

---

# Using administrative data to understand the geography of case ascertainment

---

N. Yiannakoulias (1) PhD; D. P. Schopflocher (2) PhD; L. W. Svenson (3) BSc

---

## Abstract

We examined the geographic variability of information generated from different case definitions of childhood asthma derived from administrative health data used in Alberta, Canada. Our objective was to determine if analyses based on different case ascertainment algorithms identify geographic clusters in the same region of the study area. Our study group was based on a closed cohort of asthmatic children born in 1988. We used a spatial scan statistic to identify variations in the approximate location of geographic clusters of asthma based on different case definitions. Our results indicate that the geographic patterns are not greatly affected by the case ascertainment algorithm or the source of data. For example, asthmatics identified from medical claims data showed similar clustering to asthmatics defined through hospitalization and emergency department data. However, estimates of prevalence and incidence require careful consideration and validation against other data sources.

---

**Keywords:** public health surveillance, spatial analysis, administrative health data, asthma

---

## Introduction

The growing availability of electronic health data and capacities for computer hardware to warehouse and analyse such data presents new opportunities for research on a variety of health outcomes. There is a growing body of evidence suggesting that administrative health data can be particularly important resources for research and surveillance of chronic diseases.<sup>1-4</sup> One of the most important contributions of administrative health data has been to provide information that facilitates the analysis of entire populations covering large geographic regions. In Canada, this often involves provincial-scale analysis (that is, the comparisons of regions within a province) that can be particularly useful for linking health information systems to decision making in public health.

Geographic variations between regions within a province may indicate differences in epidemiology, population attributes, availability of services, exposure to environmental hazards, diagnostic practice and a variety of other factors.

One of the challenges to using administrative health data in research and surveillance is that different methods of case ascertainment may confound group differences especially if how these data were collected and/or generated varies. For example, there may be rural/urban differences in the effectiveness of certain data to identify incident stroke events.<sup>5-7</sup> These differences could be related to delivery of care (particularly in rural areas in which acute care centres often function in primary care roles), availability of diagnostic resources, geographic variations in

physician specialty and many other factors. The most commonly discussed solution to these data problems is the development of case ascertainment algorithms that combine multiple administrative health data sources and/or multiple records within a single source.<sup>3,8</sup> With this paradigm, rather than using a single record from a single data source to identify a case, multiple records are combined with data from multiple sources. Integrating multiple sources of data may improve case ascertainment by traditional measures (such as sensitivity and specificity) but also improve the geographic uniformity in the case selection process by ensuring that geographic comparisons are not overly influenced by local or regional anomalies associated with a single data source.

Most case definition algorithms using administrative data have been validated against clinical chart reviews<sup>9,10</sup> or survey responses.<sup>4,11</sup> Although the former is an important benchmark for evaluating a case definition algorithm, it is less able to characterize the ways in which an algorithm's properties might vary between regions. This is particularly true if the validation work is based on a specific site of study, rather than a comprehensive sampling of sites over a large geographic area. While comparing administrative data to survey data can be informative, this process deals with two fallible data sources and no "gold standard."

In this study, we examined the geographic variability of information generated from different case definitions of childhood

---

## Author References

1 School of Geography and Earth Sciences, McMaster University, Hamilton, ON

2 Faculty of Nursing, University of Alberta, Edmonton, AB

3 Surveillance and Environmental Health, Alberta Health & Wellness, Edmonton, AB

**Correspondence:** Nikolaos Yiannakoulias, School of Geography and Earth Sciences, McMaster University, 1280 Main Street West, Hamilton, ON L8S 4K1, Tel.: 905-525-9140, ext. 20117, Fax: 905-546-0463, Email: yiannan@mcmaster.ca

asthma derived from an administrative health data system. Our objective was to observe how the use of different case ascertainment algorithms can affect the appearance of patterns and variability on a map. Rather than comparing case definitions to a gold standard or a survey of respondents, we compared how analytical information varies across a variety of case definition algorithms. Our specific approach was to search for geographic clusters of children identified as asthma cases. If analyses based on a variety of different case definition algorithms identify clusters in the same region of the study area, the analytical information about geographic patterns is reliable even if the case ascertainment criteria differ. On the other hand, large variations in the resulting analyses could suggest that case definition algorithms should be more geographically specific, and should be designed to take into account local or regional differences in the use of services, conformity to data standards, or other factors.

Information from this study will inform the development of case definitions for childhood asthma based on administrative data by identifying the degree to which algorithms should be geographically specific. More generally, however, we present a simple framework for evaluating the geographical robustness of administrative health data for a variety of chronic conditions.

