

	<b>PES University, Bengaluru</b> (Established under Karnataka Act No. 16 of 2013)		<b>UE20CS905</b>
	<b>: END SEMESTER ASSESSMENT (ESA)</b> <b>M TECH DATA SCIENCE AND MACHINE LEARNING_ SEMESTER I</b> <b>UE20CS905 - MACHINE LEARNING - I</b>		
Time: 3 Hrs	Answer All Questions		Max Marks: 100
<b>Instructions</b> <ol style="list-style-type: none"> <li>1. Answer all the questions.</li> <li>2. Section A should be handwritten in the answer script provided and signed at the end of the same.</li> <li>3. Section B and C are coding questions which have to be answered in the system and uploaded in Olympus Login.</li> <li>4. Avoid using GridSearchCV as it might impact the system performance drastically.</li> </ol>			

Section A (20 marks)			
1	a)	What is Machine learning? State any two types of machine learning.	4
	b)	How can you handle overfitting and underfitting?	4
	c)	State the assumptions of linear regression algorithm.	4
	d)	If $y = 2x_1 + 12x_2 + 3x_3 + 5$ is the linear regression equation, then explain how the coefficients of $x_1$ and $x_2$ affect the value of $y$ .	4
	e)	Explain any two of the data preprocessing steps.	4
Section B (40 Marks)			
2	a)	<b>DATA DESCRIPTION:</b>  Since 2008, guests and hosts have used Airbnb to expand on traveling possibilities and present more unique, personalized way of experiencing the world. This dataset describes the listing activity and metrics in NYC, NY for 2019. <ul style="list-style-type: none"> <li>● id = listing ID</li> <li>● name = name of the listing</li> <li>● host_id = host ID</li> <li>● host_name = name of the host</li> <li>● neighbourhood_group = location</li> <li>● neighbourhood = area</li> <li>● latitude = latitude coordinates</li> <li>● longitude = longitude coordinates</li> <li>● room_type = listing space type</li> <li>● price = price in dollars</li> <li>● minimum_nights = amount of nights minimum</li> <li>● number_of_reviews = number of reviews</li> <li>● last_review = latest review</li> </ul>	25

		<ul style="list-style-type: none"> <li>● reviews_per_month = number of reviews per month</li> <li>● calculated_host_listings_count = amount of listing per host</li> <li>● availability_365 = number of days when listing is available for booking</li> </ul> <p>Perform Exploratory data analysis and summarize important observations from the data set</p> <p>Some pointers which would help you, but don't be limited by these</p> <p>i. Explore the types of variables (continuous, categorical etc.) (3 marks)</p> <p>ii. Calculate five point summary for numerical variables (2 marks)</p> <p>iii. Summarize observations for categorical variables – no. of categories, % observations in each category (5 marks)</p> <p>iv. Check for defects in the data. Perform necessary actions to 'fix' these defects. (10 marks)</p> <p>check for numerical/ categorical variable( wrong representation) missing values, outlier treatment, skewness, encoding</p> <p>v. Summarize relationships among variables. (5 marks)</p>	
	b)	<p>Fit a base model. Please write your key observations</p> <p>i. What is the overall R<sup>2</sup>? Please comment on whether it is good or not. ( 2 marks)</p> <p>ii. What is the adjusted R<sup>2</sup>? Is it different from R<sup>2</sup>? Why? ( 3 marks)</p> <p>iii. Which variables are significant? (4 marks)</p> <p>iv. Is there multicollinearity? (4 marks)</p> <p>v. Which other key model output parameters do you want to look at? ( 2 marks)</p>	15
<b>Section C (40 marks)</b>			
3	a)	<p>How do you improve the accuracy of the model? Write clearly the changes that you will make before re-fitting the model. Fit the final model.</p> <p>Feature Engineering / Feature Selection Regularization Cross Validation</p> <p>Please feel free to have any number of iterations to get to the final answer. Marks are awarded based on the quality of the final model you are able to achieve.</p>	20
	b)	<p>Summarize as follows</p> <p>i. Summarize the overall fit of the model and list down the measures to prove that it is a good model</p> <p>ii. Write down a business interpretation/explanation of the model – which variables are affecting the target the most and explain the</p>	20

SRN

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

		<p>relationship. Feel free to use charts or graphs to explain.</p> <p>iii. What changes from the base model had the most effect on model performance</p> <p>iv.What are the key risks to your results and interpretation?</p> <p>Justification for selecting a model</p>	
--	--	--	--