# Data Scientist 1711 Training Syllabus – Student

**Technical Class:**

- Saturday 5-7 pm PST
- Sunday 5-7 pm PST
- Tuesday 6-7 pm PST

**Office Hour:**

- Wednesday & Friday 5-7 pm PST in first 8 weeks

**Mini Project Problem Solving Session:**

- Saturday 3:30 – 4:30 pm PST from week 2 ~ week 8

| Week | Content |
|---|---|
| **Week 1** | **Introduction to Data Application**<br>1. Data science project lifecycle<br>2. Cluster and distributed computing<br>3. Hadoop Eco-system<br>4. HDFS<br>5. Basic Linux Operation |
| | **Python Data Analytics Eco-system**<br>1. What is data scientist<br>2. Key data structures in Python & Numpy<br>3. Pandas for data analytics<br>   • Importing data into Python<br>   • Exploring dataset<br>   • Renaming the columns of a DataFrame<br>   • Filtering a Data Frame<br>   • Basic operations with a Data Frame<br>   • …<br>4. Fast data visualization in Pandas |
| | **Statistical Foundations**<br>1. Probability Distribution: Normal, Binomial, $\chi^2$…<br>2. Central Limit Theorem<br>3. Bayes' Theorem<br>4. Conditional Probability<br>5. Hypothesis Testing: Confidence interval, T-test…<br>6. Sampling: proportion sampling, t-distribition…<br>7. Statistical modeling<br>… |

| Week 2 | Mini Project 1 Session: Sberbank Data Manipulation with Pandas |
|--------|----------------------------------------------------------------|
|        | **Best Practice in Data Processing**<br>1. The importance of data quality<br>2. The data formats and types<br>3. The use of RE and BeautifulSoup to collect data from webpage<br>4. Regular Data Cleaning skills on missing data and outliers |
|        | **Python Machine Learning Eco-system**<br>1. Machine learning introduction<br>&bull; The basic concept of machine learning<br>&bull; Differences between Supervised and Unsupervised machine learning<br>&bull; What can supervised and unsupervised learning do<br>2. Full machine learning flow in Python<br>3. Scikit-learn package<br>4. Basic use of sklearn to build simple regression model<br>5. What is Cross Validation |
|        | **Machine Learning Algorithm -1**<br>**Brief Introduction to Machine Learning Algorithm**<br>1. Supervised Machine Learning vs. Unsupervised Machine Learning<br>2. Regression vs. Classification<br>3. Evaluation Methods for regression<br>4. Evaluation Methods for classification:<br>&bull; how to generate confusion matrix;<br>&bull; how to generate ROC and calculate AUC<br>5. Basic principles of linear regression and logistic regression |

| Week 3 | Mini Project 2 Session : Data Cleansing Practice on Zillow Data |
|---|---|
| | **Data Analysis using Hadoop Hive 1**<br>1. The basic hive concept:<br>   • What is hive<br>   • How hive works<br>   • Hive architecture<br>2. Basic operation of hiveQL |
| | **Supervised Learning: Classification**<br>1. Evaluation Methods of classification<br>2. Basic classification model: logistic regression, decision tree<br>3. Classification Types (how binary and multi-class works)<br>4. Ensemble model method:<br>   • Bagging<br>   • Boosting<br>   • Stacking |
| | **Machine Learning Algorithm -2**<br>**SVM Classifiers**<br>1. Basic principles of SVM<br>2. Know the procedures to derive SVM<br>3. What are Kernels and kernel tricks<br>4. Some basic Kernels such Gaussian Kernel<br>5. Some important parameters such as the slack variable<br>6. What kind of problems can be solved by SVM |

