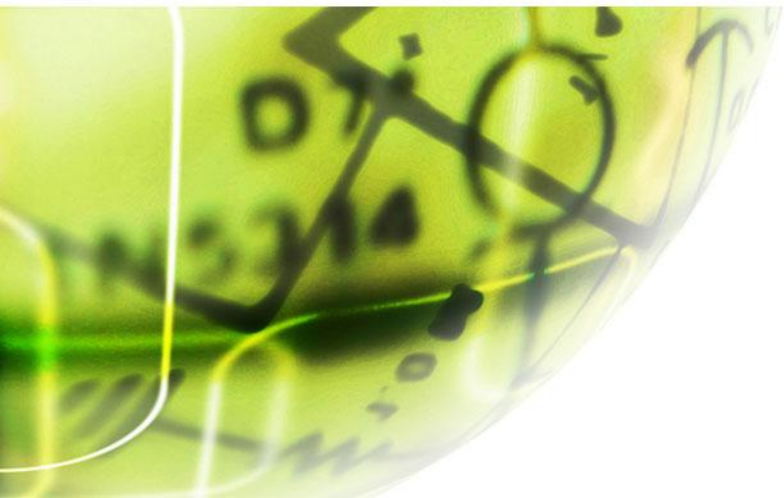




# A statistical approach to utilizing electronic health records



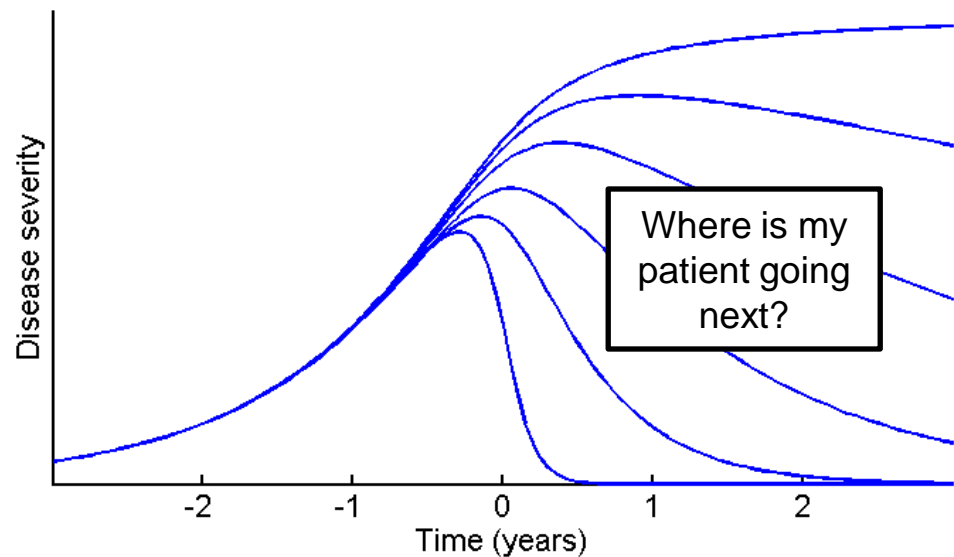
*Joseph E. Lucas, Ph. D.*



# Know your Patients



- ❑ Can a health system make use of what it knows about patients to make informed decisions?
  - Lots of data.
  - What questions to ask?
  - How to generate answers?
- ❑ Where is my patient going next?
- ❑ What will help me know?
  - EHR data, labs, genetics, “quantified self” data, genomics?



# Complicated Messy Data



- Electronic Health Records data has
  - Continuous data (labs, age, vitals)
  - Categorical data (gender, race, family history)
  - Written text (nurses' and physicians' notes, radiology reports)
  - Images (x-ray, CT, EKG)
- Important information everywhere
  - Example: A diabetic might have any or all of the following
    - Synonym of “diabetes” in a note
    - High lab values (glucose, HbA1C)
    - Relevant medications
    - Billing codes related to treatment of diabetes
    - Predisposing demographics (weight, race, family history)
    - Genetic predisposition (TCF7L2, JAZF1, HHEX, etc)
- We want to incorporate all of this information
- Don't want to be fooled by mistakes

# Approach



- ❑ Automated “Patients like me”
- ❑ Create groups of homogeneous patients
- ❑ This allows:
  - Automated generation of differential diagnosis
  - Novel comparative effectiveness studies
  - Listing of treatment options
  - Identification of Adverse drug events
  - Estimation of disease progression and prognosis
  - Assessment clinical utility of novel lab tests
- ❑ Predict probable patient type from other data
- ❑ Group patients through time

