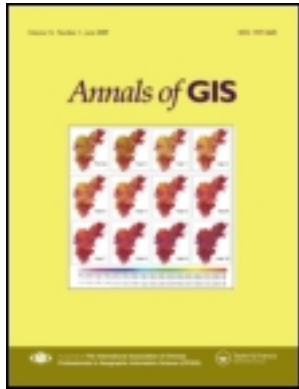


This article was downloaded by: ["Queen's University Libraries, Kingston"]

On: 23 November 2011, At: 06:40

Publisher: Taylor & Francis

Informa Ltd Registered in England and Wales Registered Number: 1072954 Registered office: Mortimer House, 37-41 Mortimer Street, London W1T 3JH, UK



## Annals of GIS

Publication details, including instructions for authors and subscription information:

<http://www.tandfonline.com/loi/tagi20>

### Spatial and temporal aberration detection methods for disease outbreaks in syndromic surveillance systems

Dongmei Chen <sup>a</sup>, John Cunningham <sup>b</sup>, Kieran Moore <sup>c</sup> & Jie Tian <sup>d</sup>

<sup>a</sup> Department of Geography, Queen's University, Kingston, ON, Canada

<sup>b</sup> Leeds, Grenville & Lanark District Health Unit, ON, Canada

<sup>c</sup> Department of Family Medicine, Queen's University, Kingston, Canada

<sup>d</sup> Department of Geology & Geography, Georgia Southern University, Statesboro, GA, USA

Available online: 23 Nov 2011

To cite this article: Dongmei Chen, John Cunningham, Kieran Moore & Jie Tian (2011): Spatial and temporal aberration detection methods for disease outbreaks in syndromic surveillance systems, *Annals of GIS*, 17:4, 211-220

To link to this article: <http://dx.doi.org/10.1080/19475683.2011.625979>

PLEASE SCROLL DOWN FOR ARTICLE

Full terms and conditions of use: <http://www.tandfonline.com/page/terms-and-conditions>

This article may be used for research, teaching, and private study purposes. Any substantial or systematic reproduction, redistribution, reselling, loan, sub-licensing, systematic supply, or distribution in any form to anyone is expressly forbidden.

The publisher does not give any warranty express or implied or make any representation that the contents will be complete or accurate or up to date. The accuracy of any instructions, formulae, and drug doses should be independently verified with primary sources. The publisher shall not be liable for any loss, actions, claims, proceedings, demand, or costs or damages whatsoever or howsoever caused arising directly or indirectly in connection with or arising out of the use of this material.

## Spatial and temporal aberration detection methods for disease outbreaks in syndromic surveillance systems

Dongmei Chen<sup>a\*</sup>, John Cunningham<sup>b</sup>, Kieran Moore<sup>c</sup> and Jie Tian<sup>d</sup>

<sup>a</sup>Department of Geography, Queen's University, Kingston, ON, Canada; <sup>b</sup>Leeds, Grenville & Lanark District Health Unit, ON, Canada;

<sup>c</sup>Department of Family Medicine, Queen's University, Kingston, Canada; <sup>d</sup>Department of Geology & Geography, Georgia Southern University, Statesboro, GA, USA

(Received 7 September 2011; final version received 18 September 2011)

Early surveillance of notifiable infectious diseases is a key element for their control by public health agencies. The goal of syndromic disease surveillance is to identify emerging infectious risks to public health in real or near real time as a method of early detection, trend monitoring, and false-alarm avoidance. This article reviews temporal, spatial, and spatial-temporal aberration detection techniques that can be used to facilitate the early detection of infectious disease outbreaks that can occur in nonrandom yet clustered distributions in geographic information systems (GIS)-based syndromic surveillance systems. The focus is on the approaches appropriate for prospective surveillance data. In addition, this article discusses the impact of data privacy, security, and data quality on detection algorithms and explores what the future GIS-based syndromic surveillance systems may hold.

**Keywords:** infectious disease; syndromic surveillance; GIS; disease outbreak detection; spatial and temporal aberration

### 1. Introduction

World events such as the potential for lethal naturally occurring outbreaks such as severe acute respiratory syndrome (SARS) in Hong Kong and Canada (Lai *et al.* 2004), the threat of bioterror as exemplified by the release of anthrax in the United States in 2001 (Nordin and Al. 2005), the heightened potential for a pandemic influenza outbreak, and the 2009 swine flu outbreak are resulting in an increased awareness among the public of the risk of a devastating infectious disease outbreak in the future. In order to monitor public health in an effective and responsible manner for disease prevention and control in a timely fashion, the development of early warning systems (or surveillance systems) for disease outbreak surveillance has become a priority in many nations. Early detection of these potential health risks requires surveillance systems based on multiple sources of data from a large proportion of the population under surveillance, as well as the capability to statistically analyze spatial and temporal data in real time or near real time (Gestland *et al.* 2003, Nordin and Al. 2005, Yan *et al.* 2008).

Traditionally, disease surveillance relies on routine and manual filing of case reports by hospital and family physicians or clinicians. Disease outbreaks are detected when observed case counts of disease morbidity or mortality

are significantly higher than what are statistically expected based on a previous season's reporting (Gestland and *et al.* 2003, Dafni and *et al.* 2004). The usefulness of this surveillance approach to early detection of infectious disease is limited by the delays that occur due to data acquisition and analysis, laboratory reporting, and physician apprehension to report a disease trend before definitive testing results are obtained (Gestland and *et al.* 2003). The term *syndromic surveillance* refers to the methods that rely on early identification of disease outbreaks by the detection of clinical case features that are apparent before confirmatory diagnoses are made by physician or laboratory analysis (Dafni *et al.* 2004, Henning 2004, Mandl 2004). A syndromic surveillance is usually assisted by computer-based data acquisition and generation of statistical alerts in real-time or near-real-time outbreaks of disease and includes a functional capacity for data collection, analysis, and dissemination linked to public health programs (Henning 2004). Syndromic surveillance systems rely on a continuous and real-time acquisition of health data via the use of automated protocols. This type of alert system aids the frontline physician in a hospital emergency department or walk-in clinic, who is only alert to individual cases rather than the overall spatial and temporal patterns of disease that patients presenting across an entire geographic region may represent (Kamel-Boulos 2004, Mandl 2004).

