


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

UPES
End Semester Examination, December 2023

Course: Probability and Statistics
Semester: III
Program: B. Tech. CSE
Course Code: CSEG 2036P

Time: 3 hrs.
Max. Marks: 100

SECTION A
(5Qx4M=20Marks)

| S. No. | | Marks | CO | | | | | | | | | | | | |
|--------|--|---------------|---------------------|-------|-------|-------|---------------|---------------|---------------|-------|---------------|---------------|---------------|----------|------------|
| Q 1 | <p>Discuss covariance of random variables. Illustrate that for random variable X and Y,</p> $\text{Cov}(aX + b, cY + d) = ac \times \text{Cov}(X, Y)$ <p>for constants a, b, c and d.</p> | 4 | CO2 | | | | | | | | | | | | |
| Q 2 | <p>Define Marginal Probability Distributions. Apply your understanding of marginal probability mass functions to evaluate that for random variables X and Y, given the joint probability mass function,</p> <table border="1" style="margin: 10px auto; border-collapse: collapse;"> <thead> <tr> <th style="width: 20%;"></th> <th style="width: 20%;">Y = 0</th> <th style="width: 20%;">Y = 1</th> <th style="width: 20%;">Y = 2</th> </tr> </thead> <tbody> <tr> <th style="text-align: center;">X = 0</th> <td style="text-align: center;">$\frac{1}{9}$</td> <td style="text-align: center;">$\frac{1}{6}$</td> <td style="text-align: center;">$\frac{1}{9}$</td> </tr> <tr> <th style="text-align: center;">X = 1</th> <td style="text-align: center;">$\frac{1}{6}$</td> <td style="text-align: center;">$\frac{1}{9}$</td> <td style="text-align: center;">$\frac{1}{3}$</td> </tr> </tbody> </table> <p>Identify if X and Y statistically independent.</p> | | Y = 0 | Y = 1 | Y = 2 | X = 0 | $\frac{1}{9}$ | $\frac{1}{6}$ | $\frac{1}{9}$ | X = 1 | $\frac{1}{6}$ | $\frac{1}{9}$ | $\frac{1}{3}$ | 4 | CO1 |
| | Y = 0 | Y = 1 | Y = 2 | | | | | | | | | | | | |
| X = 0 | $\frac{1}{9}$ | $\frac{1}{6}$ | $\frac{1}{9}$ | | | | | | | | | | | | |
| X = 1 | $\frac{1}{6}$ | $\frac{1}{9}$ | $\frac{1}{3}$ | | | | | | | | | | | | |
| Q 3 | <p>Outline what is meant by random variables. Identify c and d if we have a random variable X with the associated probability density function,</p> $f(x) = cx^{d-1}, 0 \leq x \leq 1$ <p>and if the second central moment $E[X^2]$ is 0.6.</p> | 4 | CO1, CO2 | | | | | | | | | | | | |

| | | | |
|-----|--|---|-----|
| Q 4 | <p>Define sample spaces. Identify the set expression as well as Venn diagram representation for the following cases:</p> <ol style="list-style-type: none"> At least one of the events $A, B,$ or C occurs At most two of the events $A, B,$ or C occur. <p>for a sample space S and three events A, B and C.</p> | 4 | CO1 |
|-----|--|---|-----|

| | | | |
|-----|---|---|-----|
| Q 5 | <p>Discuss correlation coefficient. Identify $Var(X'), Var(Y')$ and $r_{X'Y'}$ in terms of $Var(X), Var(Y)$ and r_{XY} respectively, if X is the height of students in a class in centimeters and Y is the weight of the students in kilograms, and we undertake a transformation to height in inches (X') and weight in pounds (Y'):</p> $X \rightarrow X' = 0.3937 \times X$ $Y \rightarrow Y' = 2.2046 \times Y$ | 4 | CO2 |
|-----|---|---|-----|

SECTION B
(4Qx10M= 40 Marks)

