

1 Extracting Urban Functional Regions from Points
2 of Interest and Human Activities on
3 Location-based Social Networks

4 Song Gao, Krzysztof Janowicz, Helen Couclelis

5 Department of Geography,
6 University of California, Santa Barbara, CA, USA

7 **Abstract:** Data about points of interest (POI) have been widely used in
8 studying urban land use types and for sensing human behaviors. However, it
9 is difficult to quantify the right mix or the spatial relations among different
10 POI types indicative of specific urban functions. In this research, we develop a
11 statistical framework to help discover semantically meaningful topics and func-
12 tional regions based on the co-occurrence patterns of POI types. The framework
13 applies the latent Dirichlet allocation (LDA) topic modeling technique and in-
14 corporates user check-in activities on location-based social networks. Using a
15 large corpus of about 100,000 Foursquare venues and user check-in behaviors
16 in the ten most populated urban areas of the United States, we demonstrate
17 the effectiveness of our proposed methodology by identifying distinctive types of
18 latent topics and further, by extracting urban functional regions using K-means
19 clustering and Delaunay triangulation spatial constraints clustering. We show
20 that a region can support multiple functions but with different probabilities,

21 while the same type of functional region can span multiple geographically non-
22 adjacent locations. Since each region can be modeled as a vector consisting of
23 multinomial topic distributions, similar regions with regard to their thematic
24 topic signatures can be identified. Compared to remote sensing images which
25 mainly uncover the physical landscape of urban environments, our popularity-
26 based POI topic modeling approach can be seen as a complementary social
27 sensing view on urban space based on human activities.

28 **Keywords:** urban functions, topic modeling, latent Dirichlet allocation,
29 LBSN, place type.

30 1 Introduction

31 Cities support a variety of functions that relate to land use types, including
32 residential, commercial, industrial, transportation, and business regions and in-
33 frastructure, while affording different types of human activities, such as living,
34 working, commuting, shopping, eating, and recreation. Rapid urbanization and
35 new construction have caused land use changes and urban expansions in many
36 areas. Remote sensing images together with spatial metrics have been widely
37 used to classify urban land use and monitor change at different spatial scales
38 (Barnsley and Barr, 1996; Herold et al., 2005; Banzhaf and Netzband, 2012).
39 However, human activities usually take place in different types of points of in-
40 terest (POIs). Remote sensing techniques perform well in extracting physical
41 characteristics, such as land surface reflectivity and texture of urban space but
42 are not good in identifying functional interaction patterns or in helping under-
43 stand socioeconomic environments (Pei et al., 2014; Liu et al., 2015). Compared
44 to other datasets and methods in remote sensing and field mapping, using POI
45 data, social media, and their associated methods can lead to a better under-
46 standing of individual-level and group-level utilization of urban space at a fine-

47 grained spatial and temporal resolution. Rich *social sensing* techniques can help
48 bridge the semantic gap between land use classification and urban functional
49 regions. The function of a place is determined by what type of activities can
50 occur there (Janowicz, 2012; Zhong et al., 2014). The same types of POIs can
51 be located in different land use types and may also support different functions.
52 For example, *restaurants* are found in residential areas and in commercial areas,
53 as well as in industrial areas. The main function of *universities* is education,
54 but they also support sports activities, music shows, and so on. Previous studies
55 have demonstrated that different POI types have distinctive semantic signatures
56 (Janowicz, 2012) (i.e., spatial, temporal, and thematic distributions) based on
57 crowd-sourced location-based social media data analysis, in analogy to spectral
58 bands in remote sensing (McKenzie et al., 2015). There is a growing trend of
59 using location-awareness sensing data (e.g., trajectories from mobile phones),
60 POI data, and social media feeds to study the spatial and social structure of ur-
61 ban environments (Pei et al., 2014; Hu et al., 2015; Jiang et al., 2015; Liu et al.,
62 2015; McKenzie et al., 2015; Steiger et al., 2016; Yao et al., 2017). However, few
63 studies have investigated the latent relationships among different types of POIs
64 and how they spatially interact with each other to support urban functions,
65 such as education, business, and shopping. In this research, we aim to develop
66 a data-driven framework to discover urban functional regions from POIs and
67 associated human activities on location-based social networks (LBSN).

68 We argue that geographic knowledge and measures of spatial distribution
69 over POI types (categories) can be employed to derive latent classification fea-
70 tures for these said types, which will then enable the detection and the abstrac-
71 tion of higher-level functional regions (i.e., semantically coherent areas of inter-
72 est) such as shopping areas, business districts, educational areas, and tourist
73 zones. To test this claim, we will study the co-occurrence patterns of different

74 POI categories as well as the associated human activities (i.e., mobility, check-
75 ins, reviews, and comments), and thus employ analytical measures to quantify
76 their differences and conduct classifications of functional regions.

77 The contributions of this research are as follows:

- 78 • We propose a novel framework to study urban functional regions by em-
79 ploying data about Points Of Interest and human activities derived from
80 social media.
- 81 • We incorporate location-based social network user check-ins into a proba-
82 bilistic topic modeling technique to discover functional co-occurrence pat-
83 terns of different POI types.
- 84 • The proposed method can support functional inferences for specific type of
85 regions and thus serve as a new heuristic to enable the search for similar
86 urban places/regions, based on their POI-type distributions and corre-
87 sponding human activities, and using natural language processing and
88 machine learning techniques.

89 The remainder of this article is structured as follows. Section 2 introduces
90 background material and related work. Next, Section 3, discusses the datasets
91 used and the selection of study areas. Section 4 introduces the methods used
92 in our framework and specifically LDA. In Section 5, we present the results of
93 topic modeling to characterize, cluster, and compare functional regions. Next,
94 we discuss the broader implications of our work in Section 6 before concluding
95 and pointing out directions of future work in Section 7.

96 **2 Related Work**

97 With the increasing popularity of travel blogs, volunteered geographic informa-
98 tion (VGI), location-based social networks (LBSN), and so forth, researchers

99 have developed a variety of place-based studies that employ datasets from these
100 various sources. For instance, Adams and Janowicz (2012) presented a topic
101 modeling methodology to estimate geographic regions from unstructured, non
102 geo-referenced text on Wikipedia and travel blogs by computing a density sur-
103 face of geo-indicative topics over the Earth’s surface. The proposed framework
104 combined natural language processing techniques, geostatistics methods, and
105 data-driven bottom-up semantics. In order to evaluate the use of topic model-
106 ing techniques on the extraction of thematic characteristics of places, Adams and
107 McKenzie (2013) applied that approach on a set of travel blog entries to iden-
108 tify the themes that are most closely associated with specific places around the
109 world. Their proposed method is capable of measuring the degree to which cer-
110 tain themes are local or global, as well as analyzing thematic changes over time.
111 POI data play an important role for human activity-based land use, transporta-
112 tion, and environmental models. Jiang et al. (2015) utilized the Yahoo online
113 POI data together with publicly available aggregated employment data from
114 census at the block group level to derive fine resolution of disaggregated land
115 use estimates (i.e., employment by category) at the city block level. For the eval-
116 uation, they first used a variety of machine learning algorithms to match and
117 cluster POI types into a labeled business establishment taxonomy, and then
118 compared it with ground-truth data from commercial business data vendors.
119 The results demonstrated that their proposed method got a better goodness of
120 fit with a lower relative mean squared error for the estimated employment pop-
121 ulation across all city blocks than that from the traditional uniform-distribution
122 disaggregation approach. As for LBSN applications, Noulas et al. (2011) pro-
123 posed a method to classify the geographical areas and LBSN users based on
124 place types and the users’ check-in statistics in Foursquare venues. The experi-
125 ments were conducted in the metropolitan cities of London and New York and

126 identified similar regions and user groups in each city. However, they didn't
127 consider the temporal pattern of user activities. Later on, Yuan et al. (2012)
128 employed both POI type information and the temporal patterns of taxi pick-
129 ups/drop-offs in segmented map regions, utilized a topic modeling method based
130 on latent Dirichlet allocation and Dirichlet multinomial regression techniques,
131 and discovered various urban functional regions in the city of Beijing. The
132 extracted region clusters were annotated as nine different groups: *diplomatic*
133 *and embassy areas, education and science areas, developed residential areas,*
134 *emerging residential areas, developed commercial/entertainment areas, develop-*
135 *ing commercial/entertainment areas, regions under construction, areas of his-*
136 *toric interests, and nature & parks.* However, such rich multiple datasets that
137 complement each other in the same city and especially high-precision mobility
138 data are usually hard to fully access. One challenging issue is how to semanti-
139 cally classify and label the regions that are found given only one data source,
140 and how to find similar places and regions across different cities. Adams (2015)
141 proposed a novel observation-to-generalization place model and employed nat-
142 ural language processing techniques to derive place attributes. The proposed
143 methods can support similar-place-search functions and the case studies were
144 conducted using over 600,000 place articles on Wikipedia as a proof of con-
145 cept. Later, Adams and Janowicz (2015) presented a novel method to enrich
146 the place information on linked knowledge graphs using thematic signatures
147 that are derived from unstructured text through topic modeling. This method
148 can also be used to clean miscategorized places on the linked data cloud. In
149 another study, Hobel et al. (2015) developed a semantic region growing algo-
150 rithm based on the density of POIs on OpenStreetMap to extract places that
151 afford certain type of human activities, e.g., shopping areas. In their model, four
152 features including the number of *banks & ATM, restaurants, tourist facilities,*

153 *and subcategories of shops* were used to identify the shopping areas/settings.
154 They then compared the similarity of shopping areas/settings based on the four
155 features in two European capital cities: Vienna and London. By incorporat-
156 ing human spatio-temporal activity data from social media, Zhou and Zhang
157 (2016) extracted the spatial distribution hotspots of six types of urban func-
158 tions (i.e., *Travel & Transport, Education Resource, Shop & Service, Nightlife*
159 *Spot, Outdoor & Recreation, Food & Restaurant*) in the cities of Boston and
160 Chicago. Zhi et al. (2016) introduced a low-rank-approximation-based model
161 to detect functional regions based on 15 million social media check-in records
162 in the city of Shanghai, China. This method discovered latent spatio-temporal
163 human activity patterns and linked these with different functional regions. Re-
164 searchers are also interested in the regional differences on discovering thematic
165 characteristics of different POI types. McKenzie and Janowicz (2017) identified
166 the most and least spatially varying place types and compared their thematic
167 signatures internationally. The ongoing trend in this research direction lies in
168 data-synthesis-driven approaches to study places and vague cognitive regions as
169 well as the semantic generalizations of urban settings (Hobel et al., 2015, 2016;
170 Gao et al., 2017).

171 In summary, there is a variety of research studying places and place types
172 from human data traces, including spatio-temporal human mobility patterns
173 that can reveal the functions of regions. However, only a few studies have simul-
174 taneously considered both POI information and human activities on location-
175 based social network to derive urban functional regions. Moreover, to the best
176 of our knowledge, there is no thorough discussion of the robustness of discov-
177 ered urban and regional functional areas using different numbers of topics and
178 clusters. There has also been no attempt to develop an urban function ontology
179 based on the structure of POIs using a bottom-up approach.

