

12 APR 2024



22MCA2052



USN 1 B Y

# BMS INSTITUTE OF TECHNOLOGY AND MANAGEMENT

(An Autonomous Institute affiliated to Visvesvaraya Technological University, Belagavi)

## SEMESTER END EXAMINATION QUESTION PAPER

### Second Semester MCA Degree Examination

Regular / Make-up / Arrears / Supplementary

### BIG DATA ANALYTICS

Time: 3 hrs.

Max. Marks: 100

Note: 1. Answer FIVE full questions, choosing ONE full question from each module.

Q. No	Module - 1	Marks	CO & RBT
1a.	Present an overview of Big Data? Explain Analytics Process model.	10	CO1, K2
b.	Big Data is characterised by various Vs. Describe 4 Vs of Big Data.	5	CO1, K2
c.	Calculate Z scores and detect the outlier for the following data. $\mu=40$ , $SD=10$ . Data = 30, 50, 10, 40, 60, 80.	5	CO1, K2
OR			
2a.	Discuss the applications of Big Data Analytics in marketing, risk management, government, web and logistics.	10	CO1, K2
b.	Explain the requirements to satisfy a good analytical model.	5	CO1, K2
c.	Construct Box plot for given data: 51,17,25,39,7,49,67,41,20,2,43,13.	5	CO1, K2
Module – 2			
3a.	Describe the scenarios in which predictive and descriptive analytics can be used.	10	CO2, K2
b.	Explain association rules as applied in descriptive analytics.	5	CO2, K2
c.	Bring out the relevance of decision trees in data analytics with an example.	5	CO2, K2
OR			
4a.	Elaborate on any 2 techniques used for classification under predictive analytics.	10	CO2, K2
b.	Describe k-means clustering with example.	5	CO2, K2
c.	Compare and contrast bagging and boosting.	5	CO2, K2
Module – 3			
5a.	Explain HDFS architecture with a neat diagram.	10	CO3, K2
b.	Explain with a diagram, MapReduce data flow with a single reduce task and multiple reduce task.	10	CO3, K2
OR			

6a.	With a neat diagram, explain the anatomy of reading data from a file in HDFS.	10	CO3, K2
b.	What is Map Reduce? Sketch a neat diagram and explain the logical data flow in Map Reduce?	10	CO3, K2
<b>Module – 4</b>			
7a.	What is the role of Apache Spark? With suitable diagram illustrate the architecture of a Spark Application.	10	CO4, K2
b.	Describe Spark's MLlib.	5	CO4, K2
c.	Discuss the difference between a single machine and dataframe data distributed over a cluster.	5	CO4, K2
<b>OR</b>			
8a.	Describe Structured API logical planning process and physical planning process.	10	CO4, K2
b.	Discuss datasets of Spark.	5	CO4, K2
c.	What is Lazy Evaluation? Compare and contrast between Transformations and Actions.	5	CO4, K2
<b>Module – 5</b>			
9a.	Illustrate the Hive modules in the Hadoop Ecosystem.	10	CO5, K2
b.	Discuss how data is inserted/loaded into Hive tables.	5	CO5, K2
c.	Describe any 3 key data types supported by Hive.	5	CO5, K2
<b>OR</b>			
10a.	Illustrate application of Aggregate functions, Mathematical functions and Table generating functions.	10	CO5, K2
b.	Discuss how dynamic partition inserts are accomplished in hive.	5	CO5, K2
c.	Discuss how internal and External Tables are created in Hive.	5	CO5, K2

**Course Outcomes (COs): At the end of the course, the student will be able to**

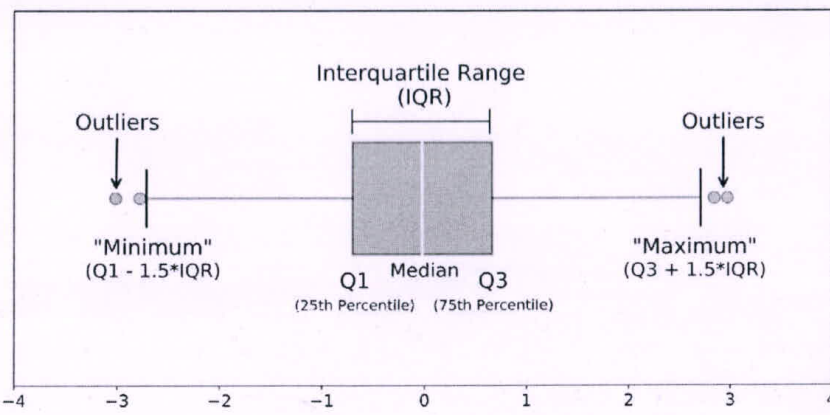
COs	Statements
CO-1	Identify the business problem for a given context and frame the objectives to solve it using data analytics tools.
CO-2	Differentiate various types of analytics algorithms and context of their application.
CO-3	Illustrate the architecture of HDFS and MapReduce.
CO-4	Explore Spark architecture and its language APIs.
CO-5	Write Hive queries against large datasets on clusters.
K1- Remembering   K2 - Understanding   K3 - Applying   K4 - Analyzing   K5 - Evaluating   K6 - Creating	

*"Success is the progressive realization of a worthy goal."*

\*\*\*\*\*



Module - 1																													
Q.No.	Questions	Marks																											
1a.	<p><b>Overview of Big Data – 2 M</b></p> <p><b>Analytical Process Model Diagram – 2M</b></p> <p><b>Explanation of Analytical Process Model – 6 M</b></p> <p><b>Big Data</b> is a collection of data that is huge in volume, yet growing exponentially with time. It is a data with so large size and complexity that none of traditional data management tools can store it or process it efficiently.</p> <p>Big Data refers to the large, diverse sets of information that grow at ever-increasing rates. It encompasses the volume of information, the velocity or speed at which it is created and collected, and the variety or scope of the data points being covered.</p> <p style="text-align: center;"><b>Overview of the Analytics Process Model</b></p> <div style="text-align: center;"> <pre> graph LR     A[Identify Business Problem] --&gt; B[Identify Data Sources]     B --&gt; C[Select the Data]     C --&gt; D[Clean the Data]     D --&gt; E[Transform the Data]     E --&gt; F[Analyze the Data]     F --&gt; G[Interpret, Evaluate, and Deploy the Model]     subgraph Preprocessing         A         B         C         D         E     end     subgraph PostAnalyticsProcessing [Post-Analytics processing]         F         G     end </pre> </div> <p><b>Steps involved in Analytics Process Model:</b></p> <ol style="list-style-type: none"> <li>1. Define the business problems to be solved</li> <li>2. All source-data need to be identified</li> <li>3. All data to be gathered as a pool</li> <li>4. Data cleaning</li> <li>5. Analytical model estimation</li> <li>6. Model interpretation and evaluation</li> </ol>	10																											
1b.	<p><b>Mention of 4 to 6 Vs – 1 M - Volume, Velocity, Value, Variability, Veracity, Variety</b></p> <p><b>Explanation of any 4 Vs – 4 M</b></p>	5																											
1c.	<p>Calculation of Z score and detection of the outlier for the following data.  <b>Mu=40, SD=10. Data = 30, 50, 10, 40, 60, 80.</b></p> <p>Z Score with +4 is an outlier in the given data.</p> <p>That means the value corresponding to +4 is 80 is an outlier.</p> <p><b>Calculation – 3 M</b></p> <p><b>Final Discussion – 2M</b></p> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th>ID</th><th>Age</th><th>Z-Score</th></tr> </thead> <tbody> <tr> <td>1</td><td>30</td><td><math>(30 - 40)/10 = -1</math></td></tr> <tr> <td>2</td><td>50</td><td><math>(50 - 40)/10 = +1</math></td></tr> <tr> <td>3</td><td>10</td><td><math>(10 - 40)/10 = -3</math></td></tr> <tr> <td>4</td><td>40</td><td><math>(40 - 40)/10 = 0</math></td></tr> <tr> <td>5</td><td>60</td><td><math>(60 - 40)/10 = +2</math></td></tr> <tr> <td>6</td><td>80</td><td><math>(80 - 40)/10 = +4</math></td></tr> <tr> <td>...</td><td>...</td><td>...</td></tr> <tr> <td></td><td><math>\mu = 40</math> <math>\sigma = 10</math></td><td><math>\mu = 0</math> <math>\sigma = 1</math></td></tr> </tbody> </table>	ID	Age	Z-Score	1	30	$(30 - 40)/10 = -1$	2	50	$(50 - 40)/10 = +1$	3	10	$(10 - 40)/10 = -3$	4	40	$(40 - 40)/10 = 0$	5	60	$(60 - 40)/10 = +2$	6	80	$(80 - 40)/10 = +4$	...	...	...		$\mu = 40$ $\sigma = 10$	$\mu = 0$ $\sigma = 1$	5
ID	Age	Z-Score																											
1	30	$(30 - 40)/10 = -1$																											
2	50	$(50 - 40)/10 = +1$																											
3	10	$(10 - 40)/10 = -3$																											
4	40	$(40 - 40)/10 = 0$																											
5	60	$(60 - 40)/10 = +2$																											
6	80	$(80 - 40)/10 = +4$																											
...	...	...																											
	$\mu = 40$ $\sigma = 10$	$\mu = 0$ $\sigma = 1$																											

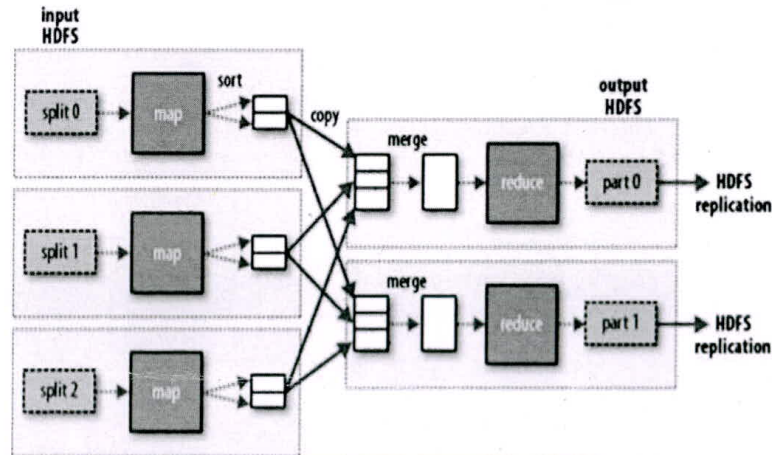
2a.	<p><b>Big Data Applications – Any Five domains - 2 M Each</b></p> <p><b>Marketing:</b> Response Modelling, Net Lift Modelling, Retention Modelling, Market Basket Analysis, Recommender Systems, Customer Segmentation</p> <p><b>Risk management:</b> Credit Risk Modelling, Market Risk Modelling, Operational Risk Modelling, Fraud Detection.</p> <p><b>Government:</b> Tax Avoidance, Social Security Fraud, Money Laundering, Terrorism Detection</p> <p><b>Web:</b> Web Analytics, Social Media Analytics, Multivariate Testing.</p> <p><b>Logistics:</b> Demand Forecasting, Supply Chain Analytics.</p> <p><b>Others:</b> Text Analytics, Business Process Analytics, Sentiment Analytics</p>	10
2b.	<p><b>Requirements to satisfy a good analytical model – Any 5 – 1 M Each</b></p> <p>Business relevance</p> <p>Statistical performance</p> <p>Interpretability</p> <p>Justifiability</p> <p>Operationally efficient</p> <p>Economical</p> <p>Local and international regulation and legislation</p>	5
2c.	<p><b>Box plot for given data: 51,17,25,39,7,49,67,41,20,2,43,13</b></p> <p><b>Box Plot Construction – 3 M</b></p> <p><b>Q1, Q2, Q3, IQR Values computation – 2 M</b></p> 	5



