# Exploratory Data Analysis with R

Matthew Renze

@matthewrenze

#KCDC15

# Motivation

*The ability to take data—to be able to understand it, to process it, to extract value from it, to visualize it, to communicate it—that's going to be a hugely important skill in the next decades, … because now we really do have essentially free and ubiquitous data. So the complimentary scarce factor is the ability to understand that data and extract value from it.*

Hal Varian, Google's Chief Economist
The McKinsey Quarterly, Jan 2009

The Economist

**The data deluge**

AND HOW TO HANDLE IT A 14-PAGE SPECIAL REPORT

The New York Times

**For Today's Graduate, Just One Word: Statistics**

By STEVE LOHR
Published: August 5, 2009

MOUNTAIN VIEW, Calif. — At Harvard, Carrie Grimes majored in anthropology and archaeology and ventured to places like Honduras, where she studied Mayan settlement patterns by mapping where artifacts were found. But she was drawn to what she calls "all the computer and math stuff" that was part of the job.

- TWITTER
- LINKEDIN
- COMMENTS (58)
- SIGN IN TO E-MAIL

**AVERAGE SALARY FOR High Paying Skills and Experience**

| SKILL | 2013 | YR/YR CHANGE |
|---|---|---|
| R | $ 115,531 | n/a |
| NoSQL | $ 114,796 | 1.6% |
| MapReduce | $ 114,396 | n/a |
| PMBok | $ 112,382 | 1.3% |
| Cassandra | $ 112,382 | n/a |
| Omnigraffle | $ 111,039 | 0.3% |
| Pig | $ 109,561 | n/a |
| SOA (Service Oriented Architecture) | $ 108,997 | -0.5% |
| Hadoop | $ 108,669 | -5.6% |
| Mongo DB | $ 107,825 | -0.4% |

Source: Dice 2014 Tech Salary Survey Results

# A Flood of Data is Coming...



## Sink     or     Swim

# Overview

- Introduction to R
- Data Munging
- Descriptive Statistics
- Data Visualization
- Beyond R and EDA

# How Does This Apply to Me?

- As a software developer, I often:
  - ☑ Perform log file analysis
  - ☑ Analyze software performance
  - ☑ Analyze code metrics for code quality
  - ☑ Detect anomalies in source data
  - ☑ Transform or clean data files to make them usable
  - ☑ Help decision makers make decisions based on data

# About Me

- Independent software consultant
- Education
  - B.S. in Computer Science (ISU)
  - B.A. in Philosophy (ISU)
- Training
  - Kimball Group – Data Warehousing
  - ESRI - ArcGIS, ArcSDE, and ArcGIS Server
  - Data Science Specialization
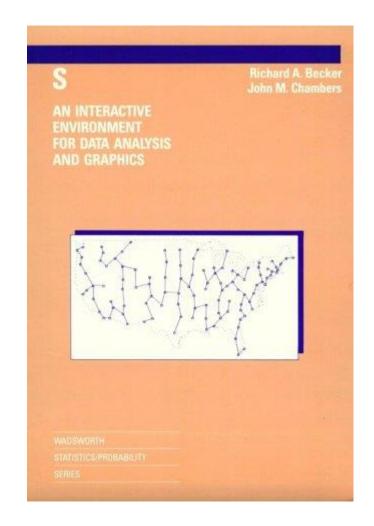    (Johns Hopkins) [In progress]

**Titanium Sponsors**

DATA BANK powered by actifio

epIQ SYSTEMS

VALOREM CONSULTING
WWW.VALOREMCONSULTING.COM
REGISTER TO WIN AN XBOX!

PAIGE TECHNOLOGIES
INTELLIGENT PAIRING. PERPETUAL SUCCESS.

**Platinum Sponsors**

perceptivesoftware from Lexmark

ADAPTIVE SOLUTIONS GROUP

KU EDWARDS CAMPUS
The University of Kansas

pinsight MEDIA+

VinSolutions
Make every connection count.

MS MULTI SERVICE
INNOVATION WHERE IT MATTERS.

**Gold Sponsors**

DST SYSTEMS

LRS Consulting Services

NETCHEMIA
Transforming the way education works

ADVANTAGE TECH

KEYHOLE SOFTWARE

Balance Innovations

stackify

GARMIN

NEW DIRECTIONS
TOGETHER IS THE WAY FORWARD

Bradford & Galt
CONSULTING SERVICES

OAKWOOD

Cerner

jack henry & ASSOCIATES INC.

