

MÉTHODE D'IDENTIFICATION DE CONCENTRATIONS LOCALES D'ÉVÉNEMENTS DANS UN SEMIS DE POINTS APPLICATION AUX ACCIDENTS DE LA ROUTE

Arnaud BANOS

Florence HUGUENIN-RICHARD

THEMA, UPRESA 6049

Université de Franche-Comté

Résumé

L'article présente une méthode d'analyse exploratoire d'un semis de points issue de la "Machine d'Analyse Géographique" (GAM) de Stan OPENSHAW. Elle permet d'identifier de manière automatisée des concentrations locales significatives d'événements localisés dans l'espace géographique. À titre d'exemple, la méthode est appliquée aux accidents de la route ayant eu lieu dans la Communauté Urbaine de Lille.

Summary

The Geographical Analysis Machine (GAM), developed by Stan OPENSHAW, is an automated exploratory spatial data analysis method of point data. It allows to identify localized geographical clustering of events in space. This method is applied to road accidents data in the urban community of Lille (France).

Mots-Clés

Accidents de la route, Agrégats spatiaux, Analyse exploratoire, Communauté urbaine de Lille, Machine d'analyse géographique, Semis de points

Key-Words

Clusters, Exploratory analysis, Geographical analysis machine, Lille (France), Point pattern analysis, Road accidents

Les systèmes d'information géographique sont relativement pauvres du point de vue des méthodes d'analyse spatiale qu'ils intègrent. Pourtant, un besoin toujours plus prégnant de telles méthodes, adaptées à la nature des données et aux problématiques spécifiques de l'analyse spatiale se fait sentir.

L'objectif de cet article est de présenter une méthode d'analyse exploratoire d'un semis de points, qui pourrait utilement être implémentée en série dans les systèmes d'information géographique.

L'application de cette méthode aux accidents de la route dans la communauté urbaine de Lille permet d'en démontrer l'intérêt.

1. Problématique

En 1996, plus de 1 500 accidents corporels ont été recensés sur le territoire de la communauté urbaine de Lille, un ensemble de 87 communes autour de Lille, Roubaix et Tourcoing. Cette information peut être représentée sous la forme d'un semis de points (Figure 1). En effet, chaque accident est connu en tant qu'événement discret, localisé de façon précise par ses coordonnées géographiques et renseigné par des attributs de nature quantitative (nombre de blessés, nombre de tués, etc.) et qualitative (type d'usager impliqué, jour de l'accident, etc.).

A partir de cette population de référence, il est possible d'extraire des sous-populations par requête attributaire: par exemple les accidents ayant impliqué au moins un piéton âgé de moins de dix ans (Figure 1). L'observation du semis de points ainsi extrait appelle deux remarques:

- la proximité géographique des accidents, notamment en milieu urbain, rend difficile l'estimation de leur nombre lorsque l'on représente à petite échelle le semis de leurs points (de nombreux points se superposent les uns sur les autres, des points proches par la distance forment des agglomérats, etc.),
- la forme du semis de points suggère l'existence de structures spatiales proches de celles de la population de référence: les zones de forte concentration d'accidents de piétons semblent correspondre aux zones les plus accidentogènes. Ce constat visuel immédiat mérite toutefois d'être précisé: la distribution de la sous-population reflète-t-elle uniquement celle de la population de référence ou au contraire existe-t-il des concentrations locales et anormales d'accidents de piétons non identifiables à l'œil nu?

