

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

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

Autonomous Institute Affiliated to VTU

April 2024 Semester End Main Examinations

Programme: B.E.

Semester: III

Branch: Dept. 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
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.	1	a)	Elaborate on the various steps of the Data Science process.			<i>CO1</i>	<i>PO1</i>	10
		b)	Distinguish between structured and unstructured data with examples.			<i>CO2</i>	<i>PO2</i>	05
		c)	<p>Examine the text data given below and classify the type of data as nominal, ordinal, interval and ratio.</p> <p>"Survey results from a group of students were collected to analyze their preferences for extracurricular activities. Each student was asked to choose their favorite activity from a list of options, rank their top three activities, rate their satisfaction on a scale from 1 to 5, and indicate the number of hours per week they spend on their chosen activity."</p>			<i>CO1</i>	<i>PO1</i>	05
			UNIT - II					
2	a)	Interpret the need for unequal random sampling. Give an example.			<i>CO2</i>	<i>PO2</i>	03	
	b)	Elaborate on the five basic steps of the hypothesis test.			<i>CO1</i>	<i>PO1</i>	07	
	c)	<p>Consider a random sample of 24 friends on Facebook and how many friends they had on Facebook as below.</p> <p>friends = [109, 1017, 1127, 418, 625, 957, 89, 950, 946, 797, 981, 125, 455, 731, 1640, 485, 1309, 472, 1132, 1773, 906, 531, 742, 621]</p> <p>Write a Python program to</p> <ol style="list-style-type: none"> Find the mean, median, range, and standard deviation(SD). Plot a graph for the obtained SD, SD+mean, and SD-mean. Find the z-score and plot its graph. 			<i>CO4</i>	<i>PO3</i>	10	
			OR					

3	a)	<p>Infer how information gain and entropy are related. Consider the dataset below:</p> <table border="1" data-bbox="335 233 1144 907"> <thead> <tr> <th>Age</th><th>Income</th><th>Student</th><th>Credit_rating</th><th>Buys_computer</th></tr> </thead> <tbody> <tr><td>Less than or equal to 30</td><td>High</td><td>No</td><td>Fair</td><td>No</td></tr> <tr><td>Less than or equal to 30</td><td>High</td><td>No</td><td>Excellent</td><td>No</td></tr> <tr><td>Between 31 to 40</td><td>High</td><td>No</td><td>Fair</td><td>Yes</td></tr> <tr><td>Greater than 40</td><td>Medium</td><td>No</td><td>Fair</td><td>Yes</td></tr> <tr><td>Greater than 40</td><td>Low</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>Greater than 40</td><td>Low</td><td>Yes</td><td>Excellent</td><td>No</td></tr> <tr><td>Between 31 to 40</td><td>Low</td><td>Yes</td><td>Excellent</td><td>Yes</td></tr> <tr><td>Less than or equal to 30</td><td>Medium</td><td>No</td><td>Fair</td><td>No</td></tr> <tr><td>Less than or equal to 30</td><td>Low</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>Greater than 40</td><td>Medium</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>Less than or equal to 30</td><td>Medium</td><td>Yes</td><td>Excellent</td><td>Yes</td></tr> <tr><td>Between 31 to 40</td><td>Medium</td><td>No</td><td>Excellent</td><td>Yes</td></tr> <tr><td>Between 31 to 40</td><td>High</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>Greater than 40</td><td>Medium</td><td>No</td><td>Excellent</td><td>No</td></tr> </tbody> </table> <p>i) Calculate the Entropy for the entire data set. ii) Calculate the Information gained for Age, Income, and Credit_rating</p>	Age	Income	Student	Credit_rating	Buys_computer	Less than or equal to 30	High	No	Fair	No	Less than or equal to 30	High	No	Excellent	No	Between 31 to 40	High	No	Fair	Yes	Greater than 40	Medium	No	Fair	Yes	Greater than 40	Low	Yes	Fair	Yes	Greater than 40	Low	Yes	Excellent	No	Between 31 to 40	Low	Yes	Excellent	Yes	Less than or equal to 30	Medium	No	Fair	No	Less than or equal to 30	Low	Yes	Fair	Yes	Greater than 40	Medium	Yes	Fair	Yes	Less than or equal to 30	Medium	Yes	Excellent	Yes	Between 31 to 40	Medium	No	Excellent	Yes	Between 31 to 40	High	Yes	Fair	Yes	Greater than 40	Medium	No	Excellent	No	CO3	PO3	08
Age	Income	Student	Credit_rating	Buys_computer																																																																												
Less than or equal to 30	High	No	Fair	No																																																																												
Less than or equal to 30	High	No	Excellent	No																																																																												
Between 31 to 40	High	No	Fair	Yes																																																																												
Greater than 40	Medium	No	Fair	Yes																																																																												
Greater than 40	Low	Yes	Fair	Yes																																																																												
Greater than 40	Low	Yes	Excellent	No																																																																												
Between 31 to 40	Low	Yes	Excellent	Yes																																																																												
Less than or equal to 30	Medium	No	Fair	No																																																																												
Less than or equal to 30	Low	Yes	Fair	Yes																																																																												
Greater than 40	Medium	Yes	Fair	Yes																																																																												
Less than or equal to 30	Medium	Yes	Excellent	Yes																																																																												
Between 31 to 40	Medium	No	Excellent	Yes																																																																												
Between 31 to 40	High	Yes	Fair	Yes																																																																												
Greater than 40	Medium	No	Excellent	No																																																																												
b)		<p>A new project assignment order is received by an IT company and the authority wants to assign the projects according to the salary package. The dataset is as below.</p> <p>Dataset for Employees</p> <table border="1" data-bbox="319 1275 1144 1641"> <thead> <tr> <th rowspan="2">Results</th> <th colspan="2">Training</th> <th rowspan="2">Total</th> </tr> <tr> <th>Without Professional Training</th> <th>With Professional Training</th> </tr> </thead> <tbody> <tr> <td>Salary package obtained by employees</td> <td>Poor salary</td> <td>05</td> <td>00</td> <td>05</td> </tr> <tr> <td></td> <td>Below-average salary</td> <td>10</td> <td>00</td> <td>10</td> </tr> <tr> <td></td> <td>Average salary</td> <td>40</td> <td>10</td> <td>50</td> </tr> <tr> <td></td> <td>Good salary</td> <td>05</td> <td>30</td> <td>35</td> </tr> <tr> <td></td> <td>Excellent salary</td> <td>00</td> <td>05</td> <td>05</td> </tr> <tr> <td></td> <td>Total</td> <td>60</td> <td>45</td> <td>105</td> </tr> </tbody> </table> <p>Identify the type and find the probability that</p> <ol style="list-style-type: none"> Employees that they have undergone professional training An employee has attended professional training and also has a good salary package An employee has a good salary package given that the employee has not undergone professional training 	Results	Training		Total	Without Professional Training	With Professional Training	Salary package obtained by employees	Poor salary	05	00	05		Below-average salary	10	00	10		Average salary	40	10	50		Good salary	05	30	35		Excellent salary	00	05	05		Total	60	45	105	CO2	PO2	06																																							
Results	Training			Total																																																																												
	Without Professional Training	With Professional Training																																																																														
Salary package obtained by employees	Poor salary	05	00	05																																																																												
	Below-average salary	10	00	10																																																																												
	Average salary	40	10	50																																																																												
	Good salary	05	30	35																																																																												
	Excellent salary	00	05	05																																																																												
	Total	60	45	105																																																																												

