

Machine Learning and Event Detection for Population Health

Daniel B. Neill, Ph.D.

**Associate Professor of Computer Science and Public Service
Associate Professor of Urban Analytics, NYU CUSP
Director, Machine Learning for Good (ML4G) Laboratory**

New York University

E-mail: daniel.neill@nyu.edu

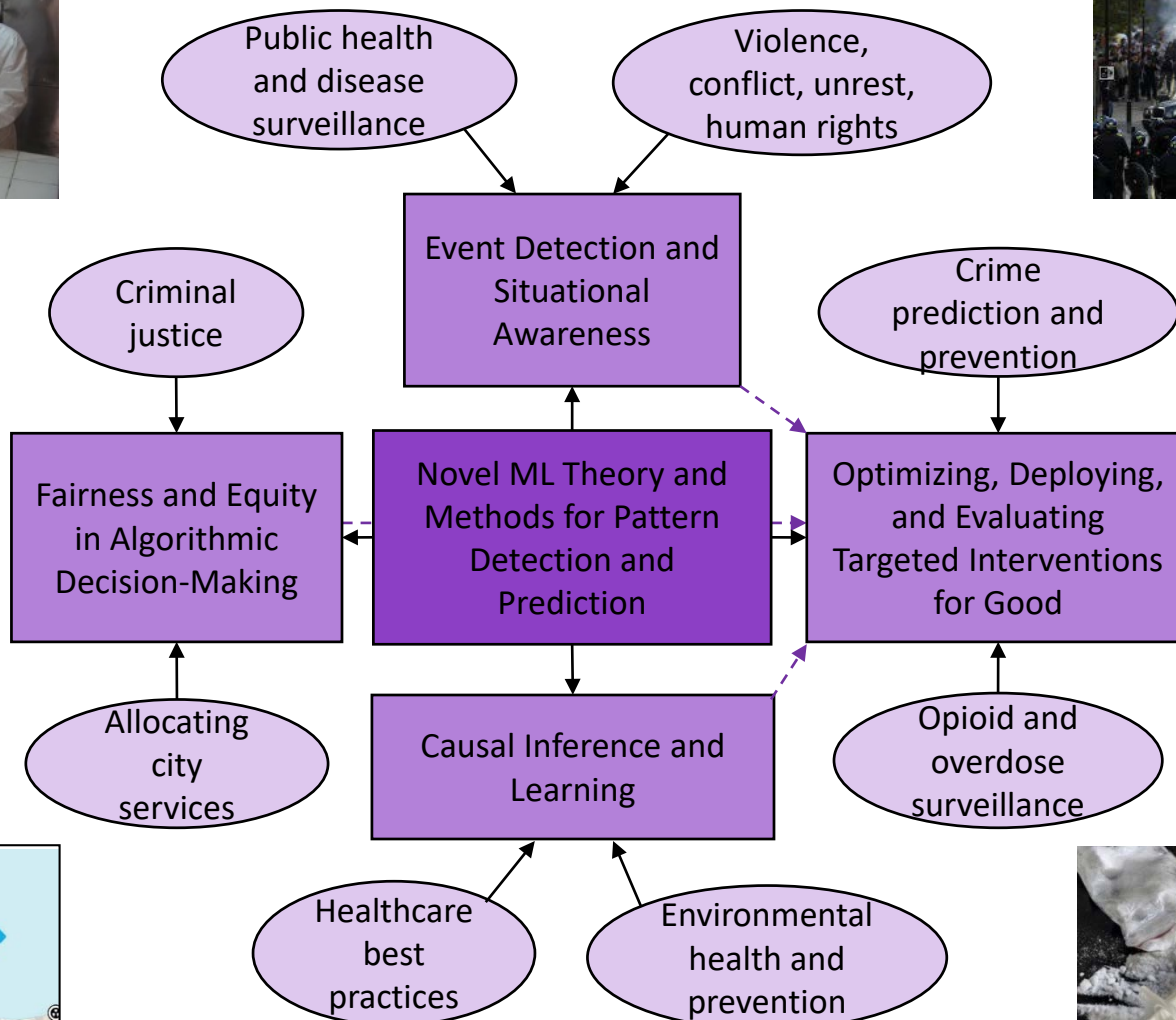
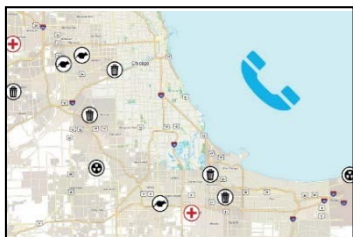
Web: <http://www.cs.nyu.edu/~neill>



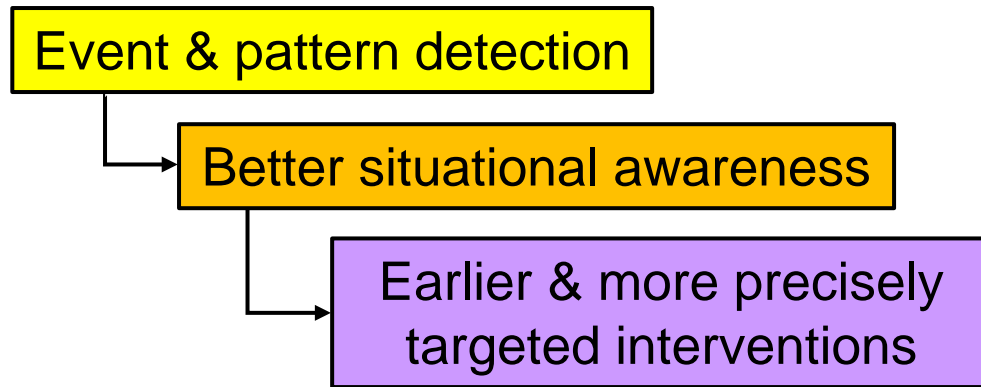
NYU

Center for Urban
Science + Progress

The Machine Learning for Good Lab @ NYU



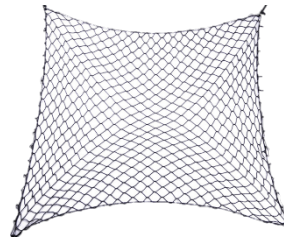
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities

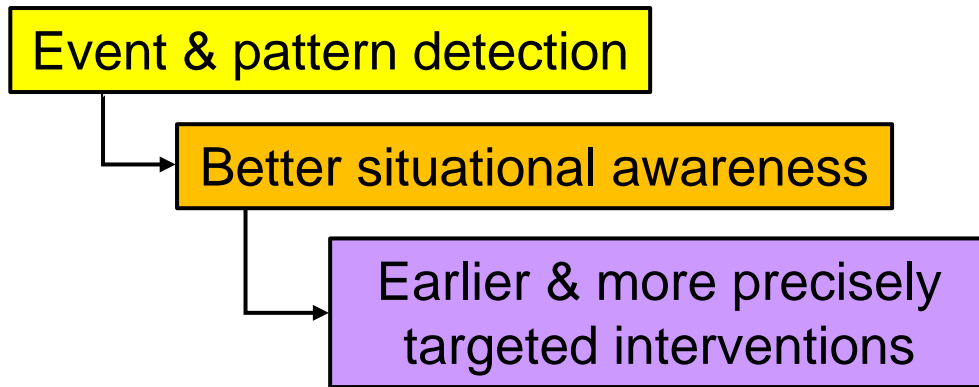


Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis

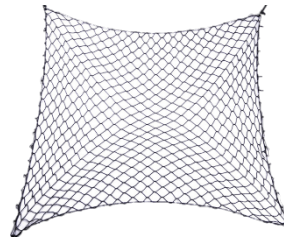
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities



Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis

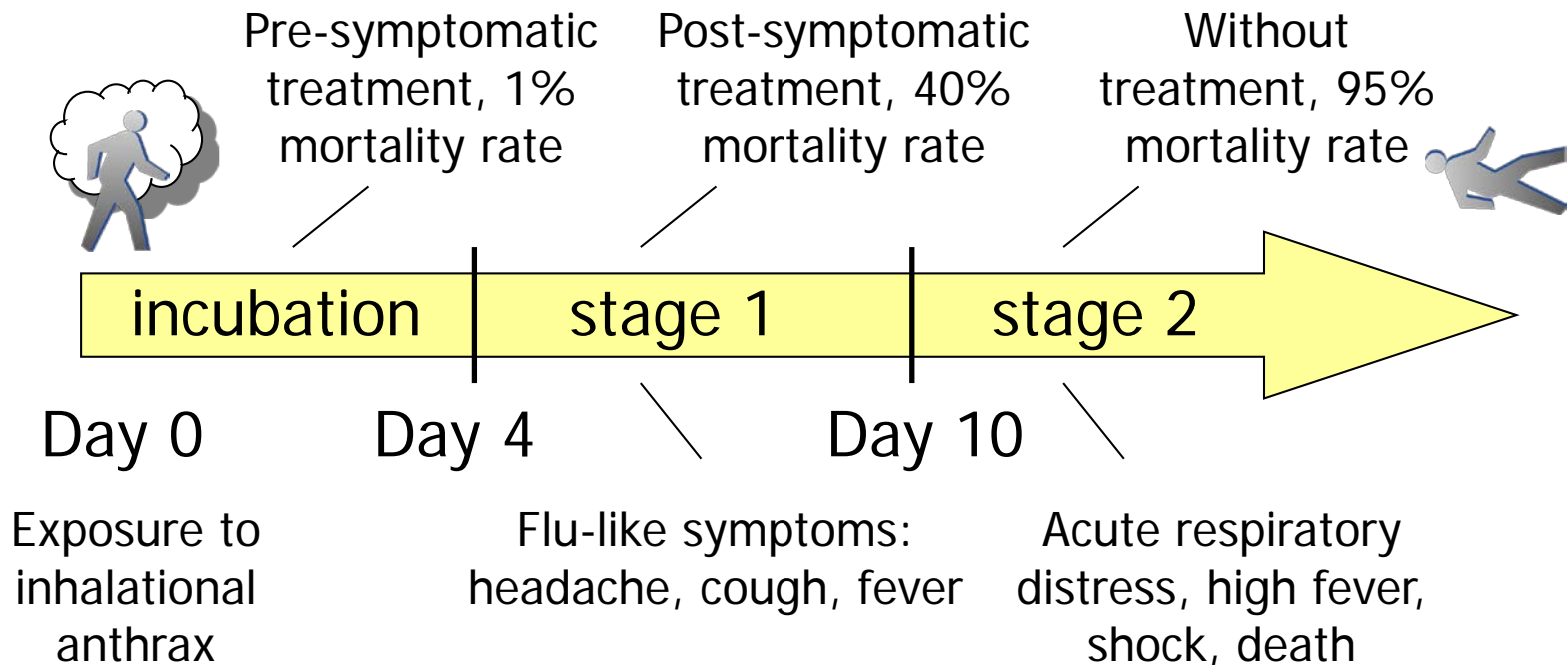
Why worry about disease outbreaks?

- Bioterrorist attacks are a very real, and scary, possibility
 - Large anthrax release over a major city could kill 1-3 million and hospitalize millions more.
- Emerging infectious diseases
 - “Conservative estimate” of 2-7 million deaths from pandemic avian influenza.
- Better response to common outbreaks and emerging public health trends.



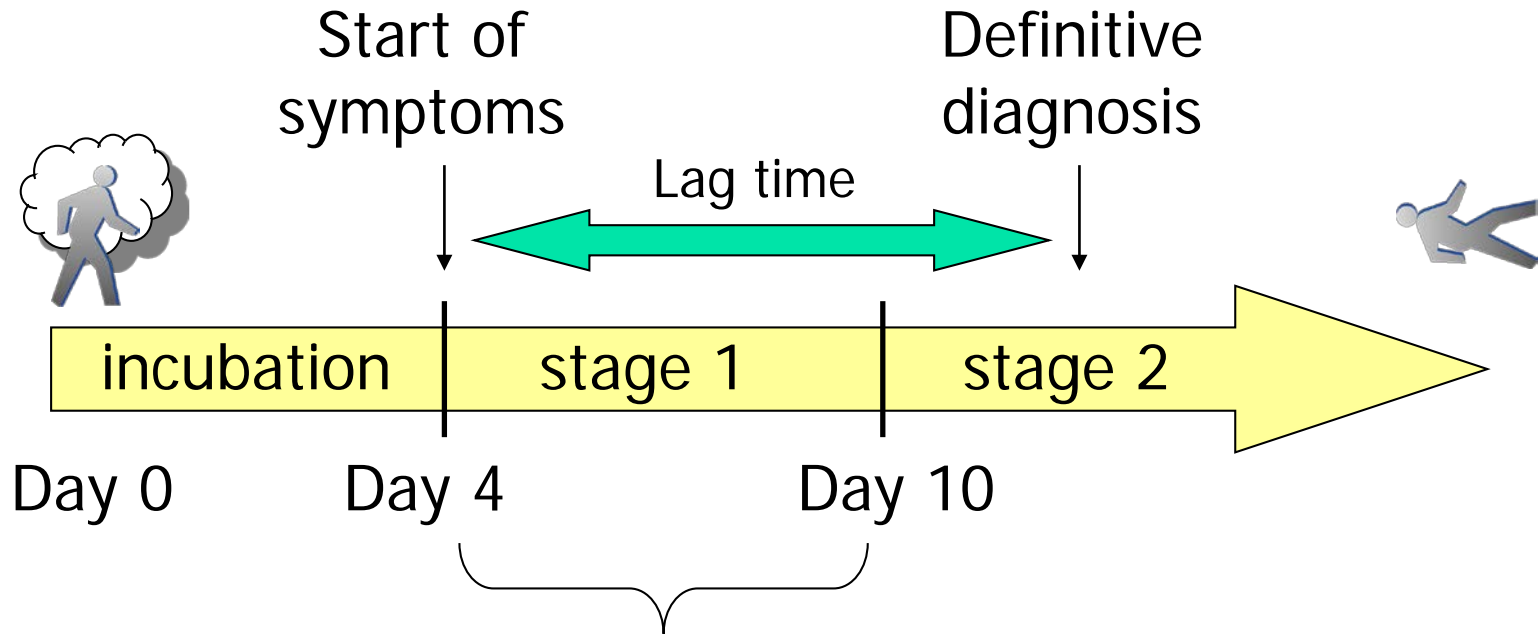
Benefits of early detection

Reduces **cost to society**, both in lives and in dollars!



