

A Binomial Model of Transients in Daily ED Visits for Detecting Infectious Disease Outbreaks

James Dunyak
Kenneth Mandl
Mojdeh Mohtashemi

The threat of biological warfare and the emergence of new infectious agents spreading at a global scale have highlighted the need for major enhancements to the public health infrastructure. Effective confrontation of these urgent crises requires rapid and accurate detection of unusual epidemiologic trends, for which our current surveillance capabilities are not adequate. Critical for real time surveillance are two components: real-time data and real-time interpretation of data. Today, most existing surveillance systems are capable of monitoring and capturing real time data. However, the state of practice for detecting temporal and spatial abnormalities in surveillance data remains inadequate. We introduce a locally stationary binomial model of early detection of epidemiologic events, applied to real historical data pertaining to the daily number of visits with respiratory syndromes to the emergency department (ED). We show that when simulated outbreaks are introduced into the respiratory data, our uniformly most powerful detection algorithm under a constant false alarm rate is capable of detecting such irregularities in the data with high sensitivity, specificity, and in a timely manner.

The threat of biological warfare and the emergence of new infectious agents spreading at a global scale have accelerated efforts to improve current public health surveillance capabilities. According to an estimate by the Centers for Disease Control and Prevention (CDC), “as of May 2003, health departments in the United States have initiated syndromic surveillance systems in approximately 100 sites throughout the country” (1). Today, most existing surveillance systems are capable of monitoring and capturing real-time data. However, the state of practice for detecting abnormalities in surveillance data remains inadequate.

In the early stages of an outbreak, case fluctuations are highly stochastic. Statistical sampling of the infected population by emergency departments leads to a different stochastic effect. Not all infected individuals will appear for measurement at a specific emergency department. Some will go to private practitioners or other hospitals; some will simply not be treated. Such statistical sampling can be shown to be a dominant source of the short-term variation, with sampling effects a dominant component of the variability over short periods. Both the stochastic nature of an epidemic at its early stages and sampling effects should be considered when developing detection techniques. While the daily rates may not be accurately predictable due to intrinsic random variation from day to day that is difficult to account for, it is conceivable that the underlying dynamics producing the daily variations are governed by transient properties that can be tracked to detect change within a short time period.

Transients in biological processes are important for several reasons. First, in biological systems it is not always possible to wait for the *eventual* behavior to emerge. Long before the limiting behavior is reached, the system is perturbed by external impacts

so that all we observe are transient trajectories. Second, there are many processes that lead to the same eventual behavior, rendering the limiting behavior useless for diagnosis. Finally, early manifestations of many biological processes—such as those of illnesses from exposure to chemical or biological agents—are similar, often resulting in misdiagnosis of the disease. Therefore, new and innovative mathematical tools are needed to enhance our understanding of the underlying processes governing the transient dynamics of such phenomena (2, 3).

In this paper, we introduce a nonstationary binomial model of transients for detecting unusual epidemiologic events, applied to historical data on daily number of visits with respiratory syndromes. We show that when simulated outbreaks are introduced into the respiratory data, our uniformly most powerful detection algorithm under a constant false alarm rate is capable of detecting such irregularities in the data with high sensitivity, specificity, and in a timely manner. Throughout this paper, we make the key assumptions that the underlying disease processes are highly infectious and manifested with non-specific (flu-like) symptoms in patients early in the development.

MATERIALS AND METHODS

This paper develops and tests a widely applicable technique for detection of temporal anomalies to improve surveillance. The technique is verified using data collected from the emergency department (ED) of a large, urban, academic pediatric hospital (CH ED) from 6/1/1992-5/31/2003. ED chief complaints were used to identify patients with infectious respiratory illness based on a triage process using a pre-defined list of 181 choices. Using a previously validated subset of codes (4), historic time series were

developed describing the number of patients in this respiratory syndrome grouping per day. Institutional board approval was obtained.

Signal plus noise model

Our analysis begins with a straightforward model of the number of patients arriving at the emergency department each day meeting our syndrome grouping criteria. The number of daily visits $Y(t)$ is the sum of two components: the number of patients $S(t)$ associated with a new outbreak, and the number of background or “noise” patients $N(t)$ in the syndrome grouping not associated with the new outbreak in question.

$$Y(t) = S(t) + N(t)$$

During periods without new outbreaks, which spans our entire dataset, $S(t) = 0$.

Figure 1 is a snapshot of part of the time series for the CH ED dataset. Note the strong seasonal variation as well as the longer-term trend towards reduced ED visits. Strong day-of-the-week dependency is another source of variability in the dataset as illustrated by Figure 2. In our study dataset Sundays followed by Saturdays, on average, have the highest number of ED visits.