\*Corresponding author. Email: chendm@queensu.ca

Many types of syndromic surveillance systems have been operating in various jurisdictions across different nations (Gestland *et al.* 2003, Lombardo *et al.* 2003, Tsui *et al.* 2003, Cooper 2007, Lombardo and Buckeridge 2007, Moore *et al.* 2008). According to a study conducted by the Center for Disease Control and Prevention (CDC) of United States, there are more than 100 sites in United States alone that implemented syndromic surveillance systems (Buehler *et al.* 2003). These systems vary widely in terms of the system architecture, information processing and management techniques, characteristics and sources of the data, methods employed to collect surveillance data, geographic coverage, disease focuses, and the analytical methodology used to determine whether a disease outbreak has occurred (Bravata *et al.* 2004, Yan *et al.* 2008). More detailed comparison and review of different syndromic surveillance systems can be found in Yan *et al.* (2008).

The application of geographic information systems (GIS) and spatial analysis to health-related research is increasing at a rapid rate. GIS is usually applied to single-issue, isolated, and time-limited etiological research that is retrospective in nature rather than to real-time prospective data associated with health planning, coordinating, or protection activities (Kamel-Boulos 2004, Moore *et al.* 2008). Geographic location is a key component in 80–90% of all data generated by the government and health sectors in the United States and Canada (Kamel-Boulos 2004, Kennedy *et al.* 2008). The advent of affordable and quick geocoding of addresses allows for surveillance on a finer spatial scale than it had ever been in the past. Therefore, the processes used to transform these data into useful information will require the systematic use of GIS and spatial analysis methods in linking, referencing, visualizing, and analyzing the various existing and yet-to-be-identified disease-related surveillance datasets (Yan *et al.* 2008).

Many different disease aberration detection algorithms are used in operating syndromic surveillance systems worldwide (Dafni *et al.* 2004, Lombardo and Buckeridge 2007). These algorithms are utilized as statistical alert mechanisms to generate an alarm whenever the observed frequency of disease occurrences exceeds what is expected and help, therefore, in the detection of a signal pattern of disease in human populations that can be distinguished from what is considered to be background noise (Tsui *et al.* 2003). The compromise with disease detection algorithms is that they provide a timely and cheap, but noisy, approximation of what could be obtained from the manual analysis of patient data that would only be available in the late stages of an epidemic disease outbreak (Tsui *et al.* 2003, Wong *et al.* 2003). Although much geographic information is collected by surveillance systems, the most commonly used aberration methods are concerned with detecting disease anomaly in time more than in space. The lack of effective spatial detection algorithms causes a deficit of the true potential and abilities of detecting outbreaks in syndromic

surveillance systems. With the increasing availability of the more GIS-based syndromic systems, effective spatial and spatial–temporal aberration detection approaches need to be tested and implemented.

To help the implementation of spatial and temporal detection of disease outbreaks in syndromic surveillance systems, this article provides a basic review on aberration detection methods that syndromic surveillance systems currently employ and discusses potential spatial and spatial–temporal methods that can be included in GIS-based syndromic systems for early detection of infectious disease outbreaks. Issues, limitations, challenges, and future research of detection methods and their potential in syndromic surveillance systems are also addressed.

## 2. Spatial–temporal aberration detection methods

Quickly detecting spatial, temporal, and spatial–temporal pattern change is the key for real-time surveillance. Aberration detection algorithms are statistical tools used as alert mechanisms to generate an alarm whenever the observed frequency of disease occurrences exceeds what is expected and, therefore, to detect a signal pattern of disease in human populations that can be distinguished from what is considered to be background noise (Tsui *et al.* 2003). These algorithms detect disease in temporal and/or spatial excess from what is normally expected and form a *cluster*. These clusters are characteristic of the distribution of a particular disease over a defined region during a specified period of time or region (Quataert *et al.* 1999, Rogerson and Yamada 2009). A temporal cluster would be an excess of disease frequency during a particular time period, whereas a spatial cluster would be an excess of disease found within a specific geographic area. A syndromic surveillance system is maximized when both spatial and temporal cluster detection algorithms work in concert to recognize emerging infectious disease at an early stage of infectivity (Kleinman *et al.* 2004).

### 2.1. Temporal aberration detection methods

Time series methods based on autoregressive integrated moving average (ARIMA) are one of the simplest methods that have been applied to detect disease outbreak (Reis and Mandl 2003). Such methods are available in commonly used statistical packages (e.g. SAS software) and are easy to use. One limitation of the ARIMA methods is their difficulty in updating model parameters when new data come in.

The *recursive least squares* (RLS) is another simple algorithm based on autoregressive linear models now being used in surveillance systems. The purpose of the RLS algorithm is to minimize the sum of the absolute squares of difference between the desired number of disease cases and the actual cases (Haykin 2002). This algorithm's strength

is that it uses recent historical data (for the past few days) rather than historical data collected over the previous year to derive its signal spikes for each syndrome category (e.g., gastrointestinal or respiratory) and can thereby be useful in situations where there is a high turnover of population through a region (Moore and Al. 2000, Das *et al.* 2003, Gestland *et al.* 2003). The RLS algorithm generates alerts when the observations of disease incidence exceed the 95% confidence interval of a predicted disease incidence.

Another common algorithm used in most surveillance system is the cumulative sum (CUSUM), which calculates the CUSUM of differences between the sequential observations of a variable from an expected mean. In surveillance systems, CUSUM compares the mean disease events observed during the previous several days or hours to a moving time window of means based on the previous 6 months to 2 years (Moore and Al 2000, Das *et al.* 2003, Moore *et al.* 2008). This allows for the correction of data for seasonal variation in disease incidence (Mandl 2004). When the number of disease cases is significantly above the average during a period of time (e.g. 7 days), the amount added to the CUSUM will be positive and the sum will steadily increase. A segment of the CUSUM will have an upward slope indicating a period in which disease appearance frequency is above average.

