# INTRO to DATA SCIENCE
## Lecture 1: Welcome

Mark Holt

Biomedical Data Scientist

# I. COURSE OVERVIEW

# II. WHAT IS DATA SCIENCE?

# III. LAB

-SETUP DEVELOPER ENVIRONMENT

-DATA EXPLORATION

# I. COURSE OVERVIEW

### CONTACT INFO

#### OFFICE HOURS

Mark Holt           mrgh@me.com             by appointment

Susan Sun           sun.w.susan@gmail.com

**Class M/W 6:30-9:30 PM** - 6/3 - 8/19
- ○ GA West (10 E. 21st St, 4th Floor), Room

## TOPICS

- **GRADIENT DESCENT & LINEAR REGRESSION**
- **LOGISTIC REGRESSION**
- **BASIC VISUALIZATION**
- **REGULARIZATION**
- **UNSUPERVISED LEARNING - CLUSTERING**
- **DIMENSIONALITY REDUCTION**
- **BAYESIAN INFERENCE & AB TESTING**
- **NATURAL LANGUAGE PROCESSING**
- **RECOMMENDATION SYSTEMS**
- **SUPPORT VECTOR MACHINES**
- **DATABASES**
- **PLUS GUEST SPEAKERS**

## SKILLS

- ▸ **PROGRAMMING**
- ▸ **PRESENTATION & COMMUNICATION**

## ASSIGNMENTS

▸ **DATAEXPLORER CHALLENGES**

▸ **TERM PROJECT**

▸ **HOW MUCH HOMEWORK????**

## TOOLS

PYTHON DATA SCIENCE STACK

- SCIKIT-LEARN (MACHINE LEARNING)
- NUMPY, SCIPY (LINEAR ALGEBRA, NUMERICAL COMPUTATION)
- MATPLOTLIB (VISUALIZATION)
- PANDAS (MODELING, EASY-TO-USE DATA STRUCTURES)
- IPYTHON (INTERACTIVE INTERFACE)

# COURSE ETIQUETTE & INFO

- There is absolutely no such thing as a stupid question
- No-one knows everything about data science so please respect your fellow students
- feel free to eat/drink during class, but try to do it without disturbing the class
- Please remember that you all come from different backgrounds, what is simple for some may be challenging for others. If you are ahead of the class, or finding it easy going, say during a lab, don't disturb the rest of the class with loud side conversations
- There will be 2 breaks during the class of about 10 mins each

## INTRODUCTIONS

Partner with someone next to you. Introduce yourselves. When we reconvene, you will introduce your partner to the group.

Please share your partner's:
▸ name
▸ occupation
▸ experience with Python and scikit-learn
▸ what s/he is most excited about learning/doing
▸ what s/he is most apprehensive about learning/doing
▸ something unusual or quirky about them to help ME remember your names

# II. WHAT IS DATA SCIENCE?

# FUN FACT:

- ‣ Every Day We Create 2.5 Quintillion Bytes of Data
- ‣ 90% of current data was collected in the past two years

Figure 1
**Data is growing at a 40 percent compound annual rate, reaching nearly 45 ZB by 2020**

Data in zettabytes (ZB)



Source: Oracle, 2012

- A set of tools and techniques used to extract useful information (aka knowledge) from data

- An interdisciplinary, problem-oriented subject

- The application of the *Scientific Method* to solving business problems
- An art

**WHO USES DATA SCIENCE?**

# WHO USES DATA SCIENCE?

- *"a data analyst who lives in California"*

- *"a business analyst who lives in New York"*

- *"a statistician who lives in San Francisco"*

Michael E. Driscoll
@medriscoll

Following

Data scientists: better statisticians than most programmers & better programmers than most statisticians bit.ly/NHmRqu @peteskomoroch

← Reply ⇄ Retweet ★ Favorite ••• More ▽ Pocket

- Software Engineering

- Machine Learning

- Domain expertise

- Communication

- Visualization

- Stats

- Math

- Intuition

# 9. Data Munging

Denoising · Feature Extraction · Binning Sparse Values · Unbiased Estimators · Handling Missing Values · Data Scrubbing · Normalization · Dimensionality & Numerosity Reduction