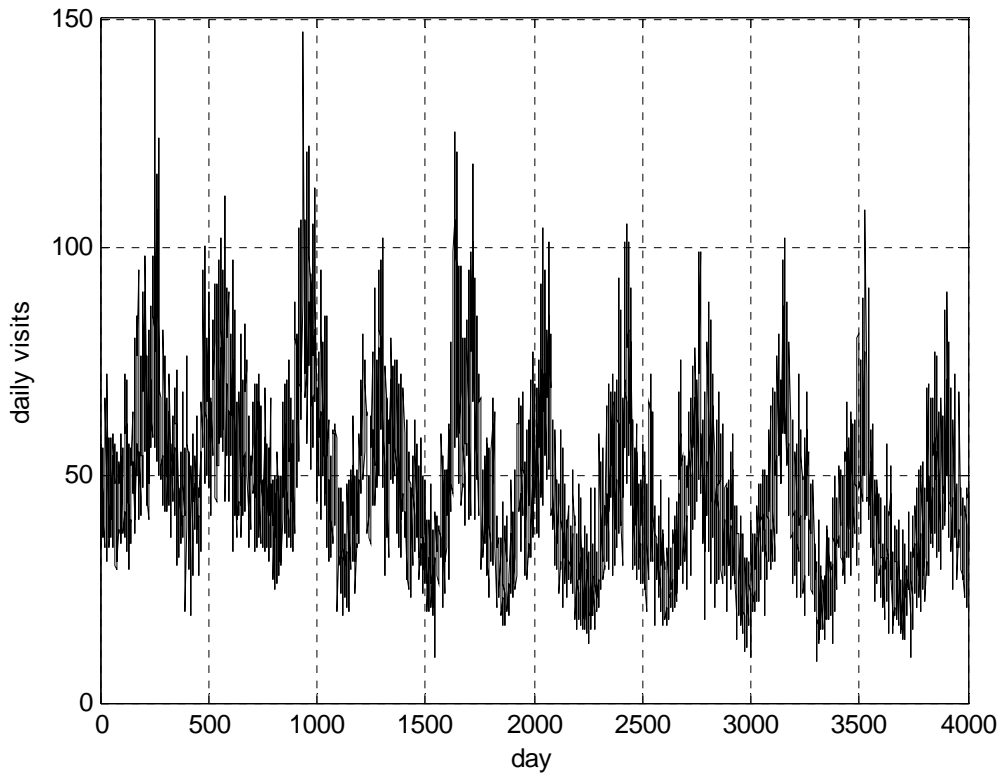


Figure 1: Daily Visits to the CH ED, June 1, 1996-May 31, 2000

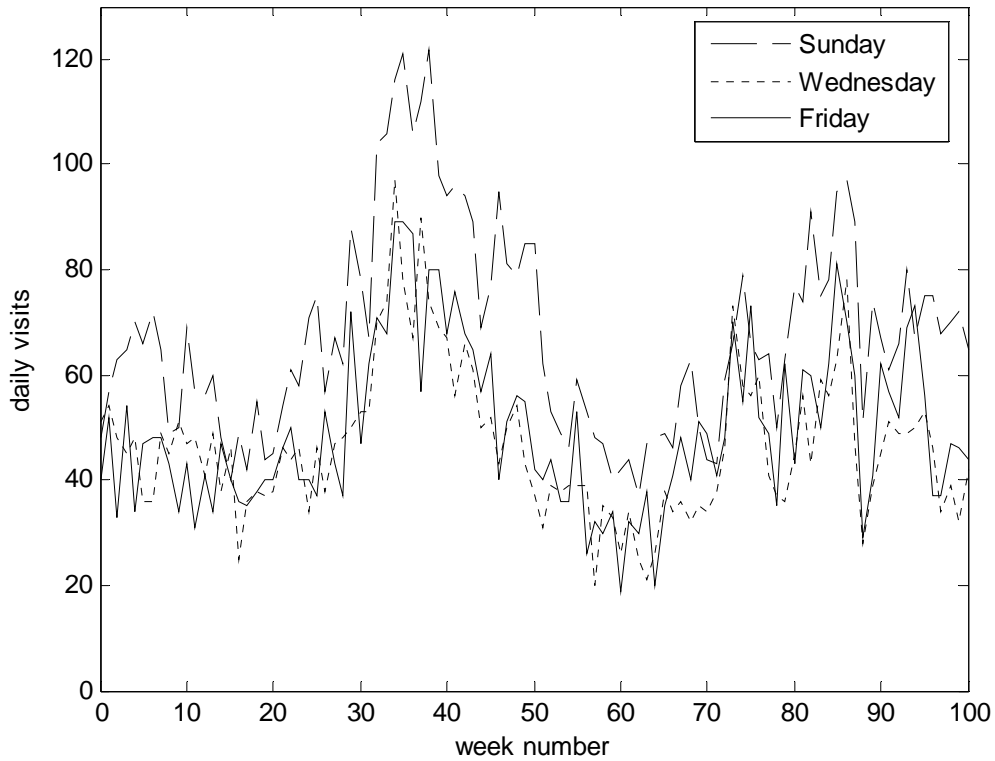


Figure 2: Day-of-the-week variability

Our goal is to identify when a new outbreak has occurred, $S(t) > 0$, in the presence of highly variable background noise $N(t)$. We do this through careful characterization and de-trending of the weekly and longer time scale variations in the background noise $N(t)$. The resulting constant false alarm rate detectors provide a uniformly most powerful test for the presence of an outbreak.

Binomial sampling models

A straightforward binomial model is developed to describe the time dynamics of $N(t)$, the background noise. Consider a catchment of size $K(t)$ that constitutes the sample

population from which the CH ED patient population is selected. The explicit dependence on time allows for a slowly evolving population size. On a particular day t , a member of the catchment population has a certain probability $p(t)$ of acquiring a respiratory infection, and thus experiencing a respiratory syndrome. The infected individual will then appear at the CH ED with a probability $q(t)$. In this case, a simple model for the number of patients arriving at the CH ED on a particular day is binomial

$$N(t) \square B(K(t), p(t)q(t)) .$$

A number of different timescales drive the evolving statistics of $N(t)$. Certainly over periods as long as the eleven years in this study the underlying population size $K(t)$ may have changed. Demographic attributes, such as an aging population, may cause $p(t)$ to vary, while time of the year and severity of the flu season will also affect $p(t)$. The parameter $q(t)$ may be influenced by changes in health care policies and insurance practices. Furthermore, the day of the week variability also impacts the probability of arriving at this specific ED given that a respiratory infection has occurred. These evolutions, however, are occurring on a slower scale, with a timescale on the order of weeks to years.

The various model parameters are not separately measurable from the data without measurement of many other variables. Therefore, we follow a highly nonparametric approach, which provides a robust technique applicable to other datasets from different institutions. The actual form of the detectors described below follow from applying the central limit theorem (CLT) to the large sample binomial distribution for $N(t)$. In this context, the time series is represented as a locally stationary Gaussian process with

$$N(t) \sim \text{Normal}\left(K(t)p(t)q(t), K(t)(p(t)q(t) - p(t)^2q(t)^2)\right).$$

Since $p(t)^2q(t)^2 \ll p(t)q(t) \ll 1$, we have

$$N(t) \sim \text{Normal}(K(t)p(t)q(t), K(t)p(t)q(t)). \quad [1]$$

We can thus estimate both the mean and variance from the time varying mean $K(t)p(t)q(t)$.