## Methods

### Data

Our study area, the province of Alberta, Canada, is well suited to this analysis because of both the availability of multiple sources of administrative health data and the existence of a population registry that can be used to identify the location of residence over time. Like other Canadian provinces, Alberta maintains a publicly funded single payer health care insurance system that covers most health services. Residents of the province who do not opt out of the provincially insured health care system are required to register with the Alberta Health Care Insurance Plan and are then provided

a unique lifetime identifier that enables the linkage of health data sources. Over 99% of the province's population is registered under this system.

Alberta Health and Wellness has designed a system for the creation of longitudinal health data profiles in support of a variety of public health surveillance activities. These profiles are based on a linkage of four databases: the fee-for-service medical claims (claims), the Ambulatory Care Classification System containing emergency department admissions (emergency), an in-patient hospital services system (in-patient) and the Alberta Health Care Insurance registry system. Data from these systems are linked based on the unique numeric identifier and then tabulated into annual counts of services associated with a particular condition (defined by ICD-9\* or ICD-10† codes) over time. When linked to the population registry system (which contains information about sex, age and place of residence) this system provides a simple method for observing changes in population estimates of incidence and/or prevalence based on different case definition algorithms.

We used a data profile for asthma-related services (identified using ICD-9 code 493 or ICD-10 J45 in the primary diagnostic field) as the primary source of data in this study. Our study period was between the 1998 and 2004 calendar years. In order to control for the variations in asthma management practice that have occurred in recent years, our study is restricted to persons born in the same year. Although the claims system has records back to the early eighties, the in-patient system has electronic records dating back to 1993 and the emergency data system has electronic records that only go back to 1998. In 1998, there were 44 651 children born in 1988 residing in the province and registered with the health care insurance plan. We restricted our study population to children born in 1988 who resided in the province continuously between 1998 and 2004; thus the subjects were old enough for reasonable asthma diagnoses to be made, but

still less than 18 years old for the study period. This gave a closed study cohort of 38 905 children. Of these, 8965 had at least one medical service recorded within one of the three health databases in which asthma was the primary diagnosis associated with the service. We assumed that the minimum threshold to identify a child as asthmatic is two or more asthma-related services (not on the same day) over the study period. We referred to this as the "baseline" asthmatic group. Subject to this definition, 5110 children in the cohort were considered baseline asthmatics as of 2004, giving an asthma prevalence of roughly 13% within this cohort.

We used the data profile system to define cases based on several different case definition algorithms (Table 1), referring to these as "test case-definitions." Definitions vary based on the number and type of services within each of the data systems. These definitions were chosen to exaggerate differences in the case identification algorithms; most algorithms would require fewer contacts to qualify persons as a case than the definitions presented here. Therefore, the interpretation of our results should be viewed as an extreme example of how different case ascertainment algorithms might present different information on the geographic distribution of disease. We avoided over-counting contacts for particular episodes by counting only one of any of these contacts in a single day. When there were multiple contacts between the data systems on a particular day, we preferentially retained in-patient records over emergency records and preferentially retained emergency records over claims records.

Using 2004 residential postal codes linked to the data profile, all data were aggregated to the level of municipality. Municipalities, consisting of cities, towns and villages, were restricted to those with at least 10 members of the cohort, making a total of 294. Smaller municipalities were joined to larger municipalities to ensure no members of the cohort were dropped from the analysis. Children living outside of municipal areas (e.g. on farms) were assigned

\* International Statistical Classification of Diseases and Related Health Problems, 9th Revision

† International Statistical Classification of Diseases and Related Health Problems, 10th Revision

to the municipality within which their residential postal code was included (typically, where they pick up their mail). As a result, some rural-dwelling children may have been assigned to a municipality that is neither the closest to their residence nor where they receive the majority of their medical services. However, any errors in geo-referencing are common to both the numerator and denominator and should not bias our results.

### Analysis

For visualization purposes, we mapped the prevalence of baseline asthmatics in Alberta. Rather than mapping crude prevalence rates, which would be greatly affected by small numbers, we used a modelled approach to estimate relative risk in a way that manages stochastic variation in the data. We used a Poisson model to predict a function of the mean number of baseline asthmatics and included population of children as an offset to control for variations in the geographic distribution of the cohort. This model also includes a random intercept effect for each municipality and estimates spatial parameters to smooth out local variations in prevalence. In simple terms, this process averages geographically neighbouring observations with each other in a way conceptually similar to a one-dimensional moving window average. This modelling process is referred to as a generalized linear mixed model (or GLMM) approach.<sup>12</sup> We used the SAS procedure PROC GLIMMIX to solve this model.<sup>13</sup> We mapped the predicted baseline asthma morbidity ratio in the cohort at the municipality level. The map produced is of polygons shaded according to relative risk, with each polygon representing the area surrounding a municipality (Figure 1).