3a.	<b>Scenarios for Predictive Analytics Application. – 2 M Each – Any Five</b>  <b>Any 5 use cases / scenarios like healthcare, marketing, supply chain, logistics, web analytics etc.</b>  <hr/>	10
3b.	<b>Association Rules: 5 M</b>  <ul style="list-style-type: none"> <li>• Association rules are used to find correlations and co-occurrences between data sets.</li> <li>• They are ideally used to explain patterns in data from seemingly independent information repositories, such as relational databases and transactional databases.</li> <li>• It is employed in Market Basket analysis, Web usage mining etc.</li> <li>• Association rules typically start from a database of transactions.</li> <li>• Each transaction consists of a transaction identifier and a set of items.</li> <li>• Association rules are usually represented in the form <math>X \rightarrow Y</math>, where X (also called rule Antecedent) and Y (also called Consequent) are disjoint item sets (i.e., disjoint conjunctions of features).</li> <li>• Ex: If a customer buys bread, he's 70% likely of buying milk." In the above association rule, bread is the antecedent and milk is the consequent.</li> <li>• Association rule shows how frequently an item set occurs in a transaction.</li> <li>• Association rules are created by analysing data for frequent if/then patterns and using the criteria support and confidence to identify the most important relationships.</li> </ul> <hr/>	5
3c.	<b>Decision Trees – Discussion – 3 M</b>  <b>Example Diagram – 2 M</b>  <hr/>	5
4a.	<b>Any 2 techniques from Logistic Regression, Decision Trees, Neural Networks, SVM, Ensemble Methods (Bagging, Boosting, Random Forest).</b>  <b>Each 5 M with Explanation and Example Diagram</b>  <hr/>	10
4b.	<b>k-means Clustering</b>  <b>Discussion – 3 M</b>  <b>Example – 2 M</b>  <hr/>	5

