

U.S.N.								
--------	--	--	--	--	--	--	--	--

# B.M.S. College of Engineering, Bengaluru-560019

Autonomous Institute Affiliated to VTU

## June 2025 Semester End Main Examinations

**Programme: B.E.**

**Semester: III**

**Branch: CSE(DS)/AI & DS**

**Duration: 3 hrs.**

**Course Code: 23DS3PCFDS**

**Max Marks: 100**

**Course: Foundations of Data Science**

**Instructions:** 1. Answer any FIVE full questions, choosing one full question from each unit.  
2. Missing data, if any, may be suitably assumed.

UNIT - I			CO	PO	Marks
1	a)	<p>Identify and classify as nominal, ordinal interval or ratio data. Justify your answer.</p> <ul style="list-style-type: none"> <li>i) A survey asks users to rate a restaurant on a scale from 1-10.</li> <li>ii) Measuring the weight of different packages in a shipping warehouse, where the weight has a true zero.</li> <li>iii) Recording the temperatures in different cities throughout the day, where the temperature scale has equal intervals.</li> </ul>	CO2	PO2	<b>6</b>
	b)	Implement a python code to calculate the standard deviation for the given temperature data 31,32,32,31,28,31,29	CO3	PO3	<b>4</b>
	c)	A company works on data science projects, one such project is to analyse the precautionary measures for the government to carry out prior to a cyclone. Grievances are collected from social media and other platforms. Discuss the steps involved in this data science process.	CO1	PO1	<b>10</b>
<b>OR</b>					
2	a)	Enumerate on Data Science Venn diagram and elaborate the role of math and statistical knowledge.	CO1	PO1	<b>8</b>
	b)	<p>Classify the following as structured and unstructured data and justify your answer:</p> <ul style="list-style-type: none"> <li>i) The data is scientific observations organized in a CSV file</li> <li>ii) The data is collected from facebook posts.</li> </ul>	CO2	PO2	<b>6</b>
	c)	<p>Illustrate transformation of data for</p> <ul style="list-style-type: none"> <li>i. Joining 2 tables</li> <li>ii. Appending data from tables</li> <li>iii. Creating a view</li> </ul>	CO1	PO1	<b>6</b>

**Important Note:** Completing your answers, compulsorily draw diagonal cross lines on the remaining blank pages. Revealing of identification, appeal to evaluator will be treated as malpractice.

		UNIT - II															
3	a)	<p>Define the term Probability. Analyse the statements below and identify the type of probability.</p> <ol style="list-style-type: none"> <li>To find the probability that the card drawn is a heart</li> <li>To find the probability that a card is 4 and red</li> </ol>	CO2	PO2	5												
	b)	<p>Consider a sample of 5 friends who read books. The frequencies of reading the book are  <math>\text{Friends\_frequency} = [6,4,2,3,5,1,2]</math></p> <p>Write the python code for</p> <ol style="list-style-type: none"> <li>Mean</li> <li>Median</li> <li>Standard deviation</li> <li>Z-Score</li> <li>happiness among friends =&gt; happiness = [0.8,0.6,0.4,0.1,0.2,0.3,0.2]. calculate correlation between happiness and Friends_frequency</li> </ol>	CO3	PO3	10												
	c)	<p>Use the prebuilt python modules to calculate confidence interval.</p> <p>Consider the sample size as 100 and confidence interval as 95 %.</p>	CO3	PO3	5												
		OR															
4	a)	<p>7000 athletes' assembly for a Commonwealth games. As part of the survey, we need to know how many athletes take a break to during the games.</p> <ol style="list-style-type: none"> <li>Apply poisson distribution to simulate 3000 employees who take 60 minutes break.</li> <li>Apply poisson distribution to simulate 4000 employees who take 120 minutes break.</li> <li>Apply sample distribution with 250 different point estimates with break times of size 90 each.</li> </ol>	CO1	PO1	8												
	b)	Elaborate the Hypothesis tests steps involved to query if the sample coffee beans vary significantly in taste from the entire population of beans. Also address the issues involved.	CO2	PO2	7												
	c)	<p>A survey conducted categorizes adult BMIs into 4 classes: Normal, over, obesity and extreme obesity with the percentages 31.2%, 33.1%, 29.4% and 6.3% respectively. A total of 500 adults were randomly sampled. Calculate the test statistic for Chi-Square Goodness Fit.</p> <p>Note – significance level = 0.05</p> <table border="1"> <thead> <tr> <th></th><th>Under/Normal</th><th>Over</th><th>Obesity</th><th>Extreme Obesity</th><th>Total</th></tr> </thead> <tbody> <tr> <td>Observed</td><td>102</td><td>178</td><td>186</td><td>34</td><td>500</td></tr> </tbody> </table>		Under/Normal	Over	Obesity	Extreme Obesity	Total	Observed	102	178	186	34	500	CO1	PO1	5
	Under/Normal	Over	Obesity	Extreme Obesity	Total												
Observed	102	178	186	34	500												
		UNIT - III															
5	a)	<p>A random sample consists of 500 students and the marks scored is recorded as x and y respectively.</p> <ol style="list-style-type: none"> <li>Determine the regression equation for the data.</li> <li>Use the regression equation to predict the marks of 5 students in Quiz who is 5 and 6.</li> <li>Calculate the coefficient of determination (R2) and comment on it.</li> </ol>	CO1	PO1	7												

