



# Hortonworks Data Platform Analyst—Data Science (EDU-PRIV-DATASCI-200) H7G69S

<b>HPE course number</b>	H7G69S
<b>Course length</b>	3 days
<b>Delivery mode</b>	ILT
<b>View schedule, local pricing, and register</b>	<a href="#">View now</a>
<b>View related courses</b>	<a href="#">View now</a>

Data Science for the Hortonworks Data Platform covers data science principles and techniques through lecture and hands-on experience. During this three-day course, students will learn the processes and practice of data science, including machine learning and natural language processing. Students will also learn the tools and programming languages used by data scientists, including Python, IPython, Mahout, Pig, NumPy, pandas, SciPy, Scikit-learn, the Natural Language Toolkit (NLTK), and Spark MLlib.

## Why HPE Education Services?

- IDC MarketScape leader 4 years running for IT education and training\*
- Recognized by IDC for leading with global coverage, unmatched technical expertise, and targeted education consulting services\*
- Key partnerships with industry leaders OpenStack®, VMware®, Linux®, Microsoft®, ITIL, PMI, CSA, and (ISC)²
- Complete continuum of training delivery options—self-paced eLearning, custom education consulting, traditional classroom, video on-demand instruction, live virtual instructor-led with hands-on lab, dedicated onsite training
- Simplified purchase option with HPE Training Credits

## Audience

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

## 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

## Course objectives

At the completion of the course, students will be able to:

- Recognize use cases for data science
- Describe the architecture of Hadoop and YARN
- Explain the differences between supervised and unsupervised learning

- List the six machine learning tasks
- Recognize use cases for clustering, outlier detection, affinity analysis, classification, regression, and recommendation
- Use Mahout to run a machine learning algorithm on Hadoop
- Write Pig scripts to transform data on Hadoop
- Use Pig to prepare data for a machine learning algorithm
- Write a Python script
- Use NumPy to analyze Big Data
- Use the data structure classes in the pandas library
- Write a Python script that invokes a SciPy machine learning algorithm
- Explain the 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

- Write a Python script that invokes a scikit-learn machine learning algorithm
- Use the k-nearest neighbor algorithm to predict values based on a training data set
- Run a machine learning algorithm on a distributed data set on Hadoop
- Describe use cases for Natural Language Processing (NLP)
- Perform sentence segmentation on a large body of text
- Perform part-of-speech tagging
- Use the Natural Language Toolkit (NLTK) for implementing NLP tasks and machine learning algorithms
- Explain the components of a Spark application
- Write a Spark application in Python
- Run machine learning algorithms on Hadoop using Spark MLlib

## **Benefits to you**

- This course will provide in depth explanation on how to perform Hadoop 2.0 application development

## Detailed course outline

---

<b>Day 1</b>	<ul style="list-style-type: none"><li>• Unit 1: Using Hadoop for Data Science</li><li>• Unit 2: Hadoop Architecture</li><li>• Unit 3: Machine Learning</li><li>• Unit 4: Introduction to Pig</li></ul>
<b>Day 2</b>	<ul style="list-style-type: none"><li>• Unit 5: Python Programming</li><li>• Unit 6: Analyzing Data with Python</li><li>• Unit 7: Running Python on Hadoop</li></ul>
<b>Day 3</b>	<ul style="list-style-type: none"><li>• Unit 8: Implementing Machine Learning</li><li>• Unit 9: Natural Language Processing</li><li>• Unit 10: Spark MLlib</li></ul>
<b>Hands-On Labs</b>	<p>Students will complete the following hands-on labs using their own 7-node Hadoop cluster (HDP 2.1) and IPython Notebook:</p> <ul style="list-style-type: none"><li>• Setting Up a Development Environment</li><li>• Using HDFS Commands</li><li>• Using Mahout for Machine Learning</li><li>• Getting Started with Pig</li><li>• Exploring Data with Pig</li><li>• Using the IPython Notebook</li><li>• Data Analysis with Python</li><li>• Interpolating Data Points</li><li>• Define a Pig UDF in Python</li><li>• Streaming Python with Pig</li><li>• K-Nearest Neighbor</li><li>• K-Means Clustering</li><li>• Using NLTK for Natural Language Processing</li><li>• Classifying Text using Naive Bayes</li><li>• Spark Programming</li><li>• Running Data Science Algorithms using Spark MLlib</li></ul>

---

Learn more at  
[hpe.com/ww/learnbigdata](http://hpe.com/ww/learnbigdata)

**Follow us:**



---

© Copyright 2015–2016 Hewlett Packard Enterprise Development LP. The information contained herein is subject to change without notice. The only warranties for Hewlett Packard Enterprise products and services are set forth in the express warranty statements accompanying such products and services. Nothing herein should be construed as constituting an additional warranty. Hewlett Packard Enterprise shall not be liable for technical or editorial errors or omissions contained herein.

Microsoft is either a registered trademark or trademark of Microsoft Corporation in the United States and/or other countries. The OpenStack Word Mark is either a registered trademark/service mark or trademark/service mark of the OpenStack Foundation, in the United States and other countries and is used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation or the OpenStack community. Pivotal and Cloud Foundry are trademarks and/or registered trademarks of Pivotal Software, Inc. in the United States and/or other countries. Linux is the registered trademark of Linus Torvalds in the U.S. and other countries. VMware is a registered trademark or trademark of VMware, Inc. in the United States and/or other jurisdictions. All other third-party trademark(s) is/are property of their respective owner(s).