



A new integrated GIS-based analysis to detect hotspots: A case study of the city of Sherbrooke



Homayoun Harirforoush*, Lynda Bellalite

Department of Applied Geomatics, University of Sherbrooke, 2500 boul. de l'Université Sherbrooke, J1K 2R1 Quebec, Canada

ARTICLE INFO

Article history:

Received 26 February 2016

Received in revised form 4 June 2016

Accepted 11 August 2016

Available online 24 August 2016

Keywords:

Network KDE

Hotspot

Crash rate

Spatial analysis

GIS

Exposure data

ABSTRACT

This paper proposes a two-step integrated method for identifying traffic accident (TA) hotspots on a roadway network. The first step includes a spatial analysis method called network kernel density estimation (KDE). The second step is a network screening method using the critical crash rate, which is described in the Highway Safety Manual (HSM). The method was examined by using three years of TAs (2011–2013) in Sherbrooke, Canada. The network KDE uses TAs to graphically display sites with a high crash density. Two different crash patterns were used for identifying these locations: (1) a crash pattern that includes three-year aggregated crash data, and (2) a crash pattern that involves three-year merged crash data. The results of the two crash patterns were evaluated based on a prediction accuracy index (PAI). It was found that the results obtained from the merged crash data outperformed the other. On the other hand, crash clustering in a site does not imply a site is hotspot and it is better to be tested by other factors. High crash density locations were then tested by the critical crash rate, which helps to create an accurate comparison of sites. The importance of the critical crash rate is that it takes several factors into account such as the amount of exposure, the type of intersection, variance in crash data, etc. We realized that the hotspots determined using the two methods reflect very problematic locations and filter out the locations that do not have a problem. This approach could help transportation authorities and safety specialists to identify and prioritize sites that require more safety attention.

© 2016 Elsevier Ltd. All rights reserved.

1. Introduction

Understanding when and where traffic accidents (TAs) occur on a road network is one of the most important concerns of transport authorities. The identification of high-risk locations on a road network can improve TA reduction efforts. Road safety is certainly one of the greatest concerns facing Canada, and indeed the world. The World Health Organization (WHO) reported that approximately 1.24 million people die every year on the world's roadways, which is the eighth leading cause of death worldwide. It is also estimated that the cost resulting from road TAs runs up to billions of dollars (Toroyan, 2013). To improve road safety and create a safe driving environment, it is crucial to identify road segments where the average crash density is relatively high compared to other segments of the network. These locations are known as hotspots, which are characterized by a high crash concentration, relative to the distribution of TAs across the whole study area (e.g., county, state,

municipal district, or downtown) (Chainy and Ratcliffe, 2013). Previous studies show that the occurrences of TAs are rarely random in time and space. Their occurrence, in reality, is determined by some important factors, including traffic volume, weather conditions, geometric design, etc.

The Poisson distribution is appropriate for the analysis of TAs occurring within a given year and at individual sites. It deals with the occurrence of some random events during a given interval (Ayyub and McCuen, 2011). In fact, it is assumed that the variance and the mean are equal, while if the variance is larger than the mean then the assumption is wrong. Therefore, the negative binomial distribution as a generalization of the Poisson distribution can be used in the case of a non-random distribution.

The spatial analysis of point events, referred to as point pattern analysis, has been widely used for analysing the distribution of a set of points (i.e., crashes) on a surface (i.e., network) (Ervin, 2015). The PPA method is divided into two main categories (Bailey and Gatrell, 1995; O'Sullivan and Unwin, 2014): (1) density-based methods (called first-order properties) and (2) distance-based methods (called second-order properties). The first group measures the intensity of point events based on the density in a region. It includes techniques such as kernel density estimation (KDE) and

* Corresponding author.

E-mail addresses: homayoun.harirforoush@usherbrooke.ca (H. Harirforoush), lynda.bellalite@usherbrooke.ca (L. Bellalite).

quadrat analysis. The second group measures the spatial dependence of point events based on the distance of points from each other. This group includes methods such as nearest neighbour distances, K-functions and Moran's I (Steenberghen et al., 2010; Xie and Yan, 2008). For instance, previous studies applied K-function methods to analyse the distribution of point patterns in a network. Their results showed that there is a significant chance of overestimating clusters of point patterns (Lu and Chen, 2007; Yamada and Thill, 2004).

In TA analysis, KDE is one of the most popular density-based methods and has been widely used for detecting dangerous road segments (Silverman, 1986; Xie and Yan, 2013). The purpose of KDE is to create a smooth density surface of point events over space by counting the number of crashes at each location as a density estimation. For each point event in the network, a kernel density surface is defined, and the density value is highest at its centre and decreases as it moves away from the centre (Silverman, 1986; Vemulapalli, 2015). This method is suitable for visualizing the crash data as a continuous surface (Chainey and Ratcliffe, 2013).

The KDE can be classified into two categories: planar KDE and network KDE. The first approach uses the Euclidean distance for estimating the density of point events. A review of previous studies indicates that planar KDE has been widely used in TA analyses, such as for hotspots (Flahaut et al., 2003; Sabel et al., 2005), highway TAs (Erdogan et al., 2008), vehicle-wildlife crashes (Krisp and Durot, 2007), fatal automobile crashes (Oris, 2011), and weather-related accidents (Khan et al., 2008). The planar KDE method estimates the crash density in a cell moving across a two-dimensional homogenous space. Crashes are weighted based on the Euclidean distance, where crashes closer to the centre contribute a higher value (Chainey and Ratcliffe, 2013).

However, this method has significant limitations: (1) in the case of crashes occurred inside a roadway network, the assumption of two-dimensional space does not hold (Xie and Yan, 2008), and (2) the density of the road network is ignored. Some cells might have the same density values, while they may include different numbers of road sections. The result (real density value) is therefore biased (Tao et al., 2011). Different studies have tried to overcome these limitations by extending the planar method into network space. In an early phase, researchers compared planar KDE and network KDE, which showed the advantages of using network KDE (Borruso, 2008; Kuo et al., 2011; Larsen, 2010; Steenberghen et al., 2004; Yamada and Thill, 2004). Recently, several studies employed the network KDE method to estimate hotspot locations (Mohaymany et al., 2013; Timothée et al., 2010; Steenberghen et al., 2010; Sugihara et al., 2010; Vemulapalli, 2015; Xie and Yan, 2013; Young and Park, 2014). Nonetheless, a major drawback of both the planar and network KDE is that there is no specific statistical approach for testing hotspots (Xie and Yan, 2008; Yao et al., 2015; Nie et al., 2015).

In addition, the network KDE method is based on the network distance and measures the density of crashes along a one-dimensional space (Timothée et al., 2010). For this purpose, an ArcGIS-based toolbox known as SANET was developed by a group of researchers (Okabe et al., 2006). They proposed a network-constrained kernel called "equal split discontinuous at nodes" to calculate the density of point events along a network (Okabe et al., 2009).

A review of previous studies shows that the length of a road segment has a great impact on the results. Some studies, like Miaou (1994), used unequal-length roadway segments and others, like (Erdogan et al., 2008; Yamada and Thill, 2010), used equal-length roadway segments (Nie et al., 2015). Their results show that the hotspot locations are varied using the unequal-length roadway segments. This variation in spatial analysis is known as the modifiable

area unit problem (MAUP). Therefore, this study used equal-length roadway segments for the network KDE calculation.

In recent years, the network KDE method has been widely used in road safety studies to detect hazardous accident locations (Larsen, 2010; Mohaymany et al., 2013; Nie et al., 2015; Oris, 2011; Vemulapalli, 2015; Xie and Yan, 2008; Young and Park, 2014; Loo and Yao, 2013). However, many of these studies neglected two main factors in their crash analysis. Firstly, they used the long-term aggregated crash data (for instance, three or more consecutive years) at each site regardless of whether these crashes occurred continuously at a particular location. In fact, they did not consider if the high crash density at a specific site was due to chance or some continuing problem (like geometry design problem at a given location). In general, crash frequency at a site varies from year to year around a steady mean value. However, due to the random variation of TA occurrence, the extreme case selected in one year may have a lower frequency in the next year. Secondly, they only used the raw crash counts for their analysis, which may result in misleading information about the most appropriate sites for treatment (FHWA, 2011). In fact, these studies failed to take into account other safety parameters such as exposure data. (Erdogan et al., 2008; Larsen, 2010; Mohaymany et al., 2013; Vemulapalli, 2015). Naturally, the more factors being considered such as exposure can improve the model.

Accordingly, we propose a method to overcome these limitations. First, the potential hotspot locations must be identified. These locations have the following properties: (1) Crash density is relatively high at a particular site, and (2) Crashes occur consecutively (at least three successive years) at a particular site. Then, the obtained potential hotspot locations must be examined by other safety parameters such as exposure data, which provides an appropriate comparison of sites. Traffic volume is the most common type of exposure data and is often used at segments and intersections. It provides a common metric to the collision data, hence sites can be compared more appropriately (FHWA, 2011). The critical crash rate method uses exposure data to define relative safety compared to other similar intersections or segments (FHWA, 2011). The method provides the expected crash rate for sites with similar characteristics (i.e., same traffic control system, same traffic volume, same number of legs). The critical crash rate method compares the crash rate at a site (road segment or intersection) with a critical crash rate specific to each site (AASHTO, 2010). If the crash rate at a site exceeds the critical crash rate, the site is then considered to be deviant due to unfavourable characteristics of the site (FHWA, 2012; PIARC, 2003).

