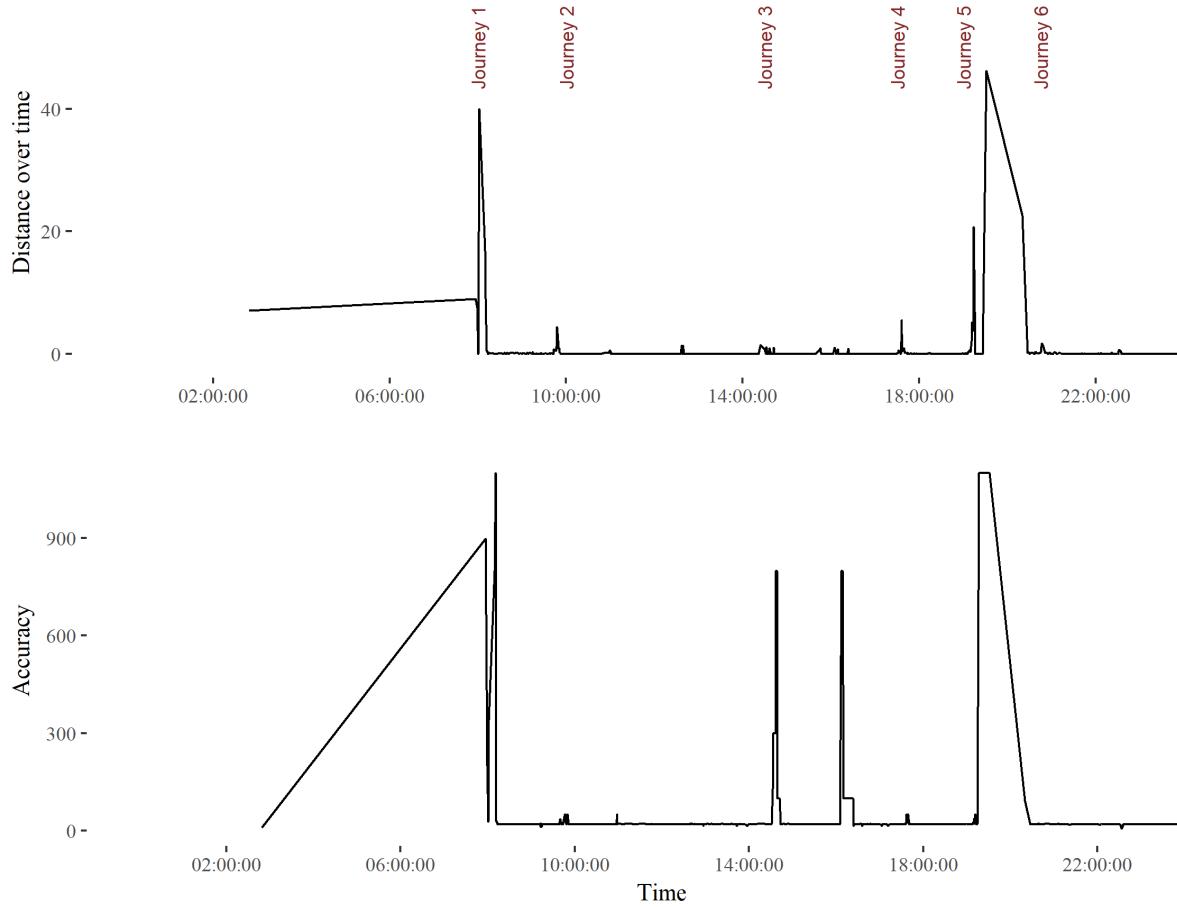
	Five minutes		One Hour		One Day	
	Mean	Median	Mean	Median	Mean	Median
Barnett & Onella	82	4	345	6	9,273	12
Palmius	43	0	497	4	NA	NA
PPMI	269	0	908	0	5,757	0
Naive Baseline	426	0	1,502	0	14,266	1,288