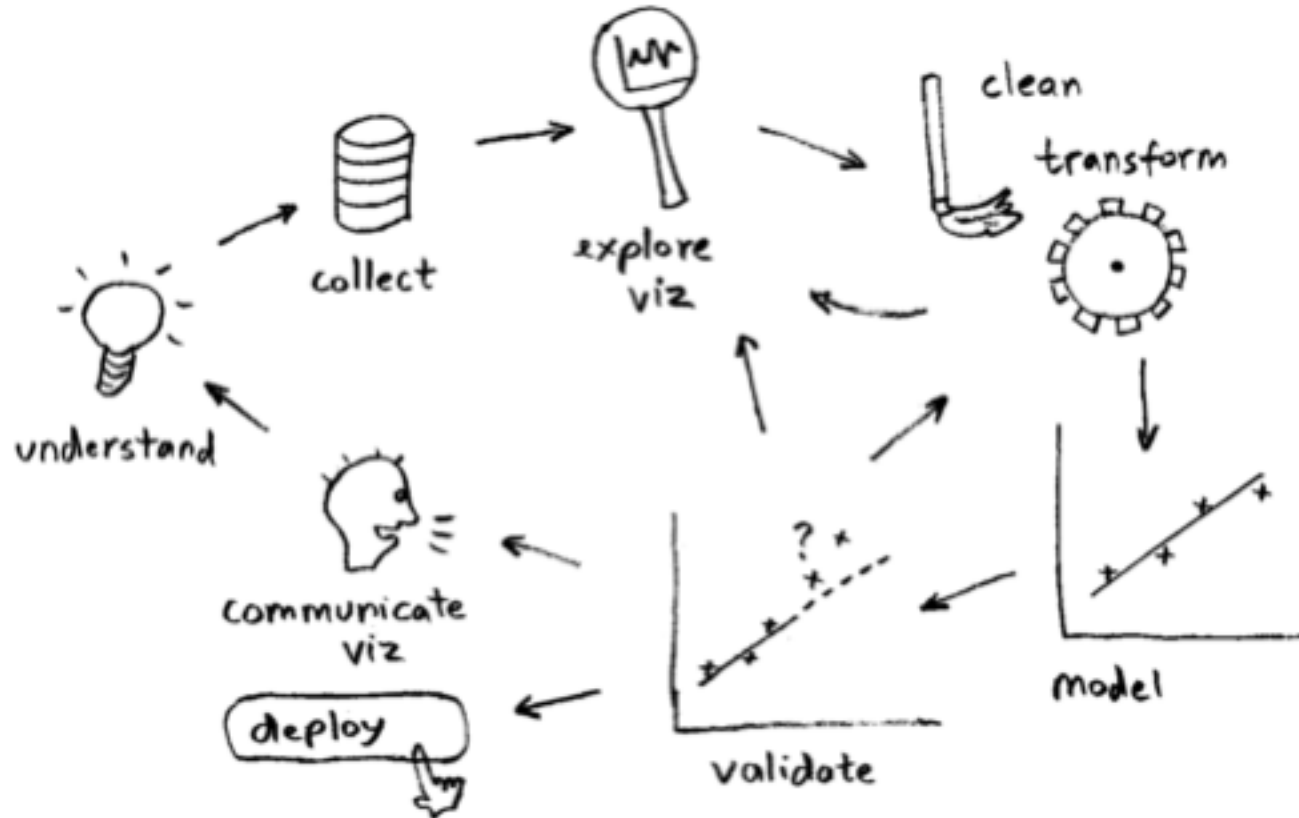Sampling · Stratified Sampling · Principal Component Analysis

Using ETL · How much Data? · Google OpenRefine · Data Survey

# 8. Data Ingestion

Denormalization & Enrichment · Data Frame · Data Integration · Data Sources & Acquisition · Data Discovery · Summary of Data Formats

# 5. Text Mining / NLP

Named Entity Recognition · Text Analysis · Term Document Matrix · Term Frequency & Weight · Support Vector Machines · Association Rules · Market Based Analysis · Feature Extraction · Using Mahout · Using Weka · Using NLTK

Clustering

Hierarchical Clustering · K-Means Clustering · Neural Network · Sentiment Analysis · Collaborative Filtering · Tagging · Vocabulary Mapping · Gensis · Topics?

