

# Exploratory Data Analysis with R

@matthewrenze

#codemash

# 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 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...



Source: <http://www.dot.gov.nt.ca/>



Source: Wikipedia

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
  - B.A. in Philosophy
- Community
  - Pluralsight Author
  - ASPInsider
  - Public Speaker
  - Open-Source Software

IOWA STATE  
UNIVERSITY



# Introduction to R



# What is R?

- Open source
- Language and environment
- Numerical and graphical analysis
- Cross platform



# What is R?

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

and best of all...

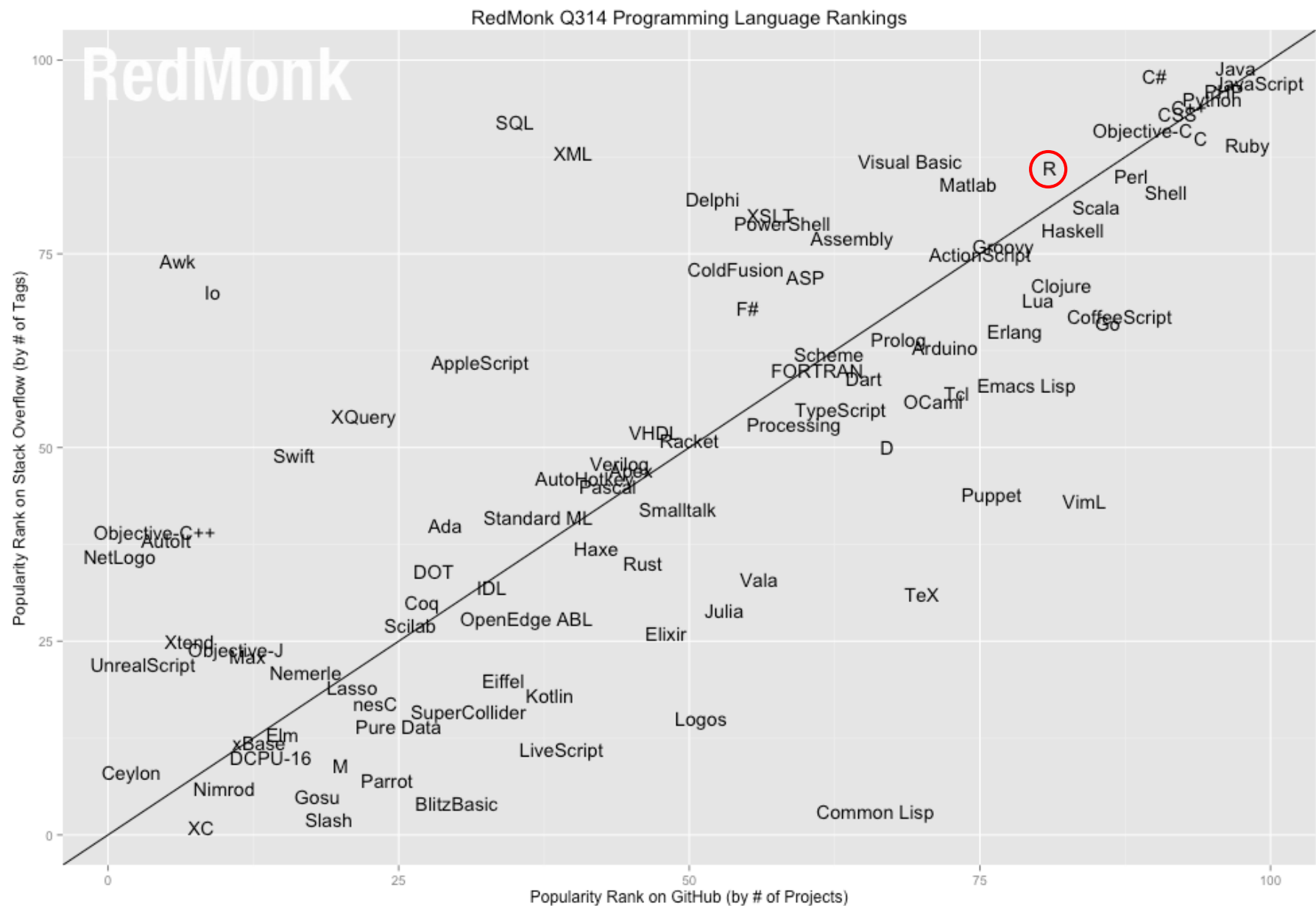




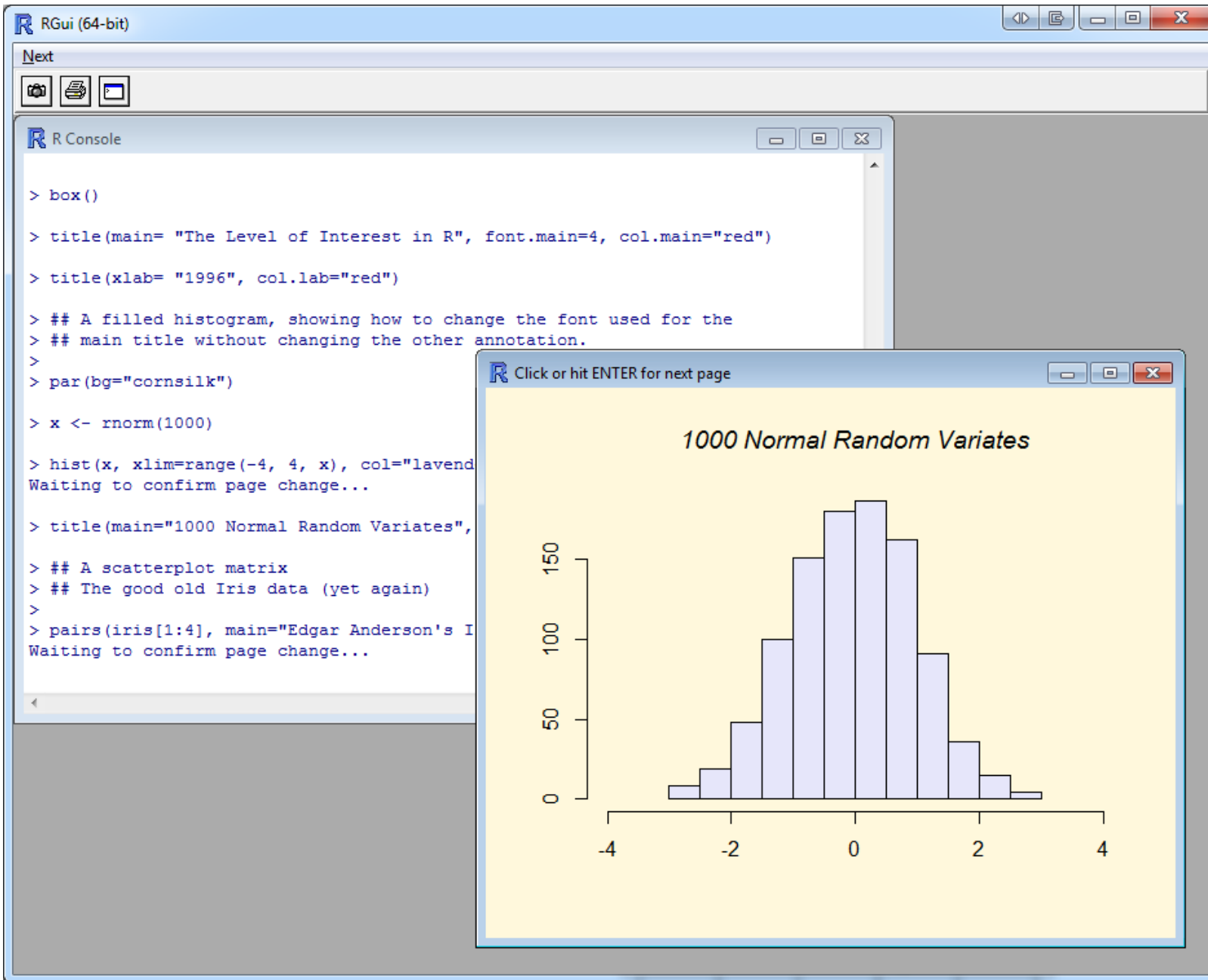
FREE

# FREE

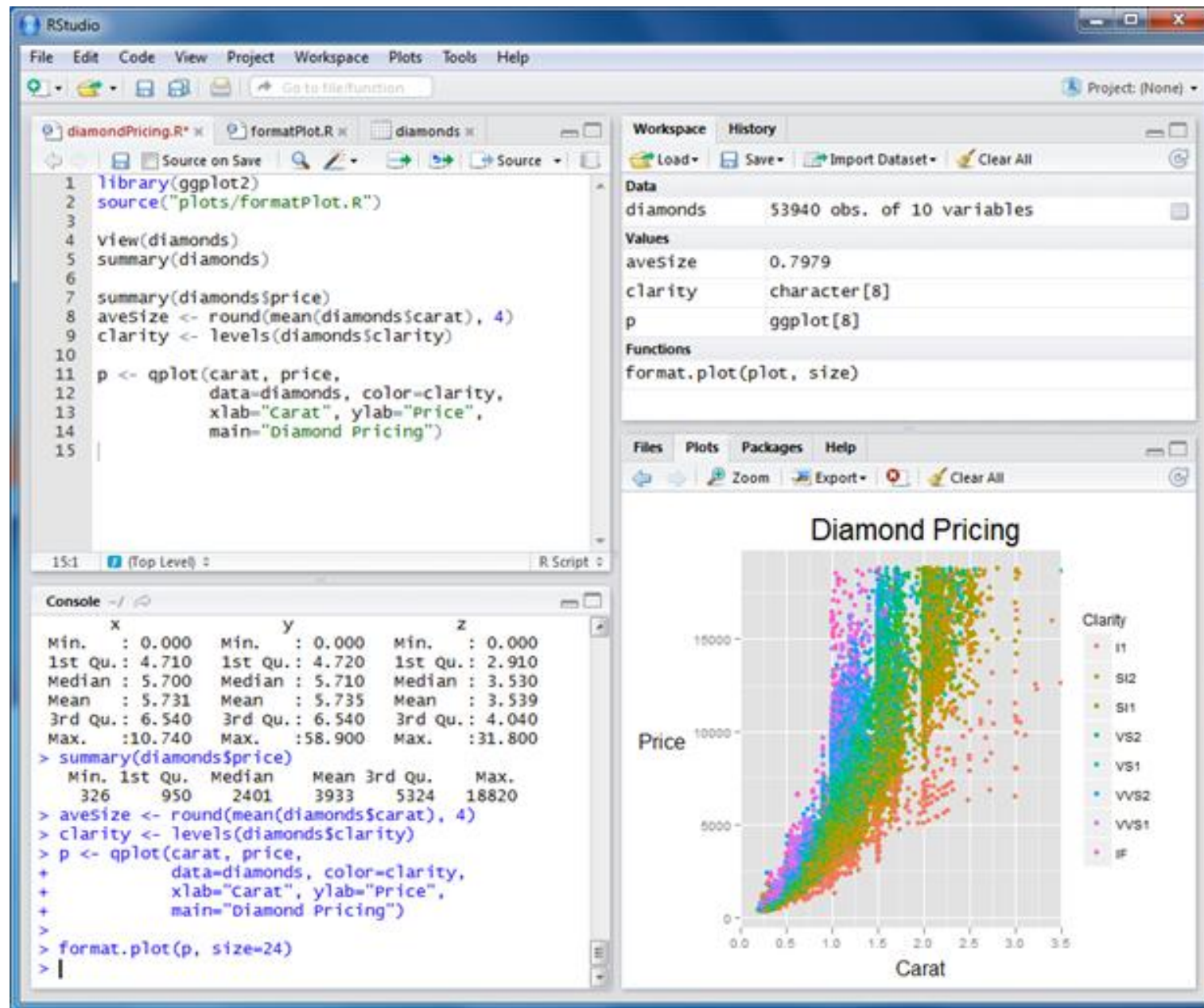




Source: <http://redmonk.com/sograzy/2014/06/13/language-rankings-6-14/>







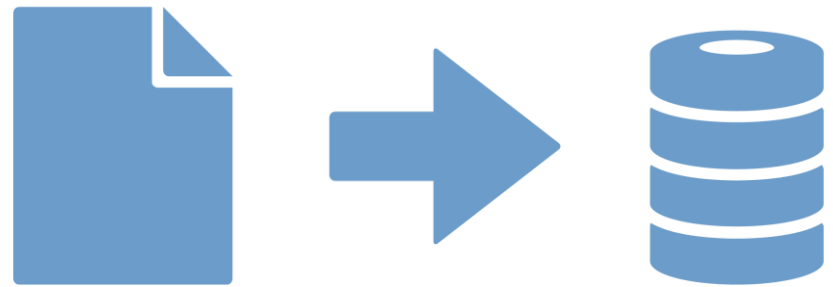
Code Demo



# Data Munging

# Data Munging

- Transforming data
