

TCSS 588 Bioinformatics

Ka Yee Yeung

kayee@uw.edu

Institute of Technology, UW-Tacoma

3/28/2018

1

Overview

- Jupyter notebooks
- Why do we care about statistics? What can we do with predictive models?
- Review on statistics
- Regression
- Building predictive models

2

Life cycle of an academic project

1. Individual exploratory work
2. Collaborative development
3. Parallel (scale up using cloud computing)
4. Publication and communication (reproducibly!)

<https://www.slideshare.net/mbussonn/jupyter-a-platform-for-data-science-at-scale>

Data science notebooks

- Mid 1980's: Start of computational notebooks: Matlab, Mathematica notebooks, Maple worksheets
 - GUI allowed for the interactive creation and editing of notebook documents that contain pretty-printed program code, formatted text
- 2010-2011: IPython web-notebooks and prototype
- 2014: Project Jupyter is a spinoff project from IPython

<https://www.datacamp.com/community/blog/ipython-jupyter#gs.wngbPCE>



The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and explanatory text.

The name "Jupyter" is actually short for "Julia, Python and R"

".ipynb" files: JSON (JavaScript Object Notation) based files embedding input and output

<http://jupyter.org/>

Kernels

- A kernel is a program that runs and introspects the user's code: it provides computation and communication with the frontend interfaces, such as notebooks.
- The Jupyter Notebook Application has three main kernels: the IPython, IRkernel and IJulia kernels.
- Community maintained kernels: Ruby, Javascript, Scala, Perl, Octave etc.
- <https://github.com/jupyter/jupyter/wiki/Jupyter-kernels>

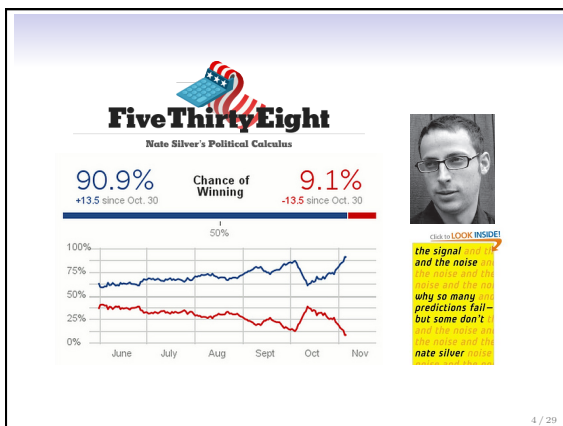
Go to your notebook now

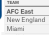
7

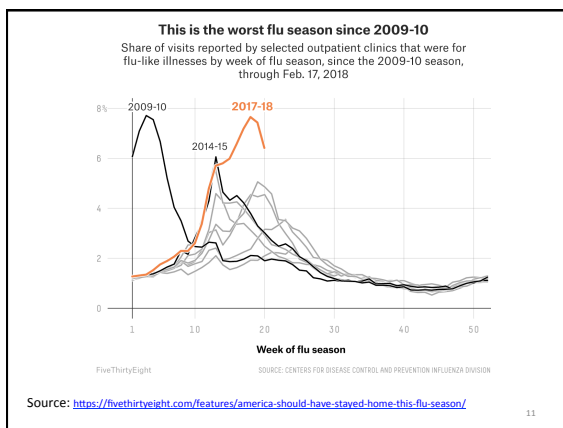
Overview

- Jupyter notebooks
- Why do we care about statistics? What can we do with predictive models?
- Review on statistics
- Regression
- Building predictive models

