

MIT Open Access Articles

*A new insight into land use classification
based on aggregated mobile phone data*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Pei, Tao, Stanislav Sobolevsky, Carlo Ratti, Shih-Lung Shaw, Ting Li, and Chenghu Zhou. "A New Insight into Land Use Classification Based on Aggregated Mobile Phone Data." *International Journal of Geographical Information Science* 28, no. 9 (May 8, 2014): 1988–2007.

As Published: <http://dx.doi.org/10.1080/13658816.2014.913794>

Publisher: Taylor & Francis

Persistent URL: <http://hdl.handle.net/1721.1/101646>

Version: Original manuscript: author's manuscript prior to formal peer review

Terms of use: Creative Commons Attribution-Noncommercial-Share Alike





senseable city lab:...

1 A New Insight into Land Use Classification Based on Aggregated Mobile Phone Data

2

3

4 TAO PEI

5 *State Key Laboratory of Resources and Environmental Information System, Institute*
6 *of Geographic Sciences and Natural Resources Research, CAS, 11A, Datun Road*
7 *Anwai, Beijing 100101, China*

8

9 *SENSEable City Lab, School of Architecture and Planning, Massachusetts Institute of*
10 *Technology, Cambridge, MA 02139, USA*

11

12

13 Stanislav Sobolevsky (Corresponding Author)

14 *Senseable City Lab, School Of Architecture And Planning, Massachusetts Institute Of*
15 *Technology, Cambridge, Ma 02139, Usa*

16

17

18 Carlo Ratti

19 *Senseable City Lab, School Of Architecture And Planning, Massachusetts Institute Of*
20 *Technology, Cambridge, Ma 02139, Usa*

21

22

23 Shih-Lung Shaw

24 *Department Of Geography, University Of Tennessee, 304 Burchfiel Geography*
25 *Building, Knoxville, Tn 37996-0925, Usa*

26 *State Key Laboratory Of Information Engineering In Surveying, Mapping And Remote*
27 *Sensing, Wuhan University, Wuhan, Hubei 430079, China*

28

29

30 Chenghu Zhou

31 *State Key Laboratory of Resources and Environmental Information System, Institute*
32 *of Geographic Sciences and Natural Resources Research, CAS, 11A, Datun Road,*
33 *Beijing 100101, China*

34

35

36

37

38

39

40

41

42

43 Abstract

44 Land use classification is essential for urban planning. Urban land use types can be
45 differentiated either by their physical characteristics (such as reflectivity and texture)
46 or social functions. Remote sensing techniques have been recognized as a vital
47 method for urban land use classification because of their ability to capture the
48 physical characteristics of land use. Although significant progress has been achieved
49 in remote sensing methods designed for urban land use classification, most techniques
50 focus on physical characteristics, whereas knowledge of social functions is not
51 adequately used. Owing to the wide usage of mobile phones, the activities of residents,
52 which can be retrieved from the mobile phone data, can be determined in order to
53 indicate the social function of land use. This could bring about the opportunity to
54 derive land use information from mobile phone data. To verify the application of this
55 new data source to urban land use classification, we first construct a time series of
56 aggregated mobile phone data to characterize land use types. This time series is
57 composed of two aspects: the hourly relative pattern, and the total call volume. A
58 semi-supervised fuzzy c-means clustering approach is then applied to infer the land
59 use types. The method is validated using mobile phone data collected in Singapore.
60 Land use is determined with a detection rate of 58.03%. An analysis of the land use
61 classification results shows that the accuracy decreases as the heterogeneity of land
62 use increases, and increases as the density of cell phone towers increases.

63 Keywords: land use; mobile phone data; classification; FCM; Singapore

64

65 1. Introduction

66 The classification of urban land use is essential for urban planning. Urban land
67 use, defined as the recognized human use of land in a city, can be differentiated either
68 by its physical characteristics (such as reflectivity and texture) or social functions (i.e.,
69 residential areas are for living whereas industrial areas are for working). Among urban
70 land use classification methods, remote sensing techniques are recognized as a vital
71 method because of their ability to capture the physical characteristics of land use.
72 Conventional land-use remote sensing methods classify land use based on spectral and
73 textural characteristics (Gong and Howarth 1990; Fisher 1997; Shaban and Dikshit
74 2001; Lu and Weng 2006). Nevertheless, because land use classes are heterogeneous
75 in both their spectral and textural characteristics, methods that rely on remote sensing
76 information and their derived characteristics are unable to differentiate between some
77 land use types (i.e., residential and commercial). Because of this, more auxiliary
78 information, such as contextual properties, field sizes and shapes, parcel information,
79 and expert knowledge, has been used to infer land use patterns (De Wit and Clevers,
80 2004; Platt and Rapoza, 2008; Wu et al. 2009; Hu and Wang, 2013). However, this
81 need for additional information not only increases the cost, but also delays the update
82 process. Although significant progress has been made in remote sensing techniques,
83 there is a tendency to focus on the utilization of information concerning physical
84 characteristics of land use, and knowledge of social functions is not adequately used
85 in the classification process.

86 Owing to the wide usage of mobile phones, the daily activities of residents in
87 various regions can be easily captured and used to indicate the social function of the
88 land use type. In other words, within different land use areas, people may demonstrate
89 different routine activities (for example, in residential areas, people usually leave

90 home for work in the morning and return in the evening, whereas in business areas the
91 opposite pattern can be found). This may allow us to derive the activities of residents,
92 and then the social functions of different land use types, from mobile phone data. As a
93 result, mobile phone data may provide a new insight into traditional urban land use
94 from the perspective of social function. The objective of this paper is to verify the
95 applicability of the potential data source for urban land use classification, and then
96 evaluate the results given by this new source of information.

97 The remainder of the paper is structured as follows. Section 2 introduces a newly
98 constructed time series, as well as the semi-supervised cluster method for urban land
99 use classification. In Section 3, the mobile phone data used in this paper are described.
100 Section 4 presents the overall procedure and the results of land use classification.
101 Section 5 validates the classification result by comparing it with that given by either
102 the call pattern or call volume alone. Section 6 discusses the factors affecting the
103 uncertainty in the classification, and Section 7 presents our conclusions and
104 suggestions for future work relating to land use classification based on mobile phone
105 data.

106

107 2. Related work

108 The retrieval of land use from mobile phone data can be divided into two stages.
109 The first is to retrieve the residents' activities based on mobile phone data. The
110 second is to infer land use from the residents' activities. Regarding the first stage,
111 recent research can be grouped into two categories. The first aims to reveal
112 individual mobility patterns using call detail record data, which consist of the
113 different base transceiver station (BTS) locations from which users have made calls
114 (Gonzalez, et al., 2008; Song et al. 2010; Calabrese et al., 2011). The second is based

115 on the aggregation of the total calling time (or numbers) at each BTS in a certain
116 temporal interval. Since our paper only uses the relationship between the mobility
117 and the aggregated mobile phone data in the inference of urban land use, the
118 literature review below will focus on the achievements of aggregated mobile phone
119 data.

120 The spatiotemporal variation regarding BTS has been extensively studied to
121 retrieve various residents' activities. Recent approaches include the description of
122 urban landscapes (i.e., the space-time structure of residents' activities in a city) (Ratti
123 et al. 2006; Pulselli et al., 2006; Sevtsuk and Ratti, 2010; Sun et al. 2011;
124 Jacobs-Crisioni and Koomen, 2012; Loibl and Peters-Anders, 2012), population
125 estimates (Vieira et al. 2010; Manfredini et al., 2011; Rubioa et al., 2013), the
126 identification of specific social groups (Vaccari et al. 2009), and the detection of
127 social events (Traag et al. 2011; Laura et al. 2012).

128 The inference of land use types in this context is dependent on their social
129 functions which can be derived from the residents' activities (namely, the overall
130 characteristics of human communication in the urban area). This contains two aspects:
131 the relative weekly calling pattern ("pattern" hereafter) and the total calling volume
132 ("volume" hereafter). The pattern is defined as the share of hourly calling volume in a
133 certain period. The calling volume of a BTS is defined as the total time (or number) of
134 calls managed by that BTS in its area of coverage over a given period of time. Unlike
135 the static residential population density, the volume is the overall characteristic of how
136 many people actually use mobile phones, indicating the activeness of their
137 communicational interactions. To identify and extract recurring patterns of mobile
138 phone usage and relate them to some land use types, Reades et al. (2009) proposed the
139 eigen-decomposition method, a process similar to factoring but suitable for complex

140 datasets. Calabrese et al. (2010) used an eigen-decomposition analysis to reveal the
141 relationship between mobile phone data and the residential and business areas.
142 Caceres et al. (2012) used a new tessellation technique to differentiate parks from
143 residential areas by detecting changes in human density retrieved from mobile phone
144 data.