- Raw data to usable data
- Data must be cleaned first



# 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



**PROD. NO.**  
**SCENE**

**TAKE**

**ROLL**









- Column with wrong name
- Rows with missing values
- Runtime column has units
- Revenue in multiple scales
- Wrong file format



Code Demo



# Descriptive Statistics

# Descriptive Statistics

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

Movie Runtime	
Statistic	Value (minutes)
Minimum	38
1 <sup>st</sup> Quartile	93
Median	101
Mean	104
3 <sup>rd</sup> Quartile	113
Maximum	219

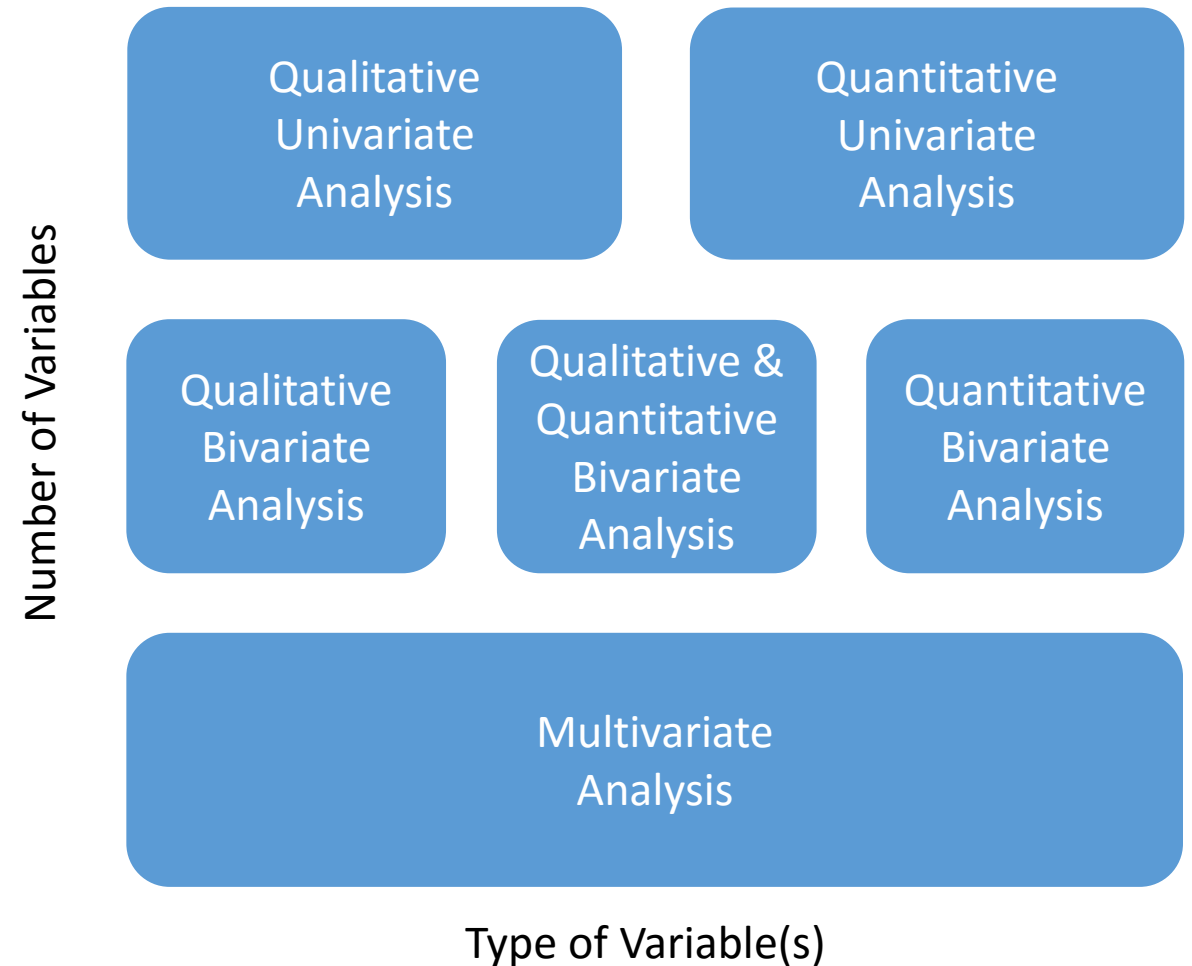
# Statistical Terms

- Observations
- Variables
- Qualitative variable
- Quantitative variable

ID	Date	Customer	Product	Quantity
1	2015-08-27	John	Pizza	2
2	2015-08-27	John	Soda	2
3	2015-08-27	Jill	Salad	1
4	2015-08-27	Jill	Milk	1
5	2015-08-28	Miko	Pizza	3
6	2015-08-28	Miko	Soda	2
7	2015-08-28	Sam	Pizza	1
8	2015-08-28	Sam	Milk	1