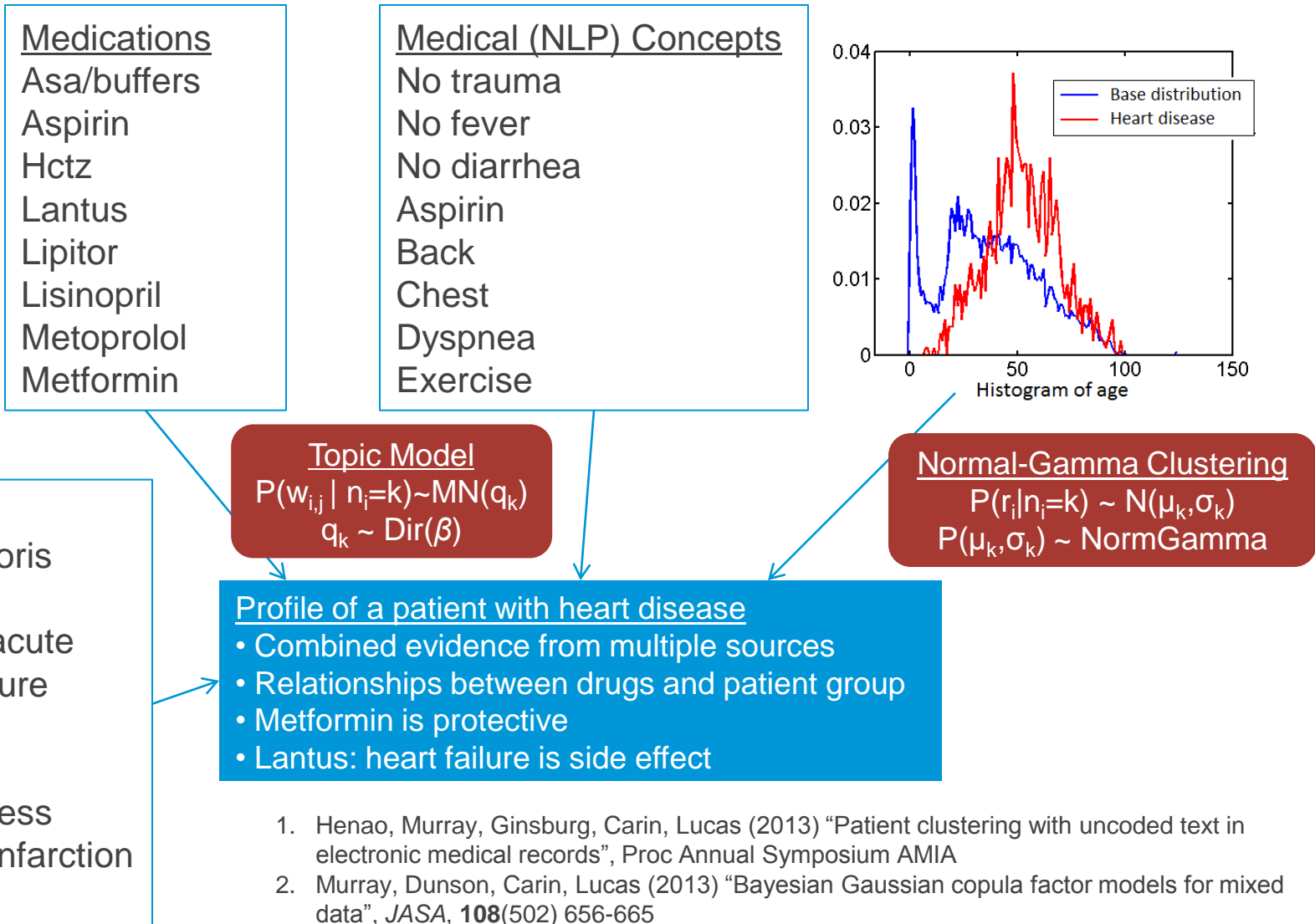
$$P(n_i = k|X) = f(m_k) \prod_l P_l(n_i = k|X)$$

- ❑ Product densities from each component
- ❑ Each data type is treated independently
- ❑ Doesn't explicitly model correlation across different data types



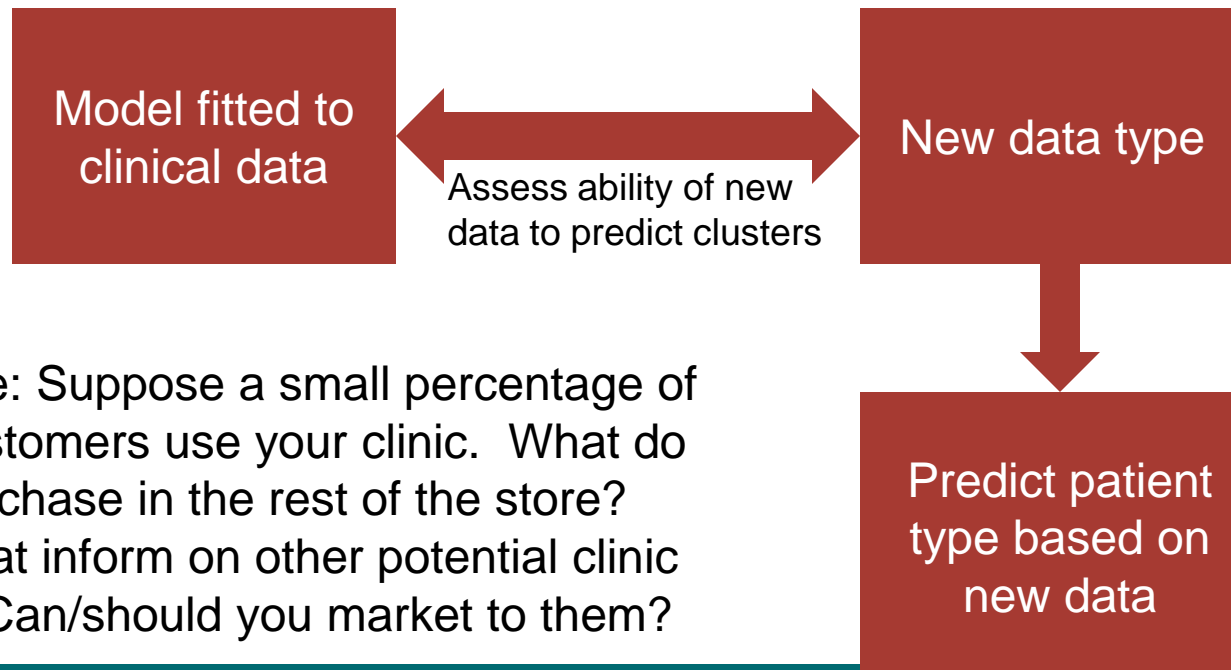
# Topic Models, NLP, Mixed data

Clustering patient visits to identify features of disease processes



# Nice features

- Very fast mixing
- MCMC without ever having to sample parameters
- Can compute MAP
- Model knows when it doesn't know
- Use to identify trends in other data