180 **3 Study Area and Datasets**

181 **3.1 Study Area**

182 Urban areas – cities for short – are the highly populated places on the planet
183 and include metropolitan regions, urban districts, towns, and suburbs. In or-
184 der to explore the thematic characteristics and semantic clusters of urban areas
185 in connection with urban functions, the ten most populated U.S. cities based
186 on the 2015 population census: *New York, Los Angeles, Chicago, Houston,*
187 *Philadelphia, Phoenix, San Antonio, San Diego, Dallas, and San Jose* and their
188 surrounding metropolitan regions were selected as our study areas. The carto-
189 graphic boundaries of those ten metropolitan areas are downloaded from the
190 U.S. Census Bureau’s TIGER geographic database¹.

191 **3.2 Points of Interest Dataset**

192 People usually go to different POIs for different kinds of activities, e.g., study-
193 ing, working, dining, shopping, and relaxing. We assume that the spatial dis-
194 tributions and interactions of different types of POIs reflect particular urban
195 functions. Location-based social networks such as Foursquare have created
196 traces of social interactions based on the physical location of users. In these
197 LBSN systems, users can check-in to a venue (i.e., a POI), rate it, and share
198 their comments or tips. As shown in Figure 1, we first randomly generated
199 200 points as search locations in each urban area and then identified the sur-
200 rounding Foursquare venues with their attribute information including *name,*
201 *location coordinates, place category, number of check-ins, number of checked*
202 *users, number of tips, and the rating score* in each search locations. Note that
203 because of the Foursquare developer API limits, we only retrieved at most 50
204 nearby venues given a random search point. The POI data were collected in De-

¹https://www.census.gov/geo/maps-data/data/cbf/cbf_msa.html

Table 1: The 95% inverse CDF distance thresholds for all the ten urban areas

City Name	95% Inverse-CDF Distance Threshold (km)
Chicago	7.389
Dallas	8.994
Houston	8.647
Los Angeles	4.474
New York	7.123
Philadelphia	8.894
Phoenix	7.969
San Antonio	7.475
San Diego	7.837
San Jose	4.822
Average	7.363

205 cember 2016 and the attribute information for all venues is a historic snapshot
 206 at that time. There is a total of 480 different POI types in our data. Figure 2
 207 shows the empirical cumulative density function (CDF) of the distance distri-
 208 bution between each Foursquare venue and the corresponding search location.
 209 Steeper curves (with larger slope values) before reaching the relatively steady
 210 state (about 95% cumulative probability) show that more POIs are closer to
 211 the search locations given the same number of Foursquare venues. In order to
 212 generate most 'nearby' POIs around each search location, we further spatially
 213 filtered out those venues outside the 95% inverse CDF distance threshold; i.e.,
 214 we only selected those venues within a relatively small search distance. The
 215 distance thresholds differ among cities as shown in Table 1.

216 4 Methods

217 4.1 Popularity-based Probabilistic Topic Model

218 Probabilistic topic models have been widely used to discover latent thematic
 219 characteristics and their structure when analyzing large sources of textual doc-

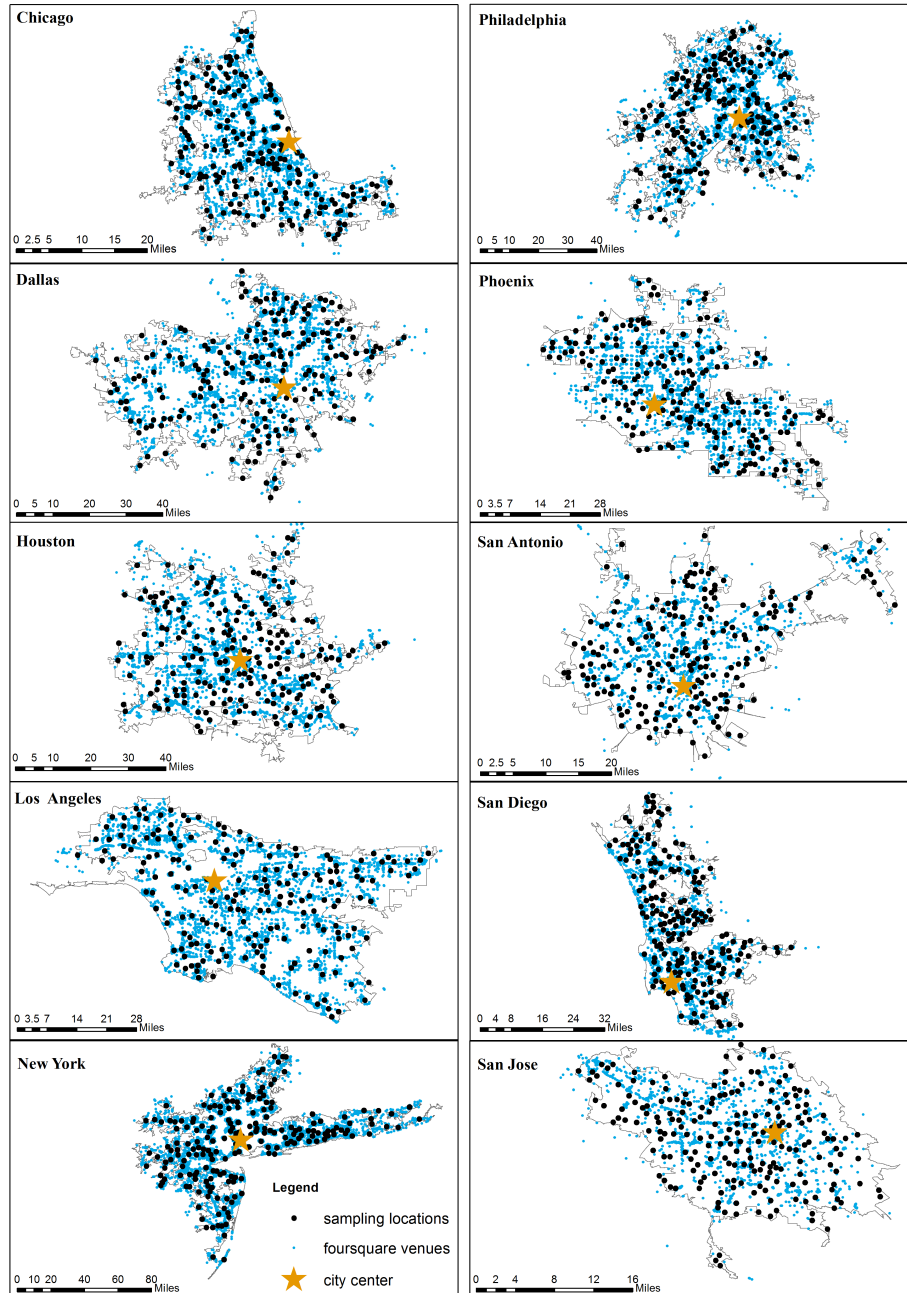


Figure 1: The spatial distributions of sampling locations and the collected Foursquare venues (POIs) in ten urban areas.

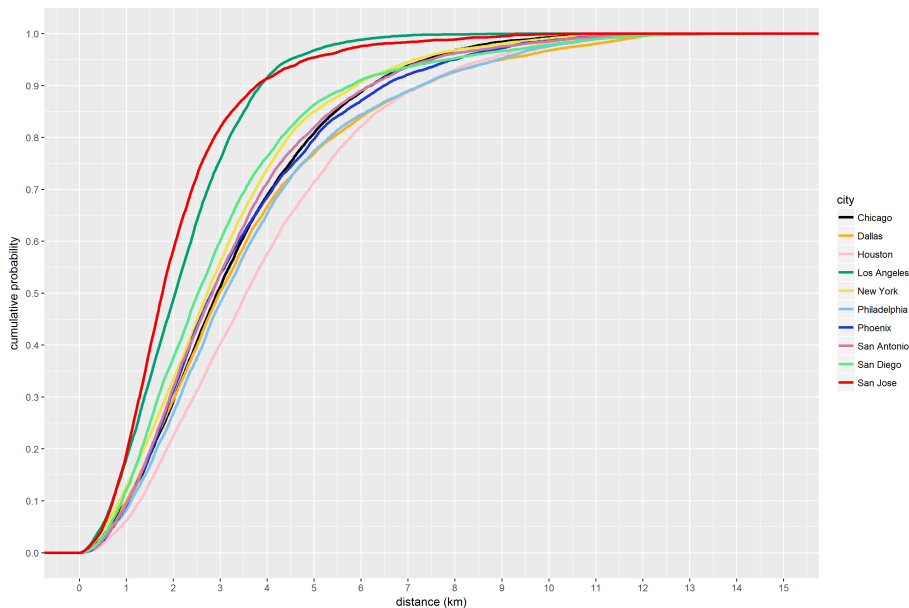


Figure 2: The cumulative density function of the distance distribution between each Foursquare venue and the corresponding search location.

220 uments (Blei et al., 2003; Steyvers and Griffiths, 2007; Blei, 2012). The *latent*
 221 *Dirichlet allocation (LDA)* is among the most popular topic modeling methods.
 222 LDA is an unsupervised generative probabilistic model that takes a bag-of-
 223 words approach (which implies that the order of words in the document does
 224 not matter) to constructing topics. The key idea of LDA is that documents can
 225 be represented as a joint probability distribution over latent topics and each
 226 topic is characterized by a distribution over words (Blei et al., 2003). Assume
 227 that there are total K number of topics associated with N words in the doc-
 228 ument corpus D , and α and η represent the prior parameters for the Dirichlet
 229 document-topic and topic-word distribution respectively. The mathematical re-
 230 lationship between the latent variables and the observed variables is described
 231 below:

$$\begin{aligned}
& p(\beta_{1:K}, \theta_{1:D}, Z_{1:D}, W_{1:D}) \\
&= \prod_{i=1}^K p(\beta_i) \prod_{i=1}^D p(\theta_d) \left(\prod_{n=1}^N p(Z_{(d,n)} | \theta_d) p(W_{(d,n)} | \beta_{1:K}, Z_{(d,n)}) \right) \quad (1)
\end{aligned}$$

232 As shown in Figure 3, the generative process can be described as follows:

233 I. Let β_k denote a probabilistic distribution over the word vocabulary for a
 234 given topic k , and draw $\beta_k \sim \text{Dir}(\eta)$;

235 II. Let θ_d represent the topic proportions for the d_{th} document, and draw θ_d
 236 $\sim \text{Dir}(\alpha)$;

237 III. Let $Z_{(d,n)}$ denote the topic assignment for the n_{th} word in document d and
 238 $W_{(d,n)}$ represent the n_{th} word in document d from a fixed vocabulary,
 239 and draw the multinomial distributions: $Z_{(d,n)} \sim \text{Multinomial}(\theta_d)$ and
 240 $W_{(d,n)} \sim \text{Multinomial}(\beta_{z_{(d,n)}})$.