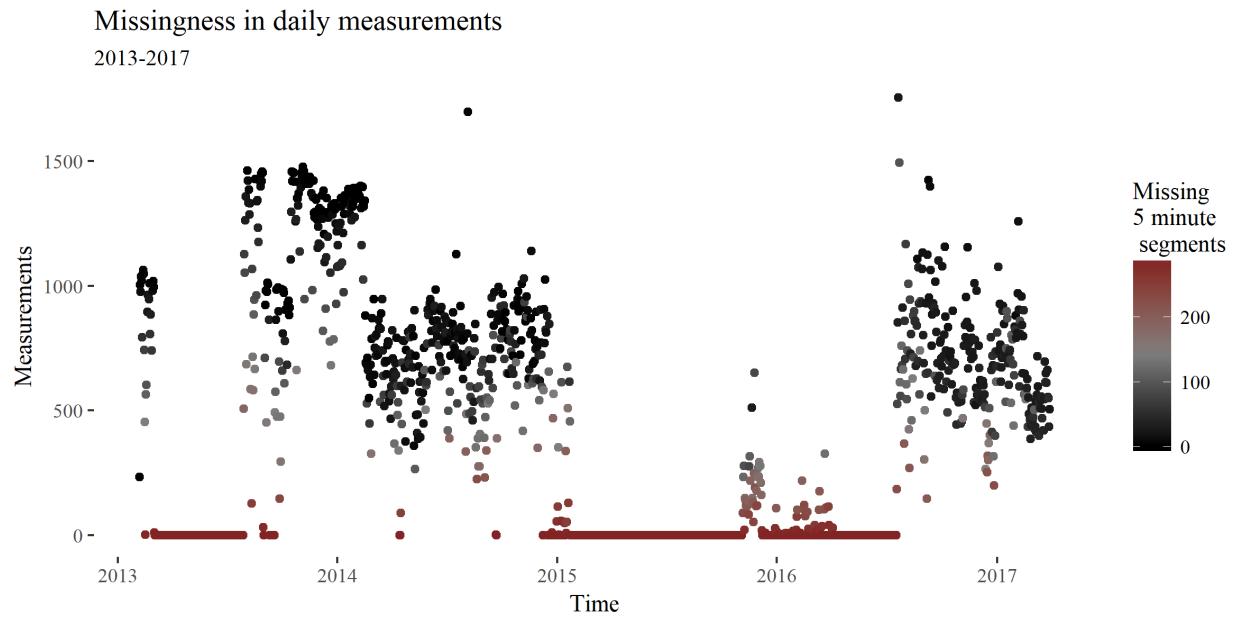
*Figure 1.* Example of missing data over the entire duration of a secondary log. The x-axis denotes time, the y-axis shows how many measurements are made and each point is a five minute window. For this day there were several periods with no information. These points lie on the x-axis.



*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 1. 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. For this day there were several periods with no information. These points are filled with red and lie on the x-axis.

