

# Hadoop course content

<b>COURSE DETAILS</b>	
<ol style="list-style-type: none"> <li>1. In-detail explanation on the concepts of HDFS &amp; MapReduce frameworks</li> <li>2. What is Hadoop 2.X Architecture &amp; How to set up Hadoop Cluster</li> <li>3. How to write complex MapReduce Programs</li> <li>4. In-detail explanation on how to load data using tools like Sqoop &amp; Flume</li> <li>5. How to perform data analysis using tools like PIG, HIVE &amp; YARN</li> <li>6. How to implement &amp; integrate HBASE &amp; MapReduce</li> <li>7. How to execute Advanced Usage and Indexing</li> <li>8. How to schedule jobs using Oozie</li> <li>9. What are the best practices for overall Hadoop development</li> <li>10. RTAs on Data Analytics</li> <li>11. What is Spark &amp; brief about its ecosystem &amp; how to work on RDD Using Spark</li> </ol>	<p><b>Programming languages:</b> Java &amp; Scala</p> <p><b>Frame works:</b> Hadoop Distributed File System (HDFS) &amp; MapReduce, spark</p> <p><b>Loading Tools:</b> Sqoop &amp; Flume</p> <p><b>Analytical Tools:</b> Pig, Hive and YARN</p> <p><b>Scheduling Tools:</b> Oozie</p>

## CURRICULUM for HADOOP 2.X

S No	Concepts	Syllabus Objectives	Topics	RTAs
1	Understanding Big Data and Hadoop	<p>The syllabus for this lecture would brief about:</p> <ol style="list-style-type: none"> <li>1. Big Data</li> <li>2. Big Data problems &amp; solutions, their limitations</li> <li>3. HADOOP's solutions that handle Big Data issue</li> <li>4. Common Hadoop Ecosystem and its Architecture</li> <li>5. Introduction to HDFS</li> <li>6. What is a file and how to write &amp; read</li> <li>7. Brief on MapReduce</li> </ol>	<ol style="list-style-type: none"> <li>1. Big Data, Limitations and Solutions of existing Data Analytics Architecture,</li> <li>2. Hadoop,</li> <li>3. Hadoop Features,</li> <li>4. Hadoop Ecosystem,</li> <li>5. Hadoop 2.x core components,</li> <li>6. Hadoop Storage: HDFS,</li> <li>7. Hadoop Processing</li> <li>8. MapReduce Framework,</li> <li>9. Hadoop Different Distributions.</li> </ol>	

		Framework and its working style.		
2	Hadoop requirements	The syllabus for this lecture would brief about: Hadoop Requirements	<ol style="list-style-type: none"> <li>1. Machine learning algorithms(field of study that gives computer the ability to learn without being explicitly programmed) R, PYTHON Mlib, Weka</li> <li>2. Large scale distributed systems, MPP columnar stores(v Machine learning algorithms(field of study that gives computer the ability to learn without being explicitly programmed) R, PYTHON Mlib, Weka</li> <li>3. large scale distributed systems, MPP columnar stores (vertica), Graph and NOSQL</li> <li>4. UNIX/LINUX, Shell Scripts</li> <li>5. java program frameworks in Agile/scrum methodology</li> <li>6. SPLUNK (BI tool), TALEND (it's an open studio for big data)</li> <li>7. Agile software development practices</li> <li>8. cloud based testing</li> <li>9. Experience with writing test strategy/plan for IT projects and test status reporting</li> <li>10. Tableau, Datameer, Azure</li> <li>11. Serialization such as JSON(JavaScript Object Notation) and BSON(Binary Script Object Notation)</li> <li>14. Apache Mahout</li> <li>15. Sentiment analysis</li> <li>16. Spark streaming and SparkQL</li> <li>17. Data pipelines, Data API's</li> <li>18. Oracle/DBA2/SQL</li> <li>19. Basic ML tools(Machine learning) (weka/R/RapidMiner/SAS/SPSS)</li> <li>20. Perl</li> <li>21. Knowledge on Data cleaning and transformations</li> </ol>	

