

Cloudera Developer Training for Apache Spark™ and Hadoop

32 hours

Course overview:

This four-day hands-on training course delivers the key concepts and expertise developers need to use Apache Spark to develop high-performance parallel applications. Participants will learn how to use Spark SQL to query structured data and Spark Streaming to perform real-time processing on streaming data from a variety of sources. Developers will also practice writing applications that use core Spark to perform ETL processing and iterative algorithms.

The course covers how to work with “big data” stored in a distributed file system, and execute Spark applications on a Hadoop cluster.

After taking this course, participants will be prepared to face real-world challenges and build applications to execute faster decisions, better decisions, and interactive analysis, applied to a wide variety of use cases, architectures, and industries.

Course Objectives:

- How the Apache Hadoop ecosystem fits in with the data processing lifecycle
- How data is distributed, stored, and processed in a Hadoop cluster
- How to write, configure, and deploy Apache Spark applications on a Hadoop cluster
- How to use the Spark shell and Spark applications to explore, process, and analyze distributed data
- How to query data using Spark SQL, DataFrames, and Datasets
- How to use Spark Streaming to process a live data stream

Who Should Attend:

This course is designed for developers and engineers who have programming experience, but prior knowledge of Hadoop and/or Spark is not required.

- Apache Spark examples and hands-on exercises are presented in Scala and Python. The ability to program in one of those languages is required.
- Basic familiarity with the Linux command line is assumed.
- Basic knowledge of SQL is helpful

Course Contents:

Introduction to Apache Hadoop and the Hadoop Ecosystem

- Apache Hadoop Overview
- Data Processing
- Introduction to the Hands-On Exercises

Apache Hadoop File Storage

- Apache Hadoop Cluster Components
- HDFS Architecture
- Using HDFS

Distributed Processing on an Apache Hadoop Cluster

- YARN Architecture
- Working With YARN

Apache Spark Basics

- What is Apache Spark?
- Starting the Spark Shell
- Using the Spark Shell
- Getting Started with Datasets and DataFrames
- DataFrame Operations

Working with DataFrames and Schemas

- Creating DataFrames from Data Sources
- Saving DataFrames to Data Sources
- DataFrame Schemas
- Eager and Lazy Execution

Analyzing Data with DataFrame Queries

- Querying DataFrames Using Column Expressions
- Grouping and Aggregation Queries
- Joining DataFrames

RDD Overview

- RDD Overview
- RDD Data Sources
- Creating and Saving RDDs
- RDD Operations

Transforming Data with RDDs

- Writing and Passing Transformation Functions
- Transformation Execution
- Converting Between RDDs and DataFrames

Aggregating Data with Pair RDDs

- Key-Value Pair RDDs
- Map-Reduce
- Other Pair RDD Operations

Querying Tables and Views with SQL

- Querying Tables in Spark Using SQL
- Querying Files and Views
- The Catalog API

Working with Datasets in Scala

- Datasets and DataFrames
- Creating Datasets
- Loading and Saving Datasets
- Dataset Operations

Writing, Configuring, and Running Apache Spark Applications

- Writing a Spark Application
- Building and Running an Application
- Application Deployment Mode
- The Spark Application Web UI
- Configuring Application Properties

Spark Distributed Processing

- Review: Apache Spark on a Cluster
- RDD Partitions
- Example: Partitioning in Queries
- Stages and Tasks
- Job Execution Planning
- Example: Catalyst Execution Plan
- Example: RDD Execution Plan

Distributed Data Persistence

- DataFrame and Dataset Persistence
- Persistence Storage Levels
- Viewing Persisted RDDs

Common Patterns in Spark Data Processing

- Common Apache Spark Use Cases
- Iterative Algorithms in Apache Spark
- Machine Learning
- Example: k-means

Introduction to Structured Streaming

- Apache Spark Streaming Overview
- Creating Streaming DataFrames
- Transforming DataFrames
- Executing Streaming Queries

Structured Streaming with Apache Kafka

- Overview
- Receiving Kafka Messages
- Sending Kafka Messages

Aggregating and Joining Streaming DataFrames

- Streaming Aggregation
- Joining Streaming DataFrames

Message Processing with Apache Kafka

- What Is Apache Kafka?
- Apache Kafka Overview
- Scaling Apache Kafka
- Apache Kafka Cluster Architecture
- Apache Kafka Command Line Tools