

Spark and Scala :- 25 days – 25 hrs

Describe Features of Apache Spark

- How Spark fits in Big Data ecosystem
- Why Spark & Hadoop fit together

Define Spark Components

- Driver Program
- o Spark Context
- Cluster Manager
- Worker
- o Executor
- o Task
- Spark RDD
- o Spark Context
- Spark Libraries

Load data into Spark

- Different data sources and formats
- o HDFS
- o Amazon S3
- o Local File System
- o Text
- o JSON
- o CSV
- o Sequence File
- Create & Use RDD, Data Frames

Apply dataset operations to Resilient Distributed Datasets

- Transformation
- Actions
- Cache Intermediate RDD
- o Lineage Graph
- o Lazy Evaluation

Use Spark DataFrames for simple queries

- Create Data Frame
- Spark Interactive shell (Scala & Python)
- Spark SQL

Define different ways to run your application

Build and launch a standalone application

- Spark Program Life Cycle
- Function of Spark Context
- Different Way to Launch Spark Application
 - o Local
 - o Standalone
 - o Hadoop YARN
 - o Apache Mesos
- Launch Spark Application
 - o Spark-Submit
 - o Monitor the Spark Job

Describe & Create pair RDD

- Key-Value pair
- Apache Spark vs Apache Hadoop MapReduce
- Create RDD from existing non-pair RDD
- Create pair RDD by loading certain formats
- Create pair RDD from in-memory collection of pairs

Apply Operations on pair RDD

- Group ByKey
- Reduce ByKey
- Other Transformations
 - o Joins

Control partitioning across nodes

- RDD Partition
- Types of Partition
 - o Hash Partitioning
 - o Range Partitioning

- Benefit of Partitioning
- Best Practices

More on Data Frames

- Explore Data in DataFrames
- Create UDFs (user define functions)
 - o UDF with Scala DSL
 - o UDF with SQL
- Repartition Data Frames.
- Infer Schema by Reflection
- DataFrame from database table
- DataFrame from JSON

Monitor Apache Spark Applications

- Spark Execution Model
- Debug and Tune Spark Applications

Identify Spark Unified Stack Components

- Spark SQL
- Spark Streaming
- Spark MLlib
- Spark GraphX

Benefits of Apache Spark over Hadoop Ecosystem

Describe Spark Data pipeline Use Cases

- Spark Streaming Architecture
- Dstream and a spark streaming application
 - o Define Use Case (Time Series Data)
 - o Basic Steps
 - o Save Data to HBase
- Operations on DStream
 - o Transformations
 - o Data Frame and SQL Operations
- Define Windowed Operation
 - o Sliding Window

- o Windowed Computation
- o Window based Transformation
- o Window Operations
- Fault tolerance of streaming applications
- o Fault Tolerance in Spark Streaming
- o Fault Tolerance in Spark RDD
- o Check pointing

Describe Graph X

Define Regular, Directed, and property graphs

Create a Property Graph

Perform Operations on Graphs

Describe Apache Spark MLib

Describe the Machine Learning Techniques

- Classifications
- Clustering
- Collaborative Filtering

Use Collaborative filtering to predict user choice

Scala

- Introduction
- A first example
- Expressions and Simple Functions
- First Class function
- Classes and Objects
- Case classes and Pattern matching
- Generic types and methods
- Lists
- For- Comprehension
- Mutable State
- Computing with Streams
- Lazy Values
- Implicit Parameters and Conversions

- Handley / Milner type Interface
- Abstraction for concurrency