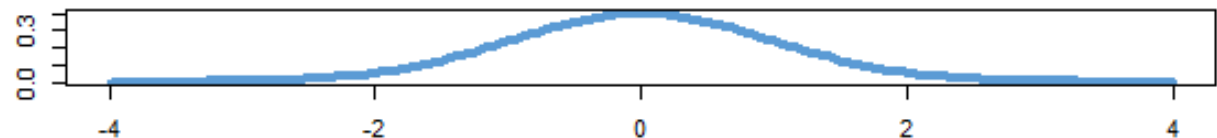
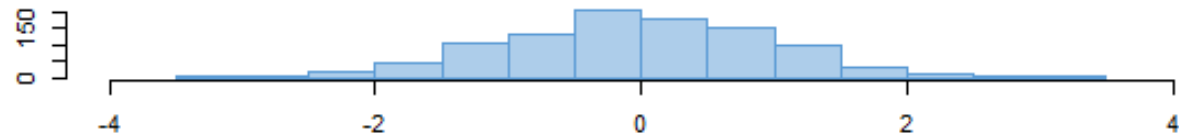
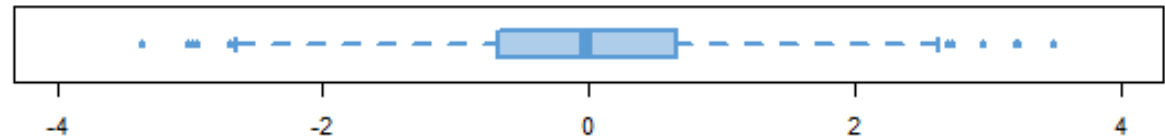
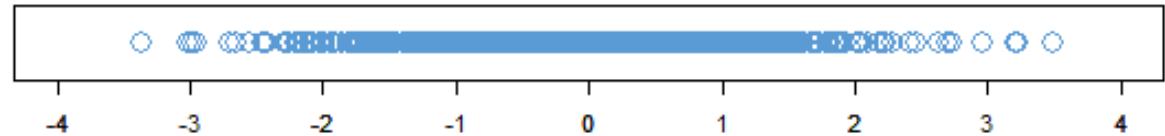
# Types of Analysis

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



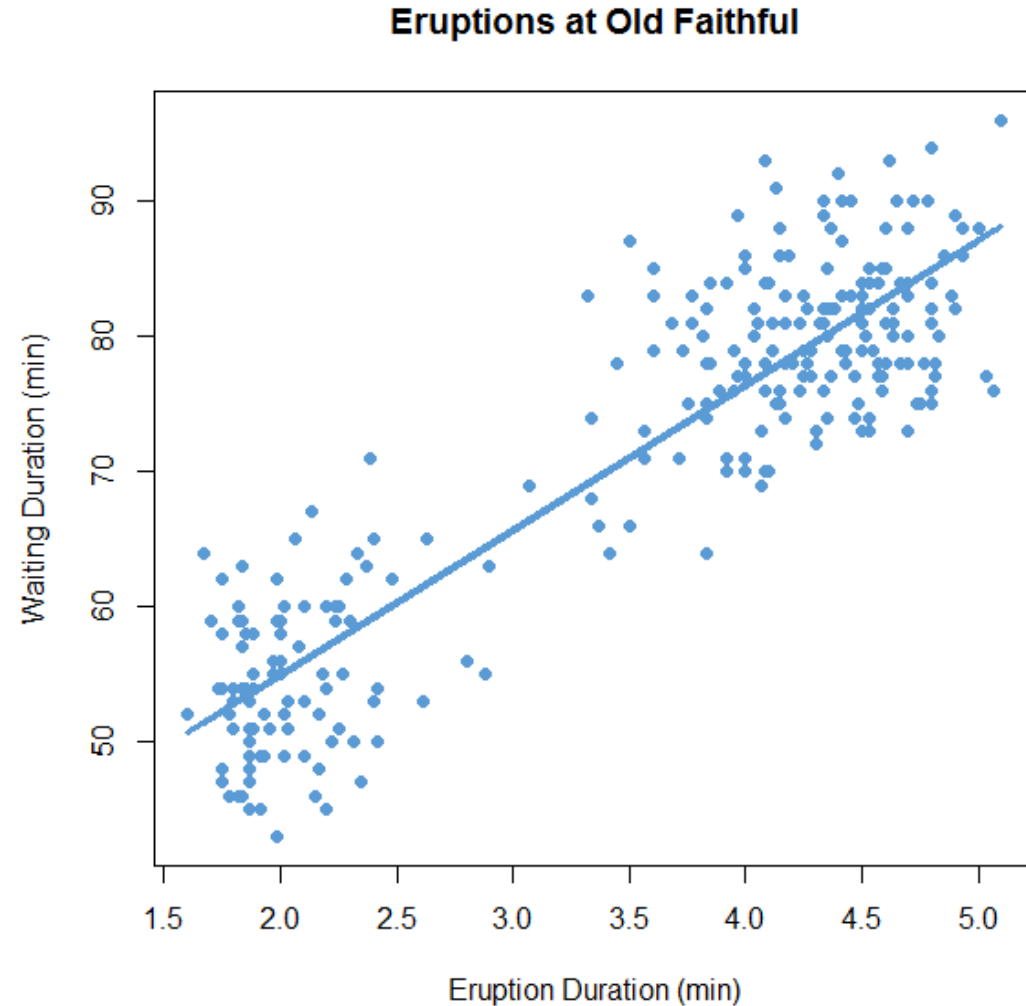
# Univariate Analysis

- One variable
- Qualitative
  - Frequency
- Quantitative
  - Central tendency
  - Dispersion



# Bivariate Analysis

- Qualitative
  - Joint frequency
- Quantitative
  - Two variables
    - Predictor
    - Outcome
  - Measures
    - Covariance
    - Correlation











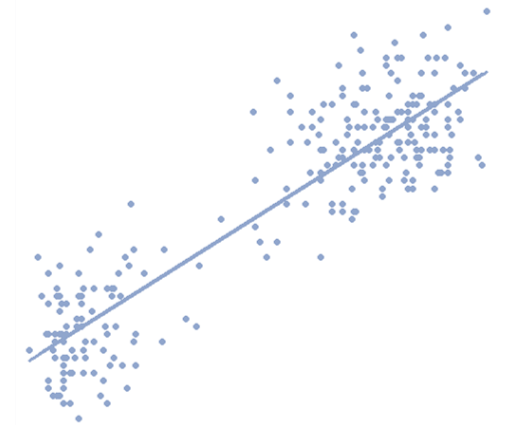
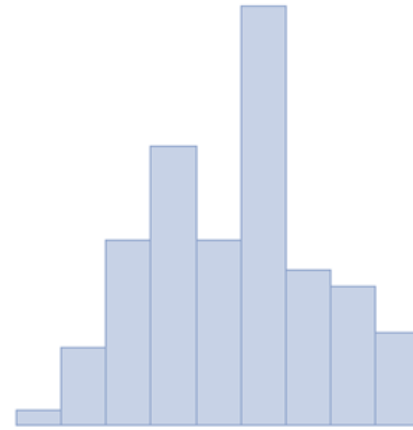
Code Demo



# Data Visualization

# Data Visualization

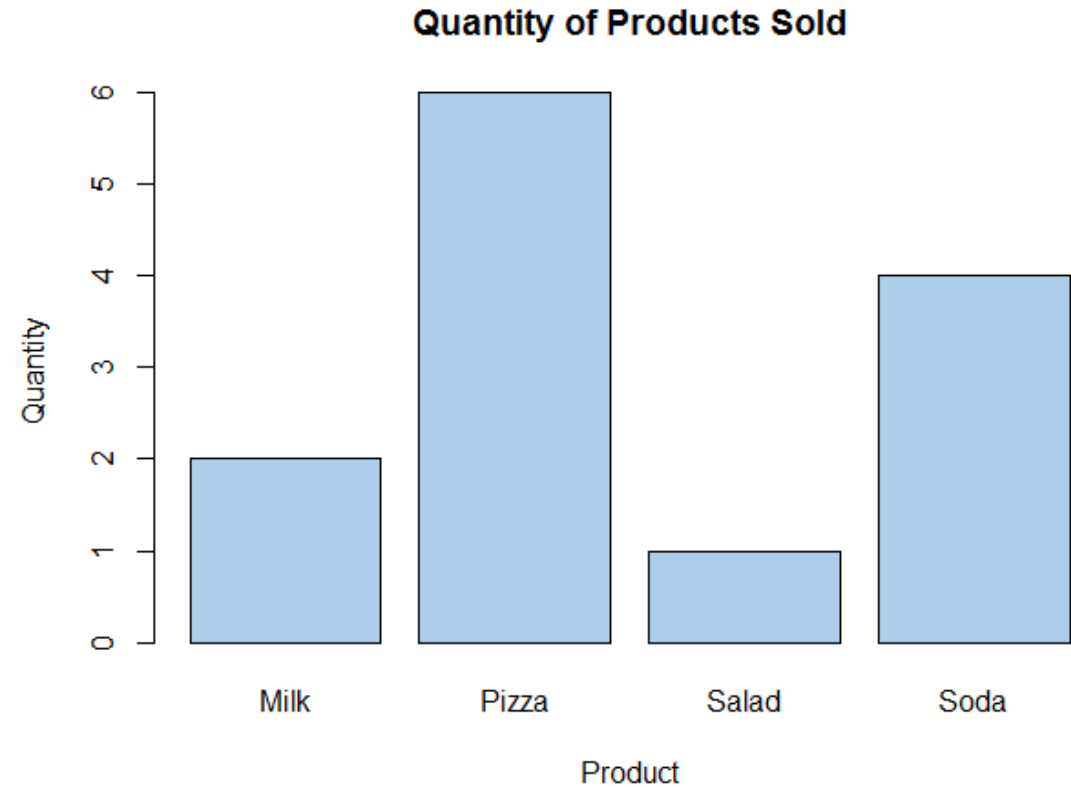
- Visual data representation
- For human pattern recognition
- Map dimensions to visual characteristics