145 Although these studies have addressed the relationship between land use and
146 mobile phone data, they have only focused on the identification of specific land use
147 types, not the classification of urban land use. In order to enhance the land use
148 classification, Soto and Frias-Martinez (2011a and 2011b) used the normalized time
149 series of the volume for a weekday and a weekend day (a time series consists of 48
150 points, each of which is the volume calculated at each hour and normalized by the
151 total volume of the 2 days) to identify the land use pattern. The same method was
152 applied to Twitter data by Frias-Martinez et al. (2012). Andrienko et al. (2013) used
153 the normalized timelines of mobile phone calls at each BTS to identify the
154 heterogeneity of the Ivory Coast at the country scale. Because the normalized data
155 only cover the temporal variation of the volume within the same BTS, the difference
156 in the total volume between BTSs was neglected. Therefore, regarding the problem of
157 heterogeneous land use (for example, downtown areas may have a variety of
158 commercial, residential, and recreational activities), methods based solely on
159 normalized patterns might fail to discern between different land use types that are not
160 homogenous.

161 To adapt the mobile phone data to urban land use classification, Toole et al.
162 (2012) proposed a supervised classification method for the data that combined the
163 normalized calling pattern and the volume (namely, “activity” in their paper). The
164 aggregated data were first converted to the residual of the Z-score normalization,

165 which reveals the flow into and out of the city center over the course of a day. The
166 random forest method, proposed by Breiman (2001), was then employed to determine
167 land use types. Although this method significantly enhanced the land use
168 classification, two aspects still need to be improved. First, the random forest, similar
169 to the neural network method, is a black box model (Berthold, 2010), which makes
170 the classification difficult to interpret. Second, only two-day pattern (an average
171 weekday and an average weekend) was used to infer the urban land use (Toole et al.,
172 2012). The difference between weekdays and that between weekends are neglected,
173 despite the fact that the significant differences exist between weekdays and between
174 weekends in terms of activities of residents (Jia and Jiang, 2012; Liu et al., 2012; Soto
175 and Frias-Martinez, 2011a).

176 Although previous studies have made substantial progresses, we think two key
177 problems should be further studied to evaluate the capability of this new data source
178 to infer urban land use. First, the time series model that represents land use type at the
179 BTS level should be improved to enhance urban land use classification. On the one
180 hand, the model should be more sophisticated and incorporate more characteristics
181 (say, the differences between weekdays and between weekends, new indices derived
182 from aggregated mobile phone data) in order to better differentiate between different
183 land use types. This is because the land use is not only dynamically changing, but is
184 often also heterogeneous in some areas. Thus, either the pattern or the volume may
185 not fully interpret the social functions of different land use types. On the other hand,
186 the model should be more transparent to allow an evaluation of the effects of different
187 characteristics on land use classification. This may help us analyze and improve the
188 classification method. Second, because mobile phone data is a new data source in
189 terms of urban planning, it is important to evaluate the uncertainties and influential

190 factors behind land use classification. These include three aspects. One is related to
191 the model, and specifically the different characteristics in the time series. The second
192 concerns the data, particularly the BTS density. The third considers the ground truth,
193 and specifically the heterogeneity of land use.

194 To overcome these key problems, we construct a new time series by generating a
195 linear combination of the four-day call pattern and volume. This time series not only
196 utilizes more characteristics of mobile phone data, but also makes the classification
197 result easier to interpret. A new semi-supervised scheme is proposed to infer the land
198 use based on this time series. Using this process, we can classify the urban land use
199 and understand the different effects imposed by the call pattern and volume on the
200 classification result. Finally, the uncertainties of land use classification are analyzed in
201 terms of the dissimilarity between land use definition and classification result, mixture
202 of land use, BTS density, and the fuzzy membership value generated by the proposed
203 method.

204

205 3. Semi-supervised fuzzy c-means (FCM) clustering method for urban land use 206 classification

207 We first construct a synthesized time series, which is the linear combination of
208 the normalized pattern and the total calling volume. The pattern part can be
209 determined by the characteristics of the mobile phone data that will be used. Then, to
210 determine different types of land use types with the synthetic time series, we use a
211 semi-supervised clustering FCM method. Thus, the effect of different parts of the time
212 series on the classification can be determined by calculating the ratios in the distance
213 between cluster centers and the time series.

214 The process of classification is divided into the following five steps. 1) Place the

215 aggregated mobile phone data from each BTS into a mesh. 2) Construct the
216 synthesized time series that combines the normalized pattern with the calling volume.
217 A coefficient (β) is introduced to weight the pattern versus the volume. 3) Determine
218 β by training samples of different land uses, which are selected based on expert
219 knowledge. 4) Cluster the time series of mobile phone data using FCM. 5)
220 Post-process the clustering result by assigning each cluster to different land use types.
221 Each of these steps is now described in detail.

222

223 3.1. Gridding the data

224 Before being used to identify urban land use, the mobile phone data, aggregated
225 hourly at the BTS level, are interpolated to generate a mesh grid for further
226 computation. The data generated by each cell on an hourly basis form a time series.
227 The procedure is divided into four stages. First, a Voronoi polygon system is
228 generated using the BTS tower locations. Next, the volume in each BTS polygon is
229 divided by its area to give the volume density. The inverse distance weighting (IDW)
230 method is then used to generate the grid at hourly intervals. Finally, the hourly values
231 generated over each BTS form a time series.

232

233 3.2. Constructing the time series of aggregated mobile data

234 The time series we use in our method consists of two parts. The first is the hourly
235 pattern of mobile phone data. The second is the total volume, given by:

$$236 \quad Z_i = [X_i \beta \cdot Y_i] \quad (1)$$

237 , where Z_i ($\{z_{i,j}, i = 1, 2, \dots, n; j = 1, 2, \dots, T\}$) is the combined time series for cell i ,

238 X_i ($\{x_{i,j}, i = 1, 2, \dots, n; j = 1, 2, \dots, T\}$) is the pattern for cell i (see equation (2)), n is the

239 number of cells in the grid, T is the number of hours considered in the pattern, and Y_i
 240 is the volume for cell i modified by the range transformation (equation (3)).

$$241 \quad X_{i,j} = \frac{b_{i,j}}{\sum_{j=1}^T b_{i,j}} \quad (i = 1,2, \dots, n; j = 1,2, \dots, T) \quad (2)$$

$$242 \quad Y_i = \frac{2[\sum_{j=1}^T b_{i,j} - \min(\sum_{j=1}^T b_{i,j})]}{\max(\sum_{j=1}^T b_{i,j}) - \min(\sum_{j=1}^T b_{i,j})} \quad (i = 1,2, \dots, n) \quad (3)$$

243 , where $b_{i,j}$ is the original hourly calling volume at cell i . Note that we multiply the
 244 numerator by 2 to ensure that Y_i has the same range as X_i . The reason we use
 245 range transform is for a comparison of the roles played by the pattern and the volume
 246 in the classification.

247 3.3. Determination of β

248 To estimate the coefficient β , we select L ($L = \sum_1^K l_k$) samples from K land
 249 use types (l_k is the number of samples for land use type k). These land use types
 250 should already be known from other information sources, e.g., points of interest (POI)
 251 in Google Earth. The center for each land use sample group
 252 ($C_k(\{c_{k,j}, k = 1,2, \dots, K; j = 1,2, \dots, T\})$) can be determined by averaging the
 253 sample time series:

$$254 \quad c_{k,j} = \frac{1}{l_k} \sum_{i=1}^{l_k} z_{i,j}^{(k)} \quad (j = 1,2, \dots, T) \quad (4)$$

255 If we define $d_{i,j}$ as the distance between sample i and cluster center j , then the
 256 land use type for sample i can be determined by locating the minimum distance
 257 between it and each cluster center.

$$258 \quad ID_i' = \text{find}(d_{i,j} == \min(d_{i,j})) \quad (i = 1,2, \dots, K; j = 1,2, \dots, T) \quad (5)$$

259 ID'_i is the land use type of sample i . We define ID_i as the true land use type of
 260 sample i for the validation. Then the value of β can be determined by minimizing
 261 the objective function:

$$262 \quad f(\beta) = \sum_i I(Z_i) \quad (i = 1, 2, \dots, L) \quad (6)$$

263 , where $I(Z_i) = \begin{cases} 0 & ID'_i = ID_i \\ 1 & ID'_i \neq ID_i \end{cases}$ is an indicator function with $I(\cdot) = 0$ when Z_i is
 264 correctly classified; otherwise, $I(\cdot) = 1$. The objective function is calculated for
 265 different values of β . The optimized value of β is that at which $f(\beta)$ reaches its
 266 minimum.

