**\*Title Page (including Author Details)**

# An examination of five spatial disease clustering methodologies for the identification of childhood cancer clusters in Alberta, Canada

M. Torabi[1], R.J. Rosychuk[2]

1. Department of Community Health Sciences, University of Manitoba, Winnipeg, Manitoba, Canada R3E 0W3

2. Department of Pediatrics, University of Alberta, Edmonton, Alberta, Canada T6G 2J3

Correspondence to:

Rhonda J. Rosychuk, PhD, PStat

Department of Pediatrics, University of Alberta

9423 Aberhart Centre, 11402 University Avenue NW

Edmonton, Alberta T6G 2J3 CANADA

e-mail: rhonda.rosychuk@ualberta.ca

phone: +1 780 492 0318 fax: +1 780 407 6435

# An examination of five spatial disease clustering methodologies for the identification of childhood cancer clusters in Alberta, Canada

## Abstract

Cluster detection is an important part of spatial epidemiology because it may help suggest potential factors associated with disease and thus, guide further investigation of the nature of diseases. Many different methods have been proposed to test for disease clusters. In this paper, we study five popular methods for detecting spatial clusters. These methods are Besag-Newell (BN), circular spatial scan statistic (CSS), flexible spatial scan statistic (FSS), Tango's maximized excess events test (MEET), and Bayesian disease mapping (BYM). We study these five different methods by analyzing a data set of malignant cancer diagnoses in children in the province of Alberta, Canada during 1983-2004. Our results show that the potential clusters are located in the south-central part of the province. Although, all methods performed very well to detect clusters, the BN and MEET methods identified local as well as general clusters.

*Keywords:* Bayesian statistic; Cancer cases; Geographic epidemiology; Spatial cluster detection

## 1. Introduction

Childhood cancers differ from adult cancers in terms of type and distribution. Leukemias, brain and other nervous system tumours, lymphomas (lymph node cancers), bone cancers, soft tissue sarcomas, kidney cancers, eye cancers, and adrenal gland cancers, are the most common types of cancer in children, while skin, prostate, breast, lung, and colorectal cancers are the most common cancers in adults [1]. Childhood cancers also differ from adult cancers based on biological, clinical and environmental features, growth rates, and treatment responses. While in most cases the causes of childhood

cancers are unknown, the causes of adult cancers are environmental, occupational and lifestyle factors such as diet, alcohol and smoking [2, 3].

In North America, childhood cancer is the most common cause of death from disease in the pediatric population (one year of age through adolescence) [4, 5]; more deaths than asthma, diabetes, cystic fibrosis and AIDS combined in Canada [4]. With such an impact, it is important to identify regions, in childhood malignant cancer diagnoses, with high ratio of cancer diagnoses. The focus of our paper is to examine geographical variations of the number of childhood cancer diagnoses during 1983 to 2004 in the western Canadian province of Alberta.

A limited region within the study regions with a high ratio of disease cases is defined as a spatial cluster [6]. The identification of a cluster of disease can help to find potential factors associated with disease and lead to improved understanding of etiology. Moreover, identification of clusters may lead to more detailed investigations to find out the association between exposures and disease interventions [7].

Statistical cluster detection methods are generally classified into two main categories, focused and general (also called as non-focused). Methods for focused cluster detection are designed to identify regions with excess number of cases in the vicinity of potential causes (e.g., toxic waste site) [8, 9]. On the other hand, methods for general clusters are designed to identify regions with excess number of cases. Typically, these models adopt extra-Poisson variability in different ways [10, 11, 12].

Methods for focused cluster detection include, but are not limited to, circular spatial scan statistic (CSS)[13], flexible spatial scan statistic (FSS)[14], and Bayesian disease mapping (BYM)[12]. The methods for general cluster detection include the Besag and Newell (BN)[15] test and the maximizing excess events test (MEET)[16]. The aim of focused tests is to test the null hypothesis of no local spatial cluster, while, the general tests are used to detect the potential clusters in the study region. In other words, for the focused tests (CSS, FSS, and BYM), the goal is to find a cluster for a specific region of interest, and consequently the test statistics are designed to capture the potential cluster. For the general tests (BN and MEET), the goal is to find any significant cluster in the study region without specifying any region of interest.

