

Forecasting Burglary Risk in Small Areas via Network Analysis of City Streets

Maria Mahfoud¹, Sandjai Bhulai², Rob van der Mei¹, Dimitry Erkin¹, and Elenna Dugundji¹

¹ CWI National Research Institute for Mathematics and Computer Science

² Vrije Universiteit Amsterdam, Faculty of Science, Department of Mathematics

Email: M.Mahfoud@cw.nl, S.Bhulai@vu.nl, R.D.van.der.Mei@cw.nl, Dimitry.Erkin@gmail.com, E.R.Dugundji@vu.nl

Abstract—Predicting residential burglary can benefit from understanding human movement patterns within an urban area. Typically, these movements occur along street networks. To take the characteristics of such networks into account, one can use two measures in the analysis: betweenness and closeness. The former measures the popularity of a particular street segment, while the latter measures the average shortest path length from one node to every other node in the network. In this paper, we study the influence of the city street network on residential burglary by including these measures in our analysis. We show that the measures of the street network help in predicting residential burglary exposing that there is a relationship between conceptions in urban design and crime.

Keywords—predictive analytics; forecasting; street network; betweenness centrality; closeness centrality; residential burglary

I. INTRODUCTION

Residential burglary is a crime with high impact for victims. Substantial academic research has accordingly been dedicated to understanding the process of residential burglary in order to prevent future burglaries. In this attempt, several studies have focused on the role of the urban configuration in shaping crime patterns; this is regarded as one of the fundamental issues in environmental criminology, e.g., [1].

According to [2], environmental criminology is based on three premises. The first premise states that the nature of the immediate environment directly influences criminal behavior, thus a crime is not only reliant on criminogenic individuals, but also on criminogenic elements in the surroundings of a crime. The second premise states that crime is non-randomly distributed in time and space, meaning that crime is always concentrated around opportunities which occur on different moments in a day or week, or different places in a given geographical area. The third premise argues that understanding the criminogenic factors within a targeted environment, and capturing patterns and particular characteristics of that area, can reduce the number of crimes within that area.

Understanding human movement patterns within an urban area is essential for determining crime patterns [3]. These movements occur along a street network consisting of roads and intersections. Throughout the city street network, various places are connected, allowing transportation from one point to the next. Within the network, a street segment can be described as the road, or edge, linking two intersections, or nodes. In their study, [4] found that crime is tightly concentrated around crime hotspots that are located at specific points within the urban area. The urban configuration influences where these hotspots are located, suggesting that it is possible to deal with

a large proportion of crime by focusing on relatively small areas. They found that crime hotspots are characterized by being stable over time, and that the hotspots are influenced by social and contextual characteristics of a specific geographical location. To be able to understand and prevent crime, it is important to examine these very small geographic areas, often as small as addresses of street segments, within the urban area. In an analysis of crime at street segment level, [5] reveal that crime trends at specific street segments were responsible for the overall observed trend in the city, emphasizing the need for understanding the development of crime at street segment level.

In urban studies, betweenness is a measure used to determine popularity or usage potential of a particular street segment for the travel movements made by the resident or ambient population through a street network [6], [7]. In criminology, betweenness represents the collective awareness spaces developed by people, including offenders, during the course of their routine activities. This metric provides a means to represent concepts, such as offender awareness, in empirical analysis [8]. Several studies have been conducted to uncover the effects of betweenness on crime. [8] investigated whether street segments that have a higher user potential measured by the network metric betweenness, have a higher risk of burglary. Also included in their research was the geometry of street segments via a measure of their linearity and different social-demographic covariates. They concluded that betweenness is a highly significant covariate when predicting burglaries at street segment level. In another study conducted by [9], a mathematical model of crime was presented that took the street network into account. The results of this study also show an evident effect of the street network.

In this research, we examine for small urban areas (4-digit postal codes: PC4) what the influence of the city street network is on residential burglary by applying betweenness as well as another centrality measure, closeness. These two centrality measures give different indications of the accessibility of an area and we study whether a more accessible area has a higher risk of residential burglary compared to a less accessible area. For comparison, we consider the same areas defined in our previous research [10]. In this earlier study, we predicted residential burglaries within different postal code areas for the district of Amsterdam-West. We extend the model of our earlier research by including the centrality measures closeness and betweenness as explanatory variables. Furthermore, we investigate which of the two centrality measures gives better outcomes, closeness or betweenness.