267

268 3.4. Determination of final land use type

269 After determining the value of β , the time series for all cells are clustered using
 270 FCM. There are two strategies to choose the number of clusters in FCM (Bezdek,
 271 1981; Nock and Nielsen, 2006). The first is to simply set the number of clusters to the
 272 number of land use types. The second determines the number of clusters from the
 273 validation index generated on each execution of FCM (Ray and Turi, 1999). In this
 274 study, we choose the second strategy, because certain land use types are the result of a
 275 simplified urban planning map, and may thus be a combination of different specific
 276 land use types. For example, an Open space may contain areas of Park, Green,
 277 Cemetery, and Water. In this context, we would rather retain the natural structure of
 278 clusters (which might be some specific land use types) for the post-process
 279 combination than generate a predefined number of clusters, which may cause some
 280 land use type is divided into different clusters.

281

282 3.5. Post-processing to assign clusters to specific land use types

283 Once the clusters have been generated, we perform post-processing to assign
284 each cluster to an appropriate land use type. A cluster is assigned to the specific land
285 use type whose center, as represented by the samples used in section 3.3, is closest to
286 the center of the cluster. If the number of clusters is greater than the number of land
287 use types, at least one land use type will be assigned more than one cluster. If there are
288 fewer clusters than land use types, then we use the number of land use types to
289 re-cluster the data.

290

291 4. Aggregated mobile phone data from Singapore

292 The mobile phone data used for the land use classification are the hourly
293 aggregated number of calls managed by each of 5500+ BTS towers in Singapore. To
294 determine land use types from mobile phone data, we use data from a whole week
295 (Monday 28 March to Sunday 3 April, 2011). Based on the timelines of mobile phone
296 data for these seven days, we use the linear combination of the normalized pattern and
297 the call volume. The pattern is a four-day mode, i.e., general weekday, Friday,
298 Saturday, and Sunday, where the general weekday is the average pattern for Monday,
299 Tuesday, Wednesday, and Thursday. To clarify our choice of the four-day mode, we
300 consider the normalized timeline (i.e., the pattern) between different days (Table 1).
301 We choose the four-day mode for two reasons. First, Monday, Tuesday, Wednesday,
302 and Thursday are similar, and can be considered as one mode. From Table 1, we can

303 see that the three closest neighbors to each of Monday, Tuesday, Wednesday, and
304 Thursday are all from these four days themselves. For example, Tuesday, Wednesday,
305 and Thursday are closer to Monday than the other three days (i.e., Friday, Saturday,
306 and Sunday) in terms of the normalized pattern distance. (Interestingly, in most cases,
307 the temporally closer are any two of these four days, the smaller the time series
308 distance between them.) Therefore, the data for Monday–Thursday are averaged to
309 represent an ordinary weekday. Second, Friday, Saturday, and Sunday show
310 significant differences, and can be considered as three separate modes. Table 1
311 indicates that each of Friday, Saturday, and Sunday are far away from all the other
312 days. As a result, we choose this four-day mode for land use classification. This
313 ordinary weekday and the remaining three days form a 96-point time series. The
314 comparison of the detection rate between the four-day mode, the two-day mode (an
315 average weekday and an average weekend) and the seven-day mode also confirms
316 that this processing generates the best classification result (see the discussion in the
317 supplementary document).

318

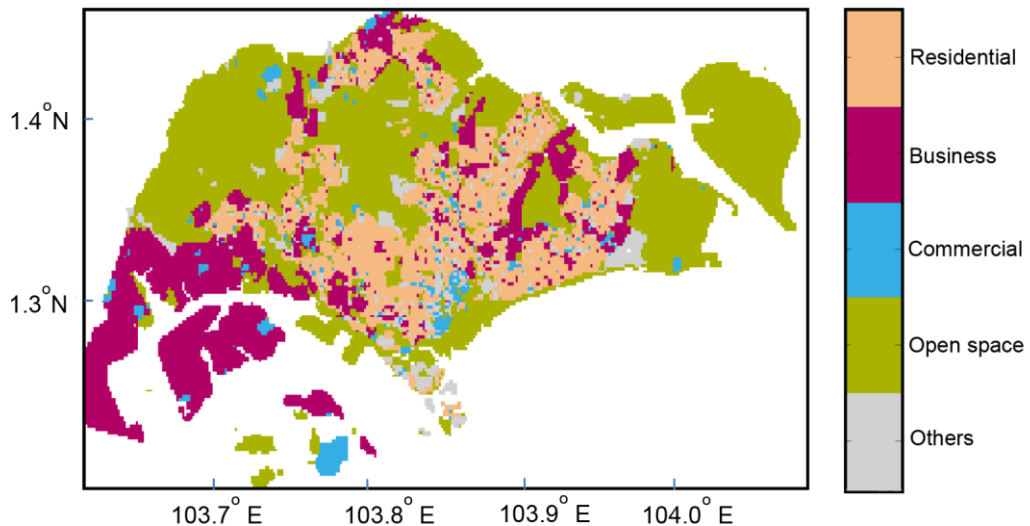
319 Table 1. Distance of normalized pattern between different days

	Mon.	Tue	Wed	Thu	Fri	Sat	Sun
Mon	0	0.0049	0.0089	0.0103	0.0175	0.0245	0.0388
Tue	0.0049	0	0.0057	0.0072	0.0137	0.0224	0.0359
Wed	0.0089	0.0057	0	0.0067	0.0099	0.0223	0.0332
Thu	0.0103	0.0072	0.0067	0	0.0113	0.0201	0.0301

Fri	0.0175	0.0137	0.0099	0.0113	0	0.0216	0.0283
Sat	0.0245	0.0224	0.0223	0.0201	0.0216	0	0.0231
Sun	0.0388	0.0359	0.0332	0.0301	0.0283	0.0231	0

320

321 In order to validate the clustering result, we use the urban planning map of
 322 Singapore, taken from the website
 323 [http://www.ur.gov.sg/uramaps/?config=config_preopen.xml&preopen=Master%20Pl](http://www.ur.gov.sg/uramaps/?config=config_preopen.xml&preopen=Master%20Plan)
 324 [an](http://www.ur.gov.sg/uramaps/?config=config_preopen.xml&preopen=Master%20Plan), and combine land use types to form the ultimate map (Figure 1). Here, we have
 325 divided Singapore into five land use types: Residential, Business, Commercial, Open
 326 space, and Others. Prior to classification, we interpolate the aggregated hourly data
 327 into a 200 m × 200 m grid using IDW, and generate 96 pattern layers and one volume
 328 layer.



329

330 Figure 1. Land use in Singapore

331

332 5. Land use classification for Singapore

333 5.1. Determination of land use types

334 After generating 97 image layers, the first 96 are transformed using equation (2) to
335 generate X_i , and the final layer is transformed using equation (3) to generate Y_i . As
336 discussed above, we combine the pattern (X_i) and the volume (Y_i) to form a new time
337 series Z_i using the coefficient β (see equation (1)). Next, we determine the value
338 of β through the following training process. First, 105 samples (allocated based on
339 the prior knowledge of the areas of different land use types: 25 samples each for
340 Residential, Business, and Open space, 20 samples for Commercial, and 10 samples
341 for Others) are chosen based on remote sensing imagery and POI data (from Google
342 Earth) as well as information provided by several residents of Singapore. To ensure
343 the samples represent their land use types, we select them according to three criteria.
344 First, samples are picked from homogeneous areas. Second, we avoid samples from
345 near the boundary between different land use types. Third, we attempt to pick samples
346 that are close to a BTS tower. The objective function $f(\beta)$ is calculated at different
347 values of β , and the results are shown in Figure 2. We can see that the minimum
348 value is acquired when β is between 0.65 and 0.80.