	c)	<p>Distinguish between disjoint and non-disjoint events. Identify the type of event for the below scenarios.</p> <ul style="list-style-type: none"> i) Coin toss ii) A used car having heated seats and a manual transmission iii) Being a first-year student and being a sophomore iv) Having blue eyes or brown hair 	CO2	PO2	06												
		UNIT - III															
4	a)	Provide a explanation of how correlation measures the strength and direction of relationships between two variables, emphasizing both positive and negative correlations.	CO1	PO1	05												
	b)	<p>The average rainfall (in mm) for the last 5 years in the state of Karnataka was recorded along with the umbrellas sold in that particular year.</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="text-align: left;">Rainfall (mm)</th> <th style="text-align: right;">Umbrellas Sold</th> </tr> </thead> <tbody> <tr> <td style="text-align: left;">121.2</td> <td style="text-align: right;">52</td> </tr> <tr> <td style="text-align: left;">152.6</td> <td style="text-align: right;">72</td> </tr> <tr> <td style="text-align: left;">98.4</td> <td style="text-align: right;">40</td> </tr> <tr> <td style="text-align: left;">171</td> <td style="text-align: right;">100</td> </tr> <tr> <td style="text-align: left;">85.6</td> <td style="text-align: right;">34</td> </tr> </tbody> </table> <p> <ol style="list-style-type: none"> 1. Find the equation of linear regression. 2. Suppose this year the state records a rainfall of 123.2 mm, what will be the expected number of umbrellas that will be sold? 3. Calculate the coefficient of determination (R^2) and comment on it </p>	Rainfall (mm)	Umbrellas Sold	121.2	52	152.6	72	98.4	40	171	100	85.6	34	CO3	PO3	10
Rainfall (mm)	Umbrellas Sold																
121.2	52																
152.6	72																
98.4	40																
171	100																
85.6	34																
	c)	Summarize Receiver Operating Characteristic(ROC).	CO1	PO1	05												
		OR															
5	a)	Depict the various stages of the multinomial logistic regression model with a neat diagram and explain.	CO1	PO1	08												
	b)	<p>Suppose, in a dataset, a feature “season” has three values – rainy, summer, and winter.</p> <p>Show how the dummy variable trap can be avoided.</p>	CO2	PO2	08												
	c)	Distinguish between linear regression and logistic regression.	CO2	PO2	04												
		UNIT - IV															
6	a)	Distinguish between single and multiple data imputation and write Python code for the same.	CO4	PO3	08												
	b)	Elaborate on the common types of inconsistent data.	CO2	PO2	06												
	c)	Using the Python’s ‘pandas’ library, write a program to perform data deduplication.	CO4	PO3	06												

UNIT - V					
7	a)	Discuss the text analytics subtasks.	<i>CO1</i>	<i>PO1</i>	08
	b)	Elaborate on the stages of natural language processing.	<i>CO2</i>	<i>PO2</i>	06
	c)	Write a Python program to create a bag of words using a counter and also remove the stop words using the ‘nltk’ python library.	<i>CO4</i>	<i>PO3</i>	06

B.M.S.C.E. - ODD SEM 2023-24