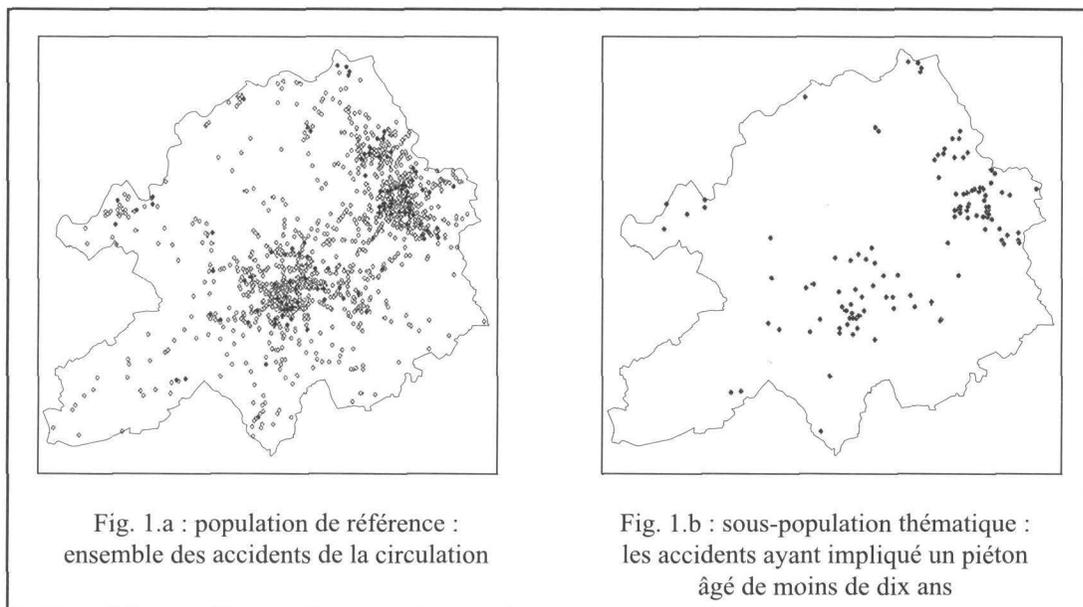


Figure 1 - Les accidents corporels dans la communauté urbaine de Lille en 1996

L'analyse d'un semis des accidents est une approche peu fréquente en accidentologie, notamment parce que les bases de données spatiales sont jusqu'ici assez rares¹.

Cependant, la recherche de concentrations d'accidents peut être rapprochée de problématiques majeures, comme:

- l'identification de sites d'accident, reposant sur l'estimation d'un indice de risque: approche déterministe basée sur la fréquence d'accidents pondérée par leur gravité ou sur des mesures d'exposition au risque (trafic, densité de population) approche probabiliste par des méthodes empiriques bayésiennes [22],
- la localisation de points noirs en milieu urbain. Un point noir est une zone d'accumulation locale d'accidents. L'hypothèse du point noir est qu'il existe sur le réseau des lieux avec un niveau anormal d'insécurité qu'il convient de modifier pour revenir à une situation plus normale. La solution la plus simple pour les détecter consiste à cartographier les accidents sur le réseau et à sélectionner les lieux où leur nombre est le plus élevé. Cependant la fréquence n'est pas un indicateur fiable dans le sens où le risque dépend fortement d'une mesure d'exposition: « la localisation informatique des accidents n'existe pas partout et, lorsqu'elle existe, les mesures de flux de véhicules - sans parler des flux de deux-roues et de piétons - ne sont pas toujours disponibles, alors qu'elles seraient nécessaires pour calculer des taux de risque » constate D. Fleury [7].

Finalement, plusieurs contraintes compliquent l'identification de concentrations d'accidents dans une approche classique: les principales méthodes sont dédiées à des données agrégées à des unités spatiales diverses (commune, tronçon de voie, carrefour), et nécessitent des informations supplémentaires relatives à des mesures d'exposition au risque très difficiles à se procurer (trafic).

Des méthodes statistiques développées pour l'analyse d'un semis de points permettent de répondre à notre problème :

- fournir des résultats de façon automatique et facilement interprétables,
- utiliser l'information brute sans opérer d'agrégation,
- détecter des concentrations anormales d'accidents sans avoir recours aux données de trafic, mais en intégrant la notion de population exposée à un risque,
- prendre en compte la dimension spatiale de la distribution des données.

Des exemples d'application de telles méthodes existent en épidémiologie et en géographie. On distinguera ici trois grandes approches: la première est basée sur des mesures de distance entre individus, la deuxième sur l'estimation de densités et la troisième est issue de la « geographical analysis machine » de Stan OPENSHAW.

La première famille de méthodes, basée sur des mesures de distance entre événements (principe général de la méthode dite des « plus proches voisins »), repose sur l'hypothèse que la population cible est issue d'un processus de Poisson hétérogène et d'intensité proportionnelle à celle de la population de référence. Il s'agit alors de comparer les distributions de ces deux populations, avec la fonction K par exemple, mise au point par RIPLEY [2]. CUSICK et EDWARDS utilisent, quant à eux, une population de contrôle, obtenue par tirage aléatoire à partir de la population de référence [6].

Ces méthodes présentent néanmoins un désavantage certain, de notre point de vue, dans la mesure où elles ne possèdent pas de traduction cartographique aisément interprétable. Deux autres approches sont à ce titre particulièrement attrayantes.