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---------|---|---------|------|------|------|------|---------|----|----|----|----|---------|----|----|----|----|-----|-----|-----|-----|-----|-----|-----|------|-----|------|------|------|------|------|------|------|----|-----|
| Q 6 | <p><i>Choice 1:</i> Define the Kruskal Wallis H Test, its null and alternate hypothesis as well as its relevant test statistic. Describe any one assumption relevant to this statistical test. Highlight how it is better than one-way ANOVA.</p> <p>Apply your understanding of the Kruskal Wallis H Test for analyzing the scores of three groups of students (Group A, Group B and Group C) with</p> <table border="1" data-bbox="318 1234 1084 1352"> <tr> <td>Group A</td> <td>73</td> <td>76</td> <td>87</td> <td>91</td> </tr> <tr> <td>Group B</td> <td>66</td> <td>72</td> <td>81</td> <td>83</td> </tr> <tr> <td>Group C</td> <td>62</td> <td>64</td> <td>71</td> <td>74</td> </tr> </table> <p><u>Given:</u> The critical value for the H test for 2 degrees of freedom and $n_1 = 4, n_2 = 4$ and $n_3 = 4$ at $\alpha = 0.05$ is 5.692.</p> <p><i>Choice 2:</i> Define regression, principle of least squares and residuals. Describe what is meant by multiple regression model.</p> <p>Apply your understanding of nonlinear regression to fit a least-square curve of the form $y = \frac{b}{x(x-a)}$ to the following data:</p> <table border="1" data-bbox="241 1713 1162 1793"> <tr> <td>x</td> <td>3.6</td> <td>4.8</td> <td>6.0</td> <td>7.2</td> <td>8.4</td> <td>9.6</td> <td>10.8</td> </tr> <tr> <td>y</td> <td>0.83</td> <td>0.31</td> <td>0.17</td> <td>0.10</td> <td>0.07</td> <td>0.05</td> <td>0.04</td> </tr> </table> | Group A | 73 | 76 | 87 | 91 | Group B | 66 | 72 | 81 | 83 | Group C | 62 | 64 | 71 | 74 | x | 3.6 | 4.8 | 6.0 | 7.2 | 8.4 | 9.6 | 10.8 | y | 0.83 | 0.31 | 0.17 | 0.10 | 0.07 | 0.05 | 0.04 | 10 | CO4 |
| Group A | 73 | 76 | 87 | 91 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Group B | 66 | 72 | 81 | 83 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Group C | 62 | 64 | 71 | 74 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| x | 3.6 | 4.8 | 6.0 | 7.2 | 8.4 | 9.6 | 10.8 | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| y | 0.83 | 0.31 | 0.17 | 0.10 | 0.07 | 0.05 | 0.04 | | | | | | | | | | | | | | | | | | | | | | | | | | | |

| Q 7 | <p>Define and illustrate the concepts of hypothesis, null hypothesis, alternative hypothesis, significance level and Type I/II errors.</p> <p>Remember and explain the generalized scheme for Hypothesis Testing.</p> <p>Apply and discuss Hypothesis Testing for Population Standard Deviations.</p> <p>Illustrate if the difference between standard deviations of samples A and B drawn from a normal distribution with population standard deviation $\sigma = 15$ is significant (for $\alpha = 0.1$) if $s_A = 8$, $s_B = 10$, $n_A = 150$ and $n_B = 200$, given</p> <table border="1" data-bbox="240 709 1162 785"> <tr> <td>z</td> <td>-2</td> <td>-1.9</td> <td>-1.8</td> <td>-1.7</td> <td>-1.6</td> <td>-1.5</td> <td>-1.4</td> </tr> <tr> <td>p</td> <td>0.04</td> <td>0.06</td> <td>0.07</td> <td>0.09</td> <td>0.11</td> <td>0.13</td> <td>0.16</td> </tr> </table> | z | -2 | -1.9 | -1.8 | -1.7 | -1.6 | -1.5 | -1.4 | p | 0.04 | 0.06 | 0.07 | 0.09 | 0.11 | 0.13 | 0.16 | 10 | CO4 | | | | | |
|--------|---|--------|--------|--------|------|------|------|------|------|-----|------|------|------|------|------|------|------|----|-----|----|---|---|----|-----|
| z | -2 | -1.9 | -1.8 | -1.7 | -1.6 | -1.5 | -1.4 | | | | | | | | | | | | | | | | | |
| p | 0.04 | 0.06 | 0.07 | 0.09 | 0.11 | 0.13 | 0.16 | | | | | | | | | | | | | | | | | |
| Q 8 | <p>Define ANOVA and the motivation behind formulating such a statistical characterization.</p> <p>Discuss the different kinds of sums of squares, degrees of freedom and mean squares.</p> <p>Identify the relevant test statistic (ANOVA coefficient) in terms of mean squares, for ANOVA.</p> <p>Apply ANOVA ($\alpha = 0.05$) for an experiment testing three types of food on separate groups of rats over a 5-week period. The main goal is to determine if there were any significant differences in the average weight (measured in grams) of the rats each week across the three food groups. The data is</p> <table border="1" data-bbox="354 1409 1052 1682"> <thead> <tr> <th>Food 1</th> <th>Food 2</th> <th>Food 3</th> </tr> </thead> <tbody> <tr> <td>8</td> <td>4</td> <td>11</td> </tr> <tr> <td>12</td> <td>5</td> <td>8</td> </tr> <tr> <td>19</td> <td>4</td> <td>7</td> </tr> <tr> <td>8</td> <td>6</td> <td>13</td> </tr> <tr> <td>6</td> <td>9</td> <td>7</td> </tr> <tr> <td>11</td> <td>7</td> <td>9</td> </tr> </tbody> </table> <p><u>Given:</u> F-statistic for degrees of freedom (2,15) at $\alpha = 0.05$ is 3.68.</p> | Food 1 | Food 2 | Food 3 | 8 | 4 | 11 | 12 | 5 | 8 | 19 | 4 | 7 | 8 | 6 | 13 | 6 | 9 | 7 | 11 | 7 | 9 | 10 | CO4 |
| Food 1 | Food 2 | Food 3 | | | | | | | | | | | | | | | | | | | | | | |
| 8 | 4 | 11 | | | | | | | | | | | | | | | | | | | | | | |
| 12 | 5 | 8 | | | | | | | | | | | | | | | | | | | | | | |
| 19 | 4 | 7 | | | | | | | | | | | | | | | | | | | | | | |
| 8 | 6 | 13 | | | | | | | | | | | | | | | | | | | | | | |
| 6 | 9 | 7 | | | | | | | | | | | | | | | | | | | | | | |
| 11 | 7 | 9 | | | | | | | | | | | | | | | | | | | | | | |
| Q 9 | <p>Define a Decision Tree. Describe what is node purity and highlight one advantage and one disadvantage of using Decision Trees.</p> | 10 | CO5 | | | | | | | | | | | | | | | | | | | | | |