Our primary analysis involved explicit hypothesis tests for geographic clusters based on the different test case-definition algorithms. We used the spatial scan statistic<sup>14</sup> to identify these geographic clusters. The spatial scan approach uses a moving and variably sized window (a circular one in this application) to search a large number of potential clusters. The method then identifies the cluster in the set that is most likely to cause the rejection of a null hypothesis of constant risk. This cluster is

referred to as the “most-likely” cluster of disease. The statistical significance of a most-likely cluster was evaluated through Monte Carlo simulation. By testing the significance of only the disease cluster most likely to cause the rejection of a null hypothesis of constant risk, the method avoids problems of multiple testing common to some other methods of local cluster detection.<sup>14</sup>

In our analysis, we investigated two general hypotheses for each of the six test case-definitions of asthma described in Table 1. The first hypothesis was that the spatial distribution of asthmatics of all definitions (including the baseline asthmatics) differed from the spatial distribution of the study cohort population. This corresponds to a null hypothesis of constant risk; that is, that there is no geographically clustered subset of municipalities that have an excess risk of asthma. We refer to this as the *constant risk null hypothesis*. Our test statistic was the Poisson model likelihood ratio,<sup>14</sup>

$$\left( \frac{c_i}{e(c_i)} \right)^{c_i} \left( \frac{C - c_i}{C - e(c_i)} \right)^{C - c_i}$$

where  $C$  is the total number of cases who are defined as asthmatic according to a particular definition,  $c_i$  is the number of these cases in municipality  $i$  and  $e(c_i)$  is the expected number of these cases in municipality  $i$ . Here we calculated  $e(c_i)$  as

$$e(c_i) = m_i g.$$

For all case definitions,  $m_i$  is the number of children in the study cohort residing in municipality  $i$  and  $g$  is the overall rate of asthmatics for a particular definition in the study cohort.

Results from the analysis above indicated whether or not to reject a null hypothesis of constant risk and approximately where there are clusters of asthmatics in the study cohort for the different test case-definitions. To determine the geographic variation in test case-definitions more explicitly, we also determined whether the geographic distribution of asthmatics according to each of the test case-definitions differed

from the distribution of baseline asthmatics. As before, we used the Poisson model likelihood-ratio test, but in this case,

$$e(c_i) = a_i h,$$

where  $a_i$  is the number of children who are asthmatic by a specified case definition residing in municipality  $i$  and  $h$  is the proportion of asthmatics according to this definition among the baseline asthmatic population. Here, a rejection of the null hypothesis for a particular test case-definition indicates that the geographic distribution of asthmatics identified by this case definition algorithm is no different from the geographic distribution of the baseline asthmatics. This corresponds to a test of a null hypothesis of constant case definition, and we refer to this as the *constant case definition null hypothesis*.

We used SaTScan v. 6.1 to search for clusters.<sup>15</sup> All clusters were bound to a size no larger than 50% of the population of Alberta, and all clusters searched were constrained to a circular shape. In all cases, a significance level of 0.05 was used to assess whether there is a most-likely cluster against the null hypotheses of constant risk and constant case definition.

### Results

Test case-definition “A” provided a prevalence estimate of 4.4%, which was less than half of the baseline group (Table 2). At the other extreme, test case-definition “F” provided prevalence estimates of less than 0.2%. For each of the test case-definitions, we also tabulated the average number of services for each child by the type of services within the medical system. For all test case-definitions, all children appeared to have frequent asthma-related contacts with the emergency system when compared to the cohort as a whole. Children in all the test case-definitions appeared to experience more contacts with medical system (for any reason) than baseline asthmatics and non-asthmatics in the cohort.

Figure 1 illustrates the model-predicted geographic distribution of baseline asthmatics in the study cohort. There is variation in the rate of asthma among the

**TABLE 1**  
**Test case-definitions**

Definition label	Test case-definition
A	6 or more services of any type, 1998-2004
B	6 or more services including a minimum of 2 or more emergency or in-patient admissions, 1998-2004
C	12 or more services of any type, 1998-2004
D	12 or more services including a minimum of 2 or more emergency or in-patient admissions, 1998-2004
E	6 or more emergency or in-patient admissions, 1998-2004
F	12 or more emergency or in-patient admissions, 1998-2004

baseline asthmatics, with the relative risk highest in the city of Calgary where it was 22% higher than the provincial average. The Edmonton area had relative risk very close to the provincial average (0.3% higher than the provincial average). Rural areas of central and northern Alberta had the lowest prevalence of asthma in the cohort.