241 In order to compute the conditional distribution of the topic structure given
 242 the word observations in documents, the expectation–maximization (EM) algo-
 243 rithm and *Gibbs sampling* are most the commonly used methods. After finishing
 244 the computation, two matrices θ and β associated with topic proportions and
 245 assignments are generated. More detailed notations, calculations, and explana-
 246 tions can be found in Blei et al. (2003).

247 In analogy with LDA’s use of textual materials, we take the type (e.g.,
 248 school, park, restaurant) of each POI as a word, the search region that contains
 249 those POIs as a document, and an urban function or a land use as a topic that
 250 represents thematic characteristics and the semantics of places. By running the
 251 LDA topic modeling technique, we can find the posterior probabilistic distri-
 252 bution of each POI type in a certain type of region conditioned on the search

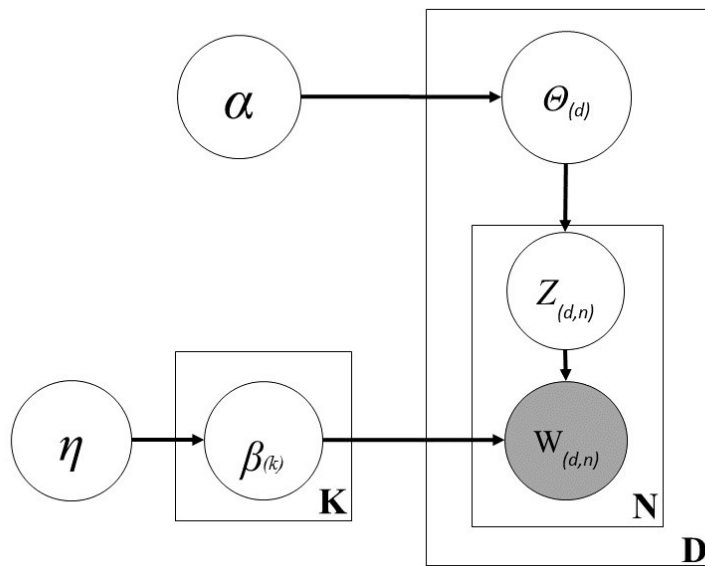


Figure 3: The graphical representation for latent Dirichlet allocation. The topic-related random variables in the generative process are the unshaded nodes while the observed words in documents are represented as a shaded node. The rectangles are "plate" notation that denotes replication.

253 region’s topic assignments. The LDA model running on POI types generates
 254 summaries of thematic place topics (e.g., beach promenades, art zones, shopping
 255 areas) with a discrete probability distribution over POI types for each topic, and
 256 infers per-search-region probability distributions over topics. For example, one
 257 would assume that a *beach promenade* topic should contain venues such as *beach*,
 258 *seafood restaurants*, and *surfing spots*; while a *shopping area* would more likely
 259 contain *clothing stores*, *cosmetics shops*, and *shoe stores*.

260 Another important concept in our method is the *popularity* of a POI as
 261 captured by its LBSN user check-in behaviors. For example, the neighbor-
 262 hood of a football stadium usually has only one instance of *stadium* surrounded
 263 by dozens of *sports bars*, *restaurants*, and *parking lots*. However, a stadium
 264 usually attracts thousands of visitors and is the dominant feature of its neigh-
 265 borhood. Thus this particular POI type makes said neighborhood distinct from
 266 other neighborhoods (e.g., nightlife zone), which also contain *cocktail bars* and
 267 *restaurants*. We need to address such a human activity effect during the gener-
 268 ation of the document-word frequency matrix. More specifically, we will rescale
 269 the POI type occurrences according to their associated POI instance check-in
 270 counts. The rescaling process can be represented as follows:

$$Freq_{(d,t)} = \sum_i \text{Log}(V_{(d,t,i)}) \quad (2)$$

271 where $Freq_{(d,t)}$ represents the rescaled occurrence for a POI type t given a
 272 search region d ; and $V_{(d,t,i)}$ is the number of unique users who have contributed
 273 their check-ins for a venue i that belongs to the same POI type t in the same
 274 search region d .

275 We then test whether such an unsupervised popularity-based LDA topic
 276 modeling technique can support the discovery of characteristic semantic regions
 277 across different U.S. cities with a similar structure of POI type mix distribution.

278 Finding the appropriate number for latent topics is important but difficult
279 given a dataset using the LDA topic model. Several metrics and methods have
280 been developed to address this issue. Griffiths and Steyvers (2004) used the
281 Gibbs sampling algorithm to obtain samples from the posterior distribution
282 over topic assignments Z at different choices of the total K number of topics,
283 and then calculated a log-likelihood $P(W|Z, k)$. The value of K at which the
284 log-likelihood gets the maximum and stabilizes after hundreds of iterations will
285 be taken as the right number of topics for a specific document corpus. With
286 the consideration of one issue that sometimes words have too many overlaps
287 across those generated latent topics, Cao et al. (2009) proposed a density-based
288 method for adaptive LDA model selection. The key idea of this algorithm is
289 to maximize the intra-topic similarity while minimize the inter-topic similar-
290 ity. They calculated the average cosine distance between pairwise topics with
291 their word assignments and then used a heuristic to find the most stable topic
292 structure given the best K value based on the topic density measure. Arun
293 et al. (2010) viewed the LDA topic model as a matrix factorization mechanism
294 and applied the symmetric Kullback–Leibler (KL) divergence (Kullback and
295 Leibler, 1951) on the distributions generated from topic-word and document-
296 topic matrices for finding the right number of topics. The best K value at which
297 the symmetric KL-divergence is the minimum would derive the most discrim-
298 inative topics and their distributions become orthogonal. In several empirical
299 geographic information studies (Adams and Janowicz, 2015; McKenzie et al.,
300 2015; Gao et al., 2017), different K numbers (e.g., 60, 100, 300) of topics have
301 been deployed to investigate place characteristics. An optimal value of K may
302 vary between different datasets and has influence on the thematic similarity of
303 POI types (McKenzie and Janowicz, 2017).

304 4.2 Functional Region Aggregation

305 After deriving those latent thematic topics by running the proposed popularity-
306 based LDA model, each region can be represented as a vector of the K -dimensional
307 POI type topics. Those regions that are semantically similar in the topic space
308 might contribute to the same urban function and can be aggregated into the
309 same cluster as a functional region. Two clustering approaches are applied
310 in this work: *K-means clustering* and *the Delaunay triangulation spatial con-*
311 *straints clustering methods*.

312 *K-means clustering* only takes the thematic characteristics of multivariate
313 topic distributions of places into consideration without any spatial constraints
314 (MacQueen et al., 1967). It is an unsupervised clustering approach in which the
315 number K needs to be predefined. The *silhouette* criterion has been widely used
316 for determining an appropriate value of K (Rousseeuw, 1987). The *silhouette*
317 value $s(i)$ quantifies how well an object i is appropriately clustered. The range of
318 *silhouette* value is between -1 and 1. A high $s(i)$ value (close to 1) indicates that
319 an object is appropriately clustered and is very dissimilar from other clusters.
320 In the region clustering process for our POI datasets, we tried different K values
321 ranging from 1 to 30 and identified the maximum average silhouette value across
322 all clusters and chose that K as the optimal K-means clustering parameter for
323 reporting the corresponding results.

324 *Delaunay triangulation spatial constraints clustering* has been introduced by
325 Assunção et al. (2006). This approach consists of three steps: (1) building a
326 connectivity graph to capture the adjacency relations between points based on
327 Delaunay triangulation spatial constraints; (2) creating a minimum spanning
328 tree (MST) (Gower and Ross, 1969) from the neighboring connectivity graph
329 with minimizing the sum of the dissimilarities over all the edges of the tree; (3)
330 partitioning the derived MST into different subtrees as spatial clusters using

331 a hierarchical division strategy to minimize the intra-cluster square deviations.
332 More implementation details about this clustering algorithm can be found at
333 Assunção et al. (2006).

334 5 Analysis and Results

335 5.1 Topic Modeling Results

336 As proposed in Section 4, before running the LDA model, we first incorporated
337 the number of visitors for each venue as a popularity score in the rescaling
338 process to generate a new document-word matrix (i.e., a search region-POI
339 type occurrence matrix) across all search regions in the ten urban areas. Next,
340 we evaluated the performance of different choices of K as the total number of
341 topics for the LDA topic model using three introduced measures. As shown in
342 Figure 4, by choosing the value of K from 5 to 200 and then running LDA topic
343 models on our POI data, we derived different topic assignment results. The
344 measure proposed by Griffiths and Steyvers (2004) aims to maximize the log-
345 likelihood of word-topic probability in the documents, while other two measures
346 (Cao et al., 2009; Arun et al., 2010) aim to minimize the proposed criteria. In
347 the ideal case, one would expect those three measures converge at the same
348 value of K . Unfortunately, in empirical studies, they do not necessarily present
349 such perfect convergence patterns. In our parameter tuning experiments, the
350 optimal K value for the “CaoJuan2009” and “Arun2010” measures is in the
351 range of 90 – 160, while the “Griffiths2004” metric gets relatively stable when
352 K reaches 130 topics.

353 Therefore, we set $K = 130$ as the total number of topics and ran 2000
354 iterations of the Gibbs sampling process to derive the posterior probabilistic
355 distribution over topic assignments. In Figure 5, we show nine of those inter-

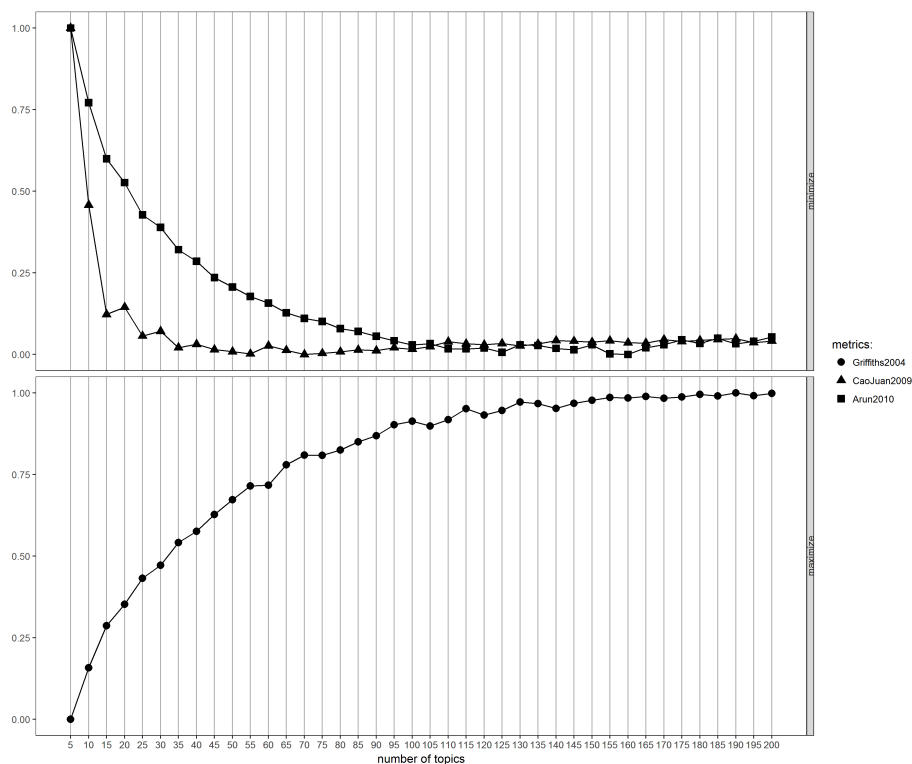


Figure 4: Find an appropriate K value for the total number of topics using three metrics.

