**1. What are the key tasks that machine learning entails? What does data pre-processing imply?**

- Key tasks in machine learning include:

  - **Data collection**: Gathering relevant data for the problem at hand.

  - **Data pre-processing**: Cleaning, transforming, and preparing the data for analysis.

  - **Feature selection/engineering:** Identifying the most informative features or creating new ones.

  - **Model selection:** Choosing an appropriate machine learning model for the problem.

  - **Model training:** Fitting the model to the training data to learn patterns.

  - **Model evaluation:** Assessing the performance of the model on unseen data.

  - **Model deployment:** Integrating the trained model into real-world applications.

- Data pre-processing implies preparing the data for analysis by applying various techniques such as:

  - **Data cleaning:** Handling missing values, removing duplicates, and correcting inconsistencies.

  - **Data transformation:** Normalizing or scaling data to a common range.

  - **Feature encoding:** Converting categorical variables into numerical representations.

  - **Feature scaling:** Ensuring features are on a similar scale to avoid dominance by certain features.

  - **Handling outliers:** Identifying and dealing with data points that deviate significantly from the norm.

  - **Data splitting:** Dividing the dataset into training, validation, and test sets for model evaluation.

**2. Describe quantitative and qualitative data in depth. Make a distinction between the two.**

- **Quantitative data:** Quantitative data represents numerical values that can be measured or counted. It involves quantities, amounts, or sizes. Examples include age, height, temperature, and income. Quantitative data can be further categorized as discrete (countable) or continuous (measurable along a continuum).

- **Qualitative data:** Qualitative data represents non-numerical values or attributes that describe characteristics, opinions, or categories. It involves qualities, properties, or attributes. Examples include gender, color, occupation, and satisfaction level. Qualitative data is often represented by labels, categories, or text.

Distinction: The key distinction between quantitative and qualitative data lies in the nature of the values they represent. Quantitative data is numerical and allows for mathematical operations,

statistical analysis, and measurement of magnitude or quantity. On the other hand, qualitative data is descriptive and categorical, providing insights into characteristics or attributes without inherent numerical significance.

**3. Create a basic data collection that includes some sample records. Have at least one attribute from each of the machine learning data types.**

Here's a basic data collection with sample records, showcasing different data types:

| ID | Name    | Age | Gender | Height | Income |
|----|---------|-----|--------|--------|--------|
| 1  | John    | 25  | Male   | 180    | 40000  |
| 2  | Emily   | 32  | Female | 165    | 60000  |
| 3  | Michael | 45  | Male   | 175    | 55000  |
| 4  | Sarah   | 28  | Female | 160    | 45000  |
| 5  | David   | 36  | Male   | 190    | 70000  |

In this example:

- ID represents a unique identifier and is typically treated as a nominal or ordinal attribute.

- Name is a nominal attribute, representing categorical data.

- Age is a discrete quantitative attribute, as it represents a countable value.

- Gender is a nominal attribute, indicating categories.

- Height is a continuous quantitative attribute, as it represents a measurable value.

- Income is a continuous quantitative attribute, representing a numerical value.2. Describe quantitative and qualitative data in depth. Make a distinction between the two.

Apologies for the inconvenience caused. Here are the answers to all the questions:

**4. What are the various causes of machine learning data issues? What are the ramifications?**

- Various causes of machine learning data issues include:

  - Missing data: Absence of values for certain attributes.

  - Outliers: Data points that significantly deviate from the normal distribution.

  - Imbalanced data: Significant disparity in the distribution of classes.

- Incorrect or inconsistent data: Errors or inconsistencies in data collection or recording.

- Noisy data: Data that contains irrelevant or misleading information.

- Ramifications of data issues:

  - Missing data can lead to **biased or inaccurate models** if not handled properly.

  - Outliers can skew the model's performance or lead to **erroneous conclusions**.

  - Imbalanced data can result in **poor predictions** for minority classes.

  - Incorrect or inconsistent data can lead to **faulty models or unreliable insights**.