349

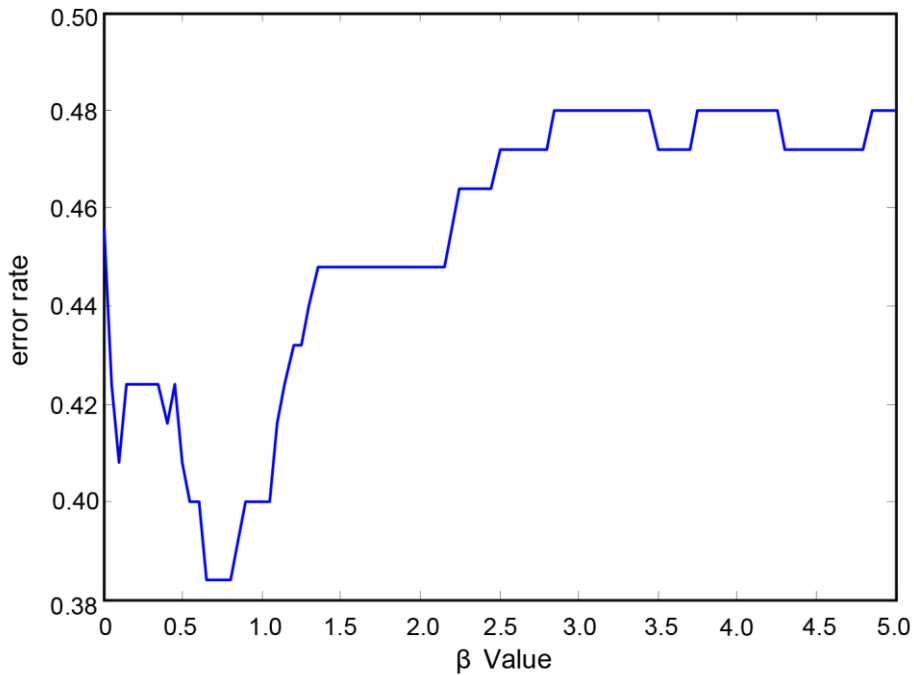


Figure 2. Error rates generated at different values of β

350

351

352

353 The sample centers of different land use types are shown in Figures 3 and 4.

354 Figure 3 shows the pattern part of the centers, each of which contains 96 points.

355 Figure 4 is a boxplot of the volume of each land use. We can see that all land use

356 types can be characterized by a combination of pattern and volume. For example,

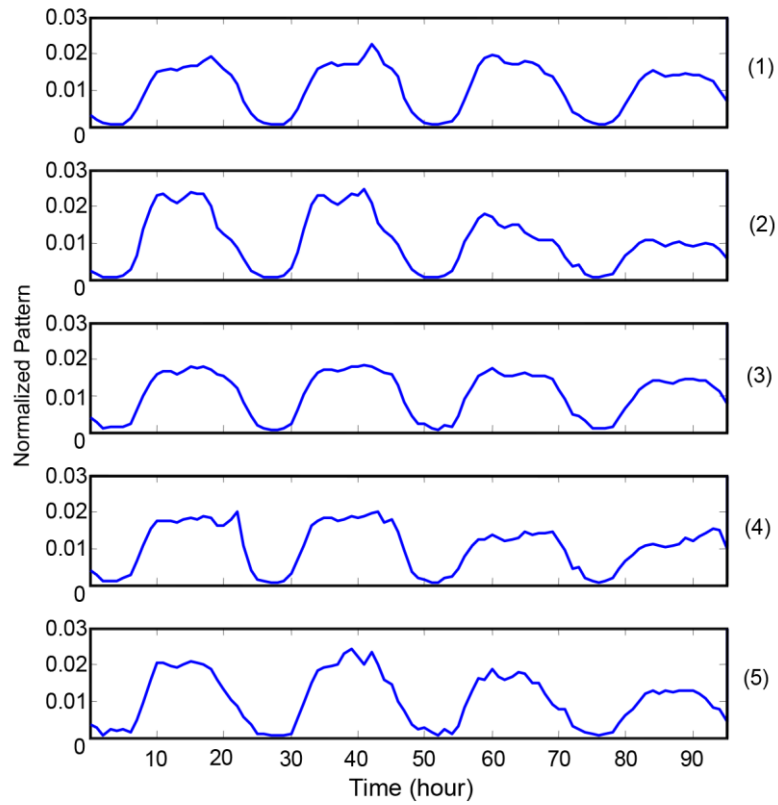
357 Residential areas are characterized by a similar size pattern for each of the four days

358 and medium volume, whereas Business areas are characterized by a high-thin pattern

359 on the ordinary weekday and Friday, a low weekend pattern, and low volume. The

360 other land use types can be similarly characterized. The characteristics of each time

361 series guarantee the classification of land use type.

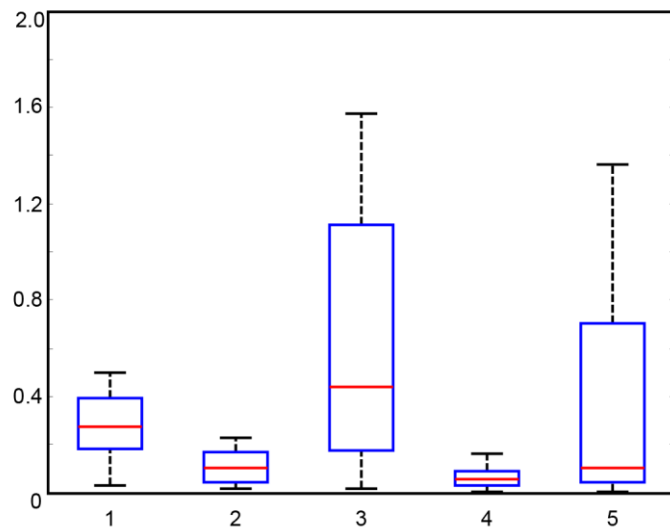


362

363 Figure 3. Patterns of centers of time series samples with $\beta = 0.75$

364 (1-Residential; 2-Business; 3-Commercial; 4-Open space; 5-Others)

365



366

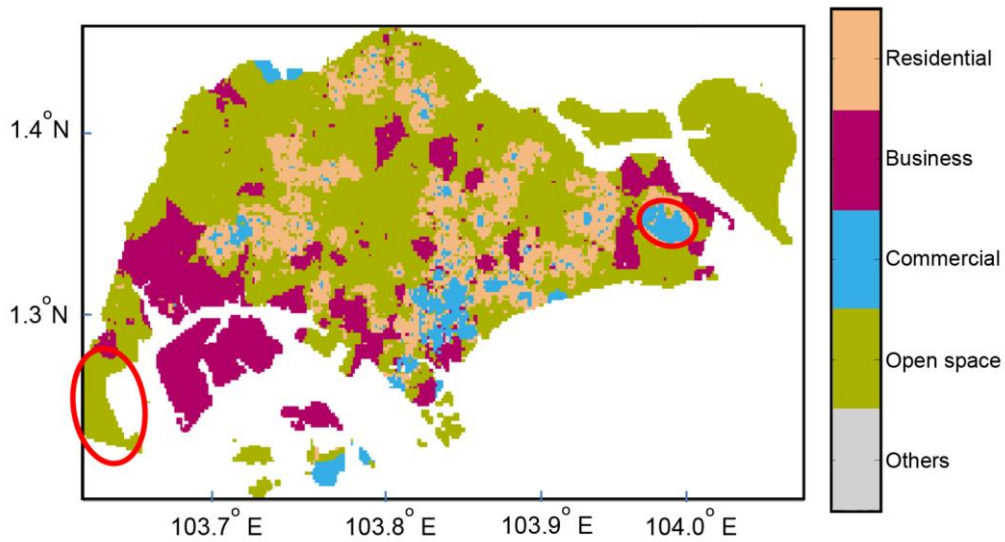
367 Figure 4. Volume of time series samples with $\beta = 0.75$

368 (1-Residential; 2-Business; 3-Commercial; 4-Open space; 5-Others)

369

370 5.2. Clustering result

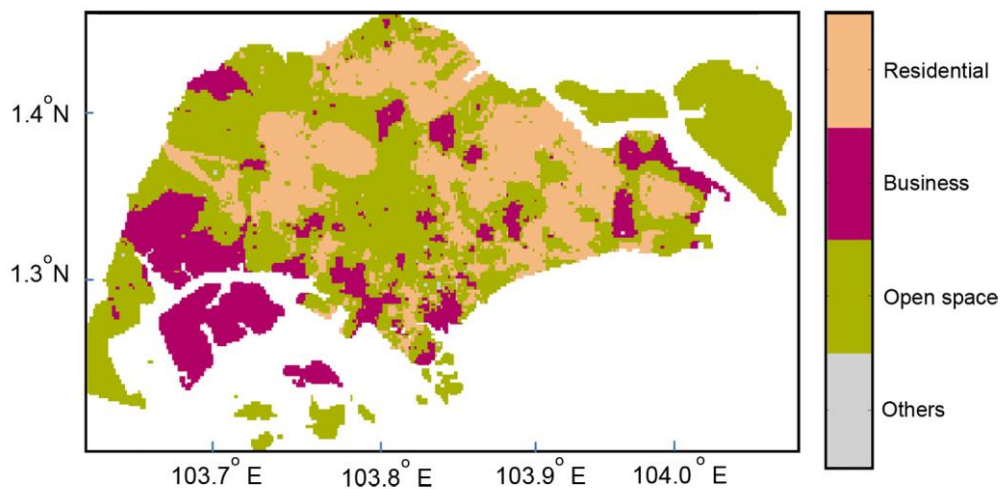
371 We use FCM to cluster the aggregated data by setting β to 0.75, based on the
372 training result. The cluster number is determined by the validity indices, which
373 indicate that the optimum cluster number is 6. After post-processing, two clusters are
374 combined and determined as Open space. Finally, we generate the land use map
375 displayed in Figure 5(a).