Statistically significant clusters under the two null hypotheses are mapped on Figures 2 and 3. Based on our null hypothesis of constant-risk, all test case-definitions of asthma with the exception of “B” were associated with a statistically significant most-likely cluster. Cluster “B”, though not statistically significant, was located in a similar region to cluster “D”. The relative risk associated with each cluster is relatively small; in most cases, the study population located inside the cluster had a 25% higher risk of asthma than the study population located outside the cluster. The one exception corresponded to case definition “F”, for which the relative risk of asthma inside the cluster is more than double the risk outside the cluster. All mapped clusters represent regions where the likelihood of rejecting the null hypothesis of constant risk was highest for each of the case definitions.

Based on our null hypothesis of constant case definition, only definition “F” reached a level of statistical significance. Definition “F” is the strictest of all test case-definitions, and the total number of cases in this cluster was very small; for this definition, only 52 of the total number of cases were found inside this cluster. Children within this cluster had an 87% higher chance of being

cases (according to definition “F”) than children outside the cluster.

### Discussion

Under the null hypothesis of constant risk, there were few apparent differences in the location of clusters across the case definitions. All but one of the clusters occurred in the southwest area of the province, though the clusters do vary considerably in geographic size; for example, clusters “A” and “C” were smaller than the other clusters. This apparent similarity was based on a qualitative assessment of a map covering a large area, and it is unclear from this map alone if these observations represent a systematically different geography in the test case-definitions when compared to each other or to the baseline asthma group. The search for clusters under the null hypothesis of constant case definition provided a more explicit test of whether or not the location of asthma clusters varied by particular case ascertainment algorithm. All but one of these searches failed to reject the null hypothesis of constant case definition. This suggests that there were relatively small differences in the geographic pattern of asthma across the different test case-definitions and that the detection of clusters was fairly robust to the precise definition selected.

For the null hypothesis of constant case definition, the only statistically significant cluster was based on definition “F”. This cluster was located in southwest Alberta, in the same region as the clusters found under the null hypothesis of constant risk. It is possible that the cluster represented a geographically concentrated region of

high-service use or serious asthmatics within the asthmatic population of the cohort. The cluster could also be related to changes in the population or health service utilization practices in the region, where there has been noteworthy population growth in recent years. Finally, it is important to be mindful of the fact that this definition corresponds to a prevalence estimate of less than 1%, and is considerably stricter than any case definition likely to be used in epidemiological or surveillance applications. Therefore, though there was a statistically significant difference in relative risk between the areas inside and outside the cluster, this amounts to a very small difference in absolute risk.

Based on these observations, it appears that comparative geographical analysis of asthma risk is not greatly affected by the case ascertainment algorithm used. More generally, our findings indicate that geographic information about relative differences in prevalence (or relative risk) may be invariant to the specific choice of case definition even if prevalence (absolute risk) varies across these case definitions. The implication of our findings, if they can be generalized to other settings or other chronic conditions, are important in applications concerned with geographic variations in illness. When using data sources that have not been validated against a gold standard, it may be more appropriate to report geographic measures of relative risk than absolute risk. This is suitable for applications in which the relationship between risk and risk factors is of primary interest. For example, an ecological correlation study of the relationship between asthma and social and environmental risk factors is likely to produce similar model coefficients across different case ascertainment algorithms. In surveillance applications, where variations and changes in absolute risk are often of interest, precise measures of prevalence and incidence remain important. Our case ascertainment algorithms produced very different estimates in prevalence of asthma in the cohort. Precise prevalence estimates are necessary to understand the actual population burden of disease, and therefore, require data that have been validated against a medically and socially acceptable case definition standard.

**TABLE 2**  
**Tabulation of asthma-related and total service utilization by all case definitions (1998-2004)**

	Baseline definition	Test case-definitions of asthma						Non-asthmatics
		A	B	C	D	E	F	
		n = 1710	n = 685	n = 638	n = 390	n = 244	n = 71	n = 33 795
Percentage of asthmatics in cohort	100.00	30.81	12.34	11.49	7.03	4.40	1.28	N/A
Percentage of the total cohort population	13.13	4.40	1.76	1.64	1.00	0.63	0.18	100.00
<b>Asthma-related services (1998-2004)</b>								
Mean claims	5.04	9.67	10.47	15.52	14.87	12.09	16.12	N/A
Median claims	3	8	8	14	13	9	13	N/A
Mean emergency	1	2.47	5.75	4.76	7.59	10.95	19.26	N/A
Median emergency	0	1	4	2	5.5	8.5	16	N/A
Mean in-patient	0.05	0.16	0.37	0.322	0.523	0.63	1.12	N/A
Median in-patient	0	0	0	0	0	0	0	N/A
<b>All services (1998-2004)</b>								
Mean claims	43.74	49.7	49.18	56.95	55.03	50.35	60.51	29.44
Median claims	36	42	41	47	46	42	46	23
Mean emergency	10.64	12.2	17.76	15.02	19.24	24.95	32.59	7.08
Median emergency	6	7	13	10	14	20	27	3
Mean in-patient	0.33	0.44	0.75	0.66	0.93	1.11	1.55	0.21
Median in-patient	0	0	0	0	0	1	1	0