This paper is organized as follows. Section II describes

the dataset and the data analysis. Section III provides the methodological framework of this research. The results of the analysis are discussed in Section IV. In Section V, conclusions and recommendations for further research are presented.

II. DATA

The data used for this research is collected from three different data sources. The first dataset is provided by the Dutch Police and ranges from the first of January 2009 to 30 April 2014. The original dataset includes all recorded incidents of residential burglaries in the city of Amsterdam recorded at a monthly level and grouped into grids of 125×125 meters resulting in 94,224 records. Next to residential burglary, the dataset includes a wide range of covariates. These covariates provide information on the geographic information of the grid such as the number of Educational Institutions (EI) in the grid. In addition to these covariates, the data includes also spatial-temporal indicators of the following crime types: violation, mugging, and robbery. These spatial-temporal indicators measure the number of times a crime type happened within a given grid cell for a given time lag. The second dataset is obtained from the Statistics Netherlands (CBS) and includes various demographic and socio-economic covariates such as the average monthly income. This data is provided on a six alphanumeric postal code level where the first two digits indicate a region and a city, the second two digits indicate a neighborhood and the last two letters indicate a range of house numbers usually a street or a segment of a street. The third dataset is an internal dataset containing different centrality measures calculated on the street network of Amsterdam.

As this research focuses on explaining and predicting residential burglaries at the four-digit postal code level (PC4), the data should be aggregated at this level. Before aggregating the data we perform some pre-processing steps. First, we check the crime records for missing postal codes: if the postal code is missing then all linked data from CBS and the street network will be missing. We observed that 309 of the total 1,812 grid cells had a missing postal code (PC6). Some of these grid cells (34) were subsequently updated manually; other grid cells referred to industrial areas, bodies of water, railroads, grasslands, and highways. As a double check, we also confirmed whether there were residential burglaries in the remaining grid cells with missing postal codes; in our case, there were indeed none. These grid cells were further removed from the dataset and the data were aggregated based on PC4 conditioning on the district as some postal codes (PC4) can cover different police districts. Discrete covariates were aggregated by taking the sum of the covariate on all PC6. For continuous covariates, this was done by taking the average on all PC6. Exploring the data is done in a similar way as discussed in [10], where an extensive data analysis is applied to the crime data and the CBS data. To analyze this data we extend the final set of covariates by the different centrality measures and repeat the same step again. The dataset was assessed for outliers and collinearity. The presence of outliers was graphically assessed by the Cleveland dot plot and analytically by the Local Outlier Factor (LOF) with 10 neighbors and a threshold of 1.3. Results of this analysis show that the training data exhibits a percentage of outliers of 7.6. The majority of these occurred in December and January. Due to the high percentage of outliers in the training set, we decided

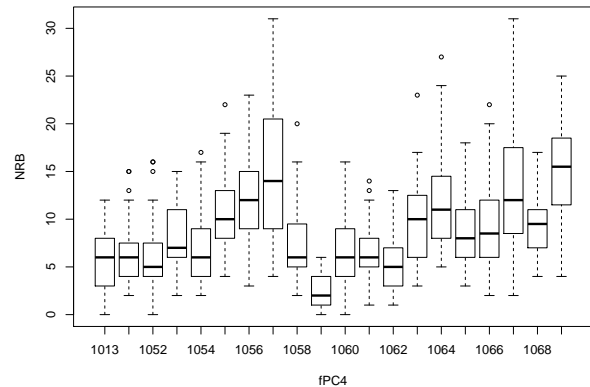


Figure 1. Boxplot of the number of burglaries conditional on the postal code indicating heterogeneity of variance in the number of burglaries within the different postal codes.

to apply the analysis initially without outliers then apply the analysis with the outliers.

