

# Detection Algorithms for Biosurveillance: A tutorial

|                |  |   |  |
|----------------|--|---|--|
| Greg Cooper    | Professor  | Computer Science and<br>RODS lab, U. Pitt | <a href="mailto:gfc@cbmi.upmc.edu">gfc@cbmi.upmc.edu</a>     |
| Bill Hogan     | Assistant Professor                                      | RODS lab, U. Pitt                         | <a href="mailto:wrh@cbmi.pitt.edu">wrh@cbmi.pitt.edu</a>     |
| Andrew Moore   | Professor and co-director of<br>Auton Lab                | Computer Science,<br>Carnegie Mellon      | <a href="mailto:awm@cs.cmu.edu">awm@cs.cmu.edu</a>           |
| Robin Sabhnani | Senior Programmer, Auton<br>Lab                          | Robotics Institute,<br>Carnegie Mellon    | <a href="mailto:sabhnani@cs.cmu.edu">sabhnani@cs.cmu.edu</a> |
| Rich Tsui      | Research Professor and<br>associate Director of RODS lab | RODS lab, U. Pitt                         | <a href="mailto:tsui@cbmi.pitt.edu">tsui@cbmi.pitt.edu</a>   |
| Mike Wagner    | Professor and Director of<br>RODS lab                    | RODS lab, U. Pitt                         | <a href="mailto:mmw@cbmi.pitt.edu">mmw@cbmi.pitt.edu</a>     |
| Weng-Keen Wong | Graduate Student, Auton Lab                              | Computer Science,<br>Carnegie Mellon      | <a href="mailto:wkw@cs.cmu.edu">wkw@cs.cmu.edu</a>           |

Tutorial slides by Andrew Moore

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials>. Comments and corrections gratefully received.

RODS: <http://www.health.pitt.edu/rods>  
Auton Lab: <http://www.autonlab.org>

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 1

## Many Methods!

| Method                    | Has Pitt/CMU tried it? | Tried but little used | Tried and used | Under development  | Multivariate signal tracking? | Spatial ? |
|---------------------------|------------------------|-----------------------|----------------|--------------------|-------------------------------|-----------|
| Time-weighted averaging   | Yes                    | Yes                   |                |                    |                               |           |
| Serfling                  | Yes                    |                       | Yes            |                    |                               |           |
| ARIMA                     | Yes                    | Yes                   |                |                    |                               |           |
| SARIMA + External Factors | Yes                    |                       | Yes            |                    |                               |           |
| Univariate HMM            | Yes                    |                       | Yes            |                    |                               |           |
| Kalman Filter             | Yes                    | Yes                   |                |                    |                               |           |
| Recursive Least Squares   | Yes                    |                       | Yes            |                    |                               |           |
| Support Vector Machine    | Yes                    | Yes                   |                |                    |                               |           |
| Neural Nets               | Yes                    | Yes                   |                |                    |                               |           |
| Randomization             | Yes                    |                       | Yes            | Yes                |                               |           |
| Spatial Scan Statistics   | Yes                    |                       |                | (w/ Howard Burkom) | Yes                           | Yes       |
| Bayesian Networks         | Yes                    |                       |                | Yes                | Yes                           |           |
| Contingency Tables        | Yes                    |                       | Yes            |                    |                               |           |
| Scalar Outlier (SOC)      | Yes                    | Yes                   |                |                    |                               |           |
| Multivariate Anomalies    | Yes                    |                       | Yes            |                    | Yes                           |           |
| Change-point statistics   | Yes                    |                       |                | Yes                |                               |           |
| FDR Tests                 | Yes                    |                       | Yes            |                    | Yes                           |           |
| WSARE (Recent patterns)   | Yes                    |                       | Yes            | Yes                | Yes                           | Yes       |
| PANDA (Causal Model)      | Yes                    |                       |                | Yes                | Yes                           | Yes       |
| FLUMOD (space/Time HMM)   |                        |                       |                | Yes                | Yes                           | Yes       |

Details of these methods and bibliography available from "Summary of Biosurveillance-relevant statistical and data mining technologies" by Moore, Cooper, Tsui and Wagner. Downloadable (PDF format) from [www.cs.cmu.edu/~awm/biosurv-methods.pdf](http://www.cs.cmu.edu/~awm/biosurv-methods.pdf)

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 2

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 3

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

These are all powerful statistical methods, which means they all have to have one thing in common...

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 4

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

These are all powerful statistical methods, which means they all have to have one thing in common...

*Boring Names.*

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 5

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

These are all powerful statistical methods, which means they all have to have one thing in common...

*Boring Names.*

Univariate Anomaly Detection

Multivariate Anomaly Detection

Spatial Scan Statistics

WSARE

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 6

## What you'll learn about

- **Noticing events in bio-event time series**
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

WSARE

Spatial Scan Statistics

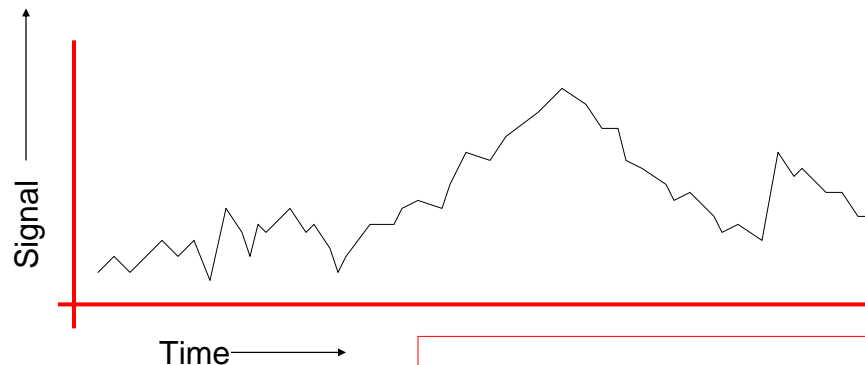
Univariate Anomaly Detection

Multivariate Anomaly Detection

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 7

## Univariate Time Series



### Example Signals:

- Number of ED visits today
- Number of ED visits this hour
- Number of Respiratory Cases Today
- School absenteeism today
- Nyquil Sales today

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 8

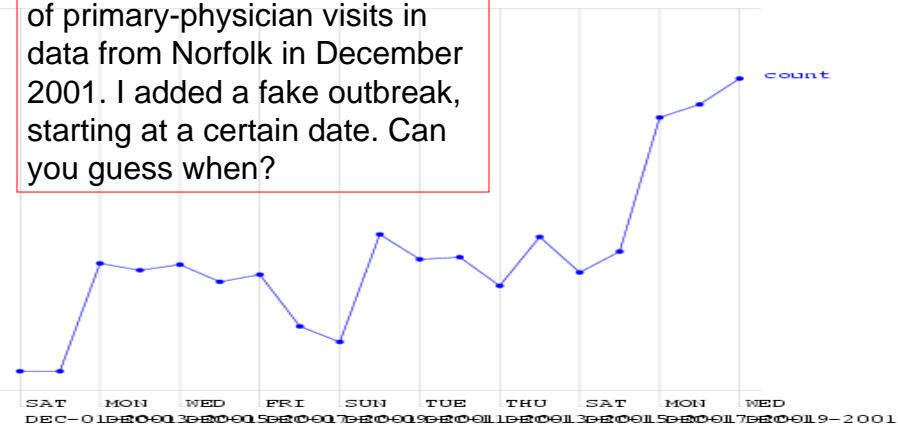
## (When) is there an anomaly?

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 9

## (When) is there an anomaly?

This is a time series of counts of primary-physician visits in data from Norfolk in December 2001. I added a fake outbreak, starting at a certain date. Can you guess when?



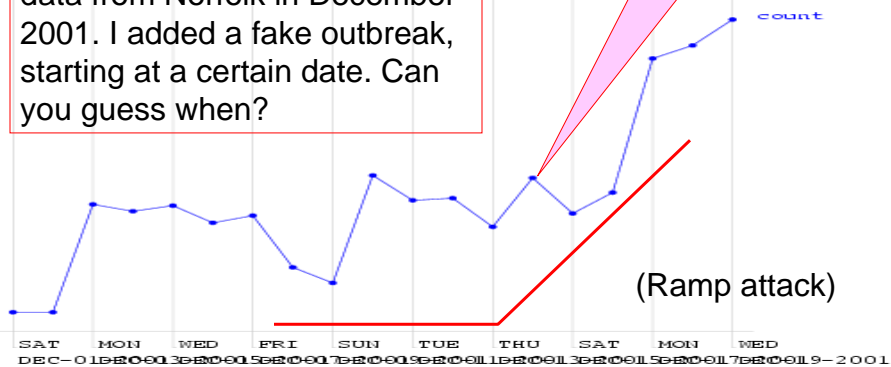
Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 10

## (When) is there an anomaly?

This is a time series of counts of primary-physician visits in data from Norfolk in December 2001. I added a fake outbreak, starting at a certain date. Can you guess when?

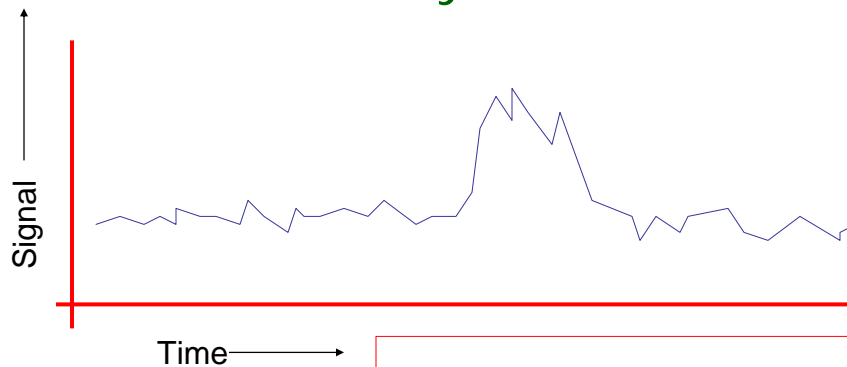
Here (much too high for a Friday)



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 11

## An easy case

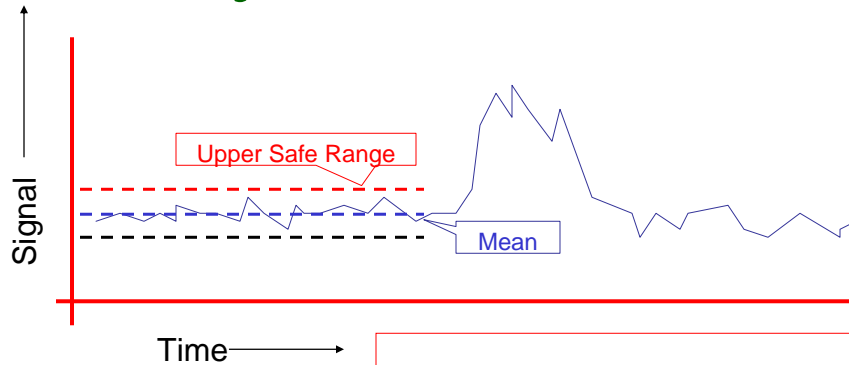


Dealt with by Statistical Quality Control  
Record the mean and standard deviation up to the current time.  
Signal an alarm if we go outside 3 sigmas

Copyright © 2002, 2003, Andrew Moore

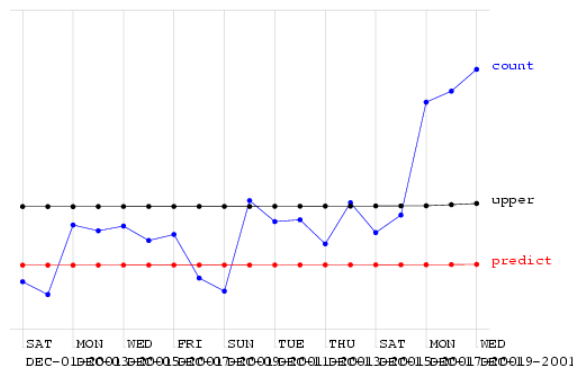
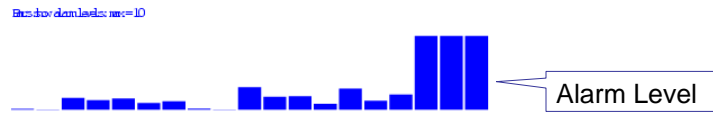
Biosurveillance Detection Algorithms: Slide 12

# An easy case: Control Charts

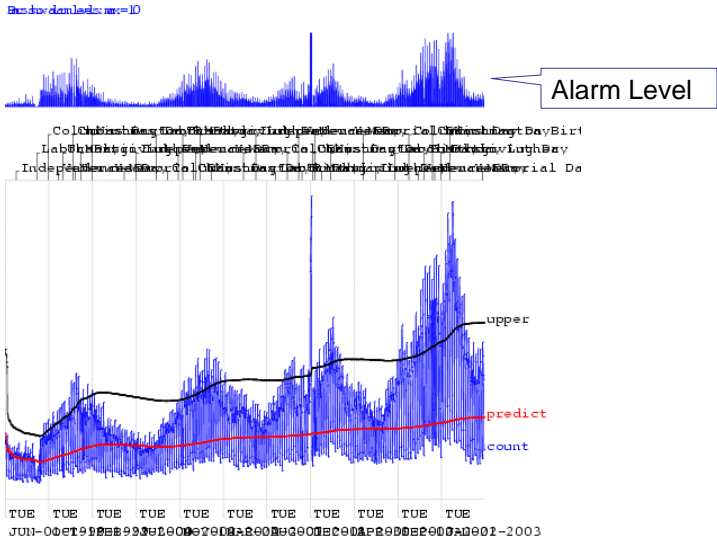


Dealt with by Statistical Quality Control  
 Record the mean and standard deviation up to the current time.  
 Signal an alarm if we go outside 3 sigmas

# Control Charts on the Norfolk Data



# Control Charts on the Norfolk Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 15

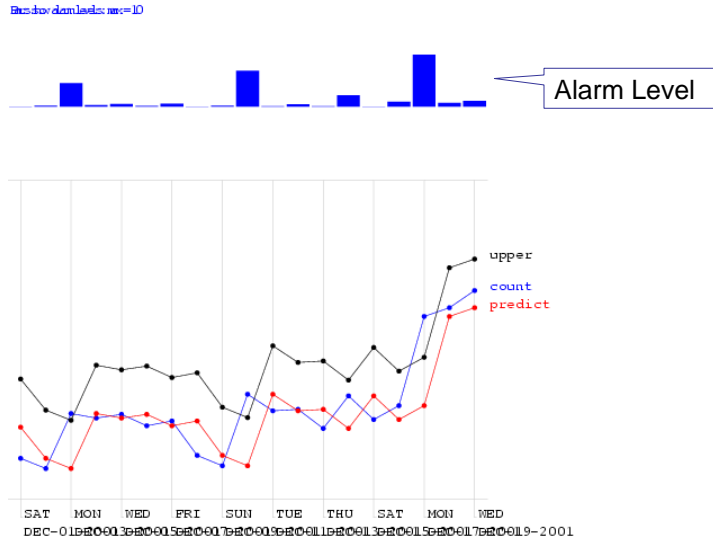
# Looking at changes from yesterday

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 16



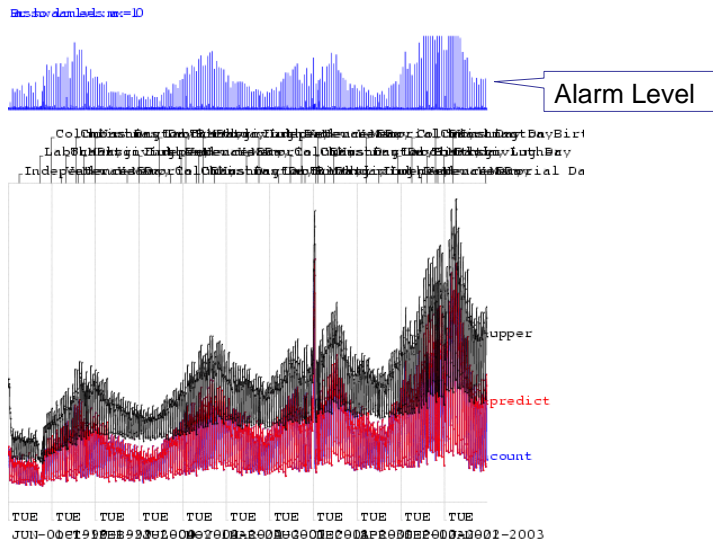
# Looking at changes from yesterday



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 17

# Looking at changes from yesterday

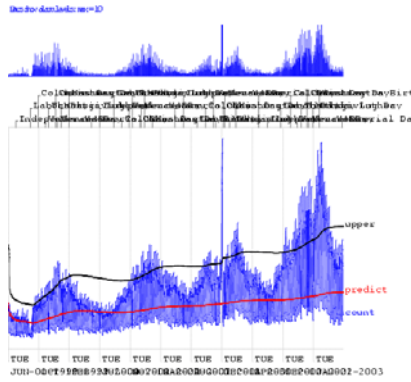


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 18

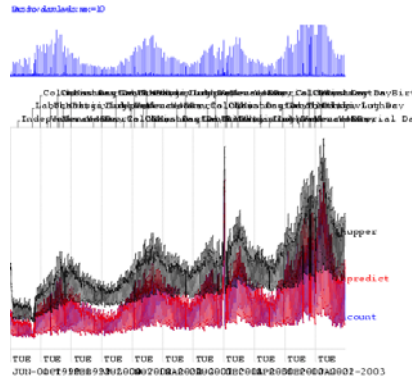
## We need a happy medium:

