

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 and transforming the data to make it suitable for analysis, including handling missing values, dealing with outliers, and normalizing or scaling the data.
- Feature engineering: Selecting or creating informative features from the available data that can be used as inputs for machine learning models.
- Model selection: Choosing the appropriate machine learning algorithm or model based on the problem and data characteristics.
- Model training: Fitting the selected model to the training data to learn patterns and relationships.
- Model evaluation: Assessing the performance of the trained model on unseen data using appropriate evaluation metrics.
- Model optimization: Fine-tuning the model's parameters to improve its performance.
- Model deployment: Integrating the trained model into a production environment for making predictions.

- Data pre-processing refers to the steps taken to clean, transform, and prepare the data for analysis. It involves tasks such as handling missing data, dealing with outliers, normalizing or scaling numerical features, encoding categorical variables, and splitting the data into training and testing sets. The goal of data pre-processing is to ensure the data is in a suitable format for the machine learning algorithms to learn from.

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

- Quantitative data: Quantitative data is numerical data that represents quantities or measurements. It can be further categorized as continuous or discrete. Continuous data can take on any value within a specific range (e.g., temperature), while discrete data only has specific, separate values (e.g., number of children). Quantitative data can be analyzed using mathematical and statistical methods.

- Qualitative data: Qualitative data, also known as categorical data, represents qualities or attributes that cannot be measured numerically. It is typically represented by labels or categories. Qualitative data can be further divided into nominal and ordinal data. Nominal data represents categories without any inherent order or ranking (e.g., color), while ordinal data represents categories with a specific order or ranking (e.g., education level). Qualitative data is analyzed using methods such as frequency counts, cross-tabulations, and visualizations specific to categorical data.

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

Example data collection:

ID	Age	Gender	Income	Product Category
1	35	Male	\$50,000	Electronics
2	42	Female	\$65,000	Clothing
3	28	Male	\$40,000	Electronics
4	45	Female	\$80,000	Books
5	31	Male	\$55,000	Electronics

- In this example, the "Age" attribute represents numerical data (quantitative).
- The "Gender" attribute represents categorical data (qualitative).
- The "Income" attribute represents numerical data (quantitative).
- The "Product Category" attribute represents categorical data (qualitative).

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, which can lead to biased or incomplete analyses.
- Outliers: Extreme values that deviate significantly from the majority of the data, which can distort the analysis and model performance.
- Imbalanced data: Significant differences in the distribution of classes or categories, which can lead to biased model predictions.
- Inconsistent data: Contradictory or conflicting information within the dataset, which can introduce errors or affect the accuracy of the analysis.
- Incorrect or noisy data: Data that contains errors, inconsistencies, or irrelevant information, which can mislead the analysis and produce inaccurate models.

- The ramifications of these data issues can include:

- Biased models: Data issues can lead to biased models that produce inaccurate or unfair predictions.
- Decreased model performance: Outliers or inconsistent data can negatively impact model performance, leading to poor predictions or low accuracy.
- Reduced generalizability: Issues like imbalanced data can affect the ability of the model to generalize well to new, unseen data.
- Misinterpretation of results: Inconsistent or noisy data can lead to incorrect interpretations and misguided conclusions from the analysis.

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

Approaches to categorical data exploration include:

- Frequency counts: Counting the occurrences of each category in a dataset.  
Example: Counting the number of male and female customers in a customer database.
- Cross-tabulation: Analyzing the relationship between two categorical variables by creating a contingency table.  
Example: Creating a cross-tabulation table to examine the distribution of product categories based on customer gender.
- Bar chart: Visualizing the distribution of categorical data using vertical bars, where the height of each bar represents the frequency or count of each category.  
Example: Creating a bar chart to visualize the number of customers in each age group.
- Pie chart: Visualizing the proportion or percentage distribution of categorical data using a circular chart divided into sectors, where each sector represents a category.  
Example: Creating a pie chart to represent the percentage of customers in each product category.

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:
  - Bias in analysis: Missing values may introduce bias in the analysis, as the patterns or relationships in the data may be distorted.
  - Inaccurate model performance: Missing values can lead to inaccurate model performance due to incomplete information and biased estimation.
  - Reduced sample size: Missing values can reduce the effective sample size available for analysis, potentially leading to reduced statistical power.
- To address missing values, several approaches can be employed:
  - Complete case analysis: Exclude records with missing values from the analysis. However, this can result in a loss of information and reduced sample size.
  - Imputation: Estimate missing values using statistical techniques, such as mean imputation (replacing missing values with the mean of the available data) or regression imputation (predicting missing values based on other variables).
  - Multiple imputation: Generate multiple imputed datasets, each with plausible values for the missing data, and analyze them separately to account for the uncertainty introduced by imputation.
  - Advanced methods: Utilize more sophisticated techniques, such as expectation-maximization algorithms or machine learning models, to impute missing values based on patterns in the data.