La seconde famille de méthodes repose sur l'estimation d'une surface de risque, à partir des semis de points originaux. Pour cela, l'intensité² de chacune des distributions est estimée à l'aide d'une fonction de densité de Kernel: il s'agit d'une fonction mobile circulaire tri-dimensionnelle qui balaye systématiquement l'espace étudié. Chaque événement rencontré est alors pondéré, en fonction de sa distance au centre [1] [5] [9]. Le ratio des deux estimations ainsi obtenues permet d'aboutir, après interpolation, à la surface de risque recherchée [4] [19]. Cette solution aurait fort bien pu être retenue ici et fait d'ailleurs l'objet d'investigations de notre part. Nous avons néanmoins préféré nous orienter vers une troisième famille de méthodes, issue des travaux de Stan OPENSHAW.

La « machine d'analyse géographique », développée au début des années 1980, teste localement et de manière répétitive, un très grand nombre d'hypothèses spatiales [17]. Cette méthode de calcul intensif originale a été largement diffusée [11] [14] [16] [21] et fait même l'objet d'un site web [18]. Des tentatives d'amélioration ont été proposées [3] [13] [8]. Nous retiendrons celles apportées récemment par FOTHERINGHAM et ZHAN, et ceci pour deux raisons: d'une part, l'algorithme est plus simple à mettre en œuvre; d'autre part, les modifications introduites sont spécialement adaptées à une information de type ponctuel³. Enfin, les résultats obtenus sont remarquablement proches de ceux obtenus par la méthode d'OPENSHAW.

2. Présentation de la méthode

La méthode repose sur la comparaison statistique de la distribution spatiale de la sous-population avec sa distribution théorique associée, construite sous l'hypothèse d'une répartition spatiale aléatoire. La significativité des écarts observés entre ces deux distributions est ensuite testée par application de la loi de Poisson.

● Dans un premier temps, la proportion moyenne d'observer au hasard un individu de la sous-population dans la zone d'étude est calculée :

$$P_{\text{moyenne}} = \frac{\text{effectif de la sous-population}}{\text{effectif de la population de référence}}$$

● On génère ensuite de façon aléatoire et itérative un nombre défini par l'utilisateur de fenêtres mobiles circulaires, dont le rayon est choisi au hasard dans un intervalle paramétré. L'objectif d'un tel procédé est de couvrir au mieux et avec un nombre minimum de fenêtres l'ensemble de la zone d'étude.

● Pour chacune de ces fenêtres, un certain nombre d'opérations sont réalisées :

- les individus de la population de référence (α) et les individus de la sous-population (β) sont dénombrés,
- le nombre théorique d'individus de la sous-population que l'on devrait observer si la répartition de ces évènements était aléatoire (λ) est calculé :

$$\lambda = \alpha * P_{\text{moyenne}}$$

– enfin, l'écart entre le nombre observé d'individus de la sous-population (β) et le nombre attendu (λ) est testé à l'aide d'une distribution de Poisson :

$$P(\beta, \lambda) = \frac{e^{-\lambda} * \lambda^{\beta}}{\beta!}$$

Ce test revient à calculer la probabilité d'observer au hasard exactement β accidents ayant, par exemple « impliqués au moins un piéton de moins de 10 ans », alors que l'on en attend λ .

● On représente sur le semis de points de la sous-population toutes les fenêtres mobiles pour lesquelles la probabilité $P(\beta, \lambda)$ est inférieure à un seuil fixé par l'utilisateur. Le chevauchement des cercles exprime directement l'intensité de la concentration locale mise à jour.

Cet algorithme a été implémenté par nos soins au sein de l'environnement de programmation statistique Xlisp-Stat [20].

3. Application à la recherche en accidentologie

Appliquée à la sous-population des accidents de piétons de moins de dix ans, la méthode met en évidence un certain nombre de concentrations locales. La figure 2 présente les résultats obtenus à partir de 5 000 fenêtres mobiles de rayon compris entre 100 et 1 000 mètres⁴. Chacun des quatre graphiques correspond à un seuil de probabilité fixé par l'utilisateur.

Le nombre de concentrations mises à jour décroît à mesure que l'on réduit la valeur du seuil de probabilité, c'est-à-dire à mesure que l'on réduit la probabilité d'obtenir chaque concentration par hasard. Par ailleurs du point de vue de l'interprétation il semble raisonnable de ne retenir que les concentrations subsistant aux seuils de significativité les plus faibles (cinq pour mille et un pour mille).

Dans le cas des accidents de piétons de moins de dix ans, deux zones de concentration locale sont identifiées, l'une sur la ville de Roubaix et la seconde sur la commune d'Halluin.

Figure 2 - Identification de concentrations locales d'accidents de piétons-enfants dans la communauté urbaine de Lille en 1996

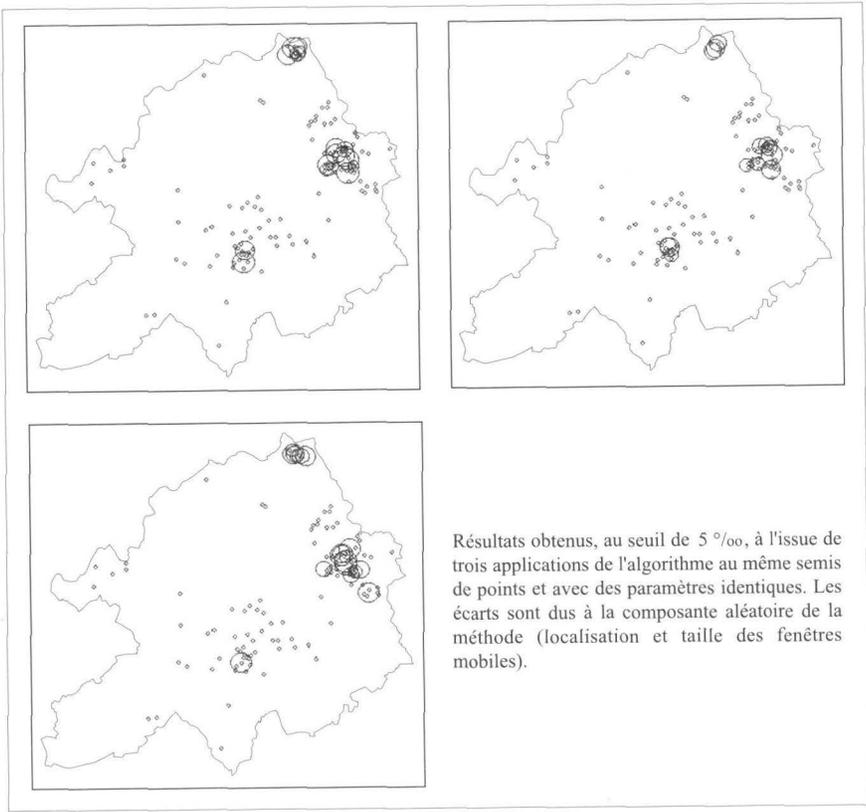
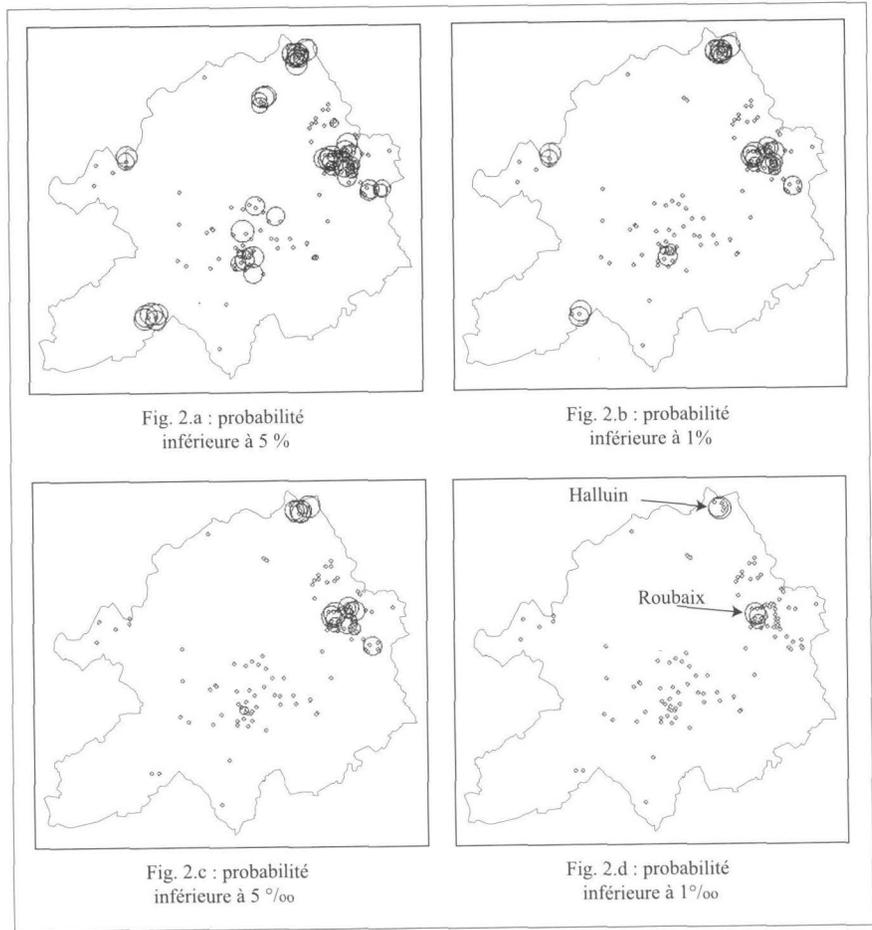


Figure 3 - Evaluation de la stabilité de la méthode

Avant d'aller plus loin dans l'interprétation de ces résultats, il semble important de déterminer la fiabilité de la méthode. Deux tests ont été appliqués.

Le premier permet d'évaluer la stabilité des résultats, c'est-à-dire la capacité de la méthode à reproduire des résultats relativement proches dans des conditions identiques. Les écarts observés ne doivent être imputables qu'à la part d'aléatoire de la méthode, résidant dans le choix de localisation du centre et la taille de chaque fenêtre mobile. La figure 3 montre ainsi les résultats de trois tests obtenus sur le même semis de points, avec des paramètres identiques et pour le seuil de probabilité de 5 pour 1000. Les mêmes concentrations d'accidents de piétons de moins de dix ans sont mises en évidence, attestant la stabilité de la méthode.

Le deuxième test revient à évaluer indirectement la validité même de la méthode. On s'assure qu'elle ne met pas en évidence des concentrations locales qui n'existent pas. Pour ce faire, une autre sous-population, de même taille que celle des accidents de piétons de moins de dix ans, a été extraite de la population totale des accidents de 1996, par un échantillonnage aléatoire. Cette nouvelle sous-population peut être interprétée comme issue d'un processus de Poisson hétérogène: les structures spatiales visibles ne sont dues qu'à la répartition même de la population de référence. L'application de la méthode, avec les mêmes paramètres, sur ce nouveau semis, ne doit donc mettre à jour aucune concentration significative. Les quelques cercles subsistant au seuil de probabilité de 5‰ (Figure 4) sont bien trop rares pour pouvoir être retenus comme des zones de concentration excessive d'accidents de piétons. Rappelons en effet que c'est la densité des cercles qui exprime l'intensité des phénomènes mis à jour.