The aim of this paper is to identify hotspot locations based on a historical crash dataset in order to improve road safety. Therefore, we integrated the network KDE and critical crash rate methods to identify crash hotspot locations. The KDE is one of the methods for examining the first-order effects of a spatial process, while the critical crash rate is from the HSM's network screening measure. Identifying hazardous locations helps transport authorities to improve road safety and to focus on the reasons behind the occurrence of these accidents.

The remainder of the study is organized as follows. Section two describes the study databases. The methods are described in section three. Section four presents the results obtained from applying the network KDE and HSM methods to select the most appropriate hotspot road locations. Finally, section five presents the discussion and conclusions of the study.

2. Study data

This study focuses on the city of Sherbrooke, in southern Quebec, in east-central Canada. Sherbrooke covers an area of approximately

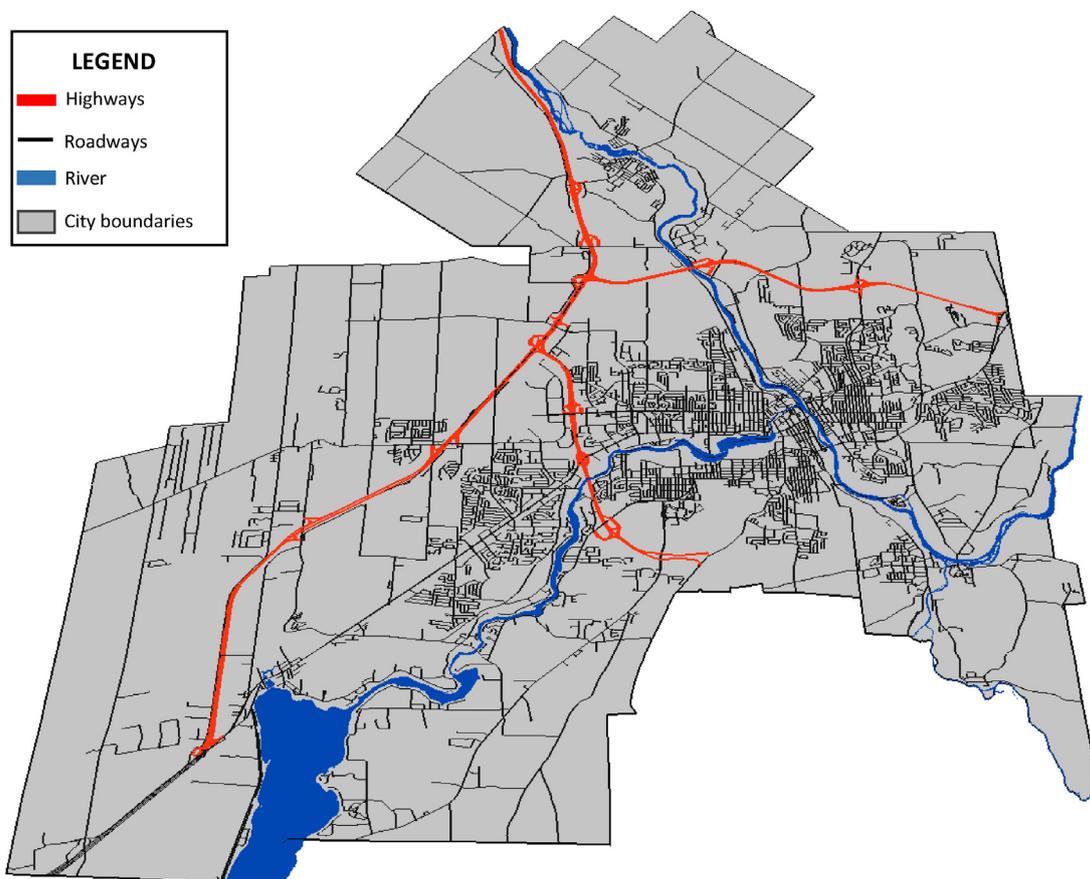


Fig. 1. The study area (Sherbrooke, Canada).

353.5 km², and its population in 2011 was about 154,600 (about 0.4% of Canada). The study only focuses on urban areas and considers all types of roadways (i.e., local, collectors, and arterial roads) within the city boundary, excluding highways. This study conducted a safety network screening of the entire network. For the study, three different databases were used from various sources. First, a roadway network base map was obtained from the “Ville de Sherbrooke”. The map was provided in a GIS shape file format, which includes roadway specifications such as shape length (segment length), road type, and speed limits. The shape file contains 8327 segments. Fig. 1 shows the study area.

Secondly, the three-year (2011–2013) TA database provided by the Société de l'assurance automobile du Québec (SAAQ). During the study period, a total of 6926 collisions were recorded on Sherbrooke's roadways. The TA database was provided in an Excel format and contains significant crash parameters, such as the date and time of an accident, accident location, age, sex, vehicle type, weather conditions, etc. These crashes were then converted into a GIS shape file and mapped using ArcGIS based on their latitude and longitude. This study only considers vehicle-to-vehicle crashes in the safety analysis, and other types of crashes such as pedestrian and cycling crashes are outside the scope of this study. Thirdly, in this study we used traffic volume data provided from different sources: (1) Quebec's Ministry of Transportation recorded the average annual daily traffic (AADT) and average daily traffic (ADT) from permanent and non-permanent stations for 44 sites throughout the city; (2) The City of Sherbrooke recorded intersection volume counts for about 1100 sites within Sherbrooke dating back to 2000. They provided ADT information for each intersection. They also provided ADT information for about 300 segments within the city for the three most recent years (2010–2013); and (3) Due to the

lack of information, a manual traffic volume count was conducted to determine the ADT. The survey was conducted for 6 h (intervals of 15 min) during the peak hours (7:00–9:00; 11:30–13:30; 16:00–18:00). Then, the total 6-h volume was converted to the 24-h volume.

Since the HSM network screening relies on the critical crash rate requires the traffic volume data for each data set (reference population). Therefore, we used the estimated traffic volume information for each corresponding site.

It should be noted that the duration of study was sufficient to limit changes in road traffic conditions, traffic volume and crash data fluctuate (Bil et al., 2013; Flahaut et al., 2003).

3. Methods

Several techniques have been suggested by researchers to identify hotspot locations. In this study, first, we applied the network KDE because it is the most common method to find a significant cluster of crashes. Then, the crash density surfaces of three-consecutive years were merged to find the potential hotspot locations. Finally, the critical crash rate method was used to identify hotspot locations throughout Sherbrooke's roadway network.

3.1. Network KDE

As stated earlier, many studies have recently used the network KDE method developed by Okabe et al. (2009) to examine the spatial correlation of point events in a network. In this study, we used the network KDE method for estimating the density of TAs on Sherbrooke's roadway network. As described by Okabe and Sugihara (2012), an unbiased kernel function must be used to

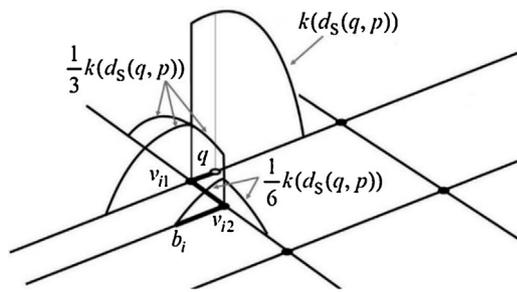


Fig. 2. Simplified example of an equal split discontinuous kernel function (modified from Okabe and Sugihara 2013).

avoid false conclusions. Hence, they formulated a kernel function named the “equal split discontinuous kernel density function”. In this approach, the network kernel function is defined for two cases: (1) kernel centre q does not coincide with a node, and (2) kernel centre q coincides with a node. In the first case, the function is defined as follows (Okabe et al., 2008):

$$K_q(p) = \begin{cases} \frac{k(d_s(q, p))}{(n_{i1} - 1)(n_{i2} - 1) \dots (n_{ik-1} - 1)} & \text{for } d_s(q, v_{ik-1}) \leq d_s(q, p) < d_s(q, v_{ik}), \\ 0 & \text{for } d_s(q, p) \geq h. \end{cases} \quad (1)$$

Where $k(x)$ is a base kernel function, y is the kernel centre, d is the shortest path distance between y and x , h is the bandwidth, and n is the degree of the node.

In this case, as shown in Fig. 2, the value of the kernel function is the same as the base kernel function as long as the kernel centre located on the shortest path does not meet a node ($0 \leq d_s(q, p) < d_s(q, v_{i1})$). Otherwise, when it reaches a node, the value of the kernel function is then equally divided into the degree of the node (Eq. (1)). This process continues until the kernel centre reaches the boundary point.

In the second case, the value of the kernel function at a vertex, say v_{i1} , divided by n_{i1} is assigned to the road links. The function for case two is defined as follows:

$$k_q(p) = \begin{cases} \frac{2k(d_s(q, p))}{n_{i1}(n_{i2} - 1) \dots (n_{ik-1} - 1)} & \text{for } d_s(q, v_{ik-1}) \leq d_s(q, p) < d_s(q, v_{ik}), \\ 0 & \text{for } d_s(q, p) \geq h. \end{cases} \quad (2)$$

As stated earlier, in the network KDE method, the road segments are divided into equal-sized sub-networks called network cells. According to Furuta et al. (2008), the procedure comprises two steps: first, the network Voronoi diagram divides the network road segments into sub-networks. Then, linear programming is used for adjusting the length of the sub-networks (split them equally). In this study, we ran the network KDE with a 10-m segment length, similar to the one suggested by Xie and Yan (2008) and Nie et al. (2015). For details on the computational process, see Okabe and Sugihara (2012).