376

377

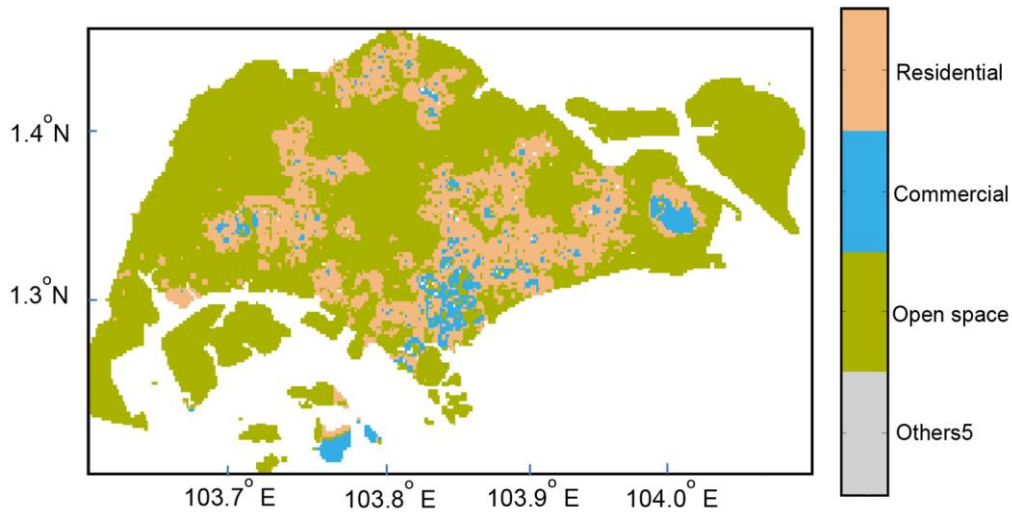
(a)



378

379

(b)



(c)

Figure 5. Clustering result for land use types in Singapore

(a) Classification generated from the synthetic time series (detection rate: 58.03%; the left red ellipse indicates the area defined as Commercial in Figure 1 is identified as Open Space; the right red ellipse indicates the area defined as Open Space in Figure 1 is identified as Commercial). (b) Classification generated from the pattern data (detection rate: 52.58%). (c) Classification generated from the volume data (detection rate: 52.68 %).

Comparing the classification result with the urban planning map (Figure 1), we find that all land use types are identified with an overall detection rate of 58.03%, which is close to that generated by Toole et al. (2012) (The detection rate is 54%). In the supplementary document, we also showed that four-day mode generates the highest detection rate compared with that for two-day mode (57.65%) and for seven-day mode (55.15%). The confusion matrix is shown in Table 2. From this table, we can see that the order in which the land use types are best detected is Open space,

397 Residential, Business, Commercial, and Others (this can be determined from the
 398 diagonal elements in the matrix, which mean the land use is correctly classified). Only
 399 Residential, Business, and Open space land use types have rates close to or above
 400 50%. The detection rates of Commercial and Others are less than 50%. In addition,
 401 some land use types have a misclassification rate of over 30%. Overall, land use is
 402 most commonly misclassified as Open space, while Others is the most likely to be
 403 misclassified.

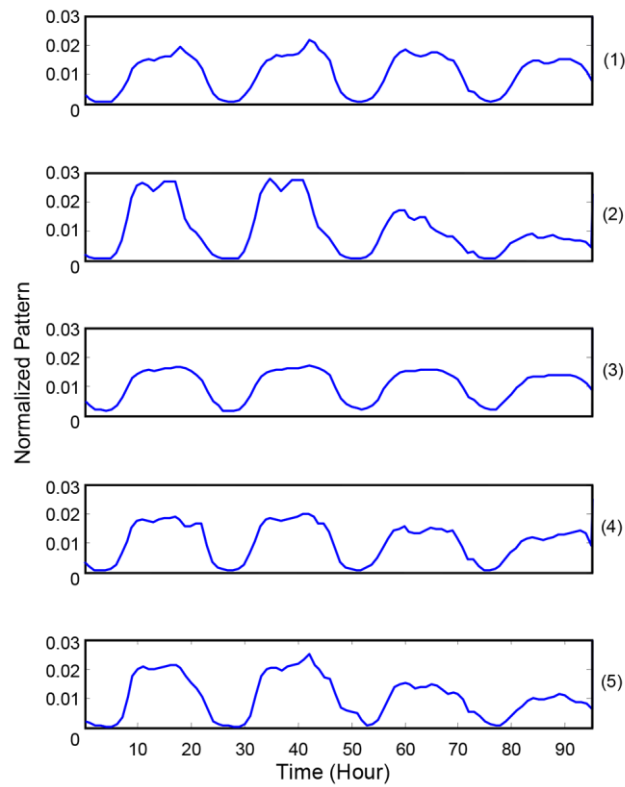
404
 405 Table 2. Confusion matrix of the classification

	Residential	Business	Commercial	Open space	Others
Residential	0.4912	0.0490	0.0658	0.3938	0.0002
Business	0.0978	0.5018	0.0174	0.3825	0.0005
Commercial	0.1612	0.1535	0.3457	0.3302	0.0093
Open space	0.0769	0.1210	0.0395	0.7622	0.0004
Others	0.0037	0.1737	0.0772	0.5026	0.2428

406
 407 To determine the reasons for this particular land use classification, we draw the
 408 center of each real land use type and that of each cluster in Figure 6. Comparing the
 409 two, we find that the Residential, Business, and Open space regions generated by our
 410 method show both a similar pattern (Figure 6a and c) and volume (Figure 6b and d) as
 411 the real land use types. Although Others in Figure 6a shows a similar pattern to the
 412 real one (“5” in Figure 6c), its volume (“5” in Figure 6b) is somewhat different
 413 (Figure 6d). The Commercial volume (“3” in Figure 6b) suggested by the clustering

414 has a larger value than the actual volume (“3” in Figure 6d), and its pattern is also
415 different (“3” in Figure 6a and c). This shows why Residential, Business, and Open
416 space have high detection rates while Commercial and Others have lower ones.

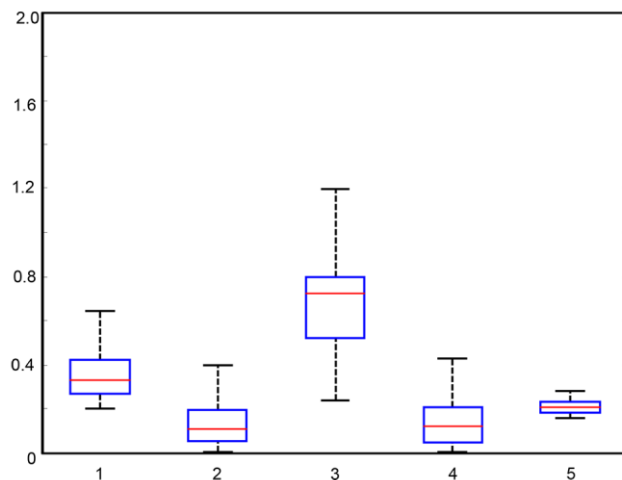
417



418

419

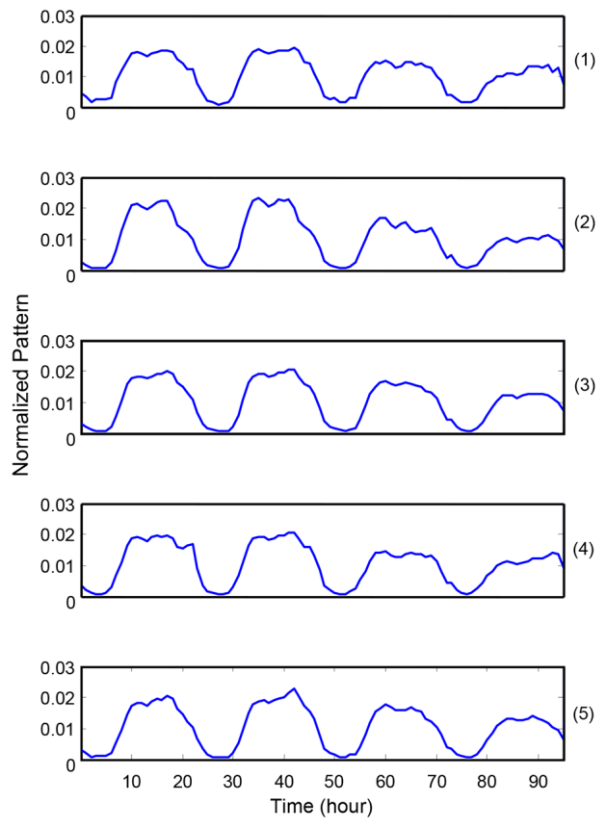
(a)



