**SRINIVAS UNIVERSITY**

**INSTITUTE OF ENGINEERING AND TECHNOLOGY**

**FIRST SEMESTER QUESTION BANK**

**COURSE CODE: BTA001**       **COURSE NAME: INTRODUCTION TO DATA SCIENCE**                **CREDITS:04**

| Q. No | Objective Type of Questions<br><br>(Single sentence answer questions, Multiple choice or fill in the blanks type)<br><br>There should be more than 50 questions for 5 Module course and<br>50 questions for 5 Module courses (Compulsory 10 questions from each module) | Module or Unit No. | Mark for each question (1mark) | CO- can be any one of the below CO1,CO2, CO3, CO4 |
|---|---|---|---|---|
| | **MODULE-1** | | | |
| 1 | Which of the following best defines Data Science?<br>A. Study of computer hardware<br>B. Extraction of knowledge and insights from structured and unstructured data<br>C. Only data visualization techniques<br>D. Database management system<br>Answer: B | 1 | 1 | CO1 |
| 2 | Which discipline is NOT a core contributor to Data Science?<br>A. Statistics | 1 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | B. Machine Learning<br>C. Philosophy<br>D. Computer Science<br>Answer: C | | | |
| 3 | The term "Data Science" was popularized during which decade?<br>A. 1970s<br>B. 1980s<br>C. 1990s<br>D. 2000s<br>Answer: D | **1** | **1** | **CO1** |
| 4 | Which component of Data Science focuses on discovering hidden patterns in data?<br>A. Data Collection<br>B. Data Storage<br>C. Data Analysis<br>D. Data Transmission<br>Answer: C | **1** | **1** | **CO1** |
| 5 | Which of the following is a key difference between Data Science and Business Analytics?<br>A. Data Science focuses only on past data<br>B. Business Analytics does not use data<br>C. Data Science emphasizes predictive and prescriptive modeling | **1** | **1** | **CO1** |

| | | | | |
|---|---|---|---|---|
| | D. Business Analytics ignores visualization<br><br>Answer: C | | | |
| 6 | **Which term refers to data that does not follow a predefined schema?**<br>A. Structured data<br>B. Semi-structured data<br>C. Unstructured data<br>D. Normalized data<br>**Answer:** C | 1 | 1 | CO1 |
| 7 | Which layer in Data Science architecture deals with data acquisition from various sources?<br><br>A. Application layer<br><br>B. Data ingestion layer<br><br>C. Analytics layer<br><br>D. Visualization layer<br><br>Answer: B | 1 | 1 | CO1 |
| 8 | Who typically occupies the top position in the Data Science hierarchy within an organization?<br><br>A. Data Analyst<br><br>B. Data Engineer<br><br>C. Chief Data Officer (CDO)<br><br>D. Business User<br><br>Answer: C | 1 | 1 | CO1 |
| 9 | Which of the following is a major challenge in Data Science projects?<br><br>A. Excessive computational power | 1 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | B. Poor data quality<br>C. Availability of tools<br>D. Too much documentation<br>**Answer:** B | | | |
| **10** | Which industry widely applies Data Science for fraud detection?<br>A. Agriculture<br>B. Healthcare<br>C. Banking and Finance<br>D. Education<br>**Answer:** C | **1** | **1** | **CO1** |
| | **MODULE 2** | | | |
| **1.** | What is the primary importance of Microsoft Excel?<br>A. Creating animations<br>B. Performing data analysis and calculations<br>C. Web browsing<br>D. Software development<br>**Answer:** B | **2** | **1** | **CO1** |
| **2** | Which feature in Excel is used to organize data in a structured format?<br>A. Pivot Table<br>B. Chart<br>C. Excel Table | **2** | **1** | **CO1** |