UnitedLex

dsi

FitBark

Commerce Bank

TRIPLE-I

BATS
Making Markets Better

LIFERAY

imodules

eccoselect
PEOPLE | PROJECTS | PERFORMANCE

CENTRIQ TRAINING

# Introduction to R

# What is R?

- R is an open source implementation of S

# What is S?

- Statistical language
- Bell Labs in 1976
- Owned by TIBCO

# A Brief History of R

- 1991 - developed by:
  - Ross Ihaka
  - Robert Gentleman
- 1995 - open source
- 2000 - v.1.0 released
- Today, R is at v.3.2.0



Source: https://www.stat.auckland.ac.nz/~ihaka/downloads/the-r-project.pdf



Source: www.auklandlifestyle.com

# What is R?

- Open source
- Implementation of S
- Language and environment
- Numerical and graphical
- Cross platform



Source: www.r-project.org

# What is R?

- Active development
- Large user community
- Modular and extensible
- 6700+ extensions


and best of all…

FREE

FREE

RedMonk Q314 Programming Language Rankings

RedMonk

Popularity Rank on Stack Overflow (by # of Tags)

Popularity Rank on GitHub (by # of Projects)

R Console

```
> box()

> title(main= "The Level of Interest in R", font.main=4, col.main="red")

> title(xlab= "1996", col.lab="red")

> ## A filled histogram, showing how to change the font used for the
> ## main title without changing the other annotation.
>
> par(bg="cornsilk")

> x <- rnorm(1000)

> hist(x, xlim=range(-4, 4, x), col="lavend
Waiting to confirm page change...

> title(main="1000 Normal Random Variates",

> ## A scatterplot matrix
> ## The good old Iris data (yet again)
>
> pairs(iris[1:4], main="Edgar Anderson's I
Waiting to confirm page change...
```

Click or hit ENTER for next page



1000 Normal Random Variates

# Code Demo

# Data Munging

# Data Munging

- Transforming data

- Raw data to usable data

- Data must be cleaned first



Source: Wikimedia

# Data Munging Tasks

- Renaming variables

- Data type conversion

- Encoding, decoding, or recoding data

- Merging data sets

- Transforming data

- Handling missing data (imputing)

- Handling anomalous values

# Loading Data in R

- File-based data

- Web-based data

- Databases

- Statistical data

- And many more...

# Cleaning Data

- This step is often the:
  - Most difficult
  - Most time consuming
- TIP: Record all steps



Source: Wikimedia

# Code Demo:
# Lending Club Dataset

- Peer-to-peer loans

- *Problem:* Data are not ready for analysis

- *Goal:* Prepare the data for analysis



Source: www.lendingclub.com

# Code Demo

# Descriptive Statistics

# Descriptive Statistics

- Describe data
- Provides a summary
- aka: Summary statistics

| Interest Rate | |
|---|---|
| Statistic | Value |
| Minimum | 5.42 |
| 1st Quartile | 10.16 |
| Median | 13.11 |
| Mean | 13.07 |
| 3rd Quartile | 15.80 |
| Maximum | 24.89 |
| | |
| Variance | 17.45 |
| Standard Deviation | 4.17 |

# Statistical Terms

- Observations
- Variables
- Qualitative variable
- Quantitative variable

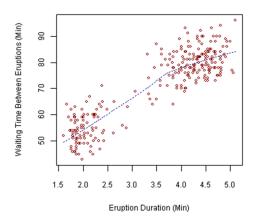| ID | Date | Customer | Product | Quantity |
|----|------|----------|---------|----------|
| 1 | 2012-10-27 | John | Pizza | 2 |
| 2 | 2012-10-27 | John | Soda | 2 |
| 3 | 2012-10-27 | Jill | Salad | 1 |
| 4 | 2012-10-27 | Bob | Milk | 1 |
| 5 | 2012-10-28 | Sue | Soda | 3 |
| 6 | 2012-10-28 | Bob | Pizza | 2 |
| 7 | 2012-10-28 | Jill | Pizza | 1 |
| 8 | 2012-10-28 | Jill | Soda | 3 |

# Types of Numerical Analysis

- Type of variables
    - Qualitative
    - Quantitative
- Number of variables
    - Univariate
    - Bivariate
    - Multivariate

# Univariate Analysis

- One variable

- Measures include:
  - Central tendency
  - Dispersion



Source: Wikipedia

# Bivariate Analysis

- Two variables
  - Predictor
  - Outcome
- Measures include
  - Covariance
  - Correlation



Source: Wikipedia

# Code Demo: Movies Data Set

- Movies from 2003
- *Goal:* What movies made the most money

Source: http://www.rossmanchance.com/iscam2/files.html

# Code Demo

# Data Visualization

# Data Visualization

- Visual data representation

- For human pattern recognition

- Map dimensions to visual characteristics

| ID | Date | Customer | Product | Quantity |
|----|------------|----------|---------|----------|
| 1 | 2012-10-27 | John | Pizza | 2 |
| 2 | 2012-10-27 | John | Soda | 2 |
| 3 | 2012-10-27 | Jill | Salad | 1 |
| 4 | 2012-10-27 | Bob | Milk | 1 |
| 5 | 2012-10-28 | Sue | Soda | 3 |
| 6 | 2012-10-28 | Bob | Pizza | 2 |
| 7 | 2012-10-28 | Jill | Pizza | 1 |
| 8 | 2012-10-28 | Jill | Soda | 3 |

