# Introduction to Big Data

## Lecture 1

Connor Dowd

March 30th 2021

# A Plea

You are not required to keep your cameras on.

Please do. It helps me immensely to be able to see how people respond as I say things.

If 15% of you are lost, I want to know about it, but with 20 cameras on, I have to notice one of only three faces being lost.

More generally though, please, jump in with questions when you have them. Use the chat if you prefer.

# I'm here to teach you.

- ▶ If you have questions,
- ▶ If you have trouble,
- ▶ If you have feedback for me,
- ▶ If you have ideas for improving the course,

Please let me know.

# Today's Class

1. Introduction
   - Goals
   - Material
   - Syllabus
2. What is Big Data?
3. Computation: R, Big Data
4. Data Viz, Statistics, Dimension Reduction
5. Testing: False Discovery

# Big Data

# Introduction

Big data **is not** an end in itself

We are doing Inference at LARGE scale!

The goal here is to learn what we can trust, how it can be used, and how to learn more.

We want to help you make good decisions with lots of messy data.

By necessity then, this class will be a mix of theory and practice.

# Theory and Practice

You need a solid foundation in statistical principles.

- ▶ We don't want to shoot ourselves in the foot

You also need a hefty dose of 'rules of thumb'.

- ▶ Inefficient procedures can be wildly inefficient
- ▶ **BUT** we don't want to reinvent the wheel.

This is hands on work.

# What is Big Data?

There are a lot of names out there for a large cluster of very similar disciplines:

- ▶ Econometrics
- ▶ Data Science
- ▶ Big Data
- ▶ Statistics
- ▶ Datamining
- ▶ Machine Learning

There are differences in focus and style here. But the similarities are larger than the distinctions.

# What is Big Data?

Big Data as a name, originates with computer scientists working with data too large for any single computer

"Big Data is one GB larger than my RAM".

But often it is associated with administrative data, where statistical notions of sampling error may start to fall apart.

More generally, it is in a nexus where a fairly tight connection between inference and prediction can start to break down.

# What is Big Data?

Big Data is focused on extracting useful truth from large datasets.

- ▶ Infer patterns in high dimensional data
- ▶ Simple and scalable algorithms
- ▶ Honest and humble model selection
- ▶ Manage conflict between "useful" and "true"
- ▶ Make some **decision**

# What is Big Data?

Big in number of observations (size $n$).

Big in number of variables (dimension $p$).

In these settings you *cannot*:

- ▶ Look at each variable and make a decision (t-tests).
- ▶ Choose from a small set of nested models (F-tests).
- ▶ Plot every variable to look for interactions and transformations.

Some of our tools are straight out of Stats 101 (regression, confidence intervals), some are close relatives of Stats 101 (PCA, MSE), some are new beasts (trees, Bagging).

# Why should you care?

- Extremely employable.

# Why should you care?

▶ Extremely employable.



We've found **68764** jobs that could be the right fit!

Your results are listed on the left.

# Why should you care?

- Extremely employable
- Avoid being hoodwinked by nutjobs with numbers.

- Avoid hoodwinking yourself with numbers. *They want to lie to you.*

# Why should you care?

- Extremely employable
- Avoid being hoodwinked by nutjobs with numbers.

- Avoid hoodwinking yourself with numbers. *They want to lie to you*.
- Its fun. Genuinely.

# Does data really want to lie to you? | YES



Figure 1: Influence of Researcher Decisions can reverse results

Admin

# Course Schedule

subject to change. . .

1. Data: Computing, plotting, principles. False Discovery. Loss.
2. Regression: Everything in two days. Linear and Logistic.
3. Model Selection: Penalties, information criteria, cross-validation
4. Treatment Effects: High Dimensional Controls, AB testing, Bootstrap.
5. Classificiation: Multinomials, KNN, sensitivity/specificity
6. Trees: CART, Random Forests, Ensembles
7. Networks?: Co-occurence, directed graphs
8. Clustering?: k-means, mixture models, association rules
9. Factors?: Latent variables, PCA, PLS

# Big Data Team

TA: Yuxiao Li – Booth PhD. She will primarily be grading your work.

Me: Booth PhD. First (and last) time teaching this course.

Office hours: Happy to meet with you. For now, email me if you want to meet. If that proves overwhelming for me, I may set formal office hours.

Structure: Today I'll see how much I get through, and then Thursday will be what we don't get to, as well as discussion of homework 1 (group assignment due next week), and more.