| | | | | |
|---|---|---|---|---|
| | D. Macro<br><br>Answer: C | | | |
| **3** | Which symbol is used for multiplication in Excel formulas?<br>A. ✕<br>B. #<br>C. *<br>D. %<br>Answer: C | **2** | **1** | **CO1** |
| **4** | To restrict the type of data entered into a cell, which Excel feature is used?<br>A. Conditional Formatting<br>B. Data Validation<br>C. Filtering<br>D. Sorting<br>Answer: B | **2** | **1** | **CO1** |
| **5** | Which option is used to arrange data in ascending or descending order?<br>A. Grouping<br>B. Subtotal<br>C. Sorting<br>D. Validation<br>Answer: C | **2** | **1** | **CO1** |
| **6** | Which Excel feature displays only records that meet specific criteria?<br>A. Sorting<br>B. Filtering<br>C. Subtotal | **2** | **1** | **CO1** |

| | D. Grouping<br>Answer: B | | | |
|---|---|---|---|---|
| **7** | Which symbol is used to start any formula in Excel?<br>A. #<br>B. $<br>C. =<br>D. @<br>Answer: C | **2** | **1** | **CO1** |
| **8** | Which logical function returns TRUE if all conditions are satisfied?<br>A. OR<br>B. NOT<br>C. IF<br>D. AND<br>Answer: D | **2** | **1** | **CO1** |
| **9** | Which chart type is best suited for showing trends over time?<br>A. Pie Chart<br>B. Column Chart<br>C. Line Chart<br>D. Bar Chart<br>Answer: C | **2** | **1** | **CO1** |
| **10** | Which file format is commonly used to import text data into Excel?<br>A. .xml | **2** | **1** | **CO1** |

| | | | | |
|---|---|---|---|---|
| | B. .xls<br>C. .csv<br>D. .accdb<br>Answer: C | | | |
| **MODULE 3** | | | | |
| 1 | Which type of machine learning uses labeled training data?<br>A. Unsupervised learning<br>B. Reinforcement learning<br>C. Supervised learning<br>D. Semi-supervised learning<br>Answer: C | 3 | 1 | CO1 |
| 2 | Which of the following is a classification algorithm?<br>A. K-Means<br>B. Decision Tree<br>C. Apriori<br>D. PCA<br>Answer: B | 3 | 1 | CO1 |
| 3 | Which algorithm is commonly used for clustering?<br>A. Logistic Regression<br>B. Naïve Bayes<br>C. K-Means | 3 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | D. Linear Regression<br><br>Answer: C | | | |
| 4 | Which technique is used for feature selection?<br><br>A. Gradient Descent<br><br>B. Principal Component Analysis (PCA)<br><br>C. K-Nearest Neighbor<br><br>D. K-Means<br><br>Answer: B | 3 | 1 | CO1 |
| 5 | Bayes' theorem is mainly used to calculate:<br>A. Mean<br>B. Conditional probability<br>C. Variance<br>D. Correlation<br>Answer: B | 3 | 1 | CO1 |
| 6 | In a Cartesian plane, how many axes are present?<br><br>A. One<br><br>B. Two<br><br>C. Three<br><br>D. Four<br><br>Answer: B | 3 | 1 | CO1 |
| 7 | Which of the following represents a linear equation?<br><br>A. $y = x^2$<br><br>B. $y = mx + c$<br><br>C. $y = \log x$ | 3 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | D.  y = e$^x$ <br><br> Answer: B | | | |
| **8** | Which SQL command is used to retrieve data from a database? <br><br> A.  INSERT <br> B.  UPDATE <br> C.  SELECT <br> D.  DELETE <br><br> Answer: C | **3** | **1** | **CO1** |
| **9** | Which SQL command category includes COMMIT and ROLLBACK? <br><br> A.  DDL <br> B.  DML <br> C.  TCL <br> D.  DCL <br><br> Answer: C | **3** | **1** | **CO1** |
| **10** | Which tool is commonly used for data science tasks? <br><br> A.  MS Paint <br> B.  Jupyter Notebook <br> C.  Notepad <br> D.  Calculator <br><br> Answer: B | **3** | **1** | **CO1** |
| **MODULE 4** | | | | |