2

Since multiple tests with different assumptions have been proposed in the literature, it is important to compare and contrast methods to examine their performance on a particular data set. The similar and diverse results provide insights on the features that the different methods can detect. In this paper, we study five different and popular methods (BN, CSS, FSS, MEET, and BYM) to detect clusters with high ratio of childhood malignant cancer cases in the province of Alberta, Canada during 1983-2004. These methods were chosen based on their popularity in the literature [17, 18]. Moreover, the relationship between the shape of clusters and the methods used to detect spatial clusters are also investigated.

## 2. Materials and methods

### 2.1. Study subjects

The study was based on a yearly data set of malignant cancer diagnoses in children (age $\leq$ 19) in the western Canadian province of Alberta during the 1983-2004 fiscal years (see http://atlas.nrcan.gc.ca/site/english/maps /reference/national/can_political_e/map.pdf for a map of Canada). During the study period, the population of Alberta increased from 2.4 million in 1983 to 3.2 million in 2004 and the average population of children numbered around 800,000. During the last study year, the province consisted of nine Regional Health Authorities that were responsible for the delivery of health care services. These regions were further sub-divided into 70 areas (called sRHAs). These non-overlapping sRHAs are the geographic unit used in our analysis and all data were linked to these geographic boundaries. In addition, a population-based centroid was provided for each sRHA and these centroids were not necessarily geographic centres. For simplicity, we call these regions $1, 2, ..., 70$. The data was aggregated over the study period 1983-2004.

The number of malignant cancer cases totaled 2,728 over the study period. The median number of yearly cases per sRHA was 1 (range 0 to 12). The distribution of gender was 45% females and 55% males and most of malignant cancer cases were in the age groups 0-4 (32%) and 15-19 years (33%). The percentages of cancer cases for other two categories, 5-9 and 10-14 years, were 16% and 19%, respectively.

3

The key data requirements for the following methods are the number of cases and the number of expected cases or the population size for each region. When the expected number of disease cases varies by important strata, such as year, age group and gender, adjustments can be made. The expected number of disease cases is then adjusted by year (1-22), age group (0-4, 5-9, 10-14, 15-19 years) and gender (male, female). Some counts are suppressed to ensure confidentiality.

### 2.2. Besag-Newell's R statistic (BN)

Besag and Newell (BN)[15] proposed a test for each region based on the number of neighbours that must be combined to contain a minimum number of cases (cluster size). The spatial relationship among the regions is characterized by calculating pairwise distances between region centroids. Denote $i_j$ as the $j$-th closest region to region $i$, $(j = 1, ..., m-1)$, $i_0 = i$ and $m$ is the number of regions. Let $N_i$ be the population at region $i$ and $N_{i:k}(= \sum_{j=1}^{k} N_{i_j})$ be the total population of the $k-$nearest neighbours of region $i$. Denote $C_i$ as the number of cases in region $i$ and $C_{i:k}(= \sum_{j=1}^{k} C_{i_j})$ is the number of cases in the $k$-th nearest neighbours of region $i$. The total population in study region is $N(= \sum_{i=1}^{m} N_i)$ and the total number of cases is $C(= \sum_{i=1}^{m} C_i)$. In this method, the cluster size (called $l$) is pre-specified. The test statistic for region $i$ is the number of regions that must be combined with region $i$, to include the nearest $l$ cases. This test statistic for region $i$ is given by

$$T_i = \min \left\{ J : l \leq \sum_{j=1}^{J} C_{i_j} \right\}. \tag{1}$$

Under the null hypothesis, $C_{i:k}$ has Poisson distribution with mean $\lambda_{i:k} = N_{i:k}C/N$, where under null hypothesis, every individual is equally likely to be a case independent of other individuals and the location of residence. For region $i$, the significance level becomes

$$P(T_i \leq t) = 1 - \sum_{z=0}^{l-1} \lambda_{i:k}^z e^{-\lambda_{i:k}} /z!. \tag{2}$$

A region is identified as cluster when the significance level is equal or less than 0.05. For calculations, we used the corresponding observed quantities (i.e., replace $C$ with the total number of observed cases in the study region).

4

This method mainly relies on a pre-determined cluster size for each test and Le et al. [19] provided a testing algorithm for the automatic selection of cluster sizes. We implement this method using R [20] code.