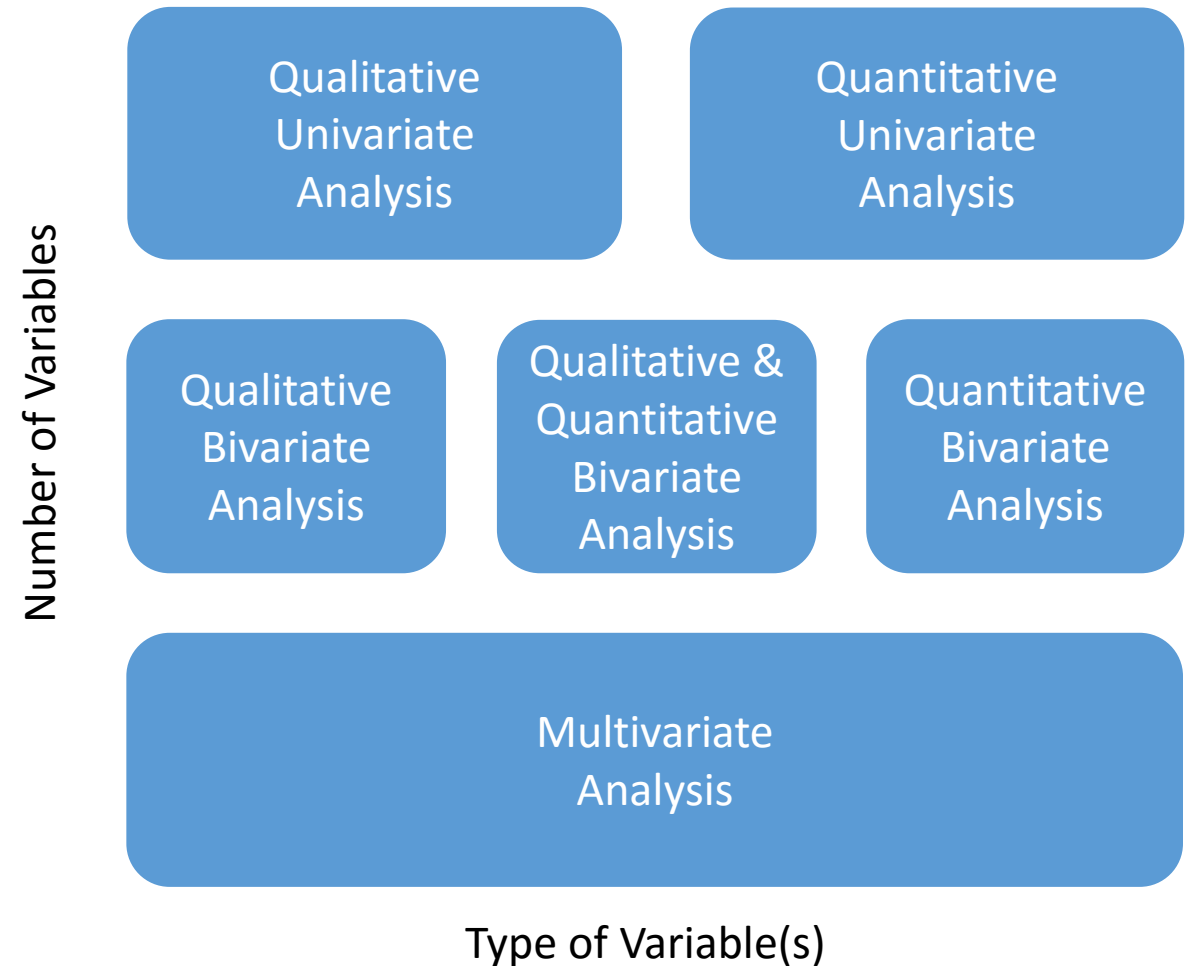
# Data Visualization

ID	Date	Customer	Product	Quantity
1	2015-08-27	John	Pizza	2
2	2015-08-27	John	Soda	2
3	2015-08-27	Jill	Salad	1
4	2015-08-27	Jill	Milk	1
5	2015-08-28	Miko	Pizza	3
6	2015-08-28	Miko	Soda	2
7	2015-08-28	Sam	Pizza	1
8	2015-08-28	Sam	Milk	1