Control Chart:  
Too insensitive to recent  
changes



Copyright © 2002, 2003, Andrew Moore

Change from yesterday:  
Too sensitive to recent  
changes



Biosurveillance Detection Algorithms: Slide 19

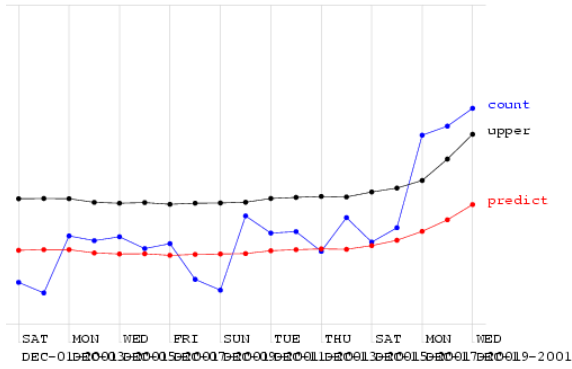
## Moving Average

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 20

# Moving Average

Enclosure: n=7387



Copyright © 2002, 2003, Andrew Moore

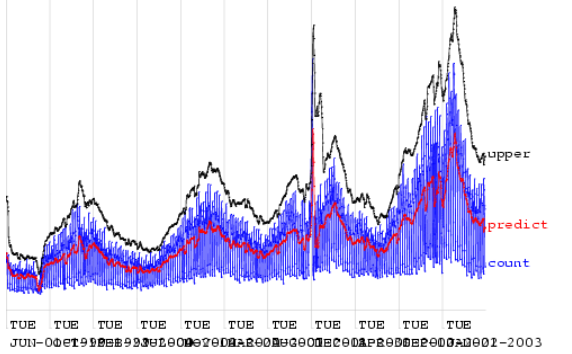
Biosurveillance Detection Algorithms: Slide 21

# Moving Average

Enclosure: n=7387



Colombia, Cuba, Czech Republic, Denmark, Germany, Greece, Hungary, India, Israel, Italy, Japan, Korea, Latvia, Lithuania, Luxembourg, Malaysia, Mexico, Netherlands, New Zealand, Norway, Poland, Portugal, Romania, Russia, Saudi Arabia, Singapore, Slovakia, Slovenia, South Africa, Spain, Sweden, Switzerland, Taiwan, Thailand, Turkey, United Kingdom, United States, Uruguay, Venezuela, Viet Nam, Zimbabwe

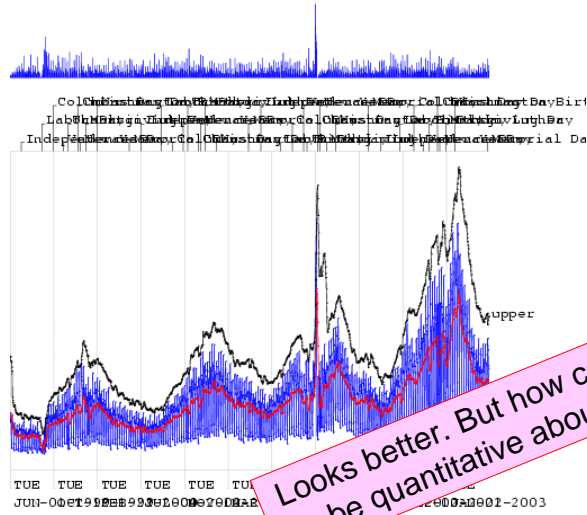


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 22

# Moving Average

https://dam.assets.nrc.gov/73807



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 23

## Algorithm Performance

Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

|                        | Allowing one False Alarm per TWO weeks... |                              | Allowing one False Alarm per SIX weeks... |                              |
|------------------------|---|------------------------------|---|------------------------------|
|                        | Fraction of spikes detected               | Days to detect a ramp attack | Fraction of spikes detected               | Days to detect a ramp attack |
| standard control chart | 0.39                                      | 3.47                         | 0.22                                      | 4.13                         |
| using yesterday        | 0.14                                      | 3.83                         | 0.1                                       | 4.7                          |

# Algorithm Performance

Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

Fraction of spikes detected  
Days to detect a ramp attack

|                        | Allowing one False Alarm per TWO weeks... |                              | Allowing one False Alarm per SIX weeks... |                              |
|------------------------|---|------------------------------|---|------------------------------|
|                        | Fraction of spikes detected               | Days to detect a ramp attack | Fraction of spikes detected               | Days to detect a ramp attack |
| standard control chart | 0.39                                      | 3.47                         | 0.22                                      | 4.13                         |
| using yesterday        | 0.14                                      | 3.83                         | 0.1                                       | 4.7                          |
| ▶ Moving Average 7     | 0.58                                      | 2.79                         | 0.51                                      | 3.31                         |

# Algorithm Performance

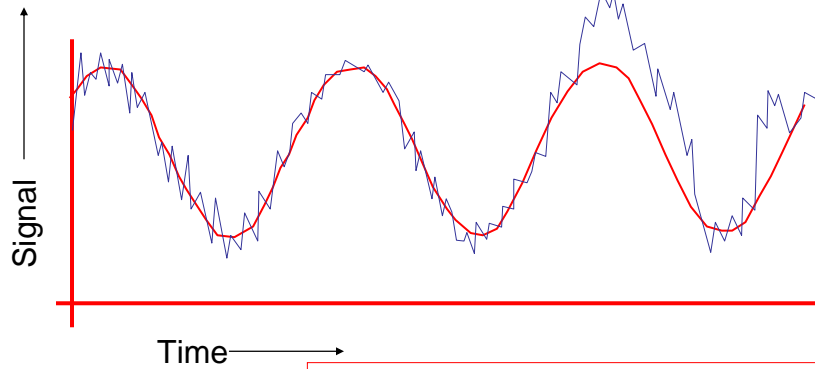
Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

Fraction of spikes detected  
Days to detect a ramp attack

|                        | Allowing one False Alarm per TWO weeks... |                              | Allowing one False Alarm per SIX weeks... |                              |
|------------------------|---|------------------------------|---|------------------------------|
|                        | Fraction of spikes detected               | Days to detect a ramp attack | Fraction of spikes detected               | Days to detect a ramp attack |
| standard control chart | 0.39                                      | 3.47                         | 0.22                                      | 4.13                         |
| using yesterday        | 0.14                                      | 3.83                         | 0.1                                       | 4.7                          |
| Moving Average 3       | 0.36                                      | 3.45                         | 0.33                                      | 3.79                         |
| Moving Average 7       | 0.58                                      | 2.79                         | 0.51                                      | 3.31                         |
| Moving Average 56      | 0.54                                      | 2.72                         | 0.44                                      | 3.54                         |

## Seasonal Effects



Fit a periodic function (e.g. sine wave) to previous data. Predict today's signal and 3-sigma confidence intervals. Signal an alarm if we're off.

Reduces False alarms from Natural outbreaks.

Different times of year deserve different thresholds.

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 27

## Algorithm Performance

Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

|                        | Allowing one False Alarm per TWO weeks... |                              | Allowing one False Alarm per SIX weeks... |                              |
|------------------------|---|------------------------------|---|------------------------------|
|                        | Fraction of spikes detected               | Days to detect a ramp attack | Fraction of spikes detected               | Days to detect a ramp attack |
| standard control chart | 0.39                                      | 3.47                         | 0.22                                      | 4.13                         |
| using yesterday        | 0.14                                      | 3.83                         | 0.1                                       | 4.7                          |
| Moving Average 3       | 0.36                                      | 3.45                         | 0.33                                      | 3.79                         |
| Moving Average 7       | 0.58                                      | 2.79                         | 0.51                                      | 3.31                         |
| Moving Average 56      | 0.54                                      | 2.72                         | 0.44                                      | 3.54                         |
| ▶ hours_of_daylight    | 0.58                                      | 2.73                         | 0.43                                      | 3.9                          |

## Day-of-week effects

Fit a day-of-week component

$$E[\text{Signal}] = a + \text{delta}_{\text{day}}$$

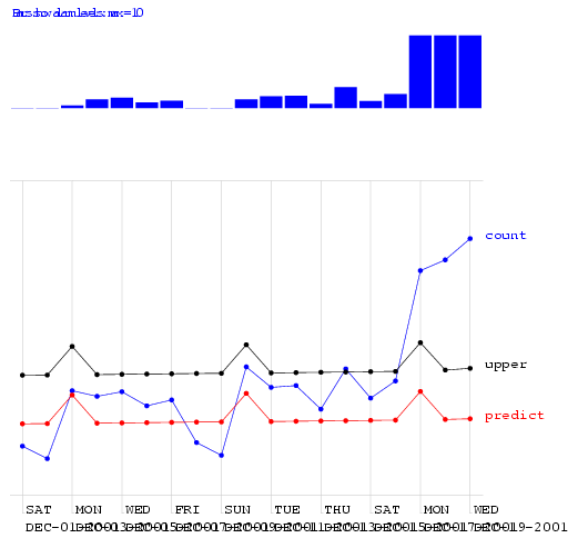
E.G:  $\text{delta}_{\text{mon}} = +5.42$ ,  $\text{delta}_{\text{tue}} = +2.20$ ,  $\text{delta}_{\text{wed}} = +3.33$ ,  $\text{delta}_{\text{thu}} = +3.10$ ,  $\text{delta}_{\text{fri}} = +4.02$ ,  $\text{delta}_{\text{sat}} = -12.2$ ,  $\text{delta}_{\text{sun}} = -23.42$

A simple form  
of ANOVA

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 29

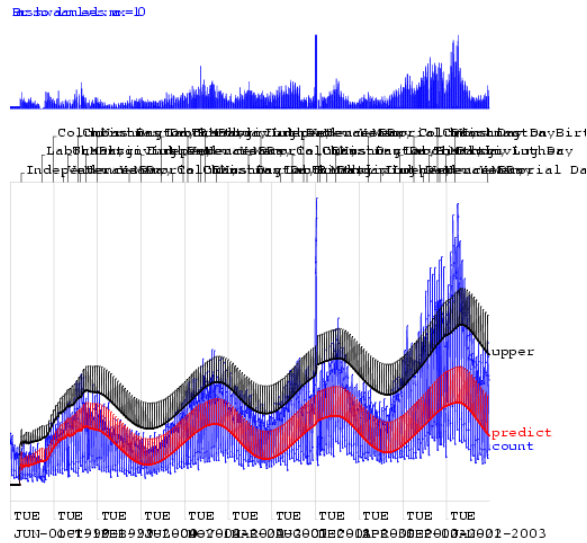
## Regression using Hours-in-day & IsMonday



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 30

## Regression using Hours-in-day & IsMonday



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 31

## Algorithm Performance

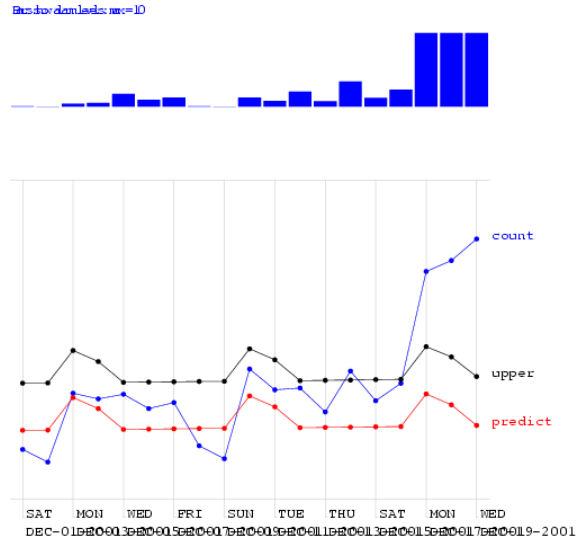
Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

|                          | Allowing one False Alarm per TWO weeks... |                              | Allowing one False Alarm per SIX weeks... |                              |
|--------------------------|---|------------------------------|---|------------------------------|
|                          | Fraction of spikes detected               | Days to detect a ramp attack | Fraction of spikes detected               | Days to detect a ramp attack |
| standard control chart   | 0.39                                      | 3.47                         | 0.22                                      | 4.13                         |
| using yesterday          | 0.14                                      | 3.83                         | 0.1                                       | 4.7                          |
| Moving Average 3         | 0.36                                      | 3.45                         | 0.33                                      | 3.79                         |
| Moving Average 7         | 0.58                                      | 2.79                         | 0.51                                      | 3.31                         |
| Moving Average 56        | 0.54                                      | 2.72                         | 0.44                                      | 3.54                         |
| hours_of_daylight        | 0.58                                      | 2.73                         | 0.43                                      | 3.9                          |
| hours_of_daylight is_mon | 0.7                                       | 2.25                         | 0.57                                      | 3.12                         |



# Regression using Mon-Tue



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 33

## Algorithm Performance

Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

Fraction of spikes detected  
Days to detect a ramp attack  
Fraction of spikes detected  
Days to detect a ramp attack

|                                     |      |      |      |      |
|-------------------------------------|------|------|------|------|
| standard control chart              | 0.39 | 3.47 | 0.22 | 4.13 |
| using yesterday                     | 0.14 | 3.83 | 0.1  | 4.7  |
| Moving Average 3                    | 0.36 | 3.45 | 0.33 | 3.79 |
| Moving Average 7                    | 0.58 | 2.79 | 0.51 | 3.31 |
| Moving Average 56                   | 0.54 | 2.72 | 0.44 | 3.54 |
| hours_of_daylight                   | 0.58 | 2.73 | 0.43 | 3.9  |
| hours_of_daylight is_mon            | 0.7  | 2.25 | 0.57 | 3.12 |
| hours_of_daylight is_mon ... is_tue | 0.72 | 1.83 | 0.57 | 3.16 |
| hours_of_daylight is_mon ... is_sat | 0.77 | 2.11 | 0.59 | 3.26 |

# CUSUM

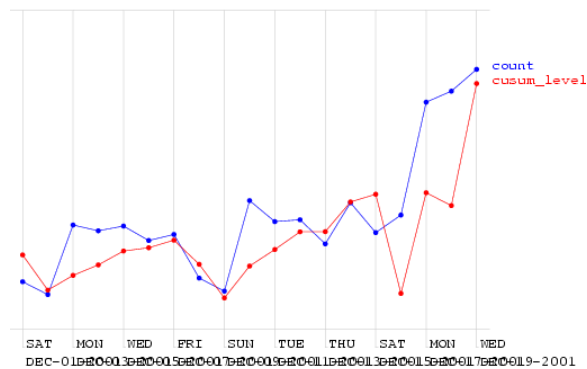
- Cumulative SUM Statistics
- Keep a running sum of “surprises”: a sum of excesses each day over the prediction
- When this sum exceeds threshold, signal alarm and reset sum

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 35

# CUSUM

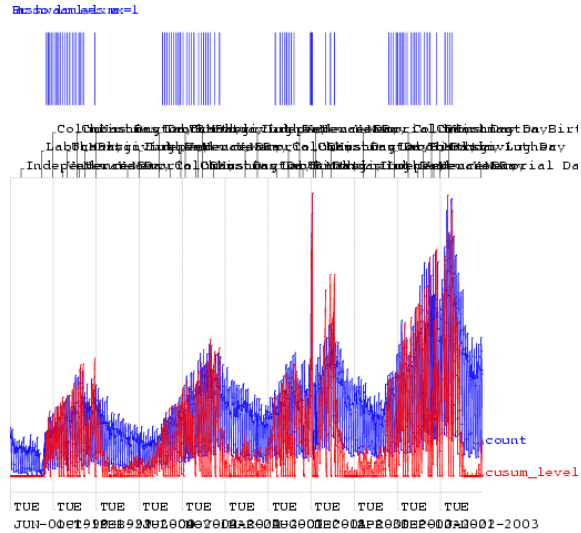
`threshold=level*nc=1`



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 36

# CUSUM



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 37

## Algorithm Performance

Allowing one False Alarm per TWO weeks...

