

# ML ASSISMENT 04

## 1. What are the key tasks involved in getting ready to work with machine learning modeling?

Sol: key Task involved are:

1. Problem Definition and Goal Setting:
  - Clearly define the problem you want to solve with machine learning.
  - Set specific goals and objectives for your modeling project. What do you want to achieve with the model?
2. Data Collection and Exploration:
  - Gather relevant and high-quality data for your project.
  - Explore and analyze the data to understand its structure, quality, and potential challenges.
  - Handle missing values, outliers, and noise in the data.
3. Data Preprocessing and Cleaning:
  - Preprocess the data by performing tasks like data normalization, scaling, and encoding categorical variables.
  - Handle data imbalances and prepare the data for model training.
4. Feature Selection and Engineering:
  - Select relevant features that will contribute to the model's performance.
  - Create new features or transform existing ones to provide additional information to the model.
5. Data Splitting:
  - Split the dataset into training, validation, and test sets.
  - The training set is used to train the model, the validation set is used to tune hyperparameters, and the test set is used to evaluate the model's generalization performance.
6. Model Selection:
  - Choose an appropriate machine learning algorithm based on the problem type (classification, regression, clustering, etc.) and the characteristics of the data.
  - Consider factors like model complexity, interpretability, and computational efficiency.
7. Model Training:
  - Train the selected model using the training data.
  - Adjust hyperparameters to achieve the best performance on the validation set.
8. Model Evaluation:
  - Evaluate the model's performance using appropriate metrics (accuracy, precision, recall, F1-score, etc.).
  - Use the validation set to fine-tune the model.
9. Model Tuning and Optimization:
  - Fine-tune hyperparameters to improve the model's performance.
10. Model Interpretation and Analysis:
  - Understand the factors that contribute to the model's predictions.
  - Use techniques like visualization to interpret the model's behavior.

11. Model Deployment:

- Prepare the model for deployment in a production environment.
- Consider factors like model serving, scalability, and monitoring.

12. Documentation:

- Document the entire process, including data sources, preprocessing steps, model architecture, hyperparameters, and results.
- This documentation aids in replicability and knowledge sharing.

Communication:

- Present your findings, insights, and results to stakeholders in a clear and understandable manner.
- Discuss the limitations and potential applications of the model.

Maintenance and Updates:

- Monitor the deployed model's performance and retrain/update it as new data becomes available.
- Keep track of changes in the data distribution and adjust the model accordingly.

**2. What are the different forms of data used in machine learning? Give a specific example for each of them.**

Sol:

1. Structure Data:

Structured data refers to data organized in a well-defined tabular format, typically represented as rows and columns.

Example: A dataset containing information about customers, such as their names, ages, genders, purchase histories, and spending amounts

2. Unstructured Data:

- Unstructured data refers to data that lacks a predefined structure, making it more challenging to analyze.
- Example: Text data from social media posts, customer reviews, or medical notes. This data can include free-form text, images, audio, and video content.

3. Time Series Data:

- Time series data consists of observations recorded over a sequence of time intervals.
- Example: Stock market price data recorded over regular time intervals.

4. Categorical Data:

- Categorical data represents different categories or labels, but there is no inherent order among the categories.
- Example: A dataset of animals categorized as "cat," "dog," "elephant," etc.

5. Numerical Data:

- Numerical data consists of continuous or discrete numerical values that can be used for mathematical calculations.
- Example: Heights of individuals, where each height measurement is a numerical value that can be compared, averaged, and analyzed mathematically.

6. Image Data:

- Image data consists of visual representations captured as pixels in a grid.
- Example: A dataset of handwritten digits (MNIST dataset), where each instance is an image of a digit (0-9) represented by pixel values.

### 7. Text Data:

- Text data includes written or typed text, often used in natural language processing tasks.
- Example: A collection of email messages, where each message can be analyzed for sentiment, topic, or other characteristics.

### 8. Audio Data:

- Audio data includes sound recordings or waveforms and is often used in speech recognition and audio analysis.
- Example: A dataset of spoken phrases in different languages, used to train a speech recognition system.

### 9. Video Data:

- Video data is a sequence of images presented over time, capturing visual information.
- Example: Video surveillance footage used to detect and track objects or activities within a scene.

### 10. Geospatial Data:

- Geospatial data includes information tied to geographical locations and is used in applications such as mapping and geolocation.
- Example: GPS coordinates of mobile devices, used to provide navigation instructions and location-based services.

