

Statistics 101C Discussion Week 1

Brice Randolph

University of California, Los Angeles

brandolph@ucla.edu

March 29, 2016

Agenda



What will we be doing in discussion section?

- Lecture Key Points
- Probability Review (Confidence intervals today)
- Interview Question(s)
- Homework questions
- Something new or an extension upon a model seen in class

Can you ditch discussion? (8 am seriously?)

- I'm going to make this discussion as useful as possible (applications & interview questions)
- Those who come to the 8 am as opposed to the 9 am get free coffee (first come first served)
- Email policy (brandolph@ucla.edu)
- Office hours (Poll?)

The Dude Recommends

Some things I have learned since graduating:



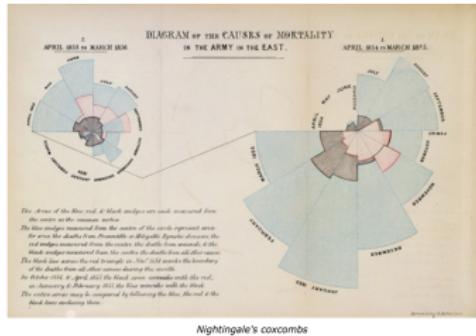
Github - take a few minutes to sign up.

Subscribe to a couple **data science blogs**

- <http://www.r-bloggers.com/>
- <http://fivethirtyeight.com/>
- <http://fastml.com/>
- <https://www.quora.com/What-are-the-best-machine-learning-blogs-or-resources-available>

The Dude Recommends

Think of your favorite example of a data visualization
(mine's the coxcomb -1856 Florence Nightingale)



Think of your favorite example of a data journalism
(Las Vegas Sun's "Do No Harm" about poor healthcare in south Nevada)

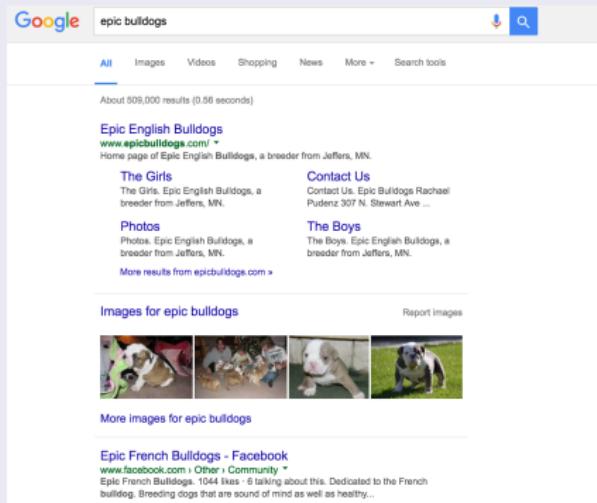
Can't think of one? Look one up and share it with your neighbor
(seriously, Go!)

Overview

What are statistical models used for?

Forecasting / Prediction:

- Stock price
- Likelihood of contracting certain diseases
- Webpage ranking



A screenshot of a Google search results page for the query "epic bulldogs". The search bar at the top shows the query. Below it, the "All" tab is selected, along with other options like Images, Videos, Shopping, News, More, and Search tools. The search results section starts with a snippet for "Epic English Bulldogs" from their website, mentioning they are a breeder from Jeffers, MN. It includes links for "The Girls", "Contact Us", and "Photos". Below this, there's a section for "The Boys" and a link to "More results from epicbulldogs.com". Further down, there's a "Images for epic bulldogs" section with four thumbnail images of bulldogs and a "Report Images" link. At the bottom of the snippet, there's a link to "More images for epic bulldogs".

Google

epic bulldogs

All Images Videos Shopping News More Search tools

About 509,000 results (0.56 seconds)

Epic English Bulldogs
www.epicbulldogs.com/ *

Home page of Epic English Bulldogs, a breeder from Jeffers, MN.

The Girls
The Girls, Epic English Bulldogs, a breeder from Jeffers, MN.

Contact Us
Contact Us, Epic Bulldogs Rachael Pudenz 307 N. Stewart Ave ...