Allowing one False Alarm per SIX weeks...

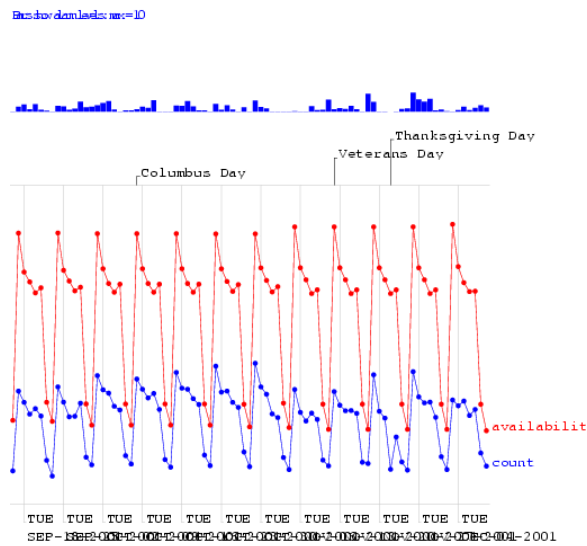
|                                     | Fraction of spikes detected | Days to detect a ramp attack | Fraction of spikes detected | Days to detect a ramp attack |
|-------------------------------------|-----------------------------|------------------------------|-----------------------------|------------------------------|
| standard control chart              | 0.39                        | 3.47                         | 0.22                        | 4.13                         |
| using yesterday                     | 0.14                        | 3.83                         | 0.1                         | 4.7                          |
| Moving Average 3                    | 0.36                        | 3.45                         | 0.33                        | 3.79                         |
| Moving Average 7                    | 0.58                        | 2.79                         | 0.51                        | 3.31                         |
| Moving Average 56                   | 0.54                        | 2.72                         | 0.44                        | 3.54                         |
| hours_of_daylight                   | 0.58                        | 2.73                         | 0.43                        | 3.9                          |
| hours_of_daylight is_mon            | 0.7                         | 2.25                         | 0.57                        | 3.12                         |
| hours_of_daylight is_mon ... is_tue | 0.72                        | 1.83                         | 0.57                        | 3.16                         |
| hours_of_daylight is_mon ... is_sat | 0.77                        | 2.11                         | 0.59                        | 3.26                         |
| ▶ CUSUM                             | 0.45                        | 2.03                         | 0.15                        | 3.55                         |

# The Sickness/Availability Model

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 39

# The Sickness/Availability Model



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 40

that the date is a Sunday. It is also possible to add local terms, such as the mean count over the previous seven days as an additional feature to allow the method to adapt better to recent local events. These additions can be very helpful, and as we shall see at the end of the chapter, regression methods that include these extra terms perform better than moving average (which had been our favorite method) on our data.

## 10. SICKNESS AVAILABILITY

Another way to deal with day-of-week variations is to use the sickness availability method to smooth the time series by removing noise due to the day-of-week effect. This algorithm transforms the daily counts in the time series into a daily sickness value, which is defined as the number of people getting sick every day irrespective of whether they seek health care or not. The term availability refers to the probability that a patient will seek health care during a specific day of the week; hence, there are a total of seven values of availability, one for each day of the week. The availability of a day can be thought of as the fraction of a weeks-worth of visits that get assigned to the given day. The sickness availability method is based on the intuitive assumptions that the expected count is the product of the true amount of sickness and the current day's availability.

We can estimate the expected availability for a specific day of week (*dow*) using the average of the availabilities on that day for the past *m* weeks. We can calculate the expected availability  $A_{dow}$  as:

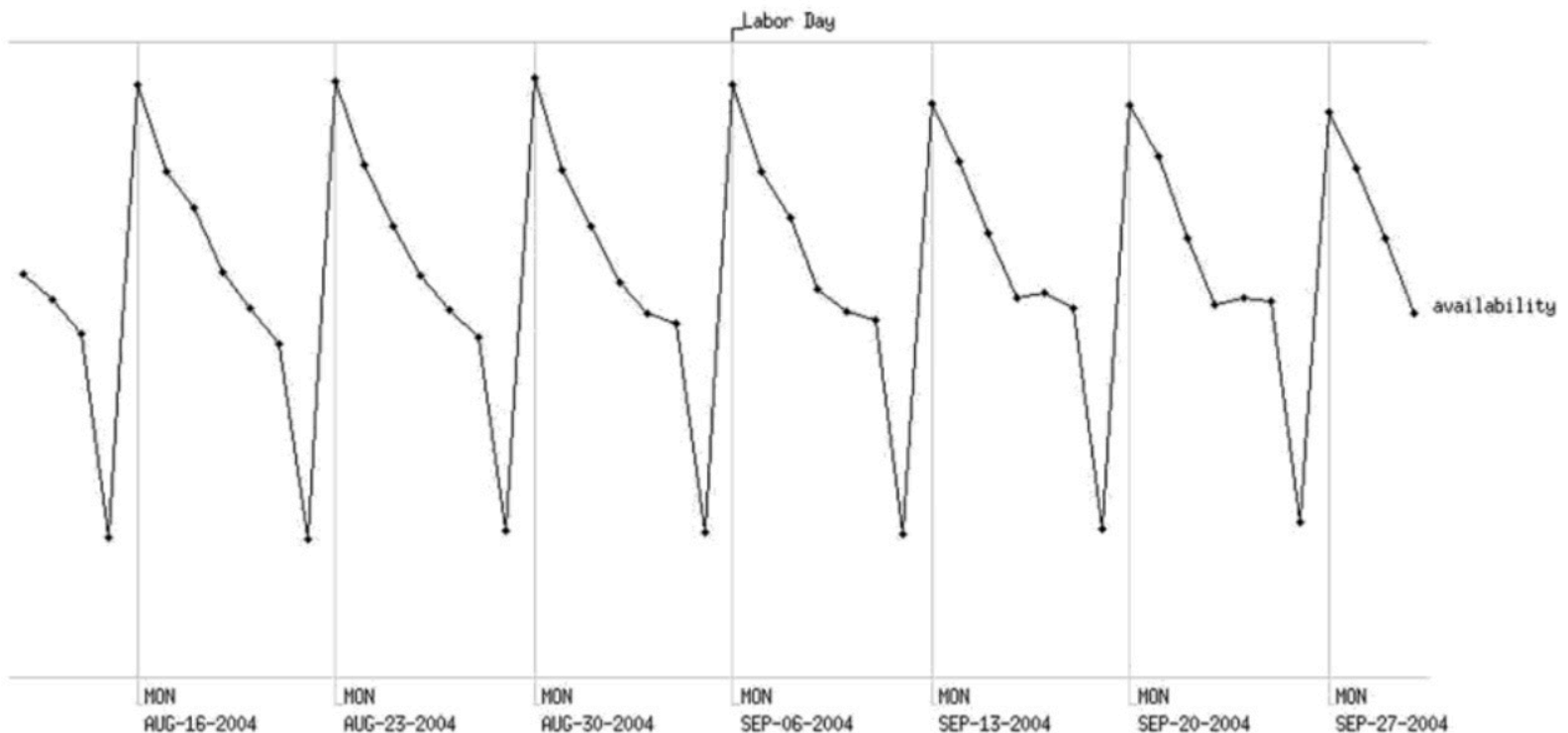
$$A_{dow} = \frac{\sum_{i=1}^m (C_{i(dow)} / \sum_{dow=0}^6 C_{i(dow)})}{m} \quad (9)$$

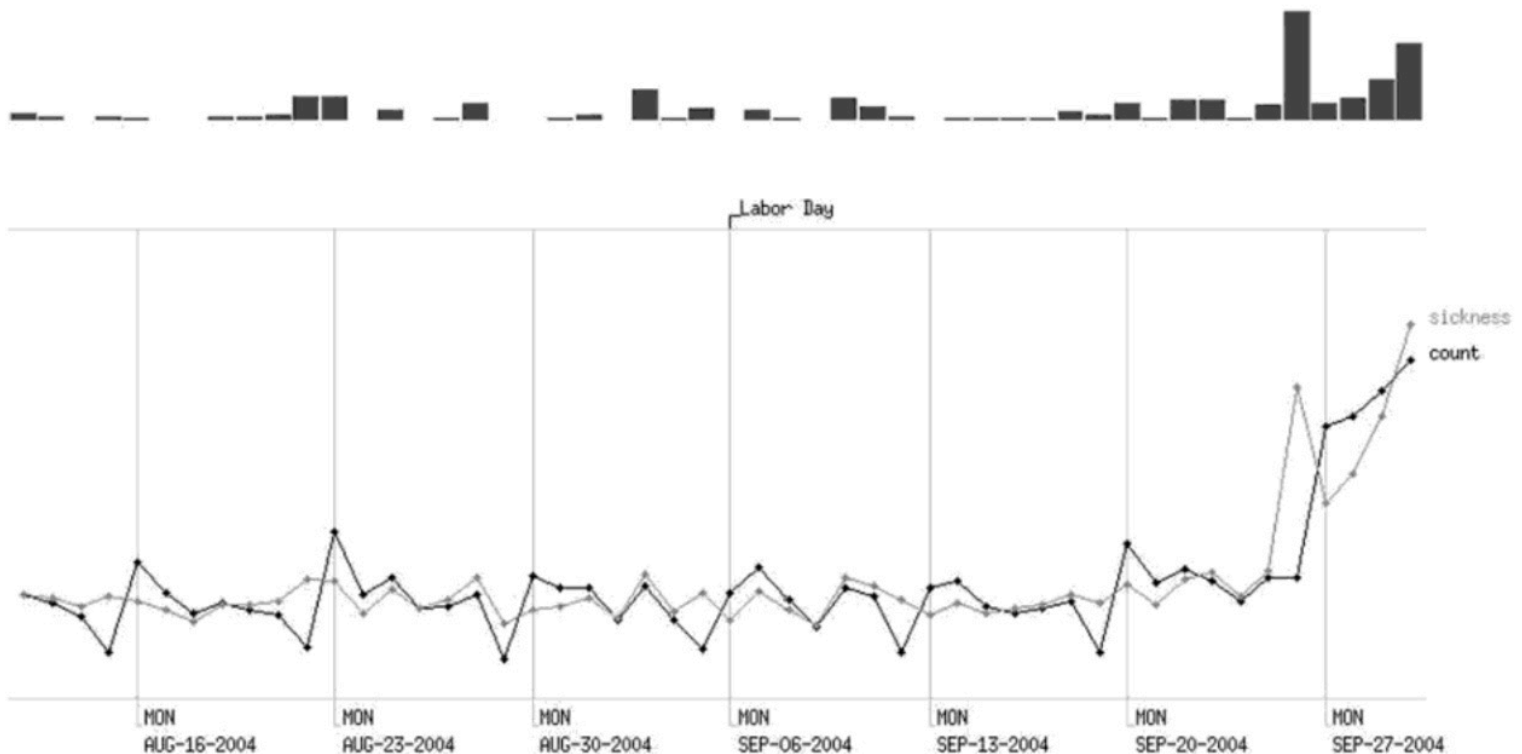
In Eq. 9,  $A_{dow}$  refers to the expected availability for the day of week specified by the variable *dow*, which takes on seven different values ranging from 0 to 6. The term  $C_{i(dow)}$  is the actual number of patients that visited the ED on the particular day of week *dow* during the *i*th week in the past. Since national holidays affect the number of patients visiting the ED, weeks containing holidays are ignored completely in the availability calculations. Finally, the parameter *m* controls the smoothness of the sickness curve.

Since sickness is defined as the total number of people in the city getting sick on a particular day, we can calculate it as follows:

$$S_{today} = \frac{C_{today}}{A_{today}} \quad (10)$$

The term  $S_{today}$  in Eq. 10 refers to the number of people who are sick on the current day while the term  $C_{today}$  refers to the number of people who visited the ED on the current day.  $A_{today}$  is the probability that a patient will visit the ED for the day of week specified by *today*. Figure 14.18 shows the availability estimated in the period leading up to the synthetic ramp outbreak: it shows a consistent picture that we usually see far more patients on Mondays than Sundays, and then the visits taper off during the rest of the week. Figure 14.19 shows both the original count and the estimated sickness (count/availability). Sickness is much more stable than the original count because day of week effects have been greatly reduced. It is now possible to run a time series algorithm on the sickness values. In this case, we chose to use moving average with a window of seven days. The resulting alarms show something that was not achieved in any





**FIGURE 14.19** The black “count” line shows the raw data. The gray “sickness” line shows the sickness after the count has been divided by the corresponding availability value from Figure 14.18. Note how the sickness time series is now smoother than the original data and decorrelated with day of week. The alarm levels are derived from moving average applied to the sickness time series.

of the previous illustrations: a strong alarm resulting from the higher-than-expected counts for the Sunday.

We would like to emphasize that the sickness availability method only smoothes the time series. It is not a detection algorithm by itself. Instead, a detection algorithm, such as the control chart, moving average, or CUSUM algorithms should be used on the smoothed data.

## 11. FURTHER COMPARISON OF THE UNIVARIATE ALGORITHMS

We insert another copy of Table 14.1, with some of the above methods added for comparison. The newly evaluated algorithms follow:

- Regression using two features: the mean count over the past week and hours of daylight. This allows the algorithm to account for seasonal variation by putting a negative coefficient in front of hours of daylight.
- Regression using the additional feature *is\_Monday*, which is set to 1 if today is Monday and 0 otherwise. This allows the algorithm to anticipate the Monday bump in physician visits and so be less prone to false positives.
- Regression using indicator variables for all days of the week except for Sunday (which would be redundant), and additionally, *hours\_of\_daylight* and *mean\_count\_over\_previous\_seven\_days*.

- Using sickness/availability to compensate day-of-week effects, and then using the approach of comparing against yesterday. This method thus looks for jumps in the day-of-week-adjusted counts.

