

COURSE CODE (CREDITS): 22M11CI112 (3)

MAX. MARKS: 35

COURSE NAME: INTRODUCTION TO DATA SCIENCE

COURSE INSTRUCTORS: Dr Nancy Singla

MAX. TIME: 2 Hours

Note: (a) All questions are compulsory.

(b) The candidate is allowed to make suitable numeric assumptions wherever required for solving problems

(c) Use of calculator is allowed.

Q. No	Question	CO	Marks																								
Q1.	<p>(a) Explain the bias–variance tradeoff in supervised learning. Why can a model with very low training error still perform poorly on unseen data?</p> <p>(b) Which of the following classifiers will be preferred while working with small and large datasets. Explain.</p> <ul style="list-style-type: none">• High bias/low variance classifiers• Low bias/high variance classifiers	CO5	[5+5]																								
Q2.	<p>(a) For the given data, compute two clusters using K-means algorithm for clustering where initial cluster centers are (1.0, 1.0) and (5.0, 7.0). Execute for two iterations.</p> <table><tr><th>Record Number</th><th>A</th><th>B</th></tr><tr><td>R1</td><td>1.0</td><td>1.0</td></tr><tr><td>R2</td><td>1.5</td><td>2.0</td></tr><tr><td>R3</td><td>3.0</td><td>4.0</td></tr><tr><td>R4</td><td>5.0</td><td>7.0</td></tr><tr><td>R5</td><td>3.5</td><td>5.0</td></tr><tr><td>R6</td><td>4.5</td><td>5.0</td></tr><tr><td>R7</td><td>3.5</td><td>4.5</td></tr></table> <p>(b) Explain how would you choose the number of clusters in k-means?</p> <p>(c) Describe a real-world scenario where k-means clustering would be useful?</p>	Record Number	A	B	R1	1.0	1.0	R2	1.5	2.0	R3	3.0	4.0	R4	5.0	7.0	R5	3.5	5.0	R6	4.5	5.0	R7	3.5	4.5	CO6	[5+3+2]
Record Number	A	B																									
R1	1.0	1.0																									
R2	1.5	2.0																									
R3	3.0	4.0																									
R4	5.0	7.0																									
R5	3.5	5.0																									
R6	4.5	5.0																									
R7	3.5	4.5																									
Q3.	<p>(a) You are given a data set consisting of variables having more than 30% missing values. Let us say, out of 50 variables, 8 variables have missing values higher than 30%. How will you deal with them?</p> <p>(b) For each of the following scenarios, state which probability distribution is most suitable and give a brief reason.</p> <ol style="list-style-type: none">A call center receives an average of 12 calls per hour. You want to model the number of calls in the next hour.The heights of adult males in a large population are being	CO3	[2+3]																								

	analyzed. iii. A company sends promotional emails. Each email has a 10% chance of being opened. You want to model the number of opens out of 200 emails.		
Q4.	<p>You are given the following five training instances:</p> <ul style="list-style-type: none"> • $x_1 = 2, x_2 = 1, y = 4$ • $x_1 = 6, x_2 = 3, y = 2$ • $x_1 = 2, x_2 = 5, y = 2$ • $x_1 = 6, x_2 = 7, y = 3$ • $x_1 = 10, x_2 = 7, y = 3$ <p>We want to model this function using the K-nearest neighbor regressor model. Predict the value of y corresponding to $(x_1, x_2) = (3, 6)$ considering k as 2 and 3.</p>	CO6	[5]
Q5.	A dataset contains product reviews stored as string objects. Describe a preprocessing pipeline to prepare this data for machine learning classification using a BoW model. Explain why each step is necessary.	CO5	[5]