

Name:	 <b>UPES</b> UNIVERSITY WITH A PURPOSE
Enrolment No:	

**UNIVERSITY OF PETROLEUM AND ENERGY STUDIES**  
**End Semester Examination, Dec 2021**

<b>Course: Statistics for Data Science</b> <b>Program: B.Tech CSE-SPZ-BD</b> <b>Course Code: CSBD3006P</b>	<b>Semester: V</b> <b>Time: 03 hrs.</b> <b>Max. Marks: 100</b>
--	--

**Instructions:**

**SECTION A**

1. Each Question will carry 4 Marks
2. Instruction: Write short answers for the following questions. (60-70 words)

S. No.		Marks	CO
Q1	Explain confounding variables and elaborate their role in correlation?	4	CO4
Q2	Mean, median and mode are three primary measures of central tendency. Explain and illustrate the effect of outliers on mean and median?	4	CO1
Q3	Explain the relation between logarithmic and exponential functions?	4	CO3
Q4	Discuss the major data types (levels of measurements) used in statistics.	2	CO2
	Explain briefly the preciousness of these data types?	2	
Q5	Discuss the differences between dependent and independent variables.	3	CO5
	Also explain the meaning of lurking variables?	1	

**SECTION B**

1. Each question will carry 10 marks.
2. Instruction: Write short / brief notes (100-150 words)
3. For question 6 choose between part a and b.
4. Attempt any one question for question 6
5. There is no such option for other questions in this section

Q6	a. Explain the basics steps in a research process in detail with a suitable example.	10	CO3
	<b>OR</b>		
	b. Serum Institute of India specializing in vaccine states that its Covishield vaccine failure rate is not more than 1%. You perform a hypothesis test to determine whether		

**Please Turn Over**

	<p>the company's claim is false. When will a type I or type II error occur? Which is more serious?</p> <ul style="list-style-type: none"> <li>• State the null and alternative hypotheses.</li> <li>• Write the possible type I and type II errors.</li> <li>• Determine which error is more serious.</li> </ul>	3																							
		3																							
		4																							
Q7	<p>A group of UPES students were given a short course in speed-reading. The instructor was curious if a monetary incentive would influence performance on a reading test taken at the end of the course.</p> <p>Half of the students were offered Rs 500 for obtaining a certain level of performance on the test, the other half were not offered money.</p> <p>Identify the independent variable and dependent variable in the study.</p>	10	CO1																						
Q8	<p>Calculate the correlation coefficient for the gross domestic products and Carbon dioxide emissions data given in the table below.</p> <table border="1" data-bbox="203 898 1172 1318"> <thead> <tr> <th>GDP (Trillions of \$), x</th> <th>CO2 Emission (Millions of Metric tons), y</th> </tr> </thead> <tbody> <tr><td>1.6</td><td>428.2</td></tr> <tr><td>3.6</td><td>828.8</td></tr> <tr><td>4.9</td><td>1214.2</td></tr> <tr><td>1.1</td><td>444.6</td></tr> <tr><td>0.9</td><td>264.0</td></tr> <tr><td>2.9</td><td>415.3</td></tr> <tr><td>2.7</td><td>571.8</td></tr> <tr><td>2.3</td><td>454.9</td></tr> <tr><td>1.6</td><td>358.7</td></tr> <tr><td>1.5</td><td>573.5</td></tr> </tbody> </table> <p>Also display the data in a scatter plot and determine whether there appears to be a positive or negative linear correlation.</p>	GDP (Trillions of \$), x	CO2 Emission (Millions of Metric tons), y	1.6	428.2	3.6	828.8	4.9	1214.2	1.1	444.6	0.9	264.0	2.9	415.3	2.7	571.8	2.3	454.9	1.6	358.7	1.5	573.5	5	CO4
GDP (Trillions of \$), x	CO2 Emission (Millions of Metric tons), y																								
1.6	428.2																								
3.6	828.8																								
4.9	1214.2																								
1.1	444.6																								
0.9	264.0																								
2.9	415.3																								
2.7	571.8																								
2.3	454.9																								
1.6	358.7																								
1.5	573.5																								
		5																							
Q9	<p>Elaborate the role of clustering in data analytics?</p> <p>Explain the types of data used in cluster analysis.</p>	3	CO5																						
		7																							

**SECTION-C**

1. Each Question carries 20 Marks.
2. Instruction: Write long answer. (Up to 350 words while explaining)
3. For question 10 choose between part a and b
4. Attempt any one question for question 10.
5. There is no such option for other question 11.

**Please Turn Over**

Q10.	<p>a. Find the mean, the median, and the mode of the sample ages of students in a class shown at the left. Which measure of central tendency best describes a typical entry of this data set? Are there any outliers?</p> <table border="1" style="margin-left: auto; margin-right: auto;"> <tr> <td style="padding: 2px;">Age of students in a class</td> <td style="padding: 2px;">20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 22, 22, 22, 23, 23, 23, 23, 24, 24, 65</td> </tr> </table> <ul style="list-style-type: none"> <li>• With the help of histogram display the distributions of data along with locations of mean, median, mode.</li> <li>• Remove the data entry 65 from the data set and then calculate the mean, median and the mode. Does the absence of the outlier change the measures? If yes, justify your answer.</li> </ul>	Age of students in a class	20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 22, 22, 22, 23, 23, 23, 23, 24, 24, 65	<b>8</b>	<b>CO2</b>
	Age of students in a class	20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 22, 22, 22, 23, 23, 23, 23, 24, 24, 65			
<p style="text-align: center;"><b>OR</b></p> <p>b. The following sample data set lists the per hour salaries of 20 employees. Construct the five-point summary of the sample data set The sample data is as follows: 30, 32, 110, 65, 55, 50, 55, 48, 43, 42, 45, 45, 33, 34, 38, 38, 34, 38, 33, 32 Analyze the data set and apply the five-point summary methods to calculate the following:</p> <ul style="list-style-type: none"> <li>• First quartile, Second quartile, Third quartile</li> <li>• Range</li> <li>• Interquartile range</li> <li>• Semi interquartile range</li> <li>• Detect the outliers if any</li> <li>• Construct a box plot for displaying the five-point summary</li> </ul>	<b>3</b>	<b>3</b>			
Q11.	<p>The regressions line always passes through the averages of data points. The formula for calculating the slop of the regression line is <math>m = r (s_y/s_x)</math> where m is the slope, r is the regression coefficient, <math>s_y</math> and <math>s_x</math> are standard deviations.</p>				
	<p>a. Using this information prove that the slope <math>m = \frac{\sum(y_i - y\text{-mean})}{\sum(x_i - x\text{-mean})}</math> where <math>y_i</math> and <math>x_i</math> are ith observation for x and y respectively.</p>	<b>12</b>	<b>CO4/ CO5</b>		
	<p>b. Use the data set in Question 8 to experimentally evaluate that <math>r (s_y/s_x) = \frac{\sum(y_i - y\text{-mean})}{\sum(x_i - x\text{-mean})}</math>.</p>	<b>8</b>			

**Please Turn Over**