

Big Data - Hadoop

Description

Data scientists build information platforms to provide deep insight and answer previously unimaginable questions. Hadoop are transforming how data scientists work by allowing interactive and iterative data analysis at scale. Learn how Hadoop enable data scientists to help companies reduce costs, increase profits, improve products, retain customers, and identify new opportunities. This course helps participants understand what data scientists do, the problems they solve, and the tools and techniques they use. Through in-class simulations, participants apply data science methods to real-world challenges in different industries and, ultimately, prepare for data scientist roles in the field.

Expectations and Goals

Through instructor-led discussion and interactive, hands-on exercises, participants will navigate the Hadoop ecosystem, and develop concrete skills such as:

- How to identify potential business use cases where data science can provide impactful results?
- How to obtain, clean and combine disparate data sources to create a coherent picture for analysis?
- What statistical methods to leverage for data exploration that will provide critical insight into your data?
- Where and when to leverage Hadoop streaming and Apache Flume for data science pipelines? What machine learning technique to use for a particular data science project?

Prerequisites

- Object Oriented Programming in Java, Exception in Java
- Knowledge of SQL Command
- Basic Command in Linux

Course Schedule

Module	Topic
Module 1	Introduction to Big Data Features of Hadoop Components in Hadoop Concept of Hadoop Ecosystem Introduction to HDFS HDFS Practical
Module 2	Collection Framework in Java (List, Map, Iterator) String Tokenizer, File Handling, String Handling Concept of Map Reduce Map Reduce Practical
Module 3	Introduction to Hive Why we need Hive Architecture of Hive Hive Data Types Hive Complex Datatypes Managed Tables External Tables Running Hive Queries Perform Joining Handling JSON Data

	<p>Handling XML Data</p> <p>Partitioned Table</p> <p>Hive UDF</p> <p>Scripting in Hive</p> <p>Performance Tuning in Hive</p> <p>Case Study in Hive based on Dataset</p>
Module 4	<p>Introduction to Pig</p> <p>Why we need Pig Technology</p> <p>Architecture of Pig</p> <p>Pig Data Types</p> <p>Different Modes in Pig</p> <p>Running Pig Command</p> <p>Tuple, Bag, Map</p> <p>Pig UDF</p> <p>Script in Pig</p> <p>Case Study in Pig Based on Dataset</p>
Module 5	<p>Introduction to SQOOP</p> <p>Importing and Exporting the RDBMS to HDFS</p> <p>Import data from RDBMS to Hive</p> <p>Export Data from Hive to RDBMS</p> <p>Assignment</p>
Module 6	<p>Introduction to Flume</p> <p>Introduction Source, Sink, Flume Agents</p> <p>Fetching Twitter Data into Solr</p> <p>Configuration to create twitter data into HDFS</p> <p>Use HiveSerde to Analyze the data</p>
Module 7	Sentiment Analysis based on Twitter Data
Module 8	<p>Concept of Impala</p> <p>Running the Queries on Impala</p> <p>Compare Impala with Hive</p> <p>Concept of HUE</p> <p>Access The Hadoop Component by means of HUE</p>
Module 9	<p>Introduction to NOSQL Database</p> <p>Compare between NOSQL and RDBMS</p> <p>Introduction to HBASE</p> <p>Why we need HBASE</p> <p>HBASE Commands</p>
Module 10	<p>Installation of Tableau</p> <p>Communicate Tableau with impala</p> <p>Plotting the Graph</p>
Module 11	Project Work and Documentation