

## Tackling The Challenges of Big Data

### Big Data Analytics

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

## Tackling The Challenges of Big Data

### Big Data Analytics

#### Applications - Medical Outcomes

Introduction

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



## My Research Group

- Technical areas
  - Machine learning and data mining
  - Algorithm design and signal processing
  - Computer vision
  - Software systems
- Application areas
  - *Medical analytics*
  - Financial analytics
  - Sports analytics
  - Mobile telemedicine
- Principal funding
  - Quanta Computer, NSF, NIH
  - NSERC, QCR, MSR, MLB.com



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Medical Analytics

- Building predictive models using machine learning
- From signals to bio-markers
- Video-based monitoring & diagnoses



MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

DATA SCIENCE

## The Good News

- Capacity to gather medically significant data growing quickly
- Better instrumentation (e.g., MRI machines, ambulatory monitors, cameras) generates more information/patient

Pictures removed due to copyright restrictions.  
We are sorry for any inconvenience this may cause.

- More storage capacity allows information to be saved
- Economic and social forces creating more aggregation of data

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

DATA SCIENCE

## Problems Posed by Increase in Available Data

- Clinicians spending ever more time
  - Studying data about their patients
    - \*And still ignoring most of the data
  - Tracking onslaught of new medical data
    - \*And not keeping up

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology

DATA SCIENCE

## Problems Posed by Increase in Available Data

- Researchers using analytical techniques that:
  - Don't scale to multi-modal data, thousands of variables, millions of patients
  - Designed to test hypotheses, not to uncover new knowledge
  - Require prohibitively expensive and time consuming clinical studies



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Using ML to Make Useful Predictions

- Predicting medical outcomes is hard
- It will be many years before biology provides answers
- But there is lots of "found data" that can be used to derive pretty good predictions
- One focus is predicting avoidable complications



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Accurate Predictions Can Help

- Better decisions about which treatments (if any) to choose
  - E.g., surgery or drugs
- Better management of medical providers
  - Match right patients with right providers, streamline utilization of resources
- Understanding what goes into making a good prediction can inform changes
  - E.g., prophylactic treatments



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## The Big Data Challenge

- Too Big: images, videos, signals
  - 10's of millions of patients, billions of bits/patient
  - Parallelism will help, but not if we use  $n^3$  algorithms
- Too Hard
  - Multiple modalities: signals, lab results, images, genomic, natural language ...
  - Always incomplete, often incorrect, outcomes ambiguous
  - Ground truth often hard to come by
  - Site-to-site variation
- One saving grace: human physiology and medical practice changes a lot more slowly than financial markets
  - What we learn likely to be of long-lasting value



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Tackling The Challenges of Big Data

### Big Data Analytics

### Applications - Medical Outcomes

Introduction

**THANK YOU**



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

---

## Tackling The Challenges of Big Data

### Module: Big Data Analytics

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

---

## Tackling The Challenges of Big Data

### Big Data Analytics

#### Applications - Medical Outcomes

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



### Computationally-Generated Biomarkers

- Biomarker (broadly defined)
  - A characteristic that can be objectively measured and provides an indicator of normal or pathological processes or expected responses to a therapeutic intervention
- Computationally-Generated Biomarker
  - A biomarker generated by applying computation to medical data
- One current project
  - Cardiovascular risk stratification



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



### Cardiovascular Risk Stratification

- Acute coronary syndrome (ACS) common: ~1.25M/year in U.S.
  - 15% - 20% of these people will suffer cardiac-related death within 4 years
- Have lots of good treatments
  - Eat better
  - Get an ICD

Pictures removed due to copyright restrictions.  
We are sorry for any inconvenience this may cause.