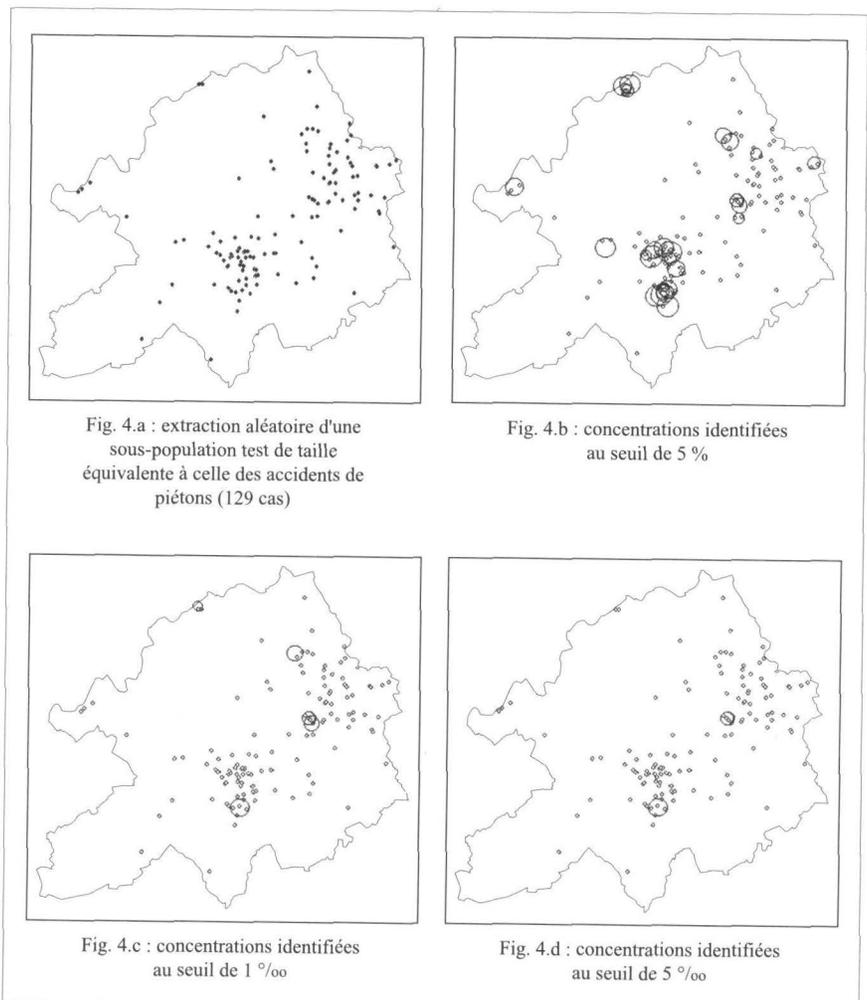


Figure 4 - Test de la méthode sur un échantillon aléatoire

L'identification de concentrations locales présente en elle-même un intérêt. Néanmoins, du point de vue de l'analyste, c'est le point de départ d'une démarche explicative visant à comprendre les processus spatiaux spécifiques susceptibles d'avoir créé de telles répartitions spatiales. Une première façon de procéder est de caractériser les accidents identifiés constitutifs des zones de concentration. Dans le cas des accidents de piétons de moins de dix ans, une série de traits communs a pu être mise en évidence en interrogeant les attributs de chaque accident concerné : essentiellement des usagers légèrement blessés ayant traversé la rue hors passage protégé ou sans précaution. Des tranches horaires spécifiques apparaissent : les heures de pointe du midi et du soir ainsi que l'heure creuse de l'après-midi.

Cette approche reste pourtant insuffisante, dans la mesure où les facteurs locaux d'insécurité que l'on cherche à identifier dépendent de l'environnement autour du lieu d'accident. Il serait alors intéressant de pouvoir intégrer ces résultats dans un système d'information géographique, gérant de l'information spatiale à plusieurs niveaux d'observation et sur plusieurs thématiques (le réseau et ses caractéristiques, des données sur l'environnement, sur la population, le trafic, etc.) [12].

4. Limites, intérêts et perspectives

La méthode présentée permet d'identifier de manière automatisée, fiable et aisément interprétable des concentrations locales d'événements. Elle offre par ailleurs le mérite d'être applicable à n'importe quel phénomène ponctuel dont la population de référence est connue. Elle mériterait amplement, à ce titre, de faire partie des outils disponibles en série au sein des SIG classiques.

Du point de vue de l'accidentologie, cette méthode peut très bien être utilisée dans la problématique plus large d'identification de sites dangereux. En effet, elle permet de cartographier rapidement des « zones à problème » et de fixer les enjeux d'une analyse plus approfondie. Elle se situe donc dans une démarche exploratoire et préliminaire de l'étude spatiale de l'insécurité routière.

Un certain nombre de réserves peuvent néanmoins être formulées. La première est liée à la nature même du phénomène étudié. L'espace des accidents de la circulation est avant tout celui du réseau viaire, à de rares exceptions près peut-être. Or la méthode utilisée ignore cette dimension essentielle, puisqu'on raisonne dans le plan. L'interprétation des résultats obtenus doit donc être menée avec prudence.

La deuxième est relative à un problème scalaire. Les événements identifiés comme constitutifs d'une même « poche » de concentration locale ne sont pas forcément très proches dans l'espace géographique. La taille des fenêtres mobiles circulaires doit en effet être gardée à l'esprit. Nous avons appliqué cette méthode à une vaste zone, la communauté urbaine de Lille, et pour une seule année, ce qui représente un peu plus de 1 500 accidents. Il serait intéressant de travailler à une échelle plus grande, comme celle d'une commune par exemple. Mais, le nombre d'accidents ayant eu lieu sur une seule année est alors insuffisant : par exemple en 1996, 389 accidents corporels ont été recensés sur la commune de Lille. Le problème ainsi soulevé peut être résolu de deux manières. Il est possible de mener des analyses sur des portions réduites du territoire mais en travaillant avec des données sur plusieurs années. Une seconde solution, plus ambitieuse, consiste à modifier l'algorithme utilisé, de façon à ne raisonner que dans l'espace du réseau, piste de recherche sur laquelle nous espérons progresser.

Une troisième réserve concerne le mode de représentation des résultats retenu. La technique de superposition graphique des fenêtres mobiles mériterait en effet d'être améliorée. Dans cette perspective, OPENSHAW a remplacé récemment l'affichage brut des fenêtres mobiles par l'estimation d'une surface de risque à partir de fonctions de densité de Kernel [18].

Enfin, la procédure de test statistique utilisée doit être considérée à sa juste valeur. Il ne s'agit pas d'un test formel de significativité, obéissant aux canons de la statistique confirmatoire. Il s'agit plutôt, comme le signale OPENSHAW, « d'un obstacle qu'une fenêtre mobile doit franchir, pour pouvoir contribuer à la représentation cartographique finale ». Par ailleurs, le choix de la distribution de Poisson est relatif. Supposer que la répartition des accidents de piétons se rapproche d'une distribution spatiale aléatoire, générée par un

processus de Poisson hétérogène, n'a en soi rien d'audacieux [15]. Néanmoins, la proportion d'accidents de piétons (de l'ordre de 10 %), nous éloigne quelque peu de la définition statistique d'un événement rare. Ce problème peut être résolu, au prix de temps de calculs décuplés, en remplaçant le test de Poisson par des procédures de type Monte Carlo ou bootstrap, ainsi que le propose OPENSHAW [18] dans sa version la plus récente de la "Geographical Analysis Machine" (GAM/K). Nous envisageons d'inclure ces alternatives dans nos prochaines investigations.

La méthode d'analyse exploratoire d'un semis de points présentée doit sans doute être affinée, avant de pouvoir être considérée comme réellement opérationnelle, tout au moins dans le domaine de l'accidentologie. Des comparaisons doivent également être menées, aussi bien entre procédures de test (Poisson, bootstrap, Monte Carlo), qu'avec d'autres méthodes, notamment celles qui estiment des surfaces de risque.

Il s'agit néanmoins, nous semble-t-il, d'une contribution utile dans le champ de l'analyse spatiale. L'intégration de ce genre de méthode au sein des systèmes d'information géographique commerciaux nous semble souhaitable. La création récente, au sein de l'équipe THEMA, d'un lien dynamique entre le système d'information géographique Arcview (ESRI) et l'environnement de programmation Xlisp-Stat nous paraît être l'occasion idéale de tester la faisabilité de cette idée.

Références bibliographiques

- [1] BANOS A., BOLOT J., 1999 : Représentation surfacique d'évènements ponctuels discrets. Comparaison méthodologique à partir de l'exemple des accidents de la route dans la CUDL, Actes de colloque, *Quatrièmes Rencontres de Théo Quant*, Besançon, à paraître
- [2] BAILEY T., GATRELL A., 1995 : *Interactive spatial data analysis*, Longman Scientific and Technical, London, 413 pages
- [3] BESAG J., NEWELL J., 1991 : The detection of clusters in rare diseases, *Journal of the Royal Statistical Society, A*, 154, Part 1, pp. 143-155
- [4] BITHELL J.-F., 1990 : An application of density estimation to geographical epidemiology, *Statistics in Medicine*, Vol. 9, pp. 691-701
- [5] BRUNSDON C., 1991 : Estimating probability surfaces in GIS: an adaptive technique, *EGIS '91, Proceedings, Second European Conference in GIS*, Brussels, Belgium, April 2-5 1991, vol. 1, EGIS Foundation, pp. 155-164
- [6] CUZICK J., EDWARDS R., 1990 : Spatial clustering for homogeneous populations, *Journal of the Royal Statistical Society, B*, 52, n° 1, pp. 73-104
- [7] FLEURY D., 1998 : *Sécurité et urbanisme. La prise en compte de la sécurité routière dans l'aménagement urbain*, Presses de l'école nationale des Ponts et chaussées, Paris, 299 pages
- [8] FOTHERINGHAM S., ZHAN B., 1996 : A comparison of three exploratory methods for cluster detection in spatial point patterns, *Geographical Analysis*, Vol. 28, n° 3, pp. 200-218
- [9] GATRELL A., 1994 : Density estimation and the visualization of point patterns, in *Visualization in geographical information systems*, John Wiley and Sons, Chichester, pp. 65-75
- [10] GATRELL A., BAILEY T., DIGGLE P., ROWLINGSON B., 1996 : Spatial point pattern analysis and its application in geographical epidemiology, *Transactions of the Institute of British Geographers*, Numéro spécial 21, pp. 256-274
- [11] GOULD P., 1992 : Epidémiologie et maladie, in *Encyclopédie de géographie*, Economica, Paris, pp. 949-969
- [12] HUGUENIN-RICHARD F., 1999 : Diagnostiquer et estimer le risque routier par une approche géographique. Premières considérations, Actes du colloque, *Quatrièmes Rencontres de Théo Quant*, Besançon, à paraître
- [13] KULLDORFF M., NAGARWALLA N., 1995 : Spatial disease clusters: detection and inference, *Statistics in medicine*, Vol. 14, pp. 799-810