This application of the CLT greatly reduces the dependency of the final detector performance on accuracy of the binomial model. At the same time, the noise is characterized by a single parameter, namely the time varying mean, which can be readily estimated in real time from the data.

Our robust algorithm is a three-step procedure applied to the observed signal+noise time series $Y(t)$:

- 1) adjust for day-of-the week variability
- 2) estimate and remove the time-varying mean
- 3) apply a matched filter detection algorithm.

The day-of-the-week adjustment is location-specific but may be easily estimated using only several months of data on temporal patterns of the emergency department use. Estimating and removing the time-varying mean is motivated by the underlying phenomenology. This mean estimate and removal is a real-time process that is not location specific. Application of the normal approximation 1 then leads directly to a uniformly most powerful detector with a constant false alarm rate (see Hypothesis test for detection and Results).

Two critical assumptions are made in the detector development. First, the time-varying mean approximates the time-varying variance, as in approximation [1]. Second,

the day-of-the-week adjustment and removal of the time-varying mean results in a white process. These are both verified below using the historic dataset.

Day-of-the-week adjustments

Day-of-the-week adjustments are made by estimating and removing day-of-the-week means. Using our historical time series $Y(t)$,

$$d(\tau) = \frac{1}{11 \times 52} \sum_{t=1}^{11 \times 52} Y(\tau + 7(t-1)), \quad \tau = 0, 1, 2, \dots, 6$$

where $d(t)$ is the day of the week specific mean. Table 1 lists the resulting day-of-the-week mean offsets. The resulting time series after removing daily variation is then $Y'(t) = Y(t) - d(t)$.

Here $d(t)$ is the obvious periodic correction for day of the week. Next, we estimate the time-varying mean in the locally stationary approximation.

Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
2	-4	-6	-6	-4	6	13

Table 1: Mean Offsets for the Days of the Week

Estimating the slowly varying mean

The time-varying mean, after removal of daily variation, is estimated through a low-pass filter with coefficients $f(\tau)$

$$\hat{N}(t) \approx \sum_{\tau=0}^{L-1} f(\tau) [Y'(t-\tau)] = \sum_{\tau=0}^{L-1} f(\tau) [Y(t-\tau) - d(t-\tau)] \quad [2]$$

The most obvious choice is a simple block average, with $f(\tau) = 1/L$ for a window of length L . Performance, however, is significantly degraded in this case due to the high side lobes associated with this block filter. Since spectral content is the main separating feature between rapid onset of outbreaks (with high frequency energy) and the slowly varying mean (with low frequency energy), low filter side lobes are important to overall performance. We address side lobe concerns, while maintaining the interpretation of a time-varying mean estimate, through use of a standard Hamming window for $f(\tau)$ (5),

$$f(\tau+1) = 0.54 - 0.46 \cos\left(2\pi \frac{\tau}{L-1}\right), \quad \tau = 0, 1, \dots, L-1.$$

The resulting filter is comparatively long in duration as compared to other low-pass filter designs, but the additional averaging contributes to the robustness of the design.

Figure 3 shows the original time series $Y(t)$ and the time-varying mean estimate $\bar{N}(t)$ from equation 3. The time series in Figure 3 is a portion of our CH ED from a period with no known outbreaks of concern. Note that the time-varying mean estimate captures a great deal of the variability of the original time series. After subtracting out the time-varying mean,

$$\tilde{Y}(t) = Y'(t) - \sum_{\tau=0}^{L-1} f(\tau)[Y'(t-\tau)] = Y(t) - d(t) - \sum_{\tau=0}^{L-1} f(\tau)[Y(t-\tau) - d(t-\tau)].$$

This may be rewritten as a high pass filter using the discrete impulse $\delta(t)$

$$\tilde{Y}(t) = \sum_{\tau=0}^{L-1} [\delta(\tau) - f(\tau)][Y'(t-\tau)] = \sum_{\tau=0}^{L-1} [\delta(\tau) - f(\tau)][Y(t-\tau) - d(t-\tau)]. \quad [3]$$

$\tilde{Y}(t)$ is the statistic we use for early detection of outbreaks. Note that the filter blocks slowly varying components while primarily passing energy associated with the outbreak.

This filter detects changes only and therefore is appropriate only for detection of the onset of an outbreak. After the outbreak has spread through the population and settled closer to a stable point, the filter will remove the contribution from the outbreak in the statistic $\tilde{Y}(t)$.

