# JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT

## TEST-2 EXAMINATION- 2025

### B.Tech-VII Semester (CSE/IT)

COURSE CODE (CREDITS): 19B1WCI731 (2)      MAX. MARKS: 25

COURSE NAME: Computational Data Analysis

COURSE INSTRUCTORS: Ekta Gandotra      MAX. TIME: 1 Hour 30 Min

*Note:*   *(a) All questions are compulsory.*
     *(b) The candidate is allowed to make Suitable numeric assumptions wherever required for solving problems*
     *(c) Calculator is allowed.*

| Q. No. | Question | CO | Marks |
|---|---|---|---|
| Q1. | Consider the following dataset: | 4 | |
| | <table><tr><th>Employee_ID</th><th>Department</th><th>Project Completed (Yes/No)</th></tr><tr><td>E101</td><td>Sales</td><td>Yes</td></tr><tr><td>E102</td><td>HR</td><td>No</td></tr><tr><td>E103</td><td>IT</td><td>Yes</td></tr><tr><td>E104</td><td>Sales</td><td>No</td></tr><tr><td>E105</td><td>IT</td><td>Yes</td></tr></table> | | |
| | a. Calculate the Information Gain of the Department feature with respect to the target variable Project Completed. | | 3 |
| | b. Evaluate why selecting Employee_ID as the top feature based on Information Gain may not be appropriate for predicting whether an employee completed a project. | | 2 |
| Q2. | A random sample of 30 students was surveyed, and each student was asked whether they attended a coaching class. The results are summarized below: | 4 | 5 |

Employee_ID / Department / Project Completed (Yes/No):

| Employee_ID | Department | Project Completed (Yes/No) |
|---|---|---|
| E101 | Sales | Yes |
| E102 | HR | No |
| E103 | IT | Yes |
| E104 | Sales | No |
| E105 | IT | Yes |

| Attended Coaching | Passed Exam | |
|---|---|---|
| | Yes | No |
| Yes | 12 | 6 |
| No | 4 | 8 |

Apply the Chi-Square ($\chi^2$) Test of Independence to check whether attending coaching is significantly associated with passing the exam. (Note: The critical value of $\chi^2$ with 1 degree of freedom is 3.841 at 5% level of significance).

| Q3. | Consider the following two-dimensional dataset consisting of four points: (2,0), (0,1), (3,4), (5,2) | 4 | |
|---|---|---|---|
| | a. Compute the first principal component (PC1) of this dataset using the PCA algorithm. | | 2 |
| | b. Project the given points onto PC1. | | 2 |
| | c. Calculate the proportion of total variance explained by PC1 and PC2. | | 2 |
| Q4. | Consider the following dataset. Apply the K-Nearest Neighbor algorithm with K = 3 to predict the rent of a house having size of 1600 Sq. Ft. and 3 occupants. Use Euclidean distance as distance metric. Perform any necessary preparation of the features before computing distances. | 2 | 5 |

| Size (Sq. Ft.) | Occupants | Rent (₹/month) |
|---|---|---|
| 550 | 1 | 8000 |
| 750 | 2 | 10000 |
| 1200 | 3 | 15000 |
| 2000 | 5 | 23000 |
| 1800 | 4 | 21000 |
| 950 | 2 | 12000 |

| Q5. | a. What is the purpose of the kernel function in SVM? How does it help the model to handle non-linearly separable data? | 2 | 1 |
|---|---|---|---|
| | b. Give two limitations of using kernel functions in SVM. | | 1 |
| | c. Explain how logistic regression handles categorical data. Illustrate your answer with an example. | | 2 |