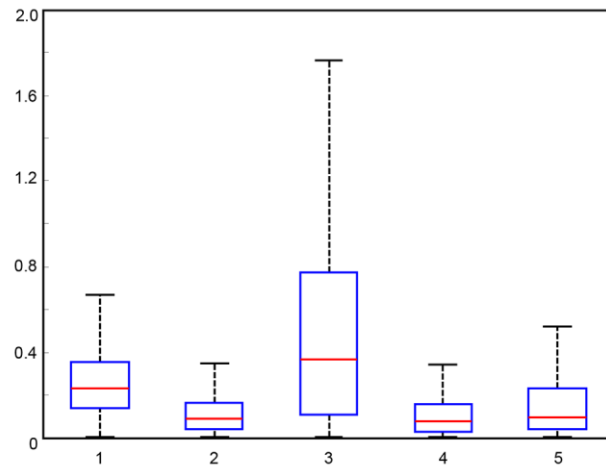
420

421

(b)



(c)



(d)

Figure 6. Centers of clusters for different land use types

(1-Residential; 2-Business; 3-Commercial; 4-Open space; 5-Others)

(a) Centers of pattern for classification; (b) Distribution of volume for classification;

(c) Centers of patterns for known land use; (d) Distribution of volume of known land

430

use

431

432 5.3. Evaluation of the effect of call pattern and volume on classification

433 We now examine how the value of β influences the detection rate. The
434 detection rate calculated for different values of β is shown in Table 3. The detection
435 rate generally increases with β until $\beta = 0.75$, then decreases for $\beta > 0.75$.

436

437 Table 3. Change in detection rate with β (four-day mode)

β value	0	0.15	0.30	0.40	0.50	0.60	0.65	0.70
Detection rate (%)	52.58	54.30	55.12	57.56	56.50	57.51	57.57	57.97
β value	0.75	0.8	0.9	1.00	1.25	1.50	2.50	∞
Detection rate (%)	58.03	57.30	55.61	55.44	54.54	54.24	54.01	52.68

438

439 As discussed in Section 2, the distance between samples and the cluster centers is
440 calculated during the FCM algorithm. The distance consists of two parts. The first (d_1)
441 is the distance between the patterns, and the second (d_2) is that between the volumes
442 weighted by β . Essentially, the value of β represents the balance between call
443 pattern and call volume, both of which are normalized. As β decreases, the weight
444 of the pattern part in the overall distance between samples and centers will increase.
445 On the contrary, as β increases, the weight of the volume part will increase. The
446 next issue is to determine which part dominates the distance (i.e., the difference in
447 discerning between land use types) in the classification generated at the optimized

448 value of β ($\beta = 0.75$). We calculated the ratio between d_1 and d_2 for all land
449 use types classified with $\beta = 0.75$. The results are given in Table 4. From this table,
450 we can see that the ratio is greater than 1 for all land use types except Commercial.
451 The average ratio is 1.6471, which indicates that the distance between the patterns is
452 generally larger than those between the weighted volumes. The ratios for different
453 land use types implies that the pattern information plays a more important role in the
454 classification for all land use types, with the exception of Commercial areas. This is
455 also consistent with the differences in the time series of different land use types,
456 which can be found in Figure 6. Specifically, Commercial has the highest volume,
457 which is significantly different from the other land use types. This causes the volume
458 to play a more important role in separating Commercial from the other types. On the
459 contrary, the other land use types show more significant differences between the
460 patterns than the volume, which leads to the larger distances between the patterns.
461 This analysis of the effect of the call pattern and volume shows that our method can
462 utilize different characteristics of mobile phone data to differentiate between land use
463 types.

464

465 Table 4. Ratio between pattern and volume for different land use types

Land use type	Residential	Business	Commercial	Open space	Others	Average
Ratio between Pattern and volume	1.1462	2.0758	0.9594	2.5467	1.5072	1.6471

466

467 6. Comparison between classifications using different information

468 To further validate the method based on the newly constructed time series, we

469 compare the classification with that generated with either the pattern or the volume.
470 The clustering validity index shows that five clusters are generated for pattern
471 information only, while four clusters are generated for the volume. The results are
472 shown in Figure 5b and c. Figure 5b indicates that the clustering based on the pattern
473 information did not identify Commercial areas, and Figure 5c indicates that the
474 clustering based on volume data did not identify the Business regions. The overall
475 detection rates are also lower (52.58% for pattern and 52.68% for volume) than that
476 based on the combination of pattern and volume.

477 The pattern information fails to identify Commercial areas because these are
478 highly mixed with Residential areas. According to the Master Plan 2008 of Singapore,
479 more than 45% of the Commercial area is either “residential with commercial on the
480 first floor” or a “mixture of commercial and residential”. This highly mixed
481 distribution causes difficulties in discerning Residential from Commercial. To
482 quantify the degree of mixing between different land use types, we can calculate the
483 posterior classification based on the pattern information, in which the land use type
484 over a cell is determined by locating the minimum distance between the pattern part
485 and the centers of known land use types. We generate the posterior confusion matrix
486 by comparing the posterior classification with the Master Plan 2008 (Table 5). This
487 shows that only 9.89% of Commercial areas are correctly classified, with 40.54%
488 mixed into Residential. This also explains why the Commercial land use type is not
489 identified from pattern information alone.

490

491

Table 5. Posterior confusion matrix of pattern information

	Residential	Business	Commercial	Open space	Others
Residence	0.6708	0.0731	0.0571	0.0138	0.1852
Business	0.1299	0.5842	0.0279	0.2285	0.0296
Commercial	0.4054	0.2679	0.0989	0.1032	0.1246
Open space	0.1645	0.3297	0.0557	0.3478	0.1024
Others	0.4640	0.2685	0.0462	0.0483	0.1729

492

493 The classification based on volume fails to detect Business land use because this
 494 volume shows no significant difference from that of Open space. The box plot of each
 495 land use type is shown in Figure 6d, indicating that Business (“2” in the figure) and
 496 Open space (“4” in the figure) have very similar median values and ranges. In this
 497 case, these two land use types cannot be separated merely by their volume, which
 498 cause only four land use types to be identified.

499

500 7. Discussion

501 In this section, we analyze the possible causes of errors generated by our
 502 method. There are four factors that may affect the error rate of the classification. The
 503 first is the difference between the definition of land use in urban planning and the
 504 function derived from the mobile phone data. The second is the degree of
 505 heterogeneity of different land use types (i.e., different land use types are mixed in the
 506 same area). The third is the precision of the information recorded, which is related to

507 the density of BTSs in each cell. The fourth is the fuzzy membership threshold (α -cut)
508 used in FCM.

509

510 7.1. Dissimilarity between definition of land use and that derived from the mobile
511 phone data

512 Previous research has found that zoned areas are not necessarily used as intended,
513 which may lead to incorrect classification (Soto and Frias-Martinez, 2011a; Toole et
514 al., 2012). However, these studies only provided some examples, without
515 summarizing all scenarios. Here, we try to list all possible situations. The first is when
516 various social activities are conducted on one land use type. As mentioned above, a
517 large portion of the residential area in Singapore is mixed with the commercial area.
518 The second is the heterogeneity of a land use type. For example, the airport is a
519 homogenous area in the Master Plan 2008, but the landing area and the terminals in
520 the airport are different in terms of social function. Thus, in the result generated by the
521 mobile phone data, the terminal is classified as Commercial, whereas the landing area
522 is classified as Open space (Figure 5a). This is because the terminal exhibits a very
523 high volume, while that of the landing area is very low. The third is that some areas
524 with specific uses are reserved for other uses in the future. For example, the western
525 part of the business area located in southwest Singapore is “misclassified” as Open
526 space by the mobile phone data (Figure 5a). In fact, this area is an empty space (this
527 can be confirmed from remote sensing images in Google Earth) that is reserved for
528 future business use.