*What is strange about recent events* (WSARE) employs a rule-based technique that compares recent disease data to those from a baseline distribution and finds recent periods changed the most (Wong *et al.* 2005). It takes advantage of the wealth of spatial, temporal, demographic, and symptomatic information in the surveillance system and investigates patient records and other patient characteristics, including geographic information, to determine whether there are statistical anomalies among patients with these characteristics on a given day when compared to patients on the same day using historical data (Wong *et al.* 2003, 2005). The WSARE algorithm incorporates a wide range of ideas, including association rules, Bayesian networks, hypothesis testing, and permutation tests to evaluate the significance of the aberration.

Many other algorithms exist but are not described here (e.g., hidden Markov model-based approaches, Wavelet, Pulsar, LOWESS, and TAD). A general description of these temporal algorithms can be found in Sonesson and Brock (2003) and Dafni *et al.* (2004). It should be noted that all the surveillance algorithms described above are based on temporal statistical methods only, with the exception of WSARE. Surveillance using purely temporal techniques is dependent on the size of the geographic unit area in which the disease surveillance counts originate. If the surveillance region is too small, the number of incidents may also be small, leading to a low algorithm sensitivity and subsequent exclusion of outbreak cases. Conversely, if the region represented in the time series is too large, early outbreak cases may be hidden in the background statistical noise.

## 2.2. Spatial and spatial-temporal detection approaches

Disease outbreaks are likely presented as clusters in space, time, or both. Most surveillance systems focus more on temporal aberration detection and GIS-based spatial mapping and visualization (Moore 2004, Moore *et al.* 2008). The continued lack of effective spatial detection algorithms creates a real deficit in the true potential and abilities of syndromic surveillance systems. Temporal statistical methods are often unable to detect the spatial clusters, especially when the number of cases is small and the surveillance region is large. Many natural disease outbreaks or biological attacks are typically started at a local spatial scale first. The statistics of clustering of health-related events have been studied for many decades (Marshall 1991, Anselin 1995, Ord and Getis 1995, Kulldorff 1997, Sonesson and Brock 2003, Lawson 2006). They are generally intended to tackle two issues: (1) Is there an overall tendency for clustering in the region? (2) If yes, where do clusters specifically occur? The first question can often be measured by global spatial statistics (Lee and Rogerson 2007), while the second is usually answered by local statistics. However, many traditional spatial clustering methods as reviewed by Lawson (2006) are developed for *retrospective* data and cannot be directly applied for spatial cluster detection in the *prospective* surveillance data, which requires an ongoing effort to detect spatial pattern change when new data become available (Rogerson 1997, Rogerson and Yamada 2004). In the following sections, our focus will be on spatial mapping and cluster detection approaches appropriate for *prospective* data although some approaches for retrospective data are discussed. Table 1 summarizes the various spatial and spatial-temporal aberration methods proposed and tested in the literature.

The earliest attempt to develop surveillance approaches for detecting spatial disease clustering can be traced back to a sequential test developed by Bartholomew (1956). Due to its complexity and underlining hypothesis, it has not received much attention (Rogerson 1997). Openshaw *et al.* (1988) proposed a geographical analysis machine (GAM) by constantly searching clusters of spatial locations with an excess of cases at a defined local significance level. Although this approach initially was not designed for surveillance data, it can be used for spatial surveillance. However, its computation demand is high when multiple testing is involved.

Rogerson (1997) was the first researcher to bring attention to spatial pattern analysis in surveillance data. He developed a CUSUM statistic that combines a modified Tango's statistic and general and focused tests of clustering to monitor changes in spatial patterns. Its shortcoming was that the approach was based on an assessment of global pattern change rather than on local changes. Rogerson suggested that CUSUM analysis of some local statistics such as local Moran's *I* or Getis's *G* may offer more promise

Table 1. The summary of spatial and spatial–temporal aberration detection methods in the literature.

Methods/algorithms	Detection type	Applications and limitations	References
A sequential test	Spatial	Tested for simulated data	Bartholomew (1956)
Geographical analysis machine	Spatial	Retrospective data; high computation cost	Openshaw <i>et al.</i> (1988)
Modified global CUSUM statistics	Spatial	Easy to implement; tested for retrospective data	Rogerson (1997, 2001), Lawson and Kleinman (2005)
Multivariate CUSUM approaches; univariate CUSUM approaches	Spatial	Tested for retrospective and prospective data; can combine multiple variables	Rogerson and Yamada (2004, 2009), Moore and AI (2000)
GLMM	Spatial	Tested for prospective surveillance data; sensitive to a small number of spatial focused cases	Lazarus <i>et al.</i> (2002); Kleinman <i>et al.</i> (2004, 2005)
SMART	Spatial	Tested for prospective surveillance data; seasonal, weekly effects and other factors can be adjusted during regression	Bradley <i>et al.</i> (2005) Yih <i>et al.</i> (2005)
M-Statistic	Spatial	Prospective hospital data	Olson <i>et al.</i> (2005)
Scan statistics	Spatial; spatial–temporal	Well tested for both retrospective and prospective data; the geometric shape of the clusters is limited; computation cost can be high	Kulldorff (1997, 2001); Das <i>et al.</i> (2003), Hefferman <i>et al.</i> (2004), Kulldorff (2011), Kleinman <i>et al.</i> (2005)
Bayesian spatial scan statistics	Spatial; spatial–temporal	Tested in prospective surveillance data; computationally efficient; can incorporate prior knowledge	Neill <i>et al.</i> (2005)
RSVC	Spatial	Tested for retrospective spatial and temporal data; can adjust the shape of clustering	Zeng <i>et al.</i> (2004), Chang <i>et al.</i> (2005)
WSARE	Temporal and spatial–temporal	Tested for prospective surveillance data; can use representative features for monitoring; baseline distribution needs to be known	Wong <i>et al.</i> (2003, 2005)
PANDA	Spatial–temporal	Tested for prospective surveillance data; can handle multiple variables; computationally extensive	Moore <i>et al.</i> (2000), Cooper <i>et al.</i> (2004)

Note: CUSUM, cumulative sum; GLMM, generalized linear mixed model; SMART, small-area regression and testing; RSVC, risk-adjusted support vector clustering; WSARE, what is strange about recent events; PANDA, population-wide anomaly detection and assessment.

than the use of global measures (Rogerson and Yamada 2009).