# Course Design

This is also a good time to note that this course is based closely off the course by the same name taught by Professor Veronika Rockova.

But some things will be different.

First time for this course with undergrads. So we don't know what pace makes sense. Feedback is welcome.

# Textbooks

The primary text for this course is "Elements of Statistical Learning" by Hastie, et al.

It is recommended as an excellent reference going forwards. PDFs are (legally) available for free at https://web.stanford.edu/~hastie/Papers/ESLII.pdf and on my website.

I will also try to recommend other textbooks as we go along for specific topics. If there is a topic you're interested in that I haven't recommended a textbook for, email me about it.

# Homeworks

Homeworks will be in small groups of 3. You may not have groups smaller than 3.

There will be ~1 homework a week, graded from 1-5.

1. Oops
2. Eh
3. Fine
4. Good. – Typical score for homework meeting my expectations. (Nothing wrong with it)
5. Nailed it. – Somehow you impressed the TA.

I don't expect there to be many 5s. Literally nothing is wrong with a 4.

There will also be a prediction competition.

# Other Grade Components

1. Participation NEW – I care a lot about asking and answering questions on the canvas discussion board.
2. Midterm – possible, curious how you all feel about a 20 question multiple choice/numeric answer midterm.
3. Take home final

# Analysis

We will be working with real data. The essence of big data is an excellent ability to make predictions using data.

To that end, we will also run a low key prediction competition. For now, there is one thing to predict.

If there is interest in the prediction competition, we will have more questions.

First prediction assignment: US covid vaccination predictions.

# Details

"How many people in the US will have had at least one dose by end of day on April 30th?"

- ▶ Predictions due before class on Thursday.
- ▶ Full details in Assignment on website
  - ▶ Along with 3 CDC datasets over time.
- ▶ You only need to provide 3 numbers:
  - ▶ A prediction
  - ▶ A 90% CI lower bound
  - ▶ A 90% CI upper bound
- ▶ Enter on canvas for full credit
- ▶ Enter on google form to compete against me
  - ▶ And register interest in more questions.

Break

Computation

# R

All of our analysis will be in R.

This is the premier data analysis platform. Other platforms are better at some things, none is better at everything.

Its free, open source, cross platform, very capable.

# R

Used by academics (in statistics, marketing, finance, genetics, etc), companies (EBAY, Google, Microsoft, Boeing, Citadel, IBM, etc), governments (RAND, DOE Labs, US Navy, etc).

R's big strength is that it is open source and has a large community contributing packages. E.g. it is very easy to install a package to do LASSO or other fancy regression tools.

But it can be unpolished. The non-standard packages are as varied as their contributors. Some are super specific, some are super general. Many are useless.

R is not flawless. But it has a large community that will help you overcome flaws.

I **strongly** recommend using the RStudio IDE. (This lecture, your homeworks, much else in this course has been made with it)

# R

The big barrier to entry is the command line interface.

▶ You type a command, R follows the command.

This can mean quite the learning curve. But it is well worth it.

▶ Once you write code for something, it becomes easy to repeat that operation.
▶ This is a critical skill for using modern statistical tools.

I will post all code online. I strongly recommend you look at my solutions to HWs, etc if you run into trouble.

# R – Introduction

Open R (or RStudio). Just like any other program.

In some sense, R is a jumped up calculator. (e.g. *,/,+,-,^,%%)

If you want to save a number, or anything else, give it a name.

```
A = 2
```

If you want to see what a variable is, just type the name:

```
A
```

```
## [1] 2
```

# R – Vectors

We can build vectors with `c()` or "`:`"

```
1:5
```

```
## [1] 1 2 3 4 5
```

```
B = c(1,5,3)
B
```

```
## [1] 1 5 3
```

# R – Vectors

We can grab subsets of vectors (and other groupings) with [ ].

```
B[2]
```

```
## [1] 5
```

And then do math with them. (";" will end a line)

```
D = A+B[2]; D
```

```
## [1] 7
```

Sometimes, its not clear what you want, and R takes a guess:

```
E = A+B
```

# R – Errors

```
E
```

```
## [1] 3 7 5
```

And sometimes you ask it to do something it can't do. So it tells
you it can't do that, and tries to explain itself.

```
G = A+b[2]
```

```
## Error in eval(expr, envir, enclos): object 'b' not found
```

If something isn't working, googling often helps. (e.g. "R error
object not found")

# R – CSVs

We are usually going to work with data in .csv files. Incredibly human-readable files, and easy to move between Excel and other programs.