529

530 7.2. Correlation between the error rate and BTS density

531 As we know, the volume of each BTS is calculated by aggregating the number of
532 calls in the polygon generated by Voronoi tessellation (Okabe et al., 2000). When the
533 BTS density is low (i.e., the area of the Voronoi polygon is large), there is a risk that
534 the volume may include calls from areas of different land use. On the contrary, when
535 the BTS density is high, calls collected in this area will have less “interference”, i.e.,
536 the signal is “purer”. In order to determine if the purity of signal affects the precision
537 of land use classification, we calculated the detection rates for different BTS densities
538 (Table 6). Note the density in this table is represented by the number of BTSs in each
539 cell. From the table, we can see that the detection rate increases with the BTS density,
540 except when the density is 0. Interestingly, the detection rate attains a relatively high
541 value (i.e., 60.56%) when the density is 0. This is because most of the cells that have a
542 density of 0 are Open space. As the signals in Open space are “purer”, the detection
543 rate in these cells is high. As a result, we can conclude that the “purer” the signal
544 recorded by a BTS (either in the homogenous and large areas with low BTS density or
545 in areas with a high BTS density), the higher the precision of the classification.

546

547 Table 6. Relationship between error rate and BTS density

Towers Density	0	1	2	3	4	5	6	7	8	11
Detection rate (%)	60.56	44.81	50.78	51.18	52.94	57.14	58.82	75.00	75.00	100.00
Number of cells	16548	2522	963	211	68	21	17	4	4	1

548

549 7.3. Relationship between error rate and mixture entropy

550 Another factor that might influence the precision is the mixture of the land use.
551 Because the resolution of Singapore's Master Plan 2008 is much higher (4 m) than
552 that of our classification (200 m), we can calculate the error rates in terms of the land
553 use entropy (En_j), which measures the randomness of the areas of different land use
554 types in each cell as:

$$555 \quad En_j = -\sum_i p_{i,j} \ln(p_{i,j}) \quad (7)$$

556 , where $p_{i,j}$ is the occupancy rate of the area of land use type i in cell j .

557 The relationship between the error rate and the land use entropy is shown in
558 Figure 7. It is interesting to see that the error rate increases with the land use entropy.
559 The reason for this is obvious. If the entropy of a cell is high, which means more land
560 use types coexist in the cell (i.e. the cell is more heterogeneous), then the error rate of
561 the classification increases. The average entropy for residential, business, commercial,
562 open space and others are 0.42, 0.18, 0.47, 0.084 and 0.57, respectively. We can see
563 that the lower the entropy of some land use type, the higher the detection rate (Table
564 2).

565

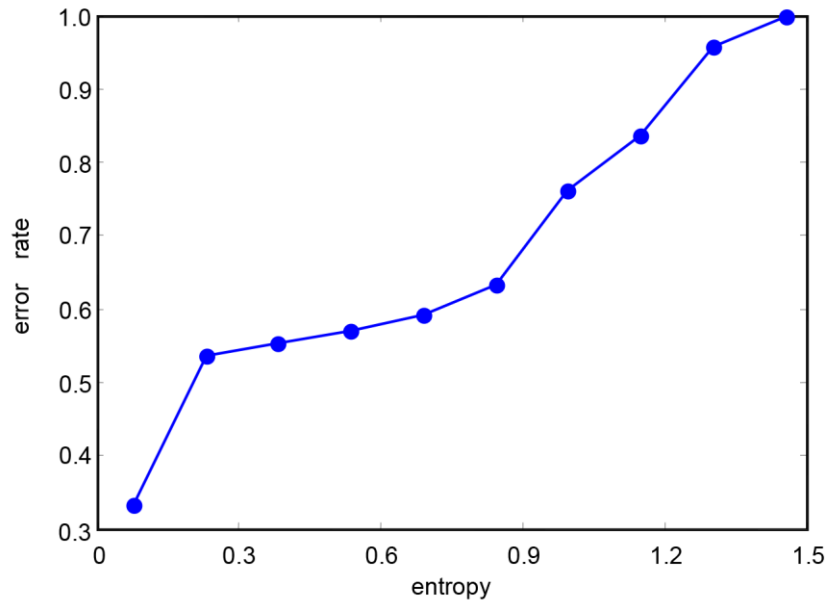


Figure 7. Relationship between land use entropy and error rate

566

567

568

569 7.4. Relationship between error rate and fuzzy membership value

570 As we know, the FCM result includes the fuzzy membership value of a sample

571 belonging to each cluster for a certain value of α -cut. Our question is: how will the

572 detection rate change if we change the value of α -cut? The detection rates obtained

573 with different α -cut values are listed in Table 7. We can see that the detection rate is

574 60.39% when α -cut is 0.5, and that 85.46% of the total area has a membership value

575 greater than 0.5. As α -cut increases to 0.8, only 45.32% of the total area attains this

576 membership value, although the detection rate increases to 72.89%. We can conclude

577 that the detection rate increases with α -cut, but must bear in mind that the area with

578 such a detection rate will decrease.

579

580 Table 7. Detection rates at different values of α -cut

Value of α -cut	0.5	0.6	0.7	0.8	0.9
Detection rate (%)	60.39	61.10	65.41	72.89	88.73
Percentage of area with membership value larger than α -cut	85.46	73.35	60.27	45.32	29.16

581

582 8. Conclusions and future work

583 In this paper, we constructed a synthesized time series of mobile phone activity
584 to identify land use types using a semi-supervised clustering method. The synthesized
585 time series was obtained as a linear combination of the (four-day) pattern and the
586 volume of aggregated data by introducing the weighting coefficient β . Our
587 classification of land use in Singapore produced a detection rate of 58.03% with β
588 set to its optimized value of 0.75, as determined by a training process. Comparisons
589 show that: (1) the data combining both the pattern and volume generate better
590 classifications than those based on either the pattern or the volume alone; (2) four-day
591 mode generates the higher detection rate than that of two-day mode and that of
592 seven-day mode. We can analyze the importance of different parts of the constructed
593 time series on the overall classification, as well as on each type of land use. The
594 results show the relative importance of ‘pattern’ over ‘volume’ in detecting most land
595 use types.

596 We also determined some factors that influence the accuracy of the land use
597 classification. First, there are substantial differences between the urban planning map
598 and the land use retrieved from mobile phone data. Second, areas of mixed land use
599 result in heterogeneous mobile phone usage, and thereby increase the error rate. Third,
600 the purity of the signal in each cell, essentially the BTS density, influences the
601 precision of classification. In general, the higher the density, the higher the precision

602 generated by the classification, except for areas where the density is 0. This indicates
603 that land use classification based on mobile phone data might generate good results in
604 areas with a high BTS density and pure land use types.

605 Our analysis shows that mobile phone data can reveal the social function of land
606 use. Nevertheless, the overall detection rate of less than 60% indicates that mobile
607 phone data alone are not adequate for urban land use classification, although in some
608 areas the data generate relatively high detection rates (e.g., areas with high BTS
609 density, pure land use, and a high fuzzy membership value). Future research can be
610 extended in the following two directions. The first is to improve the classification
611 model. One idea is to vary the parameter β over space to effectively capture the
612 characteristics of different land use types. The second is to merge more information
613 into the classification, such as remote sensing data and POI.

614
615

616 **References**

- 617 Andrienko, G., Andrienko, N. and Fuchs, G., 2013, Multi-perspective analysis of D4D
618 fine resolution data. In: Blondel V, Cordes N, Decuyper A, Deville P, Raguenez J,
619 Smoreda Z eds, Mobile phone data for development (Analysis of mobile phone
620 datasets for the development of Ivory Coast), Cambridge, MA, USA, May 1-3,
621 2013, No. 37.
- 622 Berthold, M.R., 2010, Guide to Intelligent Data Analysis. London: Springer. 394 p.
- 623 Bezdek, J.C., 1981, Pattern Recognition with Fuzzy Objective Function Algorithms.
624 Norwell: Kluwer Academic Publishers. 256 p.
- 625 Breiman, L., 2001, Random Forests. Machine Learning, 45(1), pp. 5-32.
- 626 Caceres, R., Rowland, J., Small, C. and Urbanek, S., 2012, Exploring the Use of
627 Urban Greenspace through Cellular Network Activity. In: The Second Workshop
628 on Pervasive Urban Applications (PURBA), In conjunction with Pervasive 2012,
629 Newcastle, UK, June 18-22, 2012, pp. 1-8.
- 630 Calabrese, F., Reades, J. and Ratti, C., 2010, Eigenplaces: Segmenting Space through
631 Digital Signatures. IEEE Pervasive Computing, 9(1), pp. 78-84.
- 632 Calabrese, F., Lorenzo, G.D., Liu, L. and Ratti, C., 2011, Estimating
633 Origin-Destination Flows Using Mobile Phone Location Data. IEEE Pervasive
634 Computing 10(4), pp. 36-44.
- 635 De Wit, A.J.W., and Clevers, J.G.P.W., 2004, Efficiency and Accuracy of Per-field
636 Classification for Operational Crop Mapping. International Journal of Remote

637 Sensing, 25, pp. 4091–4112.

638 Fisher, P., 1997, The Pixel: A Snare and a Delusion. *International Journal of Remote*
639 *Sensing*, 18, pp. 679–85.