| Week 4 | Mini Project 3 Session: Bank Fraud Detection (binary classification) |
|--------|----------------------------------------------------------------------|
|  | **Analysis using Hadoop Hive 2**<br>1. What is partition table<br>2. The differences between external and internal table<br>3. Advanced use of HiveQL<br>4. The basic use of HiveQL in Spark |
|  | **Supervised Learning: Regression**<br>1. Basic concept of Regression<br>2. Bias-Variance trade off<br>3. Underfitting vs. Overfitting<br>4. Linear regression analytical solution<br>5. Regularization:<br>  &bull; Lasso<br>  &bull; Ridge<br>  &bull; Elastic-Net<br>  &bull; Pros and cons of L1 and L2 regularization<br>6. Advanced techniques in regression<br>  &bull; Gradient Descendent<br>  &bull; Coordinated Descendent<br>  &bull; Stochastic Gradient Descendent<br>  &bull; Random sample consensus (RANSAC) |
|  | **Machine Learning Algorithm -3**<br>**ANN**<br>1. Basic structure of ANN<br>  &bull; Neuron<br>  &bull; Perceptron<br>2. Activation function and the common activation functions<br>3. Procedures of forward propagation and backward propagation<br>4. Derivation of ANN |

| Week 5 | Mini Project 4 Session: History Kaggle Demo: Allstate Claims Severity |
|---|---|
| | **Data Visualization with Tableau**<br>1. Hands-on data visualization & Analysis on Tableau<br>**2.** Business Insights Extraction |
| | **Advanced visualization & A/B Testing**<br>1. Basic & interactive visualization in Python<br>2. Levels of visualization<br>3. Matplotlib<br>  • Basic elements<br>  • Visualization for distribution (histogram, pie chart…)<br>  • Visualization for bi-variable relationship on continuous and categorical features<br>4. Seaborn<br>5. Exploratory analysis<br>6. A/B test and Experimentation<br>  • Business need<br>  • Design Experiment<br>  • Power Analysis<br>  • Analyzed Result |
| | **Machine Learning Algorithm -4**<br>**CNN & RNN**<br>1. What can CNN and RNN work for<br>2. Basic CNN architecture: Convolutional, RELU, Pooling and Fully Connected Layer<br>3. Basic RNN architecture: Forward Propagation and Backward Propagation in RNN<br>4. RNN Example in add operation |

| Week 6 | Mini Project 5 Session: Data Visuliazation with Duoligo User Datasets |
|---|---|
| | **Data Processing using Spark SQL and DataFrame**<br>1. Spark introduction<br> • What is Spark<br> • Why Spark is better<br>2. Spark SQL & data frame<br>3. Spark for Data Analytics<br>4. Demos to fully practice |
| | **Unsupervised Learning: Dimension Reduction**<br>1. Dimension reduction overview<br>2. Dimension reduction methods<br> • Randomized Projection<br> • Principal Component Analysis<br>  o PCA Calculation<br>  o Randomized PCA<br>  o Sparse PCA<br>3. Manifold learning<br>4. Multidimensional Scaling<br> • MDS<br> • Isomap |
| | **Machine Learning Algorithm -5**<br>**Decision Tree & Ensemble Methods**<br>1. Details in Decision Tree<br> • How Decision Tree works<br> • Measures to select the best split<br>2. Details in Ensemble Methods:<br> • Why do we need ensemble model<br> • Committees, Weighted, Predictor of Predictors<br> • Bagging and Boosting<br>3. Random Forest Tree<br>4. Gradient Boosting<br>5. Adaboost |

| Week 7 | Mini Project 6 Session: History Kaggle Demo: Airbnb New User Bookings |
|---|---|
|  | **Machine Learning using Spark MLLib** <br> 1. Relational Database and No-SQL Database <br> 2. Graph Analytics <br> • What is graph database and its applications <br> • Spark GraphX/GraphFrame <br> 3. Machine Learning in Spark <br> 4. Demo Practice by PySpark |
|  | **Unsupervised Learning: Clustering and Outlier Detection** <br> 1. Unsupervised learning introduction <br> 2. Clustering methods & techniques <br> • K-mean Algorithm <br> • Hierarchical Clustering Algorithm <br> • DBSCAN algorithm <br> 3. Outlier and anomaly detection |
|  | **Advanced Python** <br> **Basic CS Algorithm -1** <br> 1. Basic data structure <br> 2. What is Algorithm: <br> • Algorithm Analysis (Time and Space Efficiency) <br> • Theoretical Analysis and Asymptotic Notation ($\text{Big-}\Theta$, $\text{Big-}\Omega$ and $\text{Big-O}$) <br> • Master Theorem <br> 3. Search Algorithm <br> • Sequential search <br> • Binary search <br> 4. Sort Algorithm: Bubble Sort, Selection Sort, Insertion Sort, Shell Sort, Count Sort, Merge Sort <br> 5. Divide and Conquer: Quick Sort |

