



**DECEMBER 2020: END SEMESTER ASSESSMENT (ESA)  
B.TECH FIFTH SEMESTER**

UE18CS312 – Data Analytics (4 credit Elective course)

**Time: 3 Hrs**      **Answer All Questions**      **Max Marks: 100**

1	<p>a) What factors help increase attendance in the primary schools in villages? The Education Ministry has asked you to help them design a survey to gather this information.</p> <p>(i) Who would you recommend collecting the data from? (suggest any two sources or demographics that could provide a meaningful insight)</p> <p>(ii) Suggest any two variables (or attributes) that you would collect and also identify the type of data (categorical/ numeric and ordinal/ interval or ratio) for the two features.</p>	6 (2+4)																				
b)	<p>The following table contains the gross collection and budget (both variables are reported in crores of rupees) for ten movies released in 2020 and declared a ‘Success’ in the box office:</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>Collection</th><th>60</th><th>50</th><th>35</th><th>50</th><th>65</th><th>270.15</th><th>197.3</th><th>39.3</th><th>193.2</th></tr> <tr> <th>Budget</th><td>40</td><td>40</td><td>25</td><td>38</td><td>45</td><td>125</td><td>105</td><td>236</td><td>100</td></tr> </thead> </table> <p>(i) Suggest a visualization that will help us easily infer the ‘profit’ earned by the movie</p> <p>(ii) It is postulated that higher the budget, higher the box office collection. What visualization technique (or graph) will help us easily infer the relationship between collection and budget? Briefly explain.</p>	Collection	60	50	35	50	65	270.15	197.3	39.3	193.2	Budget	40	40	25	38	45	125	105	236	100	4 (2+2)
Collection	60	50	35	50	65	270.15	197.3	39.3	193.2													
Budget	40	40	25	38	45	125	105	236	100													
c)	<p>Given an example for each of the following types of error in data entry and suggest how each can be handled:</p> <p>(i) Incomplete data                   (ii) inconsistent data</p>	4																				
d)	<p>Data gathered from a batch of students is recorded below:</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th></th><th>Knows of Dragon Ball Z</th><th>Does not know of Dragon Ball Z</th></tr> </thead> <tbody> <tr> <th>Knows of Pokémon</th><td>250</td><td>200</td></tr> <tr> <th>Does not know of Pokémon</th><td>50</td><td>1000</td></tr> </tbody> </table> <p>(i) What is the chi-square (<math>\chi^2</math>) statistic for this data?</p> <p>(ii) How can this be used to determine whether knowing of Pokémon is independent of knowing of Dragon Ball-Z? (You may briefly describe the steps to be followed for this.)</p> <p>Note: You may use the following facts in your computation:</p>		Knows of Dragon Ball Z	Does not know of Dragon Ball Z	Knows of Pokémon	250	200	Does not know of Pokémon	50	1000	6 (4+2)											
	Knows of Dragon Ball Z	Does not know of Dragon Ball Z																				
Knows of Pokémon	250	200																				
Does not know of Pokémon	50	1000																				

- 2 a) The following table shows the result of multiple linear regression used to predict the number of minutes students are willing to participate in an after-school activity based on multiple factors. Answer the following questions:

	Beta	Std error	Standardized Beta	t	p
<b>Constant</b>	17.23	50.16		0.34	0.73
<b>Challenge</b>	-5.88	8.88	-0.06	-0.66	0.51
<b>Instant enjoyment</b>	-20.33	18.13	-0.14	-1.12	0.02
<b>Novelty</b>	9.95	10.78	0.09	0.92	0.001
<b>Exploration</b>	28.33	11.00	0.25	2.58	0.03

- (i) Write an equation to represent this model  
(ii) If the significance value alpha = 0.5, what are the features that would be dropped?  
(iii) In the simple linear regression model below, what does the error term,  $\epsilon$  represent?  
 $\text{Novelty} = 16.43 + 5.22 * \text{Exploration} + \epsilon$

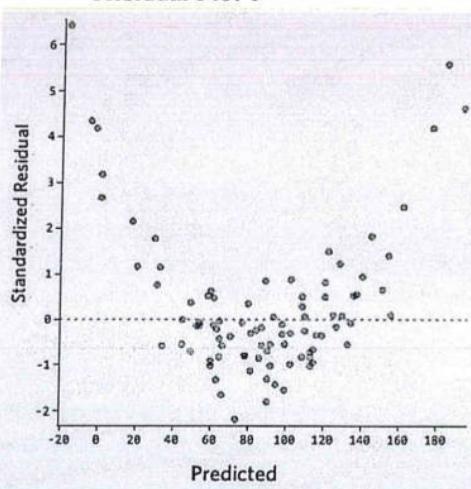
- b) In a sample of 100 engineers, 90 stay up all night the day before an exam whereas in a sample of 80 doctors, only 20 stay up all night the day before an exam.

