

SPATIAL DEPENDENCY OF URBAN SPRAWL AND THE UNDERLYING ROAD NETWORK STRUCTURE

Xianyuan Zhan¹ and Satish V. Ukkusuri²

¹JD Intelligent City Research, JD Digits. Email: zhanxianyuan@jd.com. Formerly at Lyles
School of Civil Engineering, Purdue University

²Professor, Lyles School of Civil Engineering, Purdue University. Email: sukkusur@purdue.edu

ABSTRACT

The spatial correlation between urban sprawl and the underlying road network has long been recognized in urban studies. The accessibility to road networks is often considered as an approximation for the measurement of human mobility, which is a key factor in determining potential urban sprawl in the future. Despite the close relationship between urban development and road networks, the spatial dependency of these two spatial layers has never been systematically evaluated. This paper conducts a comprehensive investigation on the spatial dependency between these two spatial layers using an urban expansion dataset between 2000 and 2010 of East Asian regions and the road network data from OpenStreetMap. Four Chinese cities, namely, Beijing, Shanghai, Chengdu and Shenzhen are selected to conduct the analysis. The spatial correlations between the urban sprawl and road networks are first quantitatively analyzed using Ripley's cross-K function. Highly significant spatial correlation have been observed in all the four tested cities. A Bayesian network model is also developed to verify the predictability of urban sprawl using the spatial and structural features extracted from the existing road networks as well as the spatial pattern of the past built-up areas. The results show an affirmative answer to the predictability of urban sprawl, by achieving an overall accuracy of 79% in classifying urban sprawl and undeveloped areas. Finally, the hidden dependencies among the urban sprawl and the extracted spatial features are interpreted and analyzed based on the Bayesian network structure learned from the data.

INTRODUCTION

The road network defines a basic template of the urban area that strongly constraints the urban development. It plays a prominent role in human mobility and activity analysis that directly impacts our understanding of urban sprawl patterns. The accessibility to road network is often considered as an important approximation for the measurement of human mobility, which is a key factor of the potential urban sprawl in a city (Obregón-Biosca et al. 2015). As a result, the co-location patterns of the existing urban settlement and the transportation infrastructures can be easily observed in many cities in the world. Such strong spatial correlation between urban built-up areas and road networks has naturally led to a hypothetical question: does strong spatial correlation exist between urban sprawl and road networks?

The evidences of the correlation between construction of new roads and future urban sprawl have long been recognized in literature (Harvey and Clark 1965; Bhatta 2010; Zischg et al. 2019). For example, the construction of highways leads to both congestion in the city and rapid outgrowth (Harvey and Clark 1965). A study by Yang (2002) on the urban sprawl of Atlanta during 1973-1999 also observed outward spread of high-density urban use along major transportation routes. Although the majority of urban sprawl studies mainly focus on analyzing macroscopic factors that impact the city growth, such as urban geometry, size relationship between cities, economic functions, social demographic and ethnic patterns, etc., many researchers have started incorporating road related features in modeling and forecasting urban sprawl (Yang and Lo 2003; Cheng and Masser 2003; Irwin and Bockstael 2007; Fregolent and Tonin 2016; Xu et al. 2014) due to the observation of existence of possible correlation between urban sprawl and road networks. For instance, Cheng and Masser (2003) developed a spatial logistic regression model to predict the urban sprawl pattern of the city of Wuhan in China. The model incorporated a set of variables that measure the distances of the given area to multiple types of road. Their study showed that the urban road infrastructure is one of the major determinants of urban growth.

Despite the growing consensus of the potential association of urban sprawl and the underlying road network, the spatial dependency patterns of these two spatial layers have never been systemat-

ically evaluated, and many important research questions still remain unanswered. For example, one simple but fundamental research question is that to what degree the urban sprawl and underlying road network are correlated. Quantifying the level of spatial correlation serves as a first step for us to fully understand the causal factors behind the coevolution of these two closely related spatial layers. Further, it helps to answer some interesting questions, such as: does a local cluster of urban sprawl always correspond to a local cluster of road networks? Are there any density threshold or other conditions exist on the structure and spatial distribution of the road network to enable the urban sprawl? If so, how can we design the road network to guide the desirable urban sprawl? This study aims to abridge these two spatial layers together by comprehensively investigating the spatial dependency between urban sprawl and road networks. Specifically, we investigate following sub-questions: (1) on what level the urban sprawl and road network are spatially correlated? (2) given the existence of such strong correlation, is urban sprawl predictable? And what are the relevant spatial features in determining the urban sprawl?

In this study, a large-scale urban expansion dataset between 2000 and 2010 of East Asia region from World Bank (Schneider et al. 2015; World Bank 2015) and the road network data extracted from OpenStreetMap (OpenStreetMap 2016) are used to conduct the analysis. Four different Chinese cities, namely, Beijing, Shanghai, Chengdu and Shenzhen are selected for the analysis. The spatial correlations between the urban sprawl and road networks are first quantitatively analyzed using cross-K function. The urban sprawl and the road network are converted into two spatial point processes with the original spatial relation information preserved to enable the utilization of cross-K function. A Bayesian network model is then developed to verify the predictability of urban sprawl using the spatial and structural features extracted from the existing road networks as well as the spatial patterns of the past built-up areas. Finally, the hidden relationships among the urban sprawl and the extracted spatial features are interpreted by analyzing the inferred structure of the Bayesian network.

The paper is organized as follows: the next section describes the data used in this study. Section 3 quantitatively evaluates the extent of spatial correlations between the urban sprawl and road

79 networks of different cities. Section 4 develops an Bayesian network model to investigate the
80 predictability of urban sprawl and further explores the contributing features associated with the
81 urban sprawl. The final section concludes the paper.

82 DATA

83 The data used in this chapter were obtained from multiple sources. The urban sprawl data were
84 obtained from a large-scale urban expansion dataset produced in a World Bank study (Schneider
85 et al. 2015; World Bank 2015), which contains the urban expansion information across the East
86 Asian region (stretching from Mongolia to the Pacific Islands) between 2000 and 2010. The data
87 are in raster map format. The map is divided into 250×250m uniform cells and a specific label is
88 assigned to each cell to indicate whether the cell is a built-up area before 2000, a new urban sprawl
89 area between 2000 and 2010, or an undeveloped area by 2010. The road network data of the year
90 2012 from OpenStreetMap OpenStreetMap (2016) were extracted to provide an approximation of
91 road network structure at the year 2010. Ideally, we want to use two sets of road network data around
92 2000 and 2010 to fully explore the spatial dependency between urban sprawl and road networks.
93 Unfortunately, high-quality road network data in China around the year 2000 are not obtainable,
94 hence we focus on analyzing the spatial dependency of urban sprawl and the road network structure
95 after the urban sprawl.

96 Four representative cities from different geographical regions of China were selected to conduct
97 the spatial dependency analysis, namely Beijing (north), Shanghai (east), Chengdu (southwest) and
98 Shenzhen (south). All the four cities have experienced rapid urban sprawl during the 2000-2010
99 period. Both the urban sprawl and the road network data of the four cities were extracted. The urban
100 sprawl pattern and the underlying road networks are illustrated in Fig.1a-1d (map data obtained
101 from OpenStreetMap), and their summary statistics are presented in Table 1. The urban sprawl
102 data in raster format were further processed and converted into a set of labeled points by placing
103 a point at the center of each 250×250m cell. Every cell is thus represented by a specific point
104 with a label that indicates the urban development condition of the area. The conversion from cell
105 to points enables the urban sprawl to be modeled as a spatial point process and allows for more

106 efficient computation in the analysis. For convenience, we refer to the points labeled as built-up
107 areas as *built-up points*, the urban sprawl area as *urban sprawl points* and the rest of the points as
108 *undeveloped points*.

109 **SPATIAL CORRELATION BETWEEN URBAN SPRAWL AND ROAD NETWORK**

110 The co-location patterns of urban sprawl and road networks are evident in many cities, since
111 roads are built to make the human settlements accessible to other parts of the city. However, such
112 correlation patterns have never been systematically investigated, and many research questions still
113 remain. For example, to what degree these two spatial layers are correlated? Does a local cluster
114 of urban sprawl always correspond to a local cluster of road networks? To adequately address
115 these questions, the spatial correlation between the urban sprawl and road network needs to be
116 quantitatively evaluated. In this section, we develop a quantitative evaluation method to examine
117 the degree of spatial correlation between urban sprawl and road networks. The approach is based
118 on a new distance measure (point-to-road distance) for evaluating the spatial proximity between a
119 point process and a spatial network, as well as the cross-K function.

120 **Model development**

121 There is a large body of research and statistical techniques on measuring the spatial correlation
122 between two spatial point processes, including the Ripley's Cross-K function with Monte Carlo
123 simulation (Cressie 1993), mean nearest-neighbor distance (Dixon 2002), and spatial regression
124 models (Chou 1997). However, there is no existing method on evaluating the spatial correlation
125 between a point process and a spatial network, as the spatial relationship between a point and a line
126 segment is far more complex than a pair of points. To utilize the well-established theoretical results
127 on the spatial correlations of point processes as well as generalize the analysis to point-network
128 analysis, we construct a new point process from the road network while preserving sufficient spatial
129 relationship between the given point process and the original road network. The key step of this
130 construction is to introduce a new set of points (referred as *access point*) and a new distance
131 measure (referred as *point-to-road distance*) to evaluate the spatial proximity between a point and
132 a road segment. Given a target point (e.g. an urban sprawl point), the corresponding access point

133 on a road segments is defined as its projection (point on the road segment that has the minimum
134 distance to the target point) on the road segment. Since the points on a curve does not necessarily
135 form a convex set, the projection operation may not lead to the unique solution. Under such cases,
136 we only pick one point in the solution set as the access point. The point-to-road distance between
137 the target point and a road segment can thus be defined as the great-circle distance (the shortest
138 distance between two points on the surface of a sphere) with Earth radius between the target point
139 and the corresponding access point of the road segment. Fig.2 presents a conceptual illustration of
140 the aforementioned access point and the point-to-road distance. There are three important features
141 of this construction:

- 142 1. Every target point has only one access point for each road segment.
- 143 2. The point-to-road distance captures the spatial proximity of the target point and a road
144 segment. As a large point-to-road distance indicates all points on the road segment are far
145 away from the target point.
- 146 3. All roads that have a point-to-road distance smaller than h to the target point will have
147 non-empty intersection with the area formed by the neighborhood of the target point with
148 radius h . Hence the set of access point in the neighborhood of the target point corresponds
149 to the same set of road segments that intersect with the neighborhood area.

150 Given the above three important features of the access points and the point-to-road distance, the
151 spatial relationship between the point process and the spatial network can be adequately captured.
152 The spatial correlation of the two spatial layers can thus be approximated by the spatial correlation
153 between the point process and a set of constructed access point processes for each target point in
154 the original point process.