### 2.3. Circular spatial scan statistic (CSS)

The spatial scan statistic has been used in a wide range of applications within the field of epidemiology [21]. The circular spatial scan statistic imposes a circular window $S$ on each region, and for any of those regions, the radius of the circle varies from zero to a pre-specified maximum distance $d$ or a pre-specified maximum number of regions $K$ to be included in the cluster. Let $S_{i:j}(j = 1, ..., J)$ denote the window composed by the *(j-1)-th* nearest neighbours to region $i$. The set of all windows to be scanned by the circular spatial scan statistic is $S_1 = \left\{ S_{i:j}; i = 1, ..., m; j = 1, ..., J \right\}$. For each circle, a likelihood ratio statistic is computed based on the number of observed and expected cases within and outside the circle. Let $L_0$ and $L_i(i = 1, ..., m)$ be likelihood under the null and alternative hypothesis, where the null hypothesis is no cluster in region $i$ and the alternative hypothesis is a cluster in region $i$ based on its *j-th* nearest neighbours. Then the likelihood ratio statistic is given by

$$\max_i \frac{L_i}{L_0} = \left(\frac{C_i}{E_i}\right)^{C_i} \left(\frac{N - C_i}{N - E_i}\right)^{N-C_i} I(C_i > E_i), \tag{3}$$

where $C_i$ and $E_i$ denote the observed and expected number of cases in a circle, respectively, and $(N - C_i)$ and $(N - E_i)$ denote the observed and expected number of cases outside the circle, respectively. Note that the indicator function $I(\cdot)$ is equal to one when $C_i > E_i$ and 0 elsewhere.

The circles with the highest likelihood ratio values are identified as potential clusters. We can implement this method using SaTScan [22] or FleXScan [23] software. In general, the $K$ is chosen to include at most 50% of population at risk. We used $K = 15$, the FleXScan default, and since our example uses aggregate data, the region centroid had to be included in the radius of the circle for the region to be part of the circle.

### 2.4. Flexible spatial scan statistic (FSS)

This method is similar to the method of CSS; however, the detected cluster is allowed to be *flexible* in shape while at the same time the cluster is confined to a relatively small neighbourhood of each region. The flexible scan statistic imposes an irregularly shaped window $S$ on each region by connecting its adjacent regions. For each region $i$, the set of irregularly shaped windows with length $j$, the $j$ connected regions including $i$, can move from 1 to the pre-specified maximum $J$. The connected regions are restricted to the subsets of the set of regions $i$ and *(J-1)-th* nearest neighbours to the region $i$, where $J$ is a pre-specified maximum length of cluster. The set of all windows to be scanned by the flexible spatial scan statistic is then $S_2 = \left\{ S_{i:j(k)}; i = 1, ..., m; j = 1, ..., J; k = 1, ..., k_{ij} \right\}$. Note that the circular spatial scan statistic considers $J$ circles for a given region $i$; however, the flexible spatial scan statistics considers $J$ circles in addition to the all sets of connected regions whose centroids are located within the *J-th* largest concentric circle. As a consequence, the size of $S_2$ is much larger than $S_1$ which is at most $mJ$. Under the Poisson assumption, the test statistic for the flexible spatial scan statistic based on the likelihood ratio test is obtained by (3), where the circle defined in (3) now refers to the $S_2$ rather than $S_1$. We implement this method with the FleXScan software, using the default setting $J = 15$. Similar to the circular spatial scan statistic, the circles with the highest likelihood ratio values are identified as potential clusters.

### 2.5. Tango's maximized excess events test (MEET)

For a given parameter $k$, the Excess Events Test statistic [24] is defined as

$$T_0(k) = \sum_i U_i(k) = \sum_i \sum_j e^{-4d_{ij}^2/k^2}(C_i - N_iC/N)(C_j - N_jC/N),$$

where $d_{ij}$ is the distance between region $i$ and $j$. However, the choice of $k$ refers to the geographical scale of clustering. A large $k$ makes the test sensitive to geographically large clusters, while a small $k$ makes the test more sensitive to small clusters. To detect clustering irrespectively of its geographical scale, Tango [16] proposed the Maximized Excess Events Test (MEET) as

$$T = \min_{0 \leq k \leq K} P\{T_0(k) > t_0(k)|H_0, k\}, \tag{4}$$

where $t_0(k)$ is the observed value of the Excess Events Test statistic for a given $k$, and $K$ is upper bound on $k$. In practice, the test uses a line search by discretization of $k$, and the MEET statistic is evaluated by Monte Carlo hypothesis testing. The null hypothesis of no clustering, $H_0$, is rejected when the test statistic is small. Given $k$ and under the null hypothesis, the test statistic $T_0(k)$ has an asymptotically chi-square distribution. If the null hypothesis of no clustering is rejected, the most likely centres of clusters may be identified by the region with maximum of $U_i(k)$. In practice, the regions with high outlying percentages will likely be the locations of clusters.