- (i) What is the odds ratio of an engineer staying up all night the day before an exam to a doctor staying up all night?  
(ii) Suggest any one approach to deal with categorical variables when designing a regression model.  
(For example, consider a month's worth of data from the postal department on zip code (the explanatory variable) and the number of letters delivered to that zip code in a day (as the predicted variable). How would we use zip code in building a regression model for this data?)

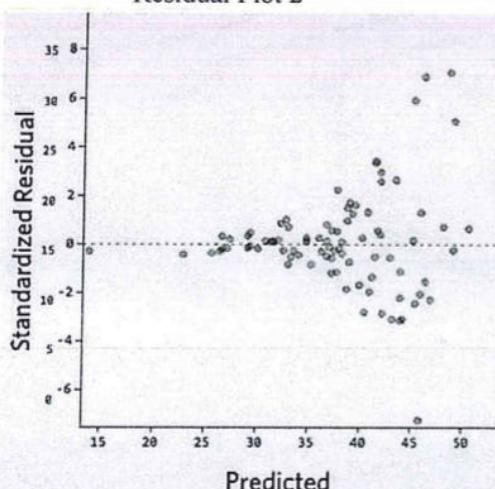
- c) Based on these residual plots from two liner regression models, can the two models be considered good? If there is any problem you can diagnose based on the plots, briefly explain what that is.

(The dotted line indicates the zero reading on the vertical axis in both plots.)

Residual Plot 1



Residual Plot 2



6  
(3\*2)

4  
(2+2)

4  
(2+2)

	d)	<p>A logistic regression model to categorize actions in a video clip yielded the following confusion matrix. Answer the questions below:</p> <table style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th colspan="2"></th><th colspan="3">Predicted class</th></tr> <tr> <th colspan="2"></th><th>Dance</th><th>Martial arts</th><th>Pilates</th></tr> </thead> <tbody> <tr> <th rowspan="3">Actual class</th><th>Dance</th><td>40</td><td>30</td><td>10</td></tr> <tr> <th>Martial arts</th><td>20</td><td>50</td><td>10</td></tr> <tr> <th>Pilates</th><td>10</td><td>10</td><td>60</td></tr> </tbody> </table> <p>(i) What is the recall (or true positive rate) for Martial arts?  (ii) What is the precision for Pilates?  (Of all the clips labeled as Pilates, how many were truly clips of Pilates?)  (iii) What is the average accuracy of this classifier?</p>			Predicted class					Dance	Martial arts	Pilates	Actual class	Dance	40	30	10	Martial arts	20	50	10	Pilates	10	10	60	6 (3*2)
		Predicted class																								
		Dance	Martial arts	Pilates																						
Actual class	Dance	40	30	10																						
	Martial arts	20	50	10																						
	Pilates	10	10	60																						
3	a)	<p>Briefly answer the questions:</p> <p>(i) With a schematic sketch of a stationary signal, list the characteristics of a typical stationary signal. How can we test whether a given time series is stationary?  (ii) What are the typical values of the parameter <math>d</math> (for differencing) in an ARIMA model and how can we determine the other two parameters <math>p</math> (order of the autoregressive component) and <math>q</math> (order of the moving average component) for an ARIMA model? Explain briefly.</p>	8 (4+4)																							
	b)	<p>What model would be most appropriate in each of the following scenarios and why?</p> <p>(i) Predicting the yield of a seasonal crop  (ii) Predicting the purchase of school supplies  (iii) Predicting the stock price of Netflix</p>	6 (2*3)																							
	c)	<p>What do the following statistics tell us about a time series model? Answer each of the questions below:</p> <p>(i) <b>Theil's coefficient – what is this statistic used for and what does it mean if <math>U=1</math>?</b></p> $U = \frac{\sum_{t=1}^n (Y_{t+1} - F_{t+1})^2}{\sum_{t=1}^n (Y_{t+1} - Y_t)^2}, \text{ where } Y \text{ is the actual value, } F \text{ is the forecast and } t \text{ is the time index}$ <p>(ii) <b>Ljung-Box test – what is this test used for and why do we need <math>p</math> and <math>q</math>, the AR and MA lags of the ARIMA model for this test?</b></p> $Q(m) = n(n + 2) \sum_{k=1}^m \frac{\rho^2}{n - k}, \text{ where } n \text{ is the number of observations in the time series, } m \text{ is the total number of lags and } \rho \text{ the autocorrelation of lag } k$ <p>(iii) <b>BIC – what is this used for and how is it different from AIC?</b>  <math>BIC = -2LL + K\ln(n),</math> where LL is the log likelihood function,  <math>K</math> is the number of parameters estimated and  <math>n</math> is the number of observations in the sample</p>	6 (2*3)																							