8



<div>  <div> FiftyThirdEight's NFL Elo Playoff Odds Dec. 23, 2014 </div> </div>									
Team	Rank	W/L	W/L %	W/L %	W/L %	W/L %	W/L %	W/L %	W/L %
AFC East									
New England	1730	12	6	2.2	0.3	0.75	100%	100%	25%
Miami	1524	8	8	7.2	0.0	0.0	0%	0%	0%
Buffalo	1393	6	2	2.0	0.0	0.0	0%	0%	0%
N.Jets	1357	2	12	0.0	0.0	-14	0%	0%	0%
AFC North									
Pittsburgh	1803	10	8	5.4	0.0	-82	61%	100%	2%
Cincinnati	1594	10	4	4.0	0.0	-27	39%	100%	2%
Baltimore	1574	9	8	6.2	0.0	-108	0%	0%	0%
Cleveland	1384	7	2	8.0	0.0	-39	0%	0%	0%
AFC South									
Indianapolis	1488	10	8	7.2	0.0	-82	100%	100%	2%
Jacksonville	1484	9	8	7.2	0.0	-71	0%	0%	0%
Jacksonville	1083	2	12	0.0	0.0	199	0%	0%	0%
Tennessee	1082	2	12	0.0	0.0	197	0%	0%	0%
AFC West									
Denver	1887	11	4	1.0	0.0	-118	100%	100%	11%
San Diego	1544	8	4	6.0	0.0	-10	0%	0%	0%
Kansas City	1320	8	6	7.0	0.0	-82	0%	0%	0%
Seattle	1211	2	12	0.0	0.0	-14	0%	0%	0%
NFC East									
Dallas	1600	11	8	4.2	0.0	-94	100%	100%	1%
Philadelphia	1536	9	8	5.0	0.0	-45	0%	0%	0%
N.Y. Giants	1487	6	5	5.0	0.0	-11	0%	0%	0%
Washington	1311	4	12	0.0	0.0	-121	0%	0%	0%
NFC North									
Green Bay	1837	11	7	4.3	0.0	-134	67%	100%	7%
Detroit	1575	11	3	7.0	0.0	-43	23%	100%	10%
Chicago	1451	8	7	8.0	0.0	-18	0%	0%	0%
Carolina	1381	5	10	3.0	0.0	105	0%	0%	0%
NFC South									
Carolina	1489	8	4	6.0	0.0	-69	43%	100%	43%
Atlanta	1479	6	8	8.0	0.0	-3	0%	0%	0%
New Orleans	1462	8	6	8.0	0.0	-22	0%	0%	0%
Tampa Bay	1424	2	14	0.0	0.0	-134	0%	0%	0%
NFC West									
Seattle	1784	11	4	1.0	0.0	-140	64%	100%	31%
Arizona	1881	11	5	4.0	0.0	-14	0%	0%	0%
San Francisco	1842	7	5	8.0	0.0	-37	0%	0%	0%
St. Louis	1811	5	8	8.0	0.0	-39	0%	0%	0%



Motivation

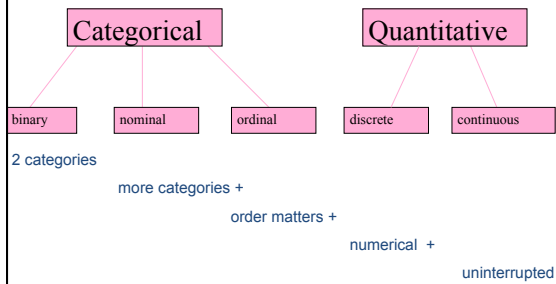
Use data to identify **relationships** among variables and use these relationships to make **predictions**.

Looking at your data

- How are the data distributed?
 - Where is the center?
 - What is the range?
 - What is the shape of the distribution (e.g., Gaussian, uniform)?
- Are there “outliers” ?
- Are there data points that don't make sense?

13

Types of Variables: Overview



14

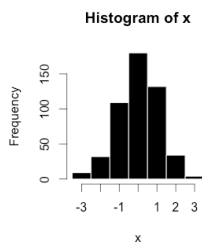
Continuous variables

- Histograms
- Box plots

15

Histogram

bin	count
-3.5	9
-2.5	32
-1.5	109
-0.5	180
0.5	132
1.5	34
2.5	4
3.5	9



Toy example from wikipedia: <http://en.wikipedia.org/wiki/Histogram>

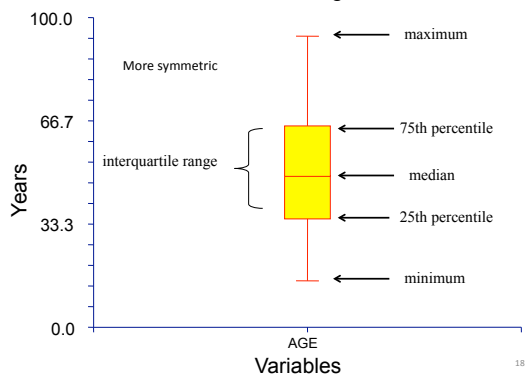
16

In-class exercise

- Predicting adverse drug reactions
<http://www.maayanlab.net/crowdsourcing/megatask1.php>
- We will re-visit this task over this course.
- Basic data types in R:
<http://www.statmethods.net/input/datatypes.html>

17

Box Plot: Age



18

Boxplots in R

- <http://www.r-bloggers.com/box-plot-with-r-tutorial/>
- <https://stat.ethz.ch/R-manual/R-devel/library/graphics/html/boxplot.html>
(advanced, more on this later)