DARPA estimate: a two-day gain in detection time and public health response could reduce fatalities by a factor of six.

Early detection is hard



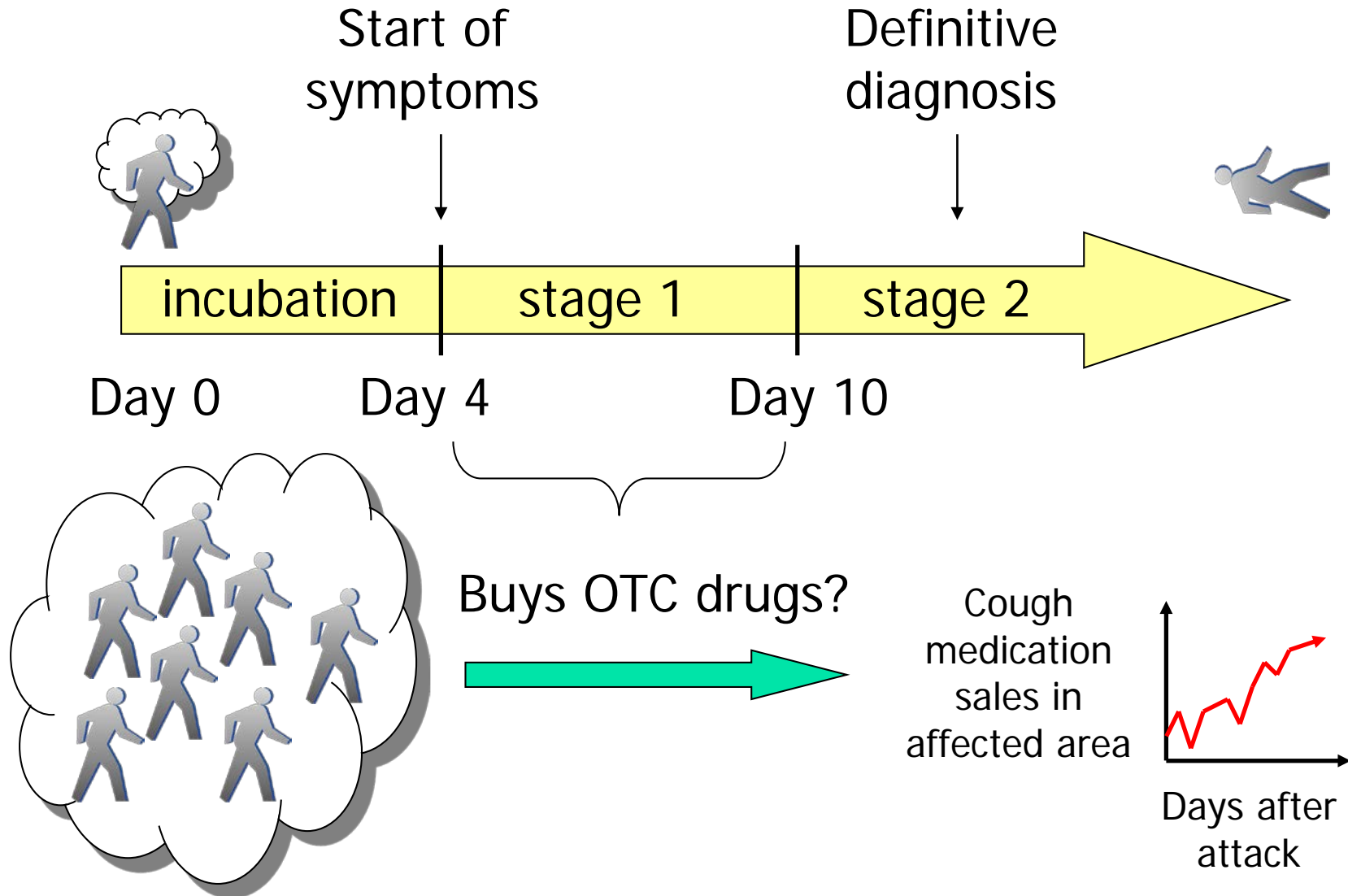
Buys OTC drugs

Skips work/school

Uses Google, Facebook, Twitter

Visits doctor/hospital/ED

Syndromic surveillance



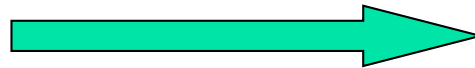
Syndromic surveillance

Start of
symptoms

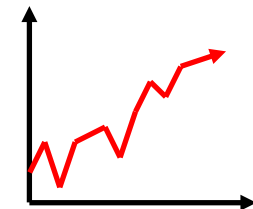
Definitive
diagnosis

We can achieve very early detection of outbreaks by gathering syndromic data, and identifying emerging spatial clusters of symptoms.

Buys OTC drugs?



Cough
medication
sales in
affected area

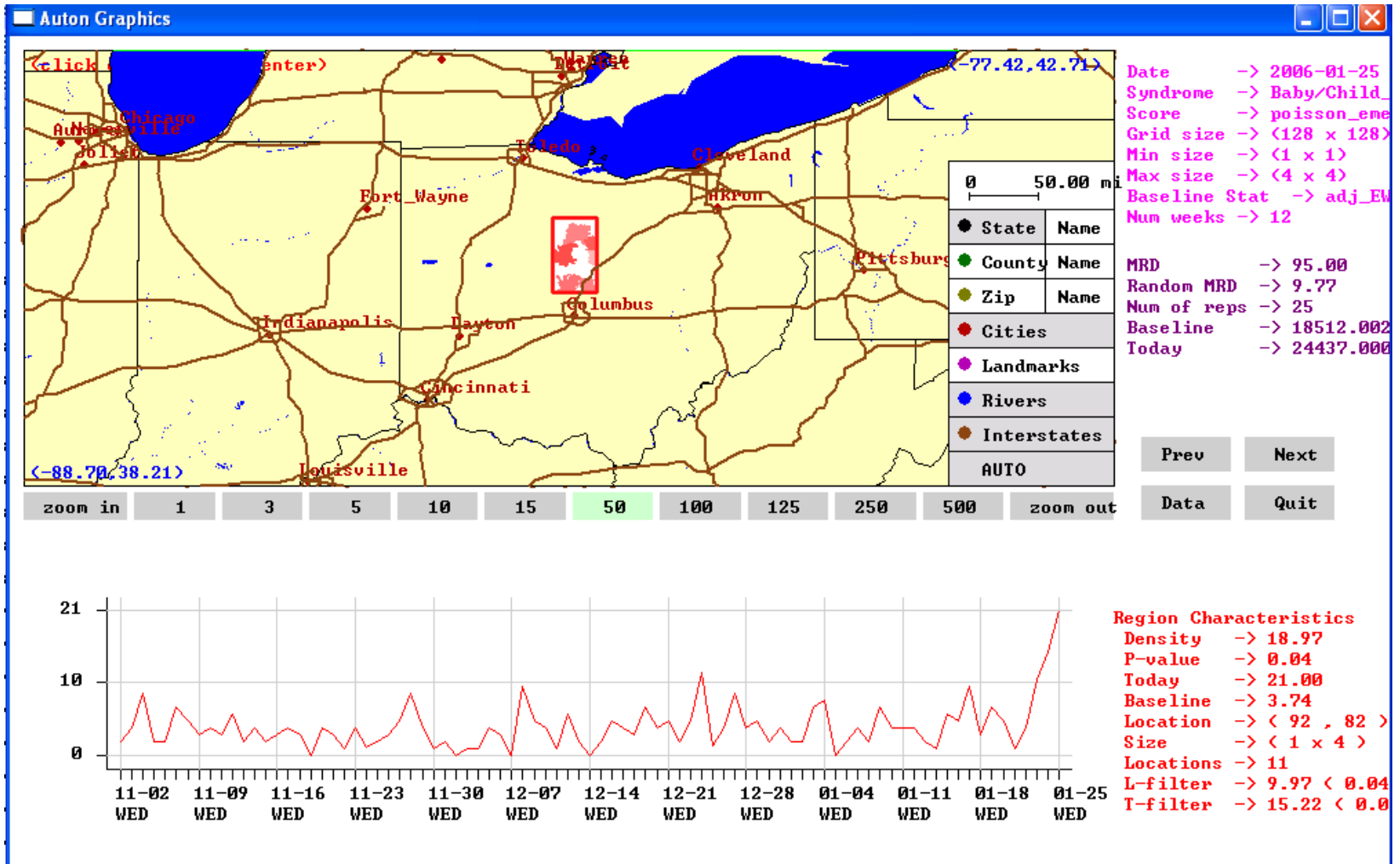


Days after
attack

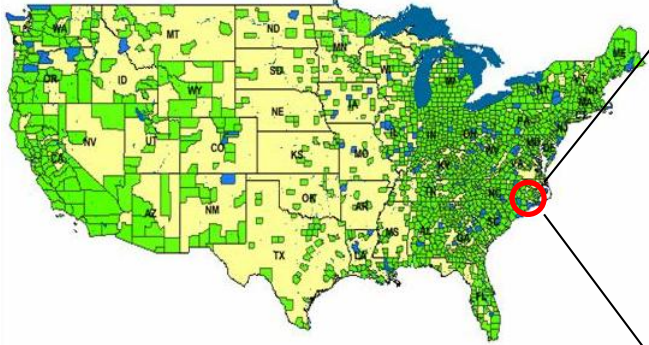


Outbreak detection example

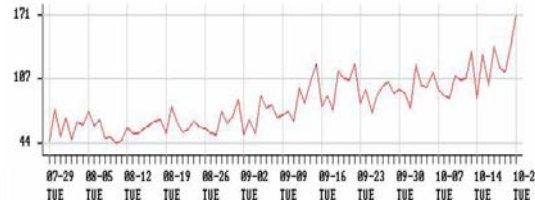
Spike in sales of pediatric electrolytes near Columbus, Ohio



Multivariate event detection



Spatial time series data from spatial locations s_i (e.g. zip codes)



Time series of counts $c_{i,m}^t$ for each zip code s_i for each data stream d_m .

Outbreak detection

- d_1 = respiratory ED
- d_2 = constitutional ED
- d_3 = OTC cough/cold
- d_4 = OTC anti-fever
(etc.)

Main goals:

- Detect** any emerging events.
- Pinpoint** the affected subset of locations and time duration.
- Characterize** the event, e.g., by identifying the affected streams.