# Types of Data Visualizations

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

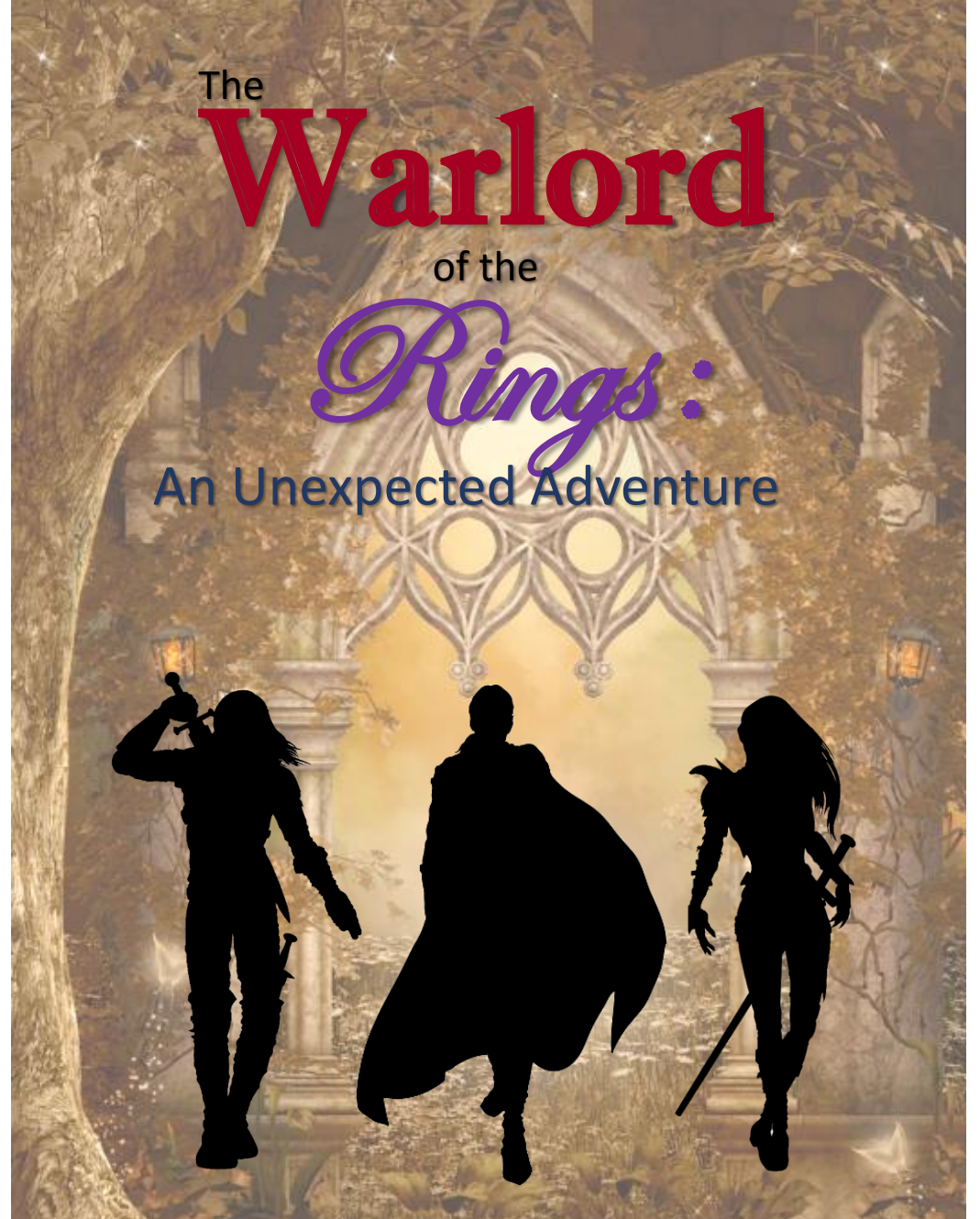






Code Demo





Feature Length

PG

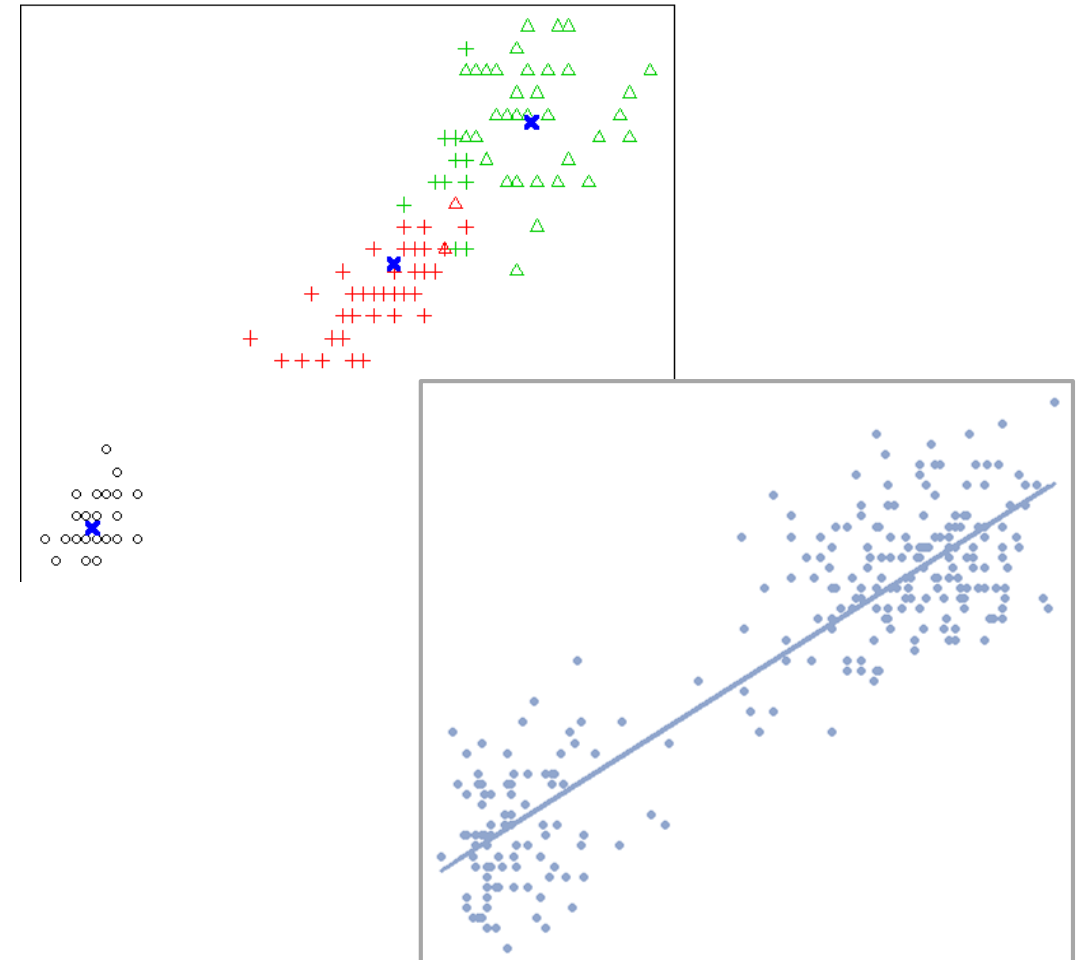
Beyond R and EDA



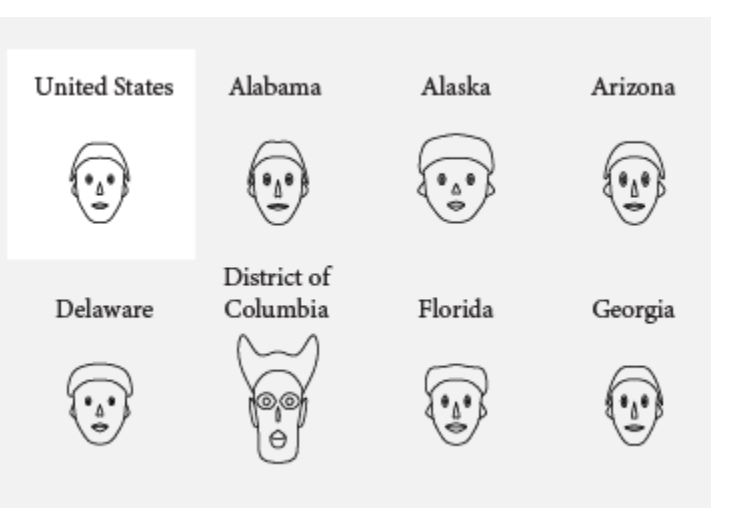
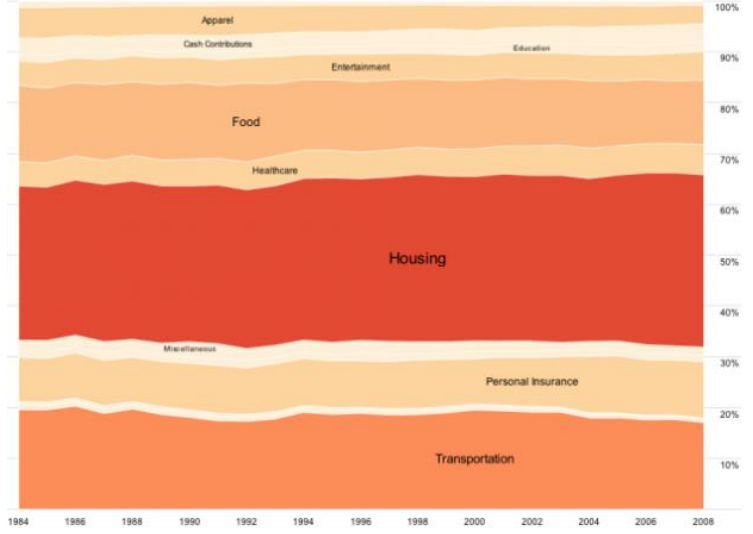
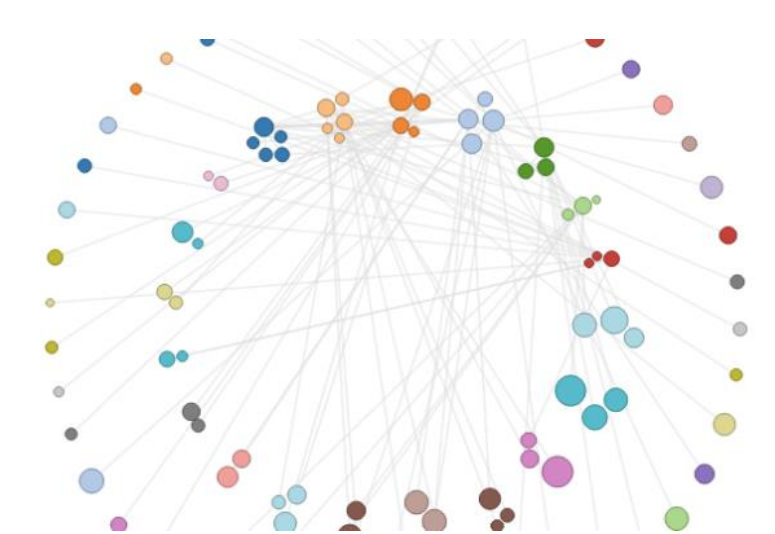
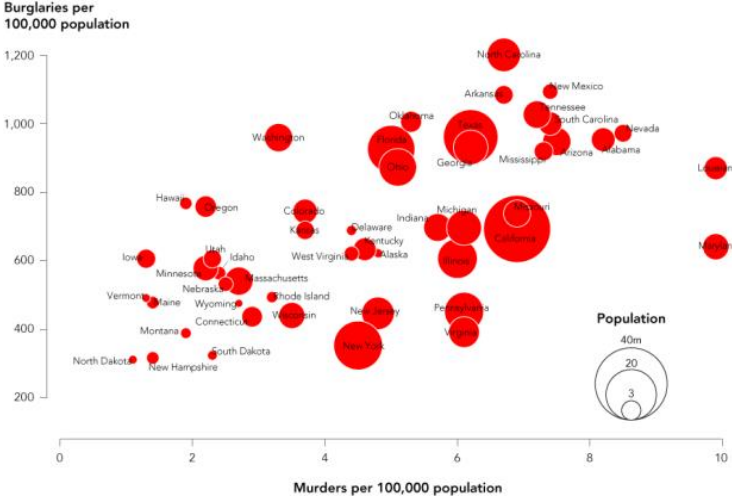
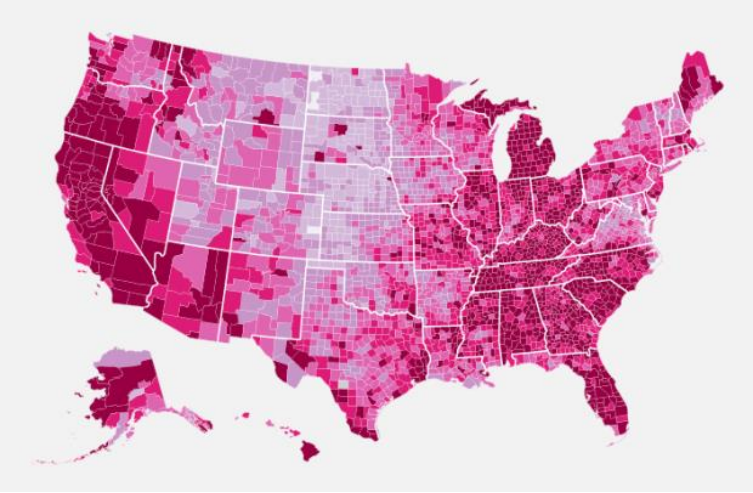
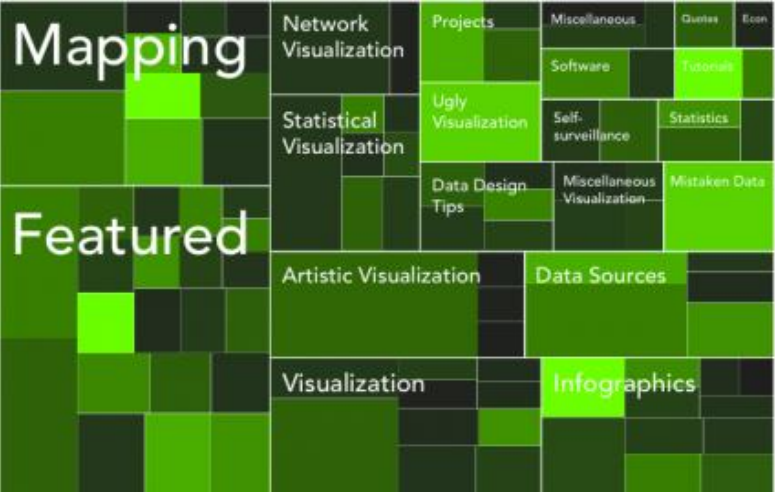
This is just the tip of the iceberg!