640 Frias-Martinez, V., Soto, V., Hohwald, H. and Frias-Martinez, E., 2012,
641 Characterizing Urban Landscapes using Geolocated Tweets, *International*
642 *Conference on Social Computing (SocialCom)*, Amsterdam, The Netherlands,
643 September 3-6, 2012, pp. 1-10.

644 Gonzalez, M., Hidalgo, C. and Barabasi, A., 2008, Understanding individual human
645 mobility patterns, *Nature*, 453, pp. 779-782.

646 Gong, P., and Howarth, P., 1990, The use of structural information for improving
647 land-cover classification accuracies at the rural-urban fringe, *Photogramm. Eng.*
648 *Remote Sens.*, 56(1), pp. 67–73.

649 Hu, S.G. and Wang, L., 2013, Automated urban land-use classification with remote
650 sensing, *International Journal of Remote Sensing*, 34(3), pp. 790-803.

651 Jacobs-Crisioni, C.G.W. and Koomen, E., 2012, Linking urban structure and activity
652 dynamics using cell phone usage data. In: *Proceedings of the workshop on*
653 *Complexity Modeling for Urban Structure and Dynamics for AGILE2012*,
654 Avignon, France, April 24-27.

655 Jia, T. and Jiang, B., 2012, Exploring Human Activity Patterns Using Taxicab Static
656 Points, *ISPRS International Journal of Geo-Information*, 1, pp. 89-107.

657 Laura, F., Marco, M. and Massimo, C., 2012, Discovering events in the city via
658 mobile network analysis. *Journal of Ambient Intelligence and Humanized*
659 *Computing*, doi: 10.1007/s12652-012-0169-0.

660 Liu, Y., Wang, F.H., Xiao, Y. and Gao, S., 2012, Urban land uses and traffic
661 ‘source-sink areas’: Evidence from GPS-enabled taxi data in Shanghai.
662 *Landscape and Urban Planning*, 106, pp. 73-87.

663 Loibl, W. and Peters-Anders, J., 2012, Mobile Phone Data as Source to Discover
664 Spatial Activity and Motion Patterns. In: *Jekel T, Car A, Strobl J, Griesebner G*
665 *(Eds.) (2012): GI_Forum 2012: Geovizualisation, Society and Learning.*
666 *Wichmann Verlag, Berlin & Offenbach*, July 1, 2012, pp. 524-532.

667 Lu, D. and Weng, Q., 2006, Use of Impervious Surface in Urban Land-Use
668 Classification. *Remote Sensing of Environment*, 102, pp. 146–60.

669 Manfredini, F., Tagliolato, P. and Rosa, C.D., 2011, Monitoring Temporary
670 Populations through Cellular Core Network Data. *Lecture Notes in Computer*
671 *Science*, 6783, pp. 151-161.

672 Nock, R. and Nielsen, F., 2006, On Weighting Clustering, *IEEE Trans. on Pattern*
673 *Analysis and Machine Intelligence*, 28 (8), pp. 1–13.

674 Okabe, A., Boots, B., Sugihara, K. and Chiu, S.N., 2000, *Spatial Tessellations –*
675 *Concepts and Applications of Voronoi Diagrams (2nd edition)*. John Wiley. 671p.

676 Platt, R.V. and Rapoza, L., 2008, An Evaluation of an Object Oriented Paradigm for
677 Land Use/Land Cover Classification. *The Professional Geographer*, 60, pp.
678 87–100.

679 Pulselli, R.M., Ratti, C. and Tiezzi, E., 2006, City Out of Chaos: Social Patterns and
680 Organization. In: *Urban Systems*. *International Journal of Ecodynamics*, 1(2), pp.

681 125-134.

682 Ratti, C., Pulselli, R. M., Williams, S. and Frenchman, D., 2006, Mobile Landscapes:
683 Using Location Data from Cell Phones for Urban Analysis. *Environment and*
684 *Planning B*, 33(5), pp. 727-748.

685 Ray, S. and Turi, R.H., 1999, Determination of number of clusters in k-means
686 clustering and application in color image segmentation. In: Pal NR, De AK, Das J
687 (eds), *Proceedings of the 4th International Conference on Advances in Pattern*
688 *Recognition and Digital Techniques (ICAPRDT'99)*, Calcutta, India, December,
689 27-29, 1999, New Delhi, India: Narosa Publishing House, pp. 137-143.

690 Reades, J., Calabrese, F. and Ratti, C., 2009, Eigenplaces: analysing cities using the
691 space-time structure of the mobile phone network. *Environ Planning B*, 36(5), pp.
692 824 – 836.

693 Rubioa, A., Sanchezb, A. and Frias-Martinez, E., 2013, Adaptive non-parametric
694 identification of dense areas using cell phone records for urban analysis,
695 *Engineering Applications of Artificial Intelligence*, 26(1), pp. 551 – 563.

696 Sevtsuk, A. and Ratti, C., 2010, Does Urban Mobility Have a Daily Routine?
697 Learning from the Aggregate Data of Mobile Networks. *Journal of Urban*
698 *Technology*, 17, pp. 41–60.

699 Shaban, M.A., and Dikshit, O., 2001, Improvement of Classification in Urban Areas
700 by the Use of Textural Features: The Case Study of Lucknow City, Uttar Pradesh.
701 *International Journal of Remote Sensing*, 22, pp. 565–93.

702 Song, C., Qu, Z., Blumm, N., and Barabasi, A.-L., 2010, Limits of Predictability in
703 Human Mobility. *Science*, 327(5968), pp. 1018-1021.

704 Soto, V. and Frias-Martinez, E., 2011a, Automated land use identification using
705 cell-phone records. In: *Proceedings of the 3rd ACM international workshop on*
706 *MobiArch, HotPlanet '11*, Bethesda, Maryland, USA, June 28-28, 2011,
707 doi:10.1145/2000172.2000179, pp. 17-22.

708 Soto, V., and Frias-Martinez, E., 2011b, Robust Land Use Characterization of Urban
709 Landscapes using Cell Phone Data, *First Workshop on Pervasive Urban*
710 *Applications*, San Francisco, USA, June 12-15, pp. 1-8.

711 Sun, J.B., Yuan, J., Wang, Y., Si, H.B. and Shan, X.M., 2011, Exploring space–time
712 structure of human mobility in urban space. *Physica A*, 390, pp. 929–942.

713 Toole, J.L., Ulm, M., González, M.C. and Bauer, D., 2012, Inferring land use from
714 mobile phone activity. In: *Proceedings of the ACM SIGKDD International*
715 *Workshop on Urban Computing*, Beijing, China, August 12-12, 2012,
716 doi:10.1145/2346496.2346498.

717 Traag, V.A., Browet, A., Calabrese, F. and Morlot, F., 2011, Social Event Detection in
718 Massive Mobile Phone Data Using Probabilistic Location Inference. In:
719 *Proceeding of IEEE SocialCom*, Boston, MA, October 9-11, pp. 1-4.

720 Vaccari, A., Gerber, A., Biderman, A. and Ratti, C., 2009, Towards estimating the
721 presence of visitors from the aggregate mobile phone network activity they
722 generate. In: *Proceedings of the 11th International Conference on Computers in*
723 *Urban Planning and Urban Management*, Hong. Kong, 16th -18th June, pp. 1-11.

724 Vieira, M.R., Frias-Martinez, V., Nuria, O. and Frias-Martinez, E., 2010,

725 Characterizing Dense Urban Areas from Mobile Phone-Call Data: Discovery and
726 Social Dynamics, In: Proceedings of IEEE International Conference on Social
727 Computing / IEEE International Conference on Privacy, Security, Risk and Trust,
728 Minneapolis, MN, USA, August 20-22, 2010, pp. 241-248.
729 Wu, S., Qiu, X., Usery, L. and Wang, L., 2009, Using Geometrical, Textural, and
730 Contextual Information of Land Parcels for Classification of Detailed Urban
731 Land Use. *Annals of the Association of American Geographers*, 99, pp. 1–23.