| 1. | Correlation measures the:<br>A. Causation between two variables<br>B. Strength and direction of relationship between variables<br>C. Difference between variables<br>D. Distribution of data<br>Answer: B | 4 | 1 | CO1 |
|---|---|---|---|---|
| 2. | Which correlation coefficient value indicates a strong positive correlation?<br>A. $-0.9$<br>B. $-0.1$<br>C. 0<br>D. $+0.9$<br>Answer: D | 4 | 1 | CO1 |
| 3. | Which regression technique is used for binary classification problems?<br>A. Linear Regression<br>B. Polynomial Regression<br>C. Logistic Regression<br>D. Ridge Regression<br>Answer: C | 4 | 1 | CO1 |
| 4. | The Gaussian distribution is also known as:<br>A. Uniform distribution<br>B. Binomial distribution<br>C. Normal distribution | 4 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | D. Poisson distribution<br>Answer: C | | | |
| 5. | Standardization converts data to have:<br>A. Mean = 1 and Variance = 0<br>B. Mean = 0 and Standard Deviation = 1<br>C. Mean = 1 and SD = 1<br>D. Mean = 0 and Variance = 0<br>Answer: B | 4 | 1 | CO1 |
| 6. | Z-score represents:<br>A. Raw data value<br>B. Probability value<br>C. Number of standard deviations from the mean<br>D. Variance of data<br>Answer: C | 4 | 1 | CO1 |
| 7. | According to the Central Limit Theorem, the sampling distribution of the mean approaches:<br>A. Uniform distribution<br>B. Binomial distribution<br>C. Normal distribution<br>D. Exponential distribution<br>Answer: C | 4 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| 8. | Markowitz Portfolio Optimization primarily focuses on:<br>A. Maximizing risk<br>B. Minimizing return<br>C. Optimizing risk and return<br>D. Eliminating variance<br>Answer: C | 4 | 1 | CO1 |
| 9. | Standardizing x and y variables in linear regression helps to:<br>A. Increase data size<br>B. Simplify coefficient interpretation<br>C. Remove correlation<br>D. Increase error<br>Answer: B | 4 | 1 | CO1 |
| 10. | Information gain in linear regression helps in:<br>A. Increasing variance<br>B. Feature selection and model improvement<br>C. Reducing data size<br>D. Data visualization<br>Answer: B | 4 | 1 | CO1 |
| **MODULE 5** | | | | |
| 1. | Which visualization technique is best suited to show the relationship between two continuous variables?<br>A. Histogram | 5 | 1 | CO1 |

| | | | | |
|---|---|---|---|---|
| | B. Scatter plot<br>C. Pie chart<br>D. Map<br>**Answer:** B | | | |
| **2.** | Which measure represents the average value of continuous data?<br>A. Median<br>B. Mode<br>C. Mean<br>D. Range<br>**Answer:** C | **5** | **1** | **CO1** |
| **3.** | Standard deviation is used to measure:<br>A. Central tendency<br>B. Frequency<br>C. Dispersion of data<br>D. Percentage<br>**Answer:** C | **5** | **1** | **CO1** |
| **4.** | Which statistical measure is most appropriate for categorical data analysis?<br>A. Mean<br>B. Standard deviation<br>C. Frequency and percentage<br>D. Variance<br>**Answer:** C | **5** | **1** | **CO1** |
| **5.** | Which Python data type is immutable?<br>A. List | **5** | **1** | **CO1** |