4	a)	<p>What is the most appropriate recommendation system for each of the following applications? Explain in a line why this is most appropriate.            (You may choose from: content-based, collaborative filtering, knowledge-based (case based, query based or constraint-based) and apriori based association rule-mining)</p> <ul style="list-style-type: none"> <li>(i) PlanMyTrip – a vacation planning website</li> <li>(ii) Facebook friend recommendations to a new user</li> <li>(iii) Recipes on a cooking site</li> <li>(iv) Predicting whether a particular poll will be answered by a customer who has been active on the news site for the past eight months.</li> </ul>	8 (2*4)																									
	b)	<p>In the context of designing a recommendation system, answer the following questions:</p> <ul style="list-style-type: none"> <li>(i) In what way is entropy a better criterion for splitting a node than either Gini index or classification error?</li> <li>(ii) If the confidence of an association rule is high, can we assume the rule would be interesting? Substantiate your answer briefly.</li> </ul>	6 (3+3)																									
	c)	<p>The Department of CSE would like to form review panels for the Capstone projects. Each team has submitted a title of their project and a brief description of the problem statement.</p> <ul style="list-style-type: none"> <li>(i) Outline the steps to cluster similar problem statements together.</li> <li>(ii) There may be multiple configurations of clusters that come up from step (i) above. Suggest an evaluation criterion that will ensure we select a configuration that has maximum inter-cluster distance and minimum intra-cluster distance. Briefly explain how the criterion you suggest can help test these conditions (of minimum intra-cluster distance and maximum inter-cluster distance).</li> </ul>	6 (4+2)																									
5	a)	<p>(i) What is the canonical form of an absorbing state Markov Chain? State what each component matrix stand for.</p> <p>(ii) The number of flights cancelled by an airline is modelled using a Markov chain. The state transition matrix between flight cancellations is shown below:            (0 – no cancellations, 1-3 represents the number of cancellations on a day):</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <th></th> <th>0</th> <th>1</th> <th>2</th> <th>3</th> </tr> <tr> <th>0</th> <td>0.45</td> <td>0.30</td> <td>0.20</td> <td>0.05</td> </tr> <tr> <th>1</th> <td>0.15</td> <td>0.60</td> <td>0.15</td> <td>0.10</td> </tr> <tr> <th>2</th> <td>0.10</td> <td>0.30</td> <td>0.40</td> <td>0.20</td> </tr> <tr> <th>3</th> <td>0</td> <td>0.10</td> <td>0.70</td> <td>0.20</td> </tr> </table> <p>If there are no cancellations initially, what is the probability there will be at least one cancellation after two days?</p>		0	1	2	3	0	0.45	0.30	0.20	0.05	1	0.15	0.60	0.15	0.10	2	0.10	0.30	0.40	0.20	3	0	0.10	0.70	0.20	8 (4+4)
	0	1	2	3																								
0	0.45	0.30	0.20	0.05																								
1	0.15	0.60	0.15	0.10																								
2	0.10	0.30	0.40	0.20																								
3	0	0.10	0.70	0.20																								
b)	<p>Check whether the matrices below are regular.</p> <p>(Note: A regular matrix is one that has all positive entries for some power of the matrix; it is not required to check beyond a power of 4.)</p> <p>(i) <table border="1" style="display: inline-table; vertical-align: middle;">0 0.5 1 0.5</table></p> <p>(ii) <table border="1" style="display: inline-table; vertical-align: middle;">0 1 1 0</table></p>	6 (3+3)																										
c)	<p>(i) Suggest a method to check whether a hidden variable is confounding.</p> <p>(ii) Microsoft Stream is considering introducing a ‘Download all’ button for all video clips in a class group. The advantage is that students do not have to click on each link to download the recorded session for each class. But the disadvantage is that the file is extremely large and students may not have the bandwidth to download the material with one click. As an expert on Data Analytics, outline the steps you would take to design this test for the MS Stream team.</p>	6 (3+3)																										