Sales by Product

# Types of Data Visualizations

- Type of variable(s)
  - Qualitative
  - Quantitative
- Number of variables
  - Univariate
  - Bivariate
  - Multivariate



**Old Faithful Eruptions**



Source: Wikipedia

# Code Demo: Movies Data Set

- *Goal:* Visualize what types of movies make the most money



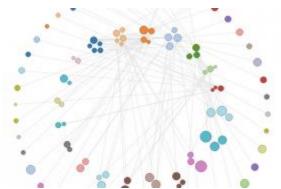Source: http://www.rossmanchance.com/iscam2/files.html

# Code Demo

# Beyond R and EDA

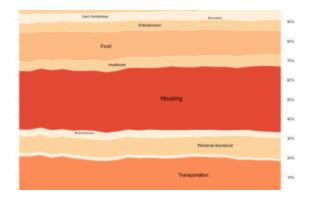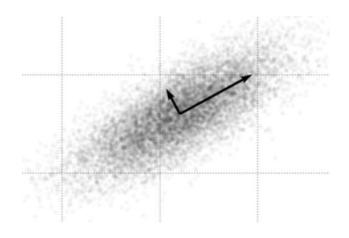This is just the tip of the iceberg!
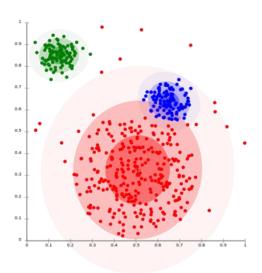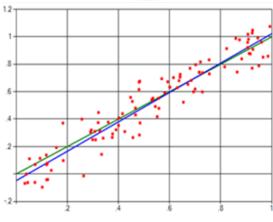
# Advanced Visualizations with R



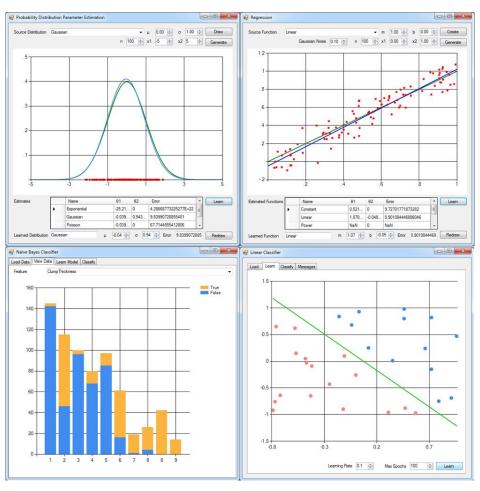Source: Flowing Data

# Advanced Data Analysis with R

- Cluster Analysis

- Statistical Modeling

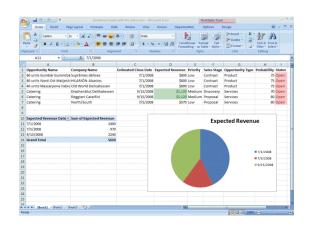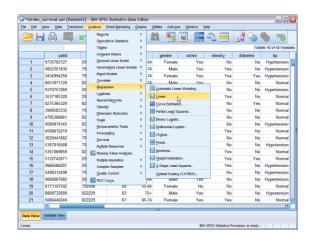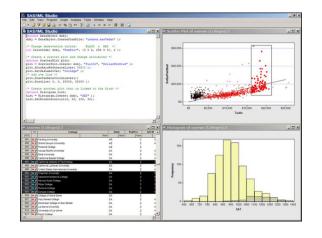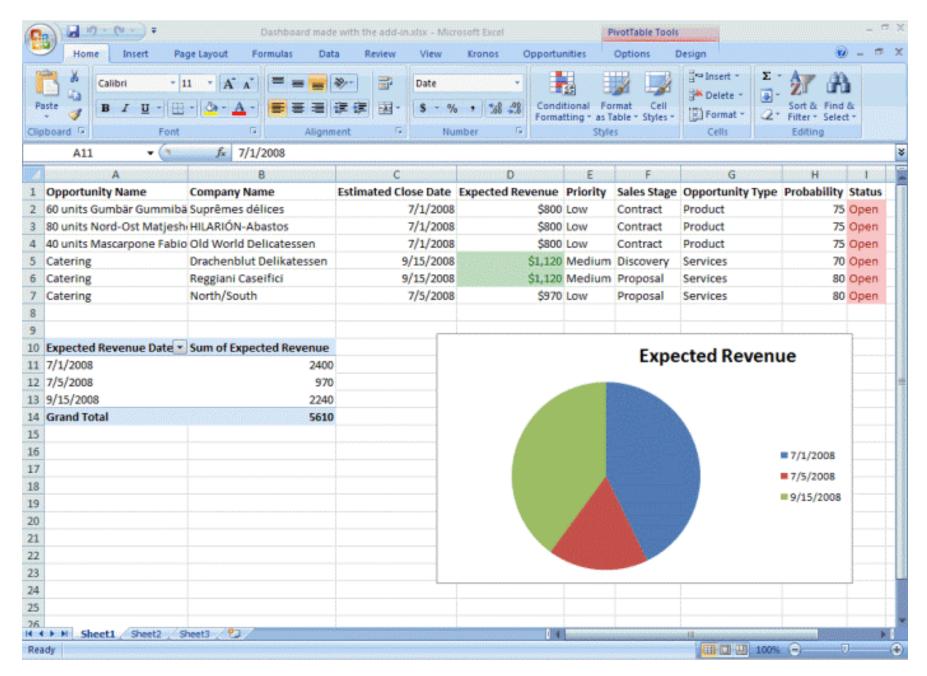- Dimensionality Reduction