  - Noisy data can introduce **unnecessary complexity** and **affect the model's generalization(over-fitting).**

**5. Demonstrate various approaches to categorical data exploration with appropriate examples.**

- Various approaches to explore categorical data include:

  - **Frequency distribution**: Counting the occurrences of each category.

    Example: Counting the number of students in each grade level (e.g., 9th, 10th, 11th, 12th).

  - Bar chart: Visualizing the frequency distribution using bars of different heights.

    Example: Plotting a bar chart to display the number of customers in different age groups (e.g., 18-24, 25-34, 35-44).

  - Pie chart: Illustrating the proportion of each category in relation to the whole.

    Example: Creating a pie chart to show the percentage of market share for different smartphone brands.

  - Cross-tabulation: Examining the relationship between two categorical variables.

    Example: Analyzing the relationship between education level and employment status.

**6. How would the learning activity be affected if certain variables have missing values? Having said that, what can be done about it?**

- If certain variables have missing values, the learning activity can be affected in several ways:

  - **Loss of information**: Missing values can result in the loss of valuable information, leading to biased or incomplete models.

- **Decreased model performance**: Missing values can disrupt the training process and negatively impact the model's performance.

  - **Incorrect conclusions**: Ignoring missing values without appropriate handling can lead to inaccurate conclusions and predictions.

- Ways to handle missing values:

  - Removal: If the missing values are negligible in quantity, the corresponding records or variables can be removed from the dataset.

  - Imputation: Missing values can be replaced with estimated or predicted values based on the available data or statistical techniques.

  - Advanced techniques: More sophisticated methods like multiple imputation or matrix completion can be employed to handle missing values.

**7. Describe the various methods for dealing with missing data values in depth.**

- Various methods for dealing with missing data values include:

  - **Mean/Median imputation**: Replace missing values with the mean or median of the available data for that attribute.

  - **Mode imputation:** Replace missing values with the mode (most frequent value) of the available data for that attribute.

  - **Hot deck imputation:** Replace missing values with values randomly selected from similar existing records.

  - **Multiple imputation:** Generate multiple imputed datasets by estimating missing values using statistical models, then combine the results for analysis.

  - **K-nearest neighbors imputation:** Predict missing values by finding the k-nearest neighbors based on other attributes and using their values as imputations.

  - **Matrix completion**: Utilize matrix factorization or other techniques to fill in missing values based on patterns in the data.

**8. What are the various data pre-processing techniques? Explain dimensionality reduction and feature selection in a few words.**

Various data pre-processing techniques include:

- **Data cleaning:** Handling missing values, duplicates, and inconsistencies in the dataset.
- **Data transformation:** Normalizing, scaling, or encoding features to ensure they are suitable for analysis.
- **Feature selection:** Identifying the most relevant subset of features that contribute to the prediction task.

- **Dimensionality reduction**: Reducing the number of features while preserving important information to simplify the model and improve efficiency.
- **Handling outliers:** Detecting and addressing data points that deviate significantly from the majority, which can impact the model's performance.
- **Balancing classes:** Addressing class imbalance by oversampling the minority class or undersampling the majority class.
- **Data integration**: Combining multiple datasets from different sources to create a unified dataset for analysis.

Dimensionality reduction: It is a process of reducing the number of features in a dataset while preserving as much important information as possible. This helps to address the curse of dimensionality, improve computational efficiency, and avoid overfitting.

**9. i) What is the IQR? What criteria are used to assess it?**

- IQR stands for Interquartile Range. It is a measure of statistical dispersion and represents the range between the first quartile (Q1) and the third quartile (Q3) in a dataset. It provides insights into the spread and variability of the middle 50% of the data. The formula for IQR is IQR = Q3 - Q1.
- Criteria used to assess the IQR:
    - **Outliers:** The IQR is commonly used in outlier detection. Data points below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR are often considered outliers.
    - **Skewness:** The IQR can indicate the skewness of the data distribution. If Q1 is close to Q3, the data is approximately symmetric. If Q1 is significantly lower or higher than Q3, the data may be skewed.

