# PRESIDENCY UNIVERSITY

## BENGALURU

### End Term Examinations –December-2025

**Date:** 15-12-2025                              **Time:** 1.00pm to 04.00pm

| | |
|---|---|
| **School:** SOCSE | **Program:** BD |
| **Course Code** : CBD3406 | **Course Name:** Introduction to Data Science and Big Data |
| **Semester**: V | **Max Marks**: 100     **Weightage**: 50% |

| CO – Levels | CO1 | CO2 | CO3 | CO4 | CO5 |
|---|---|---|---|---|---|
| **Marks** | 24 | 24 | 24 | 28 | |

**Instructions:**

    (i)  Read all questions carefully and answer accordingly.
    (ii) Do not write anything on the question paper other than roll number.

## Part A

**Answer ALL the Questions. Each question carries 2marks.**           **10Q x2M=20M**

| | | | | |
|---|---|---|---|---|
| 1. | Define data science | 2 Marks | L1 | CO1 |
| 2. | State the meaning of 'Velocity' in the 5Vs of Big Data. | 2 Marks | L1 | CO1 |
| 3. | Define standard deviation. | 2 Marks | L2 | CO2 |
| 4. | Define  outliers in a dataset. | 2 Marks | L2 | CO2 |
| 5. | List any two commonly used  distance metrics  in supervised learning. | 2 Marks | L2 | CO3 |
| 6. | List two popular classification algorithms used in machine learning. | 2 Marks | L2 | CO3 |
| 7. | Illustrate the roles of the Name  node | 2 Marks | L2 | CO4 |
| 8. | List any two sources of semi-structured data. | 2 Marks | L2 | CO4 |
| 9. | Define Spark DataFrame | 2 Marks | L1 | CO4 |
| 10. | Illustrate the difference between SQL and NoSQL databases. | 2 Marks | L2 | CO4 |

**PART-B**

**Answer all the Questions.**                    **Total Marks 80M**

| 11. | Outline the key responsibilities in data science team<br><br>**Or**<br><br>Summarize the various stages of the Data Science Workflow | 10 Marks | L1 | CO1 |
|-----|---|---|---|---|
| 12. | | | | |
| 13. | Outline the architecture and components of Hadoop Distributed File System (HDFS) with a neat diagram<br><br>**Or**<br><br>Summarize 5 VS with real world example | 10 Marks | L1 | CO1 |
| 14. | | | | |
| 15. | Identify data quality issues in the following dataset and write the cleaning steps to correct them.<br><br>**Raw Dataset**<br><br>| Cust_ID | Age | Rating | Gender | State | Purchase |<br>|---|---|---|---|---|---|<br>| 501 | 28 | 4 | Female | TN | 10-09-2024 |<br>| 502 | NA | Five | M | KL | 2024/08/15 |<br>| 503 | 999 | 3 | male | KA | Null |<br>| 504 | 40 | 2 | F | TN | 09-07-2024 |<br>| 505 | 32 | 5 | Female | Tamil | 15/08/2024 |<br><br>**Or**<br><br>Explain its purpose of descriptive statistics | 10 Marks | L2 | CO2 |
| 16. | | | | |
| 17. | Illustrate and show the image of Heat map for the following data<br><br>| Age | Avg_Session_Time | Total_Purchases | Clicks_Per_Session | Customer_Rating |<br>|---|---|---|---|---|<br>| 58 | 19.11 | 13 | 7.2 | 4 |<br>| 48 | 8.9 | 6 | 14.6 | 5 |<br>| 34 | 16.04 | 16 | 7.7 | 5 |<br>| 27 | 5.2 | 5 | 9 | 3 |<br>| 40 | 8.36 | 5 | 12.4 | 4 | | 10 Marks | L2 | CO2 |

| | | | | |
|---|---|---|---|---|
| | **Or** | | | |
| 18. | Explain with chart of a histogram to show frequency distribution of sample student marks data | | | |
| 19. | Apply Linear Regression using a simple dataset and explain the concept of correlation between dependent and independent variables. | 10 Marks | L3 | CO3 |
| 20. | **Or**<br><br>Apply Logistic regression to sample dataset | | | |
| 21. | Apply True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) with example | 10 Marks | L3 | CO3 |
| 22. | **Or**<br><br>Apply a K-NN classifier with a small data set (10 Marks) | | | |
| 23. | Describe YARN'S primary role in managing big data workloads within the Hadoop ecosystem | 10 Marks | L2 | CO4 |
| 24. | **Or**<br><br>Describe the working of the Map phase and Reduce phase in MapReduce with word count example | | | |
| | | | | |
| 25. | Illustrate about Apache S**park** | 10 Marks | L2 | CO4 |
| 26 | **Or**<br><br>Explain the architecture and advantages of using Spark SQL in data science | | | |