Although not an explicit objective of our study, our results do reveal interesting geographic patterns of paediatric asthma in Alberta. Firstly, the distribution of relative asthma risk based on Figure 1 suggests lowest risk in rural central and northern Alberta and highest risk in southern Alberta, particularly Calgary. The clusters of asthmatics based on the test case-definitions also identify the Calgary area as the region of highest risk. Together, these observations reflect high prevalence of asthma, as well as a high incidence of emergency and in-patient hospital admissions. The use of these services in particular reflects a high burden of treatment in a subset of children in the Calgary region. It may also be an indicator that asthma is more severe (and in turn, requires more emergency and hospital care) for patients in this part of the province. As noted above, this pattern may be related to the absence of primary care for children in an area of rapid population growth, it could also reflect fundamental differences in asthma epidemiology in this region. This explanation is supported by the apparently distinct pattern of higher prevalence in rural and urban southern Alberta, which

could reflect the role of environmental or meteorological conditions in the region.<sup>16</sup> Further research on the explanation for this geographic pattern, and identifying whether or not it is common to the paediatric or general population, is warranted.

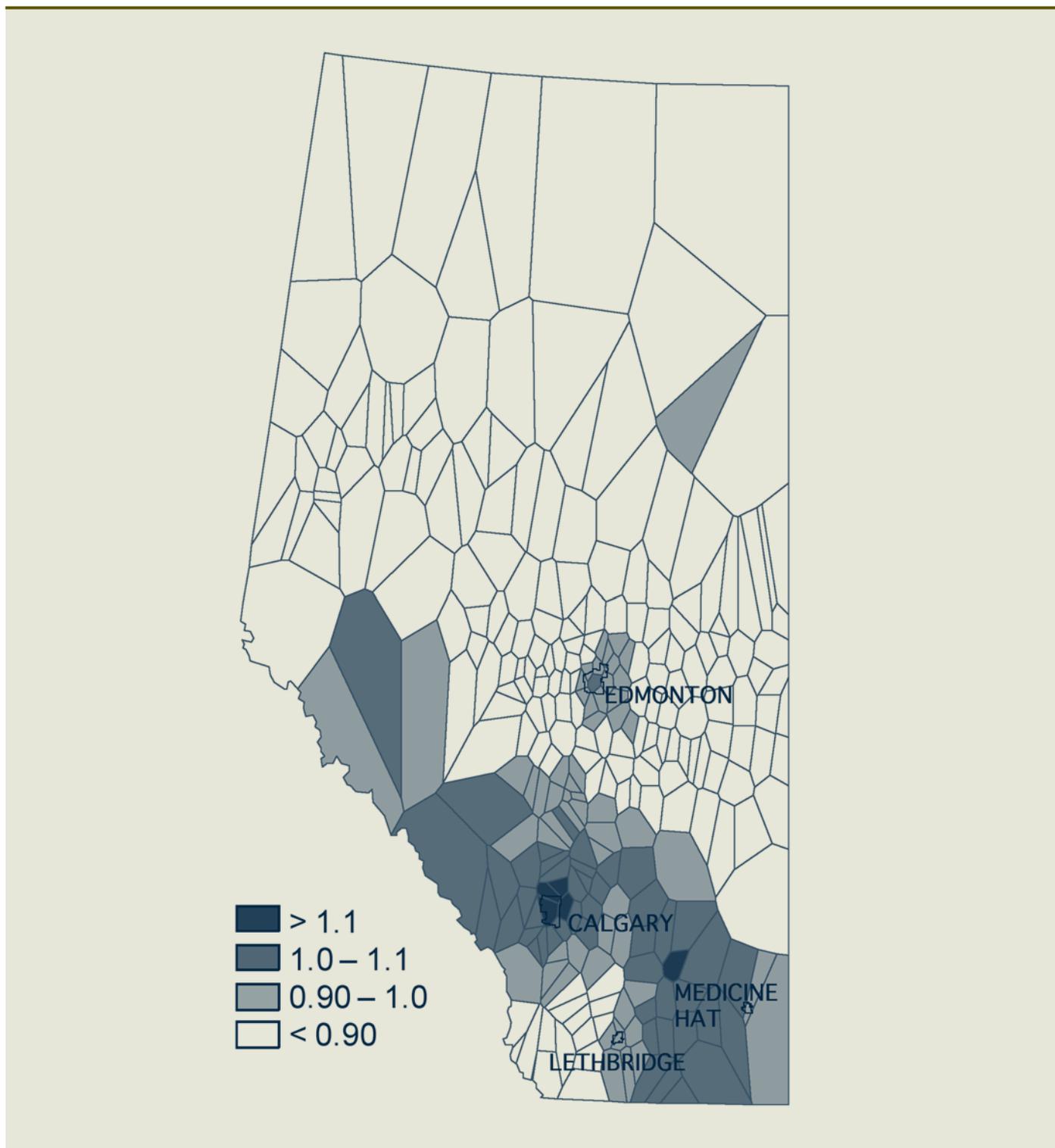
We note two potential limitations to our study. First, we excluded a large number of children immigrating to Alberta between 1998 and 2004. It is possible that children moving into the province have health utilization profiles that are considerably different from those we included in our study group. Immigration into Alberta tends to be into urban areas, and immigration is well known to affect the use of health services. Recent arrivals to the province may be more inclined to use acute-care centres for primary care. This could have resulted in systematic differences in the effectiveness of case definitions, for example, by increasing the sensitivity of definitions "E" and "F" in urban areas. Second, the ability to detect a statistically noteworthy asthma cluster is partly dependent on prevalence, and more specifically, the number of cases. All else being equal, case definitions with a higher overall prevalence are more likely