155 The most widely used spatial statistics for evaluating the spatial correlation between two point
156 processes is the cross-K function, a generalization of Ripley's K-function (Cressie 1993; Huang
157 et al. 2004; Dixon et al. 2002). The cross-K function for binary spatial features is defined as follows:

$K_{ij}(h) = \lambda_j^{-1} E[\text{number of type } j \text{ instances within distance of a randomly chosen type } i \text{ instance}]$

158 where λ_j is the density (number per unit area) of type j instances and h is the distance. Without
 159 edge effect (Cressie 1993; Dixon et al. 2002), the cross-K function can be estimated by

$$160 \quad \hat{K}_{ij}(h) = \frac{1}{\lambda_i \lambda_j A} \sum_k \sum_l I_h(d(i_k, j_l)) \quad (1)$$

161 where A is the total area of the study region, $d(i_k, j_l)$ is the distance between the k th instance
 162 of type i and the l th instance of type j ; $I_h(d(i_k, j_l))$ is the indicator function which takes value
 163 1 if $d(i_k, j_l)$ is smaller than h , 0 otherwise. The edge effect arises because points outside the
 164 study region are not counted in the numerator, even if they are within distance h of a point in
 165 the study region (Dixon et al. 2002). However, fully removal of the edge effect needs to perform
 166 the computationally expensive edge-correction, thus is often omitted (Cressie 1993). In our case,
 167 since the spatial correlation is measured between the urban sprawl points and their corresponding
 168 access point sets, above estimator needs to be modified to make it suitable for our analysis. Denote
 169 $P = \{p_i, i = 1, 2, \dots, n\}$ as the set of urban sprawl points, $Q = \{q_r^i, r = 1, 2, \dots, m\}$ as the set
 170 of access points of urban sprawl point p_i , in which q_r^i is the access point of p_i on road r , and
 171 $d_{pr}(p_i, r) = d(p_i, q_r^i)$ is the point-to-road distance between point p_i and road r . The equivalent
 172 estimator of cross-K function without edge-correction can be modified as follows

$$173 \quad \hat{K}_{ij}(h) = \frac{A}{nm} \sum_{i=1}^n \sum_{r=1}^m I_h(d_{pr}(p_i, r)) \quad (2)$$

174 The cross-K function characterizes pairwise spatial relationship between two point processes.
 175 To test the extent of spatial correlation between two point processes, a typical treatment is to
 176 compare the results against the curve of $K(h) = \pi h^2$, which indicates complete spatial randomness
 177 when edge effect is not present. However, in this study, the edge effect is present. This can be
 178 easily observed from the urban sprawl pattern of Shanghai (Fig.1b) and Shenzhen (Fig.1d), as a

179 considerable amount of urban sprawl areas are located near the boundary of the two cities. Thus
180 comparing against the curve of $K(h) = \pi h^2$ is not appropriate and may lead to biased results. In
181 this study, we adopt an alternative approach by comparing the result against a baseline cross-K
182 function produced by the urban sprawl points and a completely independent point process with
183 points randomly generated inside the study region. The random points generation is performed
184 using Monte Carlo simulation with the number of random points equal to the number of road
185 segments. We run the simulation 20 times for each city and report the 0.05-0.95 quantile of the
186 computed baseline cross-K functions.

187 **Experimental results**

188 Fig.3a-3d show the results of the cross-K function for the four cities investigated in this study.
189 In all of the four cities, the curves of cross-K functions are much higher than the 0.05-0.95 quantile
190 of the baseline cross-K functions. It suggests that in all distance scales, the average number of road
191 segments within distance h to an urban sprawl area is much higher compared with an independent
192 and randomly distributed point process with the same density. This clearly indicates that strong
193 attraction behavior exists, that the urban sprawl areas and the road segments tend to be co-located.
194 The cross-K function results for Beijing, Shanghai and Shenzhen exhibit similar patterns, where
195 their cross-K function values are about twice as much as the baseline cross-K function values,
196 which means in average, there are about twice the number of road segments located within certain
197 distance to an urban sprawl point compared with the number of points generated by an independent
198 randomly distributed spatial process. The variance of the baseline cross-K function of Shenzhen is
199 a little larger (wider 0.05-0.95 quantile range) compared with the results of Beijing and Shanghai,
200 which might caused by more significant edge effect. Different from the previous three cities,
201 Chengdu exhibits much stronger spatial correlation between urban sprawl and road networks, as its
202 cross-K function curve is significantly higher than the baseline cross-K function curve. This might
203 caused by the mountainous terrain surrounding Chengdu, that most of the urban development areas
204 and roads are concentrated in a series of clusters. The higher level of local clustering contributes
205 to a stronger spatial correlation between the two spatial layers. Despite the differences, the results

206 in all the four cities confirmed the existence of strong spatial correlation between the urban sprawl
207 and the underlying road networks.

208 **PREDICTABILITY OF URBAN SPRAWL GIVEN UNDERLYING ROAD NETWORK**

209 The previous section provides an affirmative answer to the existence of strong correlation
210 between urban sprawl and the underlying road network. Given the existence of strong spatial
211 correlation, a natural and more interesting research question to be answered is that: is urban sprawl
212 predictable using the information contained in the road network structures? And what are the
213 relevant spatial and structural features in the road network that provide discriminative information
214 on deciding whether an area is part of the future urban sprawl or remains undeveloped during
215 a specific observation period. Answering these questions will not only contribute to a better
216 understanding of the mechanism behind the co-evolution of urban sprawl and road network growth,
217 but also have important practical implications. For example, understanding of the structural impacts
218 of road networks can provide insights on guiding road network construction that lead to desirable
219 and healthy urban sprawl in the future.

220 In this section, we focus on addressing aforementioned questions by verifying the predictability
221 of the urban sprawl using spatial and structural features extracted from the road network as well as
222 the spatial pattern of the past built-up areas. The problem is cast into a binary classification and
223 prediction problem on the area label (urban sprawl or undeveloped). A Bayesian network model
224 is developed to learn the dependencies and causal relationships between extracted spatial features
225 and the area labels, and further predicts the target area label.