In addition, there are several types of kernel functions, such as quadratic, uniform, gaussian, trigonometric, etc., but the results of the network KDE are more dependent on the search bandwidth (Xie and Yan 2008; Mohaymany et al., 2013). Therefore, it is crucial to select an appropriate bandwidth. If the search bandwidth is too large, the density patterns will be too smooth; hence, it would be difficult to differentiate between local hotspot locations. On the other hand, a narrow search bandwidth may produce a very sharp density pattern and only highlight individual hotspot locations. Accordingly, the results of both cases may lead to false conclusions.

In previous studies, researchers used an iterative (trial and error) technique to obtain an optimal search bandwidth (Mohaymany et al., 2013; Plug et al., 2011; Silverman, 1986; Xie and Yan, 2008; Young and Park, 2014). In this study, we followed their sugges-

tion and tested search bandwidths from 50 to 500 m. As shown in Fig. 3, the number of clusters (hotspots) gradually becomes larger when the search bandwidth increases from 50 to 500 m. It appears that wider search bandwidths (300 and 500 m) may produce a very smooth density pattern and hazardous locations are mixed with their neighbours. Hence, it would be difficult to identify accurate hazardous locations. The large bandwidths also produce unrealistic density clusters that their density ranges are dramatically high. The small search bandwidth (50 m), on the other hand, produces a very spiky density pattern with many isolated individual clusters (shown in red). In this study, a search bandwidth of 100 m was selected for the analysis of high-density TA locations using the network KDE method.

This study also used the Natural Breaks (Jenks) classification technique to determine the best arrangement of crash densities into different classes. This technique groups similar values and maximizes the differences between classes. In other words, the boundaries between classes are set where there are relatively big differences in the data values (ESRI, 2015). In this study, the crash densities were classified into five classes.

3.2. Threshold selection

The network KDE is an appropriate approach that provides an overview of the overall distribution of collisions. However, its analytical capabilities usually end at this point. On the other hand, our results show that there are many segments with a density of zero or one. Therefore, a threshold value was selected to represent abnormal segments (hazardous locations) in the normal pattern. In this study, the threshold value was set to three standard deviations from the mean value (Larsen, 2010).

3.3. Potential hotspots

To identify potential hotspots, the network KDE must first be implemented to generate a collision density map for each year (i.e., 2011, 2012, and 2013). The determined threshold must then be applied. To extract potential hotspots, three density maps higher than the threshold should be merged and their joint spots should be selected for further processing. As Fig. 4 shows, the top maps are the collision density maps higher than the determined threshold for the years 2011, 2012 and 2013, and the bottom map shows the result of merging three consecutive years. The selected site in the bottom map of Fig. 4 shows a potential hotspot location (College/Belvedere intersection).

3.4. Prediction accuracy index

The prediction accuracy index (PAI) was first developed by Chainey et al. (2008). It was created to evaluate the performance of two methods. This approach was initially developed in crime mapping (Chainey et al., 2008; Van Patten et al., 2009; Hart and Zandbergen, 2012) and has been recently used in road safety studies. According to Thakali et al. (2015), PAI is a ratio of the proportion

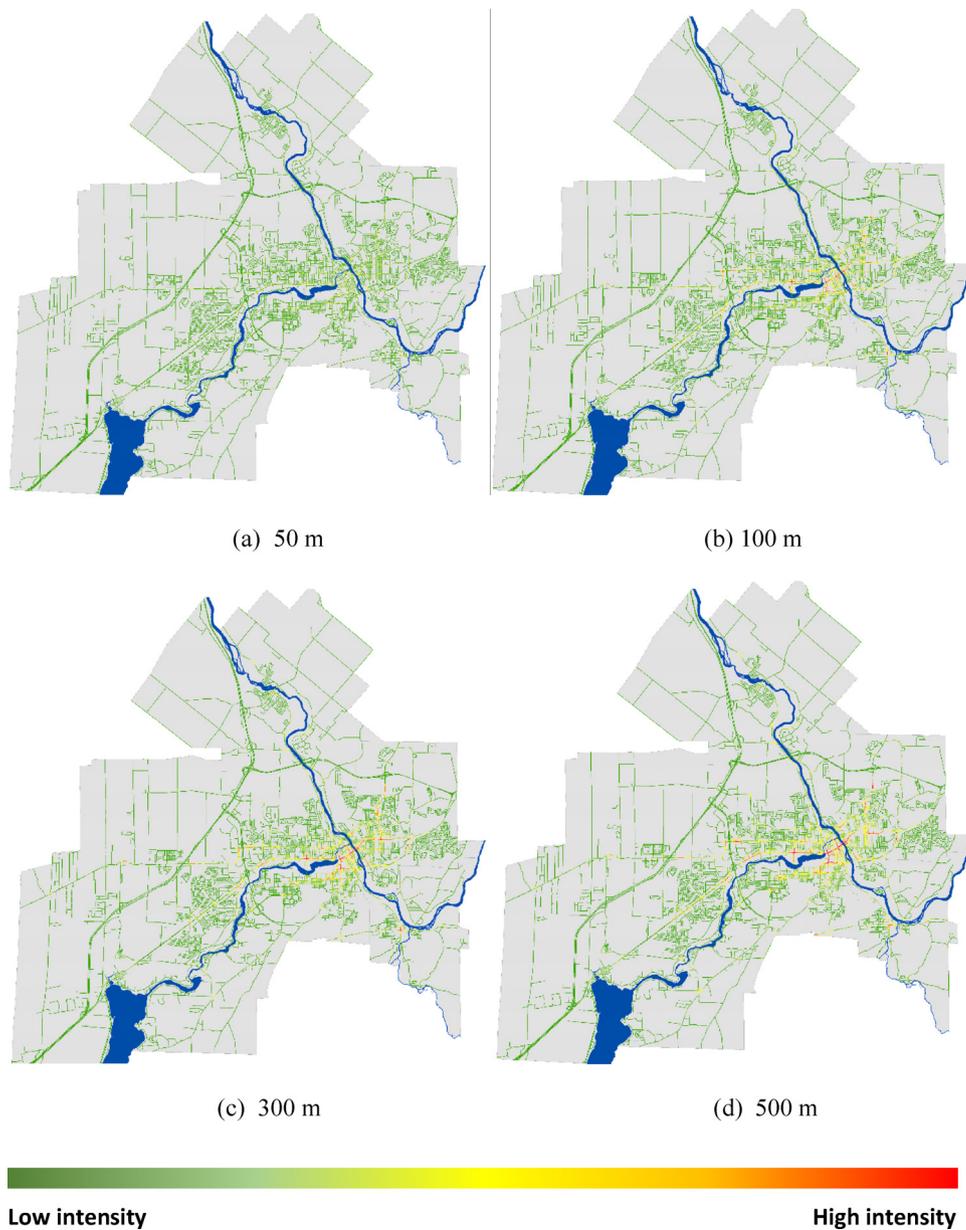


Fig. 3. Different search bandwidths (50,100, 300, and 500 m) and their impact on the density surface.

of crashes occurring within the identified hotspot to the proportion of area covered by it. The formula is given as follows (Thakali et al., 2015):

$$PAI = \frac{\text{Proportion of hotspot area}}{\text{study area}} = \frac{\frac{n}{N} \times 100}{\frac{m}{M} \times 100} \quad (3)$$

where n is the number of crashes in the hotspots, N is the total number of crashes, m is the length of highway section in the hotspots or area covered, and M is the total length of the highway section or total area covered. It should be noted that, a larger PAI value means better ability of a method to locate high potential crashes in an area.

3.5. Critical crash rate

Traffic engineers commonly use the critical crash rate as the HSM screening method. This method compares the crash rate at a location with the critical crash rate of sites with similar characteristics. The critical crash rate is a function of the average crash

rate of a reference group related to the traffic volume, the site, and a desired level of confidence (FHWA, 2011). If the value of the crash rate is higher than the critical crash rate, the division is due to unfavourable characteristics of the road segment or intersection (FHWA, 2004).

This method requires a reference population of a sufficient size and is useful when a large number of sites (reference population) is available. The HSM critical crash rate varies among intersections and segments. A reference population for intersections could be divided based on operational (i.e., signalized intersection or stop-controlled intersections) or geometric (i.e., three-leg or four-leg) characteristics. The accident characteristics of signalized and stop-controlled intersections are quite different, and they should not be mixed in the same population (Dunn et al., 2015). A reference population for segments could be divided based on roadside characteristics (i.e., lanes).

In this paper, the reference population was selected based on the potential hotspot locations (i.e., the selected results obtained