Rogerson (2001) employed a combination of CUSUM and the Knox test for space–time interactions to a group of point events. The Knox test is based on counting geographic event pairs (points) that occur within a predetermined critical distance and time. The application of these procedures to the occurrence of Burkett’s lymphoma in the West Nile region of Uganda helped researchers determine that there was evidence for significant space–time clustering of disease cases at a particular location (Rogerson 2001). This study was retrospective in nature, using data from 1961 to 1975; however, the combination of Knox and CUSUM methodologies can be used in a prospective analysis of emerging spatial patterns. Later, Rogerson and Yamada (2004) compared the univariate and multivariate CUSUM approaches to monitor changes in spatial patterns of disease using both simulated data and country-level breast cancer data in the northeastern United States. They concluded that the multivariate CUSUM method was preferred to detect spatial pattern when the degree of spatial autocorrelation was strong, while the multiple univariate CUSUM generally performed better at detecting changes

in rates that occur in a small number of regions and spatial autocorrelation was weak or was poorly estimated.

A different approach to spatial and temporal cluster detection was used by Lazarus *et al.* (2002), who employed a generalized linear mixed model (GLMM or logistic regression) of hospital visits and tele-health calls with geocoded personal addresses to detect potential clusters of bioterror-related anthrax cases in eastern Massachusetts. The main advantage of using a GLMM is its predictive ability built on historical data. The results of the model can be presented using a probability-based metric that is adjusted for covariates and unmeasured features. This allowed the model to be useful for predicting the number of expected disease cases for a particular day of the year in advance. A major recommendation of the authors was that more data sources (absenteeism, over-the-counter sales) added to the model would make an important supplement to the system in terms of predictive ability. Kleinman *et al.* (2004) also used a combined algorithm of GLMM and CUSUM to detect incident disease clusters in small areas. This small-scale approach is useful in the detection of a deliberate release of a biological agent as well as the detection of other small-scale illnesses such as food poisoning (Kleinman

*et al.* 2004). This research emphasizes the importance of both temporal and spatial cluster analyses as a more effective method in syndromic surveillance than in temporal analysis alone.

Small-area regression and testing (SMART) is also a regression analysis modified from GLMM (Bradley *et al.* 2005). The SMART model generates a predicted value and the difference between the predicted and observed values based on approximately 2 years of baseline data and parameters for day of the week, holiday, day after a holiday, sine, and cosine seasonal terms. These values are transformed into a statistical measure of the time period in which one is interested to determine the recurrence interval (RI) for the observed count. RI values are then entered into SMART score groups to determine the statistical significance of outbreaks in each small area. Past experiment study suggested that SMART gave slightly inferior results to the spatial scan statistic method although both methods achieved satisfied performances (Kleinman *et al.* 2005).

Instead of using temporal analysis and analyzing spatial data at the scale of census tract, Olson *et al.* (2005) used individual patient addresses and translated them into longitude and latitude coordinates using ArcGIS (Olson *et al.* 2005). A comparison metric known as M-Statistic was then used to determine whether discrepancies existed in an inter-point distance distribution. Here, distance values between what was expected based on a baseline historical patient spatial distribution and what were observed from real-time hospital patient data were determined. Olson *et al.* (2005) found that this methodology was sensitive in detecting spatial clustering even when the outbreak was small, but that the system would need to be adjusted for the seasonality of the disease under surveillance.

Scan statistics and its variations have been widely used in syndromic surveillance systems to detect and evaluate spatial, temporal, or spatial-temporal clusters (Kulldorff 1997, Yan *et al.* 2008). The basic principle is to gradually scan a window across time and/or space, comparing the number of observed and expected observations inside the window at each location. In theory, the scanning window can be either an interval (in time), a circle or an ellipse (in space), or a cylinder with a circular or elliptical base (in space-time) (Kulldorff 2011). Monte Carlo simulations are performed to rank the likelihood values so that a *P*-value (ranking position) can be assigned to each cluster to indicate its statistical significance. The main advantages of scan statistic methods include that it avoids preselection bias on size or location of clusters and can be easily adjusted for non-uniform population density as well as other factors. The scan statistic was originally developed for retrospective analysis of chronic disease but has since been adapted for infectious disease surveillance systems (Kulldorff 2001, Das *et al.* 2003, Hefferman *et al.* 2004). With the availability of SatScan, a freely downloadable software that implements various types of

scan statistics, these scan statistic approaches have been widely applied to different surveillance systems. Duczmal and Buckeridge (2005) proposed a modified spatial scan statistic by considering work-related factors. The results show that the modified approach achieved greater detection power than the scan statistics if not considering people movement.

Several known problems exist in spatial scan methods. First, they can only detect clusters in simple regular shape. Second, it is difficult to incorporate prior knowledge such as the size or shape of the outbreaks or the impact of infection rate. Third, local cluster search over a large region can be computationally expensive. Neill *et al.* (2005) proposed a Bayesian spatial scan statistic to address the second and third problems. The Bayesian spatial scan methods are more computational efficient and can combine a prior knowledge of an outbreak. Instead of using the Poisson method in original spatial scan statistic, a conjugate gamma-Poisson model is used to produce a spatial map showing the posterior probabilities of an outbreak in each location.

Besides the above commonly used approaches for *prospective* data, many traditional and advanced global and local spatial clustering techniques for *retrospective* data could be adopted to surveillance data. A summary of the traditional methods can be found in Lawson (2006). These methods can be used to measure the spatial association and arrangement of outbreak locations. During a communicable disease outbreak, a patient is more likely to infect people near to him/her. This relationship may be described as an autocorrelation in space and used to indicate a spatial clustering (or a hot spot) of disease outbreaks. These global and local spatial clustering methods could be used to measure the spatial association of outbreak locations and their temporal change to identify hot spots of a disease (Lee and Rogerson 2007, Rogerson and Yamada 2009).