### 2.6. Bayesian disease mapping (BYM)

A Bayesian approach using Markov chain Monte Carlo (MCMC) can also be used for cluster detection [10, 11, 25, 26]. This approach was first used by Besag et al. (BYM) [10] and the model consists of two parts. In the first part, the cases are assumed to follow a Poisson distribution with an area specific parameter $\theta_i E_i$ :

$$C_i \sim Poisson(\theta_i E_i),$$

where $C_i$ and $E_i$ are the observed and expected number of cases in region $i$, respectively. The second part of the model is obtained by

$$\log(\theta_i) = \mu + \eta_i + \phi_i,$$

where $\theta_i$ is the relative risk $(RR_i)$ in region $i$, $\mu$ is an overall mean, $\eta_i$ represents specified features of region $i$ which accommodates spatial structure, and $\phi_i$ denotes unspecified features of region $i$ which does not incorporate spatial structure. The uncorrelated component $\phi_i$ is assumed to follow a Gaussian distribution with zero mean and a common variance $\sigma_\phi^2$. The correlated component $\eta_i$ is assumed to follow an intrinsic conditionally autoregressive (ICAR) distribution depending on their neighbouring values. In particular,

$$\eta = (\eta_1, ..., \eta_m)' \sim N(0, \Sigma_\eta),$$

where $\Sigma_\eta = \sigma_\eta^2 D^{-1}$, and $\sigma_\eta^2$ is the spatial dispersion parameter. The neighbourhood matrix $D$ has its $i$-th diagonal element equal to the number of neighbours of the corresponding region, and the off-diagonal elements in each row equal -1 if the corresponding regions are neighbours and zero otherwise

[27, 28]. The parameters can be then estimated within the Bayesian framework (MCMC) using vague priors for the parameters. This produces the posterior distributions for the parameters in the model.

A cluster is defined as a region where the estimated relative risk is significantly larger than 1 (in terms of their credibility sets) [29]. To implement this method, we used WinBUGS software [30] to compute the relative risk values.

These methods have different limitations and strengths. As a limitation, we assume that the number of cases follows a Poisson distribution under the null hypothesis of no cluster for the BN and BYM methods. We also need to specify the number of regions to be included in the cluster for the CSS and FSS methods. As a strength, the CSS, FSS, and MEET methods are distribution free. Also, we do not need to specify the cluster size for the BYM and BN methods (if using a testing algorithm [19] with the BN method).

### 2.7. Specific hypotheses

Although, the specific alternative hypotheses need to be specified only for the methods CSS, FSS, and BYM, we may also want to specify the alternative hypotheses for the methods BN and MEET as well. We consider multiple alternatives that are tested separately. Further, let $RR_i$ indicate the relative risk for the $i$-th region within clusters when compared with the region outside clusters; the latter has $RR_i = 1$. For example for cluster $X$, the $RR_i$ is given by

$$RR_i = \begin{cases} 3 & i \in X \\ 1 & \text{otherwise.} \end{cases}$$

### 3. Results

We have provided the comparison of the five methods (BN, CSS, FSS, MEET, and BYM) to detect the potential clusters in our childhood malignant cancer cases for the period of 22 years (1983-2004) in the province of Alberta, Canada.

Based on the 70 regions, six different clusters were tested: (1) 19 regions from the urban south-central part of the province (called A), (2) 18 regions

from the urban central part of the province (called B), (3) two clusters of regions A and B (called AB), (4) four regions from the rural north part of the province (called D), (5) nine regions from the mixed (urban and rural) south-central part of the province (called F), and (6) a case of no clusters (called G). For G, no region was specified as a potential cluster. More precisely, the clusters are $A = \{8, 9, ..., 26\}, B = \{41, 42, ..., 58\}, D = \{59, 60, 61, 62\}$, and $F = \{27, 28, ..., 35\}$.