Compare hypotheses:

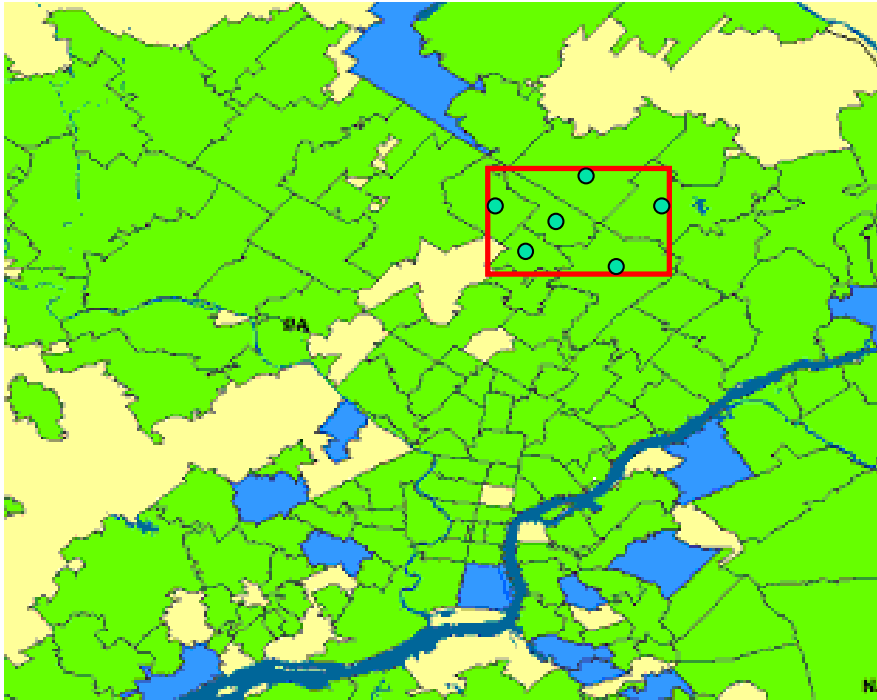
$$H_1(D, S, W)$$

- D = subset of streams
- S = subset of locations
- W = time duration

vs. H_0 : no events occurring

Expectation-based scan statistics

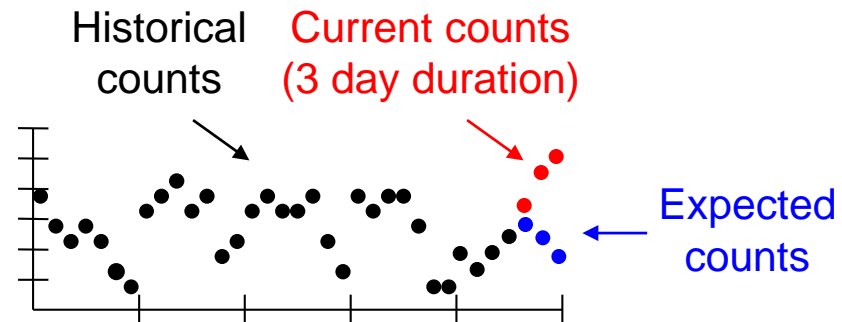
(Kulldorff, 1997; Neill and Moore, 2005)



We search for spatial regions (subsets of locations) where the recently observed counts for some subset of streams are significantly higher than expected.

We perform **time series analysis** to compute expected counts (“baselines”) for each location and stream for each recent day.

We then compare the actual and expected counts for each subset (D, S, W) under consideration.

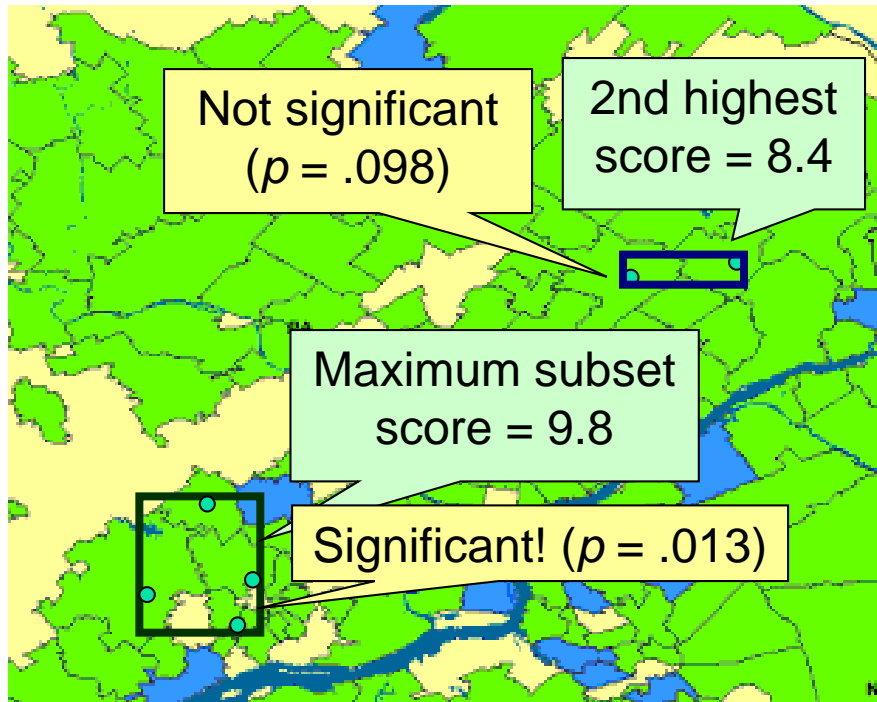


Expectation-based scan statistics

(Kulldorff, 1997; Neill and Moore, 2005)

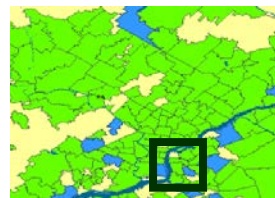
We find the subsets with highest values of a **likelihood ratio statistic**, and compute the p -value of each subset by randomization testing.

$$F(D, S, W) = \frac{\Pr(\text{Data} \mid H_1(D, S, W))}{\Pr(\text{Data} \mid H_0)}$$



To compute p-value
Compare subset score to maximum subset scores of simulated datasets under H_0 .

$F_1^* = 2.4$



$F_2^* = 9.1$



...

$F_{999}^* = 7.0$



Subset scanning for pattern detection

One key methodological idea of our work is **subset scanning**:

- We frame the pattern detection problem as a search over subsets of the data (e.g., spatial areas or subpopulations), maximizing some measure of the “interestingness” or “anomalousness” of a subset.
- This allows us to find **subtle patterns** which typical anomaly detection methods would miss (“needle in the haystack + connect the dots”)
- Once we have found the highest scoring subset, we perform a statistical test to see whether it is unlikely to have occurred by chance.

Search over subsets would be computationally infeasible but...

- Our **fast subset scan** algorithm can efficiently identify the most interesting subsets of data records **without** exhaustive search.

This enables us to solve detection problems in **milliseconds** that would previously have required **millions of years!**

Subset scanning for pattern detection

One key methodology for subset scanning:

- We can scan for patterns in data that are not necessarily contiguous in space or time.

Benefits for public health disease surveillance:

- More flexibility in defining cluster shape, leading to both improved spatial precision and higher power to detect subtle outbreaks.
- Ability to combine data from multiple streams (ED visits, OTC sales, online data, etc.)
- Analysis of massive and complex data sources such as Twitter, search queries, etc.

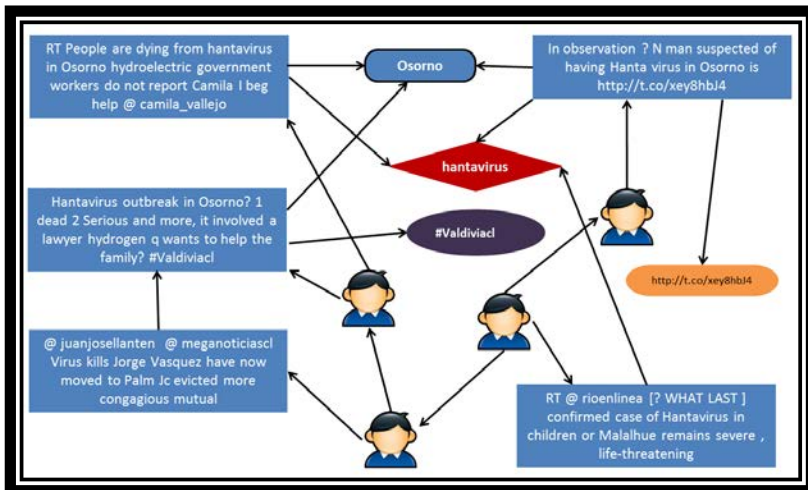
are but...

most search.

interest

This enables us to solve problems in **milliseconds** that would previously have required **millions of years!**

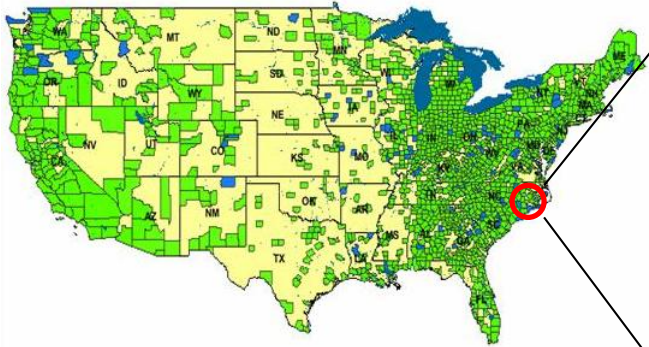
Detecting rare disease outbreaks with Twitter



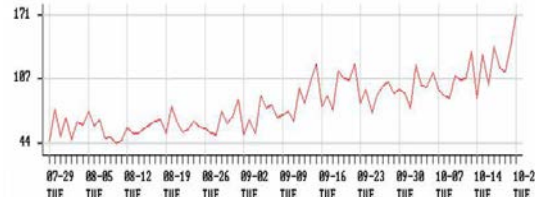
- Locations
- Users
- Keywords
- Hashtags
- Links
- Videos



Multivariate event detection



Spatial time series data from spatial locations s_i (e.g. zip codes)



Time series of counts $c_{i,m}^t$ for each zip code s_i for each data stream d_m .

Outbreak detection

- d_1 = respiratory ED
- d_2 = constitutional ED
- d_3 = OTC cough/cold
- d_4 = OTC anti-fever
(etc.)

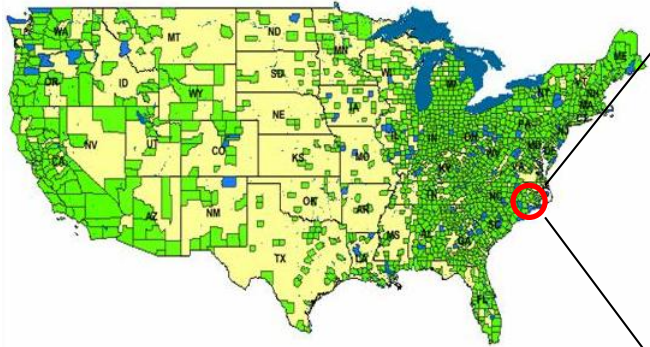
Main goals:

- Detect** any emerging events.
- Pinpoint** the affected subset of locations and time duration.
- Characterize** the event, e.g., by identifying the affected streams.

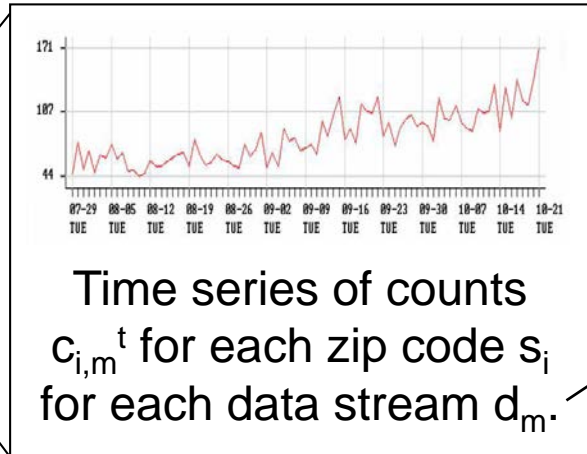
Compare hypotheses:

- $H_1(D, S, W)$
- D = subset of streams
- S = subset of locations
- W = time duration
- vs. H_0 : no events occurring

Multidimensional event detection



Spatial time series data from spatial locations s_i (e.g. zip codes)



Outbreak detection

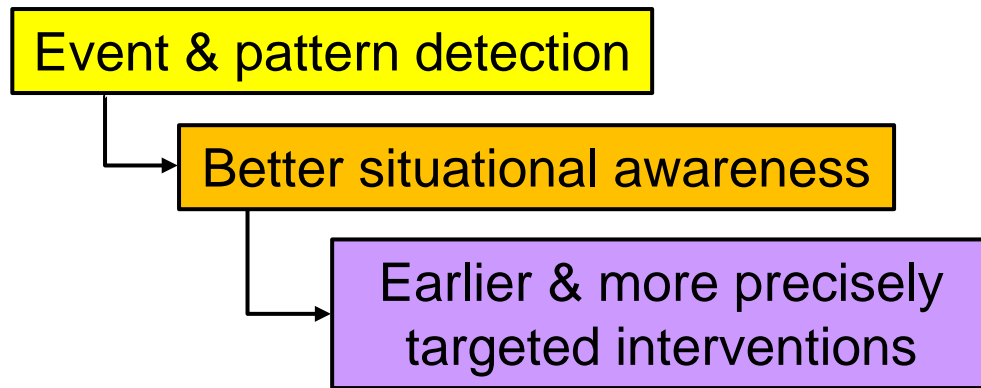
- d_1 = respiratory ED
- d_2 = constitutional ED
- d_3 = OTC cough/cold
- d_4 = OTC anti-fever
(etc.)

Additional goal: identify any differentially affected **subpopulations** P of the monitored population.