The collinearity was assessed by the calculating the variance inflation factor values (VIF) that measures the amount by which the variance of a parameter estimator is increased due to collinearity with other covariates rather than being orthogonal, e.g., [11]. A VIF threshold of 2 is used to assess collinearity [10]. This analysis results in the following set of covariates: the temporal covariate MONTH; the number of educational institutions (sEI), the number of restaurants (sRET), percentage of single-person households (aSH), the number of persons that generate income (sNPI), the total observed mugging incidents in the grid and its direct neighborhood in the last three months (sMuGL3M) and finally, the average monthly income (aAMI).

Furthermore, the relationship between residential burglaries and the categorical covariates was assessed using conditional box plots. Results show a temporal monthly effect and a spatial postal code effect on the burglaries. The effect of the postal codes on the burglaries is illustrated in Figure 1 where a clear difference in the mean and in the variance of the monthly number of burglaries is observed between the different postal codes.

III. METHODOLOGY

A. Centrality measures

Before discussing the centrality measures, we first need to introduce some important concepts of graph theory. A network represented mathematically by a graph is defined as a finite non-empty set V of vertices connected by edges E . A graph is usually written as $G = (V, E)$ where V is the set of vertices and E represents the set of edges where the number of vertices in G is called the order and the number of its edges is called the size. Two vertices u and v are said to be adjacent if there is an edge that links them together. In this case, u and v are also neighbors of each other. If two edges share one vertex then these edges are called adjacent edges. Using this concept of adjacency between all vertices represented in a matrix form

results in an adjacency matrix that summarizes all information describing a network.

Another concept for understanding centrality measures is the one of paths and shortest paths. Informally, a path is a way of traveling along edges from vertex u to vertex v without repeating any vertices [12]. Formally, a path P in a graph G is a subgraph of G whose vertices form an ordered sequence, such that every consecutive pair of vertices is connected by an edge. A path P is called a $u - v$ path in G if $P = (u = x_0, x_1, \dots, x_j = v)$ s.t. $x_0x_1, x_1x_2, \dots, x_{j-1}x_j$ are all edges of P . The number of edges in a path is called its length. The path $u - v$ with the minimum length is called the shortest path between u and v .

In the context of our analysis, a vertex represents an intersection between streets and an edge is a transport infrastructure supporting movements between the two intersections.

Paths can be considered as the key elements in defining centrality measures. In a transportation network, these centrality measures describe the flow of traffic on each particular edge of the network identifying the most important vertices in it. Some of these centrality measures that we will use in this paper are the closeness (CC) centrality and the betweenness centrality (BC).

Closeness is a very simple centrality measure to calculate. It is a geometric measure where the importance of a vertex depends on how many nodes exist at every distance. Closeness centrality can be defined as the average of the shortest path length from one node to every other node in the network and is given by:

$$CC(\nu) = \frac{1}{\sum_{d(u,\nu) < \infty} d(u,\nu)}, \quad (1)$$

where $d(u,\nu)$ is the distance between u and ν . Informally, closeness centrality measures how long it will take to spread information from node ν to all other nodes in the network and it is used to identify influential nodes in the network.

The closeness of an edge $u - v$ can be calculated by taking the average closeness values of the nodes u and v .

The betweenness centrality BC is a path-based measure that can be used to identify highly influential nodes in the flow through the network. Given a specific node ν , the intuition behind betweenness is to measure the probability that a random shortest path will pass through ν . Formally, the betweenness of node ν , $BC(\nu)$ is the percentage of shortest paths that include ν and can be calculated as follows:

$$BC(\nu) = \sum_{u \neq w \neq \nu \in V} \frac{\sigma_{u,w}(\nu)}{\sigma_{u,w}}, \quad (2)$$

where $\sigma_{u,w}$ is the total number of shortest paths between node u and w . And $\sigma_{u,w}(\nu)$ is the total number of shortest paths between node u and w that pass through ν . The betweenness of an edge e can be regarded as the degree to which an edge makes other connections possible and can be calculated in the same way by replacing the node ν by an edge e . An edge with high betweenness value forms an important bridge within the network. Removing this edge will severely hamper the flow