		<table border="1"> <thead> <tr> <th>Hours studied</th><th>Test Score</th></tr> </thead> <tbody> <tr> <td>1</td><td>50</td></tr> <tr> <td>2</td><td>55</td></tr> <tr> <td>3</td><td>65</td></tr> <tr> <td>4</td><td>70</td></tr> <tr> <td>5</td><td>75</td></tr> <tr> <td>6</td><td>85</td></tr> </tbody> </table>	Hours studied	Test Score	1	50	2	55	3	65	4	70	5	75	6	85		
Hours studied	Test Score																	
1	50																	
2	55																	
3	65																	
4	70																	
5	75																	
6	85																	
	b)	Depict the various stages of the multinomial logistic regression model with a neat diagram and elaborate each stage.	CO1	PO1 <b>8</b>														
	c)	Define correlation. A dataset contains person's height (X) and weight (Y) across 5 years. Enumerate on the correlation between X and Y, emphasizing on the type of correlations.	CO3	PO3 <b>5</b>														
<b>OR</b>																		
6	a)	Write the python code for Linear regression modelling	CO3	PO3 <b>6</b>														
	b)	Derive the equation for logistic regression and identify the category of regression models for the below data samples i) India won the toss. This also implies England lost the toss ii) At a party "Alcohol", "Soft Drinks", "Water" was being served. iii) Students' performance gradually improved from "Better" to "Best".	CO2	PO2 <b>8</b>														
	c)	Construct the dummy variable trap for the following tables below:  <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th colspan="1"><b>Major (k =4)</b></th></tr> </thead> <tbody> <tr> <td>Computer Science</td></tr> <tr> <td>Engineering</td></tr> <tr> <td>Business</td></tr> <tr> <td>Literature</td></tr> <tr> <td>Business</td></tr> <tr> <td>Engineering</td></tr> </tbody> </table> <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th colspan="1"><b>Competition Results</b></th></tr> </thead> <tbody> <tr> <td>Win</td></tr> <tr> <td>Lose</td></tr> <tr> <td>Lose</td></tr> <tr> <td>Win</td></tr> <tr> <td>Lose</td></tr> </tbody> </table>	<b>Major (k =4)</b>	Computer Science	Engineering	Business	Literature	Business	Engineering	<b>Competition Results</b>	Win	Lose	Lose	Win	Lose	CO3	PO3 <b>6</b>	
<b>Major (k =4)</b>																		
Computer Science																		
Engineering																		
Business																		
Literature																		
Business																		
Engineering																		
<b>Competition Results</b>																		
Win																		
Lose																		
Lose																		
Win																		
Lose																		
		<b>UNIT - IV</b>																
7	a)	Employees that they have undergone professional training and they have rated the course as shown below. Employees are 'A', 'B' and 'C' as shown	CO3	PO3 <b>7</b>														

		<p>‘A’ : [1,2,None,4,5],  ‘B’ : [3,None, 3,1,5],  ‘C’ : [4,5,3,None,5]</p> <p>But the rating contains missing values for few courses. Interpret with python code the approaches the data imputation technique uses to identify the missing values.</p>			
	b)	<p>A dataset contains duplicate entries as shown  Name = John, Alice, Bob, John, Age = 25,45,34,25, city = New York, Los Angeles, Los Angles, Chicago .  Write the python code to data with duplicates and de-duplicates.</p>	CO3	PO3	7
	c)	<p>Consider a movie data set and demonstrate the process of Singular Value Decomposition.</p>	CO1	PO1	6
<b>OR</b>					
8	a)	<p>What are the issues of single imputation? Implement multiple mutation using python code.</p>	CO3	PO3	8
	b)	<p>Illustrate pair wise matching by considering the records of an employee.</p>	CO1	PO1	6
	c)	<p>List any 6 common types of inconsistent data.</p>	CO1	PO1	6
<b>UNIT - V</b>					
9	a)	<p>Harry wished they wouldn’t, because he was trying to concentrate on finding his way to classes. Harry started to count the staircases. There were a hundred and forty-two staircases at Hogwarts.  Apply various text analytics tasks tweet.</p> <p>Also write python code snippets by parsing the tweet</p> <ul style="list-style-type: none"> <li>i. Tokenization</li> <li>ii. Remove punctuations in the text, using regex</li> <li>iii. Print the stop words from the text</li> <li>iv. Replace the token ‘Harry’ to ‘Harry Potter’.</li> <li>v. Create a Bag of words</li> </ul>	CO3	PO3	10
	b)	<p>Relate the role of NLP in applications such as</p> <ul style="list-style-type: none"> <li>i) Machine Translation</li> <li>ii) Virtual Assistants</li> </ul>	CO1	PO1	10
<b>OR</b>					
10	a)	<p>A business forum wants to analyze customer reviews to identify common complaints and gauge overall sentiment. Implement the text analytic steps to support the business forum.</p>	CO1	PO1	10
	b)	<p>Implement using python code</p> <ul style="list-style-type: none"> <li>i. porter stemmer</li> <li>ii. bag of words</li> </ul>	CO3	PO3	10

\*\*\*\*\*