Figure 2: CSV examples

# R – Loading Data

We can load the data into R using the `read.csv` command. It is
good practice to set a working directory first.

```
setwd("~/Github/bigdata/predictions/")
mar3 = read.csv("us_covid_vaccinations_mar3.csv",skip=3)
```

Now we can play with (some of) the data.

# R – Dataframes

mar3 is a dataframe. We can index columns with names as well as their position. Dataframes also improve on matrices by allowing different variable types in each column.

```
mar3[1,] # The first row
mar3[1:10,] # 10 rows
mar3[,2] # The second column
mar3$Total.Doses.Delivered # The second column
mar3[mar3$Total.Doses.Delivered > 1000,] #Rows w/ >1k dose
nrow(mar3)
ncol(mar3)
dim(mar3)
colnames(mar3)
```

# R – Variable types

R stores different types of data differently.

- ▶ Numeric data are just numbers.
    - ▶ Though they can take several types (integer, float, etc)
- ▶ Factor data take a small number of values.
    - ▶ Think vaccine types, or car models.
- ▶ logical/boolean values are either TRUE or FALSE
- ▶ Characters and strings are words or letters.
    - ▶ "hello". "Alaska". "aagta". "September 1st". "19".

We have many tools for converting between types and identifying
types. is.numeric, factor(), levels().

# R – Misc

Functions look like f(arg1,arg2,...)
And are incredibly useful. e.g. Create new variables:

```
mar3$log.doses = log(mar3$Total.Doses.Delivered)
```

Finding out how a function works is simple ?log will take you to
the help page.

# R – Plotting

Is straightforward for simple plots, but rapidly gets complicated.

```
plot(mar3[,10]~mar3[,14])
```

## R – Plotting

Histograms.

```
hist(mar3$Doses.Delivered.per.100K)
```

## Error in hist.default(mar3$Doses.Delivered.per.100K): 'x

```
typeof(mar3$Doses.Delivered.per.100K)
```

## [1] "character"

```
head(mar3$Doses.Delivered.per.100K)
```

## [1] "50950" "29803" "31912" "52883" "30842" "N/A"

# R – Plotting

```
mar3$Doses.Delivered.per.100K =
  as.numeric(mar3$Doses.Delivered.per.100K)
hist(mar3$Doses.Delivered.per.100K)
```

**Histogram of mar3$Doses.Delivered.per.100K**

# R – Terrible Colnames

```
colnames(mar3)[10:15]
```

```
## [1] "Percent.of.Total.Pop.with.1..Doses.by.State.of.Resi
## [2] "People.18..with.1..Doses.by.State.of.Residence"
## [3] "Percent.of.18..Pop.with.1..Doses.by.State.of.Reside
## [4] "People.with.2.Doses.by.State.of.Residence"
## [5] "Percent.of.Total.Pop.with.2.Doses.by.State.of.Resid
## [6] "People.18..with.2.Doses.by.State.of.Residence"
```

```
colnames(mar3) = c("State","Delivered","Delivered.100k",
    "X18.Delivered.100k","Administered","Administered.100k",
    "X18.Administered","X18.Administered.100k","OneDose",
    "Perc.OneDose","X18.OneDose","X18.Perc.OneDose",
    "TwoDose","Perc.TwoDose","X18.TwoDose","X18.Perc.TwoDose
```

# R – Plotting - In color?

```
plot(mar3$Delivered.100k/1000,mar3$Perc.OneDose,
     col=(mar3$Perc.TwoDose > 8.2)+1,
     main = "One Dosed vs Delivered Doses per 100 people")
```



**One Dosed vs Delivered Doses per 100 people**

# R – Scatterplots

Scatterplots are fundamental to statistics. If you can show a compelling scatterplot, you can convince people of anything.

In part, this is because they are able to contain nearly as many dimensions of information as we can interpret. We tend to max out at 3-4 dimensions.

Critical to reduce dimensions down to interpretable level for good visualization.