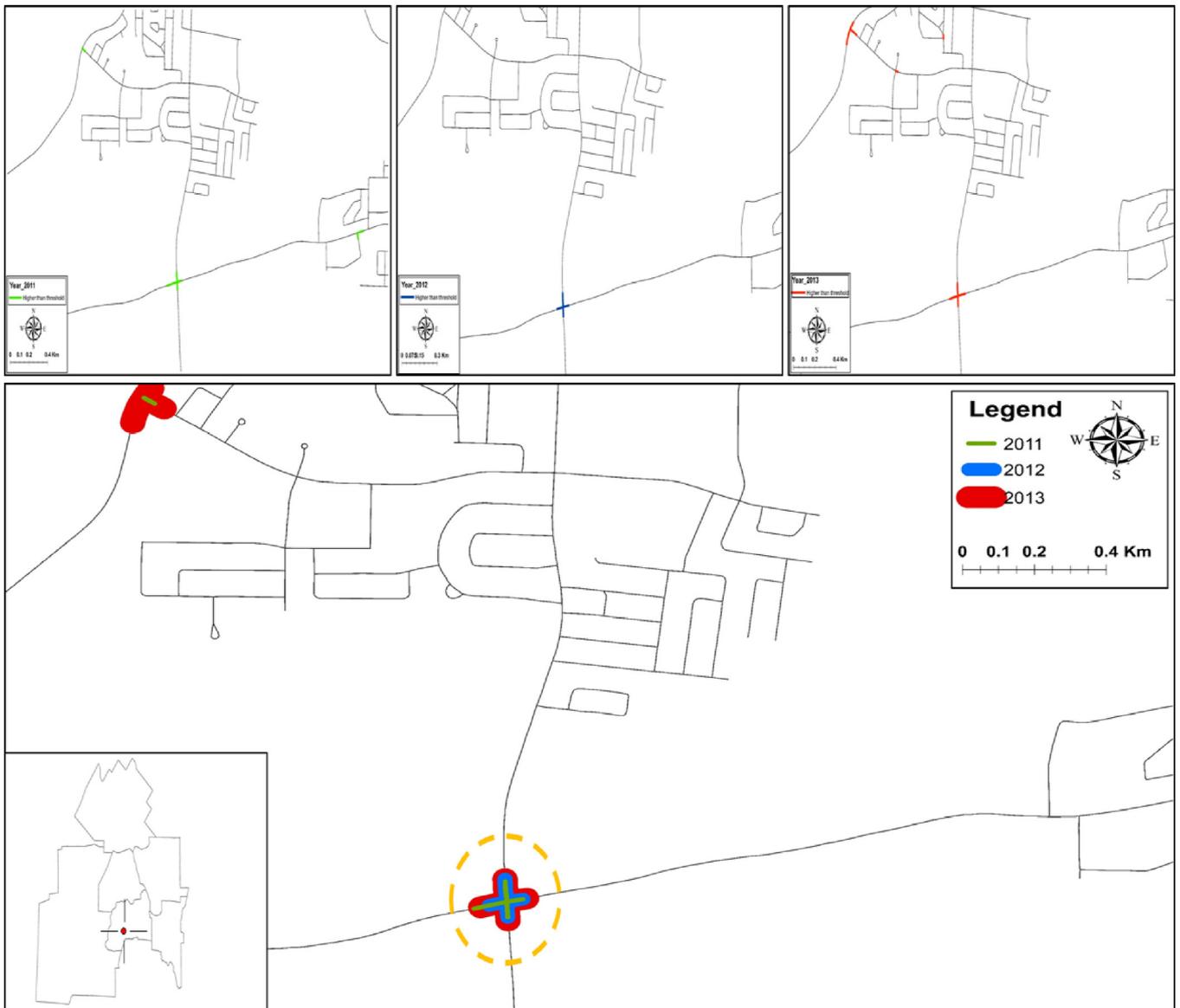


Fig. 4. Merging three collision density maps (2011, 2012, and 2013).

from the network KDE). However, we selected more intersections in each population to obtain better results. The reference population for intersections (both signalized and stop controlled) was split into three versus four-legged intersections. Therefore, for calculating the intersection critical crash rate, we selected 75 signalized intersections (including 55 four-legged and 20 three-legged) and 65 stop-controlled intersections (including 40 four-legged and 25 three-legged). On the other hand, for calculating the segment critical crash rate, we selected 24 road segments (mixed lane segments due to a lack of segments).

To obtain hazardous sites, first, the observed crash rate should be calculated for each intersection and segment using Eqs. (4) and (5) as follows (AASHTO, 2010):

$$R_i = \frac{N_{\text{observed}, i(\text{total})}}{MEV_i} \quad (4)$$

$$R_i = \frac{N_{\text{observed}, i(\text{total})}}{MVMT_i} \quad (5)$$

where R_i is the observed crash rate at intersection i , $N_{\text{observed}, i(\text{total})}$ is the total observed crashes at intersection i , MEV_i is million entering

vehicles at intersection i , and $MVMT$ is million vehicle-miles of travel.

The critical crash rate should then be calculated for each intersection and segment using Eqs. (6) and (7) as follows:

$$R_{c,i} = R_a + \left[P \times \sqrt{\frac{R_a}{MEV_i}} \right] + \left[\frac{1}{(2 \times (MEV_i))} \right] \quad (6)$$

$$R_{c,i} = R_a + \left[P \times \sqrt{\frac{R_a}{MVMT_i}} \right] + \left[\frac{1}{(2 \times (MVMT_i))} \right] \quad (7)$$

Where $R_{c,i}$ is the critical crash rate for intersection i , R_a is a weighted average crash rate for the reference population, P is a P-value corresponding to the confidence interval, MEV_i is million entering vehicles at intersection i , and $MVMT$ is million vehicle-miles of travel.

The values for the observed crash rates (i.e., Eqs. (4) and (5)) and the critical crash rates (i.e. Eqs. (6) and (7)) are then compared. If the crash rate value at the site exceeds the corresponding critical crash rate, then it is identified for further analysis. Finally, the differences between the results obtained for the crash rate and the critical crash rate are calculated and arranged in decreasing order (from largest

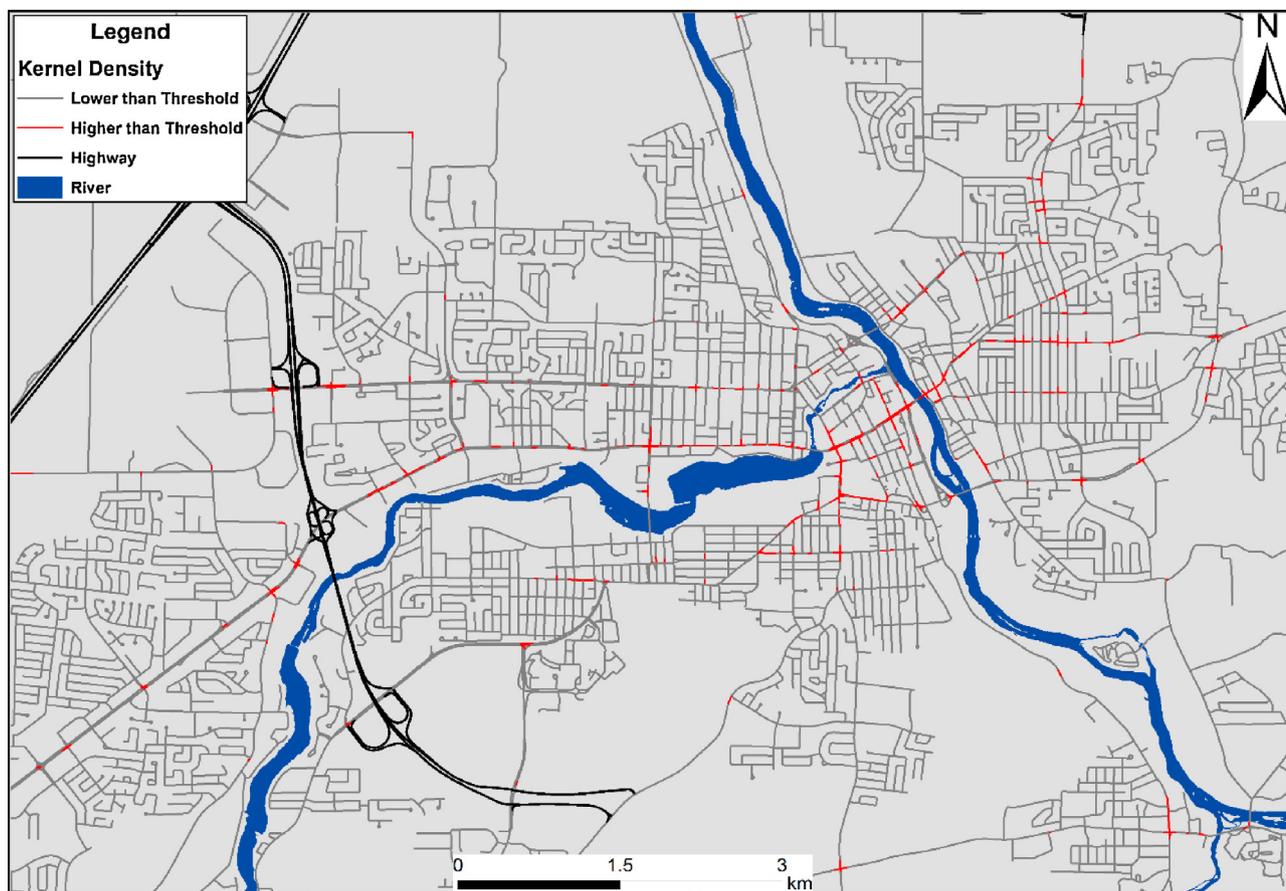


Fig. 5. Network KDE results based on three years aggregated crash data. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

to smallest). The location with the greatest difference is ranked as first, the location with the second highest difference is ranked as second, and so on (Gan et al., 2012).

4. Analysis results

The following section begins with an example of results for high crash density locations obtained from the network KDE method. Then, potential hotspot locations and a comparison with the results obtained using the network KDE method will be presented. Finally, we show how to take the exposure data into account to identify hotspot locations.

4.1. Results of three years of observed crash data using network KDE

In this section, the network KDE method used three years aggregated crash data to generate a crash density surface. It helps to identify high density sites and to examine the clustering of the collisions. The network KDE method used in this section is similar to the methods used in many prior studies (Larsen, 2010; Mohaymany et al., 2013; Plug et al., 2011).

In this method the evaluation takes place all over the study region. The results showed that there are many segments on the network which their densities are zero or near to the zero. Therefore, those segments with zero crashes or few crashes are not an interest.

The next step is finding a set of hazardous segments. Therefore, the threshold value (i.e., three standard deviation from the mean) is chosen to highlight the sites where the density are significantly

high. Fig. 5 shows those locations highlighted in red. The remaining section are not significant at chosen level. According to the results, there are 128 sites at which more crashes are taking place and their crash density are significantly high.