Example: Suppose a small percentage of your customers use your clinic. What do they purchase in the rest of the store? Does that inform on other potential clinic users? Can/should you market to them?



# Case Study: Data



54,000 records from ED

Contains

▪ Vocabularies

- Notes
- Orders
- Patient reported meds
- Diagnoses

▪ Categorical data

- Chief complaint
- Gender
- Disposition
- Zip code

▪ Continuous data

- Age
- Priority
- Vitals
- Weight

None is codified

All data subject to parsing errors

# Messy data



weakness/aching/hea  
daches

weakness/discomfort

weakness/dizziness

weakness/dizziness/re  
cent

weakness/faintness/c  
ongestion

weakness/fatigue

weakness/flaccid

weakness/heaviness

weakness/numbness

weakness/pain

weakness/shaking

weakness/tingling

weakness;

weaknessambulatory

weaknesscoughfeverl  
wbswent

weaknesscoughing

weaknessdiarrhea

weaknessdizziness

weaknesses

weaknessfalling

weaknessfatigue

weaknessn\aloc

weaknessper

weaknesss

weaknesssob

weaknesssore

weaknessunstable

weaknessvison

weaknessx

weaknss

Over 50,000 unique “words”  
with no copy editing. How to  
clean up mistakes?

# Unified Medical Language System (UMLS)



## ❑ Metathesaurus

- 10.5 million atoms in thesaurus (1.2 in SNOMED)

## ❑ Semantic Network

- 51.4 million relationships (2.9 in SNOMED)

## ❑ Lexicon

- 3.3 million spellings, inflections, properties, modifiers, abbreviations

## ❑ Together: 23 gigabytes of plaintext

- At 60 words per minute this is about 7 years of non-stop typing.

# UMLS Challenges



## ❑ Incredible number of arbitrary decisions

- What should be the semantic types?
  - Why is Mammal a semantic type but not primate
- What words to include?
- What relationships?

## ❑ Overly inclusive

- 10.5 million atoms in thesaurus(1.2 in SNOMED)
- Orders of magnitude more words in the thesaurus than are unique words in the 50k records
- How does one curate something this large?

## ❑ What do you do with it?

# Metamap Results Example



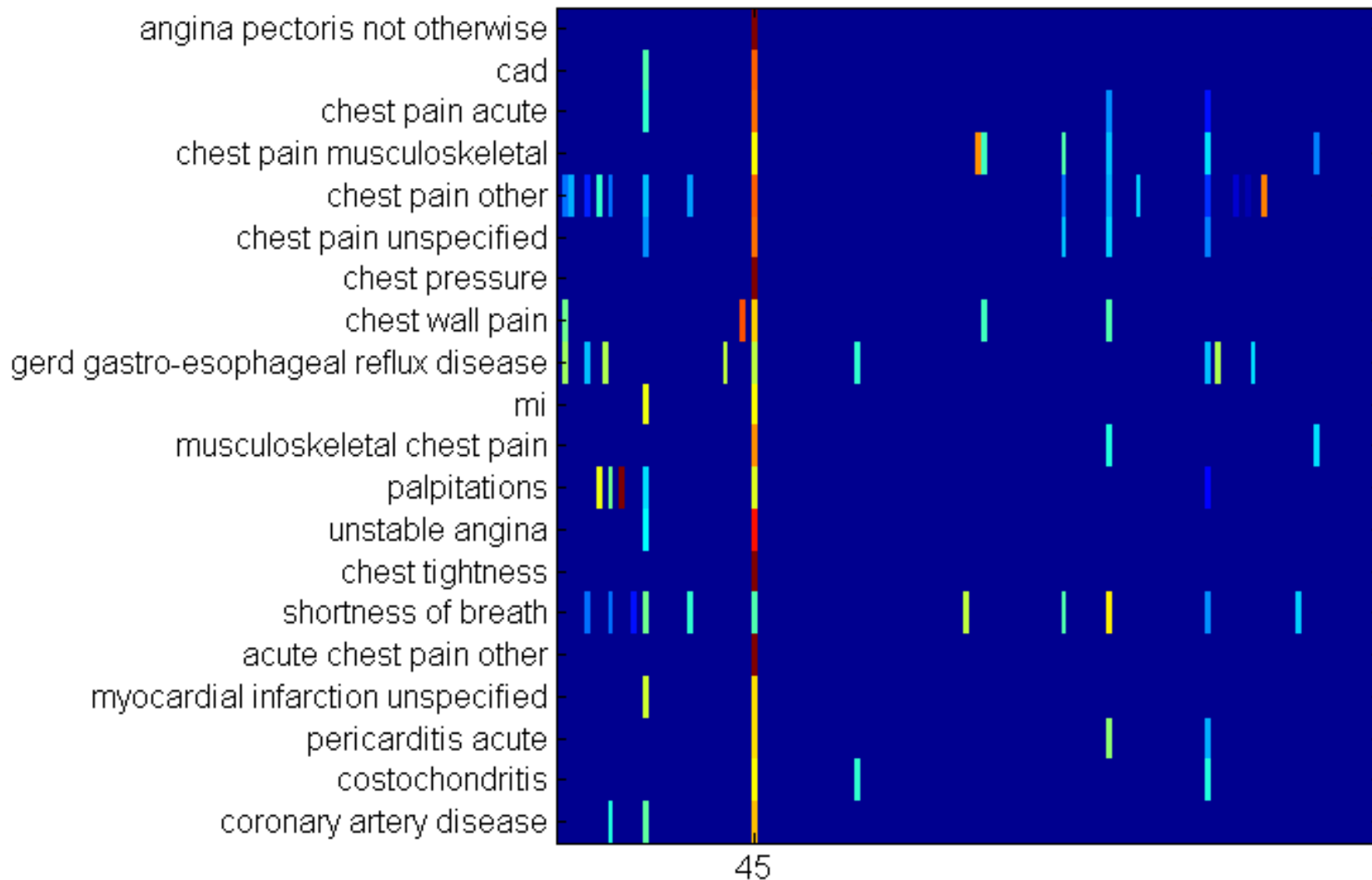
1. Pt here with c/o N/V and "shakes"
2. decreased po intake x 2 day
3. pt has pain pump which has been out of medication x 1 week
4. pt now taking PO narcotics
5. but presenting with N/V
6. pt with chronic back and neck pain

