



SOMAIYA
VIDYAVIHAR UNIVERSITY

<p align="center">Semester: January 2022 – May 2022 Examination: In-Semester Examination</p>		
<p>Programme code: 54 Programme: B.Tech. Computer Engineering (Honours in Data Science and analytics)</p>	<p>Class: SY</p>	<p>Semester: IV (SVU 2020)</p>
<p>Name of the Constituent College: K. J. Somaiya College of Engineering</p>		<p>Name of the department: COMP ENGG</p>
<p>Course Code: 116h54C401</p>	<p>Name of the Course: Applied Data Science</p>	

Question No.		Max. Marks	CO Mapped
Q1	<p>Attempt any FIVE(5) of following</p> <p>A. What are the characteristic of big data?</p> <p>B. List out skillset required for data science profile.</p> <p>C. What is need of data scrapping?</p> <p>D. List types of the exploratory data analysis</p> <p>E. Give an example of binomial distribution</p> <p>F. What is data wrangling?</p>	2 mark Each	CO1 and CO3
Q2	<p>Attempt any TWO (2) of the following</p> <p>a) Illustrate with suitable example how data science can be used to add value to business.</p> <p>b) A company packages salted peanuts in 200gm packets using the machine. A sample of 16 packets is taken from the production line at random time intervals and their contents are weighted. The mean weight of 16 packets is found as 199.5. Can we conclude that machine is working properly using statistical hypothesis testing? In the hypothesis testing clearly show all the 5 steps in the calculation, considering the significance level of 5%. (Assume value of test statistic=2.131)</p> <p>c) Are the following nominal, ordinal, interval or ratio data? Explain your answers.</p> <p>(a) Temperatures measured on the Kelvin scale.</p> <p>(b) Military ranks.</p> <p>(c) Social security numbers.</p> <p>(d) Number of passengers on buses from Delhi to Mumbai.</p> <p>(e) Code numbers given to the religion of persons married.</p>	<p>05 marks</p> <p>05 marks</p> <p>05 marks</p>	<p>CO1</p> <p>CO3</p> <p>CO1</p>
Q3	<p>a) You have been given the task to perform the data preprocessing of the data retrieved from multiple sources, before you start applying the data mining task. Identify (atleast 5) data quality issues with the sample dataset retrieved from the master data set. Suggest how do you resolve the quality issue.</p>	05 Marks	CO3

TXN-ID	NAME	AGE	HEIGHT	WEIGHT	BLOOD GROUP	COVID-19 RESULT
T001	RAMA	45	145	62kg	O+ve	Positive
T002	SEETHA	43	168	45kg	B+ve	Negative
T003	Akbar	38	172	60kg	Iam+ve	Positive
T004	BIRBAL	45	168	52kg	AB+ve	Negative
T005	THenali	22	157	78kg	B-ve	1
T006	Venkat	36	157	54kg	O-ve	Negative
T007	Rajuu	350	132	48kg	O+ve	Positive
T008	HARI	32	180	120lbs	AB-ve	Negative
T009	Inba	25		85kg	O+ve	0
T010	SysUsr789	20	165	68kg	O-ve	Negative

OR

Eleven students were asked to measure their pulses for 30 seconds and multiply by two to get their one minute pulse rate. The results were: 62,32,60,66,70,72,74,74,78,80,84. Give the five-number summary for the pulse rates and draw the boxplot

05 Marks

CO1

- b) Use min-max normalization method to normalize the following group of data by setting min=0 and max=1:
200, 300, 400, 600,1000

03 Marks

CO3

- c) What is type I and Type II error? Provide example of each.

02 Marks

CO3