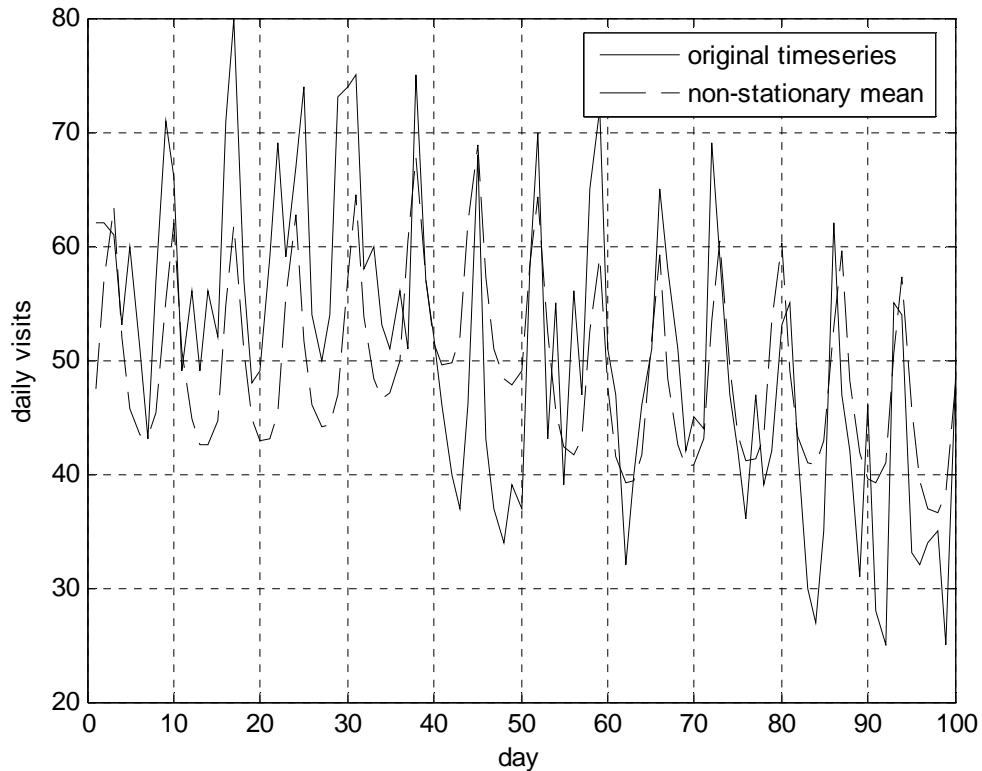


Figure 3: The Original Time series and the Estimated Non-stationary Mean

Model validation: binomial patient arrivals

We demonstrate our approach on our CH ED dataset, but the reader should note that our technique captures the population health seeking behavior and is generally applicable to different datasets from different institutions, which obviates the need to establish its validity every time a new dataset is analyzed. The binomial assumption itself is not

critical, due to the Central Limit Theorem, but the assumption that the time-varying mean approximates the time varying variance, as in approximation [1], is critical. After removing the time-varying mean via equation [3], the resulting locally stationary zero mean process has a time-varying variance. To develop a detector with constant rate of false alarm, we apply the approximation

$$std(\tilde{Y}(t)) \approx \sqrt{\bar{N}(t)} .$$

Testing this approximation is complicated because $\bar{N}(t)$ changes too rapidly to allow enough averaging for a high quality estimate. Instead, we sort the values of $\bar{N}(t)$ in ascending order resulting in $\bar{N}(i)$ for $i = 1..size(dataset)$. Since the standard deviation provides the natural scale for error probabilities, we then estimate $std(\tilde{Y}(t))$ using the ensemble of time points with nearest values of $\bar{N}(i)$. We test this approximation on our historic time series. Figure 4 shows the resulting $\bar{N}(i)$ and estimated $std(\tilde{Y}(t))$, using the sorted values of $\bar{N}(t)$. Note that $\bar{N}(i)$ is a reasonable approximation of $std(\tilde{Y}(t))$. This approximation is not perfect and would certainly fail any statistical test of fit. However, the approximation is based on a simple underlying model, the binomial sampling of patients, and as such is highly robust and widely applicable. Use of this approximation allows the development of early surveillance approaches without the expensive, arduous, and often impossible, task of collecting many years of syndromic-specific data for each location.

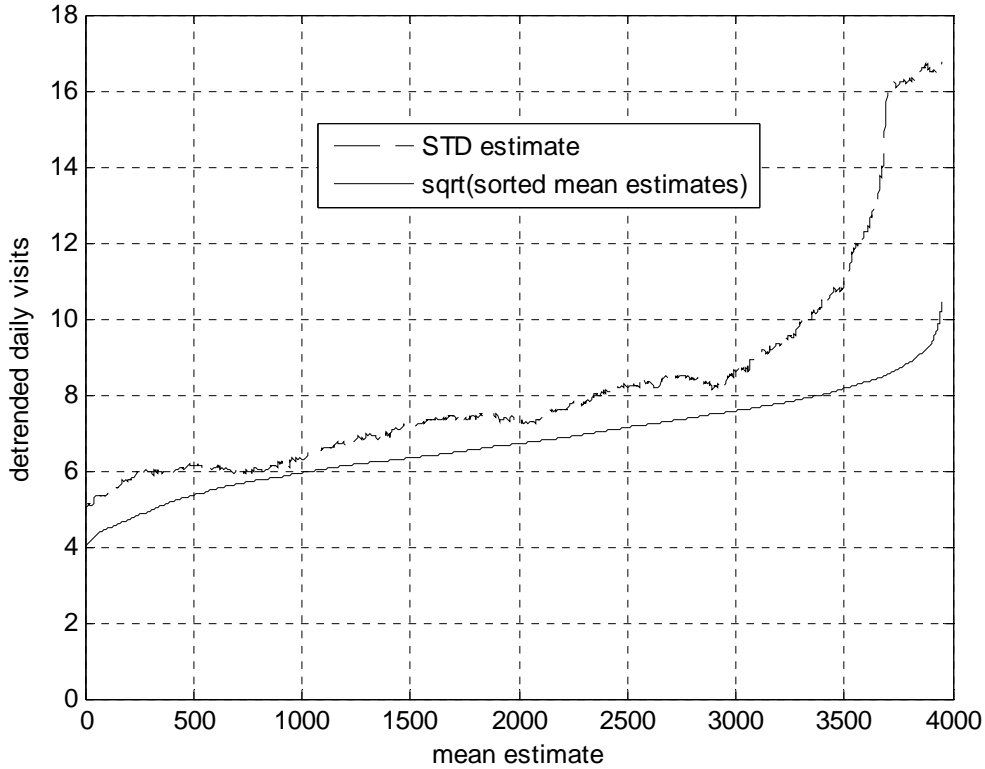


Figure 4: A Comparison of Process Standard Deviation $std(\tilde{Y}(t))$ and its Binomial

$$\text{Estimate } \sqrt{\frac{E}{N(t)}} \text{ for the CH ED Dataset}$$

Model validation: prewhitening

In the usual development of matched filter designs, we would derive a prewhitening filter to specifically match the spectrum of the data set. This approach requires large historic data sets, so we avoid it here. Instead, we demonstrate that our simple locally stationary binomial model acts as a prewhitening filter.