226 **Spatial feature extraction**

227 To begin our analysis, a set of spatial features were first extracted from the urban sprawl and the
228 road network data. A list of extracted features as well as their descriptions are presented in Table
229 2. Our target variable is *sprawlLabel* which encodes the actual state of the area, e.g. whether it
230 is an urban sprawl or undeveloped area. Two global statistics are computed, namely *minUrbDist*
231 and *minRoadDist* which are the distances of the area to the nearest built-up area and the road
232 segment (using point-to-road distance) respectively. Local spatial features within 1km radius of

233 a target area were also extracted, including the number of built-up areas (*NneighUrb*) and road
234 segments (*NneighRoad*) fall within the 1km radius; the number of road intersections with 3 or
235 more approaches (*Nintersection*); and the lengths of each road type based on corresponding type
236 information provided in OpenStreetMap (*primaryLen*, *secondaryLen*, *localLen*, *otherLen*). More
237 detailed information on the extracted features can be found in Table 2. When computing the road
238 lengths, we only calculated the part of length of the road segment within the 1km radius circle,
239 rather than the total length of the road segments. The local road density can thus be obtained by
240 summing the extracted road length of different types and divided by $\pi \times 1^2 km^2$.

241 Since the feature extraction is computationally extensive, for each city, we randomly selected
242 500 urban sprawl areas and 500 undeveloped areas as instances and extracted all the spatial features
243 for every instance. Hence, for each of the four cities, we obtain a dataset of 1000 instances. After
244 that, the dataset for each city were randomly partitioned into 2 sub-datasets, the first contains 800
245 instances which served as the training set, and the second contains 200 instances which served as the
246 testing set. In the actual implementation of the urban sprawl prediction model, the feature *otherLen*
247 was removed from the model training and testing phase. As it is correlated with *primaryLen*,
248 *secondaryLen*, *localLen* and *roadDensity* ($primaryLen + secondaryLen + localLen + otherLen \propto$
249 *roadDensity*) and does not provide additional information.

250 **Urban sprawl prediction**

251 Since all instance labels are known, we adopted the supervised learning approach in machine
252 learning to predict the area labels. Four widely used supervised learning classifiers were tested and
253 examined in this study, namely, Naïve Bayes, support vector machine (SVM), random forest and
254 Bayesian network. The Naïve Bayes is a simple probabilistic classifier based on Bayes' Theorem,
255 which assumes conditional independence among features. SVM is another popular method, which
256 classifies data by maximizing the distance between the decision boundaries defined by a set of
257 control points (support vectors) (Bishop 2006). On the other hand, the random forest (Breiman
258 2001) approach is an ensemble method, which uses a combination of randomly generated decision
259 tree classifiers to increase accuracy. The Bayesian network (Bishop 2006) is a probabilistic graphical

260 model which represents a set of random variables $U = \{x_1, x_2, \dots, x_n\}$, $n \geq 1$ and their conditional
261 dependencies via a directed acyclic graph (DAG). A Bayesian network represents a probability
262 distribution as follows

$$263 \quad P(U) = \prod_{u \in U} p(u|pa(u)) \quad (3)$$

264 where $p(u|pa(u))$ is the conditional probability table between random variable u and its parents
265 $pa(u)$. More detailed background information about the theory, learning and inference procedures
266 of Bayesian network can be found in (Bishop 2006; Koller and Friedman 2009). We used the
267 machine learning software Weka (Hall et al. 2009) to implement all of the previous four supervised
268 learning classifiers. The training dataset for all of the four cities were combined (4000 instances)
269 to train each classifier, and the overall classification accuracy on the combined test datasets (800
270 instances) were used as the final criteria for model selection. The Bayesian network was found to
271 achieve the highest overall classification accuracy of 79%, thus was selected as the final model to
272 perform the urban sprawl prediction task. Another advantage of using the Bayesian network for
273 our analysis is that Bayesian network has the ability to reflect causal relationships between between
274 variables. As is shown in Pearl (Pearl 2009), if a set of variables have causal relations, and the
275 Bayesian network is built such that arcs fully represent the causal paths between variables, then the
276 resulting Bayesian network will encode dependencies and probabilistic relations between variables.
277 This property is particularly helpful for us to identifying relevant spatial features that impact the
278 urban sprawl.

279 To implement the Bayesian network, we first discretized each continuous feature into a set
280 of discrete states. The K2 algorithm developed by Cooper et al. (Cooper and Herskovits 1992)
281 was then applied to learn the structure of the Bayesian network as well as the parameters in the
282 conditional probability tables $p(u|pa(u))$, $\forall u \in U$. In the prediction phase, using the already learned

283 conditional probability tables, the area label for area i is predicted as follows:

$$284 \quad \text{sprawlLable}_i^* = \underset{\text{sprawlLable}}{\operatorname{argmax}} P(\text{sprawlLable} | \mathbf{x}^i) \propto \prod_{u^i \in U} p(u^i | pa(u^i)) \quad (4)$$

285 where \mathbf{x}^i is the set of all observed features except *sprawlLable* for area i .

286 **Experimental result**

287 The structure of the Bayesian network learned from the training data is presented in Fig.
288 4. It can be observed that all the extracted spatial features are closely associated with the area
289 label, as there is a direct arc connecting *sprawlLable* to every other feature. This indicates that
290 the emergence of an urban sprawl area has the potential to impact all the local features of the
291 road network. It confirms our intuition that urban sprawl and road network are co-evolved with
292 each other, and the development of urban settlement will lead to the construction of new roads
293 in order to establish connectivity to other parts of the city. Among other spatial features, the
294 conditional dependency between the distance to closest road segments (*minRoadDist*) and the road
295 density (*roadDensity*) is obvious, since a location that is far away from any other road segments
296 is not possible to have high local road density. The road lengths of different types (*primaryLen*,
297 *secondaryLen*, *localLen*) are conditionally dependent on road density (*roadDensity*) and the number
298 of neighboring road segments (*NneighRoad*), which is also intuitive. The distance to the closest
299 road segment (*minRoadDist*) and the number of neighboring road segments (*NneighRoad*) are
300 conditionally dependent on the distance to the closest built-up area, which also makes sense. As
301 locations that are far away from any existing built-up areas are less likely to have dense road network
302 as well as the potential of urban sprawl. From above analysis, it can be observed that the structure
303 of Bayesian network provides rich information about the hidden relationship among the extracted
304 spatial features, and allows us to better understand the role of each spatial feature in the emergence
305 of urban sprawl.

