*Week 1*

**Data Science Toolkit – Linux, Git, Bash, and SQL**

**Data Science with R – Data Analytics – Part I**

- Linux system
  - Introduce Linux environment
  - Learn Linux commands
  - IO redirection and Pipe
  - Introduce server-side Linux usage
- Git
  - Introduce modern source code management
  - Learn common git operations
  - Setup github and personal portfolio page
- Other server related topics
  - Text editors and IDEs
  - ssh: how to communicate with a remote server
  - Linux environment variables
- SQL
  - Introduction to relational database
  - Introduction to structured query language
  - SQL major commands and examples
- Programming foundation in R I
  - Syntax
  - Data object: Vectors, Matrices, Data Frames, and Lists
  - Common functions
  - Rstudio environment and package management
  - Local data input/output
  - Introduction to R data visualization
- Programming foundation in R II
  - Data sorting and merging
  - String manipulation
  - Dates and times
  - Connecting to an external database

*Week 2*

**Data Science with R – Data Analytics – Part II**

- Data manipulation with "dplyr"
  - Tables in R
  - Join
  - Subset
  - Advanced manipulations with dplyr
- Data Visualization with "ggplot2"

- o   Histogram
- o   Point graphics
- o   Columnar graphics
- o   Line charts
- o   Pie charts
- o   Box plots
- o   Scatter plots
- o   Visualizing multivariate data
- o   Matrix-based visualizations
- o   Maps
- Introduction to Shiny
  - o   Shiny introduction
  - o   Design the User-interface
  - o   Control widgets
  - o   Build reactive output
  - o   Use data table in Shiny Apps
  - o   Use R scripts, data and packages
  - o   UI and server for the App
  - o   Make Shiny perform quickly
  - o   Matrix-based visualizations
  - o   Use reactive expressions
  - o   Share and deploy Shiny apps
- Lab: Moneyball

**Project 1 Due:  Exploratory Data Visualization**


*Week 3*
**Data Science with Python - Data Analytics – Part I**
- Python Programming Language I
  - o   Simple Values and Expressions
  - o   Functions
  - o   Lists
  - o   Conditionals
  - o   Functional programming: map, filter and reduce
- Python Programming Language II
  - o   String operations
  - o   File input/output and searching
  - o   Data Structures:
    - ▪   Mutating operations on Lists
    - ▪   Tuples, sets and dictionaries
- Python Programming Language III
  - o   Control flows

- o   Errors and exceptions
- o   Object-oriented programming
- Web scraping
  - o   Regular expression
  - o   HTML, beautiful soup and scrapy
  - o   NoSQL and MongoDB

*Week 4*

**Data Science with Python – Data Analytics – Part II**

- Numpy and Scipy
  - o   Basic data structure and operations
  - o   Matrices and linear algebra
  - o   Stats module
  - o   Random Sampling
- Pandas
  - o   Series and data frame
  - o   I/O of pandas data frame
  - o   Concatenation and merge
  - o   Arithmetic, drop, apply and describe
  - o   Selection and filter
  - o   Missing values
  - o   Grouping and aggregation
  - o   Time series
  - o   Interacting with data base
- Matplotlib and Seaborn
  - o   Basic plots
  - o   Statistical plots:
    - ▪   Scatter plots
    - ▪   Histogram
    - ▪   Boxplot
    - ▪   Barchart
  - o   Multiple figures
  - o   Advanced plots with seaborn
- Python lab: linear regression from scratch

**Project 2 Due: R Shiny Interactive Applications**

*Week 5*

**Data Science with R - Machine Learning – Part I**

- Foundations of Statistics
  - o   Descriptive Statistics

- Measures of Centrality
- Measures of Variability
- Frequency, Proportion & Contingency Tables
- Correlation
  - Hypothesis Testing
    - One Sample t-test
    - Two Sample t-test
    - F-test
    - One-way ANOVA
    - X2 Test of Independence
  - Introduction to Machine Learning
    - Supervised Learning
      - Regression
      - Classification
    - Unsupervised Learning
      - Clustering
      - Dimension Reduction
- Missingness & Imputation
  - Types of Missingness
    - MCAR
    - MAR
    - MNAR
  - Basic Methods of Imputation
    - Mean Value Imputation
    - Simple Random Imputation
    - Regression Prediction
  - K-Nearest Neighbors
    - Voronoi Tessellations
    - KNN for Classification
    - KNN for Regression
    - Distance Measures
- Linear Regression I
  - Simple Linear Regression
    - From a Mathematical Standpoint
    - Accuracy of the Coefficient Estimates
    - Performing Hypothesis Tests
    - Constructing Confidence Intervals
  - Assumptions & Diagnostics
  - Transformations
    - Power Transformation
    - Box-Cox Transformation