# Advanced Data Analysis with R

- Cluster Analysis
- Statistical Modeling
- Dimensionality Reduction
- Analysis of Variance (ANOVA)

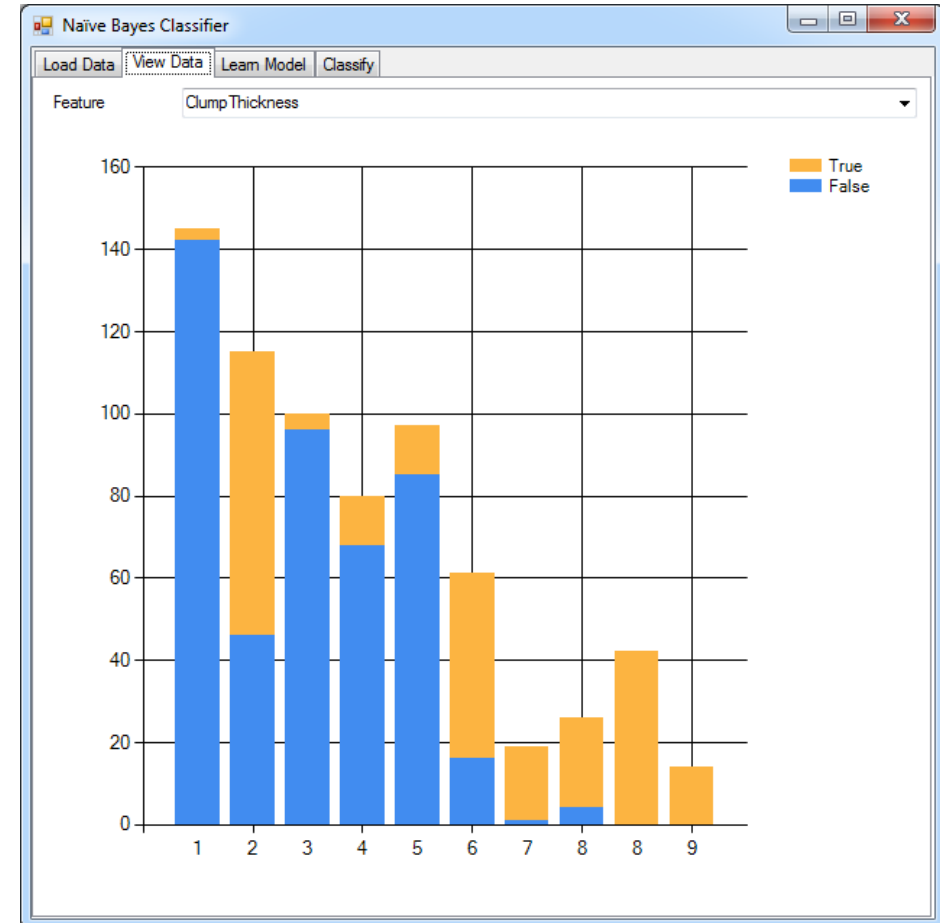
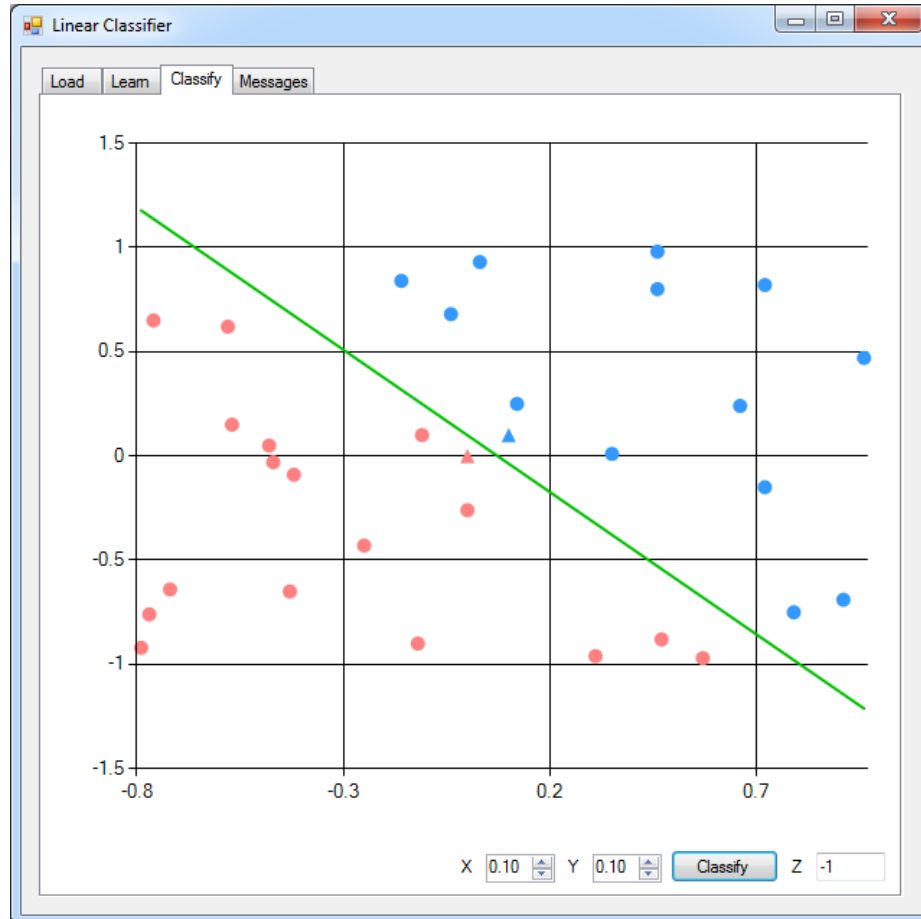






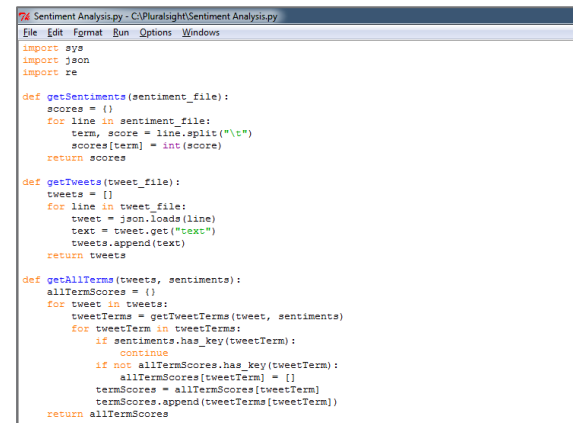
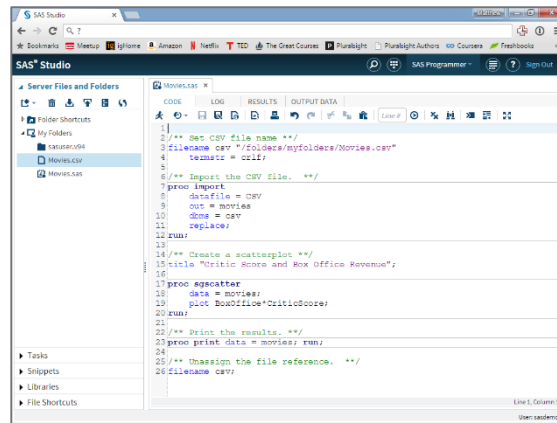
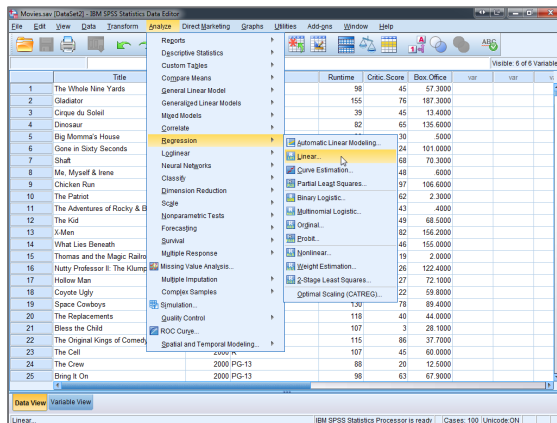
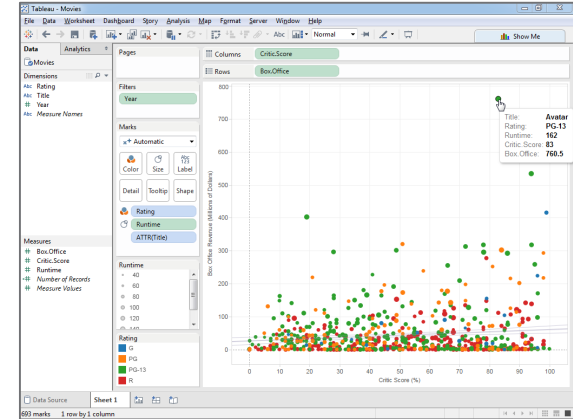
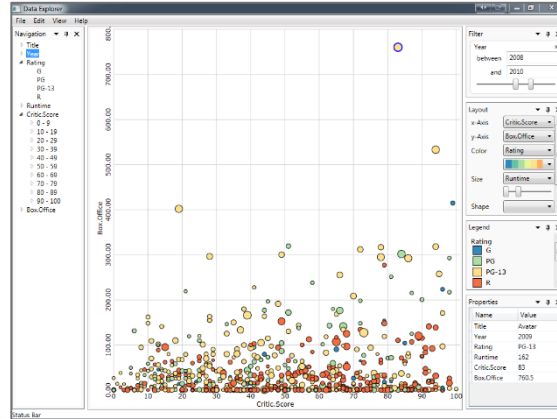
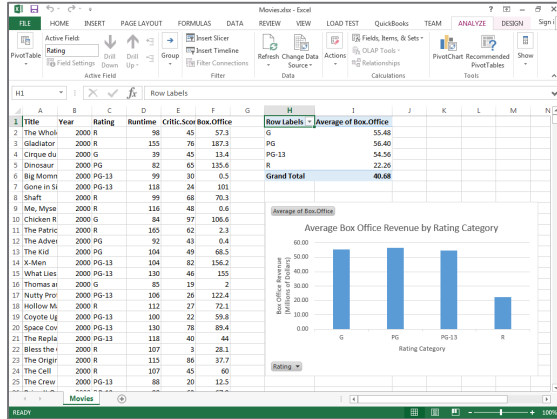
Source: Nathan Yau ([www.flowingdata.com](http://www.flowingdata.com))