- Choosing who should get which (or any) is hard
- Fine-grained risk stratification the key



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

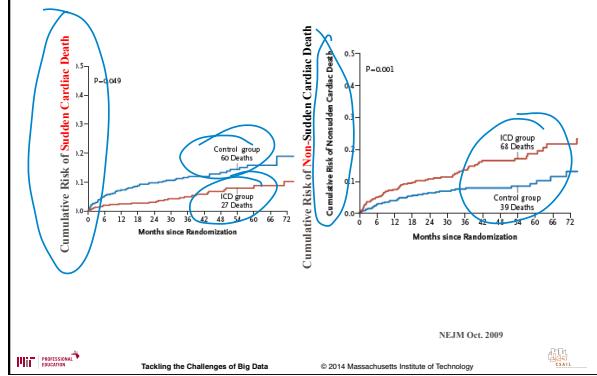
---

---

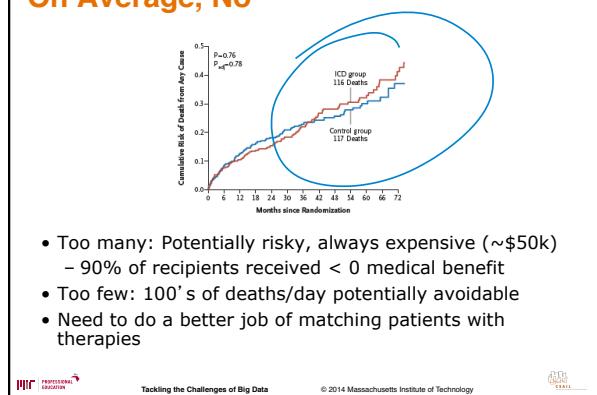
---

---

## Do ICD's Save Lives?



## On Average, No



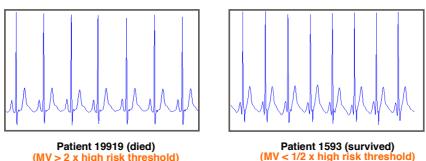
## Approaches to Identifying High Risk Cases

- Clinical characteristics
  - E.g., gender or high blood pressure
- Traditional biomarkers
  - E.g., cholesterol levels
- Echocardiography
  - E.g., LVEF
- Electrocardiography (ECG)
- We combine all of these
  - Innovation on use of ECG



## Morphological Variability

- Detect signs of small and transient electrical instability in myocardium by detecting minor differences in shape of normal appearing heart beats
- Invisible to human eye



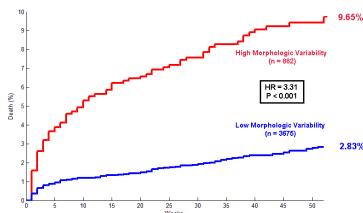
Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Data from MERLIN TIMI-36 dataset (~1B beats)

- About 4,500 ACS patients
- 1 year follow-up for cardiovascular death (193 events)



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## MV and Other Risk Variables

Multivariate Analysis of Patients with LVEF > 40%\*

Parameter	Adjusted Hazard Ratio (*)	95% Confidence Interval	P Value
<b>MV</b>	<b>2.31</b>	<b>1.28 – 4.16</b>	<b>0.005</b>
BNP	1.94	1.10 – 3.43	0.022
BMI	0.93	0.54 – 1.62	0.805
CrCl	1.48	0.86 – 2.56	0.159

\* Also adjusted for age, hypertension, diabetes, hypercholesterolemia, prior MI, prior angina, ST changes, cardiac markers

Z. Syed, C. Stultz, B. Scirica B., and J. Guttag, "Computationally generated cardiac biomarkers for risk stratification following acute coronary syndrome," *Science Translational Medicine*, September 2011.



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



**Tackling The Challenges of Big Data**  
**Big Data Analytics**  
**Applications - Medical Outcomes**

Section 2

**THANK YOU**



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Analytics**

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

**Tackling The Challenges of Big Data**  
**Big Data Analytics**  
**Applications - Medical Outcomes**

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

## Healthcare-Associated Infections (HAIs)

- 1.7 million people per year get an infection during a hospital stay
- Contribute to approx. 4% of all deaths in the US
- Estimated costs of \$20 Billion per year



MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## HAIs Remain Stubbornly Prevalent

THE WALL STREET JOURNAL..

