



DECEMBER 2018: END SEMESTER ASSESSMENT

UE16CS322: Data Analytics(SMP Section)

Time: 3 Hours (180 Minutes)

Answer All Questions

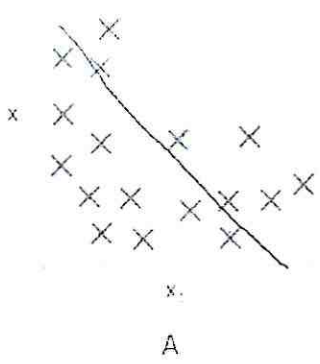
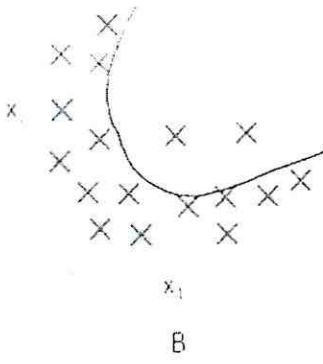
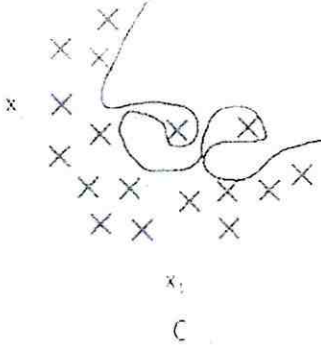
Max Marks: 100

Instructions: Only calculators are permitted. Please show your work-outs of problem clearly.

1	a	The number of customers visiting a shopping mall in Bangalore in the last 20 days is as follows: 232, 277, 261, 173, 283, 197, 251, 212, 213, 213, 229, 164, 219, 196, 186, 247, 244, 269, 216 and 272. Calculate skewness and kurtosis. What do you infer from these? How can the shopping mall use this analysis? (Hint: Kurtosis= $\sum (x_i - \bar{x})^4/n \cdot s^4$), Skewness- apply the same for 3 rd power)	6																								
	b	<p>In the above example, the age details of customers were collected and found to be as follows:</p> <table><tr><td>0-20 years</td><td>500</td></tr><tr><td>20-40</td><td>800</td></tr><tr><td>40-60</td><td>500</td></tr><tr><td>>60</td><td>200</td></tr></table> <p>The mall authorities would like to send 100 emails using clustered sampling and another 100 by stratified sampling to get feedback from people about the Mall. How do you go about constructing samples? What additional information you may need?</p>	0-20 years	500	20-40	800	40-60	500	>60	200	4																
	0-20 years	500																									
	20-40	800																									
40-60	500																										
>60	200																										
c	Explain why with a fixed number of training samples, predictive power decreases with increasing dimensions.	4																									
d	<p>Explain Principal Component Analysis taking the example of Iris data-set below. Just explain the intent and outcome expected.</p> <p>Number of Instances: 150 (50 in each of three classes)</p> <p>Number of Attributes: 4 numeric, predictive attributes and the class.</p> <p>Attribute Information:</p> <ol style="list-style-type: none">1. sepal length in cm2. sepal width in cm3. petal length in cm4. petal width in cm5. class: <p>-- Iris Setosa -- Iris Versicolour -- Iris Virginica</p>	6																									
2	a	<p>Corruption perception Index and Gini Index (measures inequalities) are given for the following countries.</p> <table><tr><td>Country</td><td>Hongkong</td><td>South Korea</td><td>China</td><td>Italy</td><td>Mongolia</td><td>Austria</td><td>Norway</td></tr><tr><td>Corruption Index</td><td>77</td><td>53</td><td>40</td><td>47</td><td>38</td><td>75</td><td>85</td></tr><tr><td>Gini</td><td>53.7</td><td>30.2</td><td>46.2</td><td>32.7</td><td>36.5</td><td>27.6</td><td>23.5</td></tr></table> <p>Draw a simple linear regression model and predict Corruption Perception for India with Gini index of 38. Note: Higher the Corruption index better the transparency. Higher the Gini Index means higher inequality. (Hint: $b_0 = \text{mean of } Y - b_1 \cdot \text{mean of } X$, $b_1 = \frac{\sum xy - n \cdot \text{mean of } x \cdot \text{mean of } y}{\sum x^2 - n \cdot (\text{mean of } x)^2}$)</p>	Country	Hongkong	South Korea	China	Italy	Mongolia	Austria	Norway	Corruption Index	77	53	40	47	38	75	85	Gini	53.7	30.2	46.2	32.7	36.5	27.6	23.5	6
Country	Hongkong	South Korea	China	Italy	Mongolia	Austria	Norway																				
Corruption Index	77	53	40	47	38	75	85																				
Gini	53.7	30.2	46.2	32.7	36.5	27.6	23.5																				

