



Data Analyst

## HDP Analyst: Data Science

### Overview

This course Provides instruction on the processes and practice of data science, including machine learning and natural language processing. Included are: tools and programming languages (Python, IPython, Mahout, Pig, NumPy, pandas, SciPy, Scikit-learn), the Natural Language Toolkit (NLTK), and Spark MLlib.

### Duration

3 days

### Target Audience

Architects, software developers, analysts and data scientists who need to apply data science and machine learning on Hadoop.

### Course Objectives

- Recognize use cases for data science on Hadoop
- Describe the Hadoop and YARN architecture
- Describe supervised and unsupervised learning differences
- Use Mahout to run a machine learning algorithm on Hadoop
- Describe the data science life cycle
- Use Pig to transform and prepare data on Hadoop
- Write a Python script
- Describe options for running Python code on a Hadoop cluster
- Write a Pig User-Defined Function in Python
- Use Pig streaming on Hadoop with a Python script
- Use machine learning algorithms
- Describe use cases for Natural Language Processing (NLP)
- Use the Natural Language Toolkit (NLTK)
- Describe the components of a Spark application
- Write a Spark application in Python
- Run machine learning algorithms using Spark MLlib
- Take data science into production

### Prerequisites

Students must have experience with at least one programming or scripting language, knowledge in statistics and/or mathematics, and a basic understanding of big data and Hadoop principles. Students new to Hadoop are encouraged to attend the *HDP Overview: Apache Hadoop Essentials* course.

### Hands-On Content

- Lab: Setting Up a Development Environment
- Demo: Block Storage
- Lab: Using HDFS Commands
- Demo: MapReduce
- Lab: Using Apache Mahout for Machine Learning
- Demo: Apache Pig
- Lab: Getting Started with Apache Pig
- Lab: Exploring Data with Pig
- Lab: Using the IPython Notebook
- Demo: The NumPy Package
- Demo: The pandas Library
- Lab: Data Analysis with Python
- Lab: Interpolating Data Points
- Lab: Defining a Pig UDF in Python
- Lab: Streaming Python with Pig
- Demo: Classification with Scikit-Learn
- Lab: Computing K-Nearest Neighbor
- Lab: Generating a K-Means Clustering
- Lab: POS Tagging Using a Decision Tree
- Lab: Using NLTK for Natural Language Processing
- Lab: Classifying Text using Naive Bayes
- Lab: Using Spark Transformations and Actions
- Lab Using Spark MLlib
- Lab: Creating a Spam Classifier with MLlib

### Format

50% Lecture/Discussion

50% Hands-on Labs

### Certification

Hortonworks offers a comprehensive certification program that identifies you as an expert in Apache Hadoop. Visit [hortonworks.com/training/certification](http://hortonworks.com/training/certification) for more information.

### Hortonworks University

Hortonworks University is your expert source for Apache Hadoop training and certification. Public and private on-site courses are available for developers, administrators, data analysts and other IT professionals involved in implementing big data solutions. Classes combine presentation material with industry-leading hands-on labs that fully prepare students for real-world Hadoop scenarios.



#### About Hortonworks

Hortonworks develops, distributes and supports the only 100 percent open source distribution of Apache Hadoop explicitly architected, built and tested for enterprise-grade deployments.

**US:** 1.855.846.7866

**International:** +1.408.916.4121

[www.hortonworks.com](http://www.hortonworks.com)

5470 Great America Parkway  
Santa Clara, CA 95054 USA