- o The Coefficient of Determination $R^2$
- Linear Regression II
  - o Multiple Linear Regression
    - ▪ From a Mathematical Standpoint
  - o Assumptions & Diagnostics
  - o Potential Problems
  - o Research Questions
  - o Variable Selection
  - o Factors
  - o Interactions
  - o Higher-Order Terms

*Week 6*
**Data Science with R - Machine Learning – Part II**
- Lab: Building Bridges
- Generalized Linear Models
  - o Logistic Regression
- The Curse of Dimensionality
  - o Ridge Regression
  - o Lasso Regression
  - o Cross-Validation
  - o Bias/Variance Tradeoff
  - o Density
  - o Principal Component Analysis
- The Curse of Dimensionality
  - o Density
  - o Principal Components Analysis
- Guest Lecture: Dataiku Part I
**Project 3 Due: Python Web Scraping**

*Week 7*
**Data Science with R - Machine Learning – Part III**
- Classification
  - o Feature Selection
  - o Support Vector Machines
  - o Decision Trees
  - o Pruning/Purity/Entropy/GINI
  - o Random Forests
  - o Bagging
  - o Boosting
- Cluster Analysis

- o   K-Means Clustering
- o   Agglomerative Clustering
- o   Hierarchical Clustering
- Neural Networks


*Week 8*
**Data Science with R - Machine Learning – Part IV**
**Introduction to Natural Language Processing**
- Case Study: Spam Detection
- Association Rules
  - o   Market Basket Analysis
- Naïve Bayes Analysis
- Introduction to Natural Language Processing
  - o   Creating corpus: stemming and lemmatization
  - o   POS tag and chunking
  - o   Text classification
- Time Series Analysis
  - o   Smoothing
  - o   Seasonal Decomposition
  - o   ARIMA
- Guest Lecture: Dataiku Part II


*Week 9*
**Data Science with Python - Machine Learning**
- Machine Learning Recap / Linear Regression
  - o   Introduction to scikit learn
  - o   Simple linear regression
  - o   Multiple linear regression
  - o   Stats module
- Classification part I
  - o   Logistic regression
  - o   Discriminant analysis
  - o   Naïve Bayes
- Model Selection
  - o   Cross-validation
  - o   Bootstrap
  - o   Feature selection
  - o   Regularization
  - o   Grid search
- Classification part II

- o   Support vector machine
- o   Decision tree
- o   Random forest
- Unsupervised learning
  - o   Principal Components Analysis
  - o   Kmeans and Hierarchical Clustering

**Project 4 Due: Machine Learning Project (It can be a Kaggle competition, a hiring partner project or a non-profit project from our partners)**

*Week 10*
**Big Data**

- Parallel processing: Introduction to Hadoop and MapReduce
  - o   HDFS
  - o   MapReduce
    - ▪   Conceptual framework
    - ▪   Streaming and Python
  - o   Examples and lab work
- MapReduce design pattern
  - o   Filtering patterns
    - ▪   Simple filtering
    - ▪   Top N
  - o   Summarization patterns
    - ▪   Numerical summarizations
    - ▪   Inverted Index summarizations
- Apache Hive:
  - o   Databases for Hadoop
  - o   Hive
    - ▪   Select
    - ▪   Joins
  - o   Compiling HiveQL to MapReduce
  - o   Technical aspects of Hive
  - o   Extending Hive with TRANSFORM
- Spark
  - o   Basics concepts
    - ▪   RDDs, transformations and actions
    - ▪   PairRDDs
  - o   Examples
    - ▪   Wordcount
    - ▪   Mean and variance

*Week 11*
**Big Data and Algorithms**

- Spark MLlib
- Amazon Web Service
- Introduction to Algorithms
  - Analysis of algorithms: big-O notation
- Sorting
  - Elementary sorts
  - Merge sorts
  - Quick sorts
- Searching
  - Linear search
  - Binary search
  - Hash tables
- Machine Learning Theory Defense Practice

*Week 12*
**Capstone Project Presentations and Review**

- Machine learning theory defense practice
- SQL code review
- R code review
- Python code review
- From the beginning of Bootcamp, you will work on hands-on projects. Now your Capstone Project lets you create your own data product that showcases your interests and talents. Students are free to use anything covered in class on this project.

**Project 5 Due: Capstone Project**