SRN

--	--	--	--	--	--	--	--	--	--	--	--

b		<p>Below are three scatter plots (A, B, C) and hand-drawn decision boundaries for logistic regression. Answer the questions that follow:</p> <div style="display: flex; justify-content: space-around; align-items: center;">    </div> <p style="text-align: center;"> A B C </p> <p>i. Which figure is over-fitting the training data the most and why?</p> <p>ii. In which model is the training error maximum and why?</p> <p>iii. Which model is more robust than the other two models and why?</p>	6																				
c		<p>Before you apply regression model to a data-set what assumptions do you need to validate? What kind of data-cleaning can one do? Explain</p>	4																				
d		<p>Explain the following:</p> <p>(i) Difference between regression and correlation</p> <p>(ii) It was found that married men earn more money by analyzing data using regression model. Can we infer that one should get married to earn money?</p>	4																				
3	a	<p>You are given the transaction data shown in the Table below from a fast food restaurant. There are 9 distinct transactions (order:1 – order:9) and each transaction involves 2-4 meal items. There are a total of 5 meal items that are involved in the transactions. Meal items are names M1-M5 for simplicity</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 30%;">Meal Item</th><th style="width: 70%;">List of Items</th></tr> </thead> <tbody> <tr><td>Order 1</td><td>M1, M2, M5</td></tr> <tr><td>Order 2</td><td>M2, M4</td></tr> <tr><td>Order 3</td><td>M2, M3</td></tr> <tr><td>Order 4</td><td>M1, M2, M4</td></tr> <tr><td>Order 5</td><td>M1, M3</td></tr> <tr><td>Order 6</td><td>M2, M3</td></tr> <tr><td>Order 7</td><td>M1, M3</td></tr> <tr><td>Order 8</td><td>M1, M2, M3, M5</td></tr> <tr><td>Order 9</td><td>M1, M2, M3</td></tr> </tbody> </table> <p>Assume minimum support is $2/9$ and the minimum confidence is $7/9$.</p> <p>i. Identify all the frequent item sets involving M5 as one of the meal items. Calculate the support for each item.</p> <p>ii. What is the support and confidence of the rules and identify which rule is stronger: $M1, M5 \rightarrow M2$ Or $M2, M5 \rightarrow M1$</p>	Meal Item	List of Items	Order 1	M1, M2, M5	Order 2	M2, M4	Order 3	M2, M3	Order 4	M1, M2, M4	Order 5	M1, M3	Order 6	M2, M3	Order 7	M1, M3	Order 8	M1, M2, M3, M5	Order 9	M1, M2, M3	8
Meal Item	List of Items																						
Order 1	M1, M2, M5																						
Order 2	M2, M4																						
Order 3	M2, M3																						
Order 4	M1, M2, M4																						
Order 5	M1, M3																						
Order 6	M2, M3																						
Order 7	M1, M3																						
Order 8	M1, M2, M3, M5																						
Order 9	M1, M2, M3																						
b		<p>Consider the 7 observations in a two-dimensional space. $\{(2,2), (4,4), (6,6), (1,4), (4,0), (5,5), (9,9)\}$. If we run the k-means clustering algorithm, to divide these data-point into 3 clusters. What would be the clusters after the first iteration for each of the data- point. Use Manhattan distance to calculate the distance between clusters. Use $(4,4)$, $(1,4)$ and $(5,5)$ as initial seeds. What would the centroids for the second iteration?</p>	6																				