| Week 8 | Mini Project 7 Session: PySpark Machine Learning |
|--------|--------------------------------------------------|
| | **Real Case Data Processing & Machine Learning in R**<br>Use the skills of what we have learned in R |
| | **Deep Learning**<br>1. Neutal Network Anatomy<br>2. Uniform approximator<br>3. CNN & RNN<br>4. LSTM (Long-Short-Term-Memory)<br>5. Use Keras to build the neural network |
| | **Advanced Python**<br>**Basic CS Algorithm -2**<br>1. Dictionary<br>2. Hashing:<br>    • Hash Code Conventions<br>    • Hash Code Design<br>    • Collision<br>    • Bucketing<br>    • Separate Chaining<br>3. Linear Probing<br>4. Quadratic Probing<br>5. Double Hashing |

**Kaggle**

| | |
|---|---|
| **Week 9** | **Kaggle Introduction**<br>1. What is Kaggle<br>2. Why we need to attend Kaggle<br>3. Tools we will use in Kaggle<br>4. The basic precedures in Kaggle (Feature Engineering, Parameter Tuning, Model Ensemble …)<br><br>**Kaggle 1 (We will cover the following topic with the data of Kaggle Topic we choose)**<br><br>1. Exploratatory Data Analysis (Data Types, Distributions, Missing Values, Correlations …)<br>2. Validation Method<br>3. Feature Engineering on:<br>   • Numetic Features (Log Transformation, Standardization …)<br>   • Categortical Features (One Hot Encoding, Label Encoding, Mean Response Encoding…)<br>   • Missing Value<br>   • Interactions<br>4. Feature Selection Methods:<br>   • Low-variance & High-correlation filters<br>   • Recursive-feature-elimination |
| **Week 10** | **Kaggle 2 (We will cover the following topic with the data of Kaggle Topic we choose)**<br>1. XGBoost highlights and its parameters<br>2. Tuning Process of XGBoost<br>3. Grid Searching and bayesian optimization |
| **Week 11** | **Kaggle 3 (We will cover the following topic with the data of Kaggle Topic we choose)**<br>1. LightGBM highlights and its parameters<br>2. Model Ensemble Details<br>3. Blending with demo code |

**Data Application Lab**

NLP

| Week 9 | **NLP 1**<br>1. Basic NLP Introducton<br>2. The Naïve Bayian Algorithm in NLP |
|---|---|
| Week 10 | **NLP 2**<br>1. Detailed Coding |
| Week 11 | **NLP 3**<br>1. Web Application Development by Flask<br>2. Enhanced NB<br>3. Negation Handling<br>4. Advanced algorithm such as RNN |

FinTech

| Week 12 | **FinTech 1**<br>3. FinTech Domain Knowledge<br>4. Introduction of Lending Club<br>5. How to request data by Lending Club API<br>6. Features at a first look<br>7. Data Preparation |
|---|---|
| Week 13 | **FinTech 2**<br>2. Feature Engineering<br>3. Baseline Model by Logisic Regression<br>4. Gradient Boost Example<br>5. Insights for this project |
| Week 14 | **FinTech 3**<br>5. Web Application Development by Flask<br>• Pickling<br>• Routing<br>• Rendering<br>• Js Basics |

Recommendation System

| Week 15 | **Recommendation System 1**<br>1.    Introduction for this project and its insights<br>2.    Understand the API and the use of Github<br>3.    Learn how to use API to crawl data from Steam<br>4.    Basic Function of Requests and BeautifulSoup |
|---|---|
| Week 16 | **Recommendation System 2**<br>1.    How to use the Github to do the version control<br>2.    The HTTP Basics<br>3.    The API Basics and request structure<br>4.    Process the raw data and set up our own database<br>5.    The knowledge of collaborative filtering<br>6.    The knowledge of content-based filtering<br>7.    The knowledge of popularity based recommendation |
| Week 17 | **Recommendation System 3**<br>1.    Build a recommender engine in Spark<br>2.    Build your own demo with Python Flask |