

[illegible]

B.E. /B.Tech / B. Arch (Full Time) - END SEMESTER EXAMINATIONS, NOV / DEC 2024

VII Semester

(Regulation 2019)

Time:3hrs

Max.Marks: 100

CO1	Identify the real world business problems and model with analytics solutions.
CO2	Solve analytics problem with relevant mathematics background knowledge.
CO3	Convert any real world decision making problem to hypothesis and apply suitable statistical testing.
CO4	Write and demonstrate simple applications involving analytics using Hadoop and MapReduce.
CO5	Use open source frameworks for modeling and storing data.
CO6	Perform data analytics and visualization using Python.

(L1-Remembering, L2-Understanding, L3-Appling, L4-Analysing, L5-Evaluating, L6-Creating)

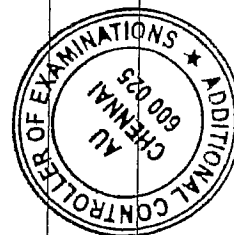
(Answer all Questions)

Q.No.	Questions	Marks	CO	BL
1	What is Big Data Analytics?	2	CO1	L1
2	List the key questions to be answered by all organization stepping into analytics.	2	CO1	L2
3	For the given univariate dataset $s=\{5,10,15,20,25,30\}$ of marks. Find mean, median, mode and standard deviation.	2	CO2	L2
4	The mean salary of the 8 people who work for a small company is 15000 rupees. When an extra worker is taken on this mean drop to 14000 rupees. How much does the new worker earn?	2	CO2	L2
5	Enumerate the drawbacks of regression analysis.	2	CO3	L1
6	Write the MongoDB syntax to find the number of documents in the student collection.	2	CO4	L1
7	List the different types of NoSQL Databases.	2	CO5	L1
8	Why NoSQL is preferred over SQL.	2	CO5	L2
9	How many rows by default the head () command retrieves?	2	CO6	L1
10	Write the correct syntax to access the row in the reverse order.	2	CO6	L1

(Restrict to a maximum of 2 subdivisions)

Q.No.	Questions	Marks	CO	BL
11 (a) i	Define big data and discuss the challenges in handling big data.	5	CO1	L4
ii	Discuss the various types of digital data, their sources, the challenges associated with handling them, and the advantages and disadvantages of utilizing them.	8	CO1	L3
OR				
11 (b) i	Define CAP theorem and discuss its implications on distributed system.	5	CO1	L4