Under a no-outbreak condition, we have developed an approximate distribution for the marginal distribution of $\tilde{Y}(t)$. We now investigate time dependence. Using a centered Hamming window, the normalized autocorrelation in our historic time series is

shown in Figure 5. Note that the time series for $\tilde{Y}(t)$ de-correlates in a single day.

Following the Central Limit Theorem, we model $\tilde{Y}(t)$ as normally distributed and uncorrelated from time sample to time sample.

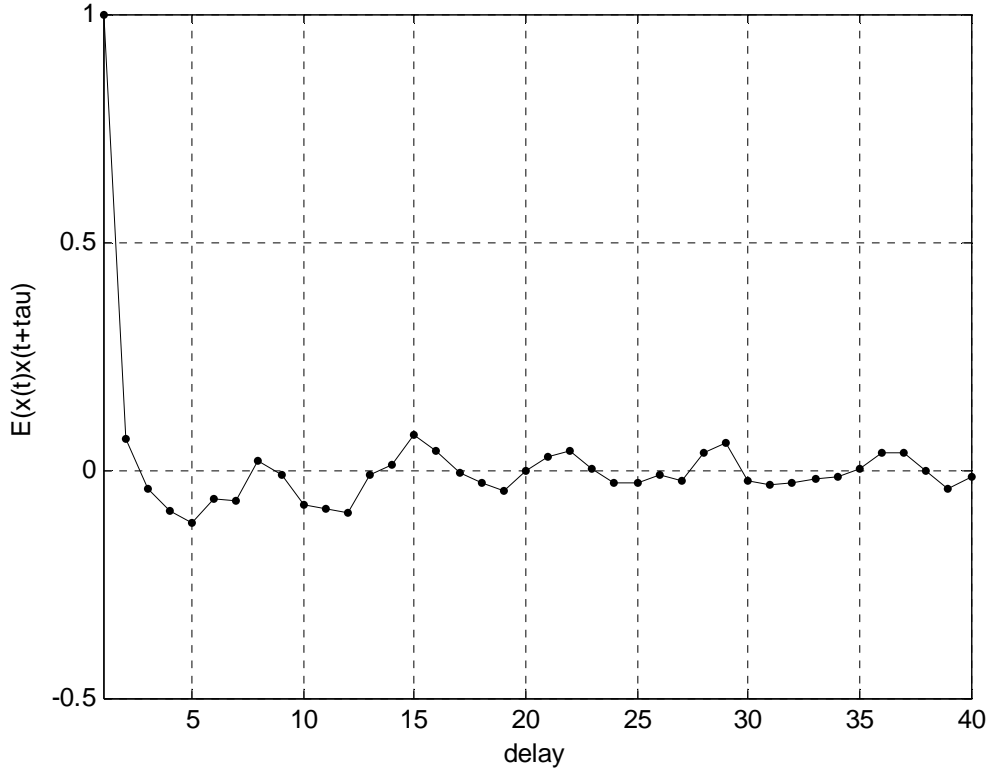


Figure 5: Normalized Autocorrelation of $\tilde{Y}(t)$

Hypothesis test for detection

We can now develop a uniformly most-powerful detection algorithm for a given class of outbreaks. Consider a class of outbreaks of the form

$$S(t) = \begin{cases} \beta s(t) & t = 0, 1, \dots, (T-1), \dots & \beta > 0 \\ 0 & t < 0 \end{cases} \quad [4]$$

where β is the class scale factor, $s(t)$ is the epidemic trajectory, and T is the detector analysis time period. Here the outbreak begins, for convenience, at time $t = 0$, Examples

of important trajectories include linear, constant where $s(t) = 1$, for all t , corresponding to infections due to environmental toxin exposure or perhaps chemical attack (6), concave up exponential functions following classic contagious outbreaks, and concave down exponential functions (7). The issue of trajectory choice is discussed in more detail in the appendix. The hypothesis testing problem for using T measurements to detect an outbreak initiated at time $t = 0$ is then

$$\begin{aligned} H_0 : Y(t) &= N(t) \\ H_1 : Y(t) &= \beta s(t) + N(t), \quad \beta > 0, \quad t = 0, 1, 2, \dots (T-1), \dots \end{aligned}$$

This test is applied as a sliding window, as discussed below. The log-likelihood test is then, for a window of length T , using the distribution $\tilde{Y}(t) \square Normal(0, \frac{\sigma^2}{N(t)})$ and independent in time,

$$\beta \sum_{t=0}^{T-1} \frac{s(t)\tilde{Y}(t)}{\frac{\sigma^2}{N(t)}} \begin{array}{l} H_1 \\ > a. \\ < \\ H_0 \end{array}$$

Note that, for a test of size α , the threshold a is a function of many parameters. However, we can use monotonicity to also write the decision in terms of another threshold a' as

$$\frac{1}{\sqrt{\sum_{t=0}^{T-1} \frac{s(t)^2}{\frac{\sigma^2}{N(t)}}}} \sum_{t=0}^{T-1} \frac{s(t)\tilde{Y}(t)}{\frac{\sigma^2}{N(t)}} \begin{array}{l} H_1 \\ > a'. \\ < \\ H_0 \end{array} \quad [5]$$

The detector in inequality [5] is a matched filter between the de-trended patient visit count $\tilde{Y}(\square)$ and the function

$$\frac{s(\square)}{\sqrt{\frac{1}{N(\square)} \sum_{\tau=0}^{T-1} \frac{s(\tau)^2}{N(\tau)}}}.$$

The matched filter provides the optimal detector in the additive white Gaussian noise case. We use this fact to motivate the choice of epidemic trajectory $s(\square)$ below.