## 3. Distinguish:

### 1. Numeric vs. categorical attributes

Numeric Attribute	Categorical Attributes
<ul style="list-style-type: none"><li>• Numeric attributes consist of numerical values that can be used for mathematical calculations.</li><li>• They can be continuous or discrete. Continuous attributes have a wide range of possible values (e.g., height, weight), while discrete attributes have a finite set of distinct values (e.g., age, number of siblings)</li></ul>	<ul style="list-style-type: none"><li>• Categorical attributes represent different categories or labels that describe the data points.</li><li>• They can be nominal or ordinal. Nominal attributes have no inherent order or ranking among categories (e.g., colors, animal types), while ordinal attributes have a specific order or ranking (e.g., education levels like "high school," "college," "graduate")</li></ul>

### 2. Feature selection vs. dimensionality reduction

Feature Selection	Dimensionality Reduction
<ul style="list-style-type: none"><li>• Feature Selection is the process of selecting a subset of relevant features from the original set of features in a dataset.</li><li>• The goal of feature selection is to choose the most informative and significant feature which contribute most predicting feature power to the model</li></ul>	<ul style="list-style-type: none"><li>• Dimensionality reduction is the process of transforming high-dimensional data into a lower-dimensional representation while preserving the most important information.</li><li>• The goal is to reduce the complexity.</li><li>• Example Techniques : PCA (principal Components Analysis).</li></ul>

**4. Make quick notes on any two of the following:**

**[ 1. The histogram 2. Use a scatter plot 3.PCA (Personal Computer Aid)]**

Sol:

- Histogram:
  - 1) histogram is a graphical representation used to display the distribution of a continuous dataset. It consists of a series of adjacent bars, where
  - 2) Each bar represents a range of values (called a "bin") and
  - 3) The height of the bar represents the frequency or count of data points within that range.
  - 4) histograms provide insights into the underlying data distribution, revealing patterns such as symmetry, skewness, and outliers. They are particularly useful for understanding the spread and central tendency of data
- Scatter Plot:
  1. A scatter plot is a two-dimensional graphical representation that displays individual data points as dots on a Cartesian plane.
  2. It is used to visualize the relationship between two continuous variables.
  3. Each dot on the scatter plot represents a single data instance, with its position determined by the values of the two variables being compared.
  4. Scatter plots help identify trends, patterns, and the strength of relationships between variables.
  5. They are especially useful for exploring correlations and detecting outliers in data

**5. Why is it necessary to investigate data? Is there a discrepancy in how qualitative and quantitative data are explored?**

Sol:

Investigating data is a critical step in the data analysis and machine learning process. It helps you understand the characteristics, patterns, and potential issues present in your dataset. Regardless of whether the data is qualitative (categorical) or quantitative (numeric).

Exploring the data dependency:

1. Data Quality Assurance: Clean and accurate data is essential for building reliable models.
2. Understanding Data Distribution: Exploring data helps you understand the distribution of values within variables.
3. Identifying Patterns and Relationships: Data exploration helps you uncover patterns, trends, and relationships between variables.
4. Feature Engineering and Selection: Exploring data aids in selecting relevant features (variables) for modeling.
5. Outlier Detection: Both qualitative and quantitative data can have outliers .
6. Assumption Checking: In statistical analysis, assumptions about the data's distribution and properties are often made. Exploring the data helps verify whether these assumptions hold true, which is essential for drawing valid conclusions.
7. Preprocessing and Transformation: This could involve scaling, normalization, or encoding categorical variables.

8. Visualization and Communication: Charts, histograms, scatter plots, and other visuals help convey insights.

## 6. What are the various histogram shapes? What exactly are 'bins'?

1. Normal Distribution (Bell Curve):
  - A symmetrical distribution with the majority of data clustered around the mean.
  - The tails of the distribution extend equally on both sides.
  - Also known as Gaussian distribution.
2. Skewed Right (Positively Skewed):
  - The tail of the distribution extends toward higher values.
  - Most data points are clustered on the left side.
  - Also called right-skewed or positively skewed distribution.
3. Skewed Left (Negatively Skewed):
  - The tail of the distribution extends toward lower values.
  - Most data points are clustered on the right side.
  - Also called left-skewed or negatively skewed distribution.
4. Bimodal Distribution:
  - The distribution has two distinct peaks.
  - Indicates the presence of two underlying modes or groups within the data.
5. Multimodal Distribution:
  - Similar to bimodal but with more than two peaks.
  - Represents data with multiple distinct groups or modes.
6. Uniform Distribution:
  - All values in the dataset have similar frequencies.
  - No particular value is more common than others.
7. Exponential Distribution:
  - A distribution that quickly drops off from the maximum value and has a long tail on one side.
  - Often seen in situations involving time intervals between events.
8. Log-Normal Distribution:
  - Data that follows a log-normal distribution results in a histogram with a skewed right shape when plotted on a linear scale. However, when plotted on a logarithmic scale, it approximates a normal distribution.

