

Data Science, Statistics, and Health

Rob Tibshirani
Departments of Biomedical Data Science & Statistics
Stanford University



Outline

1. Some general thoughts about data science and health
2. Quick introduction to supervised learning
3. **Example:** Predicting platelet usage at Stanford Hospital
4. **Example:** The Delphi project- COVID19 forecasting

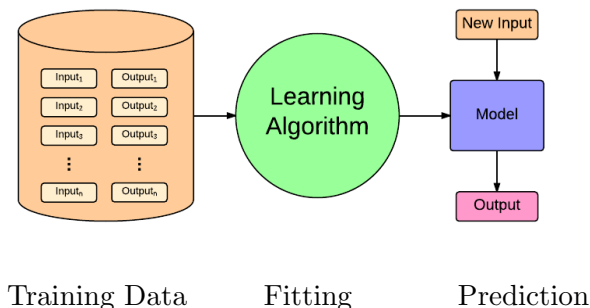
There is a lot of excitement about Data Science and health

- Artificial intelligence, predictive analytics, precision medicine are all hot areas with huge potential
- A wealth of data is now available in every area of public health and medicine, from new machines and assays, smart phones, smart watches
- Already, there have been good successes in data science in pathology, radiology, and other diagnostic specialties.

For Statisticians: 15 minutes of fame

- 2009: “ I keep saying the **sexy** job in the next ten years will be **statisticians**.” Hal Varian, Chief Economist Google
- 2012 “**Data Scientist**: The Sexiest Job of the 21st Century”
Harvard Business Review

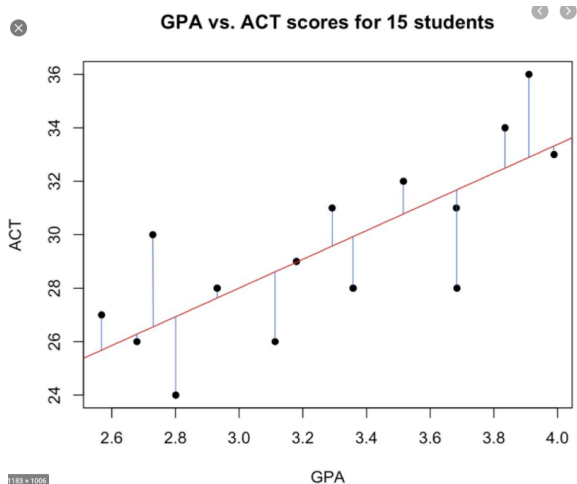
The Supervising Learning Paradigm



Traditional statistics: domain experts work for 10 years to learn good features; they bring the statistician a small clean dataset

Today's approach: we start with a large dataset with many features, and use a machine learning algorithm to find the as good ones. **A huge change.**

Simple prediction using a straight line fit



This is a form of “supervised learning” using one feature

Supervised learning via least squares regression

- We have “training data” $(x_1, y_1) \dots (x_n, y_n)$ on n individuals or units.
- Each x_i is a vector (collection) of p features measured on individual i .
- y_i is the target (outcome) that we are trying to predict.
- We assume a linear model of the form

$$y_i \approx \beta_0 + x_{i1}\beta_1 + x_{i2}\beta_2 \dots x_{ip}\beta_p$$

We estimate the unknown parameters (weights) β_j by minimizing the least squares criterion

$$\sum_i [y_i - (\beta_0 + x_{i1}\beta_1 + x_{i2}\beta_2 \dots x_{ip}\beta_p)]^2$$

This gives estimated weights $\hat{\beta}_0, \dots, \hat{\beta}_p$.

- Finally, our prediction equation for making future predictions at a new x is $\hat{y} = \hat{\beta}_0 + x_1\hat{\beta}_1 + x_2\hat{\beta}_2 \dots x_p\hat{\beta}_p$

A problem

Least squares doesn't work if the number of features p is greater than the number of observations n

Why not?

The Lasso

The **Lasso** is an estimator defined by the following optimization problem:

$$\underset{\beta_0, \beta}{\text{minimize}} \frac{1}{2} \sum_i (y_i - \beta_0 - \sum_j x_{ij} \beta_j)^2 \quad \text{subject to} \quad \sum |\beta_j| \leq s$$

- Penalty \implies sparsity (feature selection)
- Convex problem (good for computation and theory)
- Our lab has written an open-source R language package called **glmnet** for fitting lasso models (Friedman, Hastie, Simon, Tibs). Available on CRAN.
- glmnet v3.0 (just out) now features the *relaxed lasso* and other goodies (like a progress bar!)