1. Tremor [Sign or Symptom]
2. Decreased [Quantitative Concept], Oral [Spatial Concept], /day [Temporal Concept]
3. Pain [Sign or Symptom], Pump, device [Medical Device], Drugs [Pharmacologic Substance], week [Temporal Concept]
4. Take [Health Care Activity], Oral [Spatial Concept]
5. Presentation [Idea or Concept], N+ (tumor staging) [Intellectual Product]
6. Chronic [Temporal Concept], Neck pain [Sign or Symptom]

# Patients with Chest Pain



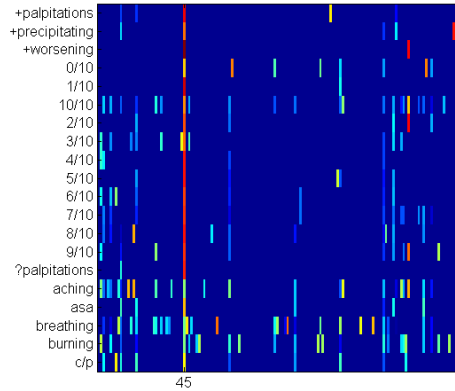
## Diagnosis



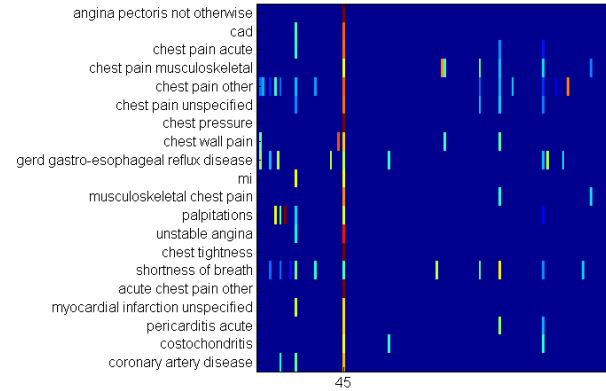
# Chest Pain



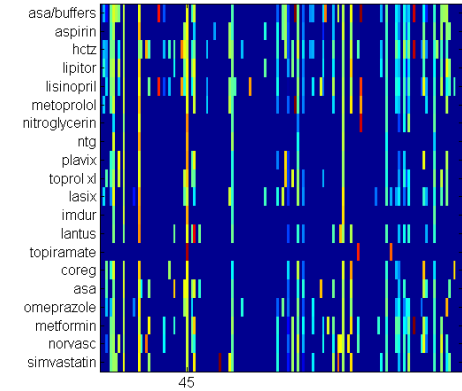
RN Notes



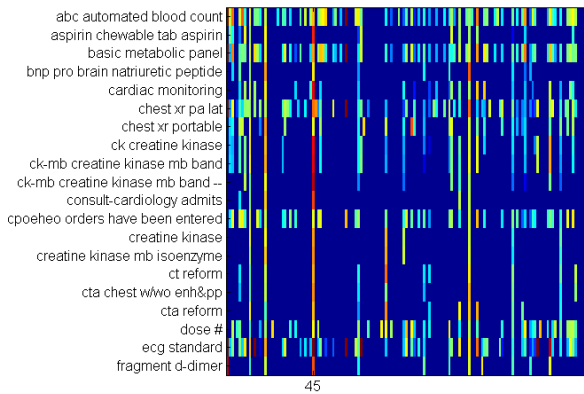
Diagnosis



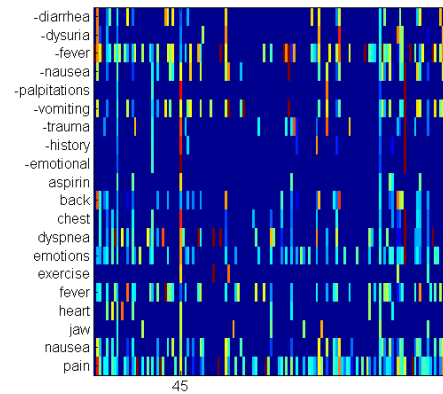
Medications



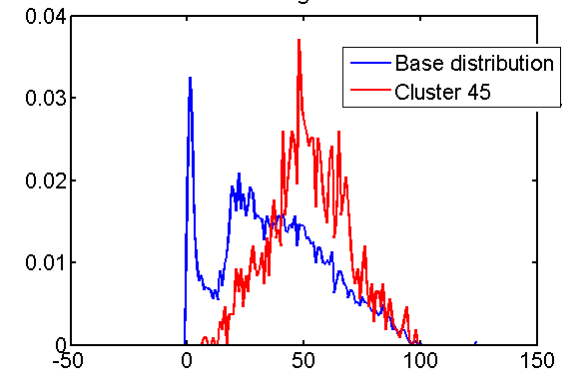
Orders



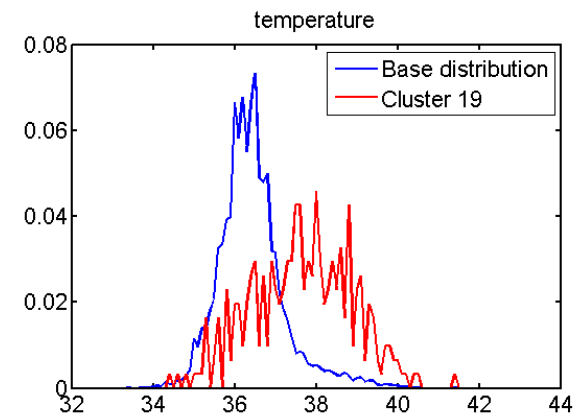
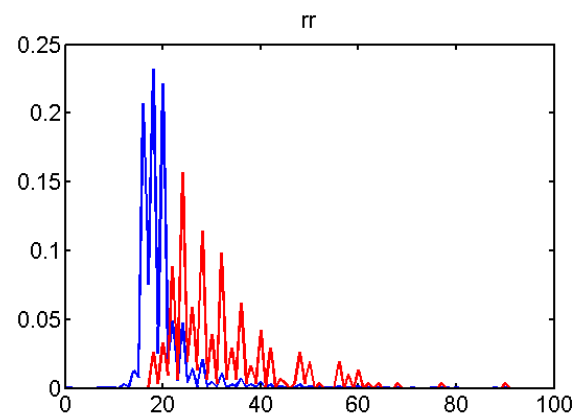
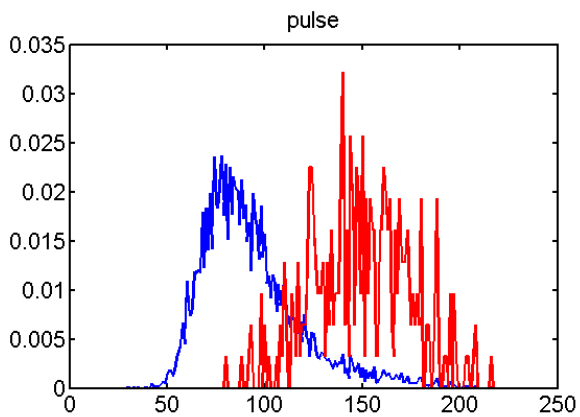
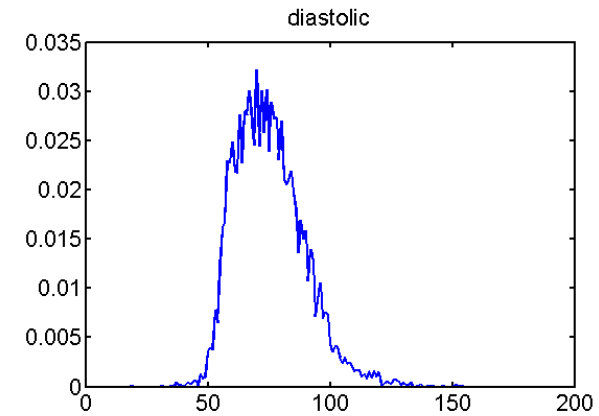
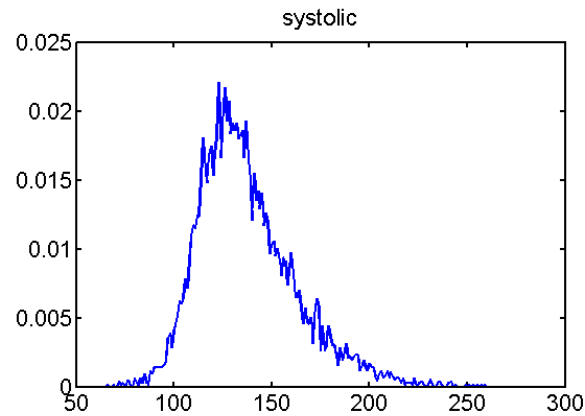
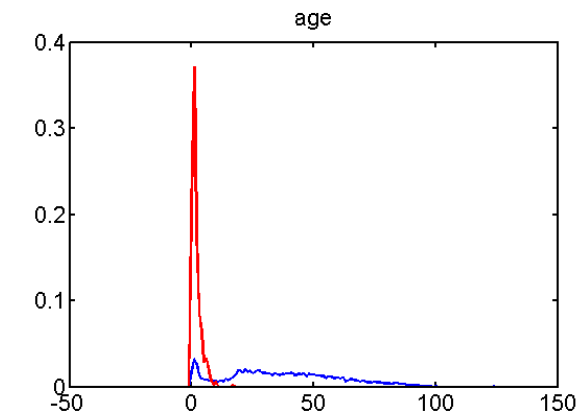
meta Notes



age



# Fever / Febrile Seizure





# Pharmacovigilance: Patients with Diabetes



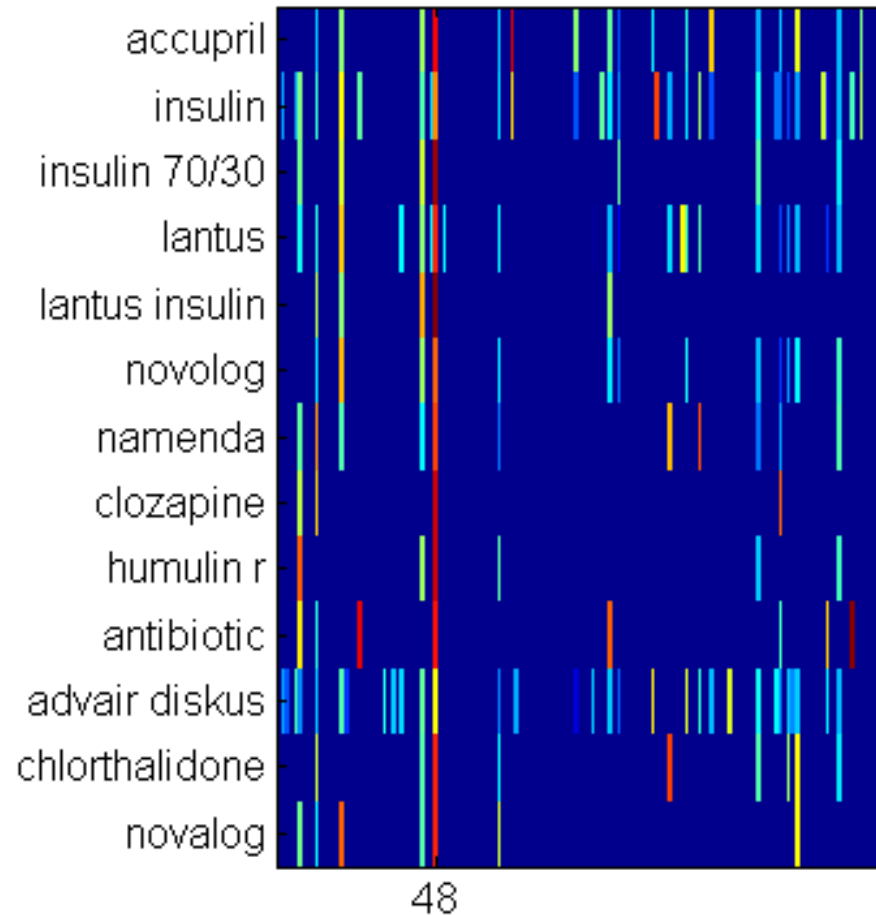
## Medications

Blood pressure

Dementia

Schizophrenia

Diuretic



# Some other associations



<u>Drug</u>	<u>Indication</u>	<u>Patient cluster</u>
Tussionex	Opioid	TIA
Altace	Blood pressure	Chest pain
Metformin	Blood sugar	Nose bleed
Tylenol	Analgesic	Rabies
Buspirone	Anxiolytic	Dog/cat bite
Cassodex	Chemo	Sickle cell
Vasotec	Blood pressure	Nose bleed
Prednisolone	Inflammation	Eye pain
Wellbutrin	Depression	Nose bleed
Zomig	Migraine	Chest pain

# Models for comparison



- Use UMLS to process RN notes
- Ignore UMLS and just throw out rare “words” from the RN notes
- Cluster with chief complaint
- Cluster with MetaMAP only