More advanced hot spot analysis approaches, including risk-adjusted support vector clustering (RSVC), prospective support vector clustering (Chang *et al.* 2005), and space-time correlation analysis (Ma *et al.* 2006), have been developed to deal with the limitation of regular clustering shape in many spatial detection approaches for disease outbreak detection. Zeng *et al.* (2004) proposed an RSVC approach that can combine the risk adjustment idea with a robust support vector clustering (SVC) method to improve the detection performance for retrospective spatial-temporal data. The RSVC method uses a dynamically adjusted Gaussian kernel function to find a hypersphere with a minimal radius to contain most of the data to indicate high-risk outbreak areas (Zeng *et al.* 2004). These approaches have been implemented in BioPortal, a biosurveillance system developed by the State of Arizona Department of Health Services and hospitals in Taiwan (Zeng *et al.* 2005, Chen and Xu 2006).

The majority of the spatial and spatial–temporal detection methods are designed to monitor spatial pattern change in individual data source. In order to take into account data from multiple sources and variables, a set of spatial–temporal detection approaches has been developed based on the theory of Bayesian network. Besides the WSARE discussed before, another example is population-wide anomaly detection and assessment (PANDA) algorithm built on a causal Bayesian network-based model to analyze the spatial–temporal probability of disease in a population as a whole (Cooper *et al.* 2004). The advantage of this algorithm is its capability of building a large set of inter-linked patient-specific probabilistic causal models to include various risk factors. A previous study showed that this algorithm can handle a population size of 1.4 million (Cooper *et al.* 2004). However, an in-depth investigation and comparison on the gain of using information from different data sources is needed.

### 3. Issues and challenges in spatial and spatial–temporal aberration detection

The key issue in spatial and spatial–temporal aberration detection is the algorithm sensitivity, specificity, and timeliness of detecting a potential outbreak (Hutwagner and Browne *et al.* 2005, Rogerson 2009). To be useful, detection algorithms need to send an alarm in a promptly fashion. The earlier a true alarm is sent, the more helpful it is to the outbreak monitor and control. Although there are many spatial and temporal detection approaches existing, there is no definitive proof that these approaches can catch all early disease outbreaks due to the issues ranging from data quality to algorithm effectiveness. The determination of appropriate alarm thresholds for each of these spatial approaches is a pivotal issue. In addition, interpretation of spatial disease distributions can be tricky as several different spatial detection algorithms may result in a similar disease pattern (Graham *et al.* 2004, Rogerson and Yamada 2009).

The use of spatial detection approaches in syndromic surveillance depends heavily on the access to useable geographic data. Most syndromic surveillance systems utilize data that were not originally created for surveillance purposes. In syndromic surveillance, it is not necessary or possible to have perfect or error-free data. Most syndromic surveillance systems generate early detection alarms from statistically noisy data. Inaccurate or incomplete geocoding of data has been addressed considerably in the literature (Hurley *et al.* 2003, Elliot and Wartenberg 2004, Moore *et al.* 2008). The result of poor data quality can be either false alarms due to an artificial amplification of a disease signal or the suppression of a disease signal when a true outbreak or biological attack occurs (Pavlin 2003). The compromise with disease detection algorithms is that they provide a timely and cheap, but noisy, approximation of what could be obtained from the manual analysis

of patient data that would only be available in the late stages of an epidemic disease outbreak (Tsui *et al.* 2003, Wong *et al.* 2003). Thus, syndromic surveillance cannot replace the requirement for epidemiologists to be on the ground tracking down clinical and exposure information. Epidemiologists must investigate every alarm to confirm whether the alarming signal is a result of valid data or misclassification. This requires the examination of individual cases from the cluster that triggered the syndromic surveillance alarm so that demographic and geographic information can be used to identify potential data misclassification (Duchin 2003).

Privacy and confidentiality concerns are always an issue impacting the accurate analysis of spatial and spatial–temporal detection methods (Kamel-Boulos 2004, Moore *et al.* 2008). There is a tension between the need to protect personal information and the requirement to use that data for the public good, as is the case in public health syndromic surveillance systems (Gestland *et al.* 2003). Although much detailed geographic data are collected in surveillance systems, some of these data must be aggregated to much coarser scale due to federal, provincial, and institutional privacy protection laws (Gestland *et al.* 2003, Moore *et al.* 2008). Data aggregation will produce ‘spatial uncertainty’ due to the masking of the outbreak signal of small clusters within aggregate data or the dilution of clusters that straddle borders between separate data aggregation areas (Jacquez 2000, Olson *et al.* 2005). This type of system will result in a loss of information and will not have the ability to identify disease clusters as readily as systems that correlate personal information and geographic place and time. The results are a lower sensitivity and specificity of the system and consequently higher likelihood to sound false alarms or fail to detect disease outbreak clusters. So far, no comprehensive tests have been conducted to evaluate the spatial and temporal sensitivity of various detection methods under different aggregation levels. Our past initial simulation showed that aggregation caused the delay in terms of spatial–temporal aberration alarming by different spatial and temporal detection algorithms in the surveillance system (Tian *et al.* 2006). More systematic tests are needed. How much impact this time delay to the infectious disease control will have depends on the disease’s incubation period and spreading speed of the disease. If it is a high contagious disease and has a short incubation period, this delay may cost valuable acting time for public health officials to control it.

The uncertainty of spatial detection algorithms caused by the variation of spatial unit sizes (such as postal/zip codes or census tracts) used in surveillance systems is also not clear. For example, in rural areas, due to the large area each postal code covers, the location of rural postal stations cannot represent the location of people living in that region at all. This would cause difficulties in outbreak detection, as there may not be a strong clustering of postal/zip codes

or hospital presentations of disease cases. Pinpointing the geographic location of the outbreak can therefore be challenging. When several random cases from the same rural postal/zip code are reported, the spatial aberration detection tends to give a false alarm. Different from urban areas, a spatial cluster size should be settled much larger when used in the rural area. This variation causes some challenges when using spatial aberration detection in the surveillance system since most spatial detection algorithms cannot automatically adjust the analysis window for different regions.

