Course Title: Data Handling and Visualization Tools
Course Code: CDCSE01

Duration: 1:30 Hours                                              Max. Marks: 25

**Note**: - Attempt all questions in the given order only. Missing data/information (if any), maybe suitably assumed & mentioned in the answer.

| Q. No. | Question | Marks | CO |
|---|---|---|---|
| 1a | Explain the stages involved in a typical data science. Illustrate each stage with an example from a real-world application. | 2.5 | 1 |
| 1b | A data scientist needs to collect and analyze data from a social media platform to study user engagement trends. Identify two potential data security issues that could arise in this context and propose solutions to mitigate these issues. | 2.5 | 1 |
| 2a | You have a dataset with 1,000 records across 5 columns: A, B, C, D, and E. The dataset is structured as follows: <br> 1. **Column A**: Contains 1,000 values. <br> 2. **Column B**: Contains 1,000 values. <br> 3. **Column C**: Contains 1,000 values. <br> 4. **Column D**: Contains 1,000 values. <br> 5. **Column E**: Contains 1,000 values. <br> Upon analysis, you find the following: <br> • 200 tuples (rows) have missing values in Column A. <br> • 150 tuples have missing values in Column B. <br> • 100 tuples have missing values in Column C. <br> • 50 tuples have missing values in Column D. <br> • 80 tuples have missing values in Column E. <br> However, some tuples have missing values in more than one column. <br> i. Determine the number of tuples with missing values in at least one column. Assume that no tuple is counted more than once, even if it has missing values in multiple columns. <br> ii. Calculate the percentage of tuples with missing values in at least one column with respect to the total number of tuples. Discuss possible strategies for handling missing data without removing the rows with missing values. Explain how each strategy can be applied and the potential impact on the dataset. | 2.5 | 2 |
| 2b | Explain the importance of data cleaning in the data pre-processing stage. Provide an example where improper data | 2.5 | 1 |

| | | | | |
|---|---|---|---|---|
| | cleaning led to misleading results. | | | |
| 3a | Given a dataset with the following data points for the variable "X": [12, 15, 13, 20, 18, 14, 17, 16, 15, 19], calculate the mean, mode, standard deviation, skewness, and kurtosis. Interpret the skewness and kurtosis values in terms of the data distribution shape. | 2.5 | 2 |
| 3b | How would you use a residual plot to evaluate the accuracy of a polynomial regression model? What specific patterns in the residual plot would indicate that the model is well-fitted or poorly fitted? | 2.5 | 3 |
| 4a | Describe the purpose of a pivot table in data analysis. Provide an example of how a pivot table can be used to summarize data for a retail store. | 2.5 | 2 |
| 4b | Discuss the interpretation of p-values in ANOVA. What does a low p-value indicate, and how should you proceed based on this result? Provide an example scenario where ANOVA might be used in a real-life situation. | 2.5 | 3 |
| 5a | What is the difference between simple regression and multiple regression? Provide a scenario where multiple regression is preferred over simple regression. | 2.5 | 3 |
| 5b | Explain a heat map and describe its typical use cases in data analysis. How does a heat map differ from other data visualization techniques like scatter plots or line graphs? | 2.5 | 3 |

me: 03 Hours

Max. Marks: 50

ote: - **Attempt all the five questions. Missing data/ information (if any), maybe suitably assumed &**
entioned in the answer.

