# Evaluation of Spatial Estimation Methods for Cluster Detection

**Jian Xing[1], Howard Burkom[2], Michael Leuze[3], James Edgerton[2], John Copeland[1], Steve Bloom[3], Jerome Tokars[1]**

*[1]Centers for Disease Control and Prevention, [2]The Johns Hopkins Applied Physics Lab, [3]Science Applications Incorporated*

## OBJECTIVE

We applied spatial scan statistics to data from CDC's BioSense system and examined the effect of the spatial prediction method on determination of anomalous disease clusters. The objectives were to decide on a reliable spatial estimation method for one BioSense data source and to establish criteria for making this decision using other sources.

## BACKGROUND

CDC's BioSense system provides near-real time situational awareness for public health monitoring through analysis of electronic health data. Determination of anomalous spatial and temporal disease clusters is a crucial part of the daily disease monitoring task. Spatial approaches depend strongly on having reliable estimated values for counts among the geographic sub-regions. If estimates are poor, algorithms will find irrelevant clusters, and clusters of importance may be missed. While many studies have focused on improved computation time and more general cluster shapes, our effort focused on finding anomalies that are correct according to available BioSense data history.

## METHODS

The study dataset included records from 32 military treatment facilities in Texas for January 1, 2004 to December 31, 2006. These records were classified into the standard 11 BioSense syndrome groups; only respiratory syndrome data is presented here. Using a newly developed computer program implementing scan statistics, we compared three spatial estimation methods using the clinic zip code. For the first estimation method, we calculated sliding baseline averages for individual zip codes with a 2-day buffer between baseline and test day. For a 7-day baseline, this method gives the same expected value for each zip code as that used for the temporal EARS C2 algorithm [1]. A second estimation method was analogous to the W2 weekday/non-weekday stratification of the C2 method. For a third estimation method, we derived expected counts by conditioning on marginal totals for both space and time, as in the SaTScan space-time permutation method [2]. This method includes the test day counts in the baseline and does not use a buffer interval. We applied each method using 1-, 4-, and 8-week baselines. We compared observed with expected counts using the Chi-square goodness-of-fit test for weekdays and weekend days.

For the entire 3 years of respiratory syndrome data, we then applied our program to determine 1-day (i.e., spatial) clusters significant at a p-value threshold value of 0.01.

## RESULTS

The data included 863,774 records (793,842 weekday, 69,932 weekend), with means of 789 records per day and 46,953 per zip code. In the sample results given in the table, the stratified mean method has the lowest cluster determination rate, especially on weekends. For the other methods, the number of clusters is larger on Sundays, when counts are much smaller, than on Mondays, and many of these clusters are likely the result of poor spatial estimation on weekend days.. The spatial goodness-of-fit statistic actually worsens when a longer baseline is used for the mean method, suggesting that seasonal baseline changes overcome the longer baseline stability. The space-time marginal and stratified mean methods both give improved goodness-of-fit, but correct management of the day-of-week effect is needed achieve a significant reduction in nuisance clusters.

| Estimation Method | Baseline (days) | Clusters $p < 0.01$ | Mon only | Sun only | $\chi^2..$ |
|---|---|---|---|---|---|
| Mean | 7 | 652 | 115 | 144 | 73.9 |
| Mean | 28 | 643 | 117 | 147 | 75.8 |
| Space-Time Marginals | 28 | 614 | 116 | 153 | 66.7 |
| Stratified Mean | 28 | 342 | 62 | 29 | 61.3 |

Table –Methods, clusters found, and $\chi^2.$ results for respiratory syndrome data.

## CONCLUSIONS

Preliminary results demonstrate substantial differences in cluster determination depending on the chosen method of spatial estimation. The overall cluster rate is high for this syndrome group because of its seasonality and noisiness in daily counts. Scan statistics will be applied using additional methods, syndromes, and effectiveness measures to determine best daily cluster detection practices for this data source.

## REFERENCES

[1] Hutwagner L, et al. The Bioterrorism Preparedness and Response Early Aberration Reporting System (EARS). J Urban Health. 2003 ;80 1:89-96.

[2] Kulldorff M, et al. (2005) A space-time permutation scan statistic for the early detection of disease outbreaks. PLoS Medicine, 2:216-224.