306 The developed Bayesian network model was validated using the test dataset of the four studied
307 cities. The testing results of the combined test set as well as the test set for each city are presented

308 in Table 3. The confusion matrix for the experiment using the combined test set is presented in
309 Table 4. In addition to the accuracy measure, F-measure (also refer as F1) is used to evaluate the
310 classification quality on the test sets, which is a commonly used accuracy measure in data mining. F-
311 measure is computed as the harmonic mean of precision (percentage of testing data that are classified
312 as positive are actually positive) and recall (percentage of positive testing data that are classified
313 as positive), which is calculated as $F\text{-measure} = 2(\text{precision} \times \text{recall})/(\text{precision} + \text{recall})$. A
314 higher F-measure suggests a better classification result.

315 It can be observed that even only use the spatial information provided in road networks as
316 well as the past built-up areas, the Bayesian network can achieve about 80% overall accuracy in
317 classifying the area labels. The accuracy of both Beijing and Shanghai are around 80%. For the
318 best case - the result of Chengdu, the accuracy even achieved 85%. As we have already observed in
319 Section 3, Chengdu exhibits highest level of spatial correlation between urban sprawl and the road
320 network. This confirmed our hypothesis that high level of spatial correlation indeed contributes to
321 the predictability of urban sprawl. The test accuracy of the experiment for Shenzhen is the lowest,
322 however, still reaches 73%. The relatively lower accuracy for Shenzhen might again associated
323 with the edge effect discussed in Section 3, which also impacts the spatial feature extraction. As
324 the training and testing instances are randomly selected, if they are located close to the boundary
325 of the study region, the local spatial features will be incomplete and lead to inaccurate information.
326 For the results of the F-measure, not surprisingly, Chengdu has the highest F-measure for both area
327 label classes among the four cities. The F-measure for urban sprawl areas are consistently higher
328 than the F-measure for areas remain undeveloped, which suggests lower misclassification error in
329 predicting urban sprawl areas, this partly indicates the urban sprawl is relatively more predictable
330 than the areas that will remain undeveloped during the observation period. All the results show that
331 the urban sprawl is highly predictable given the strong spatial correlation between urban sprawl and
332 the road networks. Furthermore, the dependency relationship among the extracted spatial features
333 revealed by the structure of Bayesian network can also be used as an input to guide the road network
334 construction that leads to a more desirable and healthy urban sprawl in the future.

335 CONCLUSION

336 This paper systematically investigates the spatial dependency between urban sprawl and the
337 underlying road networks. An urban sprawl dataset between 2000 and 2010 for East Asia region
338 and the road networks from OpenStreetMap are used in this study. Four Chinese cities, namely,
339 Beijing, Shanghai, Chengdu and Shenzhen are selected to conduct the analysis. The spatial
340 correlation between the urban sprawl and road network is first quantitatively evaluated using cross-
341 K function. Highly significant spatial correlation has been observed in all of the four tested cities.
342 A Bayesian network model is then developed to verify the predictability of urban sprawl given the
343 existence of strong spatial correlation of the two spatial layers. The results provide an affirmative
344 answer to the predictability of urban sprawl, that about 79% of overall prediction accuracy is
345 achieved by using a set of spatial and structural features extracted from the road networks as well
346 as the spatial patterns of the past built-up areas.

347 There are some limitations of this study. First, since the historical road networks around year
348 2000 is not obtainable, some of the research questions related to the co-evolution of the urban
349 sprawl and road network can not be fully explored. Also, as this study focuses on exploring the
350 spatial dependencies between the urban sprawl and the road network structure, thus the prediction
351 of urban sprawl solely rely on the spatial correlation of these two spatial layers. It is expected that
352 even higher prediction accuracy can be achieved when incorporating additional information, e.g.
353 the terrain feature of the city, demographics characteristics such as population growth, GDP growth
354 during the observation period and the prior knowledge from the local planning policies, etc. Future
355 research can be done to develop more comprehensive decision support tools for city planners for
356 more accurate urban sprawl prediction.

357 REFERENCES

- 358 Bhatta, B. (2010). *Analysis of urban growth and sprawl from remote sensing data*. Springer Science
359 & Business Media.
- 360 Bishop, C. M. (2006). *Pattern recognition and machine learning*. springer.
- 361 Breiman, L. (2001). "Random forests." *Machine learning*, 45(1), 5–32.

362 Cheng, J. and Masser, I. (2003). "Urban growth pattern modeling : a case study of Wuhan city ,
363 PR China." *Landscape and urban planning*, 62(4), 199–217.

364 Chou, Y.-H. (1997). "Exploring spatial analysis in geographic information systems." *Exploring*
365 *Spatial Analysis in Geographic Information Systems*, OnWord Press.

366 Cooper, G. F. and Herskovits, E. (1992). "A bayesian method for the induction of probabilistic
367 networks from data." *Machine learning*, 9(4), 309–347.

368 Cressie, N. A. C. (1993). *Statistics for Spatial Data*, Vol. 14 of *Wiley Series in Probability and*
369 *Statistics*. John Wiley & Sons, Inc., Hoboken, NJ, USA (September).

370 Dixon, P. M. (2002). "Nearest neighbor methods." *Encyclopedia of environmetrics*.

371 Dixon, P. M., El-shaarawi, A. H., and Piegorisch, W. W. (2002). "Ripley' s K function." *Encyclopedia*
372 *of Environmetrics*, 3, 1796–1803.

