## JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT TEST -2 EXAMINATION- 2025

M. Tech-I Semester (CSE)

COURSE CODE (CREDITS): 22M11CI112 (3)

MAX. MARKS: 25

COURSE NAME: INTRODUCTION TO DATA SCIENCE

COURSE INSTRUCTORS: Dr Nancy Singla

MAX. TIME: 1 Hour 30 Min

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.

		CO	Marks
Q. No	Question		
Q1.	You are tasked with managing semi-structured data generated from	CO3	[2+5]
, `	an online survey platform. The data includes varying fields		
1	depending on the survey type, such as user responses, timestamps,		
	device info, and optional nested metadata.		
	(a) Explain which type of database would be most suitable for		
	storing this semi-structured data?		]
	a > 1		
	(b) Using Python and the pymongo library, write code snippets to:  1. Connect to a MongoDB database called "survey_db".		]
	2. Insert a sample document representing a user response,		
	including fields: user id survey_id, response, timestamp, and		
	nested device information (type and OS).	İ	
l l	3. Update the document for a specific user_id by adding a new		
	field 'reviewed' with value True.		
	4. Delete all documents where the device type is "tablet".		
00	Your team is building a search engine that indexes articles. The	CO3	[3]
Q2.	search engine is supposed to return articles based on keywords that		
	users input. Users might search for terms like "cars," "car," and		
	"automobiles". Would you apply stemming or lemmatization to		
	reduce the word variations in the search query? Justify your choice		
in,	and explain how it would improve the search engine results.		
Č2	For each of the following scenarios, identify the most suitable type	CO4	[5]
Q3.	of data visualization to effectively represent the data and briefly		
	explain your choice:		
	a) Visualizing monthly sales trends of multiple products over the		
	past year.		
	b) Comparing the distribution and outliers of customer ages in a		
	dataset.		
	c) Showing the geographic concentration of customer purchases		
	across different cities.		
	<u> </u>		

	d) Displaying the relationship between two numerical variables to identify correlation.							
		1-score)						
	of two c							
Q4.	A custon	per hour.	CO5	[3]				
	Each cal	est.		ניין				
	(a) Calcu	ulate the probability that	eived in					
1	an hour a	are technical support requ			M			
	(b) Appr	roximate the probability of	ort calls	11/4				
	occurring in an hour.  A dataset contains the following information about hours studied (x) and corresponding even scores (x) for 5 students:							
Q5.	A datase	t contains the following i	information about hours stu	idied (x)	CO6	3+2+2]		
	and corre							
		Hours Studied (x)	Exam scores (y)		<b>%</b>			
		1	50					
		2	55					
		3	65	~1000V				
	1	4	70					
i		5	75					
		-			1			
	(a) Using the least squares method, find the linear regression							
	equation that predicts exam scores based on bours studied.							
		_						
	(b) Using your regression equation, calculate the predicted scores for				ļ			
	each student and compute the Root Mean Squared Error (RMSE) of				•			
	your model.							
	(a) E1-							
		(c) Explain the difference between Mean Squared Error (MSE) and						
		Root Mean Squared Error (RMSE). Which metric is generally preferred for evaluating regression models, and why?						
	brereued							
_ <del></del>	All .							