# Using ontologies to mine unstructured data in medicine
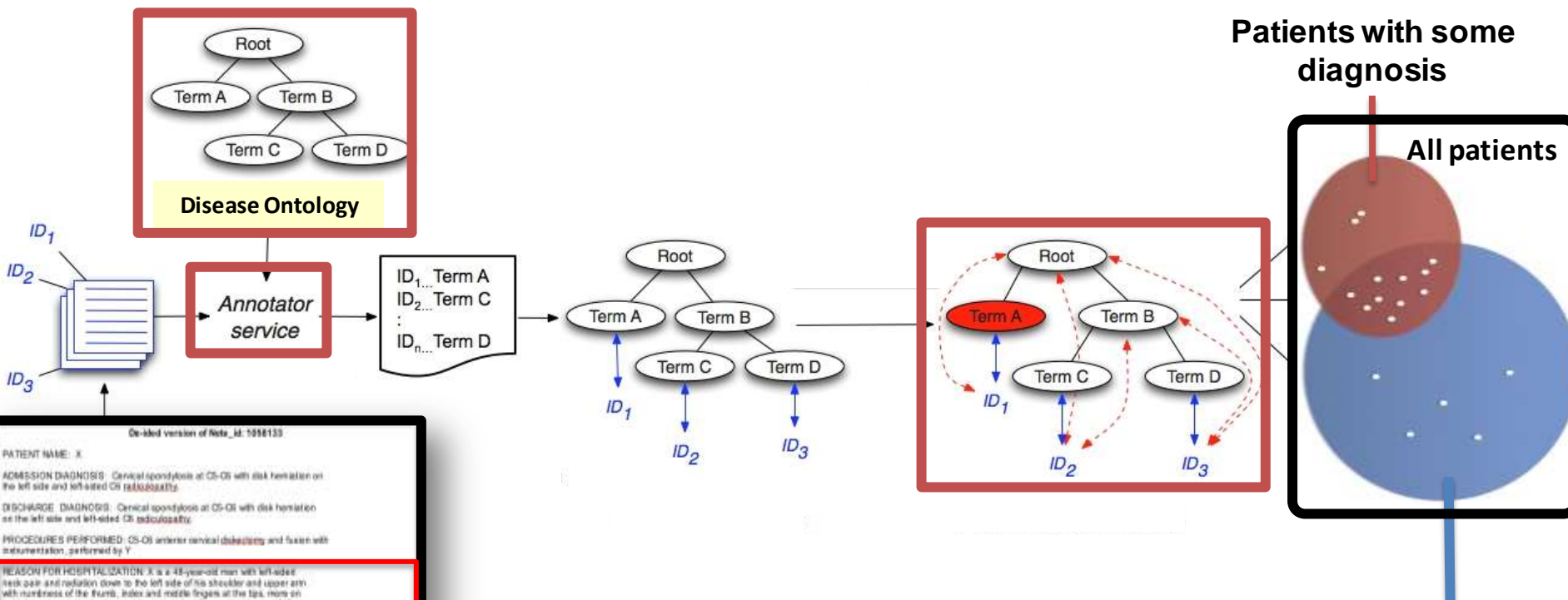
Nigam Shah, MBBS, PhD

nigam@stanford.edu

STANFORD
SCHOOL OF MEDICINE

# Profiling a patient set

# Profiling patient sets



ICD9 789.00
(*Abdominal pain, unspecified site*)

86k patient Reports

Patient records processed from U. Pittsburg NLP Repository with IRB approval.

# Associations and outcomes

|  | Gene | Disease | Drug | Device | Procedure | Environment |
|---|---|---|---|---|---|---|
| Gene |  | Gene Enrichment |  |  |  |  |
| Disease |  |  | Off-label Indications |  |  |  |
| Drug |  | Side effects |  |  |  |  |
| Device |  |  |  |  |  |  |
| Procedure |  |  |  |  |  |  |
| Environment |  |  |  |  |  |  |

What associations can we find?