373 Fregolent, L. and Tonin, S. (2016). "Local public spending and urban sprawl: Analysis of this
374 relationship in the Veneto region of Italy." *Journal of Urban Planning and Development*, 142.

375 Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., and Witten, I. H. (2009). "The
376 weka data mining software: an update." *ACM SIGKDD explorations newsletter*, 11(1), 10–18.

377 Harvey, R. O. and Clark, W. A. (1965). "The nature and economics of urban sprawl." *Land*
378 *Economics*, 1–9.

379 Huang, Y., Shekhar, S., and Xiong, H. (2004). "Discovering colocation patterns from spatial data
380 sets: a general approach." *IEEE Transactions on Knowledge and Data Engineering*, 16(12),
381 1472–1485.

382 Irwin, E. G. and Bockstael, N. E. (2007). "The evolution of urban sprawl: evidence of spatial
383 heterogeneity and increasing land fragmentation." *Proceedings of the National Academy of*
384 *Sciences*, 104(52), 20672–20677.

385 Koller, D. and Friedman, N. (2009). *Probabilistic graphical models: principles and techniques*.
386 MIT press.

387 Obregón-Biosca, S. A., Romero-Navarrete, J. A., Mendoza-Sanchez, J. F., and Betanzo-Quezada,
388 E. (2015). "Impact of mobility induced by urban sprawl: Case study of the Querétaro metropolitan

389 area.” *Journal of Urban Planning and Development*, 142(2), 05015005.

390 OpenStreetMap (2016), <<https://www.openstreetmap.org/>>.

391 Pearl, J. (2009). *Causality*. Cambridge University Press.

392 Schneider, A., Mertes, C. M., Tatem, A. J., Tan, B., Sulla-Menashe, D., Graves, S. J., Patel, N. N.,
393 Horton, J. A., Gaughan, A. E., Rollo, J. T., Schelly, I. H., Stevens, F. R., and Dastur, A. (2015).
394 “A new urban landscape in east–southeast asia, 2000-2010.” *Environmental Research Letters*,
395 10(3), 034002.

396 World Bank (2015). “World bank urban spatial dataset: Puma,
397 <<http://puma.worldbank.org/downloads>>.

398 Xu, C., Fang, S., Long, N., Teng, S., Zhang, M., and Liu, M. (2014). “Spatial patterns of dis-
399 tinct urban growth forms in relation to roads and pregrowth urban areas: Case of the nanjing
400 metropolitan region in china.” *Journal of Urban Planning and Development*, 141(1), 04014015.

401 Yang, X. (2002). “Satellite monitoring of urban spatial growth in the Atlanta metropolitan area.”
402 *Photogrammetric Engineering and Remote Sensing*, 68(7), 725–734.

403 Yang, X. and Lo, C. P. (2003). “Modelling urban growth and landscape changes in the Atlanta
404 metropolitan area.” *International Journal of Geographical Information Science*, 17(5), 463–488.

405 Zischg, J., Klinkhamer, C., Zhan, X., Rao, P. S. C., and Sitzenfrie, R. (2019). “A century of
406 topological coevolution of complex infrastructure networks in an alpine city.” *Complexity*, 2019,
407 2096749.

408 **List of Tables**

409 1 Summary of the urban sprawl in the four cities 18

410 2 Description of the extracted features 19

411 3 Testing results for the Bayesian network 20

412 4 Confusion matrix for the experiment using combined testing set 21

TABLE 1. Summary of the urban sprawl in the four cities

Summary statistics	Beijing	Shanghai	Chengdu	Shenzhen
Total area (km^2)	9,910.31	5,264.28	11,487.77	1,858.72
Total number of built-up points	28,546	17,549	11,113	9,597
Total number of urban sprawl points	13,992	16,428	8,484	4,099
Total number of undeveloped points	116,044	50,259	16,4232	16,048
Total number of road segments	46,711	11,151	47,442	18,739
Population at 2000	1633.0	1858.0	1257.9	701.2
Population at 2010	1961.2	2301.9	1404.8	1035.79
Population growth	20.10%	23.89%	11.68%	47.71%

TABLE 2. Description of the extracted features

Feature name	Description
sprawlLabel	Binary variable indicating the class of the $250 \times 250m$ area: urban sprawl (1) and undeveloped area (0)
minUrbDist	Great circle distance (km) to the nearest built-up area
minRoadDist	Point-to-road distance (km) to the nearest road segment
NneighUrb	Number of existing built-up areas within 1km radius
NneighRoad	Number of road segments within 1km radius
roadDensity	Road density within 1 km radius (km^{-1})
Nintersection	Number of intersections with at least 3 approaches within 1km radius
primaryLen	Total distance (km) of primary road (e.g. primary link, motorway, raceway)
secondaryLen	Total distance (km) of secondary road (e.g. secondary link)
localLen	Total distance (km) of local road (e.g. residential, services, pedestrian, footway)
otherLen	Total distance (km) of other road types in OpenStreetMap

TABLE 3. Testing results for the Bayesian network

Testing set	Beijing	Shanghai	Chengdu	Shenzhen	Combined
Number of test instances	200	200	200	200	800
F-measure: urban sprawl area	0.818	0.814	0.838	0.733	0.799
F-measure: undeveloped area	0.778	0.732	0.865	0.727	0.780
Overall accuracy	80%	78%	85%	73%	79%

TABLE 4. Confusion matrix for the experiment using combined testing set

Actual/Predicted	Urban sprawl	Undeveloped area	sum
Urban sprawl	334	80	414
Undeveloped area	88	298	386
Overall accuracy	79%		