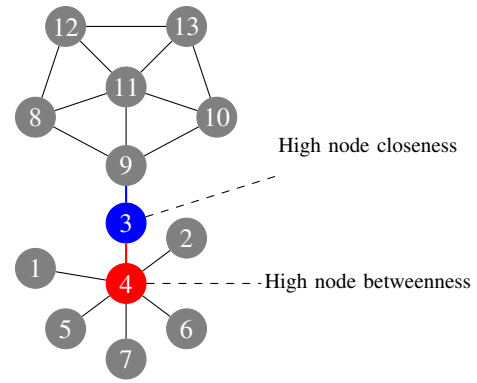


Figure 2. Illustration of high node (edge) betweenness and closeness.

of the network as it partitions the network into two large subnetworks.

High betweenness or closeness values indicate that a vertex or an edge can reach other vertices or edges, respectively, on relatively short paths. An example of a network is illustrated in Figure 2. In this example, node 3 has the highest closeness and node 4 the highest betweenness. The edge connecting the nodes 3 and 9 has the highest closeness within this network. This edge has also the highest betweenness together with the edge connecting the nodes 3 and 4.

In practice, it is almost impossible to calculate the exact betweenness or closeness scores. To make the calculations feasible, one can set a cut-off distance d and allow only paths that are at distances shorter or equal to d .

B. GAMM including centrality measures

In our paper [10], we used generalized additive mixed-effect models with different structures of the random component and showed that the one-way nested model with postal code as a random intercept has the optimal structure of the random component. Further, we showed that using the population as offset captures the most variation in the data. Moreover, the covariates month and the average monthly income seem to be the most important predictors for the number of burglaries within postal codes. In this paper, the optimal model discussed in [10] will be extended by two different centrality measures as covariates. We assess the effect of these centrality measures on explaining and forecasting the number of burglaries within the postal code. This model is given by:

$$\begin{aligned} y_{i,t} &\sim \text{Poisson}(\mu_{i,t}), \\ \mu_{i,t} &= \exp(\text{base}_{i,t} + \text{CM}_i + a_i), \\ a_i &\sim N(0, \sigma_{PC4}^2), \end{aligned} \quad (3)$$

where a_i is a random intercept for the postal code and CM_i represents the closeness CC_i or the betweenness BC_i . The $\text{base}_{i,t}$ is given by:

$$\text{base}_{i,t} = 1 + \text{sEL}_i + \text{sRET}_i + \text{aSH}_i + \text{sNPI}_i + \text{sMugL3M}_{i,t} + f_1(\text{aAMI}_i) + f_2(\text{Month}_t). \quad (4)$$

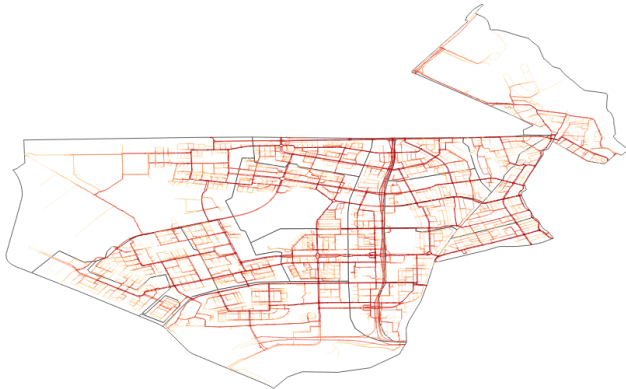


Figure 3. Betweenness of the street segments in Amsterdam West. The betweenness is calculated using the average speed on the street segment and a time threshold of four minutes.

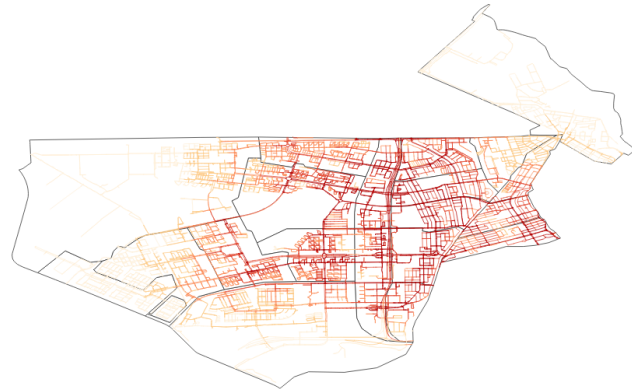


Figure 4. Closeness of the street segments in Amsterdam West. The closeness is calculated using the average speed on the street segment and a time threshold of four minutes.