"Medicare Shift Fails to Cut Hospital Infections"

-- Oct. 10<sup>th</sup> 2012

SCIENTIFIC AMERICAN

"Hospitals Fail to Take Simple Measures to Thwart Deadly Infections, Survey Says"

-- April 8<sup>th</sup>, 2013

The New York Times

"More Aggressive Action Urged to Curb Hospital Infections"

-- May 29<sup>th</sup>, 2013

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Why?

- Temporal dynamics of risk factors are not well understood
- Global risk factors do not take into account institutional differences (institutional risk factors)
- Goal: use machine learning to build data-driven hospital-specific risk stratification models for healthcare-associated infections

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## What Is Machine Learning?

- Automating automation
  - Computer programs can automatically follow rules
  - How do we determine these rules automatically in the first place?
- Machine learning focuses on getting computers to program themselves
  - Let the data do the work
  - Automatically generate programs that create the right outputs from data



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

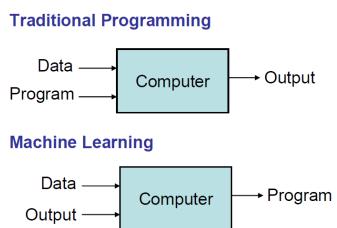
---

---

---

---

## What Is Machine Learning?



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## How Are Things Learned?

- Memorization
  - Accumulation of individual facts
  - Limited by
    - \*Time to observe facts
    - \*Memory to store facts
- Generalization
  - Deduce new facts from old facts
  - Limited by accuracy of deduction process
    - \*Essentially a predictive activity
    - \*Assumes that the past predicts the future



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Machine Learning Methods

- Many different ones, all try to learn a model that is a generalization of examples
- Supervised: given a set of feature/label pairs, find a rule that predicts the label associated with a previously unseen input
- Unsupervised: given a set of feature vectors (without labels) group them into “natural clusters”
- All have three components
  - Representation of the model
  - Metric for assessing goodness of model
  - Optimization method for learning the model



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

## More Not Always Better

- Irrelevant information can slow things down
  - Volume of search space (i.e., possible combinations of features) increases exponentially with extra dimensions
- Irrelevant information can lead to a bad model
  - Fit the noise rather than the signal
  - Especially when dimensionality large relative to number of examples



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

## Back to Infections

- Challenges
  - Temporal aspect – changes over time
  - High dimensionality – thousands of variables
  - Paucity of cases – relatively rare events
  - Institutional differences – global models fail
  - Ground truth is not available – not all patients are tested



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

## Approach

- Use data-driven models to gain insight into the temporal aspects of the problem
- Use transfer learning based techniques to
  - Leverage data from other hospitals while
  - Incorporating hospital-specific features
- Topic of next section



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

## Tackling The Challenges of Big Data Big Data Analytics Applications - Medical Outcomes

**THANK YOU**



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

---

## Tackling The Challenges of Big Data Big Data Analytics

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

---

## Tackling The Challenges of Big Data

### Big Data Analytics Applications - Medical Outcomes

**John Guttag**

Professor

Massachusetts Institute of Technology



© 2014 Massachusetts Institute of Technology



### Predicting *Clostridium difficile*

- Bacteria that takes over the gut
- Transmitted through the mouth
- Causes severe diarrhea, intestinal diseases
- Treatment: metronidazole and oral vancomycin
- 20% of cases relapse within 60-days
- 178,000/year in U.S.
  - Approx. same as invasive breast cancer



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



### Risk Factors

#### Time Invariant

- Collected at the time of admission
- e.g., admission complaint, previous admissions, home meds

#### Time Varying

- Changes during the hospitalization
- e.g., current meds, current procedures, current location, hospital conditions



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

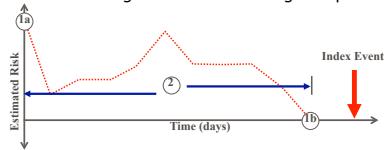
---

---

---

## Incorporating Time

- Typical Approaches in Clinical Literature
  - Calculate risk using only a snapshot of patient's state
    - a. At time of admission [Tanner et al., 2009]
    - b.  $n$  days before index event [Dubberke et al., 2011]
- Our Approach
  - Calculate risk using the entire evolving risk profile



MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## The Data

- EMR from a single large urban hospital in the US
- In-patients stays  $\geq 7$  days from a single-year
- Population:  $\sim 10,000$  hospital admissions
  - $\sim 200$  positive cases

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



## Variables Considered

### • Data available at the time of admission

- Age
- Sex
- Admission Complaint
- Admission Procedure
- Data from Hospital Admissions/Visits in last 90 days
- Medical History

### • Data collected during the hospital stay

- Medications
- Locations of the Patient
- Procedures
- Lab Results
- Vitals
- Devices
- Staff

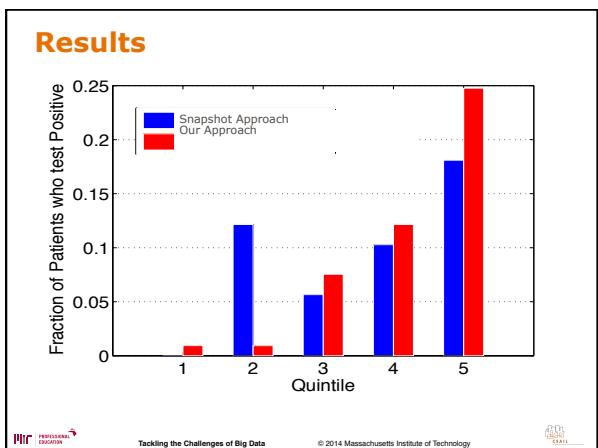
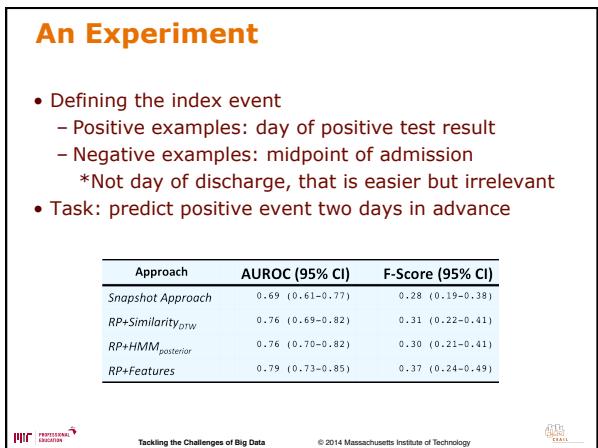
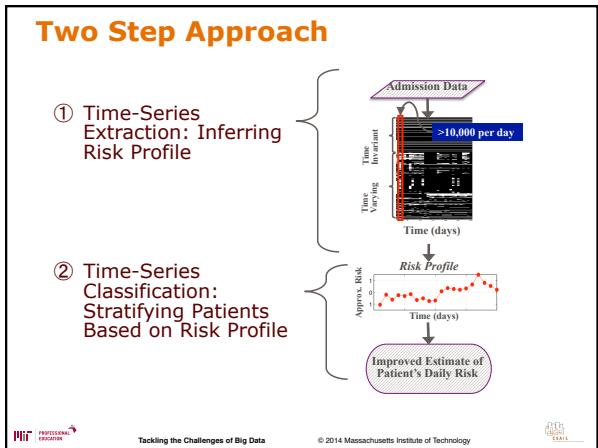
EHR Data  
Not Billing Data

MIT Professional Education

Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology





## Wrapping Up

- Data analysis has transformed many fields
  - Biology
    - \*From molecules, to genetics, to systems
  - Finance
    - \*It takes a computer to lose \$10,000,000/minute
  - Sports

Pictures removed due to copyright restrictions.  
We are sorry for any inconvenience this may cause.

- It's time to transform medicine!



Tackling the Challenges of Big Data

© 2014 Massachusetts Institute of Technology



---

---

---

---

---

---

---

---

## Tackling The Challenges of Big Data Big Data Analytics Applications - Medical Outcomes

**THANK YOU**



© 2014 Massachusetts Institute of Technology

---

---

---

---

---

---

---

---