356 esting topics related to urban functions. Note that the probability assignments
 357 for those POI types are weighted and ranked by their *term frequency–inverse*
 358 *document frequency* (i.e., POI type frequency–inverse region frequency) so that
 359 each topic can display more distinctive and meaningful POI types that are di-
 360 rectly proportional to the frequency in documents while inversely proportional
 361 to the region frequency at which a POI type occurs in the whole corpus. For
 362 instance, *coffee shop* has a very high frequency but also widely exists in most
 363 of the regions in our POI data so it plays a less important role than other cat-
 364 egories (e.g., *theme park*) in distinguishing the function of a region. Some of
 365 those meaningful topics are illustrated as follows.

366 *Topic 67* is a shopping-plaza topic that consists of various frequent occurring
367 POI types including *shopping mall*, *accessories store*, *chocolate shop*, and
368 a few restaurants. It is one of the most prominent topics across all cities.
369 *Topic 109* is a beach-related topic that consists of *beach*, *surf spot*, *island*,
370 *beach bar*, *pier*, and *so on*. In terms of spatial distribution, one would expect
371 such topics should be located in coastal or lake-side cities only. Both *Topic*
372 *21* and *Topic 25* contain *history museum* and *art museum*, but *Topic 21* is
373 more related to college/university regions since it also contains *pool*, *college*
374 *rec center*, *tourist information center*, and several other educational facilities;
375 while *Topic 25* is more likely an art museum district since it also consists of *art*
376 *gallery*, *antique shop*, *scenic lookout*, and *so on*. *Topics 6*, *117*, and *119* relate
377 to outdoor sports and leisure activity places, such as *national park*, *ski lodge*,
378 *gym*, *golf course*, *tennis court*, and various restaurants and studios. *Topic 36*
379 and *Topic 74* describe mix distribution patterns of *bar*, *restaurant*, *government*
380 *building*, *residential apartment*, and *business service*, which may suggest central-
381 city areas.

382 In order to explore the variability of the above discovered nine topics while
383 changing the total number of topics, we further investigate whether we can find
384 exactly matching or most similar topics with different values for K . Two eval-
385 uation criteria, namely *Cosine Similarity* and *Jaccard Index*, were applied for
386 this purpose (Han et al., 2011). Assume that each topic vector is a sequence of
387 probabilistic values between 0 and 1 for all the 480 POI categories. Consider-
388 ing each pair of one target topic (e.g., *Topic 6* when $K=130$) and another one
389 comparing topic (e.g., *Topic 1* when $K=10$), the cosine similarity measures the
390 cosine of the angle between two non-zero vectors defined using an inner product.
391 It is well suited to evaluate sparse vectors such as document-word matrices and
392 the topic-POI matrices in our experiments. Unlike the cosine similarity that is

393 frequently used for numeric vectors, the Jaccard index is a popular similarity
394 measure for binary and categorical data, which is defined as the cardinality of
395 the intersection divided by the cardinality of the union of two sets. We use
396 the Jaccard index to quantify the topic structure similarity for their top-fifteen
397 probabilistically ranked POI types. The larger the value, the more similar two
398 topics are, where 1 equals a perfect match while 0 indicates no overlapping top-
399 terms (i.e., POI types) in the comparison of two topics. The comparison batch
400 processing was conducted from $K=10$ to $K=150$ with a step of 10. During each
401 run with a given K , the maximum similarity values to each of the nine topics
402 were computed. As shown in Figure 6, the maximum cosine similarities for
403 *Topics 74 and 117* reach almost 1 and remained stable when the total number
404 of topics exceeded 30. As for *Topics 6, 36, 67, and 109*, we can also identify
405 most semantically similar (≥ 0.9) topics to them with K value equals to 150,
406 120, 150 or 140 respectively. This indicates the stability of identifying those
407 prominent urban functional topics related to frequently co-occurrent physical
408 facilities and services, a variety of bars and restaurants, and leisure activity
409 places. However, we cannot find very similar (≥ 0.8) topics to *Topics 21, 25,*
410 *and 119* when choosing different K values, which implies that these topics may
411 be more characteristic of a specific K value. In a similar manner, we analyze
412 the topic structure similarity using the Jaccard index. As shown in Figure 7,
413 those low similarity values illustrate the large composition variability existed in
414 the top-fifteen probabilistically ranked POI types for all discovered topics with
415 different K values. Their implications will be discussed in the section 6.

416 In short, rather than a traditional top-down approach for describing ur-
417 ban functions based on familiar compositions of POI types, we demonstrated a
418 bottom-up statistical topic-learning approach for finding underlying co-occurrence
419 relationships among different types of POIs based on data on human activities

420 extracted from location-based social networks.

421 **5.2 Searching for Similar Places**

422 Searching for similar places is an important task in geographic information re-
423 trieval and also valuable in many applications, such as tourism, real estate, and
424 immigration. People may consider many factors such as job market, affordabil-
425 ity, natural environment, and quality of life. When people consider moving into
426 new cities, they may also want to know how they will like these new places and
427 whether they can find similar neighborhoods to the ones they will be leaving.
428 Such places typically contain a mix of types of POIs that people would like
429 to visit. Fortunately, such information can be retrieved from popular location-
430 based social network platforms that have been used as a lens of social sensing to
431 capture human-place interactions. In the following, we illustrate this idea with
432 two scenarios:

433 (1) Search for similar regions given a dominant theme. We selected the city
434 of Denver as our target city, which was ranked as the best metropolitan area
435 to live in the U.S. according to a survey² from the U.S. News in 2016. It has
436 a variety of local attractions and support many activities. Here we aim to find
437 regions that are similar to those represented by *Topic 25*, which is related to
438 art districts. We collected the Foursquare POIs and user check-in data for Den-
439 ver by randomly sampling search locations and then searching for 50 nearby
440 POIs given each sample location. Based on the aforementioned data process-
441 ing procedures and the LDA topic model by incorporating the popularity score
442 based on unique Foursquare check-in users, we can infer the probabilistic com-
443 bination of different topics for a search region given its POI type co-occurrence
444 pattern. As shown in Figure 8, within this search neighborhood, we discov-
445 ered a high probabilistic topic distribution for *Topic 25*, which consists of a

²<http://realestate.usnews.com/places/rankings/best-places-to-live>

Topic 6		Topic 21		Topic 25	
Category	Prob.	Category	Prob.	Category	Prob.
gas station	0.309651	pool	0.122166	museum	0.065636
italian restaurant	0.028574	history museum	0.047285	art museum	0.047585
flower shop	0.013235	historic site	0.043474	art gallery	0.046691
national park	0.002077	college basketball court	0.015107	american restaurant	0.038677
ski lodge	0.000968	concert hall	0.012485	record shop	0.025063
jewish restaurant	0.000899	art museum	0.012412	antique shop	0.024183
auditorium	0.000847	college rec center	0.010488	building	0.002540
southern food	0.000226	park	0.007814	cycle studio	0.001103
ice cream shop	0.000123	sculpture garden	0.007216	health food store	0.000462
farmers market	0.000109	outdoor sculpture	0.005112	history museum	0.000430
club house	0.000105	college soccer field	0.004028	soup place	0.000350
bbq joint	0.000100	college cafeteria	0.003933	concert hall	0.000281
pizza place	0.000100	tourist information center	0.003731	scenic lookout	0.000264
winery	0.000072	molecular gastronomy	0.003120	animal shelter	0.000179
grocery store	0.000059	stables	0.002947	burger joint	0.000179

Topic 36		Topic 67		Topic 74	
Category	Prob.	Category	Prob.	Category	Prob.
yoga studio	0.105001	shopping mall	0.207709	bar	0.511221
science museum	0.065819	accessories store	0.056738	board shop	0.000046
boutique	0.029987	chocolate shop	0.013896	asian restaurant	0.000046
gay bar	0.015371	shoe store	0.000288	brewery	0.000036
sculpture garden	0.012688	breakfast spot	0.000282	ice cream shop	0.000030
government building	0.008197	gaming cafe	0.000196	parking	0.000025
israeli restaurant	0.004401	optical shop	0.000180	buffet	0.000025
apartment / condo	0.003005	post office	0.000114	business service	0.000019
pakistani restaurant	0.002829	bistro	0.000105	apartment / condo	0.000016
street food gathering	0.001212	dumpling restaurant	0.000096	fried chicken joint	0.000012
track stadium	0.000872	korean restaurant	0.000090	resort	0.000007
college baseball diamond	0.000602	german restaurant	0.000080	gourmet shop	0.000007
mexican restaurant	0.000542	herbs & spices store	0.000079	lighthouse	0.000006
gym / fitness center	0.000526	airport terminal	0.000078	indian restaurant	0.000006
cheese shop	0.000481	outlet store	0.000076	train	0.000006

Topic 109		Topic 117		Topic 119	
Category	Prob.	Category	Prob.	Category	Prob.
beach	0.285864	italian restaurant	0.082055	french restaurant	0.090092
surf spot	0.028952	fast food restaurant	0.000131	cocktail bar	0.072534
italian restaurant	0.015458	gym	0.000064	lounge	0.035774
island	0.005920	golf course	0.000056	tennis court	0.005389
beach bar	0.004078	sushi restaurant	0.000044	whisky bar	0.003636
board shop	0.003793	salon / barbershop	0.000043	american restaurant	0.000697
bridge	0.001484	boutique	0.000034	dry cleaner	0.000168
indie theater	0.001235	café	0.000030	pizza place	0.000118
pier	0.001187	szechuan restaurant	0.000030	café	0.000117
outdoor sculpture	0.001121	japanese restaurant	0.000030	art museum	0.000110
sri lankan restaurant	0.001040	paella restaurant	0.000027	bakery	0.000106
bistro	0.000891	men's store	0.000023	jazz club	0.000082
nature preserve	0.000851	caribbean restaurant	0.000017	chinese restaurant	0.000074
arepa restaurant	0.000751	deli / bodega	0.000016	neighborhood	0.000068
neighborhood	0.000726	massage studio	0.000014	cycle studio	0.000040

Figure 5: Nine interesting topics with their top-15 ranked POI types related to urban functions.

