

Getting Productive with Spark - 2 Days

Course Outline

Apache Spark is the next-generation successor to MapReduce. Spark is a powerful, open-source processing engine for data in the Hadoop cluster, optimized for speed, ease of use, and sophisticated analytics. The Spark framework supports streaming data processing and complex, iterative algorithms, enabling applications to run up to 100x faster than traditional Hadoop MapReduce programs.

The 2 day Spark course is aimed at developers who are encountering Spark for the first time and want to understand how to build Big Data Products with Spark. The course would enable participants to build complete, unified Big Data applications combining batch, streaming, and interactive analytics on all their data.

Developers would be able to write sophisticated parallel applications to execute faster decisions, better decisions, and real-time actions, applied to a wide variety of use cases, architectures, and industries.

The course has a practical focus, mixing presentation with in-depth hands-on labs and exercises.

Proposed Structure

Day 1 First brush <ul style="list-style-type: none">• Big Data Why and What?• Introduction to Spark• Spark shell• Programming with Spark	Day 2 Spark Streaming <ul style="list-style-type: none">• Overview• Streaming operations• Sliding window operations• Streaming Applications
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Developer to Developer Training Series

<p>RDDs</p> <ul style="list-style-type: none">• Resilient Distributed Datasets• RDD Operations• Map Reduce• Key value pair <p>Running Spark</p> <ul style="list-style-type: none">• Stand alone• Web UI• Stand alone cluster• Building and running <p>Parallel Programming</p> <ul style="list-style-type: none">• Partitions and Data Locality• Executing parallel operations <p>Caching and Persistence</p> <ul style="list-style-type: none">• Caching Overview• Distributed Persistence	<p>Good to know</p> <ul style="list-style-type: none">• Spark Context• Spark Properties• Logging• Iterative Algorithms• Graph Analysis• Machine Learning• Spark SQL <p>Spark and Hadoop</p> <ul style="list-style-type: none">• HDFS• Using HDFS with Spark• Spark and MapReduce <p>Q&A</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Course prerequisites

To benefit from this course you should have programming experience with Scala or with Python. The language of instruction is Scala. Basic Linux knowledge is expected.

For more information on the course or a discussion on your custom need, send a mail to info@knoldus.com