It should be noted that for better representation of the results and taking the details into consideration the geographical results are displayed on an enlarged scale. As shown in Fig. 5 the scale contains Sherbrooke's urban area and important roadways including local residential, collectors, main arteries, and downtown area.

As Fig. 5 shows, the high-density sites in the city of Sherbrooke are mainly located in three districts in the west, downtown, and east. In the west side, the high crash areas are mainly concentrated along two main arterial roads, which run east/west, and Jacques-Cartier Boulevard (passing over the Magog River), which runs north/south. The next high density area is Sherbrooke's downtown area (at the confluence of the Saint-François and Magog rivers) in the central business district (CBD). The traffic situation includes locals walking and driving to local shops/restaurants and buses traveling to downtown Sherbrooke. The area is bordered by three bridges to the east, which run east/west. These bridges provide the only roadway links between eastern and western Sherbrooke and traffic flow in these locations is relatively high. In the east side, high-density locations are mainly concentrated along three arterial roads (King-East Street, which runs east/west, and 12e Ave and 13e Ave, which run north/south).

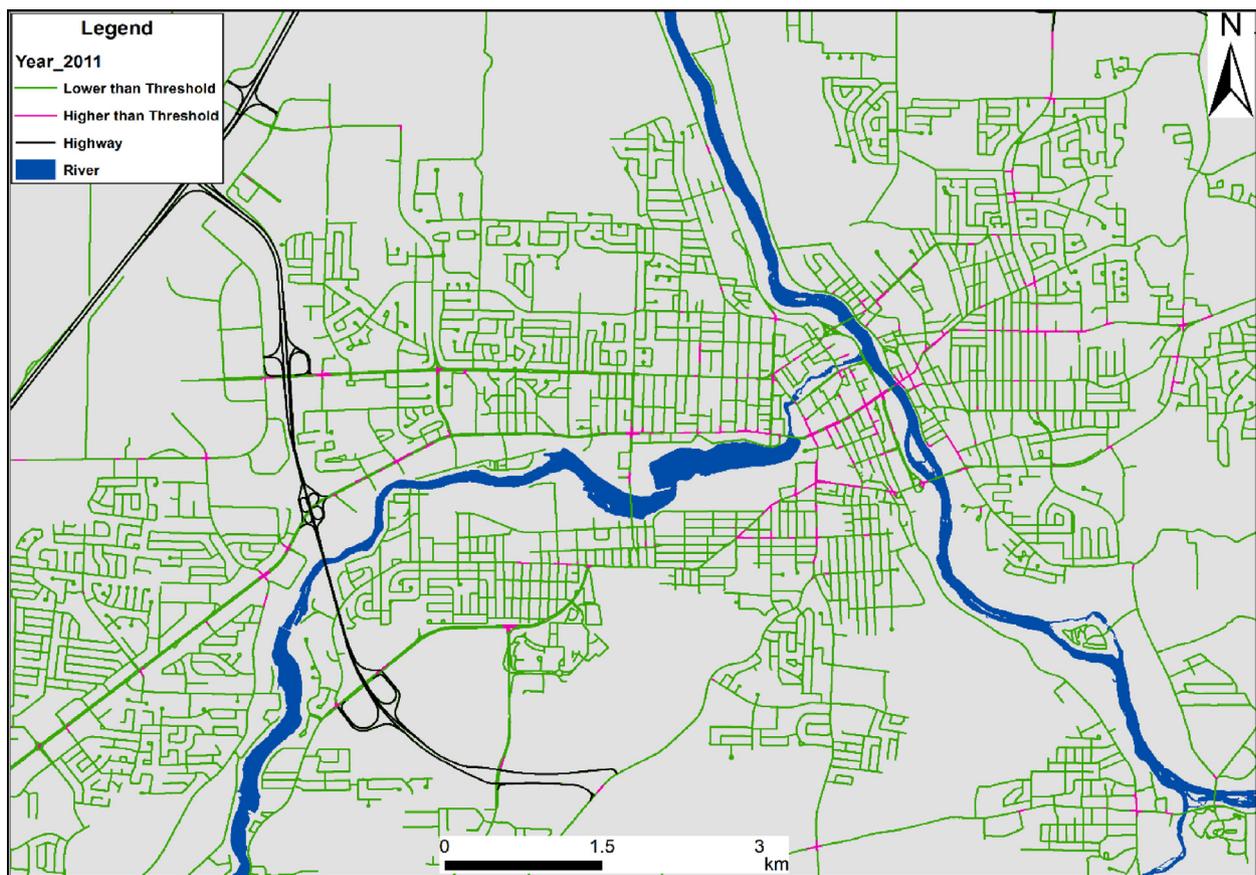
In this section we presented that the network KDE results based on aggregated crash data together with the threshold value can identify high crash density locations throughout the network. However, this approach failed to take into account whether the three years high crash density is due to chance or some continuing prob-

Table 1
Comparison between two network KDE results.

Method	No. of crashes in hotspot	Total crashes	Length of segments	Total length of segments	PAI
Network KDE (aggregated crash data)	3772	6918	46.17	1312	15.49
Network KDE (Potential hotspots)	3061	6918	35.75	1312	16.23

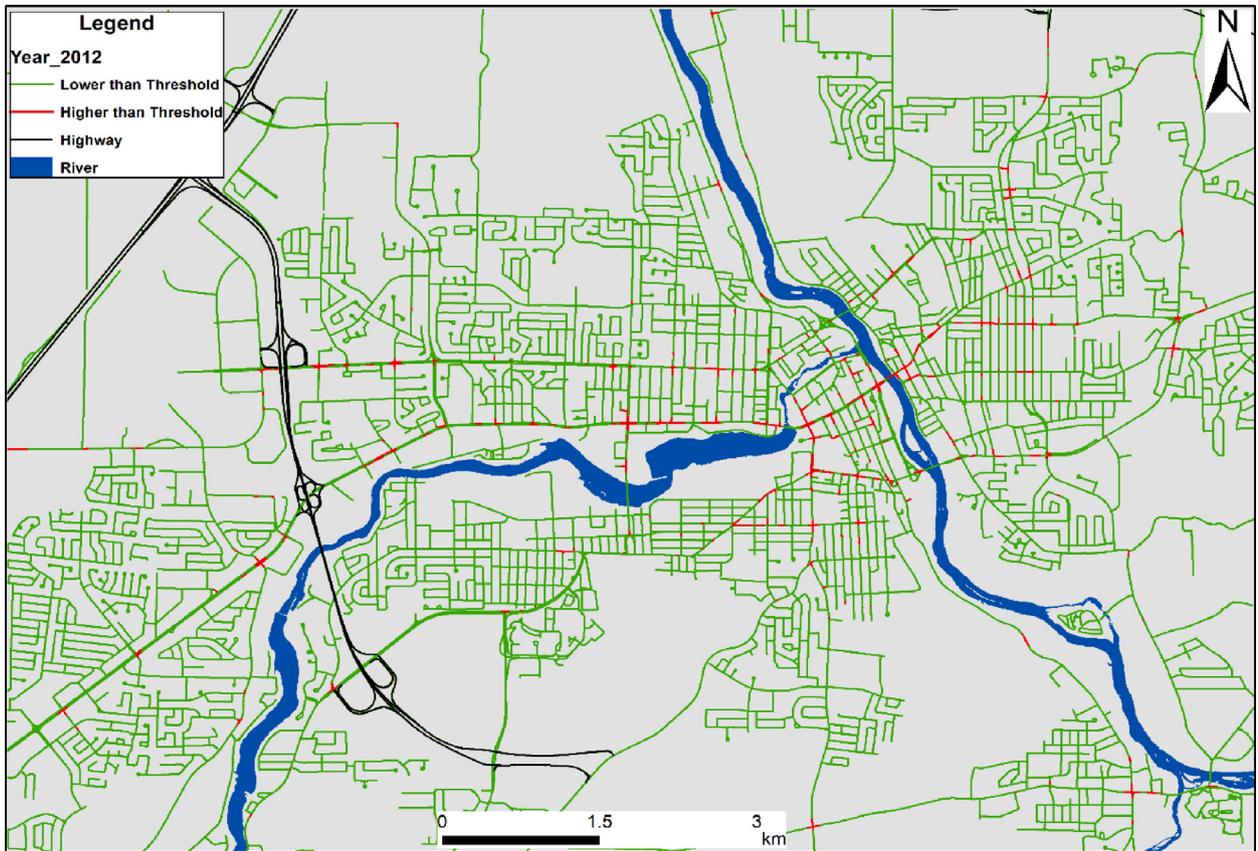
Table 2
20 hotspot locations based on the critical crash rate method.