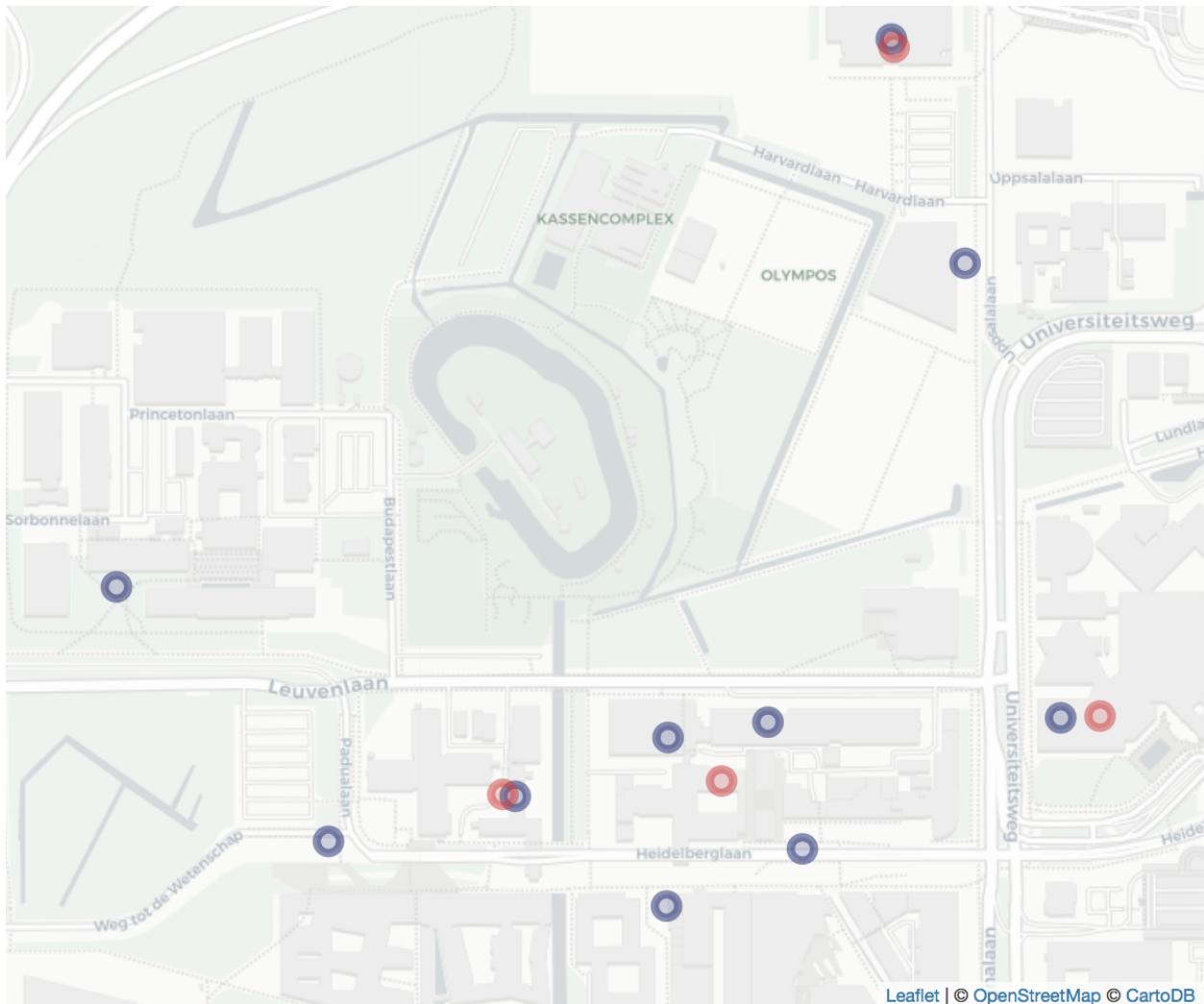


Figure 5. Example of pause locations in De Uithof university campus using 150 meters (blue) and 400 meters (red) as clustering parameters.