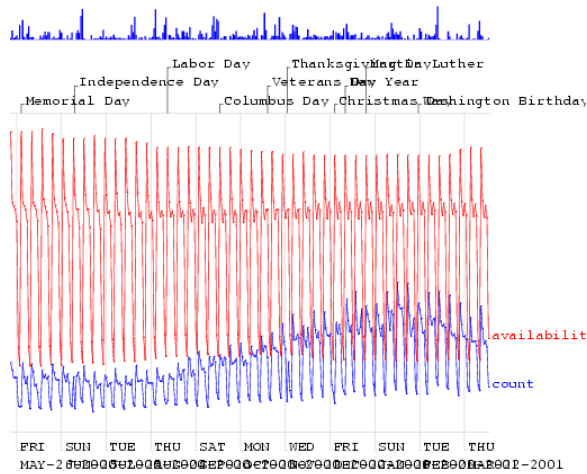
For this data set, with its seasonal components and day-of-week components, we see that sickness availability (to cope with day-of-week effects) combined with moving average (to cope with seasonal trends) performs well. Some of the regression methods perform almost equally well. We should note that this does not mean these methods are best in general: individual properties of individual data sets mean different approaches can be stronger for different data sets. Our only general advice is that in our experience relatively simple methods usually work at least as well as complex approaches. A second important note is that the numbers in Table 14.2 cannot be used as an estimate of how quickly real outbreaks are expected to be detected: the numbers are a function of many things, including the simulated magnitude of the outbreaks and the simulated noise levels.

## 12. ADDITIONAL METHODS

A complete set of approaches would require a very thick book, such as Hamilton (1994). The purpose of this chapter has been a tour of a sufficient variety of approaches to introduce the reader to some of the issues faced when choosing or implementing a time-series-based method, without going into the

# The Sickness/Availability Model

https://damianlab.com/=ID

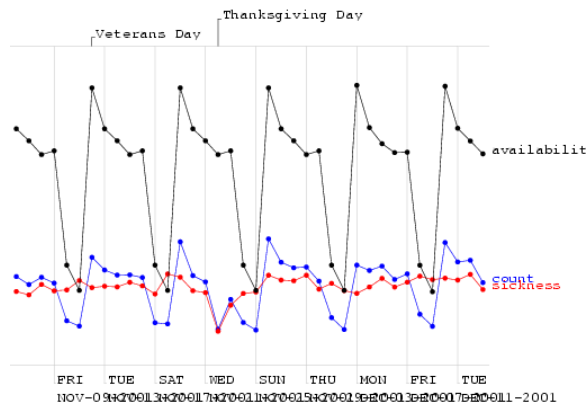


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 41

# The Sickness/Availability Model

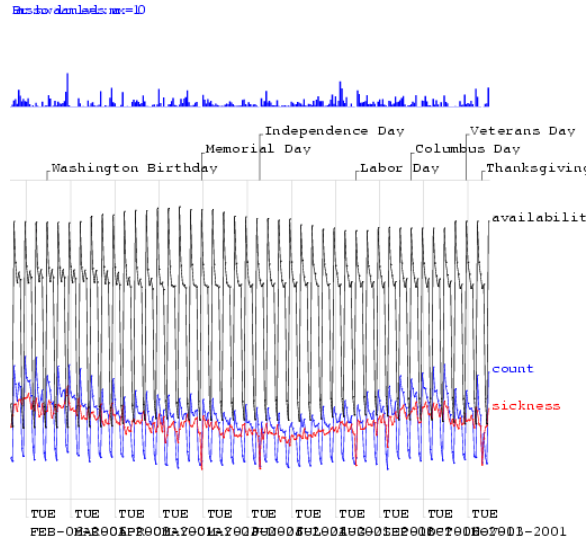
https://damianlab.com/=ID



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 42

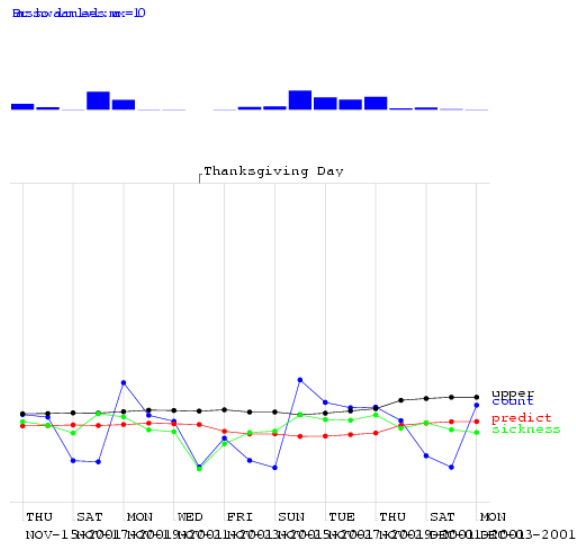
# The Sickness/Availability Model



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 43

# The Sickness/Availability Model

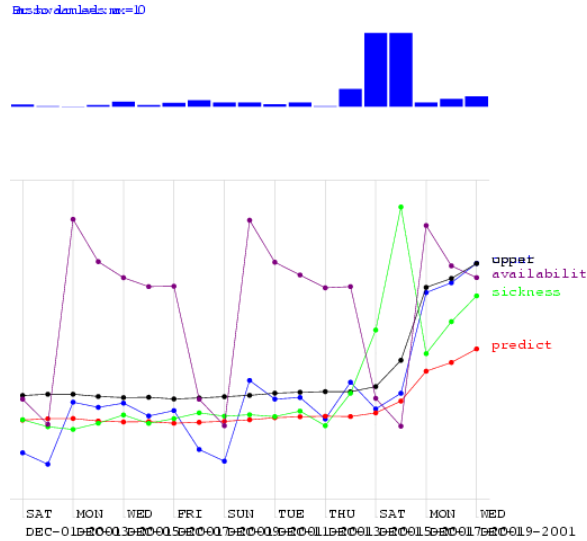


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 44



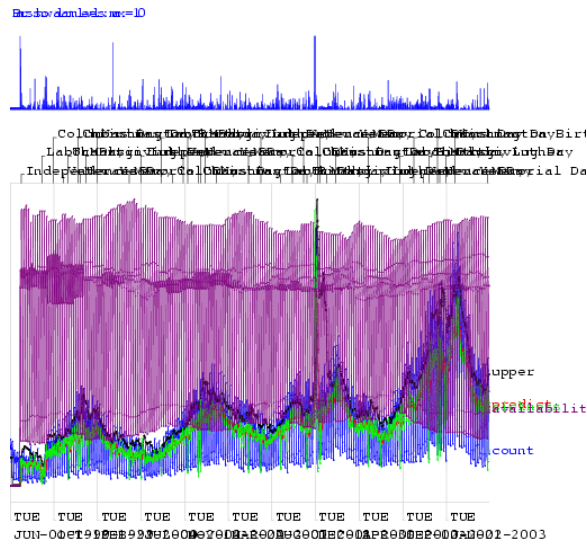
# The Sickness/Availability Model



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 45

# The Sickness/Availability Model



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 46

## Algorithm Performance

Allowing one False Alarm per TWO weeks

Allowing one False Alarm per SIX weeks...

Fraction of spikes detected  
Days to detect a ramp attack

|                                     | 0.39 | 3.47 | 0.22 | 4.13 |
|-------------------------------------|------|------|------|------|
| standard control chart              | 0.39 | 3.47 | 0.22 | 4.13 |
| using yesterday                     | 0.14 | 3.83 | 0.1  | 4.7  |
| Moving Average 3                    | 0.36 | 3.45 | 0.33 | 3.79 |
| Moving Average 7                    | 0.58 | 2.79 | 0.51 | 3.31 |
| Moving Average 56                   | 0.54 | 2.72 | 0.44 | 3.54 |
| hours_of_daylight                   | 0.58 | 2.73 | 0.43 | 3.9  |
| hours_of_daylight is_mon            | 0.7  | 2.25 | 0.57 | 3.12 |
| hours_of_daylight is_mon ... is_tue | 0.72 | 1.83 | 0.57 | 3.16 |
| hours_of_daylight is_mon ... is_sat | 0.77 | 2.11 | 0.59 | 3.26 |
| CUSUM                               | 0.45 | 2.03 | 0.15 | 3.55 |
| sa-mav-1                            | 0.86 | 1.88 | 0.74 | 2.73 |
| sa-mav-7                            | 0.87 | 1.28 | 0.83 | 1.87 |
| sa-mav-14                           | 0.86 | 1.27 | 0.82 | 1.62 |



## Algorithm Performance

Allowing one False Alarm per TWO weeks

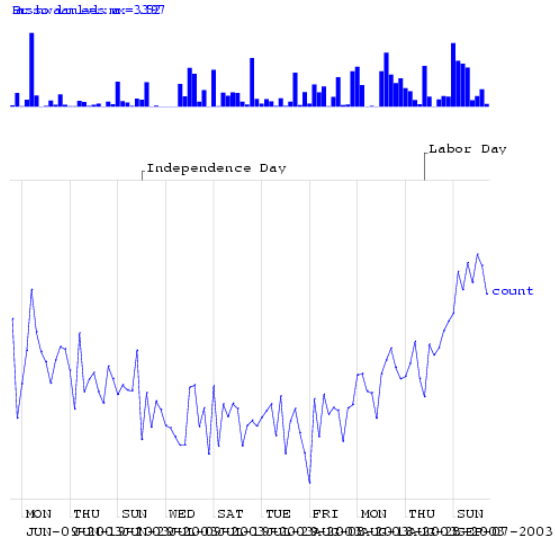
Allowing one False Alarm per SIX weeks...

Fraction of spikes detected  
Days to detect a ramp attack

|                                     | 0.39 | 3.47 | 0.22 | 4.13 |
|-------------------------------------|------|------|------|------|
| standard control chart              | 0.39 | 3.47 | 0.22 | 4.13 |
| using yesterday                     | 0.14 | 3.83 | 0.1  | 4.7  |
| Moving Average 3                    | 0.36 | 3.45 | 0.33 | 3.79 |
| Moving Average 7                    | 0.58 | 2.79 | 0.51 | 3.31 |
| Moving Average 56                   | 0.54 | 2.72 | 0.44 | 3.54 |
| hours_of_daylight                   | 0.58 | 2.73 | 0.43 | 3.9  |
| hours_of_daylight is_mon            | 0.7  | 2.25 | 0.57 | 3.12 |
| hours_of_daylight is_mon ... is_tue | 0.72 | 1.83 | 0.57 | 3.16 |
| hours_of_daylight is_mon ... is_sat | 0.77 | 2.11 | 0.59 | 3.26 |
| CUSUM                               | 0.45 | 2.03 | 0.15 | 3.55 |
| sa-mav-1                            | 0.86 | 1.88 | 0.74 | 2.73 |
| sa-mav-7                            | 0.87 | 1.28 | 0.83 | 1.87 |
| sa-mav-14                           | 0.86 | 1.27 | 0.82 | 1.62 |
| sa-regress                          | 0.73 | 1.76 | 0.67 | 2.21 |



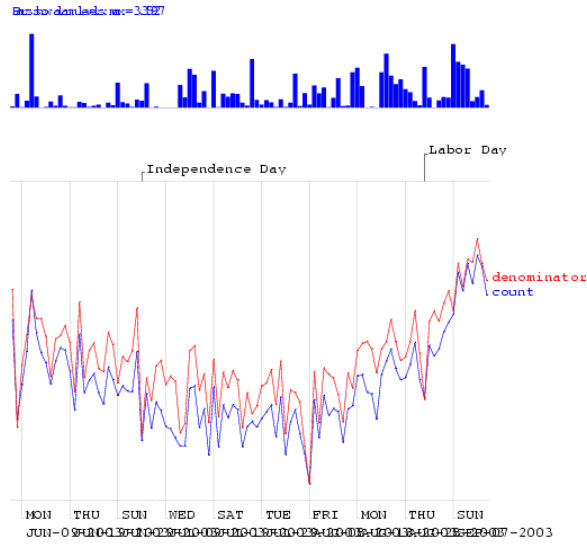
# Exploiting Denominator Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 49

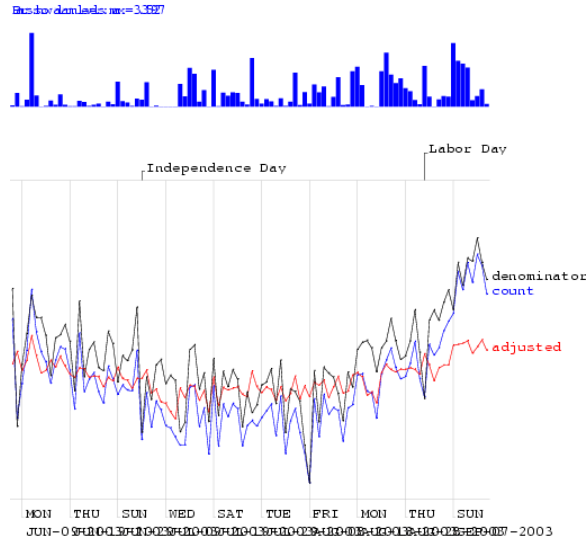
# Exploiting Denominator Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 50

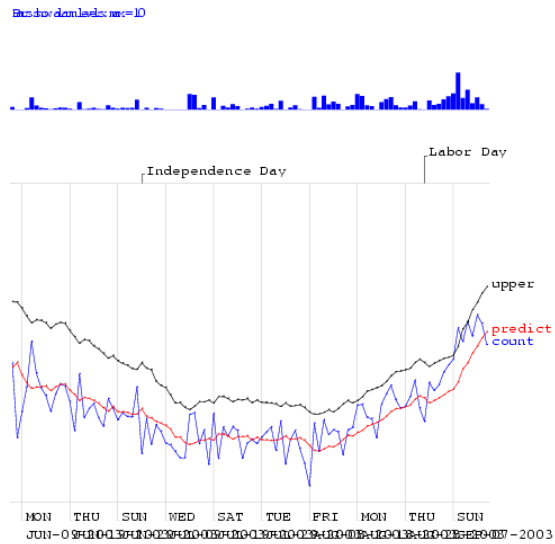
# Exploiting Denominator Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 51

# Exploiting Denominator Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 52

## Algorithm Performance

Allowing one False Alarm per TWO weeks

Allowing one False Alarm per SIX weeks...

|                                     | Fraction of spikes detected | Days to detect a ramp attack | Fraction of spikes detected | Days to detect a ramp attack |
|-------------------------------------|-----------------------------|------------------------------|-----------------------------|------------------------------|
| standard control chart              | 0.39                        | 3.47                         | 0.22                        | 4.13                         |
| using yesterday                     | 0.14                        | 3.83                         | 0.1                         | 4.7                          |
| Moving Average 3                    | 0.36                        | 3.45                         | 0.33                        | 3.79                         |
| Moving Average 7                    | 0.58                        | 2.79                         | 0.51                        | 3.31                         |
| Moving Average 56                   | 0.54                        | 2.72                         | 0.44                        | 3.54                         |
| hours_of_daylight                   | 0.58                        | 2.73                         | 0.43                        | 3.9                          |
| hours_of_daylight is_mon            | 0.7                         | 2.25                         | 0.57                        | 3.12                         |
| hours_of_daylight is_mon ... is_tue | 0.72                        | 1.83                         | 0.57                        | 3.16                         |
| hours_of_daylight is_mon ... is_sat | 0.77                        | 2.11                         | 0.59                        | 3.26                         |
| CUSUM                               | 0.45                        | 2.03                         | 0.15                        | 3.55                         |
| sa-mav-1                            | 0.86                        | 1.88                         | 0.74                        | 2.73                         |
| sa-mav-7                            | 0.87                        | 1.28                         | 0.83                        | 1.87                         |
| sa-mav-14                           | 0.86                        | 1.27                         | 0.82                        | 1.62                         |
| sa-regress                          | 0.73                        | 1.76                         | 0.67                        | 2.21                         |
| Cough with denominator              | 0.78                        | 2.15                         | 0.59                        | 2.41                         |
| Cough with MA                       | 0.65                        | 2.78                         | 0.57                        | 3.24                         |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 53

