# Machine Learning Using R, Python and Apache Spark

*Training Course Outline*

Abul Basar, abulbasar@gmail.com, +91 9591755911

**What is Machine Learning?**

Arthur Samuel, a pioneer in Artificial Intelligence (AI) defined machine learning as "field of study that gives computers the ability to learn without being explicitly programmed". The field of machine learning is concerned with the question of how to construct computer programs that automatically improve with experience. In the recent years, many machine learning applications have been developed ranging from data mining that learn to detect fraudulent credit card transactions to information filtering systems that learn reading preferences, to autonomous vehicles that learn to drive on public roads, optical character recognition systems that identify and read the text and objects from images, medical diagnostics and many more. Many researchers also think it is the best way to make progress towards human-level AI.

**Overview of the Training**

This course provides a broad introduction to machine learning, data mining, and statistical pattern recognition. Topics include: (i) Supervised learning (parametric/non-parametric algorithms, support vector machines, kernels, neural networks). (ii) Unsupervised learning (clustering, dimensionality reduction, recommender systems, deep learning). (iii) Best practices in machine learning (bias/variance theory; innovation process in machine learning and AI). You'll learn about not only the theoretical underpinnings of learning, but also gain the practical know-how needed to quickly and powerfully apply these techniques to new problems. You will spend 40% time on theory lectures and case-studies and 60% on hands on lab exercises.

**Programming Language:**

This course will be taught using **R and Python** programming languages for exploring machine learning models. A few demonstrations will be shown in Scala as well.

**Duration of the Training Course:** 5 Days

**Mode:** Instructor led classroom based training, hands on oriented.

**Prerequisite**

- Basic knowledge of statistics, mathematics and computer programming are required to take full advantage of this course.

- Basic programming knowledge in Python, R or Scala will be helpful.

# Day 1 – Introduction to Spark, Python and R for machine learning

| Topic | Description |
|---|---|
| Introduction to Machine Learning | • Machine Learning vs Statistics vs Data Science vs Data Engineering vs Data Analysis<br>• Real world scenarios to understand the perspective<br>• History, Pioneers and Modern Trends<br>• Types of Machine Learning<br>  ◦ Supervised Learning<br>  ◦ Unsupervised Learning<br>  ◦ Reinforced learning |
| Introduction to R and Python for descriptive statistics | • Inspecting Data using R and Python<br>  ◦ Examine data types – quantitative, qualitative, continuous, discrete, ordinal, nominal etc.<br>  ◦ Probability distribution<br>• Common Data Pre-processing Tasks<br>  ◦ Feature Transformation<br>  ◦ Encode Categorical features<br>  ◦ Handle Missing Data<br>  ◦ Handle Outliers<br>  ◦ Remove Problematic Features |
| Lab Exercises | • Getting started with R, R-studio and R packages<br>• Getting started with Python Anaconda distribution for data analysis – Pandas, Scikit, Numpy, Scipy, Matplotlib etc.<br>• Explore built in datasets in R and Python |
| Introduction to Spark, SparkR and Spark MLlib | • Apache Spark – a Swiss army knife for data processing at scale<br>• Apache Spark architectural overview<br>• Data Abstraction for massively parallel processing – RDD, DF, DS<br>• Building and submitting Spark application<br>• Highly available, high throughput data sources<br>• Efficient data formats for large scale data processing |
| Lab Exercise | • Working with stocks market data for S&P 500 companies |

**Learning Goals**

- Spot the use cases of machine learning

- Explore and transform data using R, Python

- Understands Spark framework for large scale, distributed data processing for machine learning

# Day 2 – Supervised Algorithms (Regression)

| Topic | Description |
|---|---|
| Linear Regression | <ul><li>Assumptions of Linear Regression</li><li>Simple Linear Regression</li><li>Multiple Linear Regression</li><li>Cross validation</li><li>Tuning and hyper parameter selection</li><li>Regularization – Ridge, Lasso, Elastic Net for robustness</li><li>Running gradient boosting techniques to solve regression problems</li></ul> |
| Lab Exercises | <ul><li>Build a regression model using on startup dataset. Measure performance, tune model.</li><li>Apply Linear Regression to predict the demand of Electricity Power Demand</li><li>[Advanced] Build regression model on housing price prediction problem from kaggle.com</li></ul> |

**Learning Goals**

- Learn to predict numeric data (for example, revenue, price etc.)

