**Pre-requisite: Basic understanding of Linux and RDBMS**

**Introduction**

What is Big Data?
- Hadoop – What, Why & Which?
  o History.
  o Hadoop components.
  o Distributions available – Apache, Horton, Cloudera, Intel, R, Greenplum, etc.
  o Installations – Pseudo and multinode cluster
- Hadoop cluster hardware & network considerations
  o Sizing standards
  o Network recommendations
  o Storage recommendations
- Rack awareness
- Concepts (horizontal scaling, replication, data locality, rack awareness)
- Nodes topologies (NameNode, Secondary NameNode, Standby NameNode, DataNode)
  o NameNode Considerations
  o Secondary vs Standby Namenode
- Hadoop Ecosystem

**HDFS & MapReduce**

- HDFS basics and high level architecture.
  o What is HDFS and it's basic architecture.
  o How is it different from other Unix OS?
  o FUSE (File system in User SpacE) and architecture
  o Configuration files, logs and directory structure
  o Reading and writing files
  o Copying data into HDFS
    ▪ Using distcp
    ▪ Using Sqoop from RDBMS
    ▪ Using Flume from external sources
  o Checking HDFS with fsck
  o Quotas and Trash
  o Cluster rebalancing
- MapReduce basics and high level architecture.
  o What is MapReduce?
  o Basic architecture, daemons & features
  o MapReduce Version 1 & 2 (YARN)
    ▪ Job tracker & task tracker
    ▪ Resource manager
    ▪ Node manager & application manager
  o Configuration and it's processing power
  o Scheduling and managing running jobs
    ▪ Failover revovery

- - - ▪ Web UI to monitor jobs
    - o Best practices to optimize MR jobs
    - o Scheduler types & configuring a Scheduler (FIFO/Fair Scheduler)

## Installation
- Pseudo cluster install using vmware
  - o Daemons running
  - o Web UI to check health of cluster
- Multi-node cluster install
  - o Installing a New Node
  - o Namenode formatting
  - o NameNode recovery options
  - o Configuring Hadoop – xml files
  - o Hadoop Ports and Web UI
  - o Compression codec
  - o Creating Users and using Quotas and Trash
  - o Log files

## Advanced Configuration, Maintenance & Monitoring
- Explicitly Including and Excluding Hosts/Nodes
- Copying Data Between Clusters
- NameNode Metadata Backup
- General System Monitoring
- Common Troubleshooting Issues
- Backup and Recovery
- Create Queues in Capacity Scheduler
- Snapshots (use cases - data backup, protection against user errors and disaster recovery)
- NFS Gateway access to HDFS
- Migration from Hadoop v1 to Hadoop v2

## High Availability & Load Balancing
- High Availability
  - o Using QJM, Zookeeper & other options
  - o Manual and auto failover
  - o Role of standby server
  - o Fencing options
  - o Split Brain Syndrome
- Federation
  - o What is Federation?
  - o Using Federation for load balancing
- Federation with HA

## Security
- Why & what is available?
- Unix accounts with standard permissions, Kerberos & LDAP.

- Securing a Hadoop Cluster with Kerberos.
- Securing a Hadoop Cluster with LDAP.

**Ecosystem Highlights**
- Hive
- Impala
- Pig
- HBase
- Sqoop
- Oozie
- Flume

**Hands-On Exercises**

1. Install a pseudo-distributed Hadoop Cluster
2. Using the Job Tracker UI to start and kill jobs
3. Install multi-node Hadoop Cluster
4. Importing data from MySQL or text file
5. Populating HDFS using Sqoop
6. Run MapReduce jobs
7. Using Fair Scheduler
8. Dead nodes and data replication
9. Adding and removing data nodes