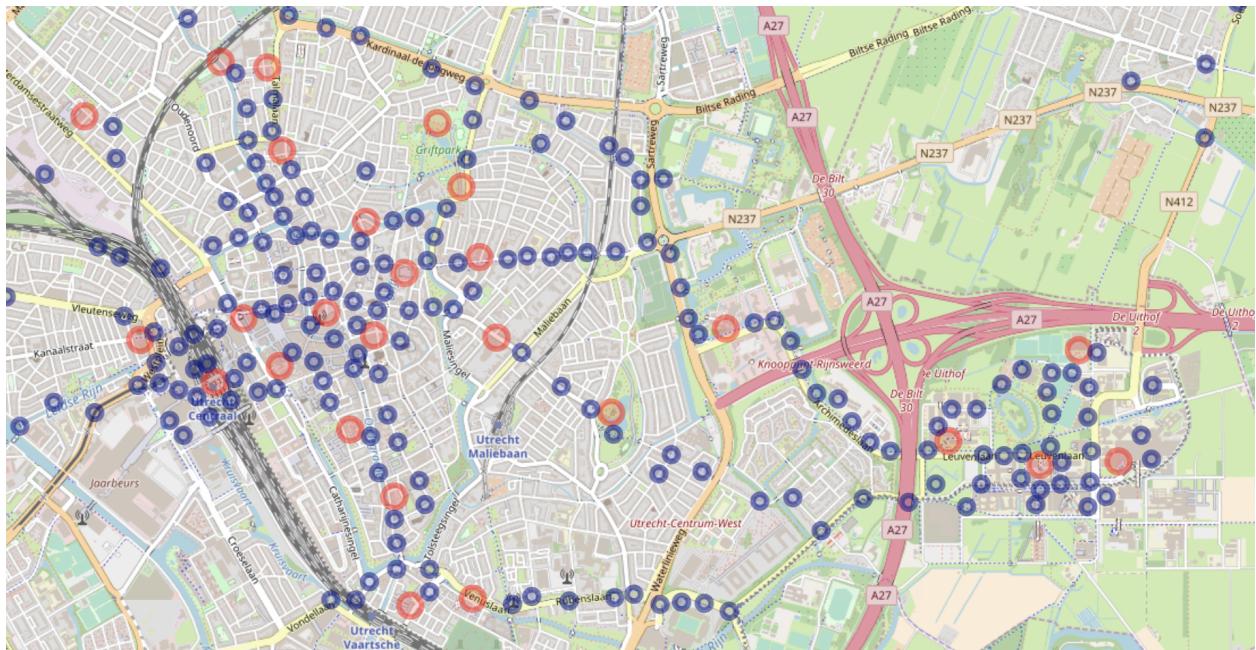


Figure 6. Excerpt of the cluster map of an individual. Red points are pause locations, blue points are path locations.