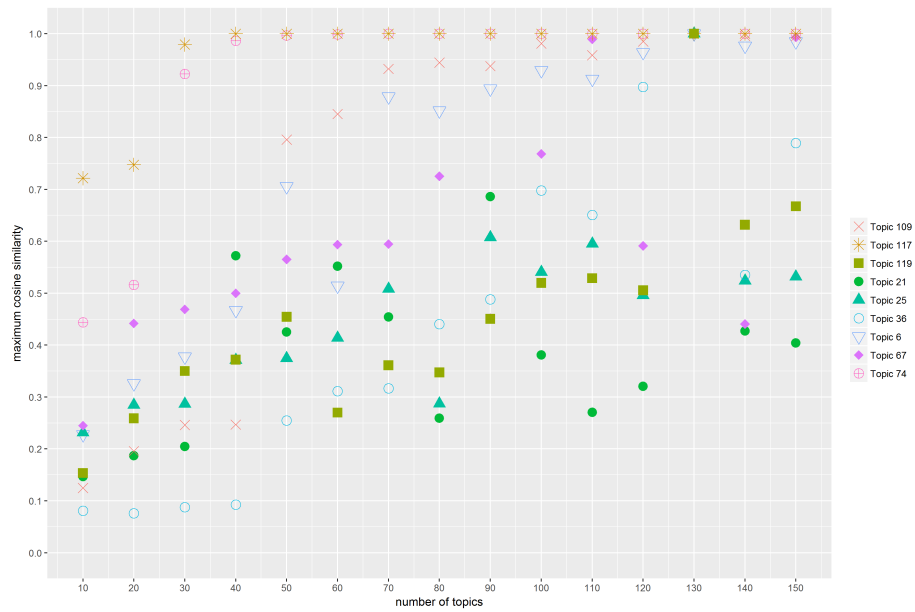


Figure 6: Maximum cosine similarity between the selected nine topics and the resulting topic models by choosing different total number of topics.

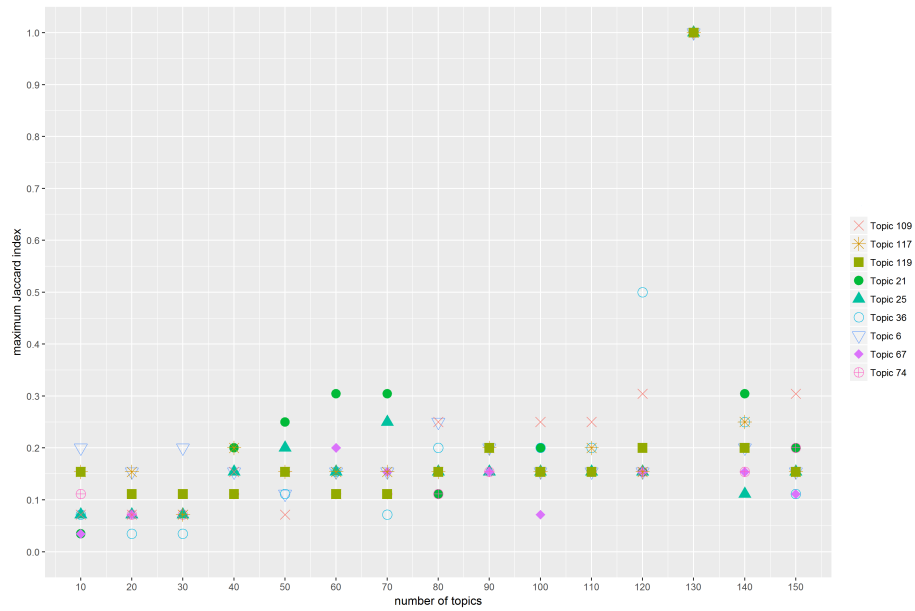


Figure 7: Maximum Jaccard similarity between the top-15 POI types of the selected nine topics and that from the resulting topic models by choosing different total number of topics.

446 variety of prominent POI types such as *art museum, art gallery, history mu-*
 447 *seum, concert hall and American restaurant.* Such a place may serve multiple
 448 functions. The second largest probabilistic topic in this search neighborhood is
 449 *Topic 121* that contains a large percentage composition of *brewery places.* By
 450 looking up other geographic background information and Web pages, we realize
 451 that local people actually also identify this region near the “*Santa Fe Dr.*” as
 452 an “Art District” in Denver, which attracts many local residents, artists and
 453 tourists³. This example illustrates the inference capability of our method to
 454 identify similar neighborhoods given certain thematic characteristics.

455 (2) Search for similar places considering all themes. After running the LDA
 456 topic model, each place can be represented as a multinomial distribution of
 457 K -dimensional POI type topics, denoted as a probability vector $[p_1, p_2, \dots, p_k]$,
 458 where all the probability values sum to one. Thus we can apply a variety of
 459 probabilistic distance or similarity measures (e.g., Hellinger distance, cosine
 460 similarity, and Jensen-Shannon divergence (JSD) (Lin, 1991) to quantify the
 461 pairwise similarity among all search regions in our POI data with regard to their
 462 POI type mix distributions. JSD is a symmetric distance measure derived from
 463 the Kullback-Leibler divergence (KLD) asymmetric distance measure between
 464 two probability distributions P and Q (Kullback and Leibler, 1951).

$$KLD(P|Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)} \quad (3)$$

$$M = \frac{P + Q}{2} \quad (4)$$

$$JSD(P|Q) = JSD(Q|P) = \frac{KLD(P|M)}{2} + \frac{KLD(Q|M)}{2} \quad (5)$$

467 The JSD is bounded by 0 and 1 if using the base 2 logarithm for the two

³<http://www.denver.org/about-denver/neighborhood-guides/artdistrict-on-santa-fe>

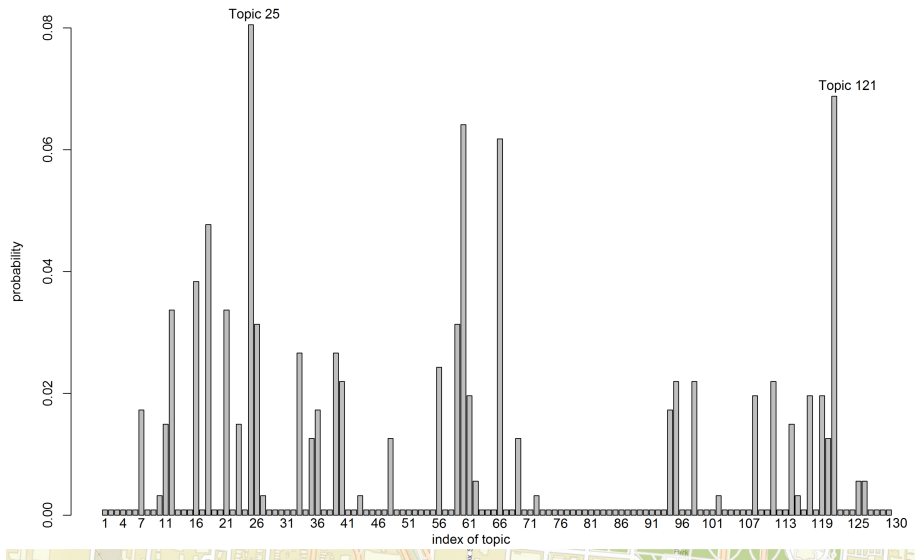


Figure 8: The topic probability distribution and the spatial distribution of Foursquare POIs around the Denver art district and museums.

468 KLD relative entropy calculation. And thus we can define a JSD-similarity
469 metric ($S_{(JSD)}$) as follows:

$$S_{(JSD)} = 1 - JSD(Q|P) \quad (6)$$

470 where the base 2 logarithm is used in the KLD and JSD calculations.

471 Therefore, according to the proposed similarity measure $S_{(JSD)}$, we can
472 analyze the pairwise similarity among our randomly selected search places that
473 contains those POIs. Figure 9 shows a JSD-similarity matrix for 200 randomly
474 selected places in Los Angeles, derived from one part of our whole dataset in
475 Section 3, and each place is represented as a vector of 130-dimensional thematic
476 topics. The similarity score in each grid is between 0 and 1. The higher the
477 value, the more similar the two places are with regard to their topic distributions.
478 The values in the diagonal are all 1. By visualizing this similarity matrix,
479 one can easily identify two anomalous red stripes (i.e., the labeled R_{th} row
480 and the C_{th} column) with relatively low similarity values across the grid cells.
481 Interestingly, as shown in Figure 10, further investigation reveals that this place
482 was sampled at a location inside the *Disneyland Resort* in the Los Angeles
483 metropolitan area, which is very different from all other randomly sampled
484 places and the dominant topic (*Topic 56*) has unusual POI types such as *theme*
485 *park*, *theme park ride/attraction*, and *gift shop*. The frequent co-occurrence of
486 those distinctive types of POIs in this region causes the very low similarity to all
487 other places. Thus, given any place, we can find the most similar or dissimilar
488 places in another geographic region based on this similarity matrix.

489 5.3 Discovering Functional Regions

490 Another goal for us is to discover urban functional regions where semantically
491 similar places group together. As described in Section 4, two clustering meth-

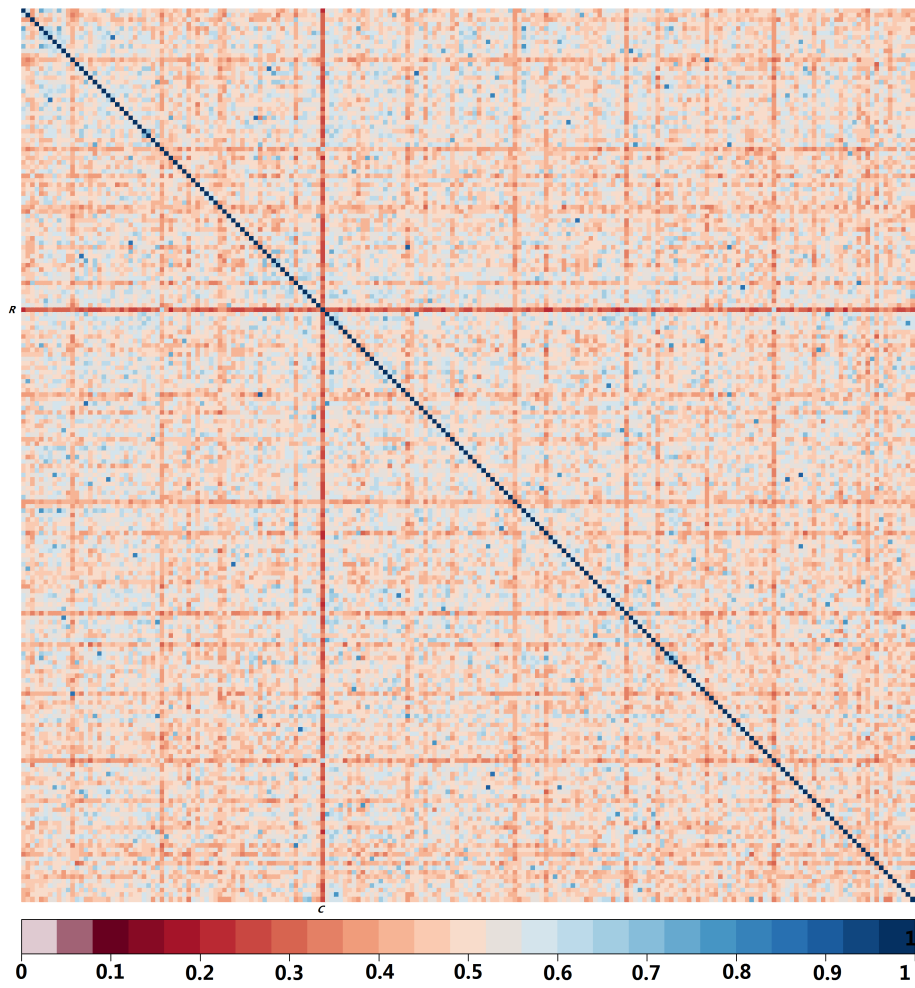


Figure 9: The JSD-similarity matrix for 200 randomly selected places in Los Angeles, where each place consists of 130 dimensional thematic topics.