- Gender (male, female, both)
- Age groups (children, adults, elderly)
- Ethnic or socio-economic groups
- Risk behaviors: e.g. intravenous drug use, multiple sexual partners

More generally, assume that we have a set of additional discrete-valued attributes $A_1..A_J$ observed for each individual case. We identify not only the affected streams, locations, and time window, but also a **subset** of values for each attribute.

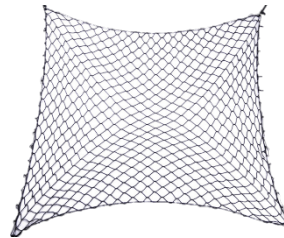
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities



Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis

Drug overdoses

- Drug overdoses are an increasingly serious problem in the United States and worldwide.
 - In 2017, more than 72,000 drug overdose deaths occurred in the U.S., more than any year in recorded history.
 - Approximately 68% of these overdose deaths involved opioids.
 - Economic costs of the crisis are estimated at \$78.5 billion annually.
- These statistics motivate public health to identify and predict emerging trends in overdoses, including geographic, demographic, and behavioral patterns, to better target interventions.
 - **Prevention** of high-risk prescribing and opioid use behaviors
 - **Treatment** of opioid addiction, e.g., medication-assisted therapy
 - **Rescue**, e.g., access to life-saving naloxone
 - **Recovery**, e.g., peer recovery coaches

Drug overdoses

- Drug overdoses are an increasingly serious problem in the United States and worldwide.
 - In 2017, more than 72,000 drug overdose deaths occurred in the U.S., more than any year in recorded history.
 - Approximately 68% of these overdose deaths involved opioids.
 - Economic costs of the crisis are estimated at \$78.5 billion annually.
- These statistics motivate public health to identify and predict emerging trends in overdoses, including geographic, demographic, and behavioral patterns, to better target interventions.
- Machine learning has potential to **save lives** by detecting subtle, emerging patterns of overdoses in their early stages and targeting an effective public health response.

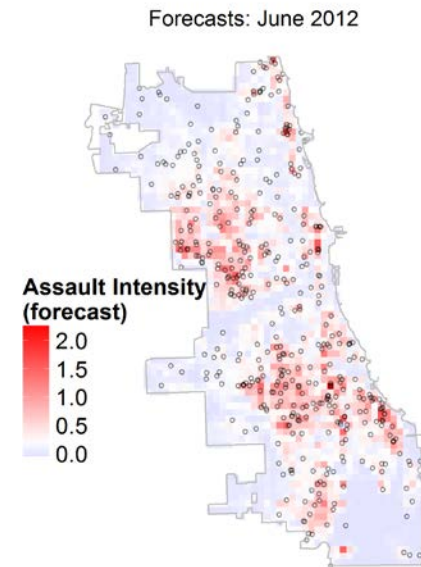
Geographic surveillance

- Answers the question, **where** should I intervene?
- Main goals: estimate predicted overdose trends in space and time; identify anomalous spikes in overdose deaths.

Useful predictors include neighborhood characteristics and recent spatio-temporal trends in overdoses and leading indicator variables (e.g., behavioral risk factors).

Gaussian processes are a useful approach for modeling correlated spatio-temporal data.

Our recent work* enables them to scale to real-world data, achieving state-of-the-art accuracy for long-term, small-area forecasting.



*SR Flaxman, AG Wilson, DB Neill, H Nickisch, AJ Smola. Fast Kronecker inference in Gaussian processes with non-Gaussian likelihoods. Proc. 32nd Intl. Conf. on Machine Learning, *PMLR* 37: 607-616, 2015.

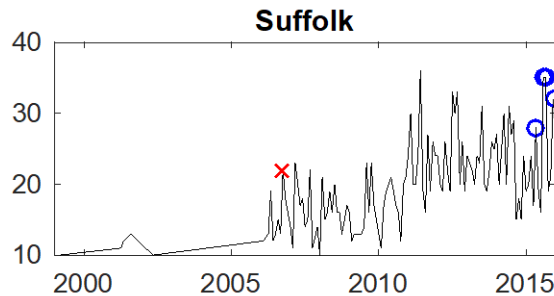
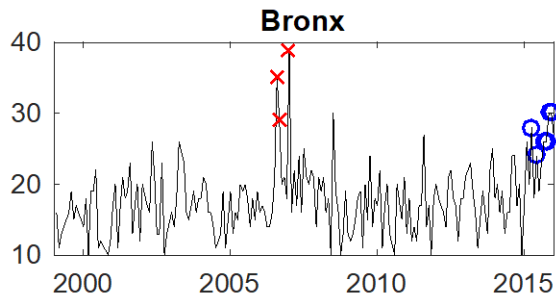
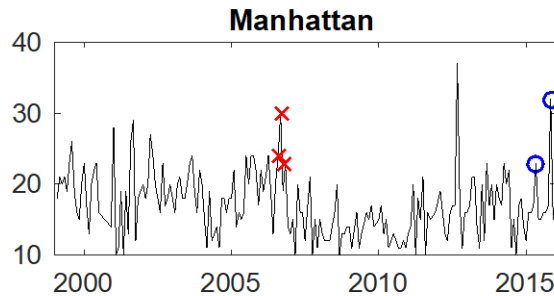
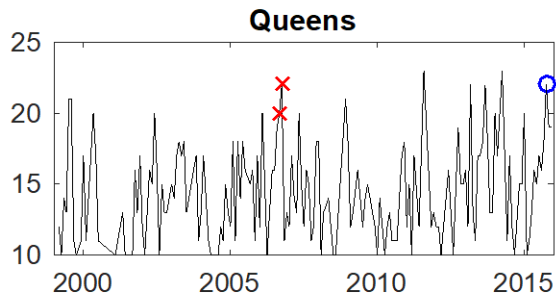
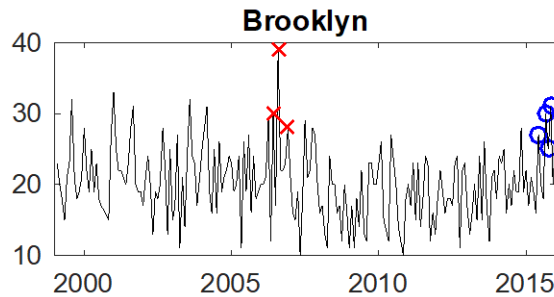
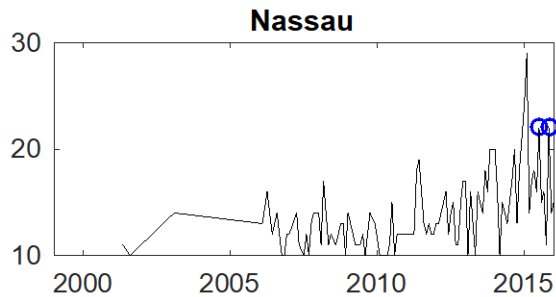
Case study: Geographic surveillance

- We analyzed **aggregate monthly counts** of fatal opioid overdoses for six New York counties from 1999-2015.
- We developed a new approach* which combines **Gaussian processes** (to model correlations) and **subset scan** (to identify the most anomalous space-time regions).
- We compared our new method to typical anomaly detection approaches on real and synthetic datasets.
 - GPSS > GP alone: nearby points matter for subtle anomalies
 - GPSS > SS alone: covariance structure matters for correlated data

*W Herlands, E McFowland III, AG Wilson, DB Neill. Gaussian process subset scanning for anomalous pattern detection in non-iid data. *Proc. 21st Intl. Conf. on Artificial Intelligence and Statistics, PMLR 84: 425-434, 2018.*

Case study: Geographic surveillance

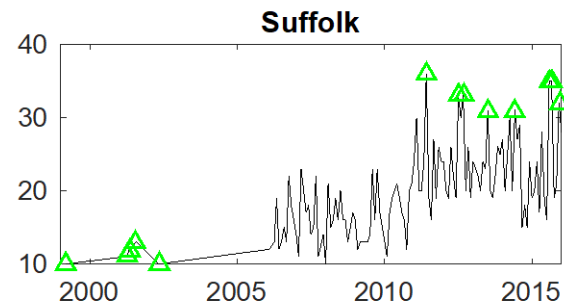
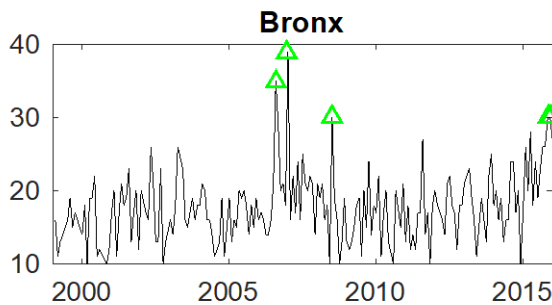
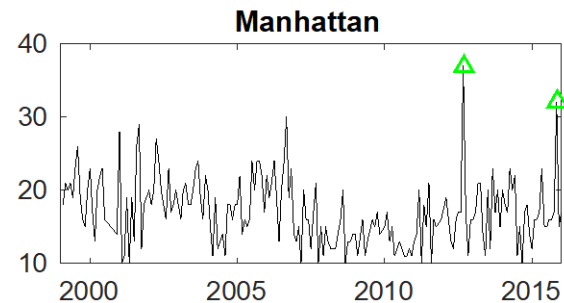
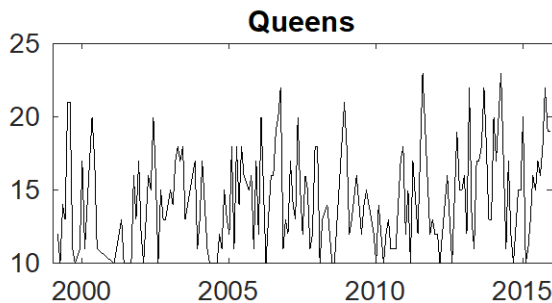
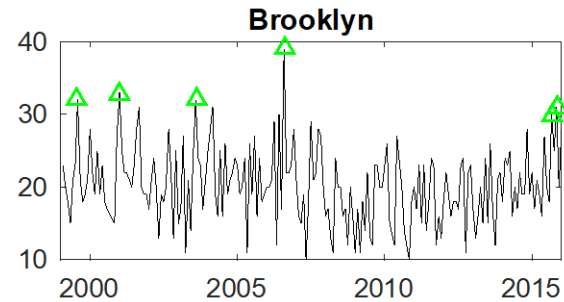
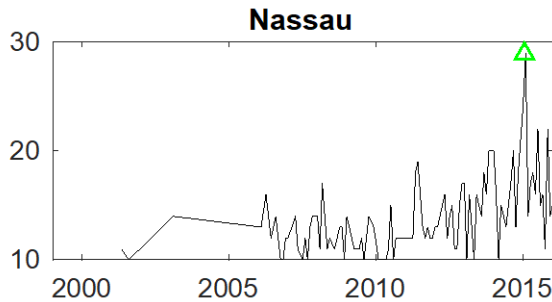
Two statistically significant spikes in overdose cases:



- ✘ Mid 2006. Just before naloxone programs.
- End of 2015. Recent surge due to fentanyl.

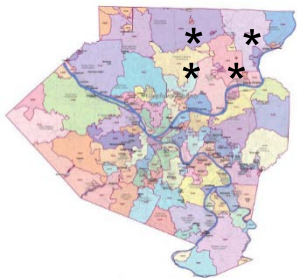
Case study: Geographic surveillance

Simpler anomaly detection methods fail to capture the relevant trends.



Subpopulation-level monitoring

- Answers the question, **for whom** should I intervene?
- Main goal: provide early warning for newly emerging subpopulation-level spikes/clusters of overdose deaths.
- We developed a novel detection method, **multi-dimensional tensor scan**, to detect emerging geographic, demographic, and behavioral patterns.
 - **Earlier detection** of emerging overdose clusters through daily surveillance runs.
 - Better characterization of **where** and **who** is affected.



X

white
males
aged
20-49

X



Multidimensional Tensor Scan

- In a nutshell: we identify **subspaces** of the attribute space (a subset of values for each attribute) with higher than expected numbers of recent case counts.
 - Spatial area (subset of locations) and time window
 - Affected genders, races, age ranges, and which drugs involved.
- We use a novel tensor decomposition approach to estimate how many counts we expect for each combination of attributes, while maintaining computational efficiency.
- Iterative conditional optimization: optimize over all subsets of values for each attribute conditional on the current subsets of values for all other attributes.
- Each conditional optimization step can be performed very efficiently, without exhaustive searching over subsets, by **fast subset scanning** (Neill, *J. Royal Stat. Soc. B*, 2012).

Overdoses in Allegheny County, PA

- We analyzed* county medical examiner data for fatal accidental drug overdoses, 2008-2015.
- ~2000 cases: for each overdose victim, we have date, location (zip), age, gender, race, and the set of drugs present in their system.
- Reduced to 30 dimensions (age decile, gender, race, presence/absence of 27 common drugs) plus space and time.
- Clusters discovered by MD-Scan were shared with Allegheny County Dept. of Human Services.

*DB Neill, W Herlands. Machine learning for drug overdose surveillance. *J. Technology in Human Services* 36(1): 8-14, 2018.

MD-Scan Overdose Results (1)



Fentanyl is a dangerous drug which has been a huge problem in western PA.

It is often mixed with white powder heroin, or sold disguised as heroin.

January 16-25, 2014:
14 deaths county-wide from fentanyl-laced heroin.

March 27 to April 21, 2015:
26 deaths county-wide from fentanyl, heroin only present in 11.

January 10 to February 7, 2015:

Cluster of 11 fentanyl-related deaths, mainly black males over 58 years of age, centered in Pittsburgh's downtown Hill District.

Very unusual demographic: common dealer / shooting gallery?

Started in the southeast suburbs of Pittsburgh and spread across the city.

Our method could have detected this pattern on **March 29**, identifying a cluster of four overdose deaths with strong geographic and demographic similarities.

Fentanyl, heroin, and combined deaths remained high through end of June (>100).

MD-Scan Overdose Results (2)

Another set of discovered overdose clusters each involved a combination of Methadone and Xanax.



The combination produces a strong high but can be deadly (~30% of methadone fatal ODs).

From 2008-2012: multiple M&X OD clusters, 3-7 cases each, localized in space and time.

Methadone: an opioid used for chronic pain relief and to treat heroin addiction, but also addictive and risk of OD.

From 2013-2015: no M&X overdose clusters; 33% and 47% drops in yearly methadone and M&X deaths respectively.



Xanax (alprazolam): a benzodiazepine prescribed for panic and anxiety disorders.

Why did these deaths cluster, when methadone and methadone + other benzo deaths did not?

What factors could explain the dramatic reduction in M&X overdose clusters?

MD-Scan Overdose Results (2)

Another set of discovered overdose clusters each involved a combination of Methadone and Xanax.



Methadone: an opioid used for chronic pain relief and to treat heroin addiction, but also addictive and risk of OD.



Xanax (alprazolam): a benzodiazepine prescribed for panic and anxiety disorders.

Increased state oversight of methadone clinics and prescribing physicians after passage of the Methadone Death and Incident Review Act (Oct 2012).

Approval of generic suboxone (buprenorphine + naloxone) in early 2013 lowered cost of suboxone treatment as an alternative to methadone clinics.

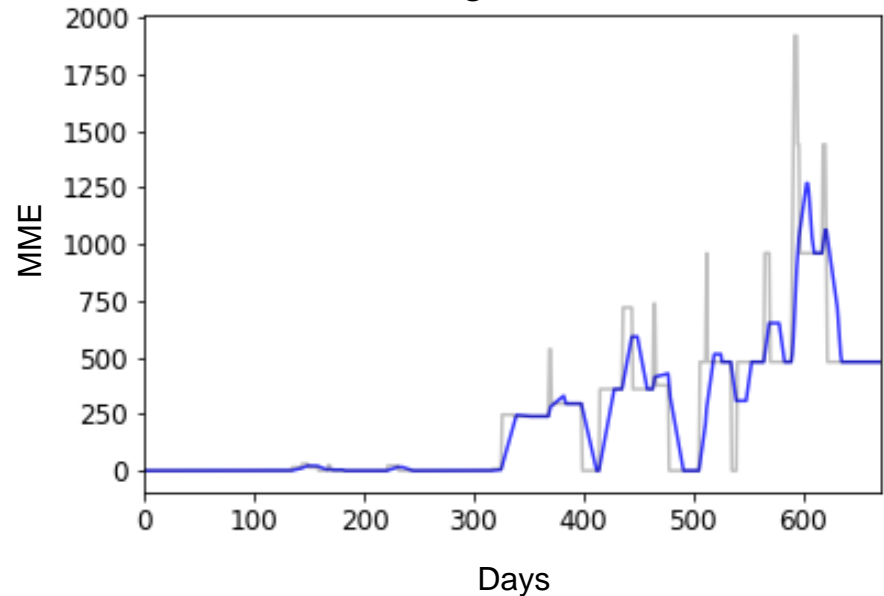
Why did these deaths cluster, when methadone and methadone + other benzo deaths did not?

What factors could explain the dramatic reduction in M&X overdose clusters?

Individual-level opioid use monitoring

- Seven years of de-identified data from over 1M individuals provided by Kansas prescription drug monitoring program (PDMP), with unique patient, prescriber, and dispensary identifiers.
- Duration and quantity of prescribed opioids are used to create timelines of morphine milligram equivalents (MME) for individual patients.
- Can we identify **early indicators** in patient MME timelines which are predictive of later opioid misuse or unsafe prescribing?

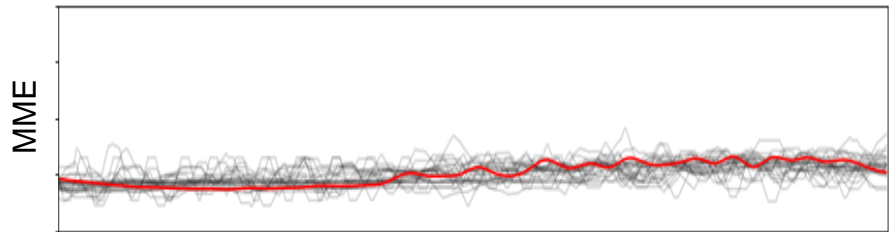
Smoothed MME Timeline for a Single Patient



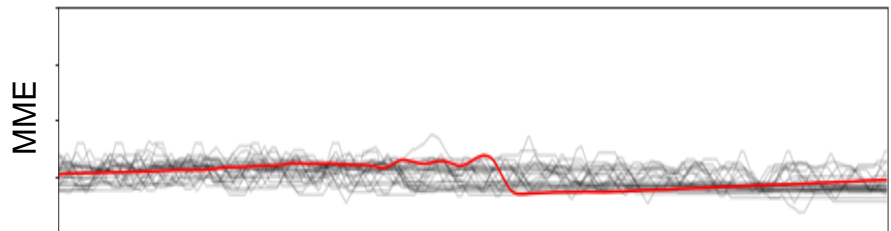
Individual-level opioid use monitoring

- Patients are clustered using the *k*-shape algorithm (Paparrizos & Gravano, 2015) to group patients with similar patterns in MME timelines.
- Are some patient clusters associated with higher risk of red flags indicating misuse or unsafe practices?
- For a new patient, can we confidently assess risk of future red flags given a partial MME timeline?