Regression
Perceptron · Linear Regression · Ranking · Logistic Regression

K-Nearest Neighbors · Naive Bayes Classifiers · Reasoning · Decision Trees · Classification Rate · Trees & Classification · Bias & Variance · Overfitting · Precision · Consistency · Modeling & Test Data · Training & Test Data · ROC Curve · Cross Validation · Supervised Learning · Unsupervised Learning · Semi-Supervised · Computational Fire · Numerical Var · True Pos · is it? · L curve Fit · Classification · Program Cluster · Correlation

Classification

# 6. Visualization

Data Exploration in R (Hist, Boxplot etc) · Uni, Bi & Multivariate Viz · ggplot2 · Histogram & Pie (Uni) · Tree & Tree Map · Scatter Plot (Bi) · Line Charts (Bi) · Spatial Charts · Survey Plot · Timeline · Decision Tree

# 4. Machine Learning

# 1. Fundamentals

Matrices & Linear Algebra Fundamentals · Hash Functions, Binary Tree, O(n) · Relational Algebra, DB Basics · Inner, Outer, Cross, Theta Join · CAP Theorem · Tabular Data · Data Frames & Series · Sharding · OLAP · Multidimensional Data Model · ETL · Reporting Vs BI Vs Analytics

JSON & XML · NoSQL · Regex · Vendor Landscape · Env Setup

# 10. Toolbox

MS Excel w/ Analysis ToolPak · Java, Python · R, R-Studio, Rattle · Weka, Knime, RapidMiner · Hadoop Dist of Choice · Spark, Storm · Flume, Scibe, Chukwa · Nutch, Talend, Scraperiki · Webscraper, Flume, Sqoop · tm, RWeka, NLTK · RHIPE · D3.js, ggplot2, Shiny · IBM Languageware · Cassandra, MongoDB

# 2. Statistics

Prob Den Fn (PDF) · ANOVA · Skewness · Continues Distributions (Normal, Poisson, Gaussian) · Cumul Dist Fn (CDF) · Random Variables · Bayes Theorem · Probability Theory · Percentiles & Outliers · Histograms · Exploratory Data Analysis · Descriptive Statistics (mean, median, range, SD, Var) · Pick a Dataset (UCI Repo)

Central Limit Theorem · Monte Carlo Method · Hypothesis Testing · P-Value · Chi² Test · Estimation · Confid Int (CI) · MLE · Kernel Density Estimate · Regression · Covariance · Correlation · Pearson Coeff

# 3. Programming

Data Frames · Reading CSV Data · Reading Raw Data · Subsetting Data · Manipulate Data Frames · Functions · Lists · Factors · Arrays · Matrices · Vectors · Variables · Expressions · R Basics · R Setup, R Studio

Rapid Miner · IBM SPSS · Factor Analysis · Install, Pkgs

# 7. Big Data

Name & Data Nodes · Setup Hadoop (IBM / Cloudera / HortonWorks) · Data Replication Principles · HDFS · Hadoop Components · Map-Reduce Fundamentals

Zookeeper Avro · Storm: Hadoop Realtime · Rhadoop, RHIPE · Cassandra · MongoDB, Neo4j

**Jeff Hammerbacher, Cloudera, Mount Sinai:**

1. Identify problem
2. Instrument data sources
3. Collect data
4. Prepare data (integrate, transform, clean, impute, filter, aggregate)
5. Build model
6. Evaluate model
7. Communicate results

# DATA SCIENCE WORKFLOW



source: DATA SCIENCE TOOLBOX SURVEY RESULTS… SURPRISE! R AND PYTHON WIN

**Problem: You are a top Internet news site competing hard in a tough environment. You have a good story but need a good headline to get readers to enter  your site to read the story. Your advertising dollars depend on your headline attracting the most readers in the shortest possible time. How would you approach solving this problem??**

In a small group, consider how you might go about solving this problem

# III. LAB: SETUP DEVENV