# Validation of Associations

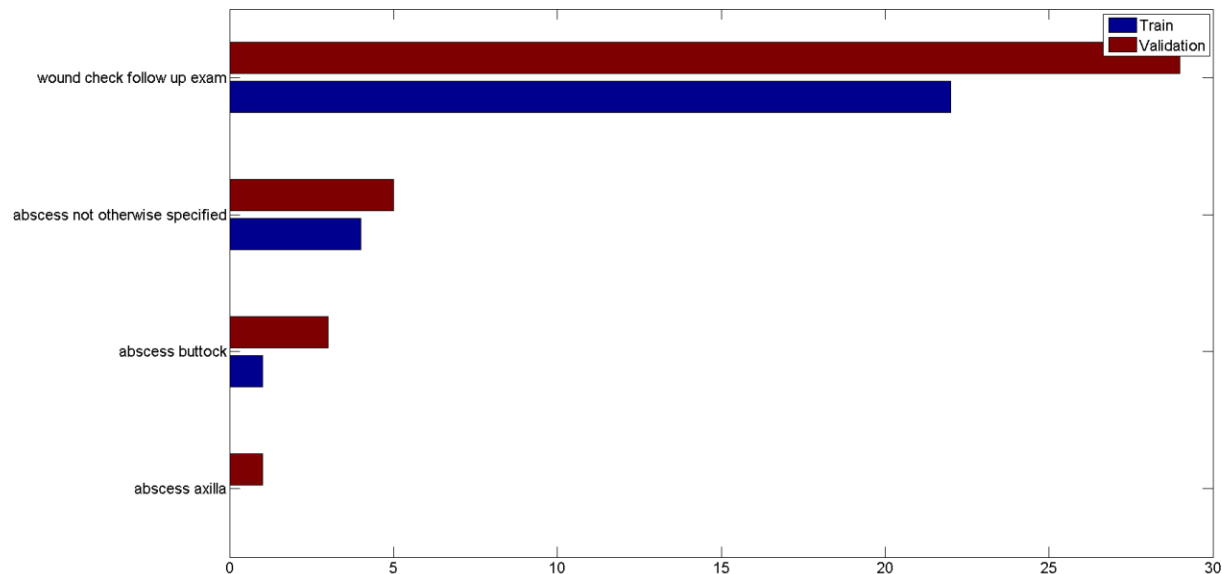


- ❑ Train the model on 20,000 samples selected uniformly at random
  - Test for association between patient clusters and drugs
- ❑ Fit the model to the test data
- ❑ Model using MetaMap
  - 4% of drug – symptom association tests have  $p\text{-value} < .05$  (3021)
  - 56% validate
- ❑ Model using terms in RN notes
  - 2% of tests significant (2258)
  - 58% validate
- ❑ Model using chief complaint
  - 3% of tests significant (9381)
  - 18% validate
- ❑ Model using MetaMap concepts only
  - 3% of tests significant (9319)
  - 22% validate

# Differential Diagnosis

## □ Average correlation (nonparametric)

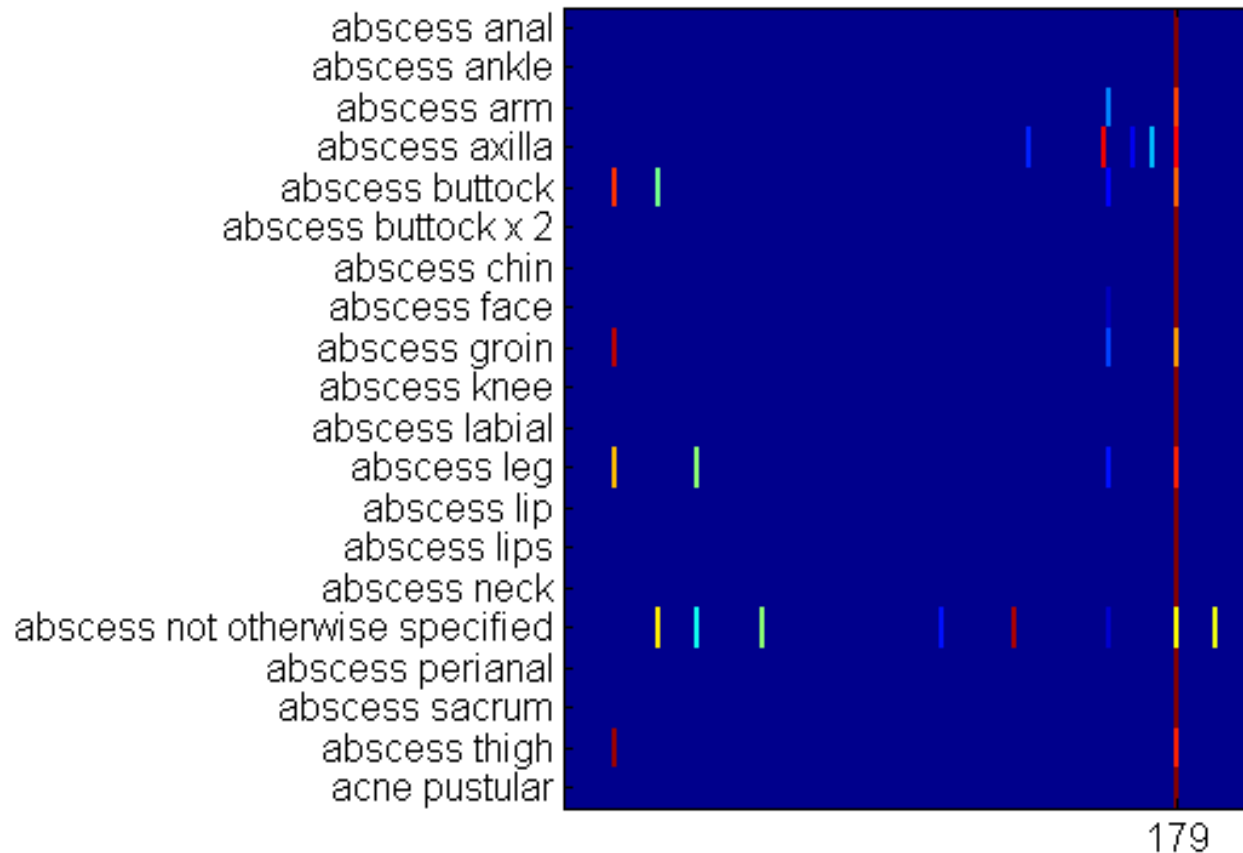
- Model with metemap, .61
- Model with RN notes, .62
- Chief complaint as classifier, .26



# Identify New Patient Populations



## Diagnosis



1.3% of patients.

