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

Getting ready to work with machine learning modeling involves several key tasks that are important to ensure the success of the modeling process. Here are some of the key tasks involved:

1. Define the problem: The first step is to clearly define the problem you are trying to solve with machine learning. This involves identifying the business problem or objective and determining how machine learning can help achieve that objective.

2. Collect and prepare data: Machine learning models require large amounts of data to train and validate. This involves collecting relevant data, cleaning and processing it to ensure it is in a suitable format for analysis.

3. Feature engineering: Feature engineering is the process of selecting and transforming the relevant features of the data to build the model. This includes selecting the most important variables, encoding categorical variables, and scaling numerical features.

4. Select and train the model: After preparing the data and selecting the relevant features, you need to select a suitable machine learning algorithm that is appropriate for the problem you are solving. This involves training the model using the data and evaluating its performance.

5. Fine-tune the model: Once you have trained the model, you need to fine-tune it to improve its performance. This involves tweaking the hyperparameters of the model and selecting the best set of hyperparameters that optimize the model's performance.

6. Deploy the model: Once the model has been fine-tuned and validated, it can be deployed in production. This involves integrating the model into the business process and ensuring it can handle real-world data.

7. Monitor and maintain the model: Machine learning models require ongoing maintenance and monitoring to ensure they continue to perform well over time. This involves monitoring the model's performance, retraining the model as needed, and updating the model as new data becomes available.

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

There are three main forms of data used in machine learning:

1. Numerical data: This type of data consists of numbers and is used in regression and classification problems. Examples include data such as age, weight, temperature, and price.

2. Categorical data: This type of data consists of categories or labels and is used in classification problems. Examples include data such as gender, color, language, and nationality.

3. Text data: This type of data consists of natural language text and is used in natural language processing (NLP) tasks such as sentiment analysis, topic modeling, and text classification. Examples include data such as customer reviews, news articles, and social media posts.

Here are specific examples for each type of data:

1. Numerical data: In a prediction problem such as predicting house prices, numerical data such as the number of bedrooms, square footage, and age of the property would be used to train the machine learning model.

2. Categorical data: In a classification problem such as predicting whether a customer will churn or not, categorical data such as the customer's gender, age group, and occupation would be used to train the model.

3. Text data: In a sentiment analysis task, text data such as customer reviews or social media posts would be used to train the model to classify the sentiment of the text as positive, negative, