# PySpark and Machine Learning   ( 3 Days )

**Goal - Make Practical Machine Learning Scalable & Easy**

---

**Objectives**

- Learn about Apache Spark and the Spark 2.0 architecture
- Build and interact with Spark DataFrames using Spark SQL
- Learn how to solve graph and deep learning problems using GraphFrames and TensorFrames respectively
- Read, transform, and understand data and use it to train machine learning models
- Build machine learning models with MLlib and ML
- Learn how to submit your applications programmatically using spark-submit
- ML Algorithms: common learning algorithms such as classification, regression, clustering, and collaborative filtering
- Featurization: feature extraction, transformation, dimensionality reduction, and selection
- Pipelines: tools for constructing, evaluating, and tuning ML Pipelines
- Persistence: saving and load algorithms, models, and Pipelines
- Utilities: linear algebra, statistics, data handling, etc.


**Who is the target audience?**
- Anyone interested in Machine Learning
- Any intermediate level people who know the basics of machine learning, including the classical algorithms like linear regression or logistic regression, but who want to learn more about it and explore all the different fields of Machine Learning.
- Any people who are not that comfortable with coding but who are interested in Machine Learning and want to apply it easily on datasets.
- Any data analysts who want to level up in Machine Learning.
- Any people who are not satisfied with their job and who want to become a Data Scientist.
- Any people who want to create added value to their business by using powerful Machine Learning tools


**Prerequisites:**

Knowledge prerequisites

- Basic Python data structures
- Basic knowledge of Pandas dataframes and SQL
- Knowledge of common data storage formats like JSON, delimiter separated files, HDFS, etc
- Entry-level machine learning

**Software Prerequisites**

- Apache Spark (Downloadable from http://spark.apache.org/downloads.html)
- A Python distribution containing IPython, Pandas and Scikit-learn
- PySpark

**Course Agenda**

| Days | Modules | Course Outline |
|---|---|---|
| 1 | Module 1 | **Introduction to Spark**<br>What is Apache Spark?<br>Spark Jobs and APIs<br>Spark 2.0 architecture<br>Installation and Configuration |
| | Module 2 | **Resilient Distributed Datasets**<br>Internal workings of an RDD<br>Creating RDDs<br>Global versus local scope<br>Transformations<br>Actions<br>Hands on Session on RDD and Spark<br>Assignments 1<br>Best Practices 1 |
| Day 2 | Module 3 | **DataFrames**<br>Python to RDD communications<br>Catalyst Optimizer refresh<br>Speeding up PySpark with DataFrames<br>Creating DataFrames<br>Simple DataFrame queries<br>Interoperating with RDDs<br>Querying with the DataFrame API<br>Hands On Session on Pandas DataFrame and PySpark<br>Assignments 2 |
| | Module 4 | **Prepare Data for Modeling**<br>Checking for duplicates, missing observations, and outliers<br>Getting familiar with your data Visualization<br>Hands on Session Data Modeling<br>Assignments 3 |

| Days | Modules | Course Outline |
|---|---|---|
| **Day 3** | **Module 5** | **Introducing MLlib**<br>Overview of the package<br>Loading and transforming the data<br>Getting to know your data<br>Creating the final dataset<br>Predicting infant survival<br>Hands on Session using PySpark MLib<br>Assignments 4 |
| | | **Introducing the ML Package**<br>Overview of the package<br>Predicting the chances of infant survival with ML<br>Parameter hyper-tuning<br>Other features of PySpark ML in action<br>Implementation of ML Algorithm<br>• Random Forest<br>• Regression<br>• K-means<br>Assignments 5 |
| | | **GraphFrames**<br>Introducing GraphFrames<br>Installing GraphFrames<br>Preparing your flights dataset<br>Building the graph<br>Executing simple queries<br>Understanding vertex degrees<br>Determining the top transfer airports<br>Understanding motifs<br>Determining airport ranking using PageRank<br>Determining the most popular non-stop flights<br>Using Breadth-First Search<br>Visualizing flights using D3<br>Assignment 6<br>**Conclusion and Summary** |