- Learn to evaluate a regression model

- Learn to validate and tune a regression model for best performance

# Day 3 – Supervised Algorithms (Classifications)

| Topic | Description |
|---|---|
| Classification | • What are classification problems and types of classification algorithms<br>• Decision Tree Classification<br>  ◦ Intuition of tree models<br>  ◦ CART regression trees<br>  ◦ Tree pruning<br>  ◦ Missing Data<br>• Classification using Random Forrest<br>• Logistic Regression<br>  ◦ Generalized linear models<br>  ◦ Interpreting coefficients in logistic regression<br>  ◦ Performance metrics – confusion matrix<br>  ◦ AUC value, ROC plot<br>• Support Vector Classifiers<br>• Compare algorithms |
| Lab Exercise | • Evaluate of credit status using customer profile<br>• Segmentation of insurance customers<br>• Wine class prediction based on the dataset from UCI ML library |

**Learning Goals**

- Learn to predict categorical data types, for example – customer segmentation, spam detection etc.

- Evaluate and tune classification models

- Hyper parameter tuning

# Day 4 – Unsupervised algorithms – clustering, PCA and Text Analytics

| Topic | Description |
|---|---|
| Clustering | <ul><li>What is clustering and use cases</li><li>Proximity Matrices – find dissimilarity between two observations<ul><li>Choice of attributes</li><li>Units of measure of attributes</li><li>Importance of scaling</li></ul></li><li>Determining number of clusters</li><li>K-means clustering</li><li>Practical Issues with K-means clustering</li><li>Measure performance of clustering</li><li>Hierarchical clustering</li><li>Practical Issues in clustering</li></ul> |
| Lab Exercises | <ul><li>Clustering movies into genres based on movie attributes</li><li>Apply clustering on images for clinical diagnostics</li></ul> |
| Feature Engineering | <ul><li>Rescaling</li><li>Feature selection techniques</li><li>Principle Component Analysis</li></ul> |
| Lab Exercise | <ul><li>Running PCA on wine classification dataset from UCI ML dataset</li></ul> |
| Text Analytics | <ul><li>Text Analysis Steps</li><li>Term Frequency – Inverse Term Frequency</li><li>POS tagging, lemmatization, stemming</li><li>Collecting raw text</li><li>Representing text</li><li>Categorize Documents by type</li><li>Determining sentiments</li><li>Gaining insights</li></ul> |
| Lab Exercises | <ul><li>Classification on IMDB comments dataset</li></ul> |

**Learning Goals**

- Apply clustering algorithms to gain insights about latent groups inside a dataset

- Identify properties of the features that impact model performance

- Transform data structure to facilitate model performance and computation speed

- Learn technique to analyze textual data and prepare it for further analysis such as classification, regression, clustering etc.

- Extract information from text in the form of word cloud, sentiment

# Day 5 – Scalable Machine Learning and Machine Learning Deployments

| Topic | Description |
|---|---|
| Scalable Machine Learning | <ul><li>Challenges with large scale machine learning</li><li>Deep dive into Spark MLlib</li><li>Spark R integrate between R with Spark</li><li>Spark ML packages for regression, classification, clustering, recommender system</li><li>Saving the trained model on persistent storage</li><li>Serving model for online prediction</li><li>[Optional] Using Prediction IO for deployment machine learning pipelines</li><li>[Optional] Machine learning for streaming application</li></ul> |
| Lab Exercise | <ul><li>Process data from stack-overflow dataset</li><li>Find community from stack-overflow dataset using tag analysis</li></ul> |

**Learning Goals**

- Learn how to build machine learning pipeline on massive data using Spark framework

- Deployment options for machine learning applications for online and offline applications.