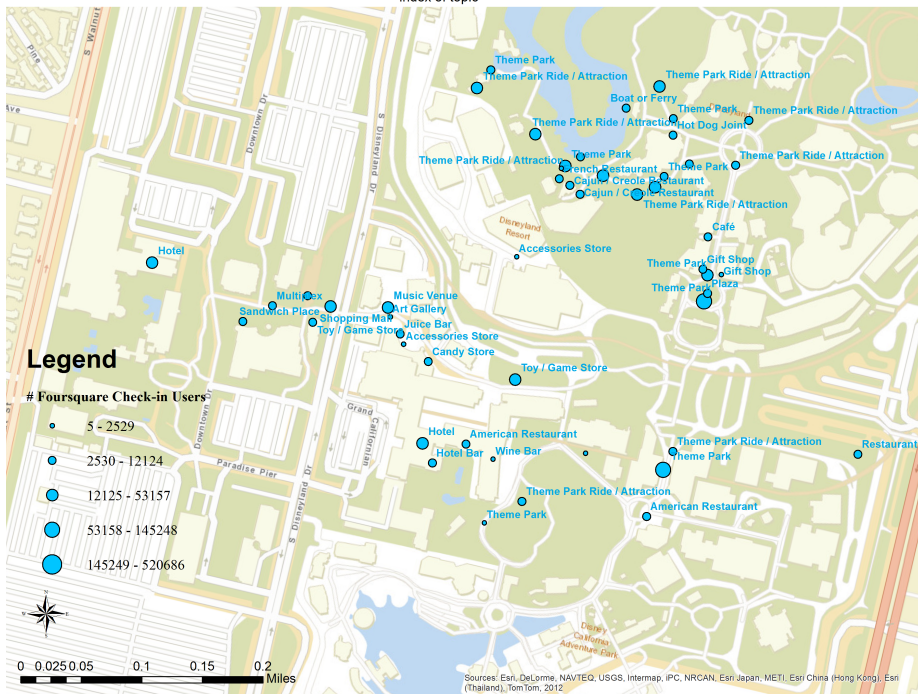
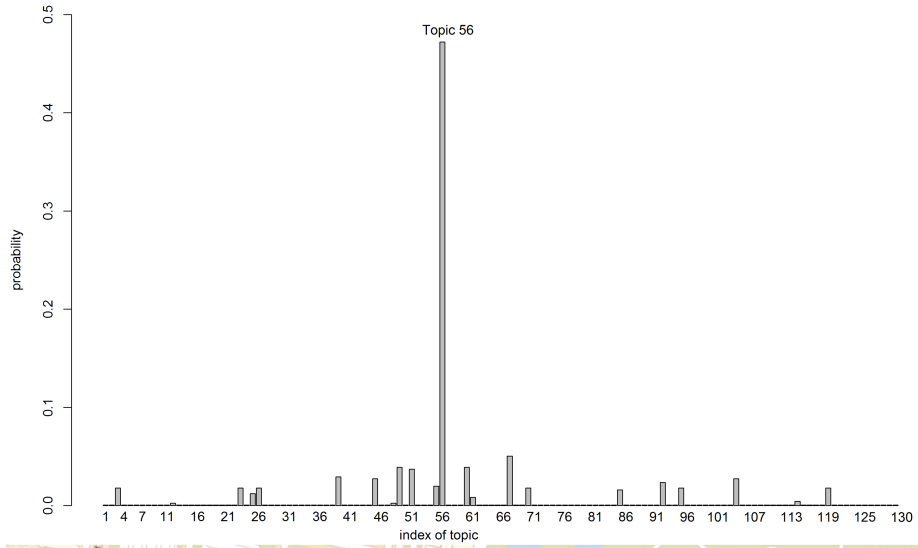


Figure 10: The topic probability distribution and the spatial distribution of Foursquare POIs in the Los Angeles Disneyland Resort.

ods are applied for aggregating similar places into functional regions. Figure 11 shows the K-means clustering result for 200 randomly sampled places in Los Angeles. The *silhouette* value for determining the optimal number of clusters in Los Angeles is 15, and thus we group those places into fifteen clusters. The circles with the same color on the map belong to the same cluster within which POI structures are more similar in their topic space of types. The numeric label on the top of each circle displays the top ranked topic that has the largest probability over the 130-dimensional thematic topic vector in this location. Note that purely K-means clustering doesn't consider any spatial constraints, and thus distant places sharing similar functions or thematic characteristics can also be grouped into the same cluster. For example, several places are dominated by food-related *Topic 30*, which contains frequent distributions of various restaurants such as *Korean restaurant*, *Mongolian restaurant*, *Portuguese restaurant*, and *Polish restaurant* are grouped into *Cluster 8*, although those places are spatially separated. However, if we take spatial constraints into considerations, only places that are semantically similar in the topic space and also located near each other can be aggregated into the same cluster. Figure 12 shows the Delaunay-triangulation-spatial-constraints clustering result for those 200 sampled places in Los Angeles. Note that we keep the same color scheme for the visualization of two clustering results, but those clusters in the same color from two maps are not identical. In the West Coast area, we can see that several places are dominated by the beach *Topic 109* and related leisure activity categories are spatially clustered together into *Cluster 7*. Although *Cluster 7* and *Cluster 3* are spatially close and share beach characteristic, *Cluster 3* tends to have another dominant POI type (*shopping plaza*) in this region, which distinguishes it from *Cluster 7* and *Cluster 10*.

By analyzing the spatial distribution of similar places and clusters, researchers

519 can have a better understanding of how certain types of POIs co-locate in order
520 to serve different urban functions from the bottom-up perspective. In addition,
521 urban planners or managers are able to further investigate the needs for comple-
522 mentary physical facilities and services related to the thematic characteristics
523 derived from the human activities on the location-based social networks. This
524 keeps in line with the human-centered and community-oriented perspectives in
525 traditional top-down urban planning and design. Furthermore, we create the
526 bounded functional regions as *convex polygons* derived from those points in the
527 same cluster (Figure 12). This can help geographic information service providers
528 develop topic-related POI search services within certain functional regions. Be-
529 cause the POI type assignments for all topics are semantically interpretable, we
530 can also select multiple dimensions of topics in geographic information queries
531 such as the *beach + shopping plaza* topics. *Cluster 3* (in Figure 12) would be a
532 good candidate since it has a mix of the dominant beach topic and the shopping
533 topic.

534 In addition, in order to test the robustness of discovered urban functional
535 areas with different probabilistic topics, we perform a series of clustering result
536 comparisons by choosing different numbers of topics ranging from 10 to 150. We
537 use their corresponding probabilistic POI type compositions as clustering fea-
538 tures and run both K-means clustering and the Delaunay-triangulation-spatial-
539 constraints clustering. Two popular metrics for comparing clustering results
540 are applied in our tests: the *Rand index* (Rand, 1971) and *normalized mutual*
541 *information (NMI)* (Strehl and Ghosh, 2002). The Rand index measures the
542 percentage of decision agreements between two clustering results X and Y. It
543 contains two types of decision agreements: (1) the number of pairs of search
544 locations within the same clustering region in X that are also in the same clus-
545 tering region in Y; (2) the number of pairs of search locations that are in different

546 clustering regions in X and also in different clustering regions in Y. The NMI
547 quantifies the mutual dependence/similarity between two clustering results us-
548 ing information theory. The detailed formula descriptions can be found in the
549 original article (Strehl and Ghosh, 2002). For both the Rand index and NMI,
550 their value range is between 0 and 1 and larger values indicate higher similarity
551 between two clustering results. Figure 13 shows the K-means clustering compar-
552 isons using the Rand index and the NMI metric between the target scenario
553 (130 topics and 15 clusters) and other scenarios with different number of topics
554 but with the same total number of clusters. Figure 14 shows the comparison
555 results for the Delaunay-triangulation-spatial-constraints clustering in a similar
556 manner. We find that the Rand index keeps a high value around 0.85 for both
557 clustering methods, which indicates a large percentage of agreements on the
558 clustering membership of those search locations and derived functional areas.
559 But the NMI values show a fluttering pattern that indicates the existence of
560 cluster membership variability. Furthermore, as for both evaluation metrics,
561 the Delaunay-triangulation-spatial-constraints clustering has a higher similar-
562 ity value in most comparison scenarios and seems to be more stable than the
563 K-means clustering results. It may imply that the spatial constraints play a role
564 in deriving the functional regions.

565 **6 Broader Implications and Discussions**

566 Based on the analysis results in this research, we show that several latent topics
567 of POI categories are spatially and semantically related to certain urban func-
568 tions. For example, the college/university topic that consists of *buildings*, *pool*,
569 *sports fields*, and *apartments*, is also co-located with several *restaurant* and *bar*
570 like topics; the *shopping plaza* topic is often also co-located with the *parking*
571 and *resort* like topics. This reveals the underlying relations of how POI cate-

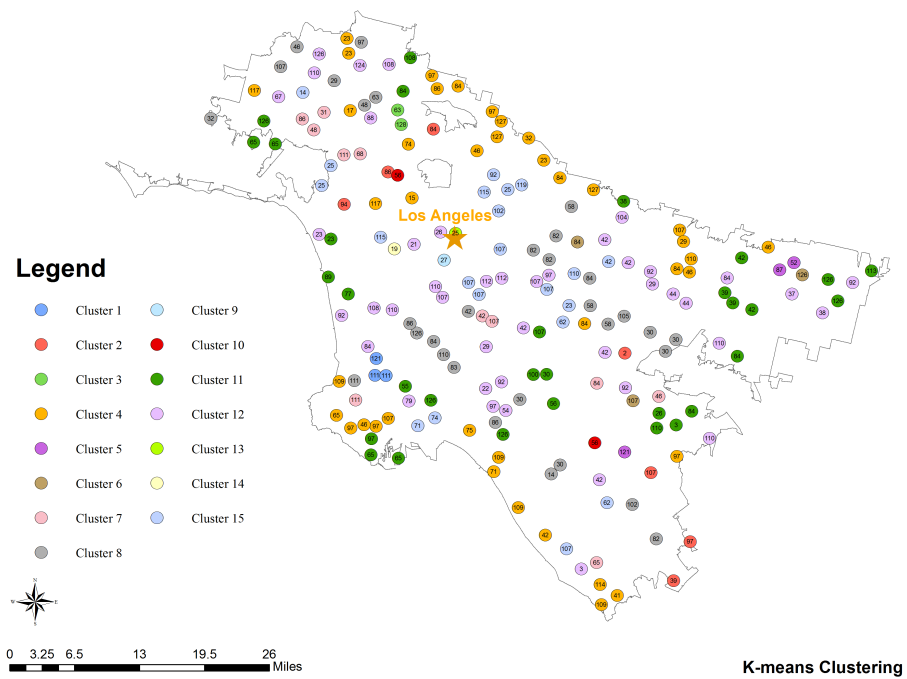


Figure 11: The K-means clustering result for 200 sampled places in Los Angeles.