The models were fitted using the Laplace approximate maximum likelihood [13]. This allows comparing the models based on the Akaike Information Criterion (AIC). All analyses were conducted using the `gamm4` package [14].

To assess the predictive performance of the models, the Root Mean Squared Error (RMSE) is calculated for an out-of-sample test. If $y_{i,t}$ denotes the realization in postal code i and in month t , and $\hat{y}_{i,t}$ denotes the forecast in the same postal code and in the same month, then the forecast error is given by $e_{i,t} = y_{i,t} - \hat{y}_{i,t}$ and the RMSE is given by:

$$RMSE = \sqrt{\frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T e_{i,t}^2}. \quad (5)$$

IV. RESULTS

In this section, we first present the results of the centrality measures. Then, we will discuss the results of the model including these centrality measures as covariates.

As discussed in Section III-A, in practice it is computationally very expensive to calculate the exact betweenness and closeness scores. In general, these can be estimated by setting up a buffer zone using a cut-off distance d and calculating these centrality measures by considering only the paths at a shorter length than d . Using historical data, the average speed per street segment was calculated and five different time cut-offs were used. Segments that are reachable within one to five minutes are used to calculate the centrality measures. Note that these averages make sense because the centrality measures are calculated for the whole city and not for each area separately.

The betweenness and the closeness on the street segment level using a cut-off of four minutes are illustrated in Figure 3 and Figure 4, respectively. The corresponding average betweenness and closeness per area are illustrated in Figure 5 and Figure 6, respectively. Figures 3 and 4 show a wide red road running from top to bottom. This road corresponds with

the A10, which is the ring road of Amsterdam. Figure 3 also shows that the roads with high betweenness correspond to the main access roads within this district. Figure 4 reveals that the roads within the areas situated on the right-hand side of the A10 have a higher closeness in general. This part of the city was built mainly before the Second World War [15] and has a higher density due to enclosed building blocks creating a more finely meshed network of roads when compared to the left-hand side of the ring road. This part was built after the Second World War and is characterized by a lower density due to more open building blocks with an emphasis on more green areas and better enclosure of the residential area via main access roads. The blank areas in the district correspond with green areas, such as parks, lakes and agricultural land.

Adding a centrality measure to the GAMM model results in a better prediction based on the RMSE. The RMSE of the GAMM model without centrality measure was about 4.5519 and as can be seen from Table I, extending the model with the betweenness or the closeness results in a generally lower RMSE. It is noteworthy that the closeness leads to better predictions when using lower thresholds (lower or equal 3 min); see Figures 7 and 8. If the threshold is four minutes or higher including the betweenness in the model results in better predictions. This can be explained by the average time an offender might need to flee from the scene of the crime on a residential street to the nearest main access road. In this case, the closeness describes the number of different routes the offender can take during his flight. Within 4 or 5 minutes, the offender can be traveling on the main access road in order to create as much distance as possible from the crime scene.

The results in the area with the postal code 1067 differ from the other areas. Including the closeness and betweenness does not improve the model, the error on the other hand increased. Taking a closer look at this area revealed that this area mainly consists of green areas with few roads. With less alternative routes available, the closeness gives a higher error.

When looking at the other areas, it is possible to say that the

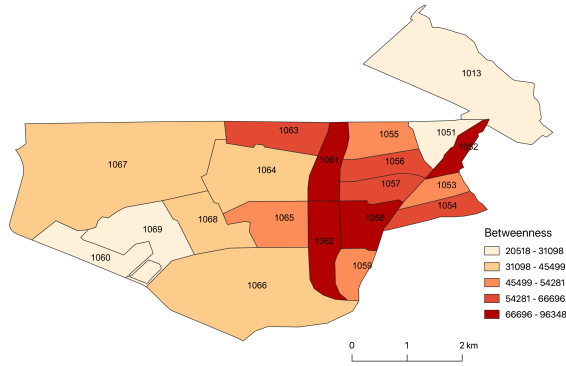


Figure 5. Average betweenness per postal code.