| | B. Dictionary<br>C. Set<br>D. Tuple<br>Answer: D | | | |
|---|---|---|---|---|
| 6. | Which Python library is mainly used for numerical computations?<br>A. Pandas<br>B. NumPy<br>C. Matplotlib<br>D. Scikit-Learn<br>Answer: B | 5 | 1 | CO1 |
| 7. | Which library is primarily used for data visualization in Python?<br>A. NumPy<br>B. Pandas<br>C. Matplotlib<br>D. Scikit-Learn<br>Answer: C | 5 | 1 | CO1 |
| 8. | Which Python data structure stores data in key–value pairs?<br>A. List<br>B. Tuple<br>C. Set<br>D. Dictionary<br>Answer: D | 5 | 1 | CO1 |

| 9. | Which library is used for data manipulation and analysis using DataFrames?<br>A. NumPy<br>B. Pandas<br>C. Matplotlib<br>D. TensorFlow<br>**Answer:** B | 5 | 1 | CO1 |
|---|---|---|---|---|
| 10. | Which Python library is widely used for implementing machine learning algorithms?<br><br>A. Matplotlib<br><br>B. NumPy<br><br>C. Pandas<br><br>D. Scikit-Learn<br><br>Answer: D | 5 | 1 | CO1 |
| **Quest ion NO.** | | | | |
| colspan="5" MODULE 1 | | | | |
| 1. | Define Data Science and explain its scope in modern data-driven organizations. | 1 | 8 | CO1, PO1 |
| 2. | Describe the history and evolution of Data Science, highlighting key technological milestones. | 1 | 8 | CO1, PO1 |
| 3. | Explain the important terminologies used in Data Science such as big data, machine learning, data mining, and artificial intelligence. | 1 | 8 | CO2, PO1 |
| 4. | Discuss the basic framework and architecture of Data Science with suitable examples. | 1 | 8 | CO2, PO1 |
| 5. | Differentiate between Data Science and Business Analytics with respect to objectives, tools, and applications. | 1 | 8 | CO3, PO2 |

| | | | | |
|---|---|---|---|---|
| 6. | Explain the importance of Data Science in today's business world with real-world use cases. | 1 | 8 | CO3, PO2 |
| 7. | Describe the primary components of Data Science and explain the role of each component. | 1 | 8 | CO2, PO1 |
| 8 | Explain the different users of Data Science in an organization and discuss the Data Science hierarchy. | 1 | 8 | CO3,PO4 |
| 9 | Provide an overview of various Data Science techniques used for data analysis and decision-making. | 1 | 8 | CO4,PO4 |
| 10 | Discuss the challenges and opportunities of Data Science in business analytics and explain its industrial applications across different sectors. | | 8 | CO4, PO2 |
| **MODULE 2** | | | | |
| 1. | Explain the importance of Microsoft Excel in data analysis and business applications. | 2 | 8 | CO1,,PO5 |
| 2. | Describe the steps involved in creating and managing Excel tables. | 2 | 8 | CO2,PO5 |
| 3. | Explain how to perform addition, subtraction, multiplication, and division in Excel with examples. | 2 | 8 | CO2,PO5 |
| 4. | Discuss Excel Data Validation and its role in maintaining data accuracy. | 2 | 8 | CO2,PO5 |
| 5. | Explain sorting, filtering, grouping, ungrouping, and subtotal operations in Excel. | 2 | 8 | CO2,PO5 |
| 6. | Introduce formulas and functions in Excel and explain their significance. | 2 | 8 | CO3,PO5 |
| 7. | Explain logical operators and conditional functions used in Excel with suitable examples. | 2 | 8 | CO4,PO5 |
| 8. | Describe different types of charts in Excel and explain how they help in data visualization. | 2 | 8 | CO4,PO5 |
| 9 | Explain the procedure to import XML, CSV (Text), and MS Access data into Excel. | 2 | 8 | CO4,PO5 |
| 10 | Discuss working with multiple worksheets and managing data across worksheets in Excel. | 2 | 8 | CO4,PO5 |
| **MODULE-3** | | | | |