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

- Listwise deletion: In this approach, records with missing values are completely excluded from the analysis. While it is simple to implement, it can result in a loss of information and reduced sample size.
- Pairwise deletion: In pairwise deletion, missing values are ignored on a variable-by-variable basis. Each analysis uses all available data for that specific analysis, resulting in varying effective sample sizes for different analyses. However, it can lead to biased results if the missingness is related to the outcome or variables of interest.

- Mean or median imputation: Missing values are replaced with the mean or median value of the available data for that variable. This method is simple to implement but can distort the distribution and relationships in the data.

- Regression imputation: Missing values are predicted based on a regression model using other variables as predictors. This method takes into account the relationships between variables but assumes the relationships hold for the missing cases.
- Multiple imputation: Multiple imputation involves creating multiple plausible imputed datasets, each with different imputed values for missing data, based on the observed data and imputation model. These datasets are then analyzed separately, and the results are combined to account for the uncertainty introduced by imputation.
- Advanced methods: There are more advanced techniques, such as expectation-maximization algorithms or machine learning-based imputation models, that can be used to impute missing values based on patterns in the data. These methods can capture complex relationships and yield more accurate imputations.

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:
  - Handling missing data: Dealing with missing values through imputation or deletion.
  - Outlier detection and treatment: Identifying and handling outliers, either by removing them or transforming them to reduce their impact.
  - Data normalization: Scaling the data to a common range, such as 0 to 1 or -1 to 1, to ensure variables with different scales have equal influence in the analysis.
  - Encoding categorical variables: Transforming categorical variables into numerical representations that can be used by machine learning algorithms, such as one-hot encoding or ordinal encoding.
  - Feature scaling: Scaling numerical features to a standard range to prevent variables with larger magnitudes from dominating the analysis.
  - Dimensionality reduction: Reducing the number of input features or variables while preserving important information. It helps to mitigate the curse of dimensionality, improve model efficiency, and remove redundant or irrelevant features.
  - Feature selection: Selecting a subset of relevant features from the original set of features to improve model performance, interpretability, and reduce computational complexity. It involves evaluating the importance or predictive power of features and retaining only the most informative ones.

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. The IQR captures the central 50% of the data.
- The IQR is assessed using the following criteria:
  - The IQR is used to identify outliers by considering observations that fall below  $Q1 - 1.5 * IQR$  or above  $Q3 + 1.5 * IQR$  as potential outliers.
  - The IQR is also used to construct box plots, where the box represents the IQR, and the whiskers extend to the data points that fall within  $Q1 - 1.5 * IQR$  and  $Q3 + 1.5 * IQR$ .

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?

- A box plot, also known as a box-and-whisker plot, visualizes the distribution of a dataset. It consists of the following components:
  - Median: The line inside the box represents the median, which is the middle value when the data is sorted.
  - Box: The box represents the IQR, with the lower edge at Q1 and the upper edge at Q3.
  - Whiskers: Lines extending from the box indicate the range of the data within  $Q1 - 1.5 * IQR$  and  $Q3 + 1.5 * IQR$ . Data points outside this range are considered potential outliers.
  - Outliers: Individual data points beyond the whiskers are marked as outliers, shown as individual points or asterisks.
- The lower whisker will surpass the upper whisker in length when the lower quartile (Q1) is larger than Q3, indicating a reversed order of the quartiles. This situation may occur if the data distribution is highly skewed.
- Box plots can be used to identify outliers by visually inspecting data points outside the whiskers. Outliers are observations that fall below  $Q1 - 1.5 * IQR$  or above  $Q3 + 1.5 * IQR$ . Box plots provide a clear representation of the central tendency, spread, and potential outliers in a dataset.

10. Make brief notes on any two of the following:

1. Data collected at regular intervals:

- Data collected at regular intervals refers to the collection of observations or measurements at consistent time intervals.
- Examples include stock market data recorded every minute, hourly weather measurements, or daily sales data.
- Regularly spaced intervals enable the analysis of trends, patterns, and seasonal variations in the data.
- Time series analysis techniques can be applied to understand and forecast future values based on historical patterns.

## 2. The gap between the quartiles:

- The gap between the quartiles, known as the interquartile range (IQR), measures the spread or dispersion of a dataset.
- It is calculated as the difference between the third quartile (Q3) and the first quartile (Q1),  $IQR = Q3 - Q1$ .
- The IQR represents the range where the central 50% of the data falls, providing a measure of the data's variability.
- It is used to identify potential outliers and construct box plots to visualize the distribution of the data.

## 3. Use a cross-tab:

- A cross-tab, also known as a contingency table, is a table that displays the joint distribution of two categorical variables.
- It presents the frequency or count of observations for each combination of categories between the two variables.
- Cross-tabs are used to analyze the relationship or association between two categorical variables and identify patterns or dependencies.
- They help understand the joint distribution, compare the frequencies across categories, and perform statistical tests, such as chi-square tests, to assess the significance of the relationship.