



**SIDDHARTH INSTITUTE OF ENGINEERING & TECHNOLOGY:: PUTTUR  
(AUTONOMOUS)**

Siddharth Nagar, Narayanavanam Road – 517583

**QUESTION BANK (DESCRIPTIVE)**

**Subject with Code: (20CS1101) INTRODUCTION TO DATA SCIENCE**

**Course & Branch: B.Tech – CSE(CAD) Regulation: R20**

**UNIT –I**

**INTRODUCTION TO DATA SCIENCE**

<b>1</b>	<b>a</b>	Define Data Science and discuss Benefits and uses of data science.	[L1][CO1]	[6M]
	<b>b</b>	Discuss the Various Processing Steps in Data Science	[L2][CO1]	[6M]
<b>2</b>		Explain in Details various data types used in Data science and Big data	[L2][CO1]	[12M]
<b>3</b>	<b>a</b>	Analyze the term: Distributed file systems	[L4][CO1]	[6M]
	<b>b</b>	How will you creating research goals in a project charter	[L1][CO1]	[6M]
<b>4</b>		Classify the term big data ecosystem	[L4][CO1]	[12M]
<b>5</b>		How will you retrieve the required data from data science	[L1][CO5]	[12M]
<b>6</b>		Discuss in detailed Data Cleaning operation in data science	[L2][CO1]	[12M]
<b>7</b>	<b>a</b>	What are various steps involved in integrating phase	[L1][CO1]	[6M]
	<b>b</b>	What is meant by exploratory data analysis	[L1][CO1]	[6M]
<b>8</b>		Examine the term: Transforming data in Data science	[L3][CO1]	[12M]
<b>9</b>	<b>a</b>	Show the various components of model building.	[L2][CO1]	[6M]
	<b>b</b>	What are the ways analyzed the data and built a well-performing model	[L2][CO1]	[6M]
<b>10</b>	<b>a</b>	How will you handling missing data in data science	[L2][CO1]	[6M]
	<b>b</b>	Examine K-nearest neighbor techniques look at the k-nearest point to make a prediction	[L4][CO1]	[6M]

**UNIT –II**  
**STATISTICAL METHODS FOR EVALUATION&ASSOCIATION RULES**

1	a	Define Hypothesis Testing	[L1][CO2]	[6M]
	b	How will you mathematically define Confidence	[L1][CO2]	[6M]
2	a	Differentiate Null Hypotheses and Alternative Hypotheses	[L4][CO2]	[6M]
	b	Examine the application property of Wilcoxon rank-sum test	[L3][CO2]	[6M]
3		Discriminate about Difference of Means	[L5][CO2]	[12M]
4		Explain the differences between BI and Data Science.	[L2][CO2]	[12M]
5		<i>Explain the following</i> a) Student's t-test b) Welch's t-test	[L2][CO2]	[12M]
6	a	What are the three characteristics of Big Data, and what are the main considerations in processing Big Data?	[L1][CO2]	[6M]
	b	How evaluation of Candidate Rules is done?	[L2][CO2]	[6M]
7	a	What is a type I error? What is a type II error? Is one always more serious than the other? Why?	[L1][CO2]	[6M]
	b	Give the difference between Validation and Testing	[L4][CO2]	[6M]
8	a	State Apriori Algorithm	[L1][CO2]	[4M]
	b	Explain Apriori Algorithm with example	[L2][CO2]	[8M]
9	a	List and discuss the four measures of significance of Association rules	[L1][CO2]	[6M]
	b	Give the Applications of Association Rules	[L1][CO2]	[6M]
10		Illustrate any five approaches to improve Apriori's efficiency when the dataset is large.	[L3][CO2]	[12M]