Inequality [5] may also be viewed as comparing the summation to a time varying threshold, resulting in the time-varying test shown in Figure 6.

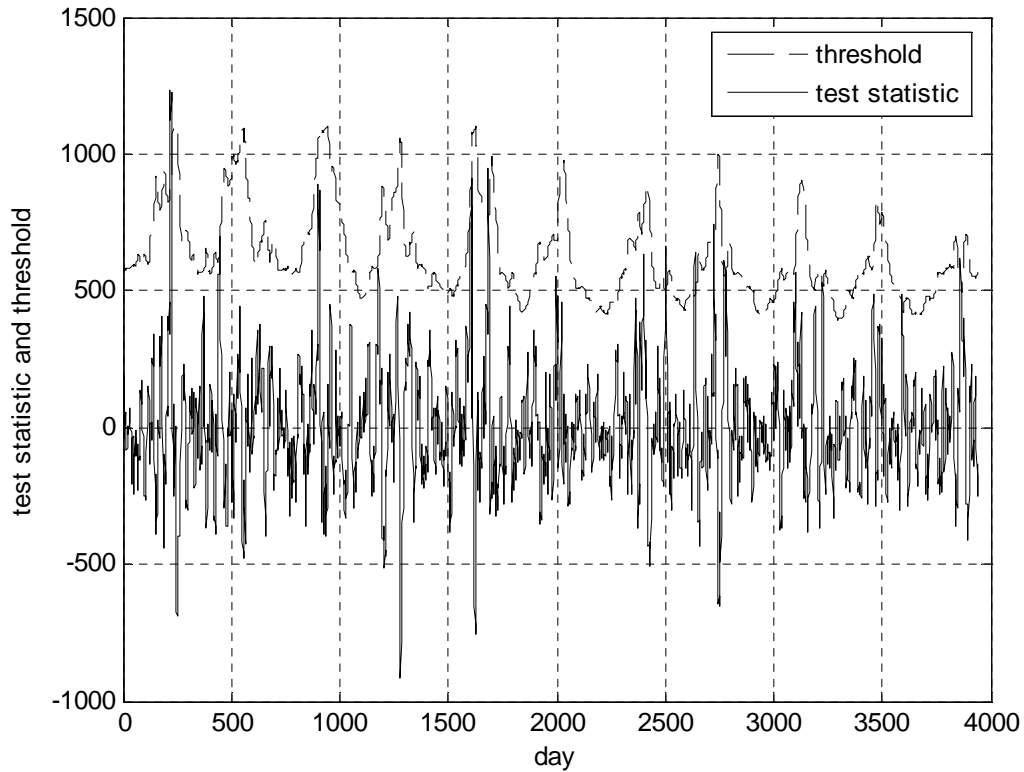


Figure 6: A Time-varying Test Statistic

The quantity on the left hand side of inequality 5 is now, under H_0 and within our approximation, a standard normal. We choose a' so that, for standard normal Z and test size α , $P(Z > a') = \alpha$.

Note that $a\phi$ is independent of the scaling factor β , making the test uniformly most powerful over the class in equation 4. This is a critical observation, since this property leads to a constant false alarm rate (CFAR) detector (5) over the epidemic class in equation [4] (see RESULTS).

The test statistic in inequality [5] is, under H_1 and within our approximation, also normally distributed

$$Normal\left(\beta\sqrt{\frac{\sum_{t=0}^{T-1} s(t)^2}{\bar{N}(t)}}, 1\right). \quad [6]$$

This allows calculation of an instantaneous probability of detection for an outbreak of size β . The reader should note, however, that both false alarm and detection events are heavily correlated over short time periods, due to the sliding window form of the detector. Note, in particular, that the probability of detection is determined by $\sum_{t=0}^{T-1} s(t)^2$. Outbreaks with increasing profiles (such as occurs in contagious disease) are more difficult to detect rapidly than outbreaks with a more constant profile. Actual performance of the test on experimental time series is analyzed below in the section on sensitivity and specificity.

Here, we confine our demonstration to linear epidemic trajectories $s(t)$ (see Appendix for a mathematical justification of linear and exponential epidemic trajectories). Exponential or polynomial shapes may also be considered, as appropriate for longer windows. Constant or even declining trajectories hold special interest for detection of outbreaks due to a constant or one-time toxin exposure. Any targeted phenomena that can be captured in a shape class are candidates for a specialized detector. However, the highly stochastic nature of the measurements will tend to overwhelm small

differences in outbreak trajectory $s(t)$. A linearly growing or constant shape will capture the generic contagious and non-contagious cases effectively.

RESULTS

Sensitivity and specificity

To address sensitivity and specificity, we randomly embedded linearly growing and constant outbreaks of various scale β , with outbreak $\beta * [1, 2, 3, \dots]$ or $\beta * [5, 5, 5, \dots]$, at thousands of locations in the CH ED dataset. Probabilities of detection were determined using a seven-day window, so an outbreak was considered detected if and only if it was detected seven days into the outbreak. (The outbreak continued after the seven days, but the detector is causal so this growth had no impact on the detection statistics.) False alarms were determined using the historic data with no simulated outbreaks. The detector, with its seven-day window, was highly correlated on adjacent days. False alarms, when they occurred, tended to persist for more than one day. False alarm events (exceeding threshold) occurring within seven days of a false alarm local maximum were classified as a single false alarm event. The reader should note, however, that most false alarm rate events contained far fewer than seven days above threshold. Our goal was to have two or fewer false alarm events per year, which is well within the range adopted in the literature (6, 7). Choosing many different values of a' in inequality [5] allows us to establish the performance of the resulting CFAR detectors as illustrated in Figure 7.