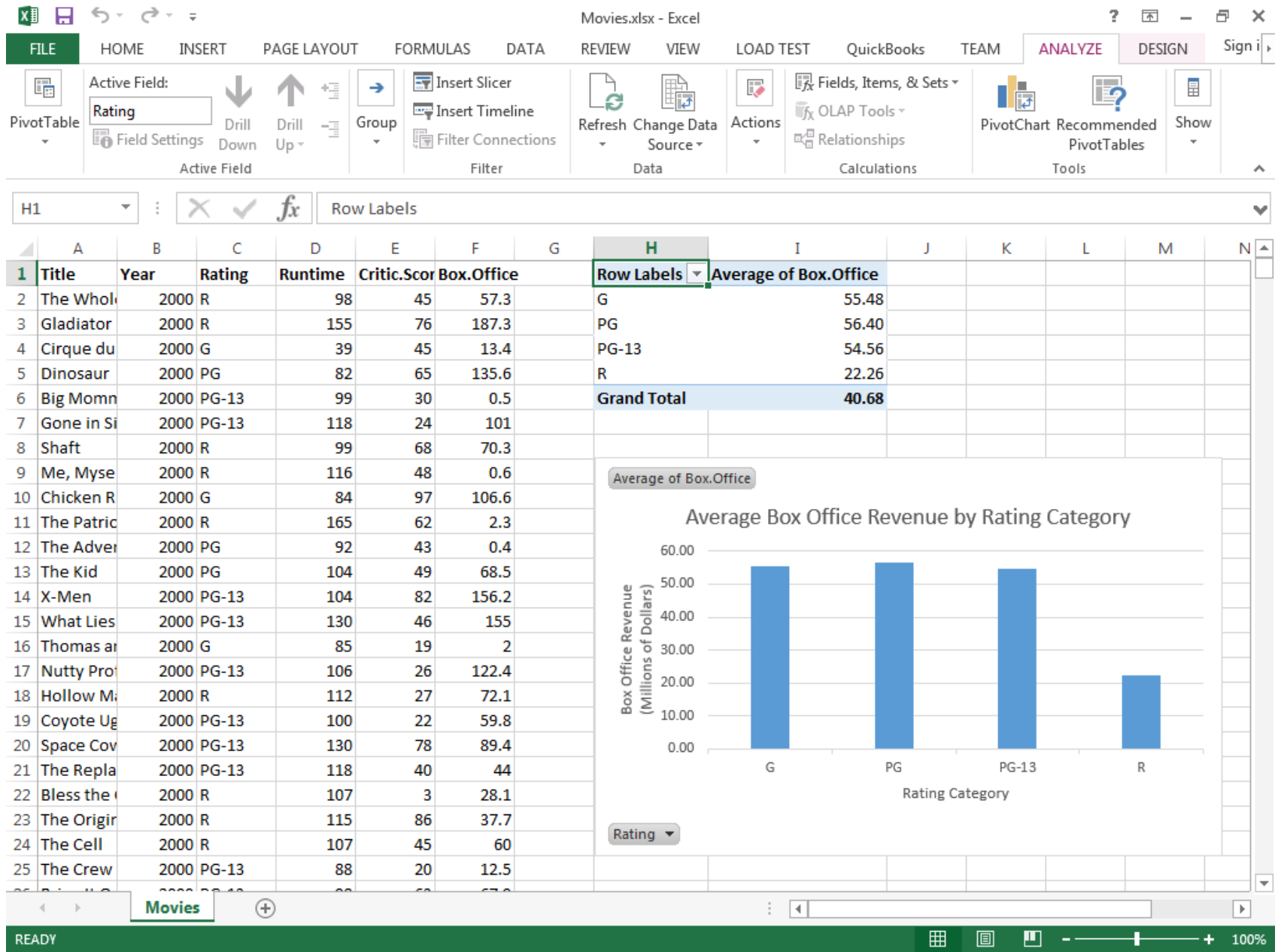
# Data Mining and Machine Learning with R