			<ul style="list-style-type: none"> <li>22. JSON transformations</li> <li>23. RStudio and R Programming language</li> <li>24. Informatica ETL products</li> <li>25. Javascript, HTML, CSS</li> <li>26. Mesos, Swarm</li> <li>27. TeraData</li> <li>28. ertica), Graph and NOSQL</li> <li>29. UNIX/LINUX, Shell Scripts</li> <li>30. java program frameworks in</li> <li>31. Agile/scrum methodology</li> <li>32. Splunk(BI tool), Talend(it's an open</li> <li>33. studio for big data)</li> <li>34. Agile software development practices</li> <li>35. cloud based testing</li> <li>36. Experience with writing test strategy/plan for IT projects and test status reporting</li> <li>37. Tableau, Datameer, Azure</li> <li>38. Serialization such as JSON(JavaScript Object Notation) and BSON(Binary Script Object Notation)</li> <li>39. Apache Mahout</li> <li>40. Sentiment analysis</li> <li>41. Spark streaming and SparkQ</li> </ul>	
3	Hadoop Architecture and HDFS	<p>The syllabus for this lecture would brief about:</p> <ul style="list-style-type: none"> <li>1. What is Hadoop Cluster Architecture</li> <li>2. What are the important Configuring files in a Hadoop Cluster</li> <li>3. What are the various Data loading techniques</li> <li>4. What are Single node and Multi nodes and their setups</li> </ul>	<ul style="list-style-type: none"> <li>1. Hadoop 2.x Cluster Architecture Federation and High Availability,</li> <li>2. A Typical Production Hadoop Cluster,</li> <li>3. Hadoop Cluster Modes,</li> <li>4. Common Hadoop Shell Commands,</li> <li>5. Hadoop 2.x Configuration Files,</li> <li>6. Single node cluster and Multi node cluster set up Hadoop Administration.</li> </ul>	
4	Hadoop MapReduce Framework	<p>The syllabus for this lecture would brief about:</p> <ul style="list-style-type: none"> <li>1. In-depth analysis on Hadoop MapReduce Framework</li> </ul>	<ul style="list-style-type: none"> <li>1. MapReduce Use Cases,</li> <li>2. Traditional way Vs MapReduce way,</li> <li>3. Why MapReduce,</li> <li>4. Hadoop 2.x MapReduce Architecture,</li> <li>5. Hadoop 2.x MapReduce Components,</li> </ul>	

	k	<ol style="list-style-type: none"> <li>2. How MapReduce works on data stored in HDFS.</li> <li>3. What are Splits, Combiner &amp; Partitioner.</li> <li>4. How to work on MapReduce using different data sets</li> </ol>	<ol style="list-style-type: none"> <li>6. YARN MR Application Execution Flow,</li> <li>7. YARN Workflow,</li> <li>8. Anatomy of MapReduce Program,</li> <li>9. Demo on MapReduce.</li> <li>10. Input Splits,</li> <li>11. Relation between Input Splits and HDFS Blocks,</li> <li>12. MapReduce Combiner &amp; Partitioner,</li> <li>13. Demo on de-identifying Health Care Data set,</li> <li>14. Demo on Weather Data set.</li> </ol>	
5	Advanced MapReduce	<p>The syllabus for this lecture would brief about:</p> <ol style="list-style-type: none"> <li>1. Advanced concepts in MapReduce</li> <li>2. What are Counters, Distributed Cache, MRUNIR, Reduce Join, Custom Input format &amp; Sequence Input Format</li> <li>3. What is XML Parsing</li> </ol>	<ol style="list-style-type: none"> <li>1.Counters,</li> <li>2.Distributed Cache,</li> <li>3.MRunit,</li> <li>4.Reduce Join,</li> <li>5.Custom Input Format,</li> <li>6.Sequence Input Format,</li> <li>7.Xml file Parsing using MapReduce.</li> </ol>	
6	Pig	<p>The syllabus for this lecture would brief about:</p> <ol style="list-style-type: none"> <li>1. What is PIG &amp; types of use, demo case</li> <li>2. How to couple PIG with MapReduce</li> <li>3. What are PIG Latin Scripting</li> <li>4. What are PIG running Modes PIG UDF, Pig Streaming, Testing PIG Scripts.</li> </ol>	<ol style="list-style-type: none"> <li>1. About Pig,</li> <li>2. MapReduce Vs Pig,</li> <li>3. Pig Use Cases,</li> <li>4. Programming Structure in Pig,</li> <li>5. Pig Running Modes,</li> <li>6. Pig components,</li> <li>7. Pig Execution,</li> <li>8. Pig Latin Program,</li> <li>9. Data Models in Pig,</li> <li>10. Pig Data Types,</li> <li>11. Shell and Utility Commands,</li> <li>12. Pig Latin Relational Operators,</li> <li>13. File Loaders,</li> <li>14. Group Operator,</li> <li>15. COGROUP Operator,</li> <li>16. Joins and COGROUP,</li> <li>17. Union,</li> <li>18. Diagnostic Operators,</li> <li>19. Specialized joins in Pig,</li> </ol>	

			<ul style="list-style-type: none"> <li>20. Built In Functions (Eval Function, Load and Store Functions, Math function, String Function, Date Function, Pig UDF, Piggybank),</li> <li>21. Parameter Substitution ( PIG macros and Pig Parameter</li> <li>22. substitution ),</li> <li>23. Pig Streaming,</li> <li>24. Testing Pig scripts with Punit,</li> <li>25. Aviation use case in PIG,</li> <li>26. Pig Demo on Healthcare Data set.</li> </ul>	
7	Hive	<p>The syllabus for this lecture would brief about:</p> <ul style="list-style-type: none"> <li>1. What are HIVE concepts</li> <li>2. What are HIVE data types</li> <li>3. What are Loading &amp; Querying in HIVE,</li> <li>4. How to run HIVE scripts</li> <li>5. What are Hive UDF</li> </ul>	<ul style="list-style-type: none"> <li>1. Hive Background,</li> <li>2. Hive Use Case,</li> <li>3. About Hive,</li> <li>4. Hive Vs Pig,</li> <li>5. Hive Architecture and Components,</li> <li>6. Metastore in Hive,</li> <li>7. Limitations of Hive,</li> <li>8. Comparison with Traditional Database,</li> <li>9. Hive Data Types and Data Models,</li> <li>10. Partitions and Buckets,</li> <li>11. Hive Tables(Managed Tables and External Tables),</li> <li>12. Importing Data,</li> <li>13. Querying Data,</li> <li>14. Managing Outputs,</li> <li>15. Hive Script,</li> <li>16. Hive UDF,</li> <li>17. Retail use case in Hive,</li> <li>18. Hive Demo on Healthcare Data set.</li> </ul>	
8	Advanced Hive and HBase	<p>The syllabus for this lecture would brief about:</p> <ul style="list-style-type: none"> <li>1. What are Advanced HIVE concepts</li> <li>2. What are UDF, Dynamic Partitioning, HIVE indexes &amp; Views</li> <li>3. What are Optimizations in HIVE</li> <li>4. In-depth analysis on HBase,</li> </ul>	<ul style="list-style-type: none"> <li>1. Hive QL: Joining Tables,</li> <li>2. Dynamic Partitioning,</li> <li>3. Custom Map/Reduce Scripts,</li> <li>4. Hive Indexes and views</li> <li>5. Hive query optimizers,</li> <li>6. Hive : Thrift Server,</li> <li>7. User Defined Functions,</li> <li>8. HBase:</li> <li>9. Introduction to NoSQL</li> <li>10. Databases and HBase,</li> </ul>	

		its Architecture, components and its running modes	11. HBase v/s RDBMS, 12. HBase Components, 13. HBase Architecture, 14. Run Modes & Configuration, 15. HBase Cluster Deployment.	
9	Advanced HBase	The syllabus for this lecture would brief about: 1. What are Advanced HBase Concepts 2. How to perform bulk loading 3. What are filters 4. What is Zookeeper and how it helps in Cluster monitoring. 5. Why HBase utilizes Zookeeper	1. HBase Data Model, 2. HBase Shell, 3. HBase Client API, 4. Data Loading Techniques, 5. ZooKeeper 6. Data Model, 7. Zookeeper Service, 8. Zookeeper, 9. Demos on Bulk Loading, 10. Getting and Inserting Data, 11. Filters in HBase.	
10	Processing Distributed Data with Apache Spark	The syllabus for this lecture would brief about: 1. What is Spark Ecosystem 2. What is Scala and its utility in Spark 3. What is SparkContext 4. How to work on RDD in Spark 5. How to run a Spark Cluster 6. Comparison of MapReduce vs Spark	1. What is Apache Spark, 2. Spark Ecosystem, 3. Spark Components, 4. History of Spark 5. Spark Versions/Releases, 6. Spark a Polyglot, 7. What is Scala?, 8. Why Scala?, 9. SparkContext, 10. RDD.	
11	Oozie	The syllabus for this lecture would brief about: 1. How multiple Hadoop ecosystem components work 2. How they should be implemented to solve Big Data Issues	1. Flume and Sqoop Demo, 2. Oozie, 3. Oozie Components, 4. Oozie Workflow, 5. Scheduling with Oozie, 6. Demo on Oozie Workflow, 7. Oozie Co-ordinator, 8. Oozie Commands, 9. Oozie Web Console, 10. Oozie for MapReduce, 11. PIG, Hive, and Sqoop, 12. Combine flow of MR, PIG, Hive in	

			Oozie	
12	Hadoop Project	<p>The syllabus for this lecture would brief about:</p> <ol style="list-style-type: none"> <li>1. What are Flume &amp; Scoop</li> <li>2. How to utilize Apache Oozie Workflow Scheduler for Hadoop Jobs</li> <li>3. What is Hadoop Tableau Integration</li> </ol>	1.Hadoop Integration with Tableau	