It is the trend in surveillance systems that more non-traditional data from multiple sources will be entered into the system. For example, analyzing patients' recent activities became quickly evident that SARS victims had recently visited Hong Kong or had been in contact with people who had done so (Lai *et al.* 2004). In the early stage of the swine flu outbreak, most patients in Canada had direct or indirect links to people who visited Mexico. As well, emergency (911) and tele-health callers and retail customers can be linked temporally and spatially through the use of their address, telephone number, or postal code (Feinberg and Shmueli 2005, Cooper 2007, Moore *et al.* 2008). Increased sales from some product groups can be good indicators of disease number and type. Absenteeism is another source of data that can be used (Lombardo and *et al.* 2003).

The combination of all these data sources can increase the strength of a signal to be detected from the background noise by surveillance algorithms. However, there are several challenges to the detection algorithms when dealing with a large volume of data from different sources. Most algorithms existing in the current surveillance systems are not designed to handle these nontraditional data and cannot utilize this information (Wong *et al.* 2005, Yan *et al.* 2008). Another issue is that current algorithms cannot handle the high level of uncertainty embedded in these nontraditional data. More robust temporal and spatial detection approaches that can integrate information from different data sources are needed in future surveillance systems. Some nontraditional approaches developed in spatial-temporal analysis of other geographic fields, such as artificial neural networks in machine learning, Bayesian spatial-temporal approaches (Chapman *et al.* 2005, Lawson 2006), or a hierarchical system of sequential alerts (Feinberg and Shmueli 2005), are recognized as being superior in extracting information from multivariate databases. In recent years, knowledge-based data mining and data-fusion approaches have also been tested in detecting and predicting disease outbreaks from multiple data sources in syndromic systems (Yan *et al.* 2008).

#### 4. Summary and conclusions

Public health syndromic surveillance allows for stakeholders and policymakers to estimate the magnitude and

distribution of a potential infectious disease outbreak in real time or near real time. Testing some surveillance systems against yearly disease cycles (e.g., influenza) has shown that the systems do validly model and predict these outbreaks (Platt and Boccino *et al.* 2003, Lombardo and Buckeridge 2007). The growing outbreak cluster can actually be recognized by a series of spatial detection with expanding time frames. Typically, syndromic surveillance systems and their associated stakeholders cut across jurisdictional boundaries (hospital, city, county, public health unit, province, federal) and, therefore, require the cooperation of these various entities (Gestland *et al.* 2003). Cooperation and facilitation are required in order for such a system to function efficiently. This is a challenge that has been gradually overcome in many jurisdictions across Canada, the United States, and Europe in recent years (Moore 2004, Rolfhamre and Grabowska *et al.* 2004, Lombardo and Buckeridge 2007).

Disease processes result in disease patterns that are both complex in nature and may operate over a range of spatial and temporal scales to produce an array of intricate patterns of disease incidence during an outbreak (Graham *et al.* 2004). Although monitoring changes in the spatial pattern of diseases is extremely useful for early, rapid detection of infectious disease outbreaks, many of the methods reviewed above have not become routines in current surveillance systems. Interpretation of results from spatial and temporal algorithms can be tricky as different algorithms may result in different conclusions for a same disease pattern. Most surveillance systems focus more on temporal aberration detection and GIS-based spatial mapping and visualization (Moore 2004, Moore *et al.* 2008). The continued lack of effective spatial detection algorithms creates a real deficit in the true potential and abilities of syndromic surveillance systems.

This article gives a review on various spatial and spatial-temporal aberration detection methods that have been used in various systems or have the potential of being used in GIS-based surveillance systems. It should be noted that the disease detection algorithms cannot guarantee an accurate detection for an epidemic disease outbreak, other than a timely, relatively cheap, and approximate alarm. Thus, syndromic surveillance cannot replace the groundwork of epidemiologists to track down clinical and exposure information.

With the increasing availability of geographic information in the syndromic systems, the spatial and spatial-temporal cluster detection should play more important roles in providing early warning of disease outbreaks. Algorithm sensitivity and specificity are ongoing issues in the use of spatial and spatial-temporal detection algorithms in the syndromic surveillance systems (Dafni *et al.* 2004). Most of these spatial and spatial-temporal techniques are sensitive to spatial and/or temporal units used in the analysis. Although there are many spatial and temporal



detection approaches existing, there is no definitive proof that these approaches can catch all early disease outbreaks due to the issues of data quality and other issues in the systems. Algorithm sensitivity needs to be considered and tested before these algorithms are implemented in GIS-based syndromic systems so that the balance between early outbreak detection and false alarms is optimized as these alarms are costly in terms of financial and human resources to investigate. More research on the impact of different algorithm parameters, spatial units, and spatial and temporal uncertainty on the performance of detection methods will be needed before they can be practically implemented into surveillance systems for disease monitoring and health policy decision.

As the surveillance system is expanded to cover larger regions, as well as to take more data sources, more robust temporal and spatial detection approaches that can integrate information from multiple different data sources are needed in the future surveillance systems. Some nontraditional approaches developed in spatial-temporal analysis of other geographic fields, such as artificial neural networks in machine learning, Bayesian spatial-temporal approaches, knowledge-based data mining, and data-fusion approaches, should be a research direction in the future.

### Acknowledgements

This research is partially supported through a project funded by the Geomatics for Informed Decision (GEOIDE), National Centre of Excellence, Canada. We thank the anonymous reviewers and Professor Ian MacLachlan for their constructive comments and suggestions for improving this article.

