

26.05.2024 (M)

Maximum Marks: 100	Semester: January 2024 - April 2024 Examination: ESE Examination	Duration: 3 Hrs.
Programme code: 54		
Programme: Honors in Data Science and Analytics	Class: SY	Semester: IV (SVU 2020)
Name of the Constituent College: K. J. Somaiya College of Engineering		Name of the department: COMPUTER
Course Code: 116h54C401	Name of the Course: Applied Data Science	
Instructions: 1) Draw neat diagrams 2) All questions are compulsory 3) Assume suitable data wherever necessary		

Que. No.	Question	Max. Marks
Q1	Solve any Four	
i)	Explain the ETL process.	20
ii)	Explain datafication with the help of any five suitable examples	5
iii)	Explain word analysis for text processing with examples.	5
iv)	Explain the 5 V's of big data with the help of suitable examples.	5
v)	With the help of an example explain stratified random sampling.	5
vi)	Differentiate between classification and clustering in data analysis. Discuss the strengths and weaknesses of each technique, and provide real-world examples of where you might use each approach.	5

Que. No.	Question	Max. Marks												
Q2 A	Solve the following	10												
i)	Suppose that survival drops off rapidly in the year following diagnosis of a certain type of advanced cancer. Suppose that the length of survival (or time-to-death) is a random variable that approximately follows an exponential distribution with parameter 2. What's the probability that a person who is diagnosed with this illness survives a year?	5												
ii)	The number of ships to arrive at a harbor on any given day is a random variable represented by x . The probability distribution for x is:	5												
	<table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td>x</td><td>10</td><td>11</td><td>12</td><td>13</td><td>14</td></tr> <tr> <td>$P(x)$</td><td>0.4</td><td>0.2</td><td>0.2</td><td>0.1</td><td>0.1</td></tr> </table> <p>a. Find the probability that at most 11 ships arrive. b. Find the variance and standard deviation for the number of ships to arrive at the harbor.</p>	x	10	11	12	13	14	$P(x)$	0.4	0.2	0.2	0.1	0.1	
x	10	11	12	13	14									
$P(x)$	0.4	0.2	0.2	0.1	0.1									
	OR													

Q2 A What is cross validation? (2 marks)
Consider the confusion matrix given below:

		Predicted	
		0	1
Actual	0	30	12
	1	8	56

Using the confusion matrix, calculate the following measures and interpret them: (2 marks each)

1. Accuracy
2. Precision
3. Recall
4. F1 score

Q2 B Solve any One

i) Consider a dataset with two features X_1 and X_2 and 15 data points:

ID	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
X_1	2	2	11	6	6	4	5	4	10	7	9	4	3	3	6
X_2	10	6	11	9	4	2	10	9	12	5	11	6	10	8	11

Apply k-means algorithm to this dataset. Assume k=3 and data points with ID 2, 7 and 15 as the initial cluster centroids.

ii) Build the decision tree using the CART algorithm with Gini Index on the following dataset:

City Size	Avg. Income	Local Investors	LOHAS Awareness	Decision
Big	High	Yes	High	Yes
Medium	Med	No	Med	No
Small	Low	Yes	Low	No
Big	High	No	High	Yes
Small	Med	Yes	High	No
Med	High	Yes	Med	Yes
Med	Med	Yes	Med	No
Big	Med	No	Med	No
Med	High	Yes	Low	No
Small	High	No	High	Yes
Small	Med	No	High	No
Med	Heigh	No	Med	No

Que. No.	Question	Max. Marks
Q3	Solve any Two	
i)	What is ensemble learning? Explain the random forest algorithm with an example.	20
ii)	Explain Hidden Markov Model (HMM) with the help of a detailed example.	10
iii)	With the help of an example illustrate the process of inference in a Bayesian Belief Network.	10

Que. No.	Question	Max. Marks																																													
Q4	Solve any Two																																														
i)	<p>Consider the following set of observations:</p> <table border="1"> <thead> <tr> <th>chills</th> <th>runny nose</th> <th>headache</th> <th>fever</th> <th>flu?</th> </tr> </thead> <tbody> <tr> <td>Y</td> <td>N</td> <td>Mild</td> <td>Y</td> <td>N</td> </tr> <tr> <td>Y</td> <td>Y</td> <td>No</td> <td>N</td> <td>Y</td> </tr> <tr> <td>Y</td> <td>N</td> <td>Strong</td> <td>Y</td> <td>Y</td> </tr> <tr> <td>N</td> <td>Y</td> <td>Mild</td> <td>Y</td> <td>Y</td> </tr> <tr> <td>N</td> <td>N</td> <td>No</td> <td>N</td> <td>N</td> </tr> <tr> <td>N</td> <td>Y</td> <td>Strong</td> <td>Y</td> <td>Y</td> </tr> <tr> <td>N</td> <td>Y</td> <td>Strong</td> <td>N</td> <td>N</td> </tr> <tr> <td>Y</td> <td>Y</td> <td>Mild</td> <td>Y</td> <td>Y</td> </tr> </tbody> </table> <p>Build a Naive Bayes classifier model with conditional probability tables.</p> <p>Given a new patient with symptoms: Runny nose=no, Headache=no, Chills=true, Fever=true, use the model to predict whether the patient has flu.</p>	chills	runny nose	headache	fever	flu?	Y	N	Mild	Y	N	Y	Y	No	N	Y	Y	N	Strong	Y	Y	N	Y	Mild	Y	Y	N	N	No	N	N	N	Y	Strong	Y	Y	N	Y	Strong	N	N	Y	Y	Mild	Y	Y	20
chills	runny nose	headache	fever	flu?																																											
Y	N	Mild	Y	N																																											
Y	Y	No	N	Y																																											
Y	N	Strong	Y	Y																																											
N	Y	Mild	Y	Y																																											
N	N	No	N	N																																											
N	Y	Strong	Y	Y																																											
N	Y	Strong	N	N																																											
Y	Y	Mild	Y	Y																																											

ii) There are five data points P_1, P_2, \dots, P_5 .

Use the distance matrix in the following table to perform single link and complete link hierarchical clustering. Show your results by drawing a dendrogram. The dendrogram should clearly show the order in which the points are merged.

	P1	P2	P3	P4	P5
P1	0.00	0.10	0.41	0.55	0.35
P2	0.10	0.00	0.64	0.47	0.98
P3	0.41	0.64	0.00	0.44	0.85
P4	0.55	0.47	0.44	0.00	0.76
P5	0.35	0.98	0.85	0.76	0.00

P.T.O.

Page 314

Que. No.	Question			Max. Marks
Q4 iii)	Consider the following dataset:			10

Height (in cms)	Weight (in kgs)	T Shirt Size
158	58	M
158	59	M
158	63	M
160	59	M
160	60	M
163	60	M
163	61	M
165	65	L
168	62	L
168	63	L
168	66	L
170	63	L
170	64	L
170	68	L

Apply k-NN algorithm to predict the T-shirt size that would fit a person with a height of 164 cm and a weight of 67 kg.

Que. No.	Question	Max. Marks
Q5	Write short notes on any four	20
i)	Logistic Regression	5
ii)	Steps in data preprocessing.	5
iii)	Gradient descent algorithm	5
iv)	SVM	5
v)	Curse of dimensionality	5
vi)	Impact of data science on business	5