**UNIT –III**  
**REGRESSION & CLASSIFICATION**

<b>1</b>	<b>a</b>	Which two basic measures does the entropy methods select the most informative attribute?	[L1][CO3]	[6M]
	<b>b</b>	Define confusion matrix	[L1][CO3]	[6M]
<b>2</b>		Explain the analytical technique Linear Regression with its model description.	[L2][CO3]	[12M]
<b>3</b>		Discuss the following with respect to linear regression a) Categorical Variables b) Confidence Intervals on the Parameters c) Confidence Interval on the Expected Outcome d) Prediction Interval on a Particular Outcome	[L2][CO3]	[12M]
<b>4</b>	<b>a</b>	Justify the usage of linear regression and logistic regression.	[L6][CO3]	[4M]
	<b>b</b>	Illustrate Logistic Regression Model.	[L3][CO3]	[8M]
<b>5</b>	<b>a</b>	Describe Decision Trees in detail with example.	[L2][CO3]	[6M]
	<b>b</b>	Difference between Alternative hypothesis and null hypothesis	[L2][CO4]	[6M]
<b>6</b>		Intercept the decision trees algorithms	[L4][CO4]	[12M]
<b>7</b>	<b>a</b>	State Bayes' Theorem	[L1][CO4]	[4M]
	<b>b</b>	Discuss Naïve Bayes classification method considering an example	[L2][CO4]	[8M]
<b>8</b>		How does one pick the most suitable method for a given classification problem?	[L2][CO4]	[12M]
<b>9</b>	<b>a</b>	Compare the C4.5 and CART algorithm of decision tree.	[L4][CO4]	[4M]
	<b>b</b>	Discriminate the way show the evaluation of decision tree is done	[L5][CO4]	[4M]
	<b>c</b>	Give the two approaches that help avoid over fitting in decision tree learning.	[L2][CO4]	[4M]
<b>10</b>		Discuss the following term: a) Accuracy b) TPR c) FPR d) FNR e) Precision	[L4][CO4]	[12M]

**UNIT –IV**  
**CLUSTERING & TIME SERIES ANALYSIS**

<b>1</b>	<b>a</b>	What is clustering?	[L1][CO5]	[6M]
	<b>b</b>	State the advantage of using PAM.	[L1][CO5]	[6M]
<b>2</b>		Illustrate the method to find k clusters from a collection of M objects with n attributes.	[L3][CO5]	[12M]
<b>3</b>	<b>a</b>	Explain any one case study for time series analysis	[L2][CO5]	[6M]
	<b>b</b>	What is forecasting in association with time series. Explain	[L1][CO6]	[6M]
<b>4</b>	<b>a</b>	Indicate when the time series $y_t$ for $t=1,2,3,\dots$ is said to be stationary time series.	[L2][CO6]	[6M]
	<b>b</b>	Express the stationary time series conditions in detail.	[L6][CO6]	[6M]
<b>5</b>		Discussion detail each part of the ARIMA model	[L2][CO5]	[12M]
<b>6</b>	<b>a</b>	List and explain time series components	[L1][CO6]	[6M]
	<b>b</b>	Discriminate the steps involved in Box-Jenkins Methodology	[L5][CO6]	[6M]
<b>7</b>	<b>a</b>	What is meant by k-means	[L1][CO5]	[4M]
	<b>b</b>	Describe k-means algorithm to find k clusters	[L2][CO5]	[8M]
<b>8</b>		Correlate ARMA and ARIMA Models	[L4][CO6]	[12M]
<b>9</b>		Express the following	[L2][CO6]	[12M]
		a) Autocorrelation Function b) Autoregressive Models		
<b>10</b>		List and describe Additional time series methods	[L2][CO6]	[12M]

**UNIT –V TEXT ANALYSIS**

<b>1</b>	<b>a</b>	Define Porter's stemming algorithm.	[L1][CO6]	[6M]
	<b>b</b>	What is Topic modeling?	[L1][CO6]	[6M]
<b>2</b>		Explain the three important steps of the text analysis	[L2][CO6]	[12M]
<b>3</b>	<b>a</b>	Sketch the flow diagram of Text analysis process	[L5][CO6]	[6M]
	<b>b</b>	Illustrate in detail the steps involved in the process of Text Analysis done by organizations	[L3][CO6]	[6M]
<b>4</b>	<b>a</b>	Define TFIDF.	[L1][CO6]	[4M]
	<b>b</b>	Describe the usage of TFIDF to compute the usefulness of each word in the texts.	[L2][CO6]	[8M]
<b>5</b>		Explain how the data science team will categorize the reviews by topics	[L2][CO6]	[12M]
<b>6</b>		Illustrate the main challenges of text analysis	[L3][CO6]	[12M]
<b>7</b>	<b>a</b>	Define Topic model. Describe LDA.	[L2][CO6]	[6M]
	<b>b</b>	Justify the process of topic modeling simplification.	[L6][CO6]	[6M]
<b>8</b>		Explain the following a) Tokenization b) Case folding	[L3][CO6]	[12M]
<b>9</b>	<b>a</b>	Explain how categorizing documents by topics is done.	[L2][CO6]	[6M]
	<b>b</b>	Interpret the procedure used in data science to gain insights into customer opinions	[L3][CO6]	[6M]
<b>10</b>	<b>a</b>	What is meant by sentiment analysis	[L1][CO6]	[4M]
	<b>b</b>	Discriminate the methods used for sentiment analysis	[L5][CO6]	[8M]

**Preparedby:**  
**Mr.G.Prasad Babu**  
**Associate Professor**