Expand on the two ways in which Decision Trees can have variable selection criterion for node allocation.

Apply your understanding of the *Gini index approach* for Decision Trees to analyze 15 students' performance in an online exam. The predictors for this data-set encompass details such as whether the student is enrolled in other online courses, their academic background and whether they are currently employed or not.

| S.No. | Target Variable | Predictor Variables | | |
|-------|-----------------|----------------------|--------------------|----------------|
| | Result | Other Online Courses | Student Background | Working Status |
| 1. | Pass | Yes | Mathematics | Not Working |
| 2. | Fail | No | Mathematics | Working |
| 3. | Fail | Yes | Mathematics | Working |
| 4. | Pass | Yes | CS | Not Working |
| 5. | Fail | No | Other | Working |
| 6. | Fail | Yes | Other | Working |
| 7. | Pass | Yes | Mathematics | Not Working |
| 8. | Pass | Yes | CS | Not Working |
| 9. | Pass | No | Mathematics | Working |
| 10. | Pass | No | CS | Working |
| 11. | Pass | Yes | CS | Working |
| 12. | Pass | No | Mathematics | Not Working |
| 13. | Fail | Yes | Other | Working |
| 14. | Fail | No | Other | Not Working |
| 15. | Fail | No | Mathematics | Working |

SECTION-C
(2Qx20M=40 Marks)

Q 10

Define a Poisson random variable $X \sim \text{Po}(\lambda)$ and **highlight** the expression for its probability distribution.

Derive the mean and variance of the Poisson distribution, considering $e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$. **Show** that this probability distribution satisfies the properties of probabilities.

Define a Gamma Function and **highlight** any two properties of the Gamma Function. **Expand** on your understanding of the Gamma Distribution $Y \sim \text{Gamma}(\alpha, \beta)$, with the expression for its probability distribution.

Derive the mean and variance of the Gamma Distribution.

20

CO3

Identify the values of $\Gamma(4)$, $\Gamma\left(\frac{7}{2}\right)$ and $\Gamma(-3)$.

Apply your understanding of cumulative distribution functions to show that $P(Y \leq \lambda) = P(X \geq \alpha)$ for $X \sim \text{Po}(\lambda)$ and $Y \sim \text{Gamma}(\alpha, \beta)$, given that the cumulative distribution function for the Poisson distribution is

$$F(x, \lambda) = \sum_{k=0}^x \frac{e^{-\lambda} \lambda^k}{k!}$$

and we take $\alpha = 2$, $\beta = 1$.

Q 11

Define a normal distribution and a standard normal table. **Derive** the points of inflection of a normal distribution.

Calculate the probability that a randomly selected student from UPES has IQ lesser than 70, given that the IQ scores of the students of UPES follow a normal distribution with a mean (μ) of 100 and a standard deviation (σ) of 15.

Given: The following segment of the standard normal table

| | | | | | | | |
|-------|-------|-------|-------|-----|-------|-------|-------|
| z | -3 | -2 | -1 | 0 | 1 | 2 | 3 |
| Value | 0.001 | 0.023 | 0.159 | 0.5 | 0.841 | 0.977 | 0.999 |

Determine all IQ scores that comprise the top 10% of the class, given that the z -score corresponding to $z \approx 1.3$ is 0.9.

Discuss sample statistics and **describe** the Method of Moments (MoM).

Highlight any two properties of a good estimator in sample statistics.

Choice 1: **Identify** the MoM estimator of the population parameters for n independent and identically distributed samples taken from a Gamma distribution.

Choice 2: **Calculate** the probability that the sample mean height of these students (for a sample of 25 students taken from the distribution mentioned above) is greater than 106.

Given: $p_{\alpha=0.05}(z = 2) = 0.977$.

20

CO3