Machine Learning in Biostatistics & Health Policy

Sherri Rose

Associate Professor Department of Health Care Policy Harvard Medical School

rose@hcp.med.harvard.edu
drsherrirose.com

July 20, 2016



Open access, freely available online

Why Most Published Research Findings **Are False** John P. A. Ioannidis

The New Hork Times

Essay

September 16, 2007

Do We Really Know What Makes Us Healthy?

By GARY TAUBES

Open access, freely available online

Why Most Published Research Findings Are False John P.A. Joannidis

The New Hork Times

nytimes.com

September 16, 2007

Do We Really Know What Makes Us Healthy?

By GARY TAUBES

4 high impact zones in statistical discovery with big data **FierceBigData**

Sentember 22, 2014

dr 4

E Like

S+1



By: Sherri Rose, Harvard University: David Durson, Duke University: Tyler McCormick, University of Washington: and Cynthia Rudin, MIT

Big data is transforming society with the help of statisticians, who possess in-depth experience and expertise in the art and science surrounding data. The American Statistical Association, or ASA, has

recently released a white paper entitled "Discovery with

Clockwise from top left: Dunson, Rudin, McCormic and Rose Transform Science and Society" (od) that highlights high-impact areas where statistical science is being applied to transformative big data research questions. Statistics is, by definition, the science of learning from data, and has had a key impact in several of the most prominent fields of discovery, including the biological sciences, health care, business analytics and recommendation systems, and the social sciences. There is a strong need to work in integrated teams comprised of domain experts, statisticians, and computer scientists in order to solve these complicated problems, which require tailored solutions using the influx of big data.

AMSTATNEWS

Statistics Ready for a Revolution

1 SEPTEMBER 2010 503 VIEWS 2 COMMENTS

Next Generation of Statisticians Must Build Tools for Massive Data Sets

Mark van der Laan, Jiann-Ping Hsu/Karl E. Peace Professor in Biostatistics and Statistics at UC Berkeley, and Sherri Rose. PhD candidate at UC Berkelev

variations

Big data and the future

At the beginning of her career Sherri Rose discusses hig data and stands amazed at its notential.







insert data



big data system

Idreos/Rose



The increasing availability of electronic health records offers a **new** resource to public health researchers.

General usefulness of this type of data to answer targeted scientific research questions is an open question.

Need **novel statistical methods** that have desirable statistical properties while remaining computationally feasible.

Machine Learning: Big Picture





Machine learning aims to

- "smooth" over the data
- make fewer assumptions



 Recent health studies have employed newer algorithms. (any mapping from data to a predictor)



- ▶ Recent health studies have employed newer algorithms.
- Researchers are then left with questions, e.g.,
 - "When should I use random forest instead of standard regression techniques?"



Recent health studies have employed newer algorithms.

- Researchers are then left with questions, e.g.,
 - "When should I use random forest instead of standard regression techniques?"

Journal of Clinical Epidemiology



Journal of Clinical Epidemiology 63 (2010) 1145-1155

Logistic regression had superior performance compared with regression trees for predicting in-hospital mortality in patients hospitalized with heart failure Peter C. Austin^{a.e.d.}, Jack V. Tu^{ab.e.d.} Douglas S. Lee^{a.d.}



Recent health studies have employed newer algorithms.

- Researchers are then left with questions, e.g.,
 - "When should I use random forest instead of standard regression techniques?"

Journal of Clinical Epidemioloav



Journal of Clinical Epidemiology 63 (2010) 1145-1155

European Journal of Neurology 2010, 17: 945-950

doi:10.1111/j.1468-1331.2010.02955 x

Logistic regression had superior trees for predicting in-hospital r

hea Random forest can predict 30-day mortality of spontaneous Peter C. Austinace*, Jack intracerebral hemorrhage with remarkable discrimination

S. -Y. Penga,b,c, Y. -C. Chuangb, T. -W. Kangb and K. -H. Tsengd *Institute of Biomedical Informatics, National Yang-Ming University, Taipei; bDepartment of Anesthesiology, Taichung Veterans General Hospital, Taichung: School of Medicine, Chung Shan Medical University, Taichung; and Department of Nephrology, Taoyuan Veterans Hospital, Taoyuan, Taiwan

Machine Learning: Ensembles



- Ensembling methods allow implementation of multiple algorithms.
- Do not need to decide beforehand which single technique to use; can use several by incorporating cross-validation.



Machine Learning: Ensembles



Super Learner: Build a collection of algorithms consisting of all weighted averages of the algorithms.

One of these weighted averages might perform better than one of the algorithms alone.



Machine Learning: Ensembles





Applications in Health Insurance

- Risk adjustment in plan payment
- 2 Effect estimation for variable importance of medical conditions

DATA

Truven MarketScan database, up to 51 million enrollees per year. Variables: age, sex, region, procedures, expenditures, etc.



Enrollment and claims from private health plans and employers.

MARKETSCAN® RESEARCH



Risk Adjustment in Plan Payment



Over 50 million people in the United States currently enrolled in an insurance program that uses risk adjustment.

- Redistributes funds based on health
- Encourages competition based on efficiency & quality
- Huge financial implications





xerox.com

Risk Adjustment in Plan Payment = Frozen



 $E[Y \mid W] = \alpha_0 + \alpha_1 W$

Potentially *\$\$\$* oversight, where it attempts to control for the impact of consumers choosing health plans.



wpb.org

Risk Adjustment in Plan Payment: Results



- **O** Super Learner had best performance.
- **②** Top 5 algorithms with reduced set of variables retained 92% of the relative efficiency of their full versions (86 variables).
 - age category 21-34
 - all five inpatient diagnoses categories
 - heart disease
 - cancer
 - diabetes
 - mental health
 - other inpatient diagnoses
 - metastatic cancer
 - stem cell transplantation/complication
 - multiple sclerosis
 - end stage renal disease

Medical Condition Variable Importance





Medical Condition Variable Importance: Results





Medical Condition Variable Importance: Results



First investigation of the impact of medical conditions on health spending as a variable importance question using double robust estimators.

Five most expensive medical conditions were

- multiple sclerosis
- 2 congestive heart failure
- Iung, brain, and other severe cancers
- **(**) major depression and bipolar disorders
- **o** chronic hepatitis.
 - Differing results compared to parametric regression.
 - What does this mean for incentives for prevention and care?

Causal Questions







Stent in Coronary Artery



HatfieldBlausen.com staff. "Blausen gallery 2014"

Hospital Profiling

Regression Only Clinical Confounders



-2.5 -1.5 -0.5 0.5 1.5 2.5 3.5 4.5

Ú. R a P. Ν м G в -2.5 -1.5 -0.5 0.5 1.5 2.5 3.5 4.5 A-IPW Full Confounders w υ ŝ R a P ò Ň м н G

-2.5 -1.5 -0.5 0.5 1.5 2.5 3.5 4.5

A-IPW Clinical Confounders

w

E

D

TMLE Clinical Confounders







Targeted Learning Methods







van der Laan & Rose, Targeted Learning: Causal Inference for Observational and Experimental Data. New York: Springer, 2011. targetedlearningbook.com

Health Policy Data Science Lab



healthpolicydatascience.org