413	List of Figures	
414	1	Illustration of urban sprawl of four Chinese cities from 2000 to 2010 and their
415		corresponding road networks 23
416	1(a)	Beijing 23
417	1(b)	Shanghai 23
418	1(c)	Chengdu 23
419	1(d)	Shenzhen 23
420	2	Conceptual illustration of the access point and the point-to-road distance 24
421	3	Cross-K function for four cities in China. Solid lines are the plots of cross-K
422		function $K(h)$. Dash lines are 0.05 and 0.95 quantiles of the baseline cross-
423		K function of the urban sprawl points and a randomly generated point process
424		estimated from 20 Monte Carlo simulations. 25
425	3(a)	Beijing 25
426	3(b)	Shanghai 25
427	3(c)	Chengdu 25
428	3(d)	Shenzhen 25
429	4	The structure of the Bayesian network learned from data 26

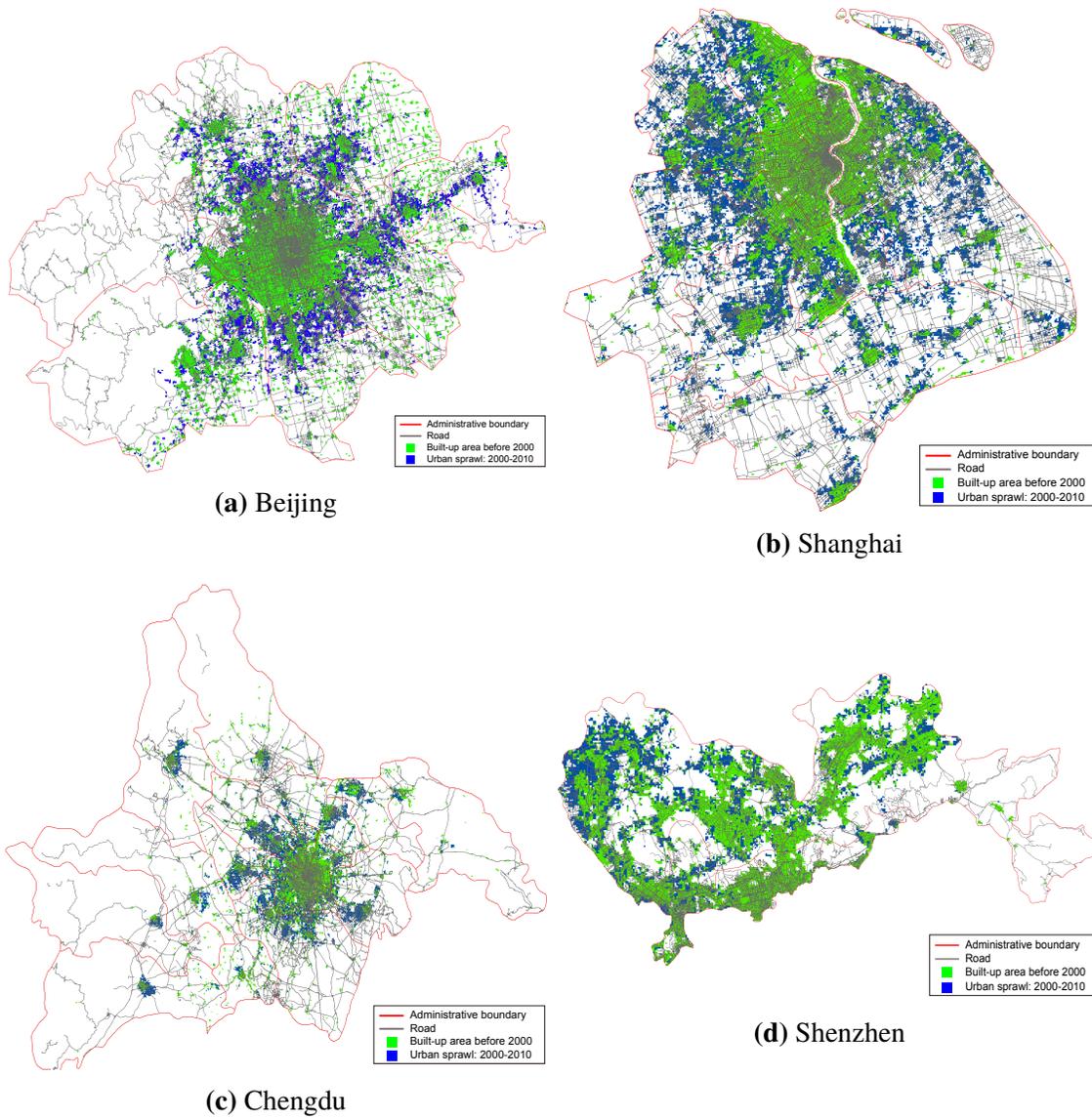


Fig. 1. Illustration of urban sprawl of four Chinese cities from 2000 to 2010 and their corresponding road networks

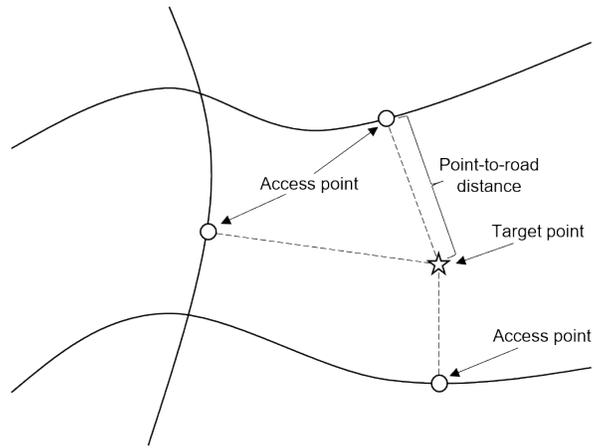
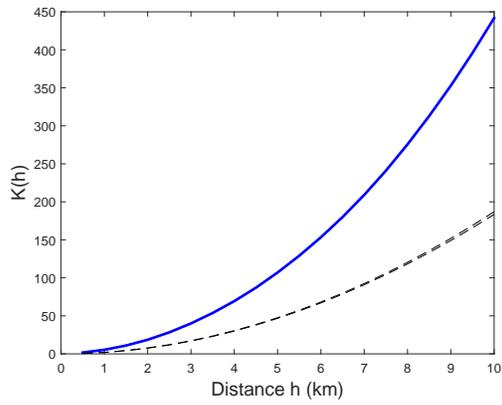
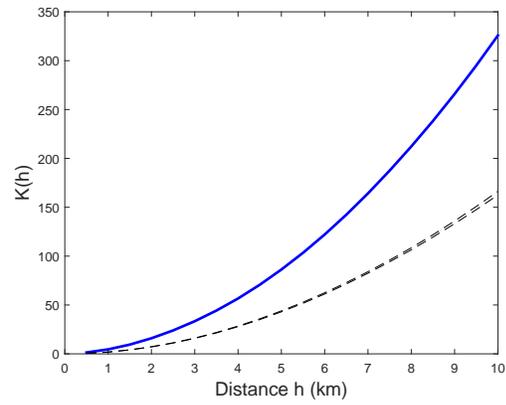