Bins also known as intervals, are the ranges into which the entire range of data is divided in histogram.

## 7. How do we deal with data outliers?

**Sol:**

- identify and visualize outliers using techniques such as box plots, scatter plots, and histograms.
- Use statistical methods like the Z-score or the interquartile range (IQR) to quantitatively identify potential outliers.
- Investigate the source of outliers and assess whether they are due to errors or represent genuine data points.

- If outliers are identified as errors or anomalies, consider removing them from the dataset if they are expected to negatively impact analysis or modeling.
- However, be cautious when removing outliers, as they might contain valuable information or insights.
- Transformations can help make the data distribution more symmetric and reduce the influence of outliers.(such as logarithmic, square root, or cube root)
- Winsorization involves replacing extreme values with the nearest non-outlier values within a predefined range. This approach mitigates the impact of outliers without completely removing them.
- Create bins or discrete categories for data values, which can help reduce the influence of individual outlier values.
- Some machine learning algorithms are less sensitive to outliers, such as tree-based algorithms (e.g., decision trees, random forests) and support vector machines.
- Consider using algorithms that are inherently robust to outliers in your analysis.
- For missing values that might appear as outliers, consider imputing them with reasonable estimates to prevent distortions in analysis.

### **8. What are the various central inclination measures? Why does mean vary too much from median in certain data sets?**

**Sol:**

Central inclination measures are statistical metrics that provide information about the center or typical value of a dataset. They help summarize the central tendency or the "average" value of the data.

- The three main central inclination measures are the mean, median, and mode.

The mean can vary significantly from the median in certain data sets, particularly when

- the data is skewed or has outliers.
  - a. Skewness refers to the asymmetry in the distribution of data.
  - b. A dataset is positively skewed (tail on the right side), the presence of a few large values (outliers) on the right side of the distribution can pull the mean to the right, making it greater than the median.
  - c. In a negatively skewed dataset (tail on the left side), outliers on the left side can pull the mean to the left, making it smaller than the median.

In contrast, the median is resistant to outliers and less affected by the shape of the distribution. It represents the middle value, which is why it's often preferred in situations where the data contains extreme values or the distribution is highly skewed.

### **9. Describe how a scatter plot can be used to investigate bivariate relationships. Is it possible to find outliers using a scatter plot?**

**Sol:** A scatter plot is a graphical representation of bivariate data, where two variables are plotted on the x and y axes. It's a powerful tool for visually exploring the relationship between two variables and identifying patterns, trends, and potential outliers.

scatter plot can be used to investigate bivariate relationships:

- Visualizing Relationship: A scatter plot allows you to see how two variables are related to each other.
- patterns and Trends: By examining the overall distribution of points on the scatter plot.
- Strength of Relationship: The tightness of the points around a trendline or curve on the scatter plot can give you an idea of the strength of the relationship between the variables.
- Outliers: outliers can be identified using the scatter plot.

## 10. Describe how cross-tabs can be used to figure out how two variables are related

Sol:

Cross-tabs is a statistical technique used to analyze the relationship between two categorical variables.

Categorical variables are those that represent different categories or groups, and they can't be ordered or measured on a continuous scale.

Cross-tabs are used to figure out how two variables are related:

- Data Collection and Organization: Example, data on a survey that includes respondents' gender (Male/Female) and their favorite type of entertainment.
- Creating a Cross-Tabulation Table: A cross-tabulation table is constructed by organizing the data into rows and columns.
- Analysis: With the cross-tabulation table in place, we can now analyze the relationship between the two variables:
  1. Frequency Analysis: You can determine how frequently each category occurs.
  2. Row Percentages and Column Percentages: Calculating percentages within rows and columns helps reveal the distribution.
  3. Chi-Square Test: A chi-square test of independence can be performed to determine whether the observed differences in the cross-tabulation table are statistically significant or just due to random chance.
  4. Visual Representation.