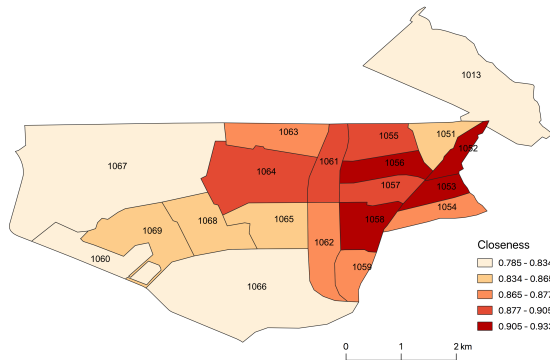


Figure 6. Average closeness per postal code using a threshold of four minutes.

building density influences the effectiveness of the centrality measures on the models. In areas with a lower density, the centrality measures have almost no influence on the outcomes, whereas in the urban areas with a high building density adding the centrality measures to the model improves the outcomes of the model.

Most studies use betweenness as a centrality measure, however, these studies focus on social networks. Given our results, we believe that the closeness is a better centrality measure for modeling crime based on small geographic areas. However, as shown there is a difference in effectiveness of this centrality measure related to the building density of the area.

Table I. ROOT MEAN SQUARED ERROR (RMSE) VALUES FROM FITTING THE GAMM MODEL WITH CLOSENESS AND BETWEENNESS USING DIFFERENT THRESHOLDS.

Model	1 min	2 min	3 min	4 min	5 min
GAMM + CC	4.5297	4.5323	4.5366	5.5437	4.5478
GAMM + BC	4.5562	4.5497	4.5405	4.5279	4.5326

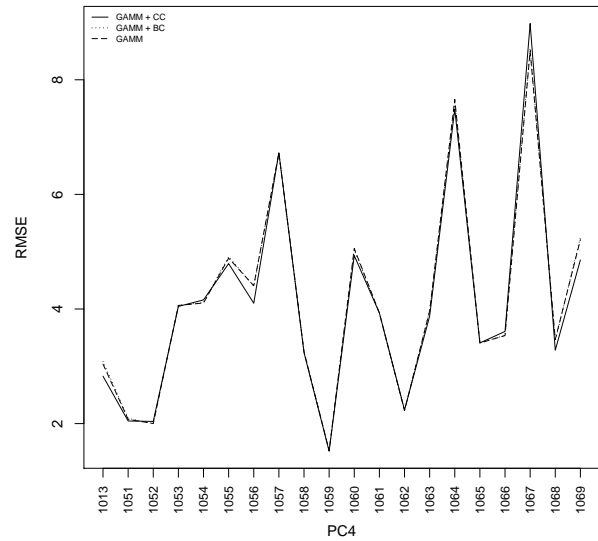


Figure 7. RMSE per PC4 base on an out-of-sample for the GAMM model, the GAMM + CC and the GAMM + BC using a threshold of 1 minute.

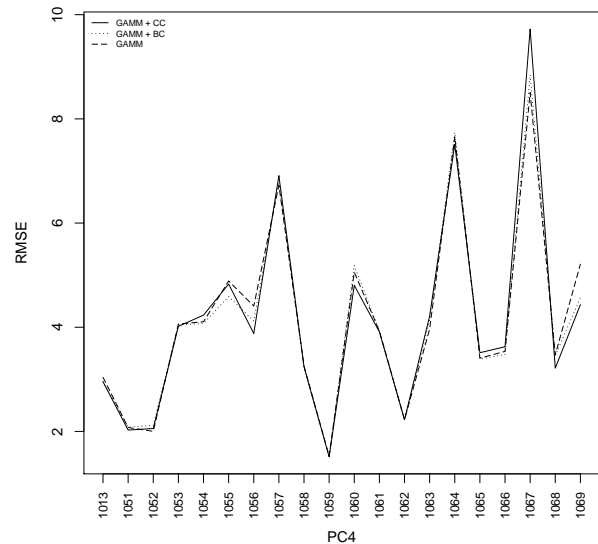


Figure 8. RMSE per PC4 base on an out-of-sample for the GAMM model, the GAMM + CC and the GAMM + BC using a threshold of 4 minutes.