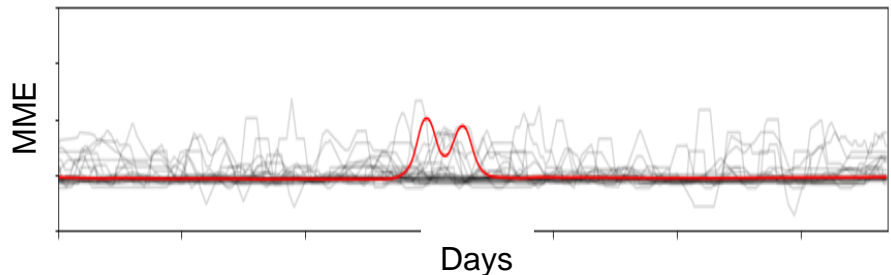
Cluster 5 of 10



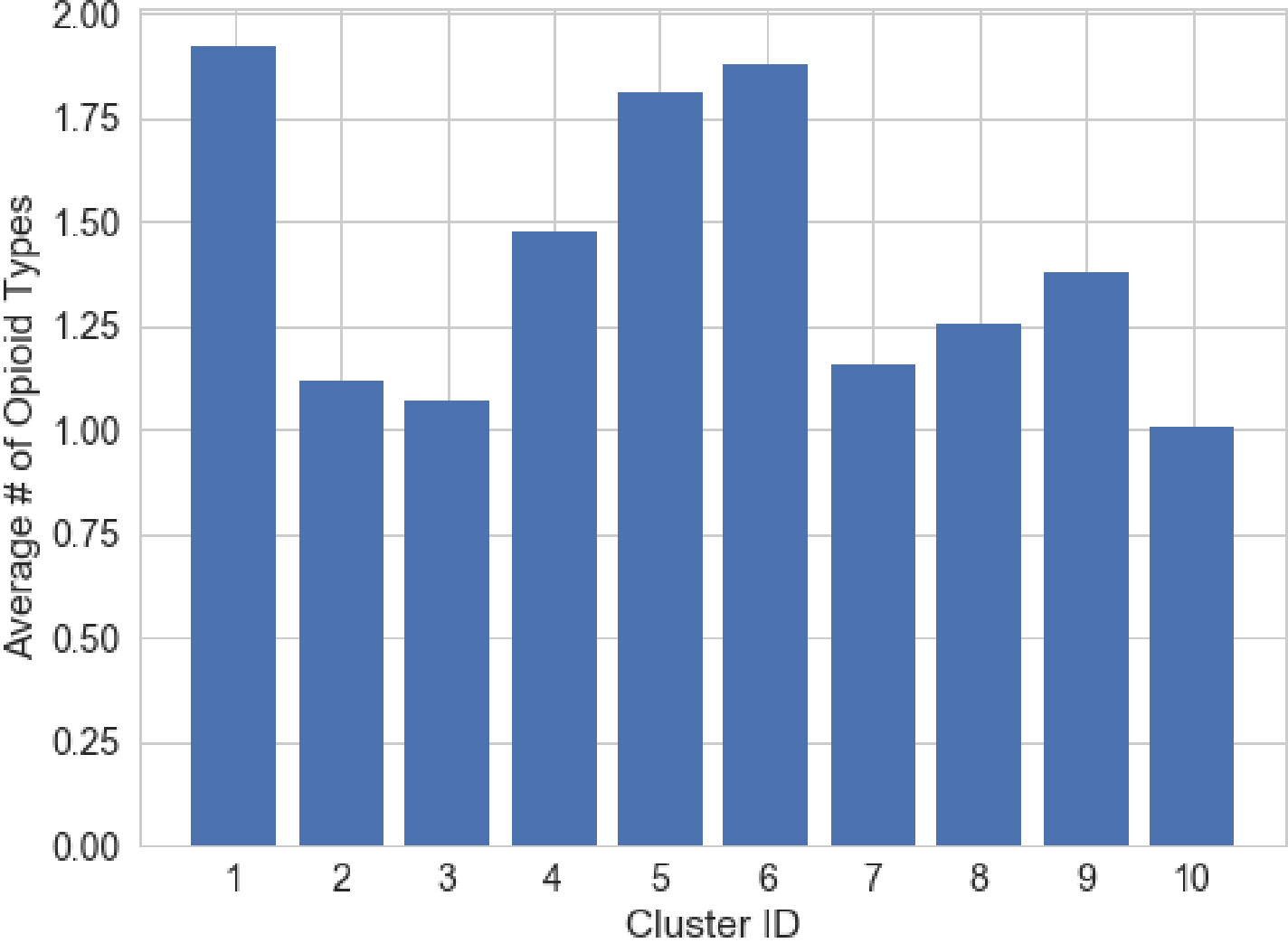
Cluster 6 of 10



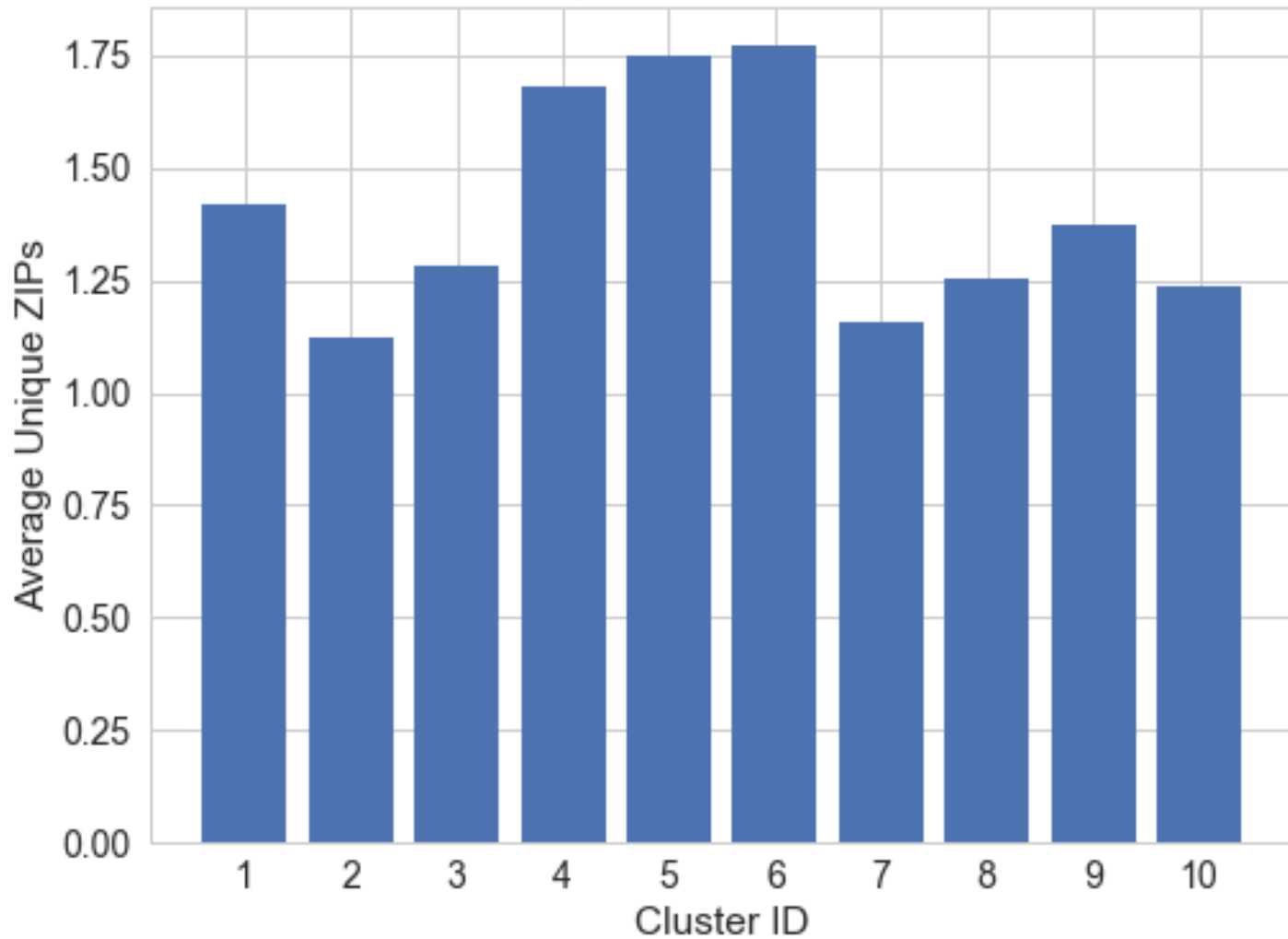
Cluster 9 of 10



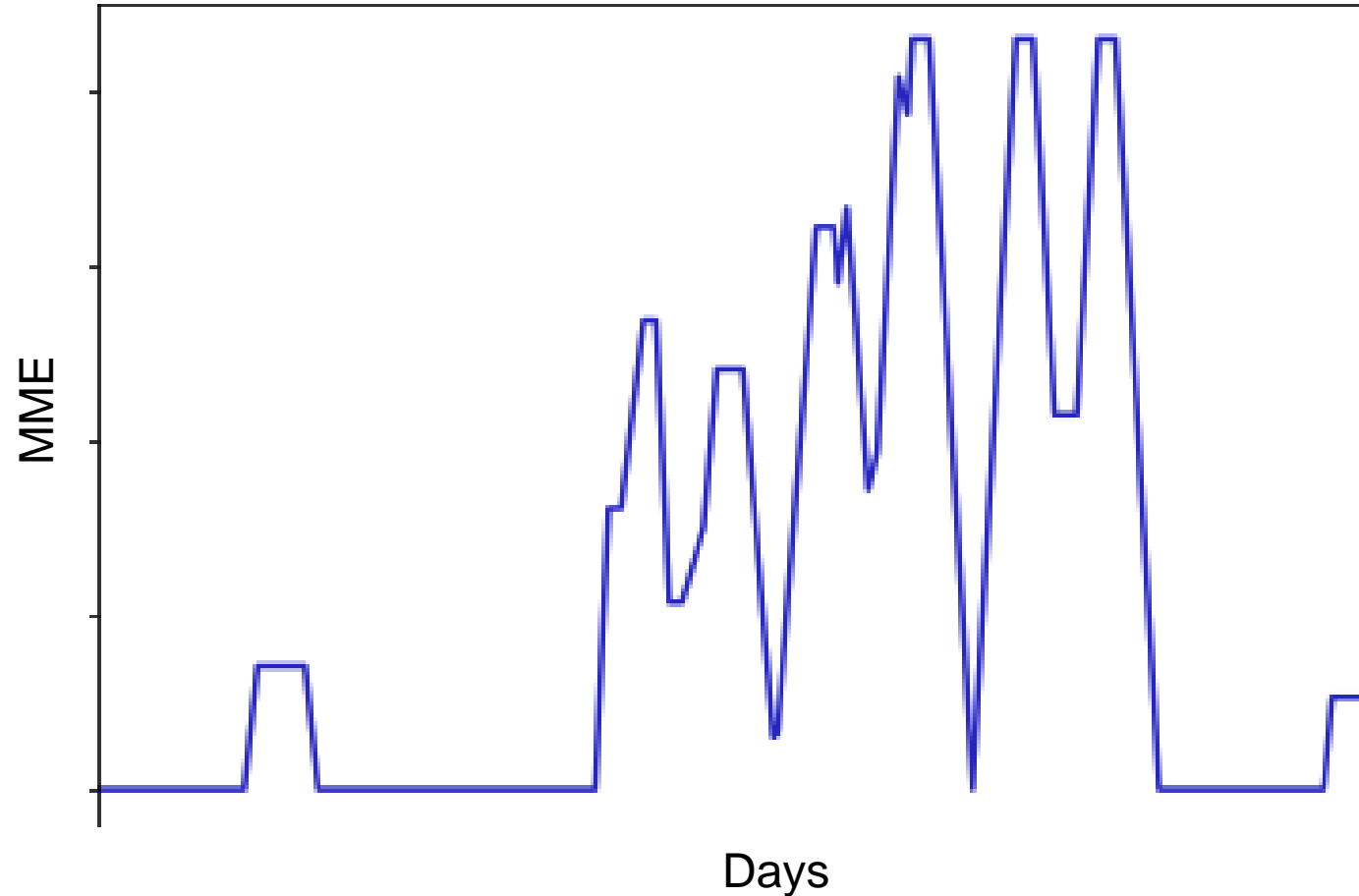
Red Flag Analysis by Cluster:
of Unique Opioid Types



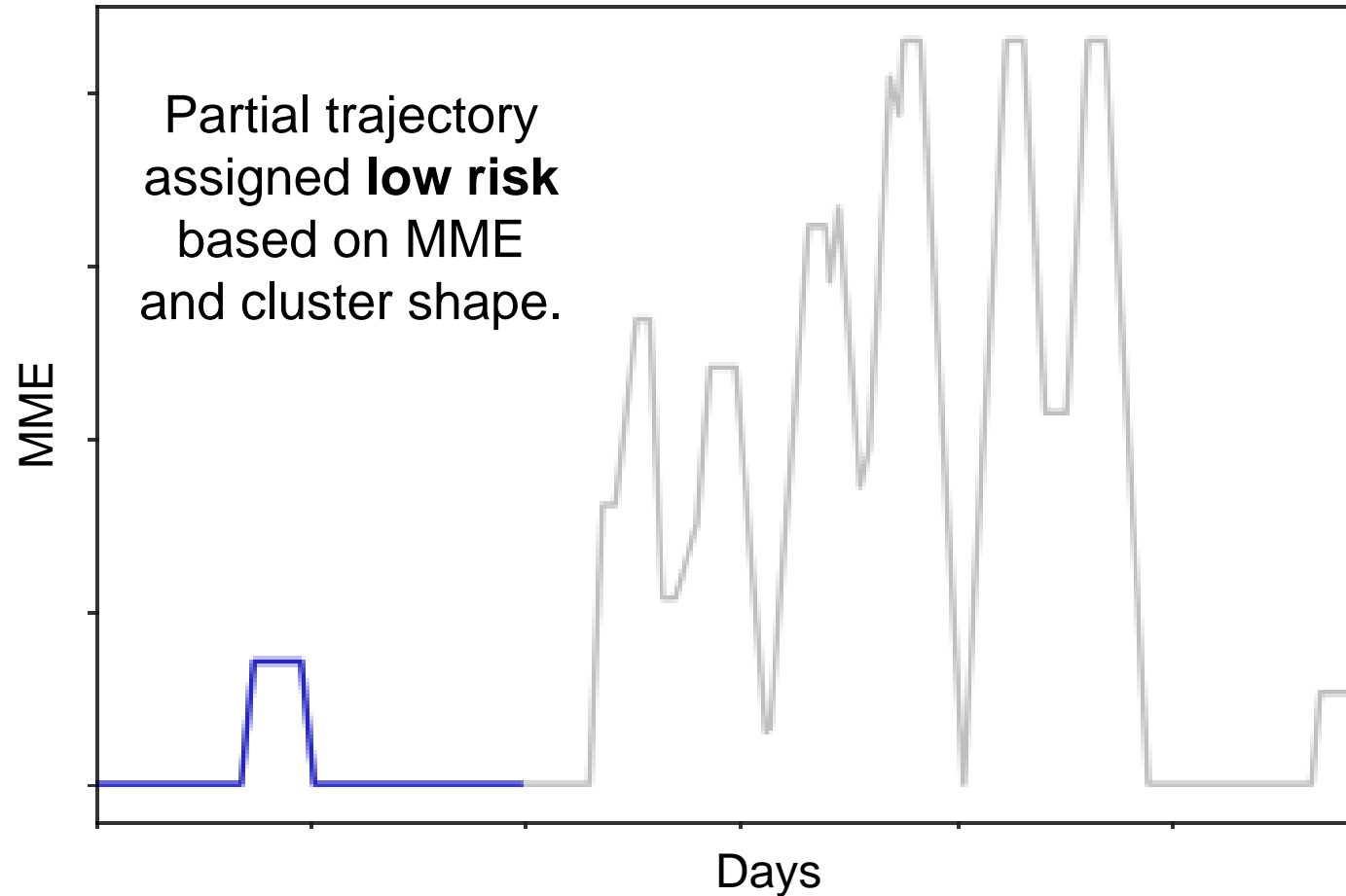
Red Flag Analysis by Cluster:
of Unique Prescriber Zip Codes



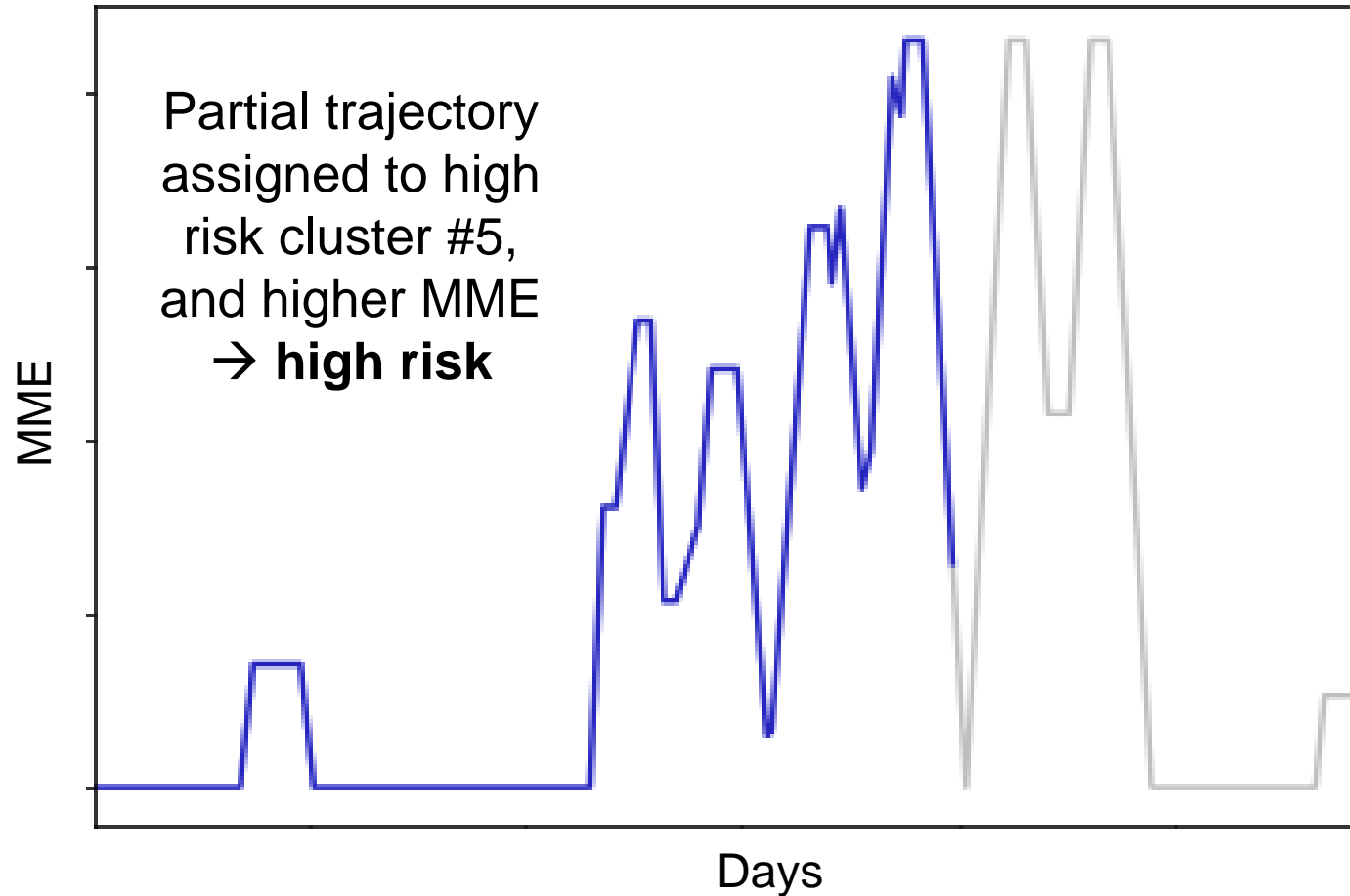
Early individual-level risk assessment by classifying partial trajectories



