| Name : | | |
|---|---|---|
| Enrolment No. : | | |

UNIVERSITY OF PETROLEUM AND ENERGY STUDIES
End Semester Examination, Dec 2024

Program Name : B.Tech CSE (All specializations)  
Course Name  : Big Data Analysis  
Course Code  : CSBD4006P  
No. of Page(s) : 2  
Instructions  : Attempt all sections.

Semester : VII  
Time : 3 hours  
Max. Marks : 100

## SECTION-A

| S. No. | Questions | Marks | CO |
|---|---|---|---|
| Q.1 | Explain the components of YARN in brief. | 4 | CO1 |
| Q.2 | In the sentence, "Amazon announced a new product in Seattle," explain how Named Entity Recognition would identify whether "Amazon" refers to the company or the river. | 4 | CO2 |
| Q.3 | Explain how does the document-based database differs from relational Database. | 4 | CO3 |
| Q.4 | Identify the learning techniques used in the following tasks.<br>a) Face Recognition<br>b) Image Segmentation<br>c) Weather Forecast<br>d) Social Network Analysis | 4 | CO4 |
| Q.5 | List the data reading procedure in Hadoop in step-by-step manner. | 4 | CO1 |

## SECTION-B

| | | | Marks | CO |
|---|---|---|---|---|
| Q.6 | A) | Consider a rapidly growing e-commerce platform that handles increasing volumes of user data and transaction requests. Explain why vertical scalability would be a viable approach if improving the performance of the existing system without adding new nodes is crucial for maintaining smooth operations. | 10 (6+4) | CO1 |
| | B) | State any two applications of horizontal scaling. | | |
| Q.7 | | You are working with a large organization that has recently adopted a Data Lake for storing massive volumes of structured and unstructured data. However, you encounter two significant issues during data management:<br>a) The small file problem and<br>b) The shared file problem.<br>Discuss how these challenges impact data storage and retrieval along with potential solutions to address them. | 10 (5+5) | CO2 |

| Q.8 | Consider the following table and a task to find the number of employees per department. Show how MapReduce would perform the operations using a diagram. | 10 | CO3 |
|---|---|---|---|

| Employee ID | Name | Department |
|---|---|---|
| 101 | Alice | HR |
| 102 | Bob | IT |
| 103 | Charlie | IT |
| 104 | Diana | HR |
| 105 | Eve | Finance |

| Q.9 | Explain the following components of Data Lakes: storage, data format, compute, and metadata. | 10 | CO4 |
|---|---|---|---|
| | OR | | |
| | Compare and Contrast Data Lakes to Data Warehouses on the following points: Modularity, Schema Enforcement, Cost and Data Format. | 10 | |

<div align="center">SECTION-C</div>

| Q.10 | Differentiate between probabilistic and deterministic classifiers. Write a short note on the following classification algorithms:<br>a) Logistic Regression<br>b) Decision Trees<br>c) Naïve Bayes | 20 (5*4) | CO1 |
|---|---|---|---|
| Q.11 | Explain with examples **any four** of the following data cleaning and transformation techniques for text data.<br>a) Sentence Segmentation<br>b) Tokenization<br>c) Stop-word Removal<br>d) Stemming<br>e) Lemmatization<br><div align="center">OR</div> | 20 (5*4) | CO3 |
| A) | Explain any four factors to consider for replication factor of data in HDFS along with a suitable example of each. | 8 | |
| B) | Explain three levels of data locality in HDFS. | 12 | CO4 |