to produce detectable patterns of clustering. It is possible that clusters were not detected for some definitions and detected for others simply based on the different total number of cases identified. Though this may be a limitation of our study design, we note that most case definitions did result in identifiable clusters of asthma in roughly the same part of the province. Furthermore, the one significant cluster found in the test of constant case definition was the least numerous of all case ascertainment algorithms.

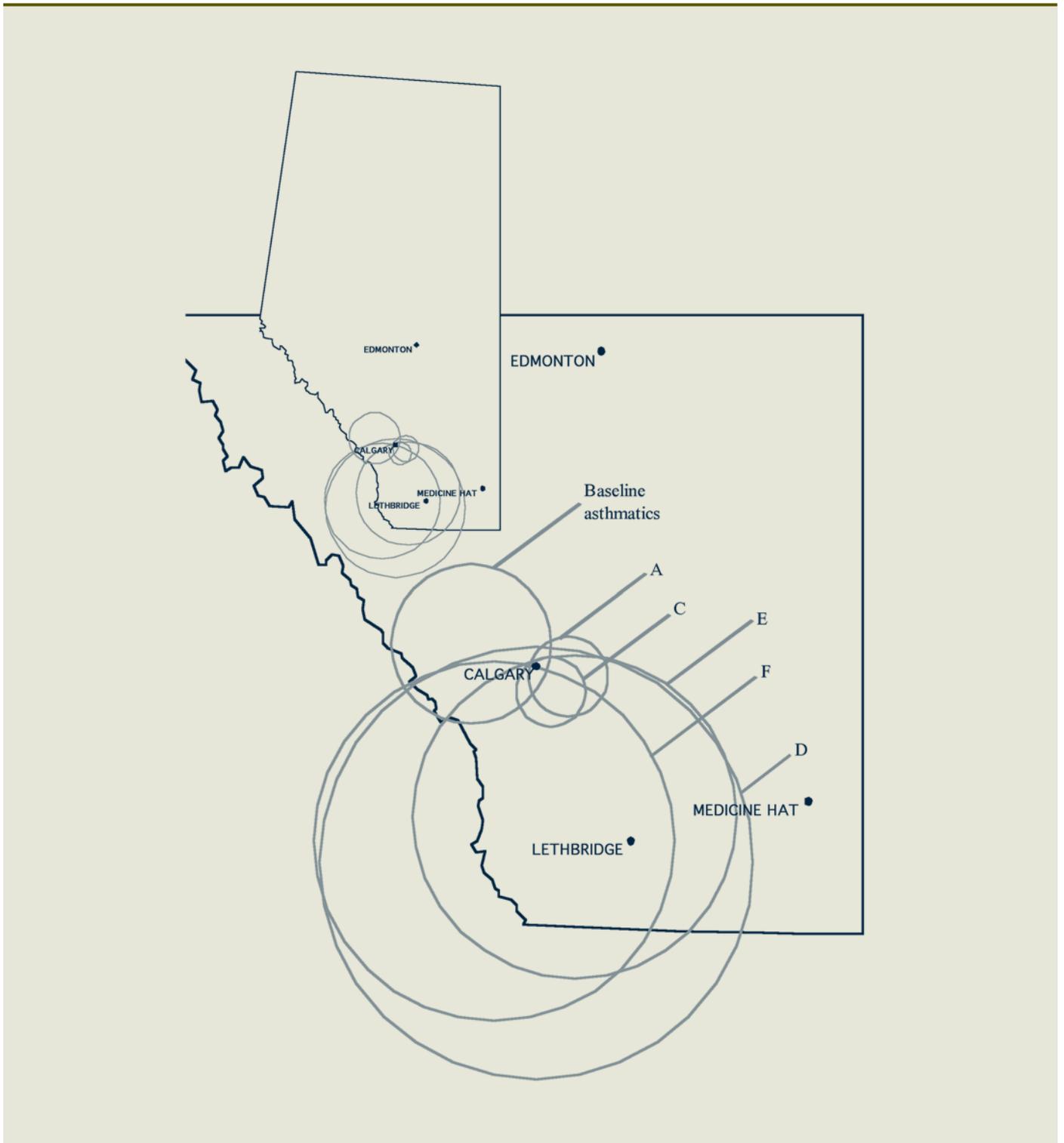