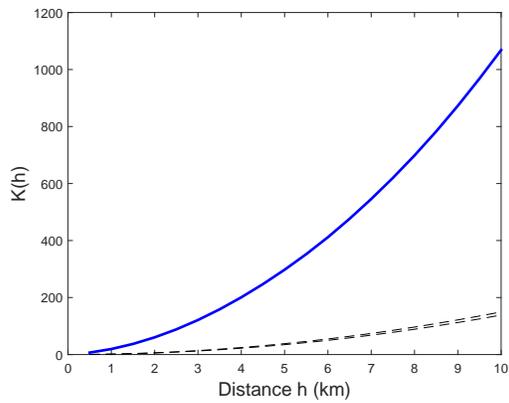
Fig. 2. Conceptual illustration of the access point and the point-to-road distance



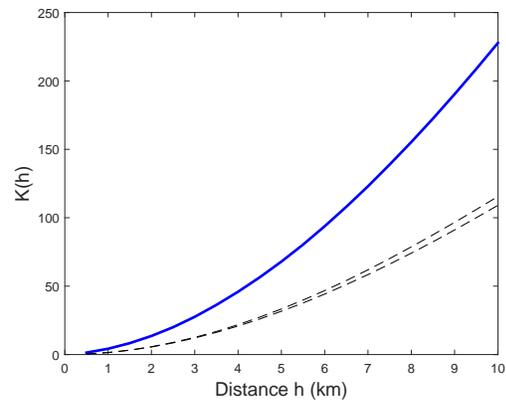
(a) Beijing



(b) Shanghai



(c) Chengdu



(d) Shenzhen

Fig. 3. Cross-K function for four cities in China. Solid lines are the plots of cross-K function $K(h)$. Dash lines are 0.05 and 0.95 quantiles of the baseline cross-K function of the urban sprawl points and a randomly generated point process estimated from 20 Monte Carlo simulations.

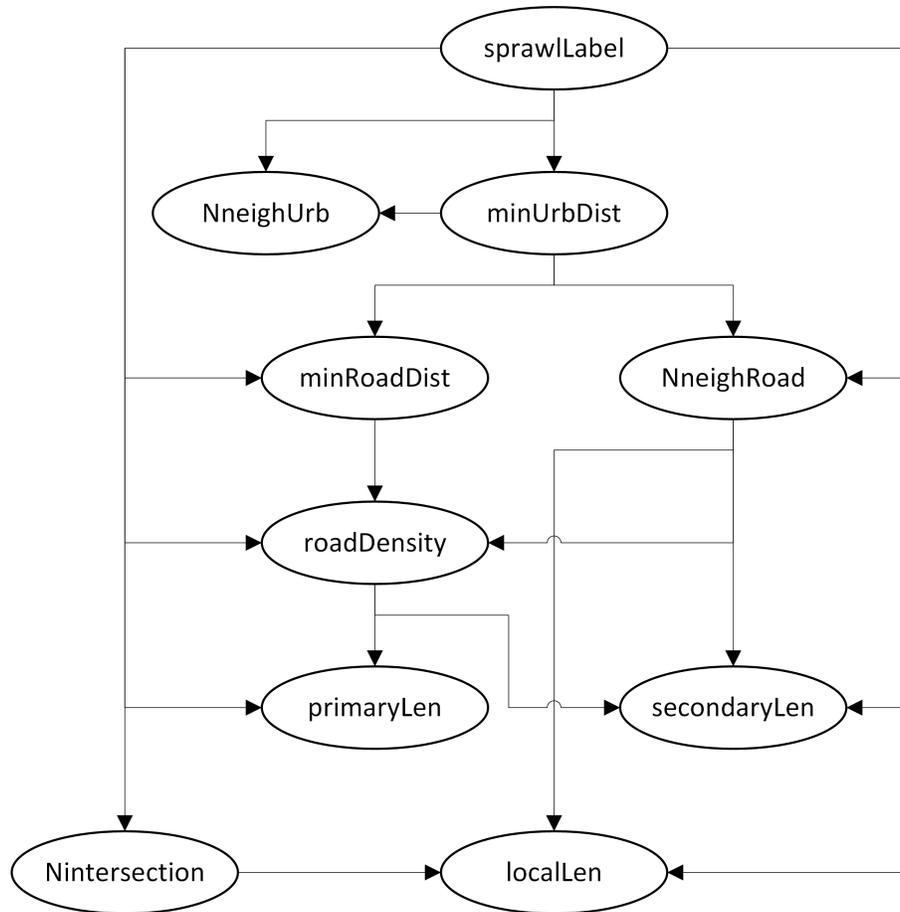


Fig. 4. The structure of the Bayesian network learned from data