Order	Intersection or Segment	Reference population name	Reference population type	Crash rate	Critical crash rate	Differences
1	INT	Galt west – Laurier	3SG	1.33	0.57	0.764
2	INT	College – Queen	4SG	1.34	0.72	0.622
3	INT	Bellvedere – College	4SG	1.35	0.82	0.525
4	INT	Terrill – Chicoyne	4ST	1.14	0.73	0.416
5	INT	Mcmanamy – Kingston	4SG	1.27	0.92	0.355
6	INT	Terrill – Rue du Cegep	4SG	1.13	0.78	0.351
7	INT	13 e Ave – Jardins Fleuris	4ST	2.13	1.82	0.319
8	INT	Jacques cartier – Tracy	3SG	0.75	0.57	0.183
9	INT	Portland – Industriel	4SG	1.29	1.11	0.177
10	INT	King west – Belvedere	4SG	1.17	1.01	0.164
11	INT	King east – Alphonse Laramee	3ST	0.56	0.43	0.132
12	INT	Sainte catherine – univeriste Blvd	3SG	0.63	0.50	0.131
13	INT	Alexandre – Ball	4SG	1.14	1.02	0.126
14	SEG	Boulevard Industriel	ART	2.47	2.35	0.119
15	INT	King west – Alexandre	3SG	0.60	0.51	0.091
16	INT	13 e Ave – Papineau	4SG	1.16	1.09	0.070
17	INT	Wellington – Aberdeen	4SG	0.91	0.86	0.057
18	INT	King west – Heneker	3ST	0.32	0.27	0.057
19	INT	Bourque – Bertrand Fabi	4SG	0.17	0.65	0.054
20	INT	King west – Jacques cartier	4SG	1.02	0.99	0.028

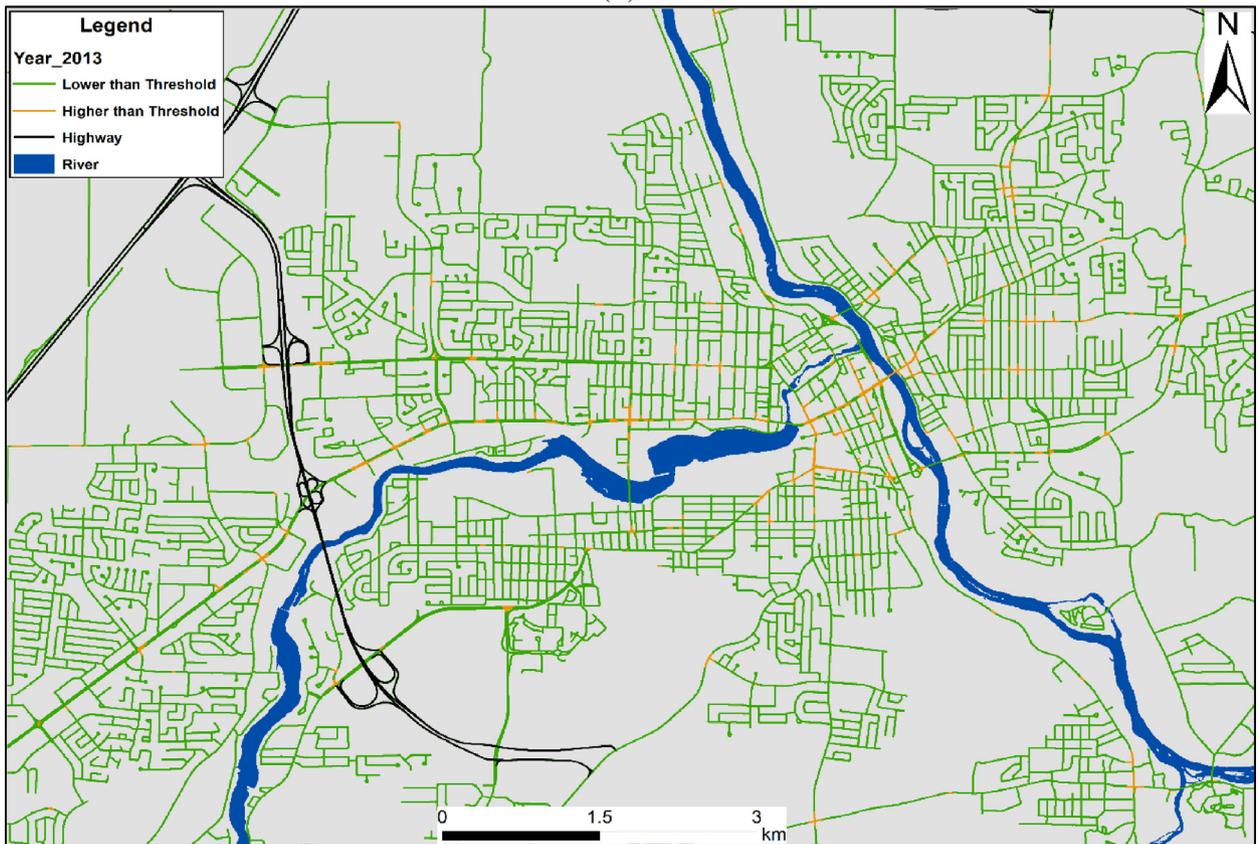


(a) 2011

Fig. 6. Collision density maps are higher than the threshold for three consecutive years: (a) 2011, (b) 2012, and (c) 2013.



(b) 2012



(c) 2013

Fig. 6. continued

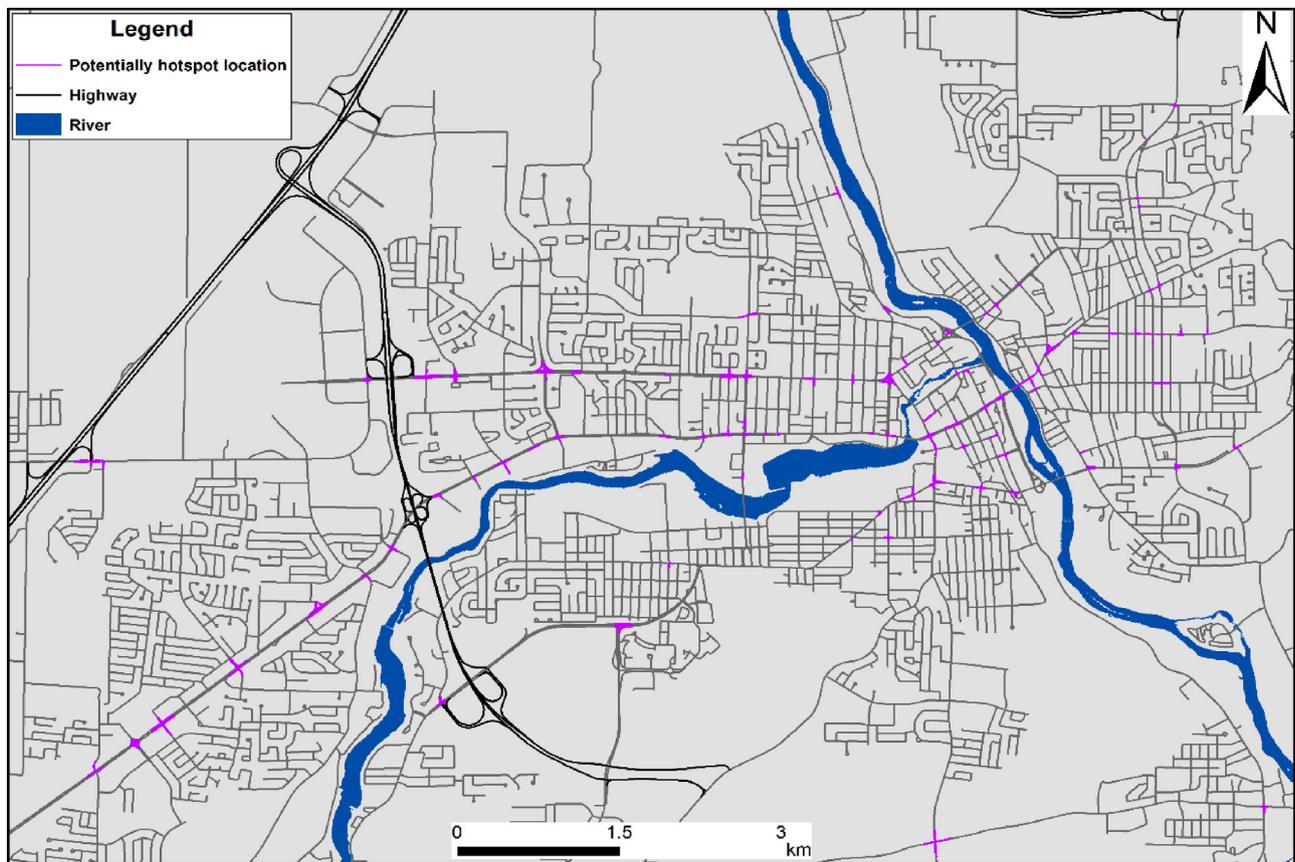


Fig. 7. Potential hotspot locations.(For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.).

lem. In order to overcome this problem a potential hotspot location proposed in the next section.

4.2. Exploring potential hotspot locations

This section discusses the method for identifying potential hotspot locations. To extract these locations, first the crash density of each three consecutive years (i.e., 2011, 2012, and 2013) must be created then the threshold value (three standard deviation from the mean) applied. Fig. 6 shows the network KDE results higher than the threshold for years 2011, 2012, and 2013. The maps clearly yielded similar patterns and the density values are significantly high in the downtown area, major arterial roadways and the entrances to the three main bridges, which run east/west. The results also show that the density surface is more continuous in 2011, followed by 2012 and 2013.

Once these three density maps are created, they must be spatially merged to extract the potential hotspot locations. Fig. 7 shows 97 potential hotspot locations in Sherbrooke (highlighted in violet), including 88 (90%) intersections and 9 segments (10%). Approximately 75% of these potential hotspot locations were at signalized intersections.

Compared to the results obtained from the three years aggregated crash data (discussed in Section 4.1), the number of high crash density locations are decreased from 128 sites to 98 sites. It means that, there are 30 locations where crashes occurred frequently in one period (for example in one year) and decreased in the next period.

The comparison also shows that the potential hotspot locations were mainly concentrated along the main arterial roadways and intersections, while the results discussed in Section 4.1 were more distributed along the network.

4.3. Methods comparison

A comparison was made between the network KDE results of three years aggregated crash data and the potential hotspots. This was performed by adopting the PAI. As Table 1 shows the PAI index for the network KDE results based on aggregated crash data was found 15.49, while a comparative value of 16.23 was found for the network KDE results based on potential hotspots. As stated earlier, larger PAI value means better ability of a method to locate high potential crashes in an area (Thakali et al., 2015).

