

TRAINING OFFERING | DEV-343

HORTONWORKS DATA PLATFORM (HDP®) DEVELOPER: Apache Spark 2 for Developers.

4 DAYS  
INTERMEDIATE DEVELOPER

This course introduces the Apache Spark distributed computing engine, and is suitable for developers, data analysts, architects, technical managers, and anyone who needs to use Spark in a hands-on manner. It is based on the Spark 2.x release. The course provides a solid technical introduction to the Spark architecture and how Spark works. It covers the basic building blocks of Spark (e.g. RDDs and the distributed compute engine), as well as higher-level constructs that provide a simpler and more capable interface. It includes in-depth coverage of Spark SQL, DataFrames, and Datasets, which are now the preferred programming API. This includes exploring possible performance issues and strategies for optimization. The course also covers more advanced capabilities such as the use of Spark Streaming to process streaming data, and integrating with the Kafka server.

PREREQUISITES

Students should be familiar with programming principles and have previous experience in software development using Scala, Java, or Python. Previous experience with data streaming, SQL, and HDP is also helpful, but not required.

TARGET AUDIENCE

Software engineers that are looking to develop in-memory applications for time sensitive and highly iterative applications in an Enterprise HDP environment.

FORMAT

50% Lecture/Discussion

50% Hands-On Labs

AGENDA SUMMARY

Day 1: Scala Ramp Up, Introduction to Spark

Day 2: RDDs and Spark Architecture, Spark SQL, DataFrames and Datasets

Day 3: Shuffling, Transformations and Performance, Performance Tuning

Day 4: Creating Standalone Applications and Spark Streaming

DAY 1 OBJECTIVES

- Scala Introduction
- Working with:
 - Variables
 - Data Types
 - Control Flow
- The Scala Interpreter
- Collections and their Standard Methods (e.g. map())
- Working with:
 - Functions
 - Methods
 - Function Literals
- Define the Following as they Relate to Scale:
 - Class
 - Object
 - Case Class
- Overview, Motivations, Spark Systems
- Spark Ecosystem
- Spark vs. Hadoop
- Acquiring and Installing Spark
- The Spark Shell, SparkContext

DAY 1 LABS

- Setting Up the Lab Environment
- Starting the Scala Interpreter
- A First Look at Spark
- A First Look at the Spark Shell

About Hortonworks

Hortonworks is a leading innovator at creating, distributing and supporting enterprise-ready open data platforms. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. Our open Connected Data Platforms power Modern Data Applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest. Along with our 1600+ partners, we provide the expertise, training and services that allows our customers to unlock the transformational value of data across any line of business. We are Powering the Future of Data™.

Contact

For further information visit
www.hortonworks.com

+1 408 675-0983
+1 855 8-HORTON
INTL: +44 (0) 20 3826 1405

DAY 2 OBJECTIVES

- RDD Concepts, Lifecycle, Lazy Evaluation
- RDD Partitioning and Transformations
- Working with RDDs Including:
 - Creating and Transforming (map, filter, etc.)
- An Overview of RDDs
- SparkSession, Loading/Saving Data, Data Formats (JSON, CSV, Parquet, text ...)
- Introducing DataFrames and DataSets (Creation and Schema Inference)
- Identify Supported Data Formats, Including:
 - JSON
 - Text
 - CSV
 - Parquet
- Working with the DataFrame (untyped) Query DSL, including:
 - Column
 - Filtering
 - Grouping
 - Aggregation
- SQL-based Queries
- Working with the Dataset (typed) API
- Mapping and Splitting (flatMap(), explode(), and split())
- Datasets vs. DataFrames vs. RDDs

DAY 2 LABS

- RDD Basics
- Operations on Multiple RDDs
- Data Formats
- Spark SQL Basics
- DataFrame Transformations
- The Dataset Typed API

About Hortonworks

Hortonworks is a leading innovator at creating, distributing and supporting enterprise-ready open data platforms. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. Our open Connected Data Platforms power Modern Data Applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest. Along with our 1600+ partners, we provide the expertise, training and services that allows our customers to unlock the transformational value of data across any line of business. We are Powering the Future of Data™.