## Other state-of-the-art methods

- Wavelets
- Change-point detection
- Kalman filters
- Hidden Markov Models

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 54

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

WSARE

Spatial Scan Statistics

Univariate Anomaly Detection

Multivariate Anomaly Detection

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 55

## Multiple Signals

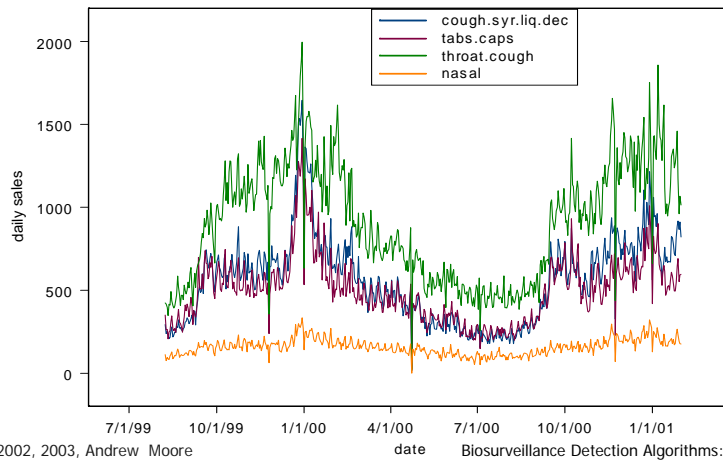


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 56

# Multivariate Signals

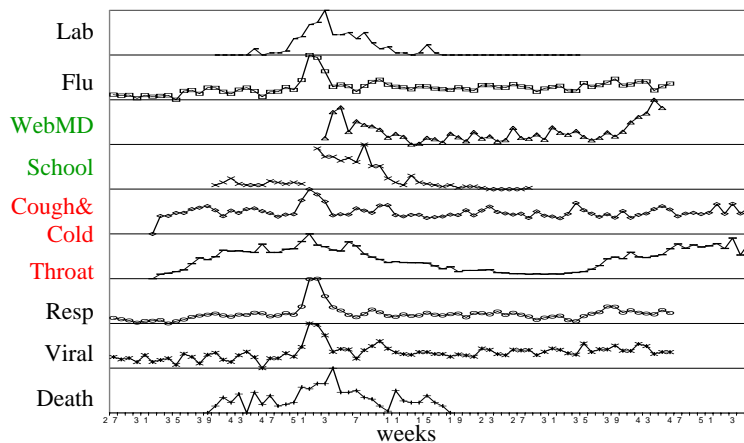
(relevant to inhalational diseases)



Copyright © 2002, 2003, Andrew Moore Biosurveillance Detection Algorithms: Slide 57

# Multi Source Signals

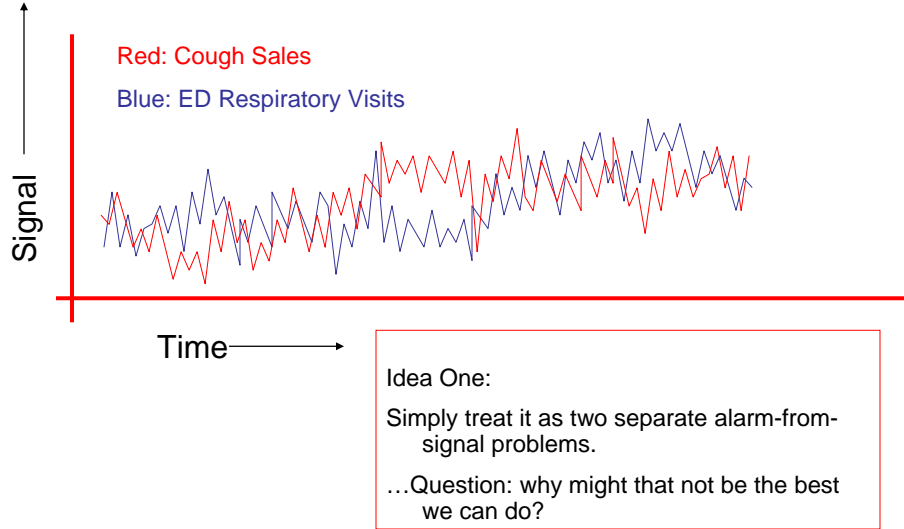
Footprint of Influenza in Routinely Collected Data



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 58

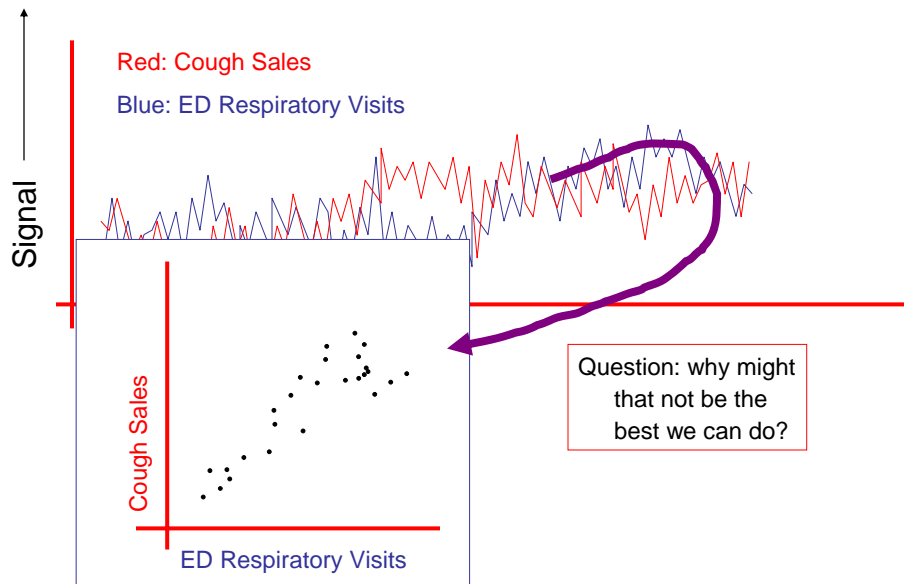
## What if you've got multiple signals?



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 59

## Another View

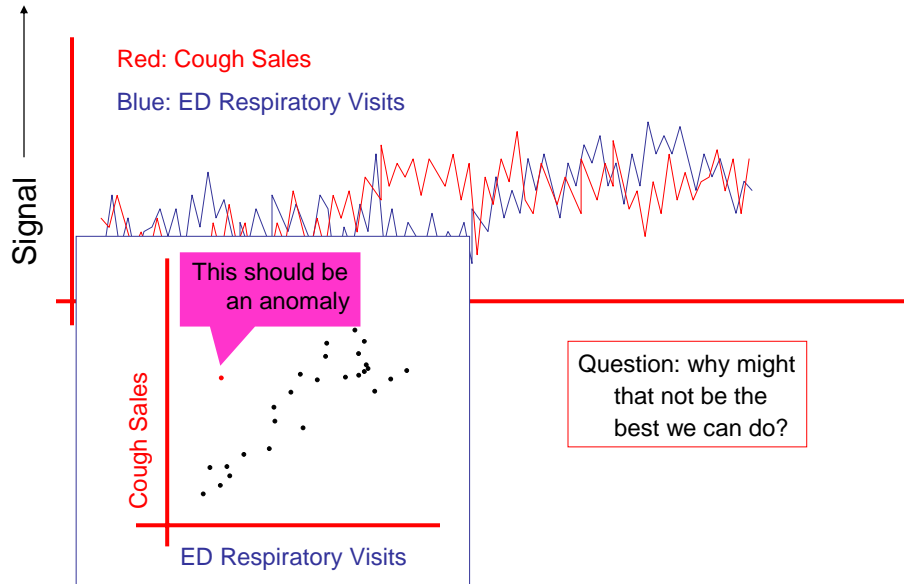


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 60



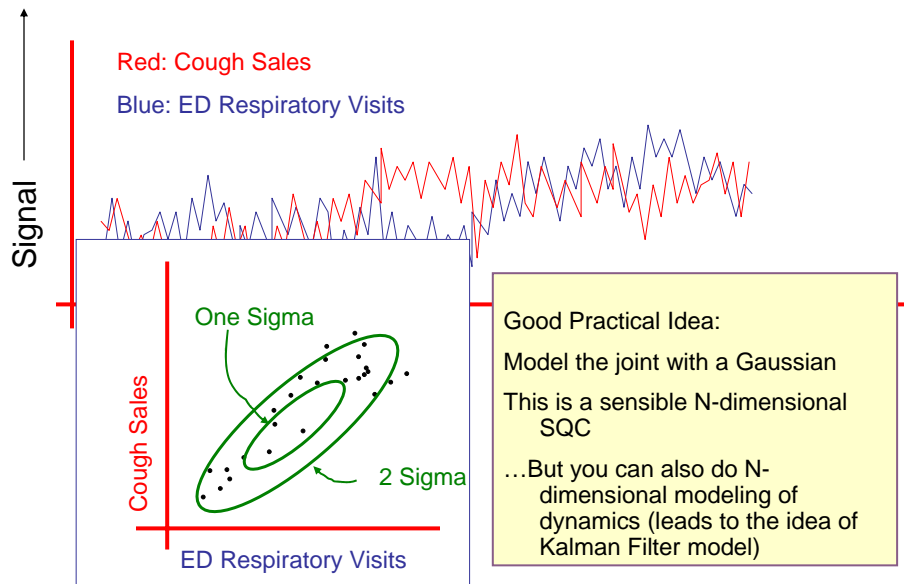
## Another View



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 61

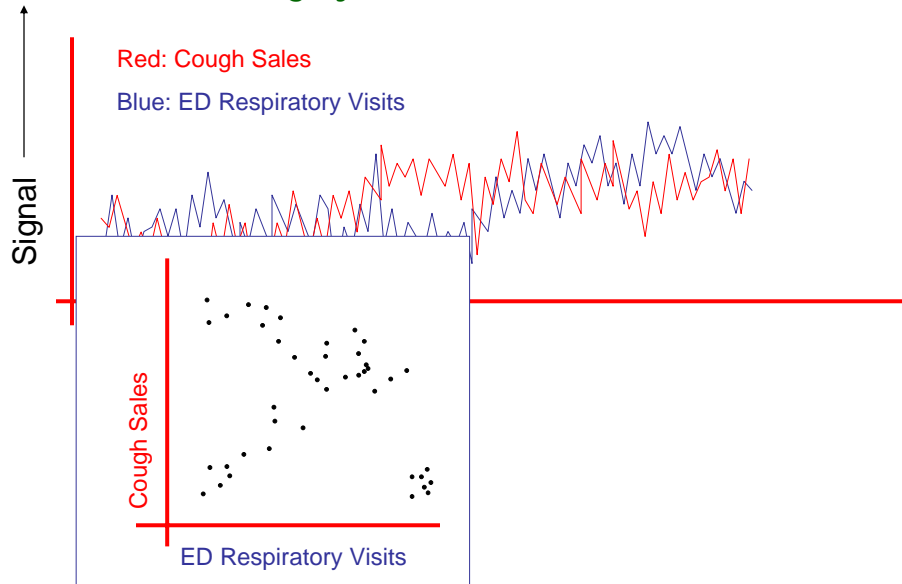
## N-dimensional Gaussian



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 62

## But what if joint N-dimensional distribution is highly non-Gaussian?



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 63

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

WSARE

Spatial Scan Statistics

Multivariate Anomaly Detection

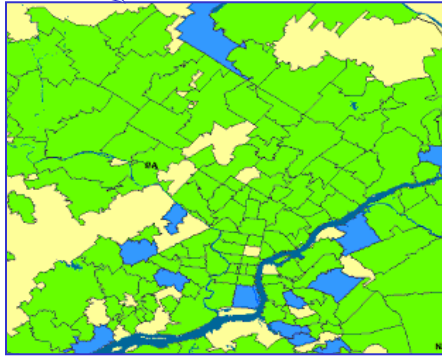
Univariate Anomaly Detection

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 64

# One Step of Spatial Scan

Entire area being scanned

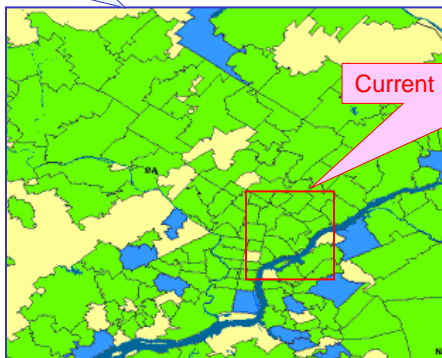


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 65

# One Step of Spatial Scan

Entire area being scanned



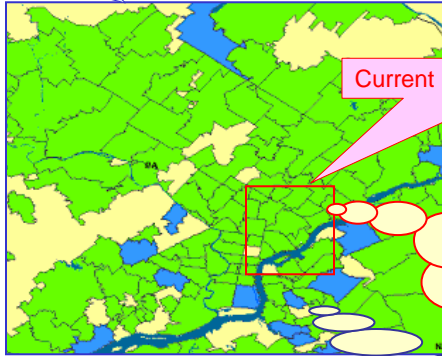
Current region being considered

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 66

# One Step of Spatial Scan

Entire area being scanned



Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

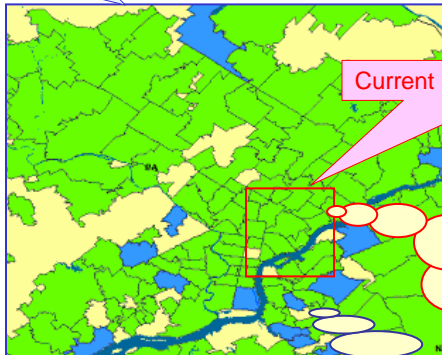
Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

Copyright © 2002, 2003, Andrew Moore

Surveillance Detection Algorithms: Slide 67

# One Step of Spatial Scan

Entire area being scanned



Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

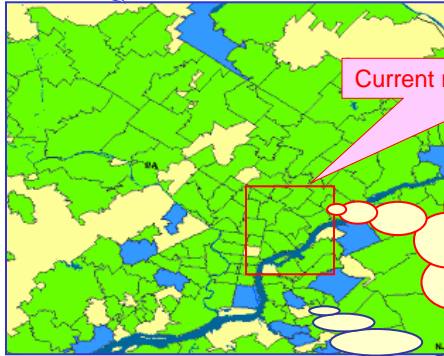
**So...** *is that a big deal?*  
Evaluated with Score function (e.g. Kulldorf's score)

Copyright © 2002, 2003, Andrew Moore

Surveillance Detection Algorithms: Slide 68

# One Step of Spatial Scan

Entire area being scanned



Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

[Score = 1.4]

Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

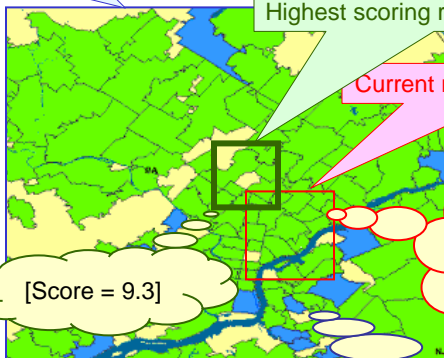
**So...** is that a big deal?  
Evaluated with Score function (e.g. Kulldorf's score)

Copyright © 2002, 2003, Andrew Moore

Surveillance Detection Algorithms: Slide 69

# Many Steps of Spatial Scan

Entire area being scanned



Highest scoring region in search so far

Current region being considered

I have a population of 5300 of whom 53 are sick (1%)

[Score = 1.4]

[Score = 9.3]

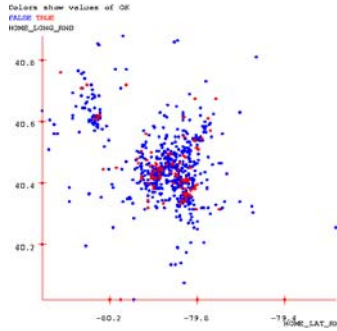
Everywhere else has a population of 2,200,000 of whom 20,000 are sick (0.9%)

**So...** is that a big deal?  
Evaluated with Score function (e.g. Kulldorf's score)

Copyright © 2002, 2003, Andrew Moore

Surveillance Detection Algorithms: Slide 70

# Scan Statistics



Standard scan statistic question:  
 Given the geographical locations of occurrences of a phenomenon, is there a region with an unusually high (low) rate of these occurrences?