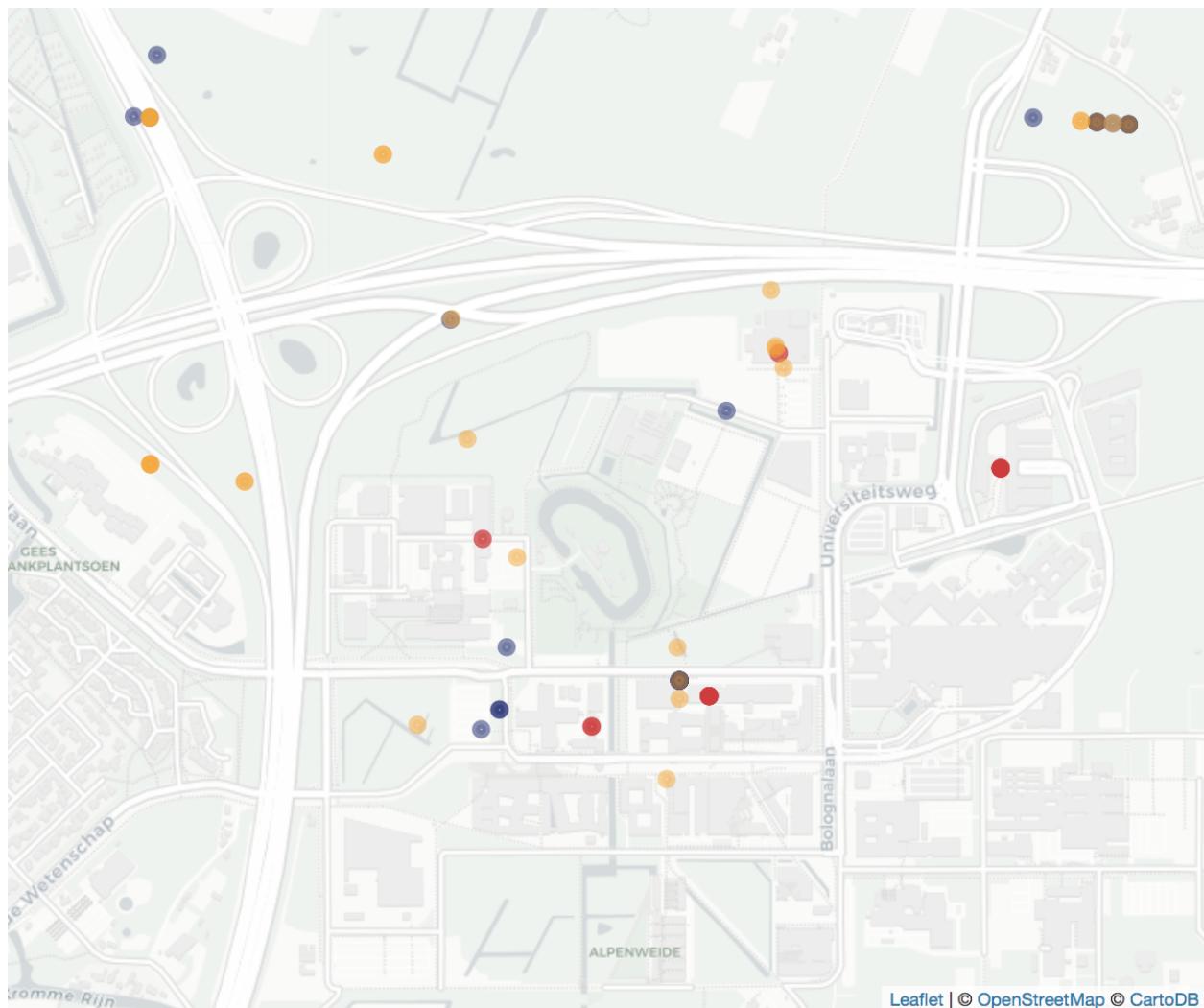
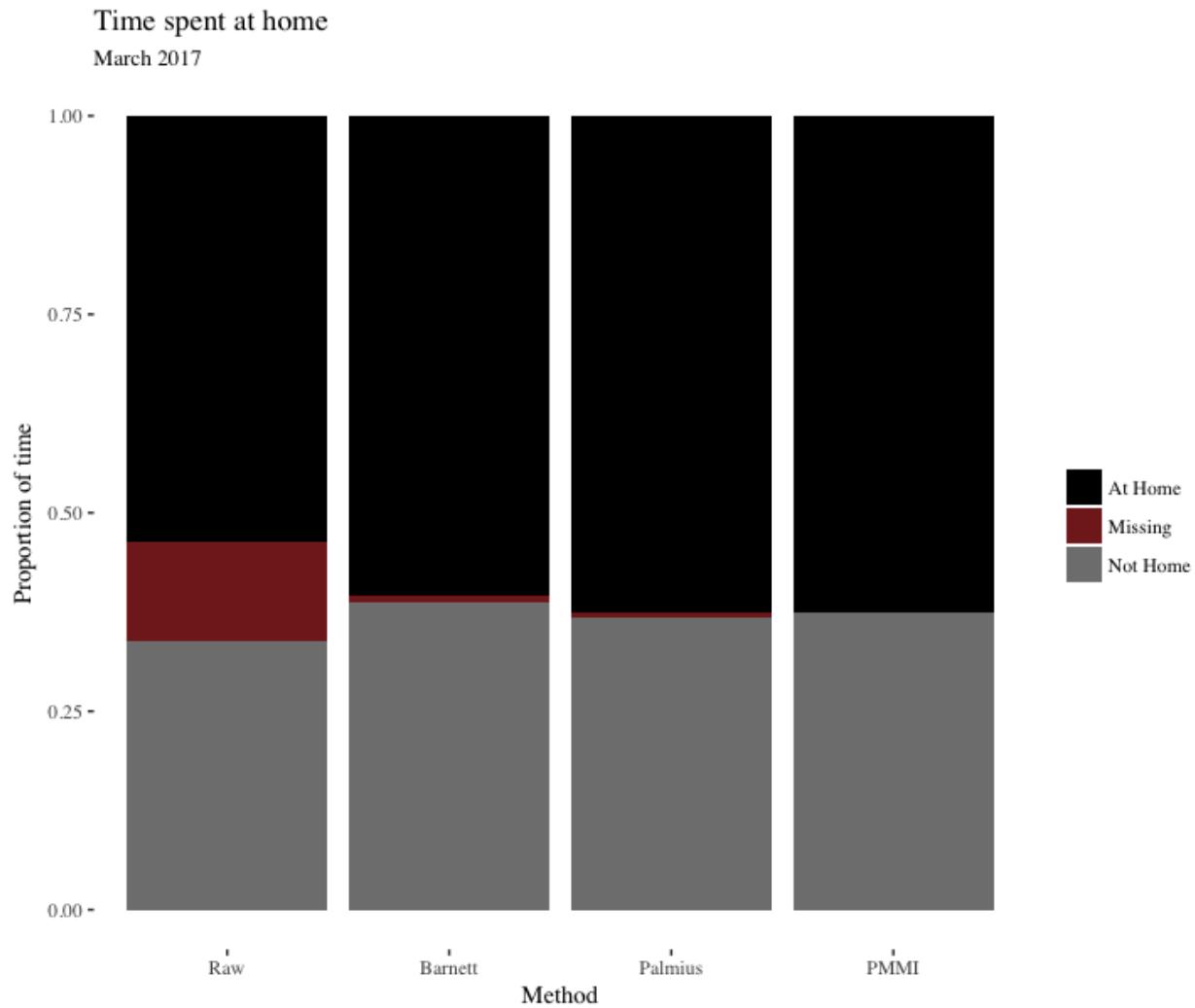


Figure 7. Excerpt. MISSING CAPTION HERE



*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.