4.4. Hotspot identification using the critical crash rate

The purpose of this section is to identify hotspot locations and minimize potential accidents in the study area. The first step was to calculate the crash rate and critical crash rate for each site. Any site with a crash rate higher than its critical crash rate will then be considered as a hotspot location. In this section, traffic volume data and crash count data for potential hotspots were used as input data.

From the analysis, 20 sites exceeded the critical crash rate. These locations are called hotspots and need to be reviewed in more depth. Table 2 shows the results of 20 hotspot locations in Sherbrooke including 19 intersections (95%) and one segment (5%). Identification of hotspot locations is crucial and is the first step in traffic safety improvement studies. Since budgets and time are limited, in many safety studies priority is given to the sites with the highest crash risk. As shown in Table 2, the sites are ordered and arranged in decreasing order according to the differences between the crash rate and critical rate, and a higher priority should be given to the greatest differences. Hence, the hotspots are classified into three priority levels: first priority level (shown in red), second

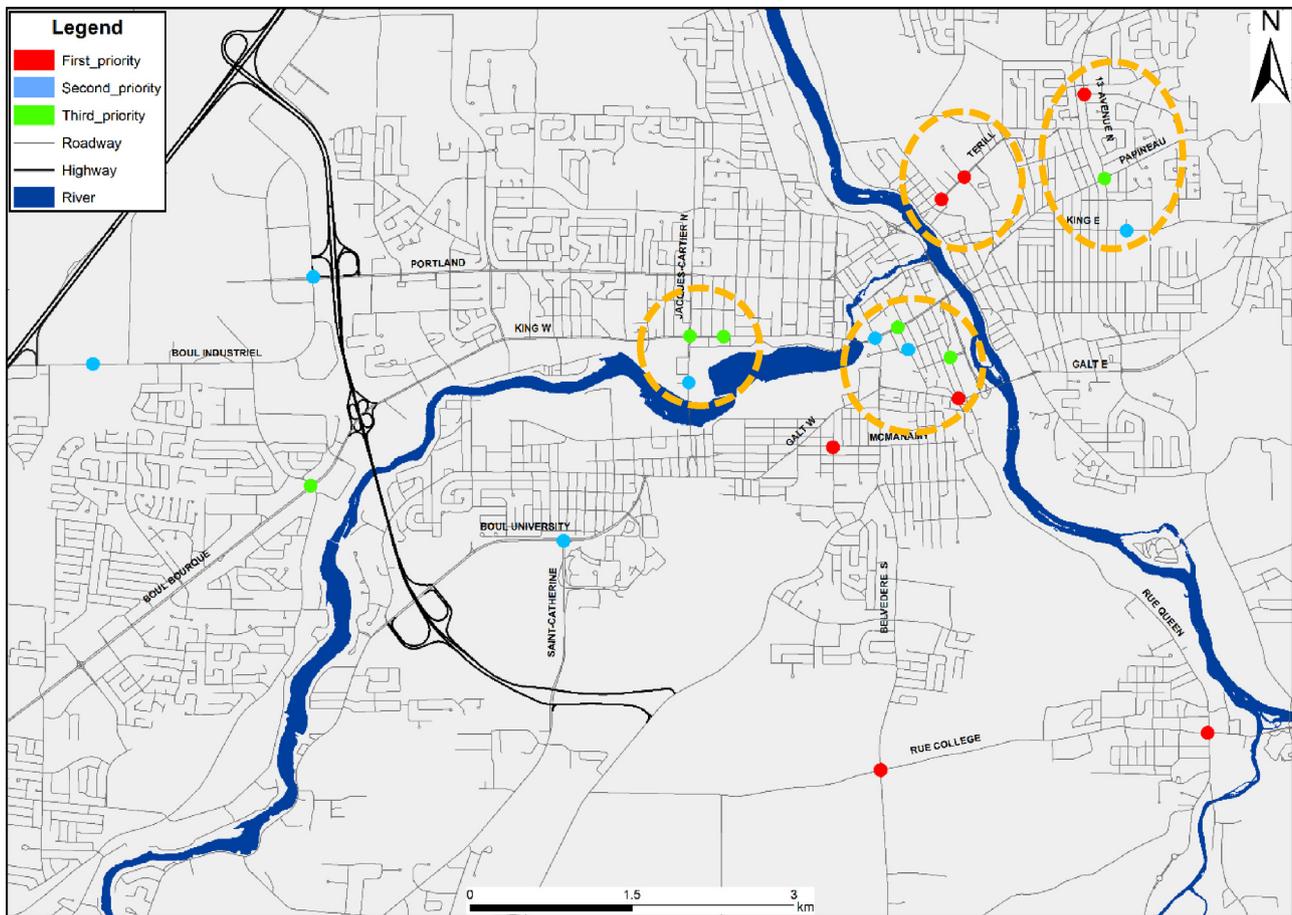


Fig. 8. Hotspot locations in the city of Sherbrooke. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

priority level (shown in blue), and third priority level (shown in green).

In the table above, INT means intersection, SEG means segment, SG means signalized intersection, ST means a stop-controlled intersection, ART means arterial roadway, and 4 means a four-leg intersection, while 3 means a three-leg intersection.

The results also show that over half of the hotspot locations (58%) were at four-leg signalized intersections followed by three-leg signalized intersections (22%), three-leg stop-controlled intersections (10%), and four-leg stop-controlled intersections (10%). This means that the signalized intersections were less safe than the stop-controlled intersections. Indeed, it is well known that signalized intersections have a considerable traffic volume.

Fig. 8 shows the distribution of hotspot locations in Sherbrooke. The map clearly shows the hotspot locations clustered in four specific areas (highlighted in orange) such as the Jacques-Cartier area (to the west), downtown area (in the centre), Terrill Street (to the east and next to the Saint-François River), and 13e Ave (further east on the Saint-François River). As shown in Fig. 8, about 25% of hazardous sites are located in the downtown area where there are many stores and restaurants as well as a high conflict between vehicles, pedestrians and cyclists. Terrill Street is dangerous, because it is located next to the Cegep de Sherbrooke, which is one of the main attraction zones (O/D zone) in Sherbrooke with about 6000 students and 8400 movements per day. The area near the Jacques-Cartier Bridge is another hazardous area since it is the only route that runs north/south to the west. 13e Ave is also risky because it is surrounded by the Quatre-Saisons shopping centre (about 2700 movements per day) and Fleurimont hospital (6700 movements per day), which are two main destination zones (O/D zones) in

Sherbrooke. It should be noted that the highlighted areas (shown in orange) are located in two main districts (Fleurimont and Jacques-Cartier), which comprise approximately 50% of the total Sherbrooke population.

5. Discussion and conclusion

This paper introduced a two-step approach to identify hotspot locations within a roadway network. The proposed approach integrated a GIS-based network KDE analysis with the HSM network screening method outlined in AASHTO (2010). It should be pointed out that this combination has a great impact on the overview of TAs from simple points to events with a spread risk. The study investigated how taking the exposure data into account affects the identification of hotspots. The approach was tested using traffic accident data from three years (2011–2013) for Sherbrooke's road network.

Unlike prior studies that used aggregated crash data, we proposed an approach to select particular sites from the network KDE results and to use them for further critical crash rate analysis. These locations, referred to as potential hotspots, followed two conditions: crashes occurred frequently (i.e., at least three consecutive years) at a particular site, and their crash density was relatively high. The advantages of using the approach are the following: (1) to filter out the locations for which the annual frequency of crashes is greater than the norm, and (2) to identify locations that have a continuing problem due to, for instance, a defect in the location (not due to chance). Likewise, a comparison was made between the network KDE results based on aggregated three-year crash data and the results based on three years of merged crash data (potential

hotspots). Their comparison, according to the PAI criterion, showed that the network KDE results based on potential hotspots are better for determining hotspots than the results based on aggregated crash data.

The critical crash rate method was selected in this study because it accounts for the significant variables that effect safety, including exposure data, the random nature of accidents and the type of intersections, as well as considering the variance in crash data. The other advantage of using the critical crash rate is that it establishes a threshold (critical rate) for comparison. If the crash rate at a site exceeds the critical crash rate, these locations are considered as hotspots. This method is very useful because it can rank and prioritize the sites in decreasing order according to the differences between the crash rate and critical rate, and a higher priority should be given to the greatest differences. Therefore, it can help transport authorities focus on the locations that have a documented problem. It should be noted that this method has a drawback, and it cannot be used in cases where the values for the traffic volume are not available for roadway segments.

The other significant advantage of this study is that it considered the whole road network (in large scale) including local, collector and artery roads, and it was not limited to one data set (i.e., only artery roads or highways). The study's results confirmed that the proposed approach can assist traffic authorities for the quick identification of the most dangerous locations within a roadway network. Finally, the method helps traffic authorities to prioritize hotspot locations more efficiently and allocate their limited budget and resources.

This study suggests that further research is needed in the following areas. First, this study used an iterative (trial and error) technique to find the most appropriate bandwidth size in the network KDE analysis. Therefore, the development of a scientific method for selecting the most appropriate bandwidth size should be considered in future research. Second, in this study, we simply used the observed crash counts and did not make a distinction between different levels of crash severity such as Property Damage Only (PDO), serious injuries and fatal crashes. Therefore, a further study is needed to show how taking the severity of crashes into consideration affects the identification of hotspots. Third, there are other factors that may affect the identification of hotspots including road geometry (e.g., road type, number of lanes), socio-economic environment (e.g., the household income) and weather conditions (e.g., rain, snow and fog). Lastly, in this study we followed the results from Xie and Yan (2008) and ran the network KDE with a 10-m length. Therefore, a sensitivity analysis is needed using other parameters.