**ii. Describe the various components of a box plot in detail? When will the lower whisker surpass the upper whisker in length? How can box plots be used to identify outliers?**

- Box plots (also known as box-and-whisker plots) represent the distribution of a dataset visually. The various components of a box plot include:
    - Median (Q2): The horizontal line within the box that represents the middle value of the dataset.
    - Box: The box represents the interquartile range (IQR), spanning from the first quartile (Q1) to the third quartile (Q3).
    - Whiskers: Vertical lines extending from the box indicate the variability outside the IQR. They can represent different ranges depending on the plot's configuration.
    - Outliers: Data points outside the whiskers are considered outliers and are plotted individually as points or asterisks.
- The lower whisker surpasses the upper whisker in length when the dataset has a highly skewed distribution or when there are outliers present on the lower end of the data. This indicates that the lower part of the dataset has a larger spread than the upper part.
- Box plots can be used to identify outliers by examining the data points outside the whiskers. Any data point that falls below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR is typically considered an outlier and can be identified visually in the plot.

**10. Brief notes on the following:**

Data collected at regular intervals:

- o Data collected at regular intervals refers to data points collected at consistent time or space intervals.
- o Examples include temperature readings recorded every hour, stock prices recorded at the end of each trading day, or patient vital signs measured every 15 minutes.
- o Regularly spaced data enables the application of time series analysis techniques, trend analysis, and forecasting models.

**The gap between the quartiles:**

- o The gap between the quartiles, also known as the interquartile range (IQR), represents the spread of the middle 50% of the dataset.
- o It is calculated as IQR = Q3 - Q1, where Q1 is the first quartile and Q3 is the third quartile.
- o The IQR provides insights into the variability and dispersion of the data distribution.
- o A larger IQR indicates a wider spread of the data, while a smaller IQR suggests a more concentrated distribution.

**Using a cross-tab:**

- o A cross-tab, short for cross-tabulation, is a table that displays the frequency distribution of two or more categorical variables.
- o It shows how the categories of one variable are distributed across the categories of another variable.
- o Cross-tabs are useful for analyzing relationships and dependencies between variables, identifying patterns, and gaining insights into the data.
- o They are commonly used in fields such as market research, social sciences, and data analysis to summarize and visualize categorical data.

**11. Comparison between:**

1. Data with nominal and ordinal values:
   - Nominal data represents categories or labels with no inherent order or ranking. Examples include colors, genders, or categories like "dog," "cat," and "bird."
   - Ordinal data represents categories with a specific order or ranking. Examples include rating scales (e.g., "low," "medium," "high"), educational levels (e.g., "elementary," "high school," "college"), or customer satisfaction levels (e.g., "very dissatisfied," "neutral," "very satisfied").
   - Nominal data can be encoded using one-hot encoding, while ordinal data can be encoded with integer values corresponding to their rank or order.

- Statistical operations like mode can be applied to nominal data, while ordinal data can also use operations like median or percentiles that preserve the ranking information.

2. Histogram and box plot:
   - Histograms and box plots are both used to visualize the distribution of **a continuous variable**.
   - A histogram represents the distribution of data by dividing it into intervals (bins) and displaying the frequency or density of data points within each bin using bars.
   - A box plot, as explained earlier, displays the quartiles, median, and outliers of a dataset using boxes and whiskers.
   - Histograms provide a more detailed view of the distribution, showing the shape, skewness, and modes, while box plots provide a summary view of the central tendency, spread, and presence of outliers.
   - Histograms are suitable for exploring the data distribution, while box plots are useful for comparing distributions, identifying outliers, and understanding the variability within groups.
3. The average and median:
   - The average (mean) is a measure of central tendency calculated by summing all the values and dividing by the total number of values.
   - The median is also a measure of central tendency and represents the middle value when the data is sorted in ascending or descending order.
   - The average is sensitive to extreme values or outliers, while the median is more robust to outliers.
   - The average can be affected by skewed distributions, while the median is less influenced by skewed data.
   - The average is suitable for symmetrical distributions, while the median is more appropriate for skewed or non-normal distributions.