## Conclusion

Administrative data represent an important resource for public health surveillance and research. Validation studies that compare case definitions based on administrative data to clinical assessments and surveys are important for understanding the strengths and weaknesses of these data, as well as determining good estimates of prevalence and incidence. Assessments of relative risk, across geography, time, age, sex, social class and other measures, are also important for a complete understanding of disease epidemiology. Our results suggest that

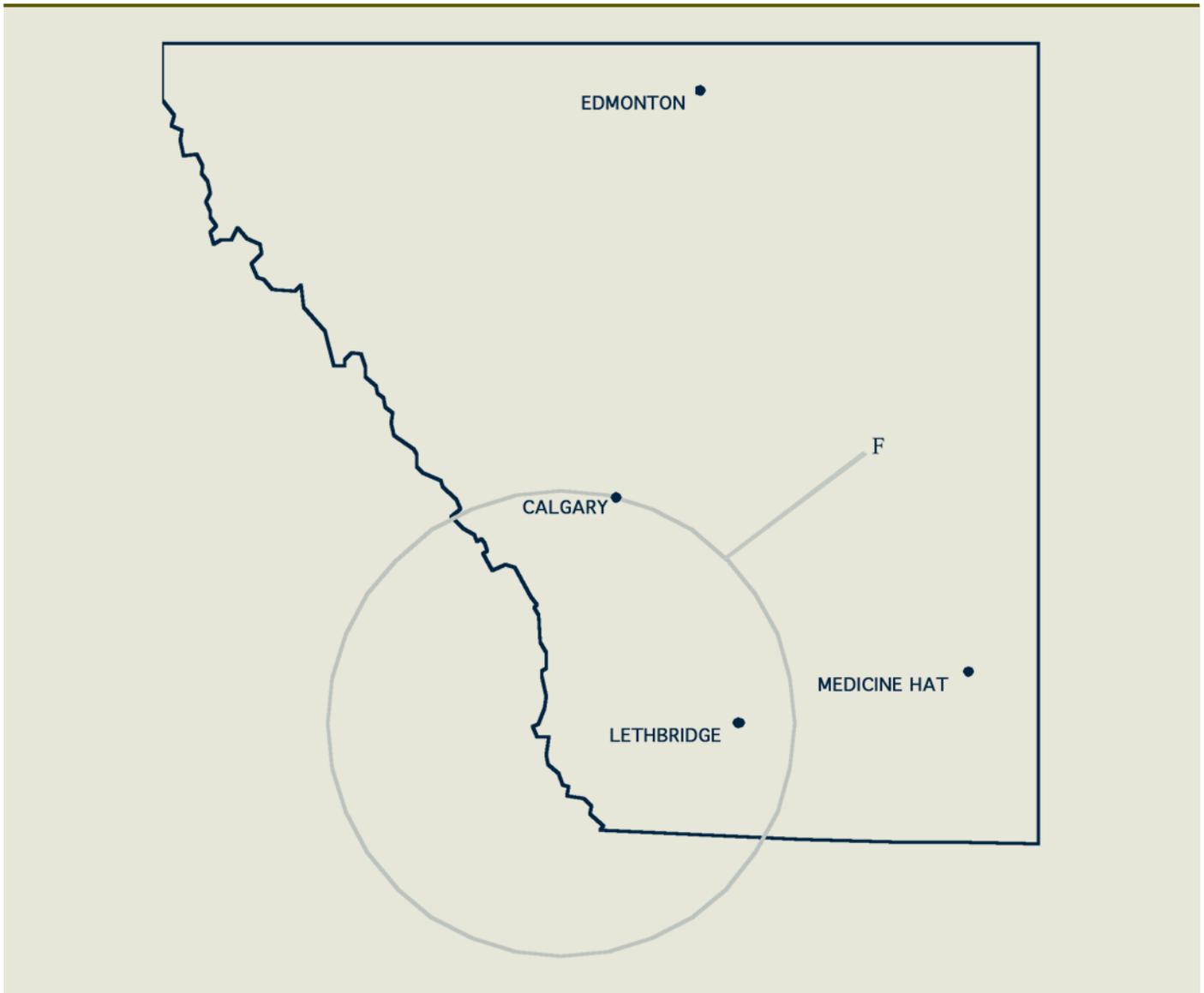
**FIGURE 1**  
**Model estimated relative risk for baseline definition asthmatics**