Early individual-level risk assessment by classifying partial trajectories



Early individual-level risk assessment by classifying partial trajectories



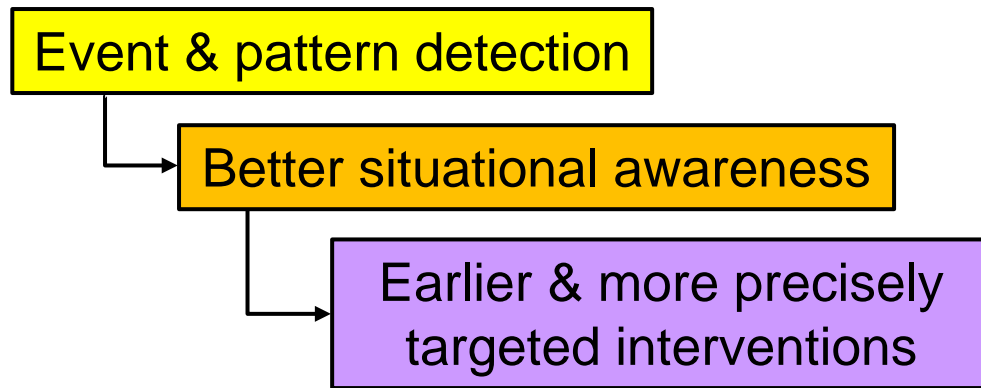
Discussion

Here we described several new methods that can be used for **early warning** and **advance forecasting** of overdoses at the geographic, subpopulation, and individual levels.

Our retrospective analyses of overdose and opioid use data from Pennsylvania, New York, and Kansas suggest high potential utility for **prospective** drug overdose surveillance systems, to facilitate targeted and effective interventions.

We are currently collaborating with an interdisciplinary team of investigators and public health practitioners, with the goals of deploying targeted interventions to prevent overdoses and evaluating their effectiveness through randomized trials.

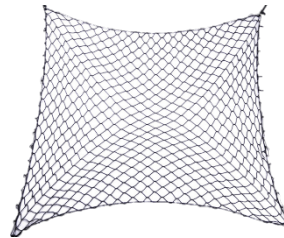
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities



Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis

Pre-syndromic surveillance

<u>Date/time</u>	<u>Hosp.</u>	<u>Age</u>	<u>Complaint</u>
Jan 1 08:00	A	19-24	runny nose
Jan 1 08:15	B	10-14	fever, chills
Jan 1 08:16	A	0-1	broken arm
Jan 2 08:20	C	65+	vomited 3x
Jan 2 08:22	A	45-64	high temp

Key challenge: A syndrome cannot be created to identify every possible cluster of potential public health significance.

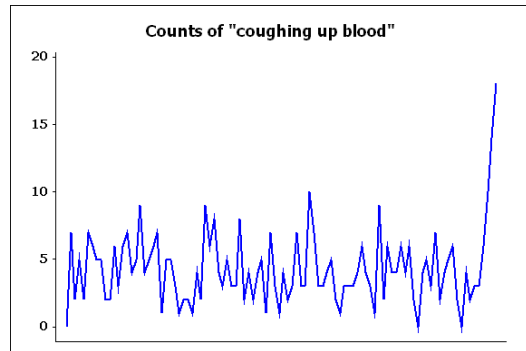
Thus a method is needed to identify relevant clusters of disease cases that do not correspond to existing syndromes.

Use case proposed by NC DOH and NYC DOHMH, solution requirements developed through a public health consultancy at the International Society for Disease Surveillance.

Where do existing methods fail?

The typical syndromic surveillance approach can effectively detect emerging outbreaks with commonly seen, general patterns of symptoms (e.g. ILI).

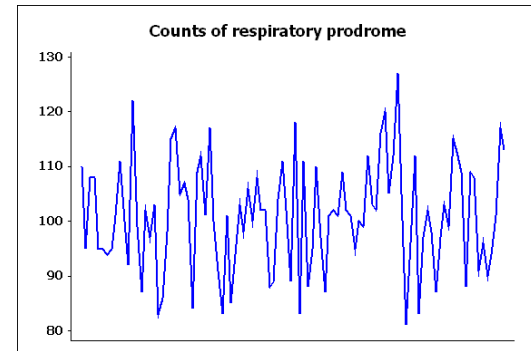
If we were monitoring these particular symptoms, it would only take a few such cases to realize that an outbreak is occurring!



What happens when something new and scary comes along?

- **More specific symptoms** ("coughing up blood")
- **Previously unseen symptoms** ("nose falls off")

Mapping specific chief complaints to a broader symptom category can dilute the outbreak signal, delaying or preventing detection.

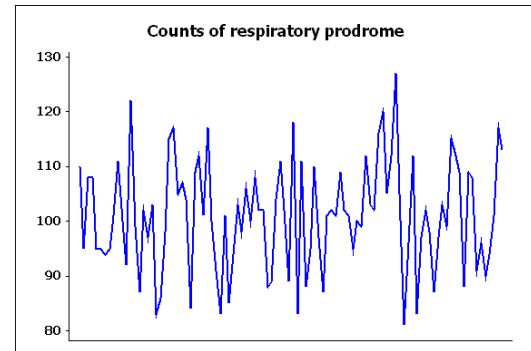
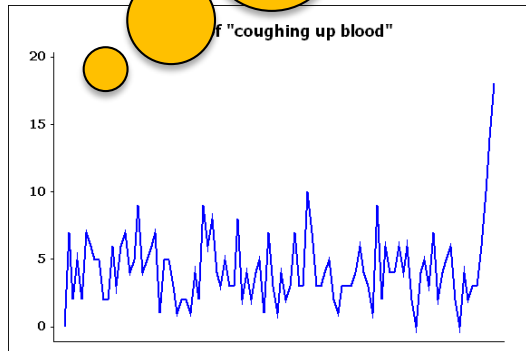


Where do existing methods fail?

The typical surveillance systems are not designed to detect something new along? symptoms (d) on off)

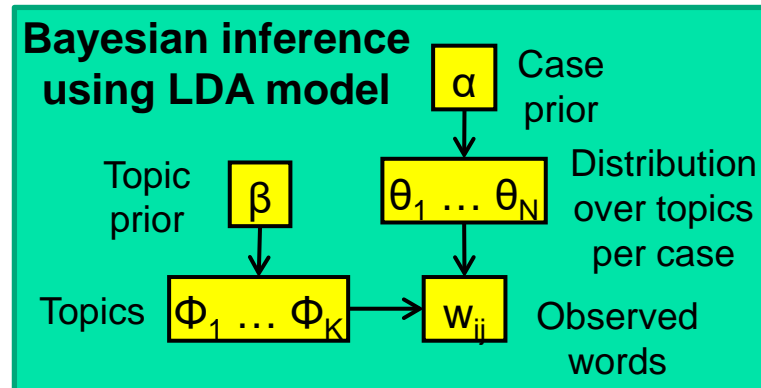
Our solution is to combine text-based (topic modeling) and event detection (multidimensional scan) approaches, to detect **emerging patterns of keywords.**

If we want to detect a particular symptom category, take a few such symptoms that an outbreak is occurring! ... symptoms or preventing detection.



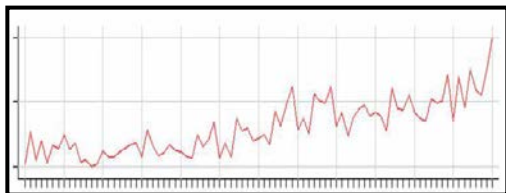
The semantic scan statistic

<u>Date/time</u>	<u>Hosp.</u>	<u>Age</u>	<u>Complaint</u>
Jan 1 08:00	A	19-24	runny nose
Jan 1 08:15	B	10-14	fever, chills
Jan 1 08:16	A	0-1	broken arm
Jan 2 08:20	C	65+	vomited 3x
Jan 2 08:22	A	45-64	high temp



ϕ_1 : vomiting, nausea, diarrhea, ...
 ϕ_2 : dizzy, lightheaded, weak, ...
 ϕ_3 : cough, throat, sore, ...

Classify cases to topics



Time series of hourly counts for each combination of hospital and age group, for each topic ϕ_j .

Now we can do a multidimensional scan, using the learned topics instead of pre-specified syndromes!

Multidimensional scanning

We consider subsets S that are a combination of a topic, time duration, set of hospitals, and age range. For each hour of data and each subset S , we compute:

Count: $C(S)$ = # of cases in that time interval matching on hospital, age range, and topic.

Baseline: $B(S)$ = expected count (28-day moving average).

Score: $F(S) = \log (\text{Pr}(\text{Data} \mid H_1(S)) / \text{Pr}(\text{Data} \mid H_0))$
 $= C \log (C/B) + B - C$, if $C > B$, or 0 otherwise

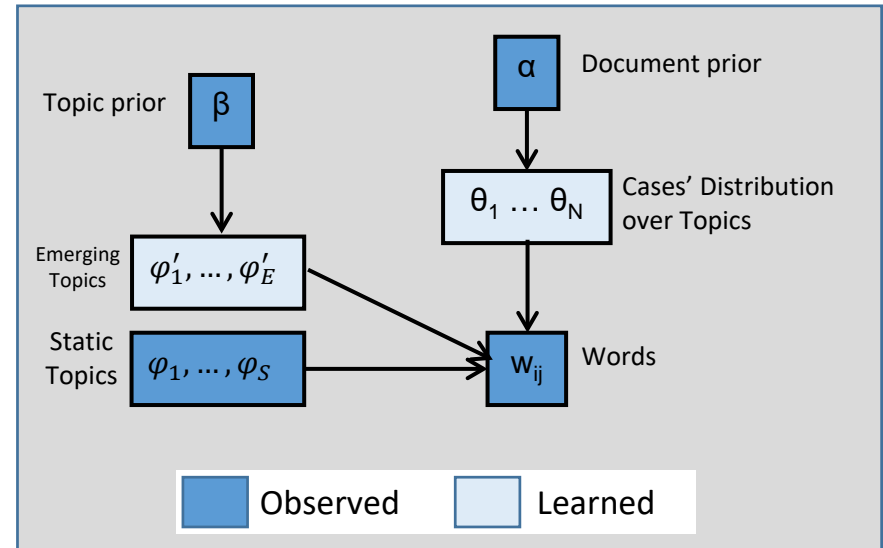
[This is the log-likelihood ratio using the **expectation-based Poisson** scan statistic, which assumes $c_i \sim \text{Poisson}(b_i)$ under H_0 , and a multiplicative increase under H_1 .]

We return cases corresponding to each top-scoring subset S .

Multidimensional Semantic Scan

Learns Two Sets of Topics

- Static Topics
 - Designed to capture common illnesses like flu.
 - Learned over a large set of historical data using a standard LDA topic model.
- Emerging Topics
 - Designed to capture rare or novel diseases that are not well explained by the static topics.
 - Learned over the most recent set of data using a new variant of LDA.