Photos
Photos, Epic English Bulldogs, a breeder from Jeffers, MN.

The Boys
The Boys, Epic English Bulldogs, a breeder from Jeffers, MN.

More results from epicbulldogs.com »

Images for epic bulldogs

Report Images



More images for epic bulldogs

Epic French Bulldogs - Facebook
www.facebook.com/OtherCommunity *

Epic French Bulldogs, 1044 likes · 6 talking about this. Dedicated to the French bulldog. Breeding dogs that are sound of mind as well as healthy...

Other Applications

Computer Vision/ Artificial Intelligence

- Object tracking
- Sentiment Analysis

Sentiment Analysis with Python NLTK Text Classification

This is a demonstration of sentiment analysis using a [NLTK 2.0.4](#) powered text classification process. It can tell you whether it thinks the text you enter below expresses **positive sentiment**, **negative sentiment**, or if it's neutral. Using **hierarchical classification**, neutrality is determined first, and sentiment polarity is determined second, but only if the text is not neutral.

Analyze Sentiment

Language: English

Enter text:

Have you been on Tinder for very long? I'm super excited to meet you and tell you about Neural Nets! I think you look great in that mirror selfie pick from the 90s.

Enter up to 50000 characters

Analyze

Sentiment Analysis Results

The text is pos.

The final sentiment is determined by looking at the classification probabilities below.

Subjectivity

- neutral: 0.1
- polar: 0.9

Polarity

- pos: 0.6
- neg: 0.4

Other Applications

Image classification/Scene Understanding

clarifai

ABOUT

PRICING

DEVELOPER

BLOG

CONTACT US

SIGN UP NC



Predicted Tags

no person ocean sea water
sand beach travel vehicle
seashore people

Similar Images



Image classification/Scene Understanding

[ABOUT](#)[PRICING](#)[DEVELOPER](#)[BLOG](#)[CONTACT US](#)

*By using the demo you agree to our [terms of service](#)



Predicted Tags

portrait people woman one
adult girl facial expression
music fashion young

Your turn

I'm going to play a song for 30 seconds and I want you to tell me who sings it.

You can use any machine learning algorithms you have on you...

Shazam

Shazam uses a smartphone or computer's built-in microphone to gather a brief sample of audio being played. It creates an acoustic fingerprint based on the sample and compares it against a central database for a match. If it finds a match, it sends information such as the artist, song title, and album back to the user.

Final Exam (Competition)

Goal: **Predict** some outcome **Given** some predictors

Data: Euthanizing Animals (Predicting if an animal will be killed?)

Kaggle: You will upload your predictions here (go check out the website)

Team with best score wins

Machine learning in one sentence

Build a model to fit some process(static or dynamic) for prediction or inference.

You have learned some **pretty good** models already. Can you name some?

Spaghetti Model from class?

Why do I say they are good when they are so simple?

Why might these be used in industry? Think audience...

General Process Summary by Galit Shmueli



In Statistics, like in Data Mining, you start with data and a goal. In statistics there is a lot of focus on inference, that is, answering population-level questions using a sample. In data mining the focus is usually prediction: you create a model from your sample (training data) in order to predict test data.

6



The process in statistics is then:

1. Explore the data using summaries and graphs - depending on how data-driven the statistician, some will be more open-minded, looking at the data from all angles, while others (especially social scientists) will look at the data through the lens of the question of interest (e.g., plot especially the variables of interest and not others)
1. Choose an appropriate statistical model family (e.g., linear regression for a continuous Y, logistic regression for a binary Y, or Poisson for count data), and perform model selection
2. Estimate the final model
3. Test model assumptions to make sure they are reasonably met (different from testing for predictive accuracy in data mining)
4. Use the model for inference -- this is the main step that differs from data mining. The word "p-value" arrives here...

Take a look at any basic stats textbook and you'll find a chapter on Exploratory Data Analysis followed by some distributions (that will help choose reasonable approximating models), then inference (confidence intervals and hypothesis tests) and regression models.