### References

- Anselin, L., 1995. Local indicators of spatial association – LISA. *Geographical Analysis*, 27, 93–115.
- Bartholomew, D.J., 1956. A sequential test of randomness for events occurring in time or space. *Biometrika*, 43, 64–78.
- Bradley, C.A., et al., 2005. BioSense: implementation of a national early event detection and situational awareness system. *MMWR (CDC)*, 54 (Suppl), 11–20.
- Bravata, D.M., et al., 2004. Systematic review: surveillance systems for early detection of bioterrorism-related diseases. *Annals of Internal Medicine*, 140, 910–922.
- Buehler, J., et al., 2003. Syndromic surveillance and bioterrorism-related epidemics. *Emerging Infectious Diseases*, 9 (1), 197–204.
- Chang, W., Zeng, D., and Chen, H., 2005. Prospective spatio-temporal data analysis for security informatics. In: *Proceedings of the 8th IEEE international conference on intelligent transportation systems*, September, Vienna, Austria. Washington, DC: IEEE Computer Society, TCD-CS-2005-45.
- Chapman, W.W., et al., 2005. Classifying free-text triage chief complaints into syndromic categories with natural language processing. *Artificial Intelligence in Medicine*, 33 (1), 31–40.
- Chen, H. and Xu, J., 2006. Intelligence and security informatics. *Annual Review of Information Science and Technology (ARIST)*, 40, 229–299.
- Cooper, D., 2007. Case study: use of tele-health data for syndromic surveillance in England and Wales. In: J.S. Lombardo and D.L. Buckeridge, eds. *Disease surveillance: a public health informatics approach*. Englewood Cliffs, NJ: Wiley-Interscience, 335–362.
- Cooper, G.F., et al., 2004. Bayesian biosurveillance of disease outbreaks. In: *Proceedings of the twentieth conference on uncertainty in artificial intelligence*. Arlington, VA: AUA Press, 94–103.
- Dafni, U.G., et al., 2004. Algorithm for statistical detection of peaks – syndromic surveillance system for the Athens 2004 Olympic games. *Morbidity and Mortality Weekly*, 53 (Suppl), 86–94.
- Das, D., et al., 2003. Enhanced drop-in syndromic surveillance in New York City following September 11, 2001. *Journal of Urban Health*, 80, i76–i88.
- Duchin, J.S., 2003. Epidemiological response to syndromic surveillance signals. *Journal of Urban Health*, 80 (2), i115–i116.
- Duczmal, L. and Buckeridge, D., 2005. Using modified spatial scan statistic to improve detection of disease outbreak when exposure occurs in workplace – Virginia, 2004. *MMWR (CDC)*, 54 (Suppl), 187.
- Elliot, P. and Wartenberg, D., 2004. Spatial epidemiology: current approaches and future challenges. *Environmental Health Perspectives*, 9, 998–1006.
- Feinberg, S.E. and Shmueli, G., 2005. Statistical issues and challenges associated with rapid detection of bio-terrorist attacks. *Statistics in Medicine*, 24, 513–529.
- Gestland, P.H., et al., 2003. Automated syndromic surveillance for the 2002 winter Olympics. *Journal of the American Medical Informatics Association*, 10, 547–554.
- Graham, A.J., Atkinson, P.M., and Danson, F.M., 2004. Spatial analysis for epidemiology. *Acta Tropica*, 91, 219–225.
- Haykin, S., 2002. *Adaptive filter theory*. Englewood Cliffs, NJ: Prentice Hall.
- Hefferman, R., et al., 2004. Syndromic surveillance in public health practice, New York City. *Emerging Infectious Disease*, 10, 858–864.
- Henning, K.J., 2004. Overview of syndromic surveillance: What is syndromic surveillance? *Morbidity and Mortality Weekly*, 53, 5–11.
- Hurley, S.E., et al., 2003. Post office box addresses: a challenge for geographic information systems-based studies. *Epidemiology*, 14 (4), 386–391.
- Hutwagner, L., et al., 2005. Comparing aberration detection methods with simulated data. *Emerging Infectious Disease*, 11 (2), 314–316.
- Jacquez, G.M., 2000. Spatial analysis in epidemiology: nascent science or a failure of GIS? *Journal of Geographical Systems*, 2, 91–97.
- Kamel-Boulos, M.N., 2004. Towards evidence-based, GIS-driven national spatial health information infrastructure and surveillance services in the United Kingdom. *International Journal of Health Geographics*, 3, 1.
- Kennedy, S., et al., 2008. The need for a national emergency health services database. *Canadian Journal of Emergency Medicine*, 10, 120–124.
- Kleinman, K., Lazarus, R., and Platt, R., 2004. A general linear mixed models approach for detecting incident clusters of disease in small areas, with an application to biological terrorism. *American Journal of Epidemiology*, 159, 217–224.