# Generation of annotated data at scale

# Detecting the Vioxx Risk Signal



ROR of 2.058, CI of [1.804, 2.349]
The $X^2$ statistic has p-value $< 10^{-7}$

ROR=1.524, CI=[0.872, 2.666] $X^2$ p-value = 0.06816.

Vioxx Patients (1,560)

Vioxx$\rightarrow$MI (339)

MI Patients (1,827)

RA Patients (14,079)

p-value $< 1.3 \times 10^{-24}$

|          | MI        | No MI      |
|----------|-----------|------------|
| **Vioxx**    | a = 339   | b = 1221   |
| **No Vioxx** | c = 1488  | d = 11031  |

We should stop acting as if our goal is to author extremely elegant theories, […] and make use of the best ally we have: the unreasonable effectiveness of data.

# A Decade of Data Mining and Still Counting

*Manfred Hauben*[1,2,3,4] and *G. Niklas Norén*[5,6]

1    Pfizer Inc., New York, New York, USA
2    New York University School of Medicine, New York, New York, USA
3    New York Medical College, Valhalla, New York, USA
4    Brunel University, West London, UK
5    Uppsala Monitoring Centre, Uppsala, Sweden
6    Stockholm University, Stockholm, Sweden

# Big Data in biomedicine

# The problem

| | On-label | Off-label |
|---|---|---|
| **Indication** | What Pharma companies get approval for | Whatever else the doctor prescribes for |
| **Side effect / Adverse effect** | Found during the pre-marketing phase | Goal of drug-safety surveillance |

- Ambulatory: 100,000 deaths and $177 billion annually

- In patient: estimated that roughly 30% of hospital stays have an adverse drug event

- 21% of prescriptions

- 73% with very little evidence

# Detecting Off-label use



PUTATIVE:
gabapentin–depression = 6.2
gabapentin–mania = 5.7
gabapentin–hypomania = 8.4
gabapentin–familial tremor = 5.0

KNOWN "accepted/limited" USE:
gabapentin–perineal neuralgia (pain) = 46.7
gabapentin–stabbing pain = 8.9
gabapentin–diffuse pain = 5.6
gabapentin–excruciating pain = 5.0

FDA–APPROVED:
gabapentin–clonic seizures = 4.9
gabapentin–rhinencephalic epilepsy = 7.9

# Detecting Adverse Events



Temporally ordered bags of terms

Patients / Cohorts of interest

Patient's Timeline

492 thousand pairs

959 thousand pairs

Adverse Events are Positives, Indications are Negatives

# Patterns worth testing (off-label usage, which is risky)

✓ **Identify off-label use**
- Find drug-indication pairs that "look like" indications

✓ **Identify which use "may be risky"**
- Use existing, known side effect databases
- Learn drug-disease associations that look like side effects

✓ **Assemble I-D-A triplets**
- Indication – Drug – Adverse effect. e.g. RA – Vioxx – MI

✓ **Test on unstructured data**

# Testing 'interesting patterns'

# Detecting Novel Associations in Large Data Sets

David N. Reshef,[1,2,3]*† Yakir A. Reshef,[2,4]*† Hilary K. Finucane,[5] Sharon R. Grossman,[2,6] Gilean McVean,[3,7] Peter J. Turnbaugh,[6] Eric S. Lander,[2,8,9] Michael Mitzenmacher,[10]‡ Pardis C. Sabeti[2,6]‡

Identifying interesting relationships between pairs of variables in large data sets is increasingly important. Here, we present a measure of dependence for two-variable relationships: the maximal information coefficient (MIC). MIC captures a wide range of associations both functional and not, and for functional relationships provides a score that roughly equals the coefficient of determination ($R^2$) of the data relative to the regression function. MIC belongs to a larger class of maximal information-based nonparametric exploration (MINE) statistics for identifying and classifying relationships. We apply MIC and MINE to data sets in global health, gene expression, major-league baseball, and the human gut microbiota and identify known and novel relationships.

The team @
www.bioontology.org/project-team


NIH Roadmap grant U54 HG004028