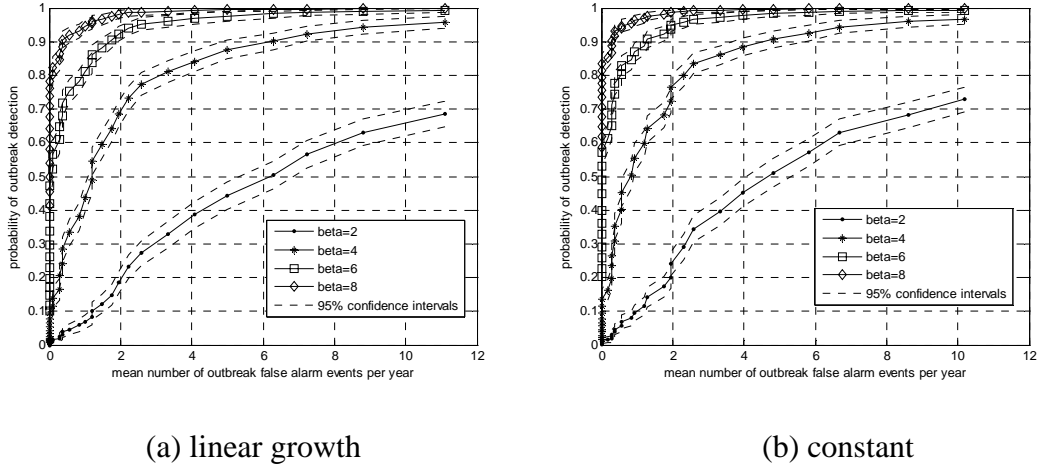
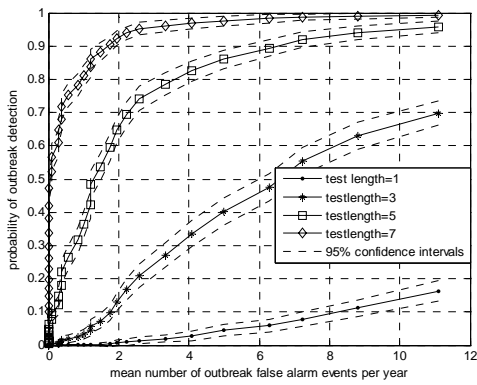


Figure 7: Probability of detection versus mean number of false alarm events per year for various epidemic sizes (a) $\beta^*[1,2,3,\dots]$ and (b) $\beta^*[5,5,5,\dots]$, using a one week window.

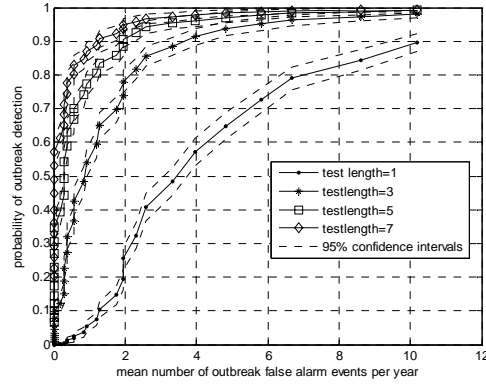
As expected, high probabilities of detection ($P_D > 0.95$) can only be achieved for larger sized outbreaks or by accepting a higher number of false alarm events. For example, an epidemic with a mean profile of $[6,12,18,24,30,36,42,\dots]$ over seven days can be detected with about 95% probability with on average two false alarm events per year. On the other hand, accepting a lower probability of detection of $P_D > 0.7$ allows detection of an epidemic with a mean profile of $[4,8,12,16,20,24,28,\dots]$ while still maintaining two false alarm events per year. A small outbreak with mean profile $[1,2,3,4,5,6,7,\dots]$ can only be detected within one week with $P_D > 0.5$ by tolerating 12 false alarm events per year or an average of one false alarm event per month. Note from Figure 7b that outbreaks following a constant profile, $\beta^*[5,5,5,\dots]$, are easier to detect. This is expected from the H_1 distribution in [6], since, even though the peak value is smaller, the total signal power in the window $\sum_{t=0}^{T-1} s(t)^2$ may be larger.

Timeliness of detection

Figure 8 explores the timeliness of detection by examining different test lengths for (a) linear $\beta^*[1,2,3,\dots]$ and (b) constant $\beta^*[5,5,5,\dots]$ outbreaks with $\beta=6$. Note that for detection window of length $T < 5$ the detector in [5] provides limited outbreak detection capability. By tolerating an average of two false alarm events per year over a five-day window the linear outbreak can be detected with probability of over 65% ($P_D > 0.65$), while the same linear outbreak can be detected with probability of about 95% if we increase the length of the detection window to seven. Constant size outbreaks are much easier to detect very earlier in the outbreak. Linear growth (with its small initial state) provides little accumulation of $\sum_{t=0}^{T-1} s(t)^2$ for small T , so the H_1 distribution in [6] shows that the detection probability will be small.



(a) linear growth



(b) constant

Figure 8: Probability of detection versus mean number of false alarm events per year for epidemic sizes (a) $6^*[1,2,3,\dots]$ and (b) $6^*[5,5,5,\dots]$ using different test lengths for detection window.

DISCUSSION

We proposed a non-parametric model of transients in the ED time series pertaining to respiratory syndromes. The benefit of the uniformly most powerful approach developed here is that the detector with constant false alarm probability is not dependent on the epidemic scaling factor β . This is critical, since the parameter β is dependent on many unknown and unobservable quantities.

The underlying technique is a simple locally stationary model of daily variations which makes the approach robust and widely applicable to different data types including non-respiratory or non-contagious disease processes. This is unlike the SEIR models of detection (7) that are specifically suited for capturing the underlying contagious transmission dynamics of the disease. Furthermore, unlike some time series surveillance techniques, including the ARMA and SEIR models (6, 7), that depend on model training and parameter estimations for detection, our proposed method does not require large historic records of patient visits in order to begin surveillance. Another advantage of the proposed method is that the detector is sensitive to the slow growth associated with the early stages of exponential shapes of epidemic trajectories that mostly resemble those of contagious infectious diseases. Constant epidemic shapes, which more closely model an environmental toxin exposure or a chemical attack, are easier to spot due to the early clustering effect that may be unavoidable under such circumstances. Constant epidemics were adopted by Reis, et al to test an ARMA model of the total ED visits (6).