In Figures 1 to 6, the areas that are statistically significant (potential clusters) are shown for each cluster and each method separately. The summary of cluster F, an area of mixed urban and rural parts of the province, is presented in Table 1. As shown, the order of significant regions of five different methods is reported for cluster F. More precisely, the regions are ordered based on which one is more significant to be as a cluster. For the BN method for example, 1 means that the regions 27, 28, 30, and 35 are most likely to constitute a significant cluster, while 8 means that the region 53 is least likely to be a significant cluster. Hence, it is easy to see which region has more contribution to constitute a cluster. For local (hot-spot) clusters (A, B, AB, D, F), almost all methods detected these areas as potential clusters.

Figure 1: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster A.

Figure 2: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster B.

Figure 3: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster AB.

Figure 4: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster D.

Figure 5: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster F.

Figure 6: Subregional health authorities (HAs) identified as potential clusters (shaded regions) for five methods (BN, CSS, FSS, MEET, and BYM) for cluster G.

For example, for cluster AB, the CSS method detected the AB cluster as a potential cluster except for regions $\{23, 25, 53, 55\}$. However, the FSS method identified AB as a potential cluster except for regions $\{23, 55\}$. The main reason for different results between CSS and FSS is due to non-circular shape of regions $\{25, 53\}$ in relation to cluster AB. Region 28 and cluster AB, without region 53, constitute a potential cluster for the BN method, while region 31 and cluster AB, without region 53, contain a potential cluster for the BYM method. The cluster AB is a potential cluster for the MEET method.

For the case of no cluster G (general test), CSS, FSS, and BYM did not identify any potential clusters. However, BN and MEET methods identified some regions as potential clusters. In BN method, the regions $\{12, 14, 15, 25, 31, 32, 51\}$ were identified as potential clusters, while regions $\{9, 15, 21, 25, 31, 34, 47, 51\}$ were potential clusters for MEET method. Note that for cluster G, we have $RR_i = 1(i = 1, ..., 70)$.

The methods CSS, FSS, and BYM detected local clusters as potential clusters. The FSS method also identified regions with a non-circular shape as a potential cluster unlike CSS method. The BN and MEET methods detected potential clusters in both scenarios (local and global clusters).

10

Table 1. The order of significant regions of five different methods for cluster $F$.

| | | | Methods | | | | |
|---|---|---|---|---|---|---|---|
| Region | $C_i$ | $E_i$ | BN | MEET | CSS | FSS | BYM |
| 12 | 67 | 64 | 5 | - | - | - | - |
| 14 | 34 | 32 | 5 | - | - | - | - |
| 15 | 64 | 43 | 6 | 9 | - | - | 10 |
| 25 | 28 | 15 | 2 | 8 | - | - | 11 |
| 27 | 19 | 6 | 1 | 4 | 3 | - | 5 |
| 28 | 17 | 6 | 1 | 6 | 3 | 1 | 8 |
| 29 | 38 | 13 | 4 | 2 | 1 | 1 | 3 |
| 30 | 17 | 6 | 1 | 6 | 1 | 1 | 6 |
| 31 | 38 | 13 | 4 | 2 | 1 | 1 | 2 |
| 32 | 104 | 35 | 3 | 1 | 1 | 1 | 1 |
| 33 | * | 3 | 5 | 7 | 1 | 1 | 9 |
| 34 | 35 | 12 | 3 | 3 | 2 | 1 | 4 |
| 35 | 18 | 6 | 1 | 5 | 2 | - | 7 |
| 51 | 58 | 44 | - | 10 | - | - | - |
| 53 | * | 8 | 8 | - | - | - | - |
| 59 | 44 | 42 | 7 | - | - | - | - |

"*" represents small count; "-" non-significant regions; $C_i$ and $E_i$ are observed and expected number of cases in region $i$; BN, MEET, CSS, FSS, and BYM are Besag-Newell's R statistic, Tango's maximized excess events test, circular spatial scan statistic, flexible spatial scan statistic, and Bayesian disease mapping methods, respectively.

## 4. Discussion

We have provided the comparison of five different methods (BN, CSS, FSS, MEET, and BYM) with potential for detecting clusters with high ratio of childhood malignant cancer cases in the province of Alberta, Canada. These five methods have been extensively used in the literature and are relatively comprehensive. These methods use different approaches (semi-parametric to parametric as well as Bayesian) to test for significant clusters.