**FIGURE 2**  
Significant clusters under a null hypothesis of constant-risk



**FIGURE 3**  
Significant clusters under a null hypothesis of constant case definition



relative geographic comparisons of disease based on case definitions from administrative data are not greatly affected by the specifics of the case ascertainment algorithm, even when the case definitions are derived from data from different sources. However, there are considerable variations in prevalence and incidence based across the different definitions, and therefore, routine surveillance requires careful consideration of the precise algorithm used.

### Acknowledgements

This project was funded through the Enhanced Chronic Disease Surveillance Grant Program of the Healthy Living and Chronic Disease Strategy, Public Health Agency of Canada.

### References

1. Tu K, Campbell NR, Chen XL, Cauch-Dudek KJ, McAlister FA. Accuracy of administrative databases in identifying patients with hypertension. *Open Med.* 2007;1:E3-5.
2. Wilchesky M, Tamblyn RM, Huang A. Validation of diagnostic codes within medical services claims. *J Clin Epidemiol.* 2004;57:131-41.
3. Hux JE, Ivis F, Flintoft V, Bica A. Diabetes in Ontario: determination of prevalence and incidence using a validated administrative data algorithm. *Diabetes Care.* 2002; 25:512-6.
4. Robinson JR, Young TK, Roos LL, Gelskey DE. Estimating the burden of disease: comparing administrative data and self-reports. *Med Care.* 1997;35:932-47.

- 
5. Liu L, Reeder B, Shuaib A, Mazagri R. Validity of stroke diagnosis on hospital discharge records in Saskatchewan, Canada: implications for stroke surveillance. *Cerebrovasc Dis.* 1999;9:224-30.
  6. Yiannakoulias N, Svenson LW, Hill MD, Schopflocher DP, James RC, Wielgosz AT, Noseworthy TW. Regional comparisons of in-patient and outpatient patterns of cerebrovascular disease diagnosis in the province of Alberta. *Chron Dis Can.* 2003; 24:9-16.
  7. Yiannakoulias N, Svenson LW, Hill MD, Schopflocher DP, Rowe BH, James RC, Wielgosz AT, Noseworthy TW. Incident cerebrovascular disease in rural and urban Alberta. *Cerebrovasc Dis.* 2004;17:72-8.
  8. Tirschwell DL, Longstreth WT Jr. Validating administrative data in stroke research. *Stroke.* 2002;33:2465-70.
  9. Donahue JG, Weiss ST, Goetsch MA, Livingston JM, Greineder DK, Platt R. Assessment of asthma using automated and full-text medical records. *J Asthma.* 1997;34:273-81.
  10. Kokotailo RA, Hill MD. Coding of stroke and stroke risk factors using international classification of diseases, revisions 9 and 10. *Stroke.* 2005;36:1776-81.
  11. Huzel L, Roos LL, Anthonisen NR, Manfreda J. Diagnosing asthma: the fit between survey and administrative database. *Can Respir J.* 2002;9:407-12.
  12. Breslow NE, Clayton DG. Approximate inference in generalized linear mixed models. *J Am Stat Assoc.* 1993;88:9-25.
  13. SAS [computer program]. Version 9.1.3. Cary, NC: SAS Institute; 2006.
  14. Kulldorff M. A spatial scan statistic. *Commun Stat Theory.* 1997;26:1481-96.
  15. SaTScan [computer program]. Version 6.1. Bethesda, MD. Kulldorff M & Information Management Services, Inc.; 2006.
  16. Verhoef MJ, Rose MS, Ramcharan S. The relationship between chinook conditions and women's physical and mental well-being. *Int J Biometeorol.* 1995;38:148-51.