NYC DOHMH dataset

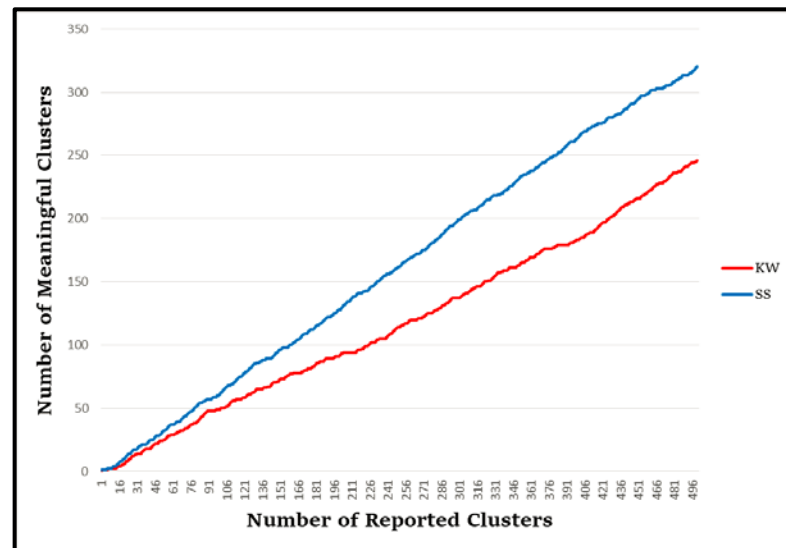
- New York City's Department of Health and Mental Hygiene, Bureau of Communicable Disease, provided us with 6 years of data (2010-2016) consisting of ~28M chief complaint cases from 53 hospitals in NYC.
- For each case, we have data on the patient's chief complaint (free text), date and time of arrival, age group, gender, and discharge ICD-9 code.
- Substantial pre-processing of the chief complaint field was necessary because of the size and messiness of the data (typos, abbreviations, etc.).

VOIMITING	VOMITINIG	VOMITINGN
VOIMITTING	VOMITINNG	VOMITINGQ
VOIMTING	VOMITIONG	VOMITINGS
VOMIITING	VOMITITING	VOMITINGT
VOMIITNG	VOMITITNG	VOMITINGX
VOMINITING	VOMITN	VOMITINGX1
VOMINTING	VOMITNG	VOMITINGX2
VOMIOTING	VOMITNIG	VOMITINGX3
VOMITE	VOMITNING	VOMITINGX4
VOMITED	VOMITO	VOMMITTING
VOMITG	VOMITOS	VOMNITING
VOMITHING	VOMITS	VOMOITING
VOMITI	VOMITT	VOMTIING
VOMITIG	VOMITTE	VOMTIN
VOMITIGN	VOMITTI	VOMTITING
VOMITIING	VOMITTING	VONMITING
VOMITIN	VOMITTING	VOOMITING
VOMITING3	VOMITUS	VOPMITING
VOMITINGA	VOMMIT	VVOMITING
VOMITINGG	VOMMITING	VOMITINGM

Variations of the words "vomit" and "vomiting" that appear > 15 times in data

Evaluation on NYC DOHMH data

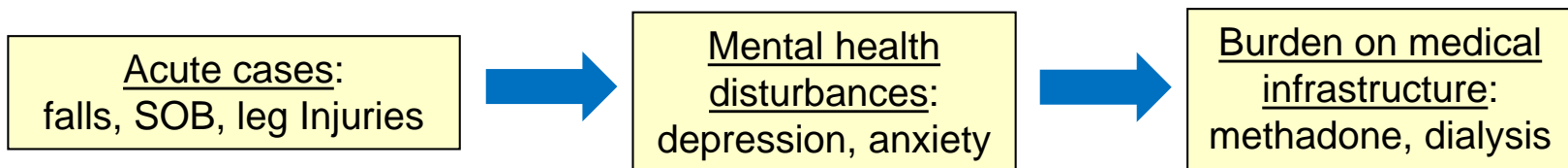
- Blinded evaluation by NYC DOHMH public health practitioners, comparing our multidimensional semantic scan approach to a state-of-the-art keyword-based scan approach.
- For each method's 500 highest scoring clusters, users indicated if the cluster is relevant, meaningful, or not of interest.



	Relevant Clusters of Interest	Meaningful Clusters of Potential Interest	Clusters Not of Interest
	Examples: bacterial meningitis, synthetic drugs use	Examples: flu, rashes, motor vehicle accidents	Examples: misspellings, non-specific words (i.e. "left")
Multidimensional Semantic Scan	53	267	180
Keyword Based Method	47	199	254

Events identified by semantic scan

The progression of detected clusters after Hurricane Sandy impacted NYC highlights the variety of strains placed on hospital emergency departments following a natural disaster:



Many other events of public health interest were identified:

Accidents
Motor vehicle
Ferry
School bus
Elevator

Contagious Diseases
Meningitis
Scabies
Ringworm
Hepatitis

Other
Drug overdoses
Smoke inhalation
Carbon monoxide poisoning
Crime related, e.g., pepper spray attacks

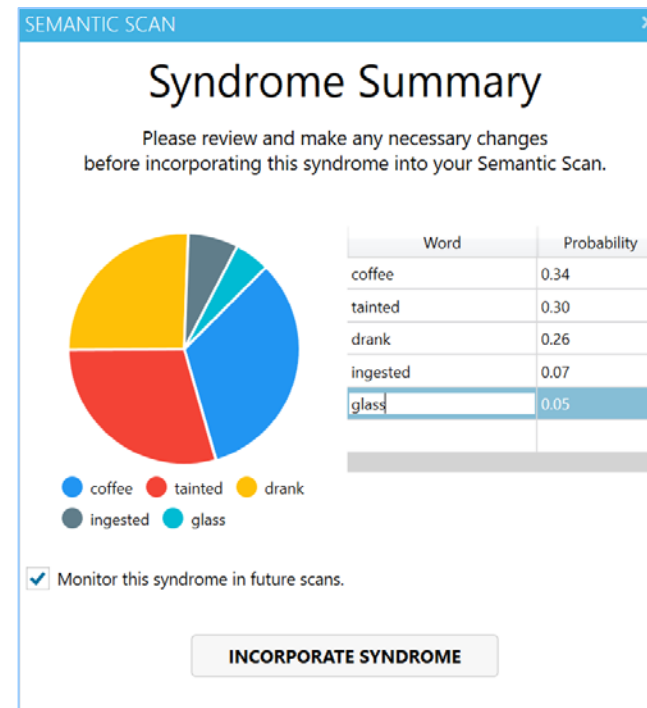
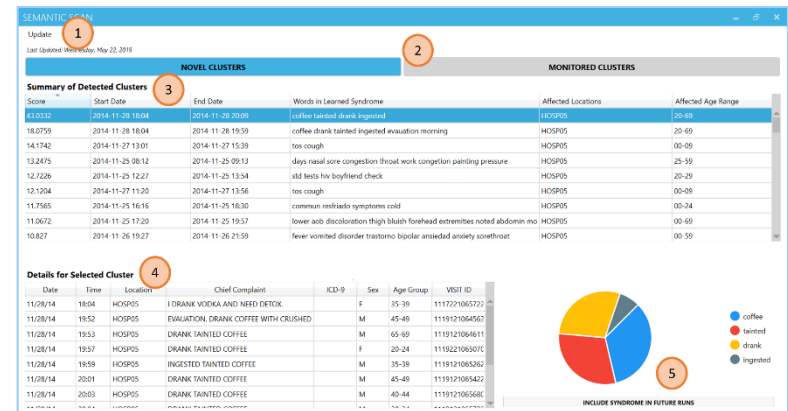
Example of a detected cluster

Arrival Date	Arrival Time	Hospital ID	Chief Complaint	Patient Sex	Patient Age
11/28/2014	7:52:00	HOSP5	EVAUATION, DRANK COFFEE WITH CRUS	M	45-49
11/28/2014	7:53:00	HOSP5	DRANK TAIANTED COFFEE	M	65-69
11/28/2014	7:57:00	HOSP5	DRANK TAIANTED COFFEE	F	20-24
11/28/2014	7:59:00	HOSP5	INGESTED TAIANTED COFFEE	M	35-39
11/28/2014	8:01:00	HOSP5	DRANK TAIANTED COFFEE	M	45-49
11/28/2014	8:03:00	HOSP5	DRANK TAIANTED COFFEE	M	40-44
11/28/2014	8:04:00	HOSP5	DRANK TAIANTED COFFEE	M	30-34
11/28/2014	8:06:00	HOSP5	DRANK TAIANTED COFFEE	M	35-39
11/28/2014	8:09:00	HOSP5	INGESTED TAIANTED COFFEE	M	25-29

This detected cluster represents 9 patients complaining of ingesting tainted coffee, and demonstrates Semantic Scan's ability to detect rare and novel events.

Incorporating user feedback

- Our system enables continual improvement of performance by including public health practitioners in the loop and incorporating their feedback.
- Our visualization interface enables users to add new syndromes (topics) and specify if they would like the system to monitor them in the future.
 - Yes → clusters corresponding to that topic appear in a separate “monitored clusters” tab.
 - No → clusters corresponding to that topic are not reported.
 - In either case, the “novel clusters” tab includes only clusters that do not correspond to any static or added topic.

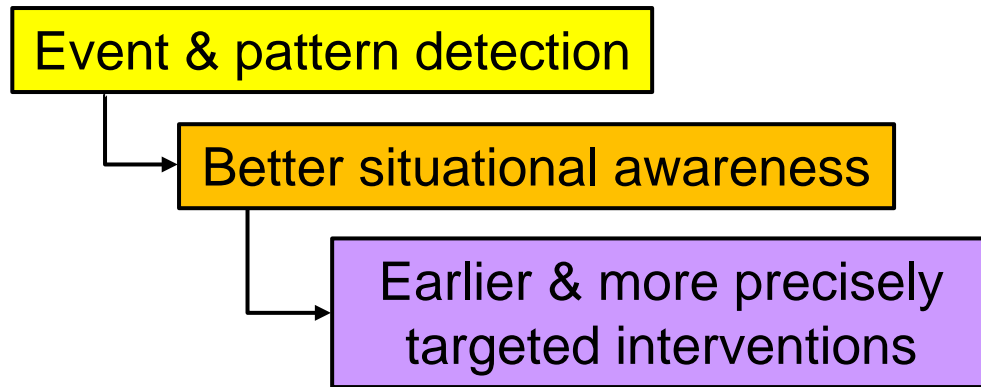


Discussion

Pre-syndromic surveillance is a **safety net** that can supplement existing ED syndromic surveillance systems by alerting public health to unusual or newly emerging threats.

Our recently proposed **semantic scan** can accurately and automatically discover pre-syndromic case clusters corresponding to novel outbreaks and other patterns of interest.

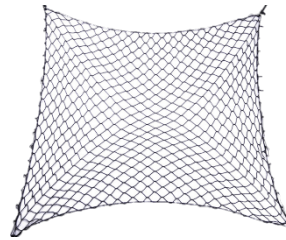
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities



Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis

Identifying causal effects of environmental exposures

We are using Medicaid data linked to detailed building characteristics in order to identify impacts of poor-quality housing on chronic health.

“Which housing conditions impact which health conditions, for which subpopulations, to what extent?”

Must adjust for known confounders, selection into treatment (exposure).

Step 1: Predictive model at building level

X = 65 diagnoses x {adult, child}

Y = building on landlord watch list?

Adult asthma and COPD

Mental health (ADHD, adjust. disorder)

Injuries (children and adults)

We have also developed an alternative scan-based approach to causal inference, based on automated discovery of natural experiments.

Key idea: treatment effects may be **heterogeneous**; use multidimensional scan to identify most affected subpopulations.

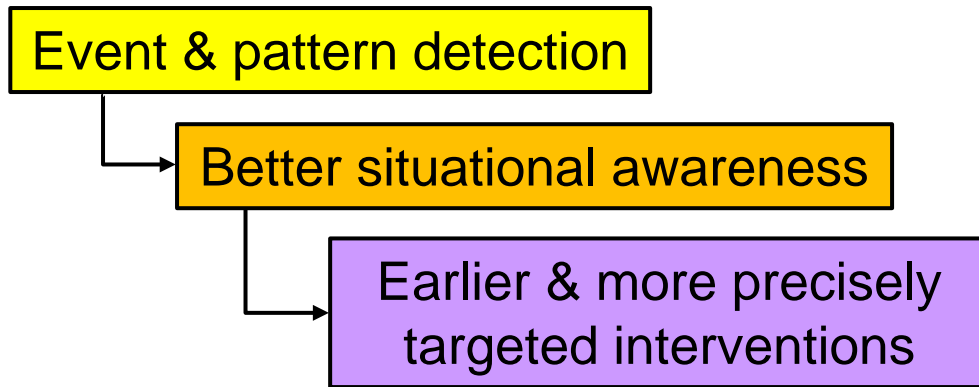
Must account for multiple hypothesis testing to bound false positive rate.

Step 2: Heterogeneous treatment effect scan

“**Crowded housing** is associated with increased respiratory conditions & injuries among Asians living in Manhattan.”



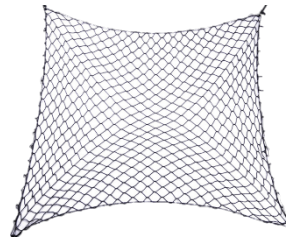
How can machine learning assist public health practitioners?



Early outbreak detection, including bioterrorism and other emerging bio-threats



Modeling and mitigating environmental causes of health disparities



Providing a safety net for novel outbreaks and unanticipated events



Interventions to combat the opioid crisis



Thanks for listening!

More details on my web site:
<http://www.cs.nyu.edu/~neill>

Or e-mail me at:
daniel.neill@nyu.edu