We considered six different alternative hypotheses, including local and global clustering, to compare the results of different methods. The CSS,

FSS, and BYM methods detected potential clusters in the scenarios of local clusters as expected, but they did not detect the global cluster (cluster G). In addition, the CSS method identified a lower number of regions combined as a potential cluster compared to FSS method, due to non-circular shape of some regions in the province of Alberta. The BN and MEET methods identified local clusters in addition to global clusters. It seems that the malignant cancer diagnoses cases tend to constitute the local clusters, particularly in south-central part of the province. Hence, we recommend using the methods of BN and MEET and particulary MEET, due to fewer assumptions compared with the BN method.

In the BYM approach, we conservatively defined a region as a cluster if the credibility set of the estimated relative risk was larger than one. One may define different decision rule where the estimated relative risk would be larger than one [31].

We adjusted our expected number of malignant cancer by three important factors age, gender, and year. We did not include any covariates in the model, since only the BN and BYM methods allow for the direct inclusion of covariates. We also note that the methods have different settings and assumptions, which motivate our comparisons. User-chosen settings are part of all cluster tests and different choices could lead to different results. The CSS, FSS, and BYM methods have been proposed for local clusters, while the BN and MEET methods have been advocated for global clusters. Under the null hypothesis, the number of cancer cases follows a Poisson distribution for the BN and BYM methods, while the test statistic for the CSS, FSS, and MEET methods has an asymptotically chi-square distribution. These features motivated us to consider these important methods and apply them to our malignant cancer cases.

In general, the potential clusters are located in the south-central part of the province (cluster G). These findings may represent real clusters or may represent different distributions of important factors that are unmeasured and unadjusted for in our modeling. Our results highlight how different methods can produce different results and further investigation may be warranted to explore these findings.

12

## References

[1] Buka I, Korantenge S, Osornio Vergas AR. Trends in childhood cancer incidence: review of environmental linkages. Pediatr Clin N Am 2007; 54: 177-203.

[2] Gochfeld M. Chronologic history of occupational medicine. J Occup Environ Med 2005; 47(2): 96-114.

[3] Kote-Jaria Z, Salmon A, Mengistu T, Copeland M, Ardern-Jones A, Locke I, Shanley S, Summersgill B, Lu YJ, Shipley J, Eeles R. Increased level of chromosomal damage after irradiation of lymphocytes from BRCA1 mutation carriers. Br J Cancer 2006; 94(2): 308-10.

[4] Borugian MJ, Spinelli JJ, Mezei G, Wilkias R, Abanto Z, McBride ML. Childhood leukemia and socioeconomic status in Canada. J Epid 2005; 16: 526-31.

[5] Robinson L. General principles of the epidemiology of childhood cancer. In: Pizzo P, Poplack D, editors. Principles and Practice of Pediatric Oncology. Philadelphia, Lippincott - Raven; 1997. p.1-10.

[6] Lawson AB. Statistical methods in spatial epidemiology. 2nd ed. London: John Wiley & Sons, Ltd; 2006.

[7] Jennings JM, Curriero FC, Celentano D, Ellen JM. Geographic identification of high gonorrhea transmission areas in Baltimore, Maryland. Am J Epid 2005; 161: 73-80.

[8] Elliott P, Briggs D, Morris S, de Hoogh C, Hurt C, Jensen TK, Maitland I, Richardson S, Wakefield J, Jarup L. Risk of adverse birth outcomes in populations living near landfill sites. Br Med J 2001; 323: 363-68.

[9] Lawson AB, Biggeri A, Williams FLR. A review of modeling approaches in health risk assessment around putative sources. In: Lawson AB, Biggeri A, Böhning D, Lesaffre E, Viel J, Bertollini R, editors. Disease Mapping and Risk Assessment for Public Health. New York, Wiley; 1999. p.231-45.

[10] Besag JE, York JC, Mollìe A. Bayesian image restoration with two applications in spatial statistics (with discussion). Ann Inst Statist Math 1991; 43: 1-59.

13

[11] Clayton D, Bernardinelli L. Bayesian methods for mapping disease risk. In: Elliott P, Cuzick J, English D, Stern R, editors. Geographical and Environmental Epidemiology: Methods for Small-Area Studies. Oxford, Oxford University Press; 1996. p.205-20.