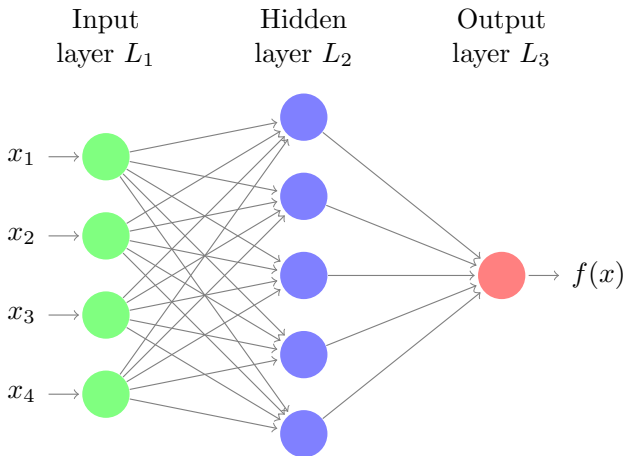
Almost 2 million downloads

The Elephant in the Room: DEEP LEARNING



Will it eat the lasso and other statistical models?

Deep Nets/Deep Learning



Neural network diagram with a single hidden layer. The hidden layer derives transformations of the inputs — nonlinear transformations of linear combinations — which are then used to model the output

How many units of platelets will the Stanford Hospital need tomorrow?



**WE WANT
YOUR GOLD.**

The stuff in your blood, not your bank.

Allison Zemek



Tho Pham



Saurabh Gombar



Leying Guan



Xiaoying Tian



Balasubramanian
Narasimhan

Big data modeling to predict platelet usage and minimize wastage in a tertiary care system

Leying Guan^{a,1}, Xiaoying Tian^{a,1}, Saurabh Gombar^b, Allison J. Zemek^b, Gomathi Krishnan^c, Robert Scott^d, Balasubramanian Narasimhan^a, Robert J. Tibshirani^{a,e,2}, and Tho D. Pham^{b,d,f,2}

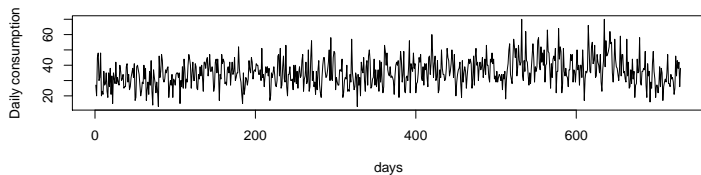
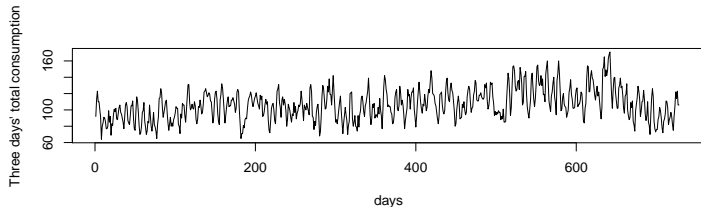
^aDepartment of Statistics, Stanford University, Stanford, CA 94305; ^bDepartment of Pathology, Stanford University, Stanford, CA 94305; ^cStanford for Clinical Informatics, Stanford University, Stanford, CA 94305; ^dStanford Hospital Transfusion Service, Stanford Medicine, Stanford, CA 94305; ^eDepartment of Biomedical Data Science, Stanford University, Stanford, CA 94305; and ^fStanford Blood Center, Stanford Medicine, Stanford, CA

Contributed by Robert J. Tibshirani, August 10, 2017 (sent for review June 25, 2017; reviewed by James Burner, Pearl Toy, and Minh-Ha Tran)

Background

- Each day Stanford hospital orders some number of units (bags) of platelets from Stanford blood center, based on the estimated need (roughly 45 units)
- The daily needs are estimated “manually”
- Platelets have just 5 days of shelf-life; they are safety-tested for 2 days. Hence are **usable for just 3 days**.
- Currently about **1400** units (bags) are wasted each year. That's about **8%** of the total number ordered.
- There's rarely any shortage (shortage is bad but not catastrophic)
- Can we do better?

Data overview



Data description

Daily platelet use from 2/8/2013 - 2/8/2015.

- Response: number of platelet transfusions on a given day.
- Covariates:
 1. **Complete blood count (CBC) data:** Platelet count, White blood cell count, Red blood cell count, Hemoglobin concentration, number of lymphocytes, ...
 2. **Census data:** location of the patient, admission date, discharge date, ...
 3. **Surgery schedule data:** scheduled surgery date, type of surgical services, ...
 4. ...