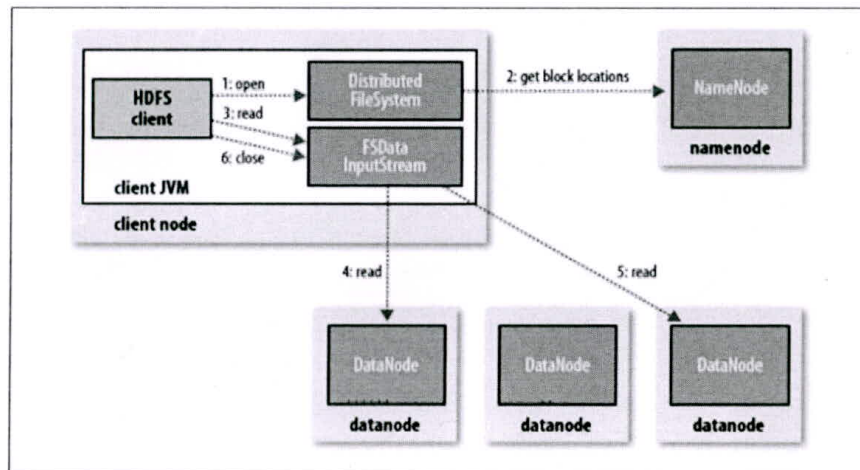
4c.	Bagging vs Boosting – Any 5 Parameters – 1 M Each	5																		
<table><tr><th>Bagging</th><th>Boosting</th></tr><tr><td>The simplest way of combining predictions that belong to the same type.</td><td>A way of combining predictions that belong to the different types.</td></tr><tr><td>Each model receives equal weight.</td><td>Models are weighted according to their performance.</td></tr><tr><td>Each model is built independently.</td><td>New models are influenced by the performance of previously built models.</td></tr><tr><td>Different training data subsets are selected using row sampling with replacement and random sampling methods from the entire training dataset.</td><td>Every new subset contains the elements that were misclassified by previous models.</td></tr><tr><td>If the classifier is unstable (high variance), then apply bagging.</td><td>If the classifier is stable and simple (high bias) the apply boosting.</td></tr><tr><td>In bagging base classifiers are trained parallelly.</td><td>In boosting base classifiers are trained sequentially.</td></tr><tr><td>Aims at decreasing variance (Overfitting leads to high variance – Overfitting means that the model performs well on the training data but does not perform accurately in the evaluation set.)</td><td>Aims at decreasing bias (Bias is a phenomenon that skews the result of an algorithm either in favour or against an idea).</td></tr><tr><td><b>Example:</b> The <b>Random forest</b> model uses <b>Bagging</b>.</td><td><b>Example:</b> The <b>AdaBoost</b> algorithm uses <b>Boosting</b>.</td></tr></table>			Bagging	Boosting	The simplest way of combining predictions that belong to the same type.	A way of combining predictions that belong to the different types.	Each model receives equal weight.	Models are weighted according to their performance.	Each model is built independently.	New models are influenced by the performance of previously built models.	Different training data subsets are selected using row sampling with replacement and random sampling methods from the entire training dataset.	Every new subset contains the elements that were misclassified by previous models.	If the classifier is unstable (high variance), then apply bagging.	If the classifier is stable and simple (high bias) the apply boosting.	In bagging base classifiers are trained parallelly.	In boosting base classifiers are trained sequentially.	Aims at decreasing variance (Overfitting leads to high variance – Overfitting means that the model performs well on the training data but does not perform accurately in the evaluation set.)	Aims at decreasing bias (Bias is a phenomenon that skews the result of an algorithm either in favour or against an idea).	<b>Example:</b> The <b>Random forest</b> model uses <b>Bagging</b> .	<b>Example:</b> The <b>AdaBoost</b> algorithm uses <b>Boosting</b> .
Bagging	Boosting																			
The simplest way of combining predictions that belong to the same type.	A way of combining predictions that belong to the different types.																			
Each model receives equal weight.	Models are weighted according to their performance.																			
Each model is built independently.	New models are influenced by the performance of previously built models.																			
Different training data subsets are selected using row sampling with replacement and random sampling methods from the entire training dataset.	Every new subset contains the elements that were misclassified by previous models.																			
If the classifier is unstable (high variance), then apply bagging.	If the classifier is stable and simple (high bias) the apply boosting.																			
In bagging base classifiers are trained parallelly.	In boosting base classifiers are trained sequentially.																			
Aims at decreasing variance (Overfitting leads to high variance – Overfitting means that the model performs well on the training data but does not perform accurately in the evaluation set.)	Aims at decreasing bias (Bias is a phenomenon that skews the result of an algorithm either in favour or against an idea).																			
<b>Example:</b> The <b>Random forest</b> model uses <b>Bagging</b> .	<b>Example:</b> The <b>AdaBoost</b> algorithm uses <b>Boosting</b> .																			
5a.	<b>Hadoop Architecture – Diagram – 4 M</b> <b>Discussion on HDFS, YARN, MapReduce – 2 M each</b>	10																		
5b.	MapReduce processing involves multiple mappers in general. However, when it comes to reducer, it may have one in general or increase it carefully as the need be but not in more numbers. Elaborate the MapReduce data flow with a no reduce task, single reduce task and multiple reduce tasks.  Single Reduce Task – Diagram, Discussion – 5M Multiple Reduce Tasks – Diagram, Discussion – 5 M	10																		
<pre>graph LR     subgraph Input_HDFS [input HDFS]         S0[split 0]         S1[split 1]         S2[split 2]     end     S0 --&gt; M0[map]     S1 --&gt; M1[map]     S2 --&gt; M2[map]     M0 -- copy --&gt; Merge[merge]     M1 -- copy --&gt; Merge     M2 -- copy --&gt; Merge     Merge --&gt; R[reduce]     R --&gt; P0[part 0]     P0 --&gt; Out_HDFS[output HDFS]     Out_HDFS --&gt; Rep[replication]</pre>																				





- 6a. Anatomy of reading Data from a file in HDFS.  
 Diagram – 4 M  
 Explanation of Steps involved – 6 M

10



- 6b. MapReduce Logical DataFlow  
 Diagram – 4 M  
 Explanation with Example – 6 M