- Kleinman, K., *et al.*, 2005. A model-adjusted spacetime scan statistic with an application to syndromic surveillance. *Epidemiological Infection*, 119, 409–419.
- Kulldorff, M., 1997. A spatial scan statistic. *Communications in Statistics: Theory and Methods*, 26, 1481–1496.
- Kulldorff, M., 2001. Prospective time periodic geographical disease surveillance using a scan statistic. *Journal of the Royal Statistical Society, Series A*, 164, 61–72.
- Kulldorff, M., 2011. SaTScan User Guide v9 www.satscan.org
- Lai, P.C., *et al.*, 2004. Understanding spatial clustering of severe acute respiratory syndrome (SARS) in Hong Kong. *Environmental Health Perspectives*, 112, 1550–1556.
- Lawson, A., 2006. *Statistical methods in spatial epidemiology*. New York: Wiley.
- Lawson, A.B. and Kleinman, K., 2005. *Spatial & syndromic surveillance for public health*. Chichester, UK: John Wiley & Sons.
- Lazarus, R., Kleinman, K., Dashevsky, I., Adams, C., Kludt, P., DeMaria, A., and Platt, R. 2002. Use of automated ambulatory-care encounter records for detection of acute illness clusters, including potential bioterrorism events. *Emerging Infectious Disease*, 8 (8), 753–760.
- Lee, G. and Rogerson, P.A., 2007. Monitoring global spatial statistics. *Journal of Stochastic Environmental Research and Risk Assessment*, 21, 545–553.
- Lombardo, J., *et al.*, 2003. A systems overview of the electronic surveillance system for the early notification of community-based epidemics (ESSENCE II). *Journal of Urban Health*, 80, i32–i42.
- Lombardo, J.S., and Buckeridge, D.L., eds., 2007. *Disease surveillance: a public health informatics approach*. New York: Wiley-Interscience.
- Ma, J., Zeng, D., and Chen, H., 2006. Spatial-temporal cross-correlation analysis: a new measure and a case study in infectious disease informatics. In: *Proceedings of the IEEE Intelligence and Security Informatics Conference 2006*, 23–24 May, San Diego, CA. Lecture Notes in Computer Science. Berlin: Springer, 542–547.
- Mandl, K.D., 2004. Implementing syndromic surveillance: a practical guide informed by early experience. *Journal of the American Medical Informatics Association*, 11, 141–150.
- Marshall, R.C., 1991. A review of the statistical analysis of spatial patterns of disease. *Journal of the Royal Statistical Society, Series A*, 154, 421–441.
- Moore, K., 2004. Real-time syndromic surveillance in Ontario, Canada: the potential use of emergency departments and telehealth. *European Journal of Emergency Medicine*, 11, 3–11.
- Moore, K. and Al, E., 2000. *Real-time outbreak and disease surveillance (RODS) for use as an Ontario syndromic surveillance tool: a technical overview*. Kingston, Frontenac, Lennox: Emergency Department Syndromic Surveillance Team and Addington Health Unit, 1–15.
- Moore, K.M., Edgar, B.L., and McGuinness, D., 2008. Implementation of an automated, real-time public health surveillance system linking emergency departments and health units: rationale and methodology. *Canadian Journal of Emergency Medicine*, 10, 114–119.
- Neill, D., Moore, A., and Cooper, G., 2005. A Bayesian spatial scan statistic. *Neural Information Processing Systems*, 18, 1003–1010.
- Nordin, J.D. and Al, E., 2005. Simulated anthrax attacks and syndromic surveillance. *Emerging Infectious Disease*, 11, 1394–1398.
- Olson, K.L., *et al.*, 2005. Real time spatial cluster detection using interpoint distances among precise patient locations. *Medical Informatics and Decision Making*, 5, 19.
- Openshaw, S., *et al.*, 1988. An investigation of leukaemia clusters by use of a geographical analysis machine. *The Lancet*, 1, 272–273.
- Ord, J.K. and Getis, A., 1995. Local spatial autocorrelation statistics: distribution issues and an application. *Geographical Analysis*, 27 (4), 286–306.
- Pavlin, J.A., 2003. Investigation of disease outbreaks detected by “syndromic” surveillance systems. *Journal of Urban Health*, 80 (2), i107–i114.
- Platt, R., *et al.*, 2003. Syndromic surveillance using minimum transfer of identifiable data: the example of the national bioterrorism syndromic surveillance demonstration program. *Journal of Urban Health*, 80 (2), i25–i31.
- Quataert, P.K.M., *et al.*, 1999. Methodological problems and the role of statistics in cluster response studies: a framework. *European Journal of Epidemiology*, 15, 821–831.
- Reis, B. and Mandl, K., 2003. Time series modeling for syndromic surveillance. *BMC Medical Informatics and Decision Making*, 3, 2.
- Rogerson, P.A., 1997. Surveillance systems for monitoring the development of spatial patterns. *Statistics in Medicine*, 16, 2081–2093.
- Rogerson, P.A., 2001. Monitoring point patterns for the development of space-time clusters. *Journal of the Royal Statistical Society*, 164, 87–96.
- Rogerson, P.A. and Yamada, I., 2004. Monitoring change in spatial patterns of disease: comparing univariate and multivariate cumulative sum approaches. *Statistics in Medicine*, 23, 2195–2214.
- Rogerson, P.A. and Yamada, I., 2009. *Statistical detection and surveillance of geographic clusters*. Boca Raton, FL: Chapman & Hall/CRC.
- Rolfhamre, P., *et al.*, 2004. Implementing a public web based GIS service for feedback of surveillance data on communicable diseases in Sweden. *Infectious Diseases*, 4, 1471–2334.
- Sonesson, C. and Brock, D., 2003. A review and discussion of prospective statistical surveillance in public health. *Journal of the Royal Statistical Society, Series A*, 166, 5–21.
- Tian, J., Chen, D., and McGuinness, D., 2006. *Spatial-temporal cluster detection approaches in the syndromic surveillance system* Reported to Queen’s University Emergency Syndromic Surveillance Team, KFL & A Public Health Unit, Ontario, Canada.
- Tsui, F., Espino, J.U., and Dato, V.M., 2003. Technical description of RODS: a real-time public health surveillance system. *Journal of the American Medical Informatics Association*, 10, 399–408.
- Wong, W., *et al.*, 2003. WSARE: what’s strange about recent events? *Journal of Urban Health*, 80, i66–i75.
- Wong, W., *et al.*, 2005. What’s strange about recent events (WSARE): an algorithm for early detection of disease outbreaks. *Journal of Machine Learning Research*, 6, 1961–1998.
- Yan, P.Y., Chen, H., and Zeng, D., 2008. Syndromic surveillance systems: public health and biodefense. *Annual Review of Information Science and Technology*, 42, 425–495.
- Yih, W.K., *et al.*, 2005. Ambulatory-care diagnoses as potential indicators of outbreaks of gastrointestinal illness: Minnesota. *Morbidity and Mortality Weekly Report*, 54, 157–162.

Zeng, D., Chang, W., and Chen, H., 2004. A comparative study of spatio-temporal hotspot analysis techniques in security informatics. In *Proceedings of the 7th IEEE transactions on intelligent transportation systems*, 3–8 October. Washington, DC: IEEE Computer Society, 106–111.

Zeng, D., *et al.*, 2005. Bio Portal: a case study in infectious disease informatics. In: *Proceedings of the 5th ACM/IEEE-CS Joint Conference on Digital Libraries*, 7–11 June, Denver, CO, 418.