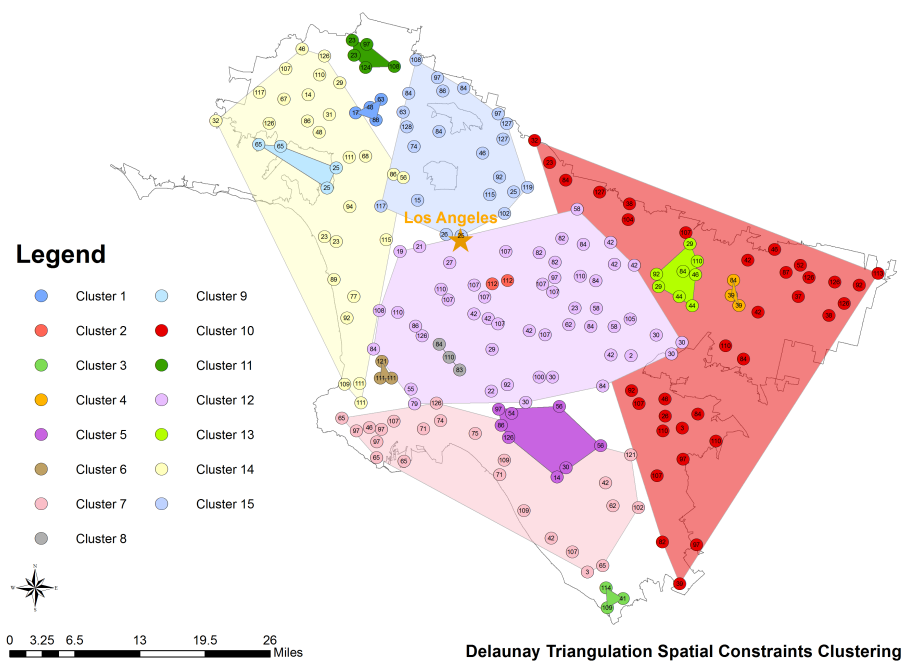


Figure 12: The Delaunay triangulation spatial constraints clustering and convex polygon generation result for 200 sampled places in Los Angeles.

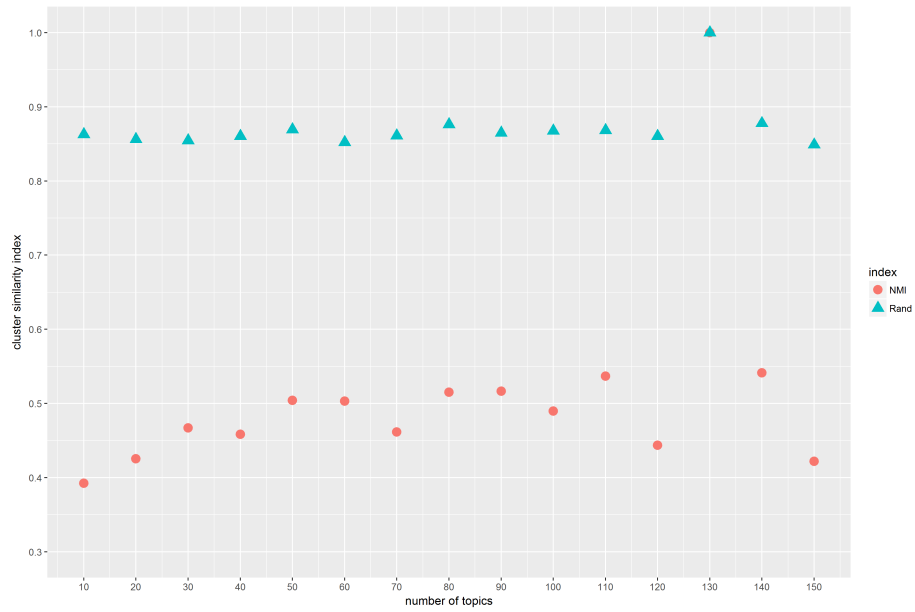


Figure 13: K-means clustering similarity evaluation using the NMI and Rand metrics with different number of topics.

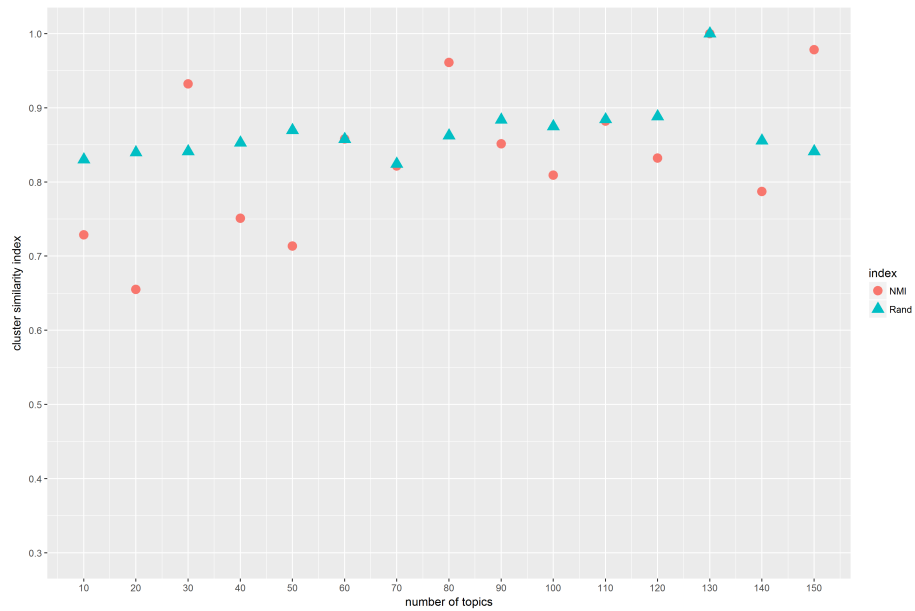


Figure 14: Delaunay triangulation spatial constraints clustering similarity evaluation using the NMI and Rand metrics with different number of topics.

572 gories function in geographic settings. Occasionally a topic may contain a less
573 meaningful or even an outlier POI category such as the *neighborhood* in *topic*
574 *109*. Data cleaning or post-processing may help to eliminate or reduce the noise.
575 This is the nature of crowdsourced data or user-generated content in which the
576 information is not validated by any authority. Also, the coverage and the accu-
577 racy of POI data in different cities may vary and the POI categories might also
578 change over time. We need to pay attention to those issues when interpreting
579 findings. In addition, we have also discovered various urban functional regions
580 as clusters of multinomial topic distributions over POI categories. However,
581 one limitation is that we cannot systematically evaluate the accuracy of those
582 derived functional regions without labeled ground truth data or the detailed
583 urban land-use GIS data. But we can test the intrinsic robustness of identify-
584 ing functional topics with different parameter settings. The variability analyses
585 were carried out at two levels: *the topic-level* and *the cluster-level*. At the topic
586 level, we found the stability in identifying prominent urban functional topics
587 related to frequently co-occurrent physical facilities and services, a variety of
588 bars and restaurants, and leisure activity places regardless of the total number
589 of topics. But the topic composition of top-ranked POI categories varies in dif-
590 ferent scenarios. It implies the variability of the semantic structure of functional
591 topics. Although choosing an optimal K in topic modeling can either maximize
592 the log-likelihood of the term-topic probability in the training document corpus,
593 or minimize the inter-topic similarity, we may miss the opportunity for discov-
594 ering some interesting topic composition structures that can only be identified
595 with a different K value or with other model parameter settings. At the cluster
596 level, a series of clustering result comparisons by choosing different numbers of
597 topics were evaluated using the Rand index and the NMI metric. We found a
598 large percentage of agreements on the clustering membership of those search

599 locations with their surrounding POIs and derived functional areas that can be
600 supported by the mix types of POIs.

601 One broader question is whether we can automatically identify those topolog-
602 ical and hierarchical relations in order to support the development of an ontology
603 for urban functional regions. As shown in Figure 15, we applied the Ward hier-
604 archical clustering method (Ward Jr, 1963) on those 130-topics derived from the
605 aforementioned LDA topic model. Each topic is a 480-dimensional probabilistic
606 vector over all POI categories in our datasets. Those semantically related topics
607 are grouped together in each step by minimizing the increment of within-cluster
608 variance after merging. This process repeats until all topic vectors merge into
609 the same group. This tree diagram is derived from a bottom-up approach and
610 can be used as a starting point with regard to constructing the urban functional
611 region ontology. However, it does not yet include the spatial relationships nor
612 the dichotomous relationship among POIs. It may be more promising to com-
613 bine this bottom-up approach with the top-down approach of the expert urban
614 geographers or planners to develop a more holistic ontology in the future.

615 **7 Conclusions and Future Work**

616 In this work, we develop a statistical framework that applies the LDA topic
617 modeling technique and incorporates user check-ins on LBSN in order to help
618 discover semantically meaningful topics and functional regions based on co-
619 occurrence patterns of POI types. The “functions” derived from probabilistic
620 topic modeling techniques can reveal the latent structure of POI mixtures and
621 the semantics of places. Based on a large corpus of about 100,000 Foursquare
622 venues and check-in behavior in the ten most populated urban areas in the
623 U.S., we demonstrate the effectiveness of proposed methodology by identifying
624 distinctive types of latent topics and further, by extracting urban functional

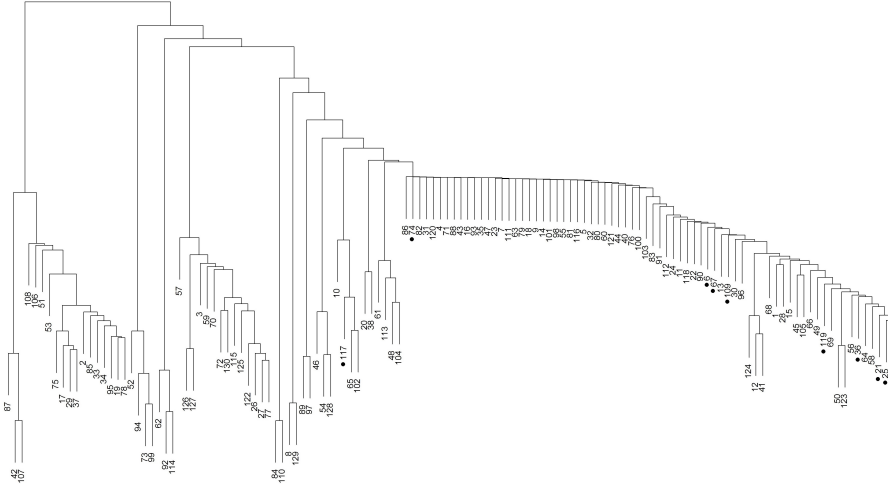


Figure 15: The dendrogram for the hierarchal clustering result on the 130-LDA topics using the Ward clustering method. The topics highlighted with a black filled-in circle are those mentioned in section 5.