10

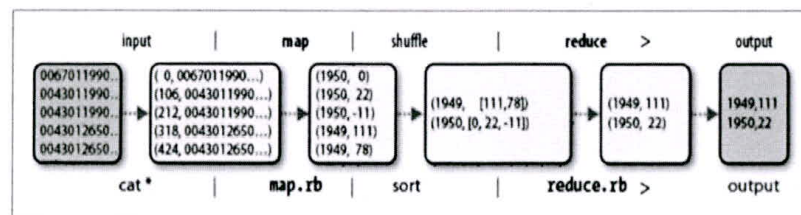
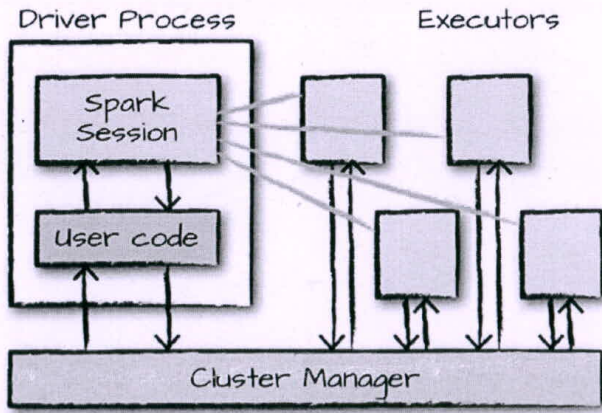
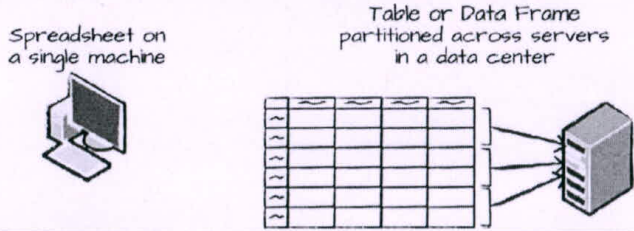
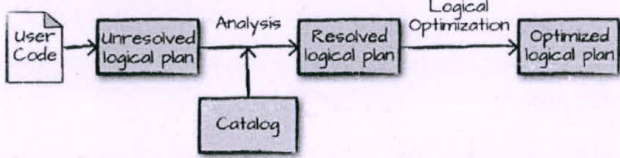
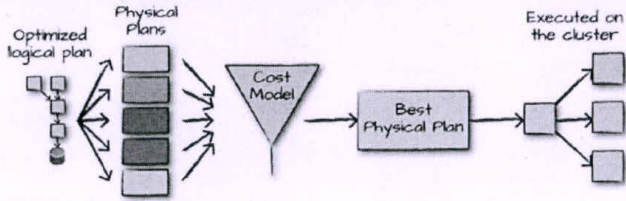
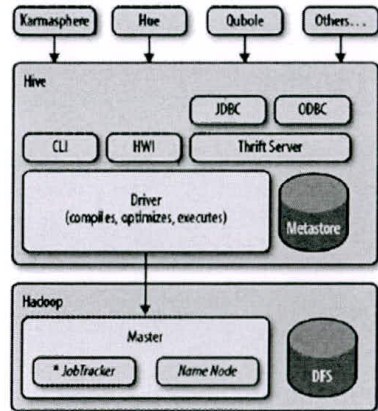


Figure 2-1. MapReduce logical data flow

7a.	<p><b>About Apache Spark – 2 M</b>  <b>Architecture of Spark – Diagram – 3 M</b>  <b>Explanation of Spark's Architecture – 5 M</b></p>  <p>The diagram illustrates the Spark architecture. On the left, a box labeled 'Driver Process' contains a 'Spark Session' box and a 'User code' box, connected by bidirectional arrows. On the right, there are four boxes representing 'Executors'. At the bottom is a 'Cluster Manager' box. Arrows show the 'Spark Session' connecting to all 'Executors'. The 'Cluster Manager' connects to the 'Driver Process' and each 'Executor'.</p>	10
7b.	<p><b>Note on Spark MLib – 5 M</b></p>	5
7c.	<p><b>Difference between a single machine and dataframe data distributed over a cluster. Diagram – 2 M</b>  <b>Discussion – 3 M</b></p>  <p>The diagram compares two data storage methods. On the left, 'Spreadsheet on a single machine' is shown with a computer icon. On the right, 'Table or Data Frame partitioned across servers in a data center' is shown with a grid of data blocks connected to a server rack icon.</p>	5
8a.	<p><b>Structured API Logical Planning Process – Diagram, Explanation - 5 M</b></p>  <p>The flowchart shows the logical planning process: 'User Code' leads to 'Unresolved logical plan', which goes through 'Analysis' (involving a 'Catalog') to become a 'Resolved logical plan'. This then undergoes 'Logical Optimization' to become an 'Optimized logical plan'.</p> <p><b>Structured API Physical Planning Process – Diagram, Explanation - 5 M</b></p>  <p>The flowchart shows the physical planning process: an 'Optimized logical plan' is converted into 'Physical Plans' (represented by a stack of blocks). These are evaluated by a 'Cost Model' (funnel shape) to select the 'Best Physical Plan', which is then 'Executed on the cluster'.</p>	10



8b.	<b>Note on Datasets – Type-safe Structured APIs</b>	5
8c.	<p><b>Lazy Evaluation – 1 M</b></p> <p>Lazy evaluation means that Spark will wait until the very last moment to execute the graph of computation instructions. In Spark, instead of modifying the data immediately when you express some operation, you build up a <i>plan</i> of transformations that you would like to apply to your source data. By waiting until the last minute to execute the code, Spark compiles this plan from your raw DataFrame transformations to a streamlined physical plan that will run as efficiently as possible across the cluster.</p> <p><b>Difference between TRANSFORMATIONS and ACTIONS – 2 M Each</b></p> <p>Transformations allow us to build up our logical transformation plan. To trigger the computation, we run an action. An action instructs Spark to compute a result from a series of transformations.</p>	5
9a.	<p><b>Hive Modules in Hadoop Ecosystem</b></p> <p><b>Diagram – 4 M</b> <b>Explanation – 6M</b></p> 	10
9b.	<p><b>Usage of CASE...WHEN...THEN Statements.</b></p> <p>The CASE ... WHEN ... THEN clauses are like if statements for individual columns in query results.</p> <p><b>Explanation – 2 M</b> <b>Example – 3 M</b></p> <pre> hive&gt; SELECT name, salary, &gt; CASE &gt;   WHEN salary &lt; 50000.0 THEN 'low' &gt;   WHEN salary &gt;= 50000.0 AND salary &lt; 70000.0 THEN 'middle' &gt;   WHEN salary &gt;= 70000.0 AND salary &lt; 100000.0 THEN 'high' &gt;   ELSE 'very high' &gt; END AS bracket FROM employees; John Doe      100000.0    very high Mary Smith    80000.0     high Todd Jones    70000.0     high Bill King     60000.0     middle Boss Man      200000.0    very high Fred Finance  150000.0    very high Stacy Accountant 60000.0     middle ... </pre> <p><i>Data inserted / loaded in 2 ways</i></p> <p><i>272 M each</i></p>	5

9c.	<p><b>Partitioned-Managed Tables of Hive.</b></p> <p><b>Explanation – 3 M</b></p> <p><b>Example – 2 M</b></p>	<p>3 key date types</p> <p>Month – 1M</p> <p>Expt any 2 – 4M.</p>	5
10a.	<p><b>Aggregate functions, Mathematical functions and Table generating functions.</b></p> <p><b>Write up on these functions - 1 M</b></p> <p><b>Discussion of any three from each group – 3 M Each</b></p>		10
10b.	<p><b>Dynamic Partition in Hive - Hive supports a dynamic partition feature, where it can infer the partitions to create based on query parameters.</b></p> <p><b>Explanation – 3 M</b></p> <p><b>Example – 2 M</b></p>		5
10c.	<p><b>Internal and External Tables are created in Hive.</b></p> <p><b>Internal Tables/Managed Tables – Explanation, Example – 2 ½ M</b></p> <p><b>External Tables – Explanation, Example – 2 ½ M</b></p> <p>*****</p>		5