# Other possible groups to add



- Minor trauma requiring x-ray, cast, splint
  - 13% of patients
- Chronic pain, headache, sickle cell crisis
  - 5.5%
- Cellulitis, local infection
  - 4%
- >50% of patients on either hydroxyurea or procrin presented with sickle cell crisis

# ID Patients by Drug



Ritonavir	metabolic imbalance (kidney disease)
Hydroxyurea	sickle cell
Procrit	sickle cell
Levodopa	Altered mental status



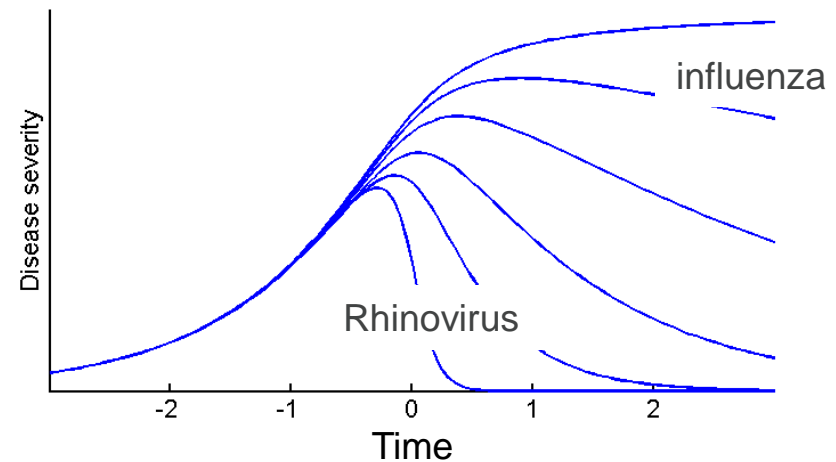
# Disease Processes over Time

*Learn a parametric model that describes the likelihood of observing particular diseases through time*

- Probability of disease over time

$$y = \text{disease severity}$$
$$\log(y) = -a(d^2+x^2)^{1/2}+bx$$

- ‘a’ and ‘b’ give control of rate of change of disease severity
- ‘d’ gives control of peak width
- Additional parameters allow control of height and temporal shift

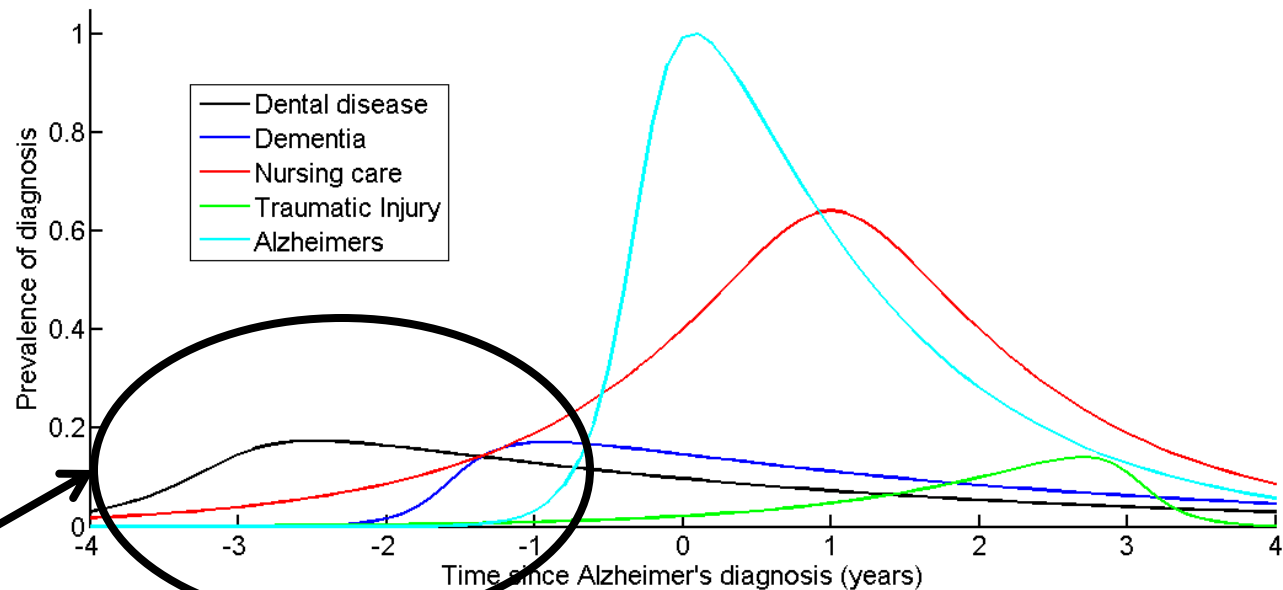


Combination of clustering model and temporal disease severity model allows identification and tracking of disease through time.

# Disease Trajectory

- ❑ New Data set: 3.5 million patients, 7 years
- ❑ Why is it important to predict future disease
  - Identification of high risk patients for early intervention
- ❑ Early intervention might be:
  - Standard therapy
  - Health coaching (telephone or in person)
  - Enrollment in clinical trial

## Trajectory of a patient population with Alzheimer's



Early signs that patient is developing Alzheimer's disease

Utilize all types of data to track and predict important features in the evolution of a patient's disease.

# Acknowledgments



## Quintiles Payer / Provider

- Brian Kelly
- Jon Morris

## Duke University

- Ricardo Henao
- Larry Carin
- Geoff Ginsburg



<http://EHRanalytics.blogspot.com>

[joseph.lucas@quintiles.com](mailto:joseph.lucas@quintiles.com)

[joe@stat.duke.edu](mailto:joe@stat.duke.edu)