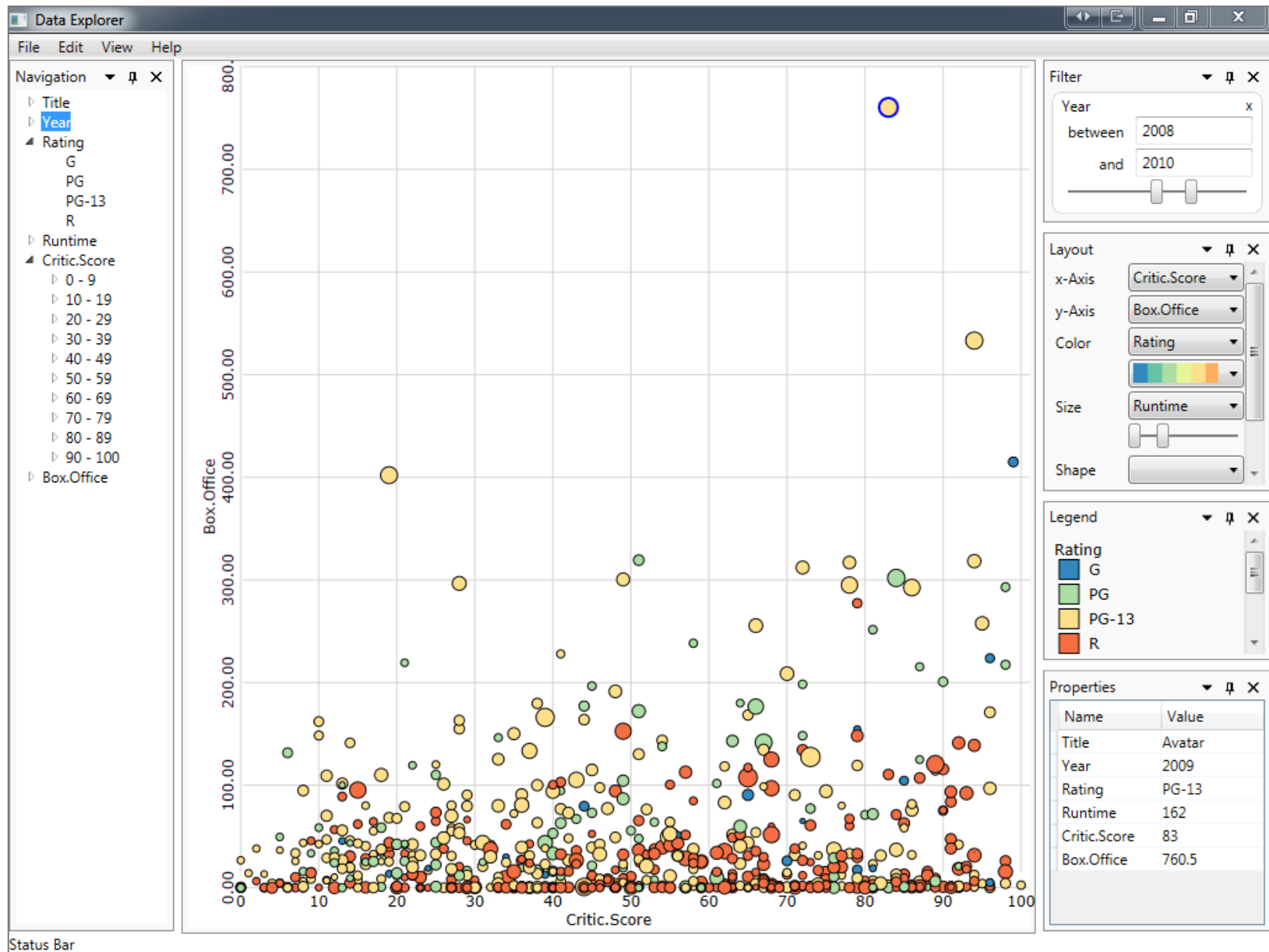


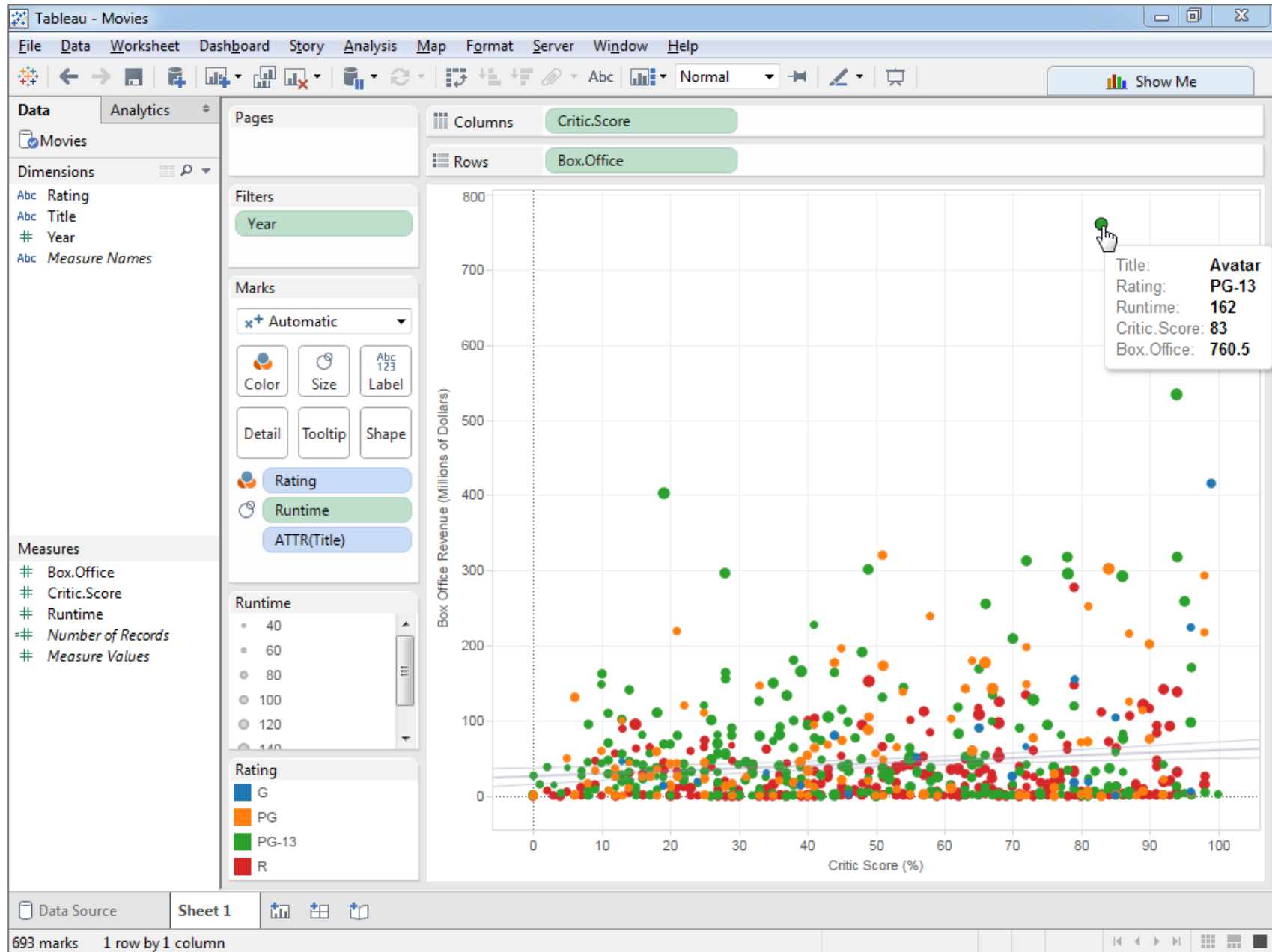


# Alternatives to R for EDA









IBM SPSS Statistics Data Editor window showing the 'Analyze' menu and the 'Linear...' regression dialog box.

**File Edit View Data Transform Analyze Direct Marketing Graphs Utilities Add-ons Window Help**

**Reports**  
**Descriptive Statistics**  
**Custom Tables**  
**Compare Means**  
**General Linear Model**  
**Generalized Linear Models**  
**Mixed Models**  
**Correlate**  
**Regression** (selected)  
**Loglinear**  
**Neural Networks**  
**Classify**  
**Dimension Reduction**  
**Scale**  
**Nonparametric Tests**  
**Forecasting**  
**Survival**  
**Multiple Response**  
**Missing Value Analysis...**  
**Multiple Imputation**  
**Complex Samples**  
**Simulation...**  
**Quality Control**  
**ROC Curve...**  
**Spatial and Temporal Modeling...**

**Automatic Linear Modeling...**  
**Linear...** (selected)  
**Curve Estimation...**  
**Partial Least Squares...**  
**Binary Logistic...**  
**Multinomial Logistic...**  
**Ordinal...**  
**Probit...**  
**Nonlinear...**  
**Weight Estimation...**  
**2-Stage Least Squares...**  
**Optimal Scaling (CATREG)...**

Visible: 6 of 6 Variables

	Title	Runtime	Critic.Score	Box.Office	var	var	var
1	The Whole Nine Yards	98	45	57.3000			
2	Gladiator	155	76	187.3000			
3	Cirque du Soleil	39	45	13.4000			
4	Dinosaur	82	65	135.6000			
5	Big Momma's House	30		.5000			
6	Gone in Sixty Seconds	24		101.0000			
7	Shaft	68		70.3000			
8	Me, Myself & Irene	48		.6000			
9	Chicken Run	97		106.6000			
10	The Patriot	62		2.3000			
11	The Adventures of Rocky & Bullwinkle	43		.4000			
12	The Kid	49		68.5000			
13	X-Men	82		156.2000			
14	What Lies Beneath	46		155.0000			
15	Thomas and the Magic Railroad	19		2.0000			
16	Nutty Professor II: The Klump	26		122.4000			
17	Hollow Man	27		72.1000			
18	Coyote Ugly	22		59.8000			
19	Space Cowboys	130		89.4000			
20	The Replacements	118		44.0000			
21	Bless the Child	107		28.1000			
22	The Original Kings of Comedy	115		37.7000			
23	The Cell	107		60.0000			
24	The Crew	88		12.5000			
25	Bring It On	98		67.9000			

**Data View** **Variable View**

Linear... IBM SPSS Statistics Processor is ready Cases: 100 Unicode:ON



File Edit Format Run Options Windows

```
import sys
import json
import re

def getSentiments(sentiment_file):
    scores = {}
    for line in sentiment_file:
        term, score = line.split("\t")
        scores[term] = int(score)
    return scores

def getTweets(tweet_file):
    tweets = []
    for line in tweet_file:
        tweet = json.loads(line)
        text = tweet.get("text")
        tweets.append(text)
    return tweets

def getAllTerms(tweets, sentiments):
    allTermScores = {}
    for tweet in tweets:
        tweetTerms = getTweetTerms(tweet, sentiments)
        for tweetTerm in tweetTerms:
            if sentiments.has_key(tweetTerm):
                continue
            if not allTermScores.has_key(tweetTerm):
                allTermScores[tweetTerm] = []
            termScores = allTermScores[tweetTerm]
            termScores.append(tweetTerms[tweetTerm])
    return allTermScores
```



# Where to Go Next...

- R website: <http://www.cran.r-project.org>
- R Studio: <http://www.rstudio.com>
- Revolutions: <http://blog.revolutionanalytics.com>
- Flowing Data: <http://flowingdata.com>
- R-Blogger: <http://www.r-bloggers.com>

# Exploratory Data Analysis with R



**Matthew Renze**

@matthewrenze | [www.matthewrenze.com](http://www.matthewrenze.com)



[www.pluralsight.com/courses/r-data-analysis](http://www.pluralsight.com/courses/r-data-analysis)

# 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

- Twitter: [@matthewrenze](https://twitter.com/matthewrenze)
- Email: [matthew@renzeconsulting.com](mailto:matthew@renzeconsulting.com)
- Website: [www.matthewrenze.com](http://www.matthewrenze.com)