Notation

y_i : actual PLT usage in day i .

x_i : amount of new PLT that arrives at day i .

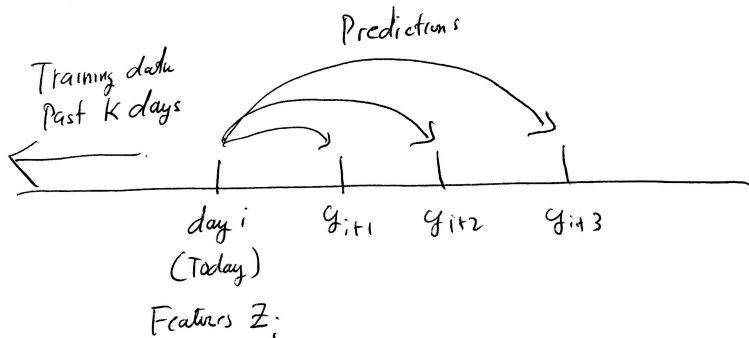
$r_i(k)$: remaining PLT which can be used in the following k days, $k = 1, 2$

w_i : PLT wasted in day i .

s_i : PLT shortage in day i .

- **Overall objective:** waste as little as possible, with little or no shortage

Our first approach



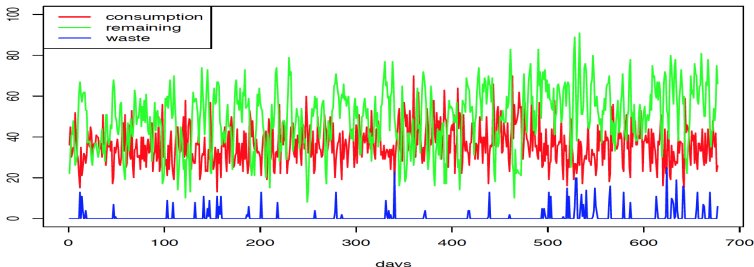
Our approach

- Build a supervised learning model (via lasso) to predict use y_i for next three days
- Use the estimates \hat{y}_i to estimate how many units x_i to order. Add a buffer to predictions to ensure there is no shortage. Do this in a “rolling manner”.

Results

Chose sensible features- previous platelet use, day of week, # patients in key wards.

Over 2 years of backtesting: no shortage, reduces waste from **1400** bags/ year (8%) to just **339** bags/year (1.9%)



Corresponds to a predicted direct savings at Stanford of \$350,000/year. If implemented nationally could result in approximately \$110 million in savings.

Moving forward

- System has just been deployed at the Stanford Blood center (R Shiny app). **But yet not in production**
- We are distributing the software around the world, for other centers to train and deploy
- see Platelet inventory R package
<https://bnaras.github.io/pip/>
- Recently this week we learned that the predictions have been way off for the past month! Data problem? Model problem? Not sure yet.
- **A remaining challenge:** How can we detect if/when the system is no longer working well?

Covidcast: A map of Real-time COVID-19 Indicators

- For the past 5 months, I have been working with Roni Rosenfeld (Chair of ML), Ryan Tibshirani (Statistics+ML) and their **Delphi** flu prediction team at Carnegie Mellon University.
- The Delphi group has been doing influenza forecasting for the CDC for the past five years, as part of the CDC national influenza forecasting challenge.
- They have done well- finishing first in the CDC national influenza forecasting challenge 3 of the past 4 years.
- They were awarded a CDC Center of Excellence in September 2019, along with U. Mass.

Then covid arrived

- In March the CDC asked Delphi (and other groups) to make covid-19 predictions, and this led to the current effort involving about 25 software engineers and statisticians
- The team includes, professors, research assistants and current grad students.
- We launched a national covid-19 indicators map in May and work continues on Phase 2: **forecasting cases and hospitalization usage**



Ryan Tibshirani

Lead Researcher, Delphi COVID-19 Response Team

Associate Professor, Department of Statistics and Machine Learning Department
Carnegie Mellon University



Roni Rosenfeld

Lead Researcher, Delphi COVID-19 Response Team

Professor and Head, Machine Learning Department
School of Computer Science
Carnegie Mellon University



Delphi COVID-19 Response Team

Getting good data is the key

- We needed data on covid-19 symptoms (not just confirmed cases) in order to forecast hospitalization needs
- Ryan approached Google, Facebook, Amazon and other companies, asking them to conduct surveys. He has spent about a month in negotiations.
- Main problem: legal concerns about health data privacy. FB solution: rather than having the survey run by FB on their site, they put just a link to an external survey at CMU.
- With this model, Amazon and Google joined. We have good data so far from Google Surveys (600K/day) and FB (100K/day).