Acknowledgements

We would like to thank three anonymous reviewers and the editor whose comments substantially increased the quality of the paper.

References

- AASHTO (American Association of State Highway and Transportation Officials), 2010. *Highway Safety Manual, first ed.* AASHTO, Washington, DC.
- Ayyub, B.M., McCuen, R.H., 2011. *Probability, Statistics, and Reliability for Engineers and Scientists*. CRC Press.
- Bil, M., Andrášik, R., Janoška, Z., 2013. Identification of hazardous road locations of traffic accidents by means of kernel density estimation and cluster significance evaluation. *Accid. Anal. Prev.* 55, 265–273.
- Bailey, T.C., Gatrell, A.C., 1995. *Interactive spatial data analysis*. In: Harlow Essex, England: Longman Scientific & Technical. J. Wiley.
- Borruso, G., 2008. Network density estimation: a GIS approach for analysing point patterns in a network space. *Trans. GIS* 12 (3), 377–402.
- Chainey, S., Ratcliffe, J., 2013. *GIS and Crime Mapping*. John Wiley & Sons.
- Chainey, S., Tompson, L., Uhlig, S., 2008. The utility of hotspot mapping for predicting spatial patterns of crime. *Secur. J.* 21 (1–2), 4–28.
- Dunn, B., McDaniel-wilson, C., Iii, J.M., Appanaitis, G., Associates, D.K.S., Batten, A.M., Brooks, A., Stabler, B., Fine, D., Cathcart, R., Crownover, D., 2015. *Anal. Proced. Manual*.
- ESRI, 2015. Data Classification Methods—ArcGIS Pro | ArcGIS for Desktop [WWW Document] (accessed 1.18.16) URL <http://pro.arcgis.com/en/pro-app/help/mapping/symbols-and-styles/data-classification-methods.htm>.
- Erdogan, S., Yilmaz, I., Baybura, T., Gullu, M., 2008. Geographical information systems aided traffic accident analysis system case study: city of Afyonkarahisar. *Accid. Anal. Prev.* 40 (1), 174–181.
- Ervin, D., 2015. Advanced Spatial Analysis [WWW Document] (accessed 12.11.15) URL <http://gispopsci.org/point-pattern-analysis/>.
- FHWA, 2004. Signalized intersections: Informational guide. Federal Highway Administration. U.S. Department of Transportation, Washington, DC, Available from <https://www.fhwa.dot.gov/publications/research/safety/04091/04091.pdf>.
- FHWA, 2012. Chapter 6 – Signalized Intersections: Informational Guide, – FHWA-HRT-04-091 [WWW Document] (accessed 12.21.15) URL <https://www.fhwa.dot.gov/publications/research/safety/04091/06.cfm>.
- FHWA, 2011. Roadway Safety Information Analysis – Safety | Federal Highway Administration [WWW Document] (accessed 12.20.15) URL http://safety.fhwa.dot.gov/local_rural/training/fhwasaxx1210/s3.cfm.
- Flahaut, B., Mouchart, M., San Martin, E., Thomas, I., 2003. The local spatial autocorrelation and the kernel method for identifying black zones: a comparative approach. *Accid. Anal. Prev.* 35 (6), 991–1004.
- Furuta, T., Suzuki, A., Okabe, A., 2008. A voronoi heuristic approach to dividing networks into equal-sized sub-networks. *Forma* 23 (2), 73–79.
- Gan, A., Haleem, K., Alluri, P., Saha, D., 2012. Standardization of Crash Analysis in Florida. Lehman Center for Transportation Research, Miami.
- Hart, T.C., Zandbergen, P.A., 2012. Effects of data quality on predictive hotspot mapping. *National Criminal Justice Reference Service*.
- Khan, G., Qin, X., Noyce, D.A., 2008. Spatial analysis of weather crash patterns. *J. Transp. Eng.* 134 (5), 191–202.
- Krisp, J.M., Durot, S., 2007. Segmentation of lines based on point densities—an optimisation of wildlife warning sign placement in southern Finland. *Accid. Anal. Prev.* 39 (1), 38–46.
- Kuo, P.F., Zeng, X., Lord, D., 2011. Guidelines for choosing hot-spot analysis tools based on data characteristics, network restrictions, and time distributions. In: *Proceedings of the 91 Annual Meeting of the Transportation Research Board, January*, pp. 22–26.
- Larsen, M., 2010. Philadelphia traffic accident cluster analysis using GIS and SANET. In: *Master of Urban Spatial Analytics*. School of Design, University of Pennsylvania.
- Loo, B.P., Yao, S., 2013. The identification of traffic crash hot zones under the link-attribute and event-based approaches in a network-constrained environment. *Comput. Environ. Urban Syst.* 41, 249–261.
- Lu, Y., Chen, X., 2007. On the false alarm of planar K-function when analyzing urban crime distributed along streets. *Soc. Sci. Res.* 36 (2), 611–632.
- Miaou, S.P., 1994. The relationship between truck accidents and geometric design of road sections: poisson versus negative binomial regressions. *Accid. Anal. Prev.* 26 (4), p.471–482.
- Mohaymany, A.S., Shahri, M., Mirbagheri, B., 2013. GIS-based method for detecting high-crash-risk road segments using network kernel density estimation. *Geo-spatial Inf. Sci.* 16 (2), 113–119.
- Nie, K., Wang, Z., Du, Q., Ren, F., Tian, Q., 2015. A network-constrained integrated method for detecting spatial cluster and risk location of traffic crash: a case study from Wuhan, China. *Sustainability* 7 (3), 2662–2677.
- O'Sullivan, D., Unwin, D., 2014. *Geographic Information Analysis*. John Wiley & Sons.
- Okabe, A., Sugihara, K., 2012. *Spatial Analysis Along Networks: Statistical and Computational Methods*. John Wiley & Sons.
- Okabe, A., Okunuki, K.I., Shiode, S., 2006. SANET: a toolbox for spatial analysis on a network. *Geogr. Anal.* 38 (1), 57–66.
- Okabe, A., Satoh, T., Sugihara, K., 2009. A kernel density estimation method for networks, its computational method and a GIS-based tool. *Int. J. Geogr. Inf. Sci.* 23 (1), 7–32.
- Oris, W.N., 2011. *Spatial Analysis of Fatal Automobile Crashes in Kentucky*. Road Safety Manual: Recommendations from the World Road Association (PIARC). Route 2 Market, 2003.
- Plug, C., Xia, J.C., Caulfield, C., 2011. Spatial and temporal visualisation techniques for crash analysis. *Accid. Anal. Prev.* 43 (6), 1937–1946.
- Sabel, C.E., Kingham, S., Nicholson, A., Bartie, P., 2005. Road traffic accident simulation modelling—a kernel estimation approach. In: *The 17th Annual Colloquium of the Spatial Information Research Centre University of Otago, Dunedin, New Zealand*, pp. 67–75.
- Silverman, B.W., 1986. *Density Estimation for Statistics and Data Analysis*, 26. CRC Press.
- Steenberghen, T., Dufays, T., Thomas, I., Flahaut, B., 2004. Intra-urban location and clustering of road accidents using GIS: a Belgian example. *Int. J. Geogr. Inf. Sci.* 18 (2), 169–181.
- Steenberghen, T., Aerts, K., Thomas, I., 2010. Spatial clustering of events on a network. *J. Trans. Geogr.* 18 (3), 411–418.
- Sugihara, K., Satoh, T., Okabe, A., 2010. Simple and unbiased kernel function for network analysis. In: *Communications and Information Technologies (ISCIT), 2010 International Symposium*, pp. 827–832, IEEE.

- Thakali, L., Kwon, T.J., Fu, L., 2015. Identification of crash hotspots using kernel density estimation and kriging methods: a comparison. *J. Mod. Transp.* 23 (2), 93–106.
- Tao, D.P.Y., Nascimento, K.M.M.A., Shekhar, M.M.S., Huang, Y., 2011. *Advances in Spatial and Temporal Databases*.
- Timotheé, P., Nicolas, L.B., Emanuele, S., Sergio, P., Stéphane, J., 2010. A network based kernel density estimator applied to Barcelona economic activities. In: *International Conference on Computational Science and Its Applications* (pp. 32–45), Springer Berlin Heidelberg.
- Toroyan, T., 2013. *Global Status Report on Road Safety 2015. Supporting a Decade of Action*. World Health Organization, Department of Violence and Injury Prevention and Disability, Geneva.
- Van Patten, I.T., McKeldin-Coner, J., Cox, D., 2009. A microspatial analysis of robbery: prospective hot spotting in a small city. *Crime Mapping. J. Res. Pract.* 1 (1), 7–32.
- Vemulapalli, S.S., 2015. GIS-based spatial and temporal analysis of aging-involved crashes in Florida (Doctoral dissertation, The Florida State University).
- Xie, Z., Yan, J., 2008. Kernel density estimation of traffic accidents in a network space. *Computers. Environ. Urban Syst.* 32 (5), 396–406.
- Xie, Z., Yan, J., 2013. Detecting traffic accident clusters with network kernel density estimation and local spatial statistics: an integrated approach. *J. Transp. Geogr.* 31, 64–71.
- Yamada, I., Thill, J.C., 2004. Comparison of planar and network K-functions in traffic accident analysis. *J. Transp. Geogr.* 12 (2), 149–158.
- Yao, S., Loo, B.P., Yang, B.Z., 2015. Traffic collisions in space: four decades of advancement in applied GIS. *Ann. GIS*, 1–14.
- Young, J., Park, P.Y., 2014. Hotzone identification with GIS-based post-network screening analysis. *J. Transp. Geogr.* 34, 106–120.