Some of the limitations of the proposed method need to be addressed. One such limitation, shared by all surveillance techniques that are based on syndromic data, can be viewed as lack of adequate power to provide timely warning. This is evident in Figure 8 where a longer choice for the detection window provides higher probability of detection.

However, this is mostly due to inherent non-specific properties of syndromic data. More symptom-specific data are needed to achieve more detection power. Another constraint in the model is due to the binomial approximation of the process variance with the locally stationary mean. As demonstrated in Figure 4, the approximation is less valid for larger values of mean. This can potentially reduce the power of the detector or limit its utility when the mean background noise is very large.

One critical issue for public health surveillance is the absence of a uniform approach for evaluating and comparing surveillance techniques. While additional work in developing effective surveillance techniques are needed, it is quite conceivable that different detection techniques may perform differently under various outbreak conditions or datasets. Some algorithms may be better suited for capturing the early transmission dynamics of contagious disease while others respond to aggregate levels of different disease processes. Finally, some techniques may be region-specific and sensitive to localized clustering of disease incidents in time and space while others detect elevated numbers across an entire area. Relative advantages and disadvantages of different surveillance techniques cannot be systematically addressed until a uniform evaluation approach is adopted.

REFERENCES

1. Buehler JW, Berkelman R, Hartley DM, Peters CJ. Syndromic surveillance and bioterrorism-related epidemics. *Emerging Infectious Diseases* 2003; 9 (10):1197-1204.

2. Mohtashemi M, Levins R. The early dynamics and diagnostics of malaria infection. Presented at the New Challenges in Tropical Medicine and Parasitology, Oxford, UK, September 2000.
3. Mohtashemi M, Levins R. Transient dynamics and early diagnostics in infectious disease. *J Math Bio* 2001; 43, 446-70.
4. Beitel AJ, Olson KL, Reis BY, et al. Use of emergency department chief complaint and diagnostic codes for identifying respiratory illness in a pediatric population. *Pediatric Emergency Care*. (in press)
5. Van Trees HL. *Detection, Estimation, and Modulation Theory, Part I*. Wiley-Interscience, 2001.
6. Reis BY, Pagano M, Mandl KD. Using temporal context to improve biosurveillance. *Proc Natl Acad Sci USA* 2003; 100, 1961-5.
7. Mohtashemi M, Szolovits P, Mandl KD. A mathematical model of early detection of outbreaks of contagious infectious disease. (in prep)
8. Tsui F-C, Espino JU, Dato, VM et al. Technical description of RODS: a real-time public health surveillance system. *J Am Med Inform Assoc* 2003; 10(5):399-408.
9. Lewis MD, Pavlin JA, Mansfield JL, et al. Disease outbreak detection system using syndromic data in the greater Washington DC area. *Am J Prevent Med* 2002; 23(3):180-6.
10. Lober WB, Karras BT, Wagner MM, et al. Roundtable on bioterrorism detection: information system-based surveillance. *J Am Med Inform Assoc* 2002; 9(2):105-15.

11. Greenko J, Mostashari F, Fine A, et al. Clinical evaluation of the Emergency Medical Services (EMS) ambulance dispatch-based syndromic surveillance system, New York City. *J Urban Health* 2003; 80(2 Suppl 1):i50-6.
12. Goldenberg A, Shmueli G, Caruana RA, et al. Early statistical detection of anthrax outbreaks by tracking over-the-counter medication sales. *Proc Natl Acad Sci USA* 2002; 99, 5237-40.
13. Kleinman K, Lazarus R, Platt R. A generalized linear mixed models approach for detecting incident clusters of disease in small areas, with an application to biological terrorism. *Am J Epidemiol* 2004; 159(3):217-24.

APPENDIX

Epidemic trajectories The choice of trajectory $s(\square)$ in equation (4) is a critical issue in the development of the detector. Consider a simple stochastic equation

$$I(n+1) = I(n) + \Delta_{\gamma}I(n) - \Delta_{\delta}I(n) .$$

Here $I(n)$ represent the proportion of infected individuals at time n ; $\Delta_{\gamma}I(\square)$ is the stochastic increments associated with new infections and $\Delta_{\delta}I(\square)$ is the stochastic decrements associated with removal of the infected. Under almost all circumstances, early in the outbreak we would expect the conditional mean of the increments to be proportional to the number infected. This follows from the dynamics of individuals spreading the infection and recovering independently. In this case, for proportionality constants γ and δ , we have

$$E(\Delta_{\gamma}I(n) | I(n)) = \gamma I(n) \quad E(\Delta_{\delta}I(n) | I(n)) = \delta I(n)$$

Thus the conditional expected value of the proportion infected at time $t = n + m$ is

$$E(I(n+m) | I(n)) = (1 + \gamma - \delta)^m I(n).$$

Using the approximation $(1+x)^m \approx 1+mx$ for $x \ll 1$, we can rewrite the above equation for short time periods as

$$E(I(n+m) | I(n)) = I(n) + mI(n)(\gamma - \delta)$$

Note that $I(n)$ is the proportion of the population who are infected, and thus we have $\gamma - \delta \ll 1$. A linearization of the dynamics is, not surprisingly, a good approximation for short periods during early stages of an outbreak. Finally, application of the result

$e^x = \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n$ suggests a linear or exponential model of the epidemic trajectory, but in

this case it is the conditional mean that grows approximately linearly or exponentially.

Deviations from the conditional mean $I(n+m) - E(I(n+m) | I(n))$ are then viewed as

“noise” in the hypothesis testing problem. The mean of this distribution is again

proportional to $I(n)$ and the deviation from mean is also viewed as noise.