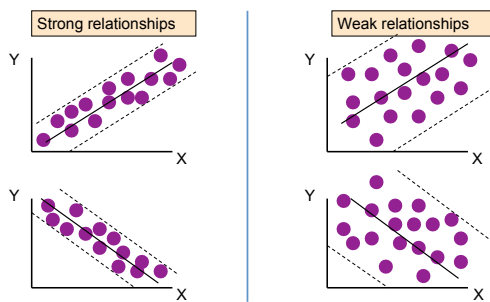
19

Correlation

- Measures the relative strength of the *linear* relationship between two variables
- Ranges between -1 and 1
- The closer to -1 , the stronger the negative linear relationship
- The closer to 1 , the stronger the positive linear relationship
- The closer to 0 , the weaker any positive linear relationship

20

Linear Correlation



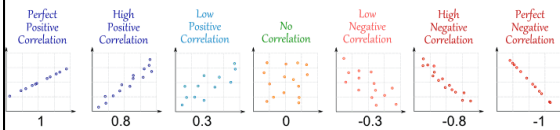
Slide from: Statistics for Managers Using Microsoft® Excel 4th Edition, 2004 Prentice-Hall

21

Pearson's Correlation Coefficient

$$\rho = \frac{\sum_{j=1}^n (A[j] - \mu_A)(B[j] - \mu_B)}{\sqrt{\sum_{j=1}^n (A[j] - \mu_A)^2 \sum_{j=1}^n (B[j] - \mu_B)^2}}$$

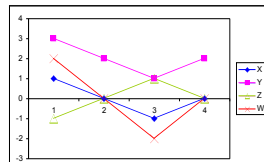
$$\mu_A = \frac{\sum_{j=1}^n A[j]}{n}$$



22

Examples

X	1	0	-1	0
Y	3	2	1	2
Z	-1	0	1	0
W	2	0	-2	0



Correlation (X,Y) = 1

Correlation (X,Z) = -1

Correlation (X,W) = 1

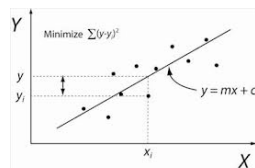
In-class exercise.

23

Linear regression

Regression analysis describes the relationship between two (or more) variables.

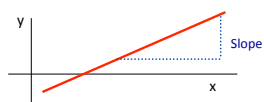
In correlation, the two variables are treated as equals. In regression, one variable is considered **independent** (=predictor) variable (X) and the other the **dependent** (=outcome, response) variable Y.



24

Simple linear regression

- Relation between 2 continuous variables



$$y = \alpha + \beta_1 x_1$$

- Regression coefficient β_1
 - Measures association between y and x
 - Amount by which y changes on average when x changes by one unit
 - Least squares method
- Residual = actual data point – fitted value

25

Predicted value for an individual...

$$\hat{y}_i = \underbrace{\alpha + \beta x_i}_{\text{Fixed – exactly on the line}} + \underbrace{\text{random error}_i}_{\text{Follows a normal distribution}}$$

Fixed –
exactly
on the
line

Follows a normal
distribution

26

Assumptions (or the fine print)

- Linear regression assumes that...
 - The relationship between X and Y is linear
 - Y is distributed normally at each value of X
 - The variance of Y at every value of X is the same (homogeneity of variances)
 - The observations are independent

27

Simple linear regression using a single predictor X .

- We assume a model

$$Y = \beta_0 + \beta_1 X + \epsilon,$$

where β_0 and β_1 are two unknown constants that represent the *intercept* and *slope*, also known as *coefficients* or *parameters*, and ϵ is the error term.

- Given some estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ for the model coefficients, we predict future sales using

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x,$$

where \hat{y} indicates a prediction of Y on the basis of $X = x$. The *hat* symbol denotes an estimated value.

4 / 48

Estimation of the parameters by least squares

- Let $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ be the prediction for Y based on the i th value of X . Then $e_i = y_i - \hat{y}_i$ represents the i th *residual*
- We define the *residual sum of squares* (RSS) as

$$\text{RSS} = e_1^2 + e_2^2 + \dots + e_n^2,$$

or equivalently as

$$\text{RSS} = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2.$$

- The least squares approach chooses $\hat{\beta}_0$ and $\hat{\beta}_1$ to minimize the RSS. The minimizing values can be shown to be

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2},$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where $\bar{y} \equiv \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{x} \equiv \frac{1}{n} \sum_{i=1}^n x_i$ are the sample means.

5 / 48

Credits

- Textbook Ch. 3
- Inspired by slides from Trevor Hastie, Robert Tibshirani, Kristin Sainani, Steven Buechler
- Data: <http://www.maayanlab.net/crowdsourcing/megatask1.php>

30