- Analysis of Variance (ANOVA)

# Data Mining and Machine Learning with R

# Alternatives to R for EDA

Source: Microsoft

Source: Tableau

Source: IBM SPSS

Source: SAS

```python
1  import sys
2  import json
3  import re
4
5  def getSentiments(sentiment_file):
6      scores = {}
7      for line in sentiment_file:
8          term, score = line.split("\t")
9          scores[term] = int(score)
10     return scores
11
12 def getTweets(tweet_file):
13     tweets = []
14     for line in tweet_file:
15   tweet = json.loads(line)
16   text = tweet.get("text")
17   tweets.append(text)
18     return tweets
19
20 def getAllTerms(tweets, sentiments):
21     allTermScores = {}
22     for tweet in tweets:
23         tweetTerms = getTweetTerms(tweet, sentiments)
24         for tweetTerm in tweetTerms:
25             if sentiments.has_key(tweetTerm):
26                 continue
27             if not allTermScores.has_key(tweetTerm):
28                 allTermScores[tweetTerm] = []
29             termScores = allTermScores[tweetTerm]
30             termScores.append(tweetTerms[tweetTerm])
31     return allTermScores
32
```
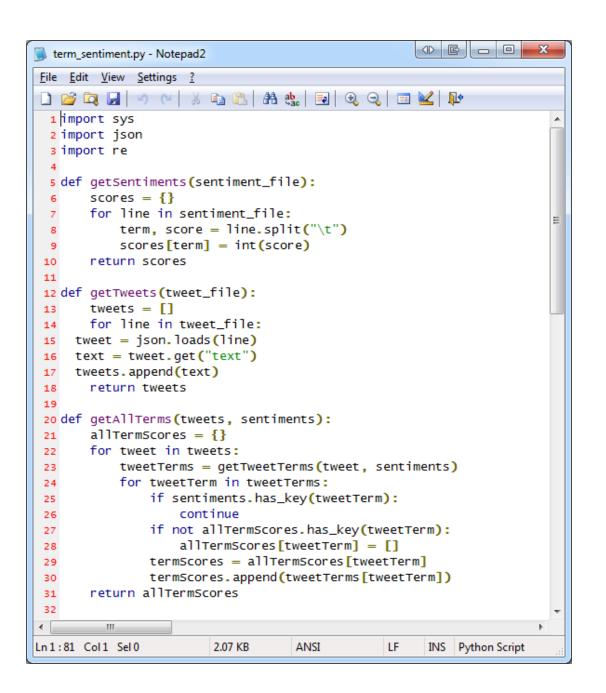
# Code Demo

# Where to Go Next…

- R website: http://www.cran.r-project.org
- R Studio: http://www.rstudio.com
- Pluralsight: http://www.pluralsight.com
- Coursera: https://www.coursera.org
- Revolutions: http://blog.revolutionanalytics.com
- Flowing Data: http://flowingdata.com
- R-Blogger: http://www.r-bloggers.com

# Conclusion

# Conclusion

- Introduction to R
- Data munging
- Descriptive statistics
- Data visualization
- Beyond R & EDA

# Feedback

- Feedback is very important to me
- One thing you liked?
- One thing I could improve?

# Contact Info

Matthew Renze
@matthewrenze
matthew@renzeconsulting.com


Renze Consulting
www.renzeconsulting.com


Data Explorer
http://www.data-explorer.com