[12] Clayton D, Kaldor J. Empirical Bayes estimates of age-standardized relative risks for use in disease mapping. Biometrics 1987; 43: 671-81.

[13] Kulldorff M. A spatial scan statistics. Comm Statist: Theor Meth 1997; 26: 1481-96.

[14] Tango T, Takahashi K. A flexibly shaped spatial scan statistic for detecting clusters. Int J Health Geogr 2005; 4:11.

[15] Besag JE, Newell J. The detection of clusters in rare diseases. J Roy Statist Soc Ser A 1991; 154: 143-55.

[16] Tango T. A test for spatial disease clustering adjusted for multiple testing. Statist Med 2000; 19: 191-204.

[17] Kulldorff M, Tango T, Park PJ. Power comparisons for disease clustering tests. Comput Statist Data An 2003; 42: 665-84.

[18] Song C, Kulldorff M. Power evaluation of disease clustering tests. Int J Health Geogr 2003; 2:9.

[19] Le ND, Petkau AJ, Rosychuk RJ. Surveillance of clustering near point sources. Statist Med 1996; 15: 727-40.

[20] R: A language and Environment for Statistical Computing [http://www.R-project.org].

[21] Fukuda Y, Umezaki M, Nakamura K, Takano T. Variations in social characteristics of spatial disease clusters: examples of colon, lung and breast cancer in Japan. Int J Health Geogr 2005; 4:16.

[22] Kulldorff M, Rand K, Gherman G, Williams G, DeFrancesco D. SaTScan V2.1: Software for the spatial and space-time scan statistics. National Centre Institute, Bethesda; 1998.

[23] Takahashi K, Yokoyama T, Tango T. FleXScan: Software for the Flexible Scan Statistic. National Institute of Public Health, Japan; 2006.

14

[24] Tango T. A class of tests for detecting "general" and "focused" clustering of rare diseases. Statist Med 1995; 14: 2323-34.

[25] Bernardinelli L, Montomoli C. Empirical Bayes versus fully Bayesian analysis of geographical variation in disease risk. Statist Med 1992; 11: 983-1007.

[26] Gilks WR, Richardson S, Spielhalter DJ. (eds). Markov chain Monte Carlo in practice. Chapman and Hall/CRC; 1995.

[27] MacNab YC, Dean CB. Parametric bootstrap and penalized quasi-likelihood inference in conditional autoregressive models. Statist Med 2000; 19: (17/18): 2421-35.

[28] MacNab YC, Dean CB. Autoregressive spatial smoothing and temporal spline smoothing for mapping rates. Biometrics 2001; 57: 949-56.

[29] Aamodt G, Samuelsen SO, Skrondal A. A simulated study of three methods for detecting disease clusters. Int J Health Geogr 2006; 5:15.

[30] Lunn DJ, Thomas A, Best N, Spiegelhalter DJ. WinBUGS- a Bayesian modelling framework: concepts, structure, and extensibility. Statist Comput 2000; 10: 325-37.

[31] Richrdson S, Thomson A, Best N, Elliott P. Interpreting posterior risk estimates in disease-mapping studies. Environ Health Persp 2004; 112(9): 1016-25.

Figure 1 BN



BN

Figure 1 BYM



BYM

Figure 1 CSS

CSS

**Figure 1 FSS**



FSS

**Figure 1 MEET**

MEET

Figure 2 BN



BN

**Figure 2 BYM**



BYM

Figure 2 CSS



CSS

Figure 2 FSS



FSS

Figure 2 MEET



MEET

Figure 3 BN



BN

Figure 3 BYM



BYM

Figure 3 CSS



CSS

Figure 3 FSS



FSS

Figure 3 MEET

MEET

Figure 4 BN



BN

Figure 4 BYM



BYM

Figure 4 CSS



CSS

**Figure 4 FSS**

FSS

Figure 4 MEET

MEET

Figure 5 BN



BN

Figure 5 BYM



BYM

**Figure 5 CSS**



CSS

Figure 5 FSS



FSS

Figure 5 MEET

MEET

Figure 6 BN



BN

**Figure 6 BYM**



BYM

Figure 6 CSS



CSS

Figure 6 FSS

FSS

Figure 6 MEET



MEET