Contact

For further information visit
www.hortonworks.com

+1 408 675-0983
+1 855 8-HORTON
INTL: +44 (0) 20 3826 1405

- Splitting Up Data

DAY 3 OBJECTIVES

- Working with:
 - Grouping
 - Reducing
 - Joining
- Shuffling, Narrow vs. Wide Dependencies, and Performance Implications
- Exploring the Catalyst Query Optimizer [explain(), Query Plans, Issues with lambdas]
- The Tungsten Optimizer (Binary Format, Cache Awareness, Whole-Stage Code Gen)
- Discuss Caching, Including:
 - Concepts
 - Storage Type
 - Guidelines
- Minimizing Shuffling for Increased Performance
- Using Broadcast Variables and Accumulators
- General Performance Guidelines
 - Using the Spark UI
 - Efficient Transformations
 - Data Storage
 - Monitoring

DAY 3 LABS

- Exploring Group Shuffling
- Seeing Catalyst at Work
- Seeing Tungsten at Work
- Working with Caching, Joins, Shuffles, Broadcasts, Accumulators
- Broadcast General Guidelines

About Hortonworks

Hortonworks is a leading innovator at creating, distributing and supporting enterprise-ready open data platforms. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. Our open Connected Data Platforms power Modern Data Applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest. Along with our 1600+ partners, we provide the expertise, training and services that allows our customers to unlock the transformational value of data across any line of business. We are Powering the Future of Data™.

Contact

For further information visit
www.hortonworks.com

+1 408 675-0983
+1 855 8-HORTON
INTL: +44 (0) 20 3826 1405

DAY 4 OBJECTIVES

- Core API, SparkSession.Builder
- Configuring and Creating a SparkSession
- Building and Running Applications - sbt/build.sbt and spark-submit
- Application Lifecycle (Driver, Executors, and Tasks)
- Cluster Managers (Standalone, YARN, Mesos)
- Logging and Debugging
- Introduction and Streaming Basics
- Spark Streaming (Spark 1.0+)
 - DStreams, Receivers, Batching
 - Stateless Transformation
 - Windowed Transformation
 - Stateful Transformation
- Structured Streaming (Spark 2+)
 - Continuous Applications
 - Table Paradigm, Result Table
 - Steps for Structured Streaming
 - Sources and Sinks
- Consuming Kafka Data
 - Kafka Overview
 - Structured Streaming - "kafka" Format
 - Processing the Stream
- OPTIONAL: Structured Streaming Sessionization (example Scala Session)
- OPTIONAL: Time Series Analysis (example PySpark SQL Window function, moving/rolling average)
- OPTIONAL: Use Case Discovery (language agnostic)

About Hortonworks

Hortonworks is a leading innovator at creating, distributing and supporting enterprise-ready open data platforms. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. Our open Connected Data Platforms power Modern Data Applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest. Along with our 1600+ partners, we provide the expertise, training and services that allows our customers to unlock the transformational value of data across any line of business. We are Powering the Future of Data™.

Contact

For further information visit
www.hortonworks.com

+1 408 675-0983
+1 855 8-HORTON
INTL: +44 (0) 20 3826 1405

DAY 4 LABS

- Spark Job Submission
- Additional Spark Capabilities
- Spark Streaming
- Spark Structured Streaming
- Spark Structured Streaming with Kafka
- OPTIONAL: Structured Streaming Sessionization (*Scala only for now*)
- OPTIONAL: Time Series Analysis (*PySpark only for now*)
- OPTIONAL: Use Case Discovery workshop (*Agile-Scrum centric, language agnostic*)

Revised 12/19/2018

LCW

About Hortonworks

Hortonworks is a leading innovator at creating, distributing and supporting enterprise-ready open data platforms. Our mission is to manage the world's data. We have a single-minded focus on driving innovation in open source communities such as Apache Hadoop, NiFi, and Spark. Our open Connected Data Platforms power Modern Data Applications that deliver actionable intelligence from all data: data-in-motion and data-at-rest. Along with our 1600+ partners, we provide the expertise, training and services that allows our customers to unlock the transformational value of data across any line of business. We are Powering the Future of Data™.

Contact

For further information visit
www.hortonworks.com

+1 408 675-0983
+1 855 8-HORTON
INTL: +44 (0) 20 3826 1405

© 2011-2019 Hortonworks Inc. All Rights Reserved.
[Privacy Policy](#) | [Terms of Service](#)