| Q. No. | Question | Marks | CO |
|---|---|---|---|
| Q 1 | **Attempt any 2 parts of the following.** | | |
| 1a | Outline the typical stages in a data science project lifecycle. i. Explain the importance of each stage, from problem definition to deployment. ii. What challenges might arise at each stage, and how can they be addressed? | 5 | 1 |
| 1b | i. Compare and contrast different data collection strategies. What are the advantages and disadvantages of each method in terms of data quality and relevance? ii. Discuss how data integration and transformation can enhance the usability of data from multiple sources. | 5 | 1 |
| 1c | Analyse the ethical implications of data science applications. i. How can data privacy concerns be addressed while leveraging data for decision-making? ii. What frameworks or guidelines should organizations follow to ensure ethical use of data? | 5 | 1 |
| Q 2 | **Attempt any 2 parts of the following.** | | |
| 2a | Discuss the concept of in-sample evaluation and its importance in model development. i. What measures (e.g., R-squared, adjusted R-squared) can be used to evaluate the performance of a regression model? ii. How do you interpret these metrics in the context of your model? | 5 | 2 |
| 2b | Explain the ANOVA (Analysis of Variance) and describe a real-time scenario where it would be appropriate to use ANOVA. i. How does ANOVA help in comparing means across multiple groups? ii. What assumptions must be met for ANOVA to be valid? | 5 | 2 |
| 2c | Given a dataset with the following values: [12, 14, 15, 15, 16, 18, 19, 20, 22, 30], calculate the skewness and kurtosis. i. Interpret the results: What do the values of skewness and kurtosis suggest about the distribution? ii. How would you visually confirm your findings using a histogram or a box plot? | 5 | 2 |
| Q 3 | **Attempt any 2 parts of the following.** | 5 | 3 |
| 3a | Define what is meant by "integrity in visualization." i. Why is integrity important in the context of data visualization? ii. Provide an example of a visualization that lacks integrity and discuss the potential implications of misleading visualizations. | | |

| | | | | |
|---|---|---|---|---|
| 3b | List and describe the seven common quantitative relationships in graphs.<br>i. Explain the significance of each relationship in data analysis.<br>ii. Provide an example of a dataset that would be appropriate for illustrating each relationship. | 5 | | 3 |
| 3c | Explain how improper visualization techniques can mislead viewers.<br>i. What are some common pitfalls in data visualization that can distort information?<br>ii. Describe how intentional or unintentional biases in visualization can impact decision-making. | 5 | | 3 |
| Q 4 | Attempt any 2 parts of the following. | | | |
| 4a | Given two datasets, one containing sales data and another containing customer information, explain the steps to merge them in Python using Pandas. Then, demonstrate reshaping techniques to pivot this data based on customer regions. | 5 | | 4 |
| 4b | Explain what regular expressions are and provide an example of how they can be used for string matching in Python.<br>i. Describe how to use the *remodule* in Python to search for a pattern in a string.<br>ii. What is the significance of special characters in regular expressions? Provide examples of at least three special characters and their meanings. | 5 | | 4 |
| 4c | Explain the concept of data aggregation and its importance in data analysis.<br>i. How does the *GroupBy* operation work in pandas, and what is its purpose?<br>ii. Provide a scenario where grouping data would be beneficial, including an example of a *DataFrame* you might use. | 5 | | 4 |
| Q 5 | Attempt any 2 parts of the following. | | | |
| 5a | Identify and explain at least three factors that can lead to overfitting in a machine learning model.<br>i. How does the tradeoff between bias and variance relate to overfitting and underfitting?<br>ii. How can the size of the training dataset contribute to overfitting?<br>iii. Describe the role of noise in the training data and its effect on overfitting. | 5 | | 5 |
| 5b | Suppose you have trained a binary classifier to predict whether emails are spam (positive class) or not spam (negative class). After evaluating your classifier on a test set of 100 emails, you obtained the following confusion matrix: | 5 | | 5 |

|  | Predicted not Spam | Predicted Spam |
|---|---|---|
| Actually not Spam | 65 | 5 |
| Actually Spam | 10 | 20 |

| | | | | |
|---|---|---|---|---|
| | Using this confusion matrix, calculate the following performance metrics for the classifier: <br> i. Accuracy <br> ii. Precision for the spam class <br> iii. Recall for the spam class <br> iv. F1-score for the spam class | | | |
| 5c | Define cross-validation and explain how it can help to prevent overfitting. Use the concept of Ridge Regression to show how regularization can improve model generalization. | 5 | 5 |