| 1. | Explain the different types of machine learning with suitable examples. | 3 | 8 | CO1,PO1 |
|---|---|---|---|---|
| 2. | List and explain machine learning algorithms used for classification, clustering, and feature selection. | 3 | 8 | CO1,PO2 |
| 3. | Explain probability theory and derive Bayes' theorem with an example. | 3 | 8 | CO1,PO1 |
| 4. | Define Bayes probability and explain its role in machine learning. | 3 | 8 | CO1,PO2 |
| 5. | Explain the Cartesian plane and equations of straight lines with graphical representation. | 3 | 8 | CO2,PO1 |
| 6 | Explain the concept of exponents and their importance in data science computations. | 3 | 8 | CO2,PO1 |
| 7 | Describe commonly used tools for data science and explain their applications. | 3 | 8 | CO4,PO1 |
| 8 | Explain SQL and describe different SQL command categories: DDL, DML, DCL, TCL, and DQL with examples. | 3 | 8 | CO3,PO5 |
| 9 | Demonstrate the use of SELECT, INSERT, UPDATE, and DELETE commands in SQL. | 3 | 8 | CO3,PO5 |
| 10 | Explain the procedure to import SQL database data into Microsoft Excel. | 3 | 8 | CO4,PO5 |
| **MODULE-4** | | | | |
| 1. | Define correlation and explain its types with suitable examples. | 4 | 8 | CO1, PO1 |
| 2. | Describe linear regression and explain its assumptions and applications. | 4 | 8 | CO1, PO1 |
| 3. | Explain logistic regression and compare it with linear regression. | 4 | 8 | CO1, PO2 |
| 4. | Explain the Gaussian (normal) distribution and discuss its properties. | 4 | 8 | CO1, PO1 |
| 5. | Explain the concept of standardization and its importance in data analysis. | 4 | 8 | CO2, PO2 |
| 6 | Explain the standard normal probability distribution and demonstrate probability calculation using Excel. | 4 | 8 | CO2,PO5 |

| 7 | Explain Z-scores and describe how probabilities are calculated using Z-score tables. | 4 | 8 | CO2, PO2 |
|---|---|---|---|---|
| 8 | State and explain the Central Limit Theorem and its significance in statistics. | 4 | 8 | CO2, PO4 |
| 9 | Explain Markowitz Portfolio Optimization and discuss the role of Gaussian algebra in finance. | 4 | 8 | CO3,PO3 |
| 10 | Explain how standardization simplifies linear regression, including modeling error and information gain. | 4 | 8 | CO4,PO4 |
| **MODULE-5** | | | | |
| 1. | Explain the importance of data visualization and describe scatter plots, charts, graphs, histograms, and maps with suitable examples. | 5 | 8 | CO1, PO1 |
| 2. | Explain descriptive statistics for continuous data, focusing on mean and standard deviation, and discuss their significance. | 5 | 8 | CO2,PO1 |
| 3. | Explain how frequency and percentage are used for analyzing categorical data with examples. | 5 | 8 | CO2,PO2 |
| 4. | Describe the basic concepts and features of Python that make it suitable for data science. | 5 | 8 | CO2,PO2 |
| 5. | Explain Python strings and lists, including their operations and applications in data analysis. | 5 | 8 | CO3,PO1 |
| 6 | Describe tuples, sets, and dictionaries in Python and explain their differences with examples. | 5 | 8 | CO4,PO5 |
| 7 | Explain the role of NumPy in data science and discuss its important features. | 5 | 8 | CO4,PO5 |
| 8 | Describe the use of Pandas and Matplotlib libraries for data analysis and visualization. | 5 | 8 | CO3,PO1 |
| 9 | Explain the purpose of the Scikit-Learn library and discuss its role in implementing machine learning models. | 5 | 8 | CO4,PO5 |
| 10 | Describe the general steps involved in implementing a machine learning model using Python libraries. | 5 | 8 | CO4,PO5 |

Signature of the subject teacher          Signature of the Course Co-Ordinator          Signature of the Dean