## Standard approach:

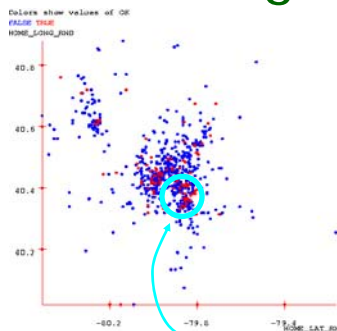
1. Compute the likelihood of the data given the hypothesis that the rate of occurrence is uniform everywhere,  $L_0$
2. For some geographical region,  $W$ , compute the likelihood that the rate of occurrence is uniform at one level inside the region and uniform at another level outside the region,  $L(W)$ .
3. Compute the likelihood ratio,  $L(W)/L_0$
4. Repeat for all regions, and find the largest likelihood ratio. This is the scan statistic,  $S^*_W$
5. Report the region,  $W$ , which yielded the max,  $S^*_W$

See [Glaz and Balakrishnan, 99] for details

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 71

# Significance testing



Given that region  $W$  is the most likely to be abnormal, is it significantly abnormal?

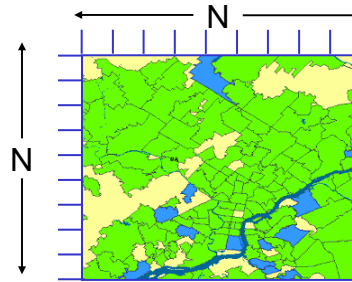
## Standard approach:

1. Generate many randomized versions of the data set by shuffling the labels (positive instance of the phenomenon or not).
2. Compute  $S^*_W$  for each randomized data set. This forms a baseline distribution for  $S^*_W$  if the null hypothesis holds.
3. Compare the observed value of  $S^*_W$  against the baseline distribution to determine a p-value.

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 72

## Fast squares speedup

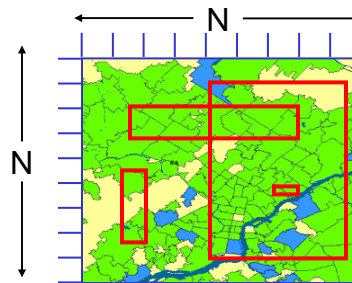


- Theoretical complexity of fast squares:  $O(N^2)$  (as opposed to naive  $N^3$ ), if maximum density region sufficiently dense.  
*If not, we can use several other speedup tricks.*
- In practice: 10-200x speedups on real and artificially generated datasets.  
*Emergency Dept. dataset (600K records): 20 minutes, versus 66 hours with naive approach.*

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 73

## Fast rectangles speedup



Work in progress

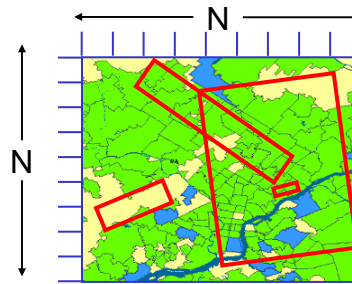
- Theoretical complexity of fast rectangles:  $O(N^2 \log N)$  (as opposed to naive  $N^4$ )

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 74

## Fast oriented rectangles speedup

Work in progress



- Theoretical complexity of fast rectangles:  $18N^2 \log N$  (as opposed to naive  $18N^4$ )

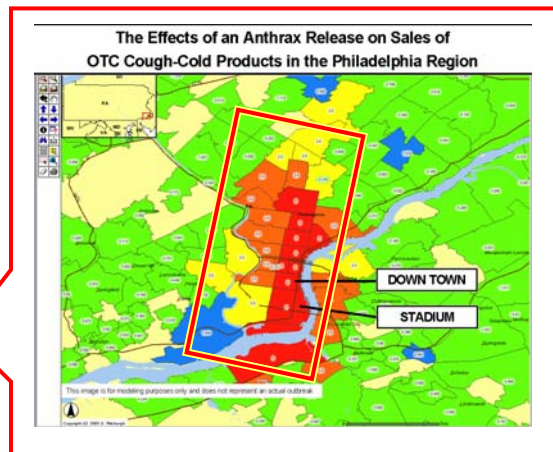
*(Angles discretized to 5 degree buckets)*

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 75

## Why the Scan Statistic speed obsession?

- Traditional Scan Statistics very expensive, especially with Randomization tests
- New "Historical Model" Scan Statistics
- Proposed new WSARE/Scan Statistic hybrid



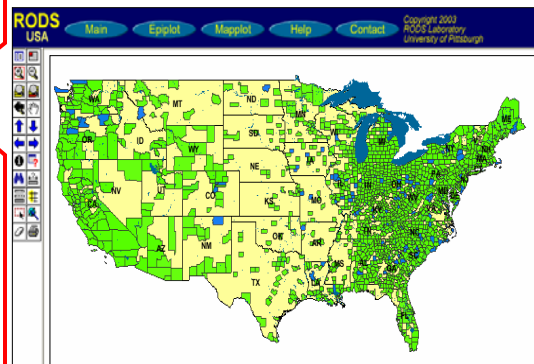
Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 76



## Why the Scan Statistic speed obsession?

- Traditional Scan Statistics very expensive, especially with Randomization tests
- New "Historical Model" Scan Statistics
- Proposed new WSARE/Scan Statistic hybrid

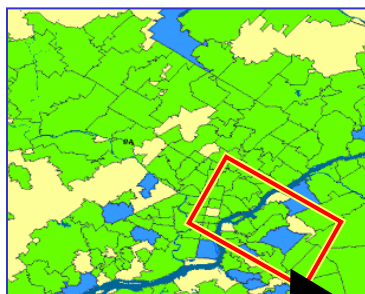


Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 77

## Why the Scan Statistic speed obsession?

- Traditional Scan Statistics very expensive, especially with Randomization tests
- New "Historical Model" Scan Statistics
- Proposed new WSARE/Scan Statistic hybrid



This is the strangest region because the age distribution of respiratory cases has changed dramatically for no reason that can be explained by known background changes

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 78

## What you'll learn about

- Noticing events in bio-event time series
- Tracking many series at once
- Detecting geographic hotspots
- Finding emerging new patterns

WSARE

Spatial Scan Statistics

Multivariate Anomaly Detection

Univariate Anomaly Detection

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 79

## But there's potentially more data than aggregates

Suppose we know that today in the ED we had...

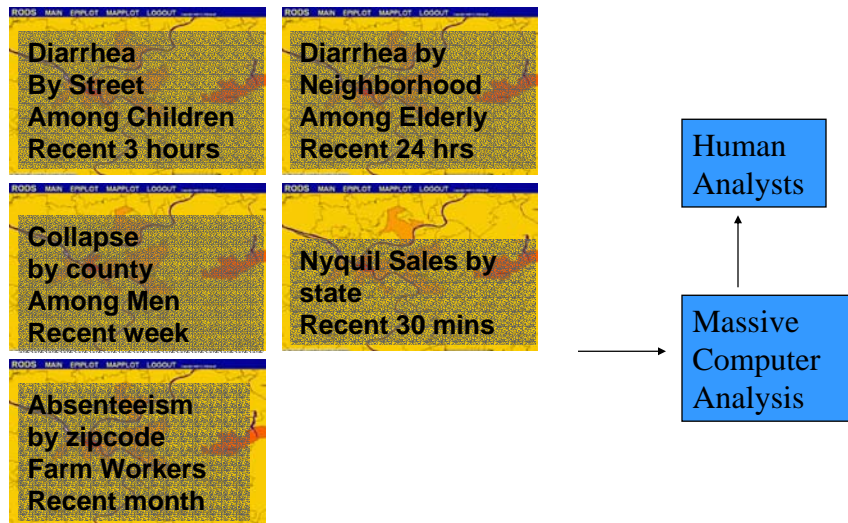
- 421 Cases
  - 78 Respiratory Cases
  - 190 Males
  - 32 Children
  - 21 from North Suburbs
  - 2 Postal workers
- (etc etc etc)

Have we made best use of all possible information?

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 80

## There are so many things to look at



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 81

## WSARE v2.0

- What's Strange About Recent Events?
- Designed to be easily applicable to any date/time-indexed biosurveillance-relevant data stream.

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 82

# WSARE v2.0

- Inputs:
  - 1. Date/time-indexed biosurveillance-relevant data stream
  - 2. Time Window Length
  - 3. Which attributes to use?

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 83

# WSARE v2.0

- Input s:
    - 1. Date/time-indexed biosurveillance-relevant data stream
    - 2. Time Window Length
    - 3. Which attributes to use?
- Example
- "last 24 hours"
- "ignore key and weather"

| Primary Key | Date  | Time  | Hospital   | ICD9 | Prodrome    | Gender | Age | Home        |              |            | Work        |              |            | Recent Flu Levels | Recent Weather | (Many more...) |
|-------------|-------|-------|------------|------|-------------|--------|-----|-------------|--------------|------------|-------------|--------------|------------|-------------------|----------------|----------------|
|             |       |       |            |      |             |        |     | Large Scale | Medium Scale | Fine Scale | Large Scale | Medium Scale | Fine Scale |                   |                |                |
| h6r32       | 6/2/2 | 14:12 | Downtown   | 781  | Fever       | M      | 20s | NE          | 15217        | A5         | NW          | 15213        | B8         | 2%                | 70R            | ...            |
| t3q15       | 6/2/2 | 14:15 | Riverside  | 717  | Respiratory | M      | 60s | NE          | 15222        | J3         | NE          | 15222        | J3         | 2%                | 70R            | ...            |
| t5hh5       | 6/2/2 | 14:15 | Smithfield | 622  | Respiratory | F      | 80s | SE          | 15210        | K9         | SE          | 15210        | K9         | 2%                | 70R            | ...            |
| :           | :     | :     | :          | :    | :           | :      | :   | :           | :            | :          | :           | :            | :          | :                 | :              | :              |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 84

## WSARE v2.0

- Inputs:
  - 1. Date/time-indexed biosurveillance-relevant data stream
  - 2. Time Window Length
  - 3. Which attributes to use?
- Outputs:
  - 1. Here are the records that most surprise me
  - 2. Here's why
  - 3. And here's how seriously you should take it

| Primary Key | Date  | Time  | Hospital    | ICD9 | Prodrome    | Gender | Age | Home        |              |            | Work        |              |            | Recent Flu Levels | Recent Weather | (Many more...) |
|-------------|-------|-------|-------------|------|-------------|--------|-----|-------------|--------------|------------|-------------|--------------|------------|-------------------|----------------|----------------|
|             |       |       |             |      |             |        |     | Large Scale | Medium Scale | Fine Scale | Large Scale | Medium Scale | Fine Scale |                   |                |                |
| h6r32       | 6/2/2 | 14:12 | Down-town   | 781  | Fever       | M      | 20s | NE          | 15217        | A5         | NW          | 15213        | B8         | 2%                | 70R            | ...            |
| t3q15       | 6/2/2 | 14:15 | River-side  | 717  | Respiratory | M      | 60s | NE          | 15222        | J3         | NE          | 15222        | J3         | 2%                | 70R            | ...            |
| t5hh5       | 6/2/2 | 14:15 | Smith-field | 622  | Respiratory | F      | 80s | SE          | 15210        | K9         | SE          | 15210        | K9         | 2%                | 70R            | ...            |
| :           | :     | :     | :           | :    | :           | :      | :   | :           | :            | :          | :           | :            | :          | :                 | :              | :              |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 85

## Simple WSARE

- Given 500 day's worth of ER cases at 15 hospitals...

| Date           | Cases              |
|----------------|--------------------|
| Thu 5/22/2000  | C1, C2, C3, C4 ... |
| Fri 5/23/2000  | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Sat 12/9/2000  | C1, C2, C3, C4 ... |
| Sun 12/10/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/16/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/23/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Fri 9/14/2001  | C1, C2, C3, C4 ... |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 86

## Simple WSARE

- Given 500 day's worth of ER cases at 15 hospitals...
- For each day...
  - Take today's cases

| Date           | Cases              |
|----------------|--------------------|
| Thu 5/22/2000  | C1, C2, C3, C4 ... |
| Fri 5/23/2000  | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Sat 12/9/2000  | C1, C2, C3, C4 ... |
| Sun 12/10/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/16/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/23/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Fri 9/14/2001  | C1, C2, C3, C4 ... |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 87

## Simple WSARE

- Given 500 day's worth of ER cases at 15 hospitals...
- For each day...
  - Take today's cases
  - The cases one week ago
  - The cases two weeks ago

| Date           | Cases              |
|----------------|--------------------|
| Thu 5/22/2000  | C1, C2, C3, C4 ... |
| Fri 5/23/2000  | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Sat 12/9/2000  | C1, C2, C3, C4 ... |
| Sun 12/10/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/16/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| Sat 12/23/2000 | C1, C2, C3, C4 ... |
| :              | :                  |
| :              | :                  |
| Fri 9/14/2001  | C1, C2, C3, C4 ... |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 88

# Simple WSARE

- Given 500 day's worth of ER cases at 15 hospitals...
- For each day...
  - Take today's cases
  - The cases one week ago
  - The cases two weeks ago
- Ask: "What's different about today?"

| DATE_AD  | ICD9   | PRODROM | GENDER | place2 | ... | ... |
|----------|--------|---------|--------|--------|-----|-----|
| 12/9/00  | 786.05 |         | 3 F    | s-e    | ... | ... |
| 12/9/00  | 789    |         | 1 F    | s-e    | ... | ... |
| 12/9/00  | 789    |         | 1 M    | n-w    | ... | ... |
| 12/9/00  | 786.05 |         | 3 M    | s-e    | ... | ... |
| :        | :      | :       | :      | :      | ... | ... |
| 12/16/00 | 787.02 |         | 2 M    | n-e    | ... | ... |
| 12/16/00 | 782.1  |         | 4 F    | s-w    | ... | ... |
| 12/16/00 | 789    |         | 1 M    | s-e    | ... | ... |
| 12/16/00 | 786.09 |         | 3 M    | n-w    | ... | ... |
| 12/23/00 | 789.09 |         | 1 M    | s-w    | ... | ... |
| 12/23/00 | 789.09 |         | 1 F    | s-w    | ... | ... |
| 12/23/00 | 782.1  |         | 4 M    | n-w    | ... | ... |
| :        | :      | :       | :      | :      | ... | ... |
| 12/23/00 | 786.09 |         | 3 M    | s-e    | ... | ... |
| 12/23/00 | 786.09 |         | 3 M    | s-e    | ... | ... |
| 12/23/00 | 780.9  |         | 2 F    | n-w    | ... | ... |
| 12/23/00 | V40.9  |         | 7 M    | s-w    | ... | ... |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 89

# Simple WSARE

- Given 500 day's worth of ER cases at 15 hospitals...
- For each day...

| DATE_AD  | ICD9   | PRODROM | GENDER | place2 | ... | ... |
|----------|--------|---------|--------|--------|-----|-----|
| 12/9/00  | 786.05 |         | 3 F    | s-e    | ... | ... |
| 12/9/00  | 789    |         | 1 F    | s-e    | ... | ... |
| 12/9/00  | 789    |         | 1 M    | n-w    | ... | ... |
| 12/9/00  | 786.05 |         | 3 M    | s-e    | ... | ... |
| :        | :      | :       | :      | :      | ... | ... |
| 12/16/00 | 787.02 |         | 2 M    | n-e    | ... | ... |
| 12/16/00 | 782.1  |         | 4 F    | s-w    | ... | ... |
| 12/16/00 | 789    |         | 1 M    | s-e    | ... | ... |
| 12/16/00 | 786.09 |         | 3 M    | n-w    | ... | ... |
| 12/23/00 | 789.09 |         | 1 M    | s-w    | ... | ... |
| 12/23/00 | 789.09 |         | 1 F    | s-w    | ... | ... |
| :        | :      | :       | :      | :      | ... | ... |
| 12/23/00 | 786.09 |         | 3 M    | s-e    | ... | ... |
| 12/23/00 | 786.09 |         | 3 M    | s-e    | ... | ... |
| 12/23/00 | 780.9  |         | 2 F    | n-w    | ... | ... |
| 12/23/00 | V40.9  |         | 7 M    | s-w    | ... | ... |

Fields we use:  
 Date, Time of Day, Prodrome, ICD9,  
**Symptoms**, Age, Gender, Coarse Location,  
 Fine Location, **ICD9 Derived Features**,  
**Census Block Derived Features**, **Work  
 Details**, **Colocation Details**

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 90

# Example

Sat 12-23-2001 (daynum 36882, dayindex 239)

35.8% ( 48/134) of today's cases have  $30 \leq \text{age} < 40$

17.0% ( 45/265) of other cases have  $30 \leq \text{age} < 40$

# Example

Sat 12-23-2001 (daynum 36882, dayindex 239)

FISHER\_PVALUE = 0.000051

35.8% ( 48/134) of today's cases have  $30 \leq \text{age} < 40$

17.0% ( 45/265) of other cases have  $30 \leq \text{age} < 40$

Table 1: A sample 2x2 Contingency Table

|                      | $C_{today}$ | $C_{other}$ |
|----------------------|-------------|-------------|
| $Age\_Decile = 3$    | 48          | 45          |
| $Age\_Decile \neq 3$ | 86          | 220         |



## Searching for the best score...

- Try ICD9 = x for each value of x
- Try Gender=M, Gender=F
- Try CoarseRegion=NE, =NW, SE, SW..
- Try FineRegion=AA,AB,AC, ... DD (4x4 Grid)
- Try Hospital=x, TimeofDay=x, Prodrome=X, ...
- [In future... features of census

**Overfitting Alert!**

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 93

## Example

```
Sat 12-23-2001 (daynum 36882, dayindex 239)
FISHER_PVALUE = 0.000051 RANDOMIZATION_PVALUE = 0.031
35.8% ( 48/134) of today's cases have 30 <= age < 40
17.0% ( 45/265) of other cases have 30 <= age < 40
```

Table 1: A sample 2x2 Contingency Table

|                      | $C_{today}$ | $C_{other}$ |
|----------------------|-------------|-------------|
| $Age\_Decile = 3$    | 48          | 45          |
| $Age\_Decile \neq 3$ | 86          | 220         |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 94

## Multiple component rules

- We would like to be able to find rules like:
  - There are a surprisingly large number of children with respiratory problems today
- or
- There are too many skin complaints among people from the affluent neighborhoods
- These are things that would be missed by casual screening
- **BUT**
  - The danger of overfitting could be much worse
  - It's very computationally demanding
  - How can we be sure the entire rule is meaningful?