625 regions using the K-means clustering and the Delaunay triangulation spatial
 626 constraints clustering methods. A region can have multiple functions but with
 627 different probabilities, while the same type of functional region can span multiple
 628 geographically non-adjacent locations. Compared with the remote sensing im-
 629 ages that mainly uncover the physical landscape of urban environments, results
 630 derived from the popularity-based POI topic model can be seen as a comple-
 631 mentary social sensing view of urban space based on human activities and the
 632 place settings of urban functions. However, there may exist gaps between the
 633 real-world business establishments and the online available POI information.
 634 Data-fusion and cross-validation relying on multiple sources may help reduce
 635 such gaps.

636 Although we have successfully identified several types of semantically mean-
 637 ingful urban functional topics, LDA topic modeling is an unsupervised approach
 638 that has certain limitation with respect to discovering plausible urban functions.
 639 In the future, we plan to investigate additional semantic signatures such as those

640 incorporating the spatial patterns of POI distributions and using supervised-
641 versions of probabilistic topic models to compare the performance of two fami-
642 lies of topic models (unsupervised or supervised) in discovering urban functional
643 regions. Last but not least, we also aim at developing a functional region on-
644 tology by combing the data-driven approach as outlined in this work with the
645 top-down knowledge engineering approach based on our understanding of urban
646 functional regions from human geography and urban planning.

647 **Acknowledgement**

648 We would like to thank Gengchen Mai for his help and discussion on the evalu-
649 ation of different clustering results. We also want to thank the editor and three
650 anonymous referees for their valuable comments and suggestions.

651 **References**

- 652 Adams, B., 2015. Finding similar places using the observation-to-generalization
653 place model. *Journal of Geographical Systems* 17 (2), 137–156.
- 654 Adams, B., Janowicz, K., 2012. On the geo-indicativeness of non-georeferenced
655 text. In: ICWSM. pp. 375–378.
- 656 Adams, B., Janowicz, K., 2015. Thematic signatures for cleansing and enriching
657 place-related linked data. *International Journal of Geographical Information*
658 *Science* 29 (4), 556–579.
- 659 Adams, B., McKenzie, G., 2013. Inferring thematic places from spatially refer-
660 enced natural language descriptions. In: *Crowdsourcing Geographic Knowl-*
661 *edge*. Springer, pp. 201–221.

- 662 Arun, R., Suresh, V., Madhavan, C. V., Murthy, M. N., 2010. On finding the nat-
663 ural number of topics with latent dirichlet allocation: Some observations. In:
664 Pacific-Asia Conference on Knowledge Discovery and Data Mining. Springer,
665 pp. 391–402.
- 666 Assunção, R. M., Neves, M. C., Câmara, G., da Costa Freitas, C., 2006. Effi-
667 cient regionalization techniques for socio-economic geographical units using
668 minimum spanning trees. *International Journal of Geographical Information*
669 *Science* 20 (7), 797–811.
- 670 Banzhaf, E., Netzband, M., 2012. Monitoring urban land use changes with re-
671 mote sensing techniques. *Applied Urban Ecology: A Global Framework* ,
672 18–32.
- 673 Barnsley, M. J., Barr, S. L., 1996. Inferring urban land use from satellite sensor
674 images using kernel-based spatial reclassification. *Photogrammetric Engineer-*
675 *ing and Remote Sensing* 62 (8), 949–958.
- 676 Blei, D. M., 2012. Probabilistic topic models. *Communications of the ACM*
677 55 (4), 77–84.
- 678 Blei, D. M., Ng, A. Y., Jordan, M. I., 2003. Latent dirichlet allocation. *Journal*
679 *of Machine Learning Research* 3 (Jan), 993–1022.
- 680 Cao, J., Xia, T., Li, J., Zhang, Y., Tang, S., 2009. A density-based method for
681 adaptive lda model selection. *Neurocomputing* 72 (7), 1775–1781.
- 682 Gao, S., Janowicz, K., Montello, D. R., Hu, Y., Yang, J.-A., McKenzie, G.,
683 Ju, Y., Gong, L., Adams, B., Yan, B., 2017. A data-synthesis-driven method
684 for detecting and extracting vague cognitive regions. *International Journal of*
685 *Geographical Information Science* , 1245–1271.

- 686 Gower, J. C., Ross, G., 1969. Minimum spanning trees and single linkage cluster
687 analysis. *Applied Statistics* , 54–64.
- 688 Griffiths, T. L., Steyvers, M., 2004. Finding scientific topics. *Proceedings of the*
689 *National Academy of Sciences* 101 (suppl 1), 5228–5235.
- 690 Han, J., Kamber, M., Pei, J., 2011. *Data mining: concepts and techniques*
691 (Third Edition). Elsevier, Waltham, MA.
- 692 Herold, M., Couclelis, H., Clarke, K. C., 2005. The role of spatial metrics in the
693 analysis and modeling of urban land use change. *Computers, Environment*
694 *and Urban Systems* 29 (4), 369–399.
- 695 Hobel, H., Abdalla, A., Fogliaroni, P., Frank, A. U., 2015. A semantic region
696 growing algorithm: extraction of urban settings. In: *AGILE 2015*. Springer,
697 pp. 19–33.
- 698 Hobel, H., Fogliaroni, P., Frank, A. U., 2016. Deriving the geographic footprint
699 of cognitive regions. In: *Geospatial Data in a Changing World*. Springer, pp.
700 67–84.
- 701 Hu, Y., Gao, S., Janowicz, K., Yu, B., Li, W., Prasad, S., 2015. Extracting and
702 understanding urban areas of interest using geotagged photos. *Computers,*
703 *Environment and Urban Systems* 54, 240–254.
- 704 Janowicz, K., 2012. Observation-driven geo-ontology engineering. *Transactions*
705 *in GIS* 16 (3), 351–374.
- 706 Jiang, S., Alves, A., Rodrigues, F., Ferreira, J., Pereira, F. C., 2015. Mining
707 point-of-interest data from social networks for urban land use classification
708 and disaggregation. *Computers, Environment and Urban Systems* 53, 36–46.
- 709 Kullback, S., Leibler, R. A., 1951. On information and sufficiency. *The Annals*
710 *of Mathematical Statistics* 22 (1), 79–86.

- 711 Lin, J., 1991. Divergence measures based on the shannon entropy. *IEEE Trans-*
712 *actions on Information Theory* 37 (1), 145–151.
- 713 Liu, Y., Liu, X., Gao, S., Gong, L., Kang, C., Zhi, Y., Chi, G., Shi, L., 2015.
714 Social sensing: A new approach to understanding our socioeconomic environ-
715 ments. *Annals of the Association of American Geographers* 105 (3), 512–530.
- 716 MacQueen, J., et al., 1967. Some methods for classification and analysis of
717 multivariate observations. In: *Proceedings of the Fifth Berkeley Symposium*
718 *on Mathematical Statistics and Probability*. Vol. 1. Oakland, CA, USA., pp.
719 281–297.
- 720 McKenzie, G., Janowicz, K., 2017. The effect of regional variation and resolution
721 on geosocial thematic signatures for points of interest. In: *Proceedings of the*
722 *2017 AGILE Conference*. Springer, pp. 237–256.
- 723 McKenzie, G., Janowicz, K., Gao, S., Yang, J.-A., Hu, Y., 2015. Poi pulse: A
724 multi-granular, semantic signature-based information observatory for the in-
725 teractive visualization of big geosocial data. *Cartographica: The International*
726 *Journal for Geographic Information and Geovisualization* 50 (2), 71–85.
- 727 Noulas, A., Scellato, S., Mascolo, C., Pontil, M., 2011. Exploiting semantic
728 annotations for clustering geographic areas and users in location-based social
729 networks. *The Social Mobile Web* 11, 02.
- 730 Pei, T., Sobolevsky, S., Ratti, C., Shaw, S.-L., Li, T., Zhou, C., 2014. A new
731 insight into land use classification based on aggregated mobile phone data.
732 *International Journal of Geographical Information Science* 28 (9), 1988–2007.
- 733 Rand, W. M., 1971. Objective criteria for the evaluation of clustering methods.
734 *Journal of the American Statistical association* 66 (336), 846–850.

- 735 Rousseeuw, P. J., 1987. Silhouettes: a graphical aid to the interpretation and
736 validation of cluster analysis. *Journal of Computational and Applied Mathe-*
737 *matics* 20, 53–65.
- 738 Steiger, E., Westerholt, R., Zipf, A., 2016. Research on social media feeds—
739 a giscience perspective. *European Handbook of Crowdsourced Geographic*
740 *Information* , 237.
- 741 Steyvers, M., Griffiths, T., 2007. Probabilistic topic models. *Handbook of Latent*
742 *Semantic Analysis* 427 (7), 424–440.
- 743 Strehl, A., Ghosh, J., 2002. Cluster ensembles—a knowledge reuse framework for
744 combining multiple partitions. *Journal of Machine Learning Research* 3 (Dec),
745 583–617.
- 746 Ward Jr, J. H., 1963. Hierarchical grouping to optimize an objective function.
747 *Journal of the American statistical association* 58 (301), 236–244.
- 748 Yao, Y., Li, X., Liu, X., Liu, P., Liang, Z., Zhang, J., Mai, K., 2017. Sensing
749 spatial distribution of urban land use by integrating points-of-interest and
750 google word2vec model. *International Journal of Geographical Information*
751 *Science* 31 (4), 825–848.
- 752 Yuan, J., Zheng, Y., Xie, X., 2012. Discovering regions of different functions
753 in a city using human mobility and pois. In: *Proceedings of the 18th ACM*
754 *SIGKDD International Conference on Knowledge Discovery and Data Mining.*
755 *ACM*, pp. 186–194.
- 756 Zhi, Y., Li, H., Wang, D., Deng, M., Wang, S., Gao, J., Duan, Z., Liu, Y.,
757 2016. Latent spatio-temporal activity structures: a new approach to inferring
758 intra-urban functional regions via social media check-in data. *Geo-spatial In-*
759 *formation Science* 19 (2), 94–105.

- 760 Zhong, C., Huang, X., Arisona, S. M., Schmitt, G., Batty, M., 2014. Infer-
761 ring building functions from a probabilistic model using public transportation
762 data. *Computers, Environment and Urban Systems* 48, 124–137.
- 763 Zhou, X., Zhang, L., 2016. Crowdsourcing functions of the living city from twit-
764 ter and foursquare data. *Cartography and Geographic Information Science*
765 43 (5), 393–404.