I described to you the classic statistical process. However, I have many issues with it. The focus on inference has completely dominated the fields, while prediction (which is extremely important and useful) has been nearly neglected. Moreover, if you look at how social scientists use statistics for inference, you'll find that they use it quite differently! You can check out more about this [here](#)

share improve this answer

answered Nov 10 '10 at 14:38

community wiki

Galit Shmueli



ggplot2

Please download **ggplot2** if you haven't already

```
install.packages("ggplot2")
```

I will also be using *ggmap* *dplyr* *ggvis* and *shiny* at times

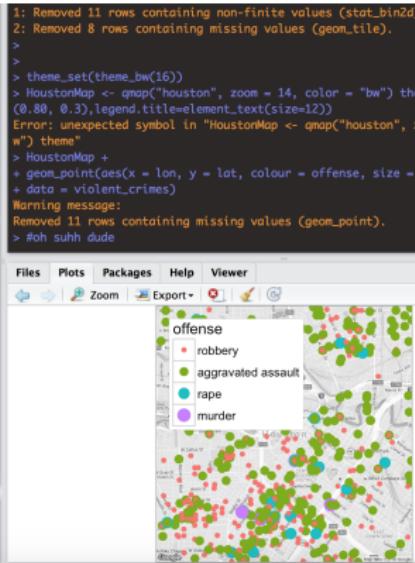
Graphics

Only recently have graphical tools been easy to create.

They had only been used at the end of an analysis to summarize the report.

Take advantage of it!

```
107 library(ggplot2)
108 library(ggmap)
109 library(dplyr)
110 # find a reasonable spatial extent
111 qmap("houston", zoom = 13)
112 gglocator(2)
113
114 str(crime)
115 # only violent crimes
116 violent_crimes <- subset(crime,
117   offense != "auto theft" & offense != "theft" & offense != "vandalism")
118 # order violent crimes
119 violent_crimes$offense <- factor(violent_crimes$offense,
120   levels = c("robbery", "aggravated assault",
121             "murder", "rape"))
122 # restrict to downtown
123 violent_crimes <- subset(violent_crimes, -95.39681 <= lon & lon <= -95.34188
124           29.73631 <= lat & lat <= 29.78400)
125
126 theme_set(theme_bw(16))
127 HoustonMap <- qmap("houston", zoom = 14, color = "bw") theme(legend.position =
128   "bottom")
129 HoustonMap +
130   geom_point(aes(x = lon, y = lat, colour = offense, size = offense),
131             data = violent_crimes)
132 HoustonMap +
133   stat_bin2d(
134     aes(x = lon, y = lat, colour = offense, fill = offense),
135     size = .5, bins = 30, alpha = 1/2,
136     data = violent_crimes
137   )
135:4 ggplot2 / ggmap example 3 R Script
Environment History
To Console To Source
data = violent_crimes)
```



Visualizations

Check this out: one of the more advanced machine learning tools

<http://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

Confidence intervals

Can someone tell me what a confidence interval might be used for?

Confidence intervals in the wild

Given the observed error (accuracy) of a model over a limited sample of data, how well does this error characterize its accuracy over additional instances?

Suppose we have

- a learned model h
- a test set S containing n instances drawn independently of one another and independent of h
- $n \geq 30$
- h makes r errors over the n instances

our best estimate of the error of h is

$$\text{error}_S(h) = \frac{r}{n}$$

Confidence intervals on error

With approximately $N\%$ probability, the true error lies in the interval

$$\text{error}_s(h) \pm z_N \sqrt{\frac{\text{error}_s(h)(1 - \text{error}_s(h))}{n}}$$

where z_N is a constant that depends on N (e.g. for 95% confidence, $z_N = 1.96$)

Look familiar?

These last two examples came from here:

<http://pages.cs.wisc.edu/~dpage/cs760/evaluating.pdf>

Practice Interview Questions

Use R: Write some code that will simulate flipping a fair coin 100 times.

What are the basic assumptions used in linear regression? How do you empirically assess if errors are normally distributed? What are the most common estimation techniques for linear regression?

Any final questions/concerns/music recommendations?