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 95

## Checking two component rules

Table 2: 2x2 Contingency Table 1 for a two component rule

|  |  |
|--|--|
| Records from Today matching $C_0$ and $C_1$              | Records from Other matching $C_0$ and $C_1$              |
| Records from Today matching $C_1$ and differing on $C_0$ | Records from Other matching $C_1$ and differing on $C_0$ |

Table 3: 2x2 Contingency Table 2 for a two component rule

|  |  |
|--|--|
| Records from Today matching $C_0$ and $C_1$              | Records from Other matching $C_0$ and $C_1$              |
| Records from Today matching $C_0$ and differing on $C_1$ | Records from Other matching $C_0$ and differing on $C_1$ |

- Must pass both tests to be allowed to live.

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 96

# WSARE v2.0

- Inputs:
  - 1. Date/time-indexed biosurveillance-relevant data stream
  - 2. Time Window Length
  - 3. Which attributes to use?
- Outputs:
  - 1. Here are the records that most surprise me
  - 2. Here's why
  - 3. And here's how seriously you should take it

| Primary Key | Date  | Time  | Hospital    | ICD9 | Prodrome    | Gender | Age | Home        |              |            | Work        |              |            | Recent Flu Levels | Recent Weather | (Many more...) |
|-------------|-------|-------|-------------|------|-------------|--------|-----|-------------|--------------|------------|-------------|--------------|------------|-------------------|----------------|----------------|
|             |       |       |             |      |             |        |     | Large Scale | Medium Scale | Fine Scale | Large Scale | Medium Scale | Fine Scale |                   |                |                |
| h6r32       | 6/2/2 | 14:12 | Down-town   | 781  | Fever       | M      | 20s | NE          | 15217        | A5         | NW          | 15213        | B8         | 2%                | 70R            | ...            |
| t3q15       | 6/2/2 | 14:15 | River-side  | 717  | Respiratory | M      | 60s | NE          | 15222        | J3         | NE          | 15222        | J3         | 2%                | 70R            | ...            |
| t5hh5       | 6/2/2 | 14:15 | Smith-field | 622  | Respiratory | F      | 80s | SE          | 15210        | K9         | SE          | 15210        | K9         | 2%                | 70R            | ...            |
| :           | :     | :     | :           | :    | :           | :      | :   | :           | :            | :          | :           | :            | :          | :                 | :              | :              |

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 97

# WSARE v2.0

- Input s:
  - 1. Date/time-indexed biosurveillance-relevant data stream
  - 2. Time Window Length
  - 3. Which attributes to use?
- Output s:
  - 1. Here are the records that most surprise me
  - 2. Here's why
  - 3. And here's how seriously you should take it

| Primary Key | Date  | Time  | Hospital    | ICD9 | Prodrome    | Gender | Age | Home        |              |            | Work        |              |            | Recent Flu Levels | Recent Weather | (Many more...) |
|-------------|-------|-------|-------------|------|-------------|--------|-----|-------------|--------------|------------|-------------|--------------|------------|-------------------|----------------|----------------|
|             |       |       |             |      |             |        |     | Large Scale | Medium Scale | Fine Scale | Large Scale | Medium Scale | Fine Scale |                   |                |                |
| h6r32       |       |       |             |      |             |        |     | NE          | 15217        | A5         | NW          | 15213        | B8         | 2%                | 70R            | ...            |
| t3q15       |       |       | side        |      | Respiratory | M      | 60s | NE          | 15222        | J3         | NE          | 15222        | J3         | 2%                | 70R            | ...            |
| t5hh5       | 6/2/2 | 14:15 | Smith-field | 622  | Respiratory | F      | 80s | SE          | 15210        | K9         | SE          | 15210        | K9         | 2%                | 70R            | ...            |
| :           | :     | :     | :           | :    | :           | :      | :   | :           | :            | :          | :           | :            | :          | :                 | :              | :              |

Normally, 8% of cases in the East are over-50s with respiratory problems.  
But today it's been 15%

Don't be too impressed!  
Taking into account all the patterns I've been searching over, there's a 20% chance I'd have found a rule this dramatic just by chance

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 98

# WSARE on recent Utah Data

Saturday June 1st in Utah:

The most surprising thing about recent records is:

Normally:

0.8% of records (50/6205) have time before 2pm and prodrome = Hemorrhagic

But recently:

2.1% of records (19/907) have time before 2pm and prodrome = Hemorrhagic

Pvalue = 0.0484042

Which means that in a world where nothing changes we'd

expect to have a result this significant about once

every 20 times we ran the program

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 99

# Results on Emergency Dept Data

```
### Rule 1: Tue 05-16-2000 (daynum 36661, dayindex 18)
SCORE = -0.00000000 PVALUE = 0.00000000
32.84% ( 44/134) of today's cases have Time Of Day4 after 6:00 pm
90.00% ( 27/30) of other cases have Time Of Day4 after 6:00 pm
```

```
### Rule 2: Fri 06-30-2000 (daynum 36706, dayindex 63)
SCORE = -0.00000000 PVALUE = 0.00000000
19.40% ( 26/134) of today's cases have Place2 = NE and Lat4 = d
5.71% ( 16/280) of other cases have Place2 = NE and Lat4 = d
```

```
### Rule 3: Wed 09-06-2000 (daynum 36774, dayindex 131)
SCORE = -0.00000000 PVALUE = 0.00000000
17.16% ( 23/134) of today's cases have Prodrome = Respiratory
and age2 less than 40
4.53% ( 12/265) of other cases have Prodrome = Respiratory
and age2 less than 40
```

```
### Rule 4: Fri 12-01-2000 (daynum 36860, dayindex 217)
SCORE = -0.00000000 PVALUE = 0.00000000
22.88% ( 27/118) of today's cases have Time Of Day4
after 6:00 pm and Lat2 = s
8.10% ( 20/247) of other cases have Time Of Day4
after 6:00 pm and Lat2 = s
```

```
### Rule 5: Sat 12-23-2000 (daynum 36882, dayindex 239)
SCORE = -0.00000000 PVALUE = 0.00000000
18.25% ( 25/137) of today's cases have ICD9 = shortness of breath
and Time Of Day2 before 3:00 pm
5.12% ( 15/293) of other cases have ICD9 = shortness of breath
and Time Of Day2 before 3:00 pm
```

```
### Rule 6: Fri 09-14-2001 (daynum 37147, dayindex 504)
SCORE = -0.00000000 PVALUE = 0.00000000
66.67% ( 30/45) of today's cases have Time Of Day4 before 10:00 am
18.42% ( 42/228) of other cases have Time Of Day4 before 10:00 am
```

Copyright © 2002, 2003, Andrew Moore

## WSARE 3.0

- “Taking into account recent flu levels...”
- “Taking into account that today is a public holiday...”
- “Taking into account that this is Spring...”
- “Taking into account recent heatwave...”
- “Taking into account that there’s a known natural Food-borne outbreak in progress...”

Bonus: More  
efficient use of  
historical data

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 101

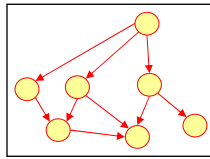
## Analysis of variance

- **Good news:**  
If you’re tracking a daily aggregate (e.g. number of flu cases in your ED, or Nyquil Sales)...then ANOVA can take care of many of these effects.
- **But...**  
What if you’re tracking a whole joint distribution of transactional events?

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 102

## Idea: Bayesian Networks



“Patients from West Park Hospital are less likely to be young”

“On Cold Tuesday Mornings the folks coming in from the North part of the city are more likely to have respiratory problems”

“The Viral prodrome is more likely to co-occur with a Rash prodrome than Botulinic”

“On the day after a major holiday, expect a boost in the morning followed by a lull in the afternoon”

Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 103

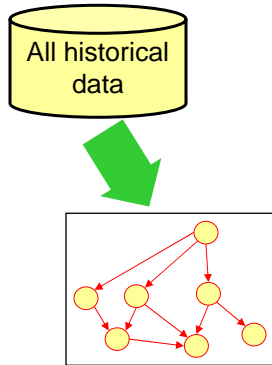
## WSARE 3.0



Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 104

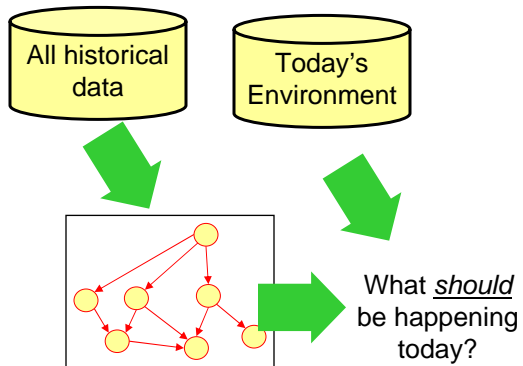
# WSARE 3.0



Copyright © 2002, 2003, Andrew Moore

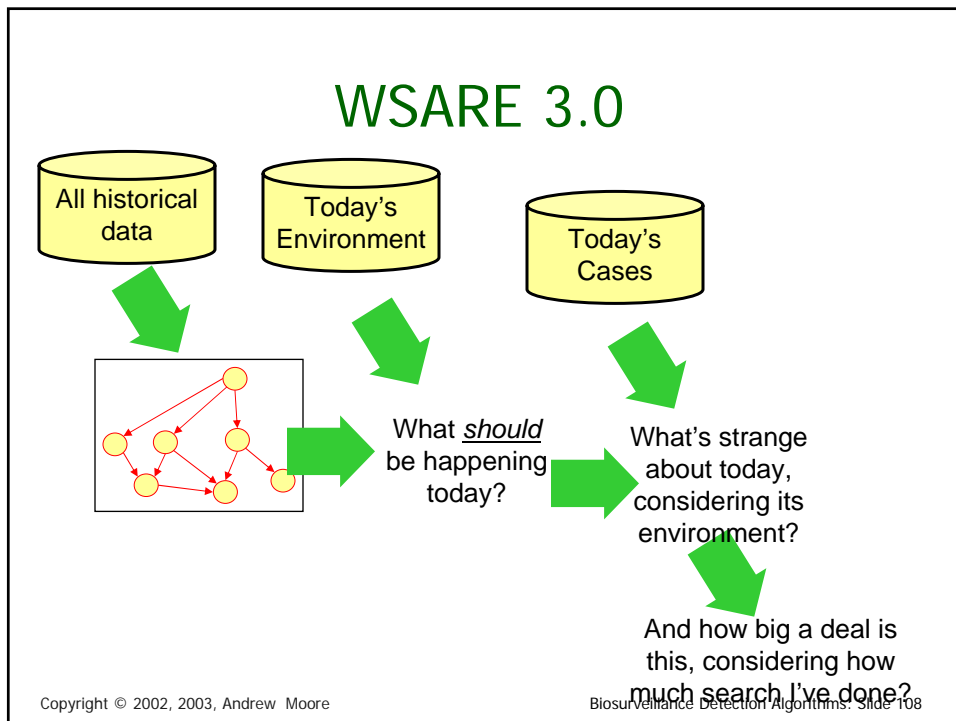
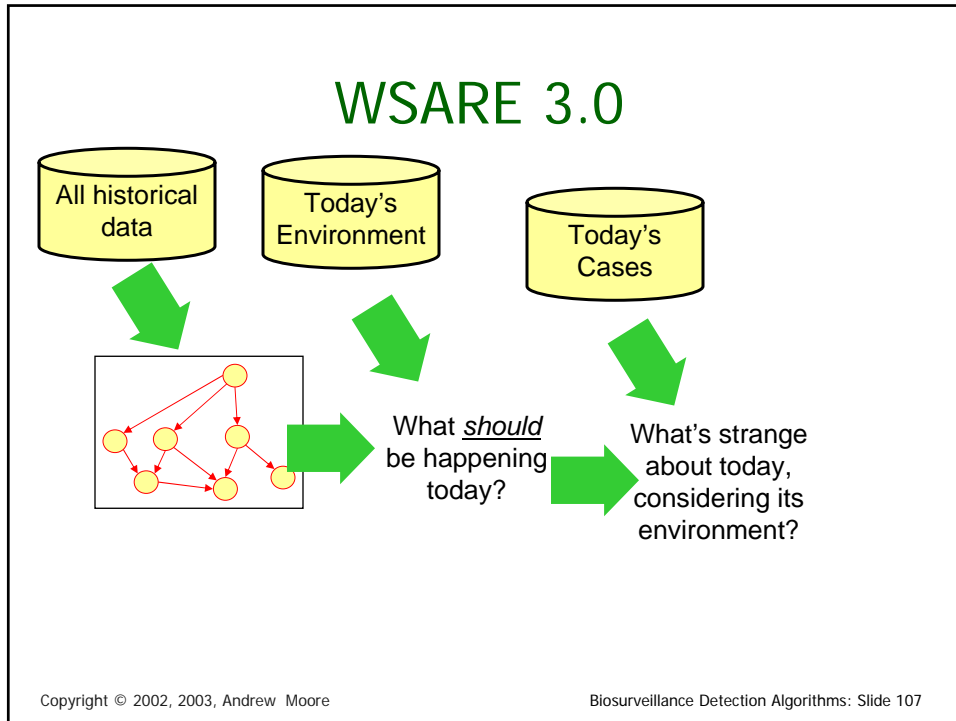
Biosurveillance Detection Algorithms: Slide 105