- [14] MARSHALL R., 1991 : A review of methods for the statistical analysis of spatial patterns of disease, *Journal of the Royal Statistical Society*, A, 154, Part 3, pp. 421-441
- [15] NICHOLSON A., WONG Y.D., 1993 : Are accidents poisson distributed? A statistical test, *Accident Analysis and Prevention*, Vol. 25, n° 1, pp. 91-97
- [16] OPENSHAW S., 1995 : Developing automated and smart spatial pattern exploration tools for geographical information systems applications, *The Statistician*, Vol. 44, n° 1, pp. 3-16
- [17] OPENSHAW S., CHARLTON M., WYMER C., CRAFT A., 1987 : A mark 1 geographical analysis machine for the automated analysis of point data sets, *International Journal of Geographical Information Systems*, Vol. 1, n° 4, pp. 335-358
- [18] OPENSHAW S., TURTON I., 1999 : Geographical analysis machine on the internet, <http://www/ccg.leeds.ac.uk/smart/gam/gam.html>
- [19] SABEL C., GATRELL A., 1997 : Exploratory spatial data analysis of motor neurone disease in North West England : beyond the address at diagnosis, GEOMED '97 Paper, <http://www.lancs.ac.uk/postgrad/sabell/rostock.html>, 8 pages
- [20] TIERNEY L., 1990 : *Lisp-Stat, an object-oriented environment for statistical computing and dynamic graphics*, John Wiley and Sons, New York, 397 pages
- [21] THOMAS R.W., 1991 : Quantitative methods : clines, hot spots and cancer clusters, *Progress in Human Geography*, 15, 4, pp. 445-455
- [22] VANDERSMISSEN M.H., POULIOT M., MORIN D., 1996 : Comment estimer l'insécurité d'un site d'accident : état de la question, *Recherche-Transport-Sécurité*, n° 51, pp. 49-60

Notes

- 1 - Le codage de la localisation précise d'un accident est une question difficile [7]. Ce problème technique est en voie de résolution par la diffusion d'un applicatif du SIG MapInfo (Concerto) par le Ministère des transports, qui propose entre autre des procédures de géocodage automatisé par adresse postale des accidents à partir des renseignements fournis dans les Bordereaux d'Analyse des Accidents Corporels (la source d'information nationale). Certaines villes avaient déjà mis au point des saisies spécifiques, comme la Communauté Urbaine de Lille, qui a créé depuis le milieu des années 1980 sa propre base de données en repérant le lieu exact de l'accident sur un référentiel informatisé de son territoire après lecture de son procès-verbal. Nous utilisons des données extraites de cette base de données
- 2 - L'intensité correspond au nombre moyen d'évènements par unité de surface.
- 3 - L'algorithme original intègre des informations de type surfacique, relatives à l'estimation de la population de référence.
- 4 - Le choix de la taille des fenêtres mobiles est crucial, mais difficilement automatisable. Il n'existe pas de critère « objectif » réellement efficace, qui déciderait à la place de l'analyste. Une démarche de type « essai-erreur » est souvent adoptée, comme dans tout problème impliquant des fenêtres mobiles. Dans le cadre de cette application, nous avons commencé à réfléchir à des outils d'aide à la décision souples, susceptibles d'aider l'utilisateur dans ses investigations.