ii	Explain the following technologies used in modern data processing and analytics: in-memory analytics, in-database processing, massively parallel processing and shared nothing architecture.	8	CO1	L3																																																																																										
12 (a)	<p>Using Navie Bayes classifier predict whether the person Buys Computer or not whose instance are: age = ≤ 30, income = Medium, Student = Yes, Credit-Rating = Fair based on the following observation.</p> <table border="1"> <thead> <tr> <th>S.No</th> <th>Age</th> <th>Income</th> <th>Student</th> <th>Credit - Rating</th> <th>Buys Computer</th> </tr> </thead> <tbody> <tr><td>1</td><td>≤ 30</td><td>High</td><td>No</td><td>Fair</td><td>No</td></tr> <tr><td>2</td><td>≤ 30</td><td>High</td><td>No</td><td>Excellent</td><td>No</td></tr> <tr><td>3</td><td>31-40</td><td>High</td><td>No</td><td>Fair</td><td>Yes</td></tr> <tr><td>4</td><td>> 40</td><td>Medium</td><td>No</td><td>Fair</td><td>Yes</td></tr> <tr><td>5</td><td>> 40</td><td>Low</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>6</td><td>> 40</td><td>Low</td><td>Yes</td><td>Excellent</td><td>No</td></tr> <tr><td>7</td><td>31-40</td><td>Low</td><td>Yes</td><td>Excellent</td><td>Yes</td></tr> <tr><td>8</td><td>≤ 30</td><td>Medium</td><td>No</td><td>Fair</td><td>No</td></tr> <tr><td>9</td><td>≤ 30</td><td>Low</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>10</td><td>> 40</td><td>Medium</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>11</td><td>≤ 30</td><td>Medium</td><td>Yes</td><td>Excellent</td><td>Yes</td></tr> <tr><td>12</td><td>31-40</td><td>Medium</td><td>No</td><td>Excellent</td><td>Yes</td></tr> <tr><td>13</td><td>31-40</td><td>High</td><td>Yes</td><td>Fair</td><td>Yes</td></tr> <tr><td>14</td><td>> 40</td><td>Medium</td><td>No</td><td>Excellent</td><td>No</td></tr> </tbody> </table>	S.No	Age	Income	Student	Credit - Rating	Buys Computer	1	≤ 30	High	No	Fair	No	2	≤ 30	High	No	Excellent	No	3	31-40	High	No	Fair	Yes	4	> 40	Medium	No	Fair	Yes	5	> 40	Low	Yes	Fair	Yes	6	> 40	Low	Yes	Excellent	No	7	31-40	Low	Yes	Excellent	Yes	8	≤ 30	Medium	No	Fair	No	9	≤ 30	Low	Yes	Fair	Yes	10	> 40	Medium	Yes	Fair	Yes	11	≤ 30	Medium	Yes	Excellent	Yes	12	31-40	Medium	No	Excellent	Yes	13	31-40	High	Yes	Fair	Yes	14	> 40	Medium	No	Excellent	No	13	CO2	L4
S.No	Age	Income	Student	Credit - Rating	Buys Computer																																																																																									
1	≤ 30	High	No	Fair	No																																																																																									
2	≤ 30	High	No	Excellent	No																																																																																									
3	31-40	High	No	Fair	Yes																																																																																									
4	> 40	Medium	No	Fair	Yes																																																																																									
5	> 40	Low	Yes	Fair	Yes																																																																																									
6	> 40	Low	Yes	Excellent	No																																																																																									
7	31-40	Low	Yes	Excellent	Yes																																																																																									
8	≤ 30	Medium	No	Fair	No																																																																																									
9	≤ 30	Low	Yes	Fair	Yes																																																																																									
10	> 40	Medium	Yes	Fair	Yes																																																																																									
11	≤ 30	Medium	Yes	Excellent	Yes																																																																																									
12	31-40	Medium	No	Excellent	Yes																																																																																									
13	31-40	High	Yes	Fair	Yes																																																																																									
14	> 40	Medium	No	Excellent	No																																																																																									
OR																																																																																														
12 (b)	<p>The data given is related to gram plant dry weight Y, soil organic matter X1, and kilograms of supplemental soil nitrogen added per 100 square meters X2. Predict the plant dry weight given soil organic matter of 5 and soil nitrogen 4</p> <table border="1"> <thead> <tr> <th>Y</th> <th>X1</th> <th>X2</th> </tr> </thead> <tbody> <tr><td>78.5</td><td>7</td><td>2.6</td></tr> <tr><td>74.3</td><td>1</td><td>2.9</td></tr> <tr><td>104.3</td><td>11</td><td>5.6</td></tr> <tr><td>87.6</td><td>11</td><td>3.1</td></tr> <tr><td>95.9</td><td>7</td><td>5.2</td></tr> <tr><td>109.2</td><td>11</td><td>5.5</td></tr> <tr><td>102.7</td><td>3</td><td>7.1</td></tr> </tbody> </table>	Y	X1	X2	78.5	7	2.6	74.3	1	2.9	104.3	11	5.6	87.6	11	3.1	95.9	7	5.2	109.2	11	5.5	102.7	3	7.1	13	CO2	L4																																																																		
Y	X1	X2																																																																																												
78.5	7	2.6																																																																																												
74.3	1	2.9																																																																																												
104.3	11	5.6																																																																																												
87.6	11	3.1																																																																																												
95.9	7	5.2																																																																																												
109.2	11	5.5																																																																																												
102.7	3	7.1																																																																																												
13 (a)	Draw the hyperplane that separates the points into positive and negative class using SVM classifier. Points belongs to positive class are (4,1), (4,-1) and (5,0) and points belonging to negative class are (2,0), (0,2) and (0,-2).	13	CO3	L3																																																																																										
OR																																																																																														
13 (b)	<p>Apply K-Means Clustering Algorithm to cluster the data points:</p> <table border="1"> <thead> <tr> <th>Data points</th> <th></th> <th></th> </tr> </thead> <tbody> <tr><td>A1</td><td>185</td><td>72</td></tr> <tr><td>A2</td><td>170</td><td>56</td></tr> <tr><td>A3</td><td>168</td><td>60</td></tr> <tr><td>A4</td><td>179</td><td>68</td></tr> <tr><td>A5</td><td>182</td><td>72</td></tr> <tr><td>B6</td><td>188</td><td>77</td></tr> </tbody> </table>	Data points			A1	185	72	A2	170	56	A3	168	60	A4	179	68	A5	182	72	B6	188	77	13	CO3	L3																																																																					
Data points																																																																																														
A1	185	72																																																																																												
A2	170	56																																																																																												
A3	168	60																																																																																												
A4	179	68																																																																																												
A5	182	72																																																																																												
B6	188	77																																																																																												



	B7	180	71				
	B8	183	84				
	B9	181	88				
	B10	177	76				
	into 2 clusters C1 and C2. The initial centroids of are C1= (185,72) and C2 = (170,56).						
14 (a)	Write in brief about Hadoop Ecosystem, including its core components, associated tools, features, applications, advantages and challenges.				13	CO4	L4
OR							
14 (b)	Explain how MapReduce is used in Matrix-Vector Multiplication, Relational-Algebra Operation, Union, Intersection and Difference Operation				13	CO4	L4
15 (a)	Explain in details about Properties, Indexing and Slicing Operation, Arithmetic Operation of NumPy.				13	CO6	L3
OR							
15 (b)	Explain in details about Data Visualization using Pandas and Matplotlib.				13	CO6	L3

PART- C(1x 15=15Marks)
(Q.No.16 is compulsory)

Q.No.	Questions	Marks	CO	BL																																	
16.	Using PCA reduce the following two-dimension data into a single dimension data. <table><tr><td>Feat ures</td><td>Ex1</td><td>Ex2</td><td>Ex3</td><td>Ex4</td><td>Ex5</td><td>Ex6</td><td>Ex7</td><td>Ex8</td><td>Ex9</td><td>Ex10</td></tr><tr><td>X</td><td>2.5</td><td>0.5</td><td>2.2</td><td>1.9</td><td>3.1</td><td>2.3</td><td>2</td><td>1</td><td>1.5</td><td>1.1</td></tr><tr><td>Y</td><td>2.4</td><td>0.7</td><td>2.9</td><td>2.2</td><td>3.0</td><td>2.7</td><td>1.6</td><td>1.1</td><td>1.6</td><td>0.9</td></tr></table>	Feat ures	Ex1	Ex2	Ex3	Ex4	Ex5	Ex6	Ex7	Ex8	Ex9	Ex10	X	2.5	0.5	2.2	1.9	3.1	2.3	2	1	1.5	1.1	Y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9	15	CO5	L5
Feat ures	Ex1	Ex2	Ex3	Ex4	Ex5	Ex6	Ex7	Ex8	Ex9	Ex10																											
X	2.5	0.5	2.2	1.9	3.1	2.3	2	1	1.5	1.1																											
Y	2.4	0.7	2.9	2.2	3.0	2.7	1.6	1.1	1.6	0.9																											