# WSARE 3.0

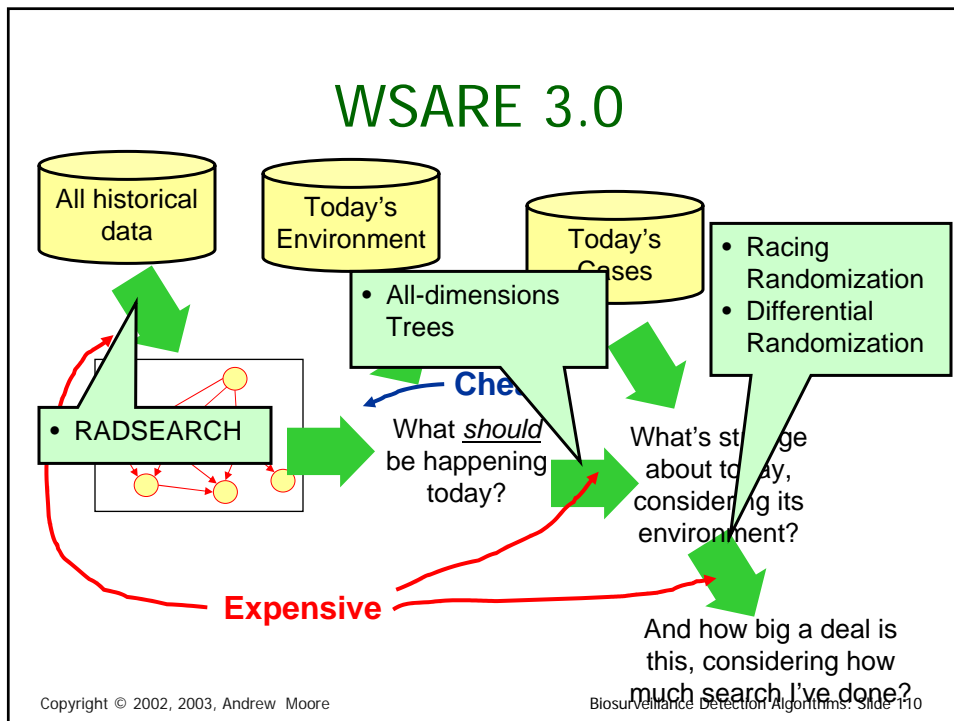
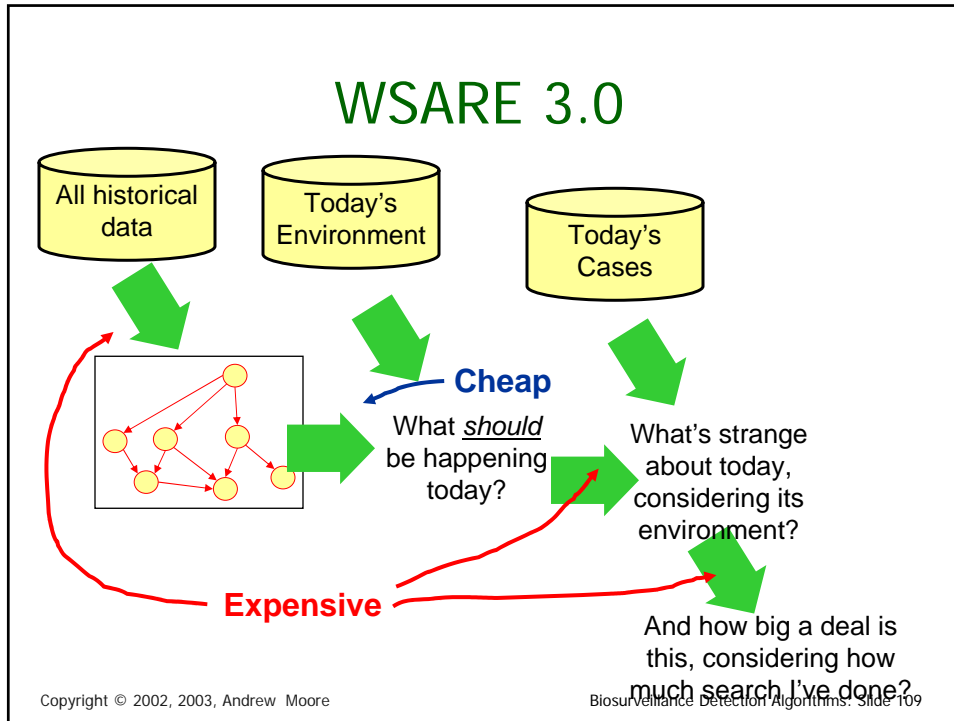


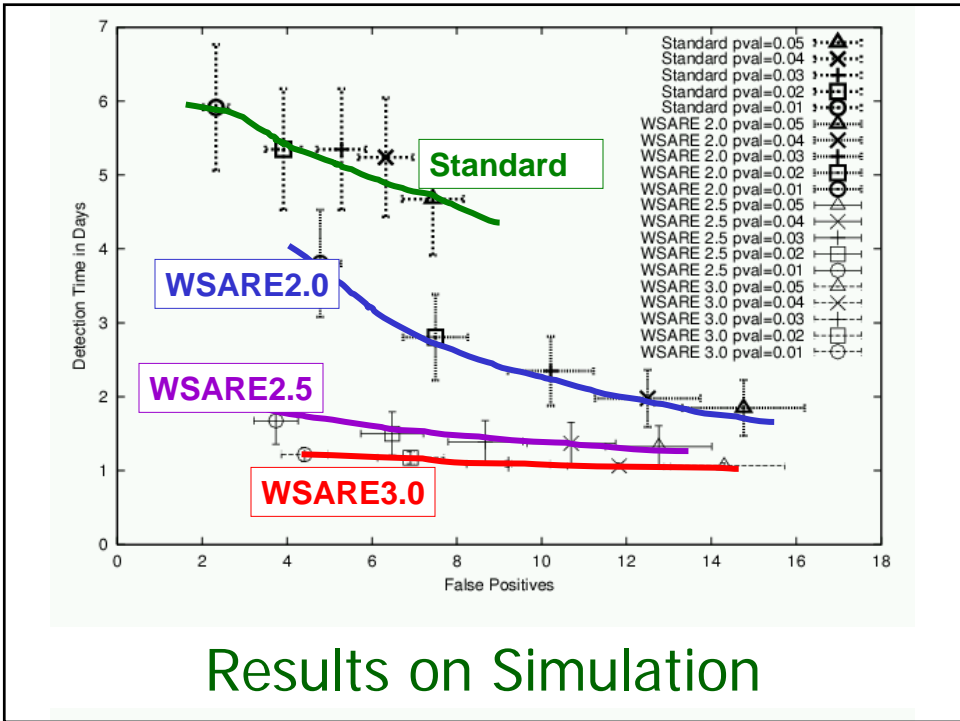
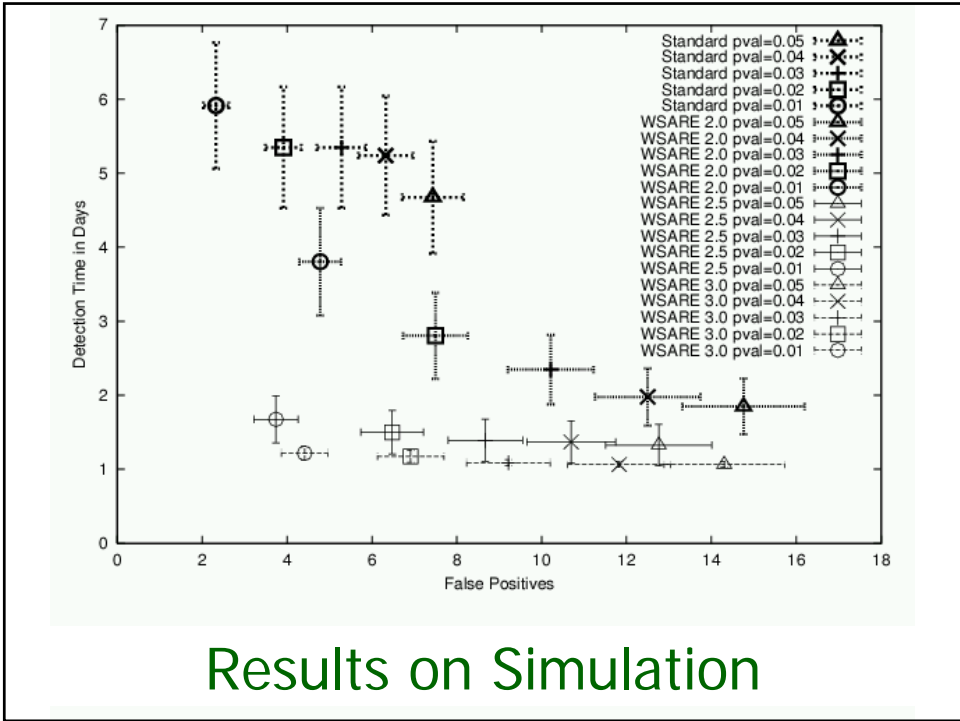
Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 106









## Conclusion

- One approach to biosurveillance: one algorithm monitoring millions of signals derived from multivariate data instead of Hundreds of univariate detectors
- Modeling historical data with Bayesian Networks to allow conditioning on unique features of today
- Computationally intense unless we're tricky!

- Searching over thousands of contingency tables on a large database...
- ...only we have to do it 10,000 times on the replicas during randomization
- ...we also need to learn Bayes Nets from databases with millions of records...
- ...and keep relearning them as data arrives online...
- ...in the end we typically search about a billion alternative Bayes net structures for modeling 800,000 records in 10 minutes
- Modeling historical data with Bayesian Networks to allow conditioning on unique features of today
- Computationally intense unless we're tricky!

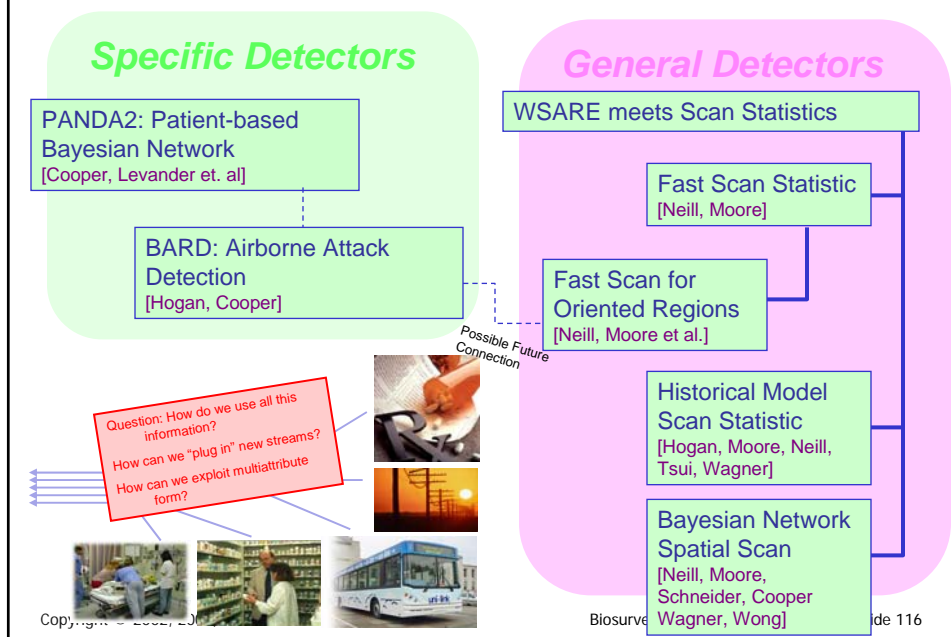
## Conclusion

- One approach to biosurveillance: one algorithm monitoring millions of signals derived from multivariate data  
instead of  
Hundreds of univariate detectors
- Modeling historical data with Bayesian Networks to allow conditioning on unique features of today
- Computationally intense unless we're tricky!
- WSARE 2.0 Deployed during the past year
- WSARE 3.0 about to go online
- WSARE now being extended to additionally exploit over the counter medicine sales

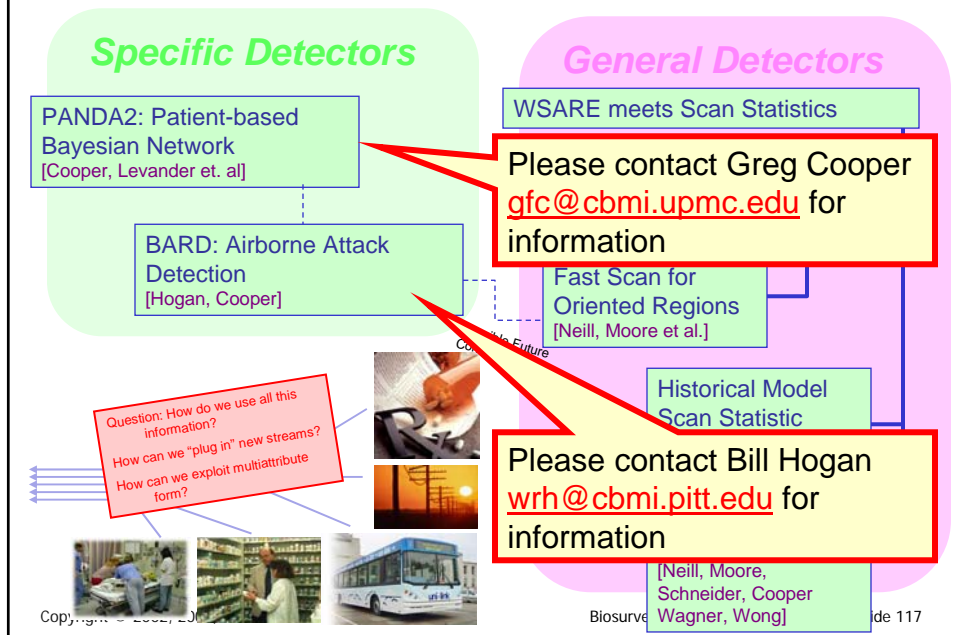
Copyright © 2002, 2003, Andrew Moore

Biosurveillance Detection Algorithms: Slide 115

## Other New Algorithmic Developments



# Other New Algorithmic Developments



## For further info

- Papers on these and other anti-terror applications:  
[www.cs.cmu.edu/~awm/antiterror](http://www.cs.cmu.edu/~awm/antiterror)
- Papers on scaling up many of these analysis methods:  
[www.cs.cmu.edu/~awm/papers.html](http://www.cs.cmu.edu/~awm/papers.html)
- Software implementing the above:  
[www.autonlab.org](http://www.autonlab.org)
- Copies of 18 lectures on 25 statistical data mining topics:  
[www.cs.cmu.edu/~awm/781](http://www.cs.cmu.edu/~awm/781)
- CD-ROM, powerpoint-synchronized video/audio recordings of the above lectures: [awm@cs.cmu.edu](mailto:awm@cs.cmu.edu)

Information Gain, Decision Trees  
 Probabilistic Reasoning, Bayes Classifiers, Density Estimation  
 Probability Densities in Data Mining  
 Gaussians in Data Mining  
 Maximum Likelihood Estimation  
 Gaussian Bayes Classifiers  
 Regression, Neural Nets  
 Overfitting: detection and avoidance  
 The many approaches to cross-validation  
 Locally Weighted Learning  
 Bayes Net, Bayes Net Structure Learning, Anomaly Detection  
 Andrew's Top 8 Favorite Regression Algorithms (Regression Trees, Cascade Correlation, Group Method Data Handling (GMDH), Multivariate Adaptive Regression Splines (MARS), Multilinear Interpolation, Radial Basis Functions, Robust Regression, Cascade Correlation + Projection Pursuit  
 Clustering, Mixture Models, Model Selection  
 K-means clustering and hierarchical clustering  
 Vapnik-Chervonenkis (VC) Dimensionality and Structural Risk Minimization  
 PAC Learning  
 Support Vector Machines  
 Time Series Analysis with Hidden Markov Models

## References

1. WSARE 3.0 : Bayesian Network based Anomaly Pattern Detection  
Wong, Moore, Cooper and Wagner [ICML/KDD 2003]
2. Fast Grid Based Computation of Spatial Scan Statistics  
Neill and Moore [NIPS 2003]

These and other Biosurveillance algorithms papers and free software available from

<http://www.autonlab.org/>

See also: <http://www.health.pitt.edu/rods>