The questions

FaceBook survey

1. In the past 24 hours, have you or anyone in your household had (yes/no for each):
 - a. Fever (of 100 degrees or higher)
 - b. Sore throat
 - c. Cough
 - d. Shortness of breath
 - e. Difficulty breathing
2. How many people in your household (including yourself) are sick (fever, along with at least one other symptom from the above list)?
3. How many people are there in your household in total (including yourself)?
4. What is your current ZIP code?

From this, covid positive is defined as fever and at least one of (cough, shortness of breath, difficulty breathing)

We also measure influenza positive as fever and at least one of (cough, sore throat)

Google asks : *how many people in your community that you know are sick with fever, along with one of sore throat, shortness of breath, cough or difficulty breathing?*

Other data sources

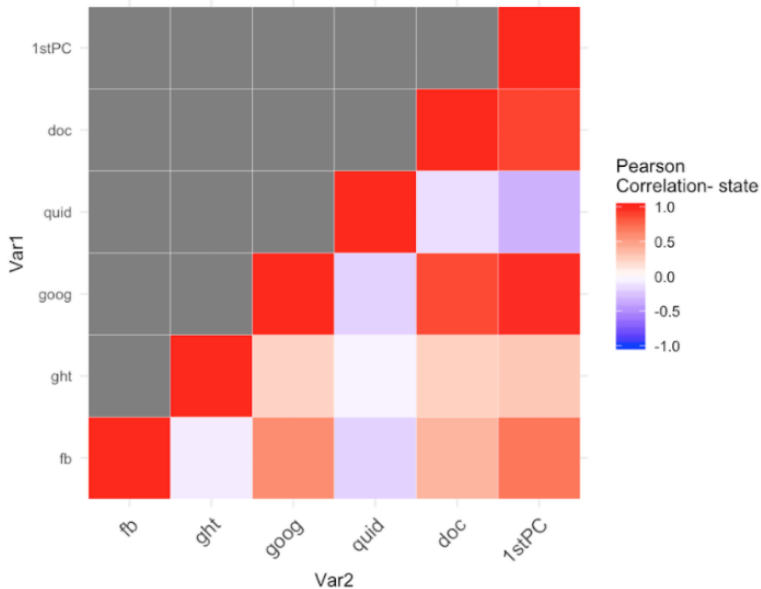
- A major medical/hospital provider is giving us real time data on doctors visits, hospital admissions, and ICU admissions
- Google Health trends is providing estimates of the percentage of Google Searches that relate to covid symptoms
- Quidel (diagnostic healthcare manufacturer) is giving us data on the number of lab tests for influenza

The experience for me

- **Exciting!**: like being in a startup.
- daily zoom meetings, github, slack, a lot of coding
- I had no idea the amount of work involved in such a project; statistical (survey sampling, estimating trends, estimating uncertainty) and lots of software engineering
- The data is complicated and at many resolutions: state, metropolitan area, hospital referral region and county.
- We launched thursday morning. At around 10pm wed, the entire map/site was broken; eventually we figured out it was because someone had decided to change a file name without telling anyone!

SHOW THE MAP <https://covidcast.cmu.edu/>

Correlation and consensus



Things I've learned

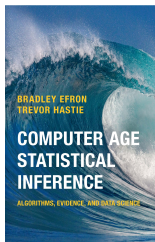
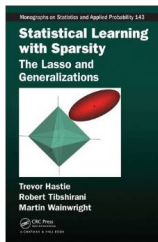
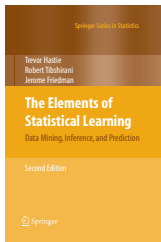
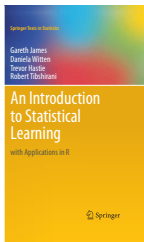
- An greater appreciation for R, R markdown, github and CRAN. Also: we are making substantial use `glmnet`.
- BUT: many CRAN libraries need improvement- numerics, defaults, documentation
- Spend time on this: it's really, really important

The next stage

- Use current signals, as well as data on hospital admissions to forecast hospitalization needs (in each Hospital referral region) a few weeks ahead.
- Try to assess the effects of interventions like Shelter at Home
- Help inform public health officials, to make better decisions. We are working with the California Dept of Public Health and starting to work with Covid Act Now, to supply their forecasts

For further reading

Many of the methods used are described in detail in our books on Statistical Learning: (last one by Efron & Hastie)



All available online for free