# R – Tidyverse

```
library(tidyverse)
ggplot(mar3,aes(x=Delivered.100k/1000,y = Perc.OneDose,col=
```

# Statistics

## R – Regression

```
fit = lm(Perc.OneDose~Delivered.100k+Perc.TwoDose, data=mar
summary(fit)

##
## Call:
## lm(formula = Perc.OneDose ~ Delivered.100k + Perc.TwoDos
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -3.3031 -1.1608 -0.3072  0.9775  3.9143
##
## Coefficients:
##                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)     6.469e+00  8.593e-01   7.528 4.20e-10 ***
## Delivered.100k  9.691e-06  3.525e-05   0.275    0.784
## Perc.TwoDose    1.138e+00  1.468e-01   7.752 1.77e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

# Regression

The basic model is as follows:
$$E[Perc.OneDose] = \beta_0 + \beta_1 Delivered.100k + \beta_2 Perc.TwoDose$$

Or, similarly:

$$Perc.OneDose = \beta_0 + \beta_1 Delivered.100k + \beta_2 Perc.TwoDose + \epsilon$$
Where $E[\epsilon] = 0$.

# Hypothesis Testing

Is "Delivered.100k" important?

H0: $\beta_1 = 0$.

```
summary(fit)
```

```
##
## Call:
## lm(formula = Perc.OneDose ~ Delivered.100k + Perc.TwoDos
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -3.3031 -1.1608 -0.3072  0.9775  3.9143
##
## Coefficients:
##                  Estimate Std. Error t value Pr(>|t|)
## (Intercept)     6.469e+00  8.593e-01   7.528 4.20e-10 ***
## Delivered.100k  9.691e-06  3.525e-05   0.275    0.784
## Perc.TwoDose    1.138e+00  1.468e-01   7.752 1.77e-10 ***
##
```

# A Single Test

The "p-value" is the probability of observing an event at least as extreme as the observed event, assuming that the null hypothesis is true.

**Testing Procedure** Choose a cutoff $\alpha$. Conclude that there is a real effect if $p < \alpha$.

This procedure tries to give only $\alpha$ probability of a false positive.

$$P(p < \alpha | H0) = \alpha$$

# Challenges

Big Data situations pose two different problems for p-values.

1. With a big enough sample (i.e. large $n$), we wind up rejecting the null even for effects that are of no practical significance, but which have a small positive $\beta$.
2. With many parameters (i.e. large $p$), we are performing many tests, and need to account for the number of tests we perform.

# Large Scale Testing

We wish to test $p$ simultaneous null hypothesis:

$$H0_1, H0_2, ..., H0_p$$

Out of the $p$ null hypothesis, $N_0$ are true nulls and $N_1 = p - N_0$ are false – i.e. there is an effect.

| | | Decision | | |
|---|---|---|---|---|
| | | Fail to Reject | Reject | Sum |
| **Truth** | Noise | Real non-Discovery (TN) | False Discovery (FD) | N_0 |
| | Signal | Missed Discovery. (FN) | Real Discovery (TD) | N_1 |
| | Sum | p-R | R | |

# Multiplicity

$\alpha$ is typically for a single test. If we repeat $p$ tests, we expect $p\alpha$ significant even when all the null hypotheses are true.

Suppose we have 100 regression coefficents, of which 5 are true effects. Now suppose we test all 100 coefficients, all 5 true effects come up (if we're lucky), as well as $95\alpha \approx 5$ more. So of the significant variables, 50% are fake.

This is the "false discovery proportion".

It depends closely on $\alpha$, the underlying truth, and how good our methods are. But it can be really big.

# False Discovery Rate

Big data is about making choices.
Instead of single tests, we'll consider:

FD Proportion = False positives / Significant = $\frac{FD}{R}$

FDP is a property of the fitted model and the truth. We don't know it.

We can control its expectation though: False Discovery Rate, $FDR = E[FDP]$.

This is the multivariate analogue of $\alpha$.

If all tests are tested at $\alpha$ level, we have $\alpha = E[FD/N_0]$, whereas $FDR = E[FD/R]$

# FDR Control

Suppose we want to know that $FDR \leq q = 0.1$.

Benjamini + Hochberg Algorithm

1. Rank your p-values smallest to largest.
2. Set p-value cutoff as $\alpha^* = max\{p_{(k)} : p_{(k)} \leq q\frac{k}{p}\}$

Then $FDR \leq q$ – assuming approximate independence between tests.

Wrap up

# Things to do

Before Thursday:

- ▶ Group up on canvas
- ▶ Submit predictions
- ▶ Play with R, tidyverse, ggplot2

# Rehash

- I'm Connor Dowd
- Big Data brings new problems
- R is important
- Graphs are important
- False Discovery Rates can be controlled

Bye!