V. CONCLUSION AND FUTURE WORK

During this research, we have tried to determine the influence of accessibility of the street network within small urban areas on residential burglary by applying the centrality measures closeness and betweenness. We have found that adding the centrality measures as a variable to our model has improved the performance of this model as can be concluded from the lower RMSE. Furthermore we have shown that there is a relation between the different conceptions in urban design

over time and residential burglary. Our results show that the pre-world War II neighborhoods suffer from more residential burglary than the neighborhoods built after the Second World War. Also, differences in the performance of the two centrality measures were found. Closeness as a centrality measure gives better predictions when taking into consideration a threshold smaller than 4 minutes. If the threshold is 4 minutes or larger, the betweenness gives better predictions. We can also conclude that the centrality measures perform better when applied to geographic areas with a high density, for example, a city center.

Our study has shown that there is a relationship between the conceptions in urban design and crime. Neighborhoods built under a certain conception of urban design tend to have a higher risk of residential burglary, which can be explained by how the public space is designed. Further research is necessary to confirm this hypothesis.

REFERENCES

- [1] P. J. Brantingham, and P. L. Brantingham, *Environmental criminology*. Sage Publications Beverly Hills, CA, 1981.
- [2] R. Wortley and M. Townsley, *Environmental criminology and crime analysis*. Taylor & Francis, 2016, vol. 18.
- [3] P. Brantingham and P. Brantingham, "Crime pattern theory," in *Environmental criminology and crime analysis*. Willan, 2013, pp. 100–116.
- [4] D. Weisburd, E. R. Groff, and S.-M. Yang, *The criminology of place: Street segments and our understanding of the crime problem*. Oxford University Press, 2012.
- [5] D. Weisburd, S. Bushway, C. Lum, and S.-M. Yang, "Trajectories of crime at places: A longitudinal study of street segments in the city of Seattle," *Criminology*, vol. 42, no. 2, 2004, pp. 283–322.
- [6] L. C. Freeman, "Centrality in social networks conceptual clarification," *Social networks*, vol. 1, no. 3, 1978, pp. 215–239.
- [7] P. Crucitti, V. Latora, and S. Porta, "Centrality measures in spatial networks of urban streets," *Physical Review E*, vol. 73, no. 3, 2006, p. 036125.
- [8] T. Davies and S. D. Johnson, "Examining the relationship between road structure and Quantitative Criminology," vol. 31, no. 3, 2015, pp. 481–507.
- [9] T. P. Davies and S. R. Bishop, "Modelling patterns of burglary on street networks," *Crime Science*, vol. 2, no. 1, 2013, p. 10.
- [10] M. Mahfoud, S. Bhulai, and R. van der Mei, "Spatio-temporal modeling for residential burglary," in *Proceedings of the 6th International Conference on Data Analytics*. IARIA, 2018, pp. 59–64.
- [11] D. Liao and R. Valliant, "Variance inflation factors in the analysis of complex survey data," *Survey Methodology*, vol. 38, no. 1, 2012, pp. 53–62.
- [12] A. Benjamin, G. Chartrand, and P. Zhang, *The fascinating world of graph theory*. Princeton University Press, 2015.
- [13] R. H. Baayen, D. J. Davidson, and D. M. Bates, "Mixed-effects modeling with crossed random effects for subjects and items," *Journal of memory and language*, vol. 59, no. 4, 2008, pp. 390–412.
- [14] S. Wood and F. Scheipl, *GAMM4: Generalized additive mixed models using mgcv and lme4*, 2014, r package version 0.2-3. [Online]. Available: <http://CRAN.R-project.org/package=gamm4>
- [15] G. Amsterdam. *De groei van Amsterdam*. [Online]. Available: <https://maps.amsterdam.nl/bouwjaar/?LANG=nl> (2018)