

MID-SEMESTER EXAMINATION, FEB, 2024

Course Code- CDCSC11

Course Title- **Big Data Analytics**

Time- 1.5 Hours

Max. Marks- 15

Note: - Attempt all questions. Missing data/ information (if any), may be suitably assumed & mentioned in the answer.

Q. No.	Question	Marks	CO
✓ 1a	Discuss the integration of cloud computing with big data and its advantages.	1.5	2
✓ 1b	What are the challenges associated with inter and trans firewall analytics?	1.5	
✓ 2a	How does Hadoop enable parallel processing, and what is the significance of data locality?	2	1
✓ 2b	Write a pseudocode for implementing word-count program using Map-reduce	1	
✓ 3a	Provide a detailed comparison of different types of NoSQL databases	2	2
✓ 3b	Provide examples of scenarios where materialized views in NoSQL databases prove beneficial.	1	
✓ 4a	How can you insert, update, and delete documents in a MongoDB collection?	2	3
✓ 4b	How do <u>aggregate data models</u> address the challenges of handling large datasets in NoSQL databases?	1	
✓ 5a	Explain HDFS architecture in detail. How HDFS ensures fault tolerance?	2	2
✓ 5b	How does the MapReduce programming model enable parallel processing?	1	

Course Code- CDCSC11

Course Title- Big Data Analytics

Time- 3 Hours

Max. Marks- 40

Note: - Attempt all questions in the given order only. Missing data/information (if any), maybe suitably assumed & mentioned in the answer.

Q. No.	Question	Marks	CO
Q1	Attempt any two parts of the following		
1a	What drives organizations to transition from traditional data warehouse tools to smarter data hubs rooted in Hadoop ecosystems? Furthermore, what differentiates the regular FileSystem from HDFS?	4	
1b	What is Crowdsourcing and how it helps in big data applications. Explain the role of big data analytics and crowdsourcing in following applications: <ul style="list-style-type: none"> • Credit Card Risk prediction • Healthcare and Medicine 	4	
1c	What are some key trends contributing to the generation and management of unstructured data? How can big data analytics be applied in fraud detection and prevention across industries such as banking and insurance?	4	
Q2	Attempt any two parts of the following		
2a	<u>Why Hadoop</u> Explain various Hadoop daemons and their roles in a Hadoop cluster. Write Hadoop commands for following: <ul style="list-style-type: none"> • Creating a file in HDFS • Copying a file from local file system to HDFS 	4	
2b	Does MongoDB support ACID transaction management and locking functionalities? How do I perform the SQL JOIN equivalent in MongoDB?	4	
2c	Where are the two types of metadata that NameNode server stores? If a file of <u>1 GB</u> has to be written on HDFS, Explain with diagram, how it will be stored in datanodes.	4	
Q3	Attempt any two parts of the following		
3a	How do you select among the Different File Formats for Storing and Processing Data? What is Avro Serialization in Hadoop?	4	
3b	State the reason why we can't perform "aggregation" (addition) in mapper? Why do we need the "reducer" for this?	4	

	You have a file that contains 200 billion URLs. How will you find the first unique URL using Hadoop MapReduce?																				
3c	Explain the architecture of YARN and how it allocates various resources to applications?	4																			
Q4	Attempt any two parts of the following																				
4a	What is the role of MasterServer and Zookeeper in HBase. What is lazy evaluation in Spark?	4																			
4b	What is the difference between logical and physical plans? Explain COGROUP in Pig with example?	4																			
4c	Why does Hive not store metadata information in HDFS? Write command to create a table student in Hive and create partitions based on semester and Elective: Student(Name, Roll No., Elective, semester, marks)	4																			
Q5	Attempt any two parts of the following																				
5a	Define Multivariate Analysis and how is it used in big data analytics. Fit a simple linear regression model to predict the price of a house based on its size.	4																			
	<table border="1"> <tr> <td>Size (in sq. feet)</td> <td>1000</td> <td>1500</td> <td>2000</td> <td>2500</td> <td>3000</td> </tr> <tr> <td>Price (in \$1000s)</td> <td>300</td> <td>450</td> <td>600</td> <td>750</td> <td>900</td> </tr> </table>	Size (in sq. feet)	1000	1500	2000	2500	3000	Price (in \$1000s)	300	450	600	750	900								
Size (in sq. feet)	1000	1500	2000	2500	3000																
Price (in \$1000s)	300	450	600	750	900																
	Using the regression model, predict the price of a house with a size of 1800 sq. feet.																				
5b	Explain how a weather prediction problem can be solved using fuzzy logic. Given $U = \{1,2,3,4,5,6,7\}$ $A = \{(3,0.7), (5,1), (6,0.8)\}$ What will be the value of $\sim A$ (where $\sim \rightarrow$ complement)	4																			
	Explain Bayesian Networks with example. Below, given the frequency table, what is the probability, that <u>players will play</u> if it is sunny weather.	4																			
	<table border="1"> <tr> <th colspan="3">Frequency Table</th> </tr> <tr> <th>Weather</th> <th>No</th> <th>Yes</th> </tr> <tr> <td>Overcast</td> <td>1</td> <td>4</td> </tr> <tr> <td>Rainy</td> <td>3</td> <td>2</td> </tr> <tr> <td>Sunny</td> <td>2</td> <td>3</td> </tr> <tr> <td>Total</td> <td>5</td> <td>9</td> </tr> </table>	Frequency Table			Weather	No	Yes	Overcast	1	4	Rainy	3	2	Sunny	2	3	Total	5	9		
Frequency Table																					
Weather	No	Yes																			
Overcast	1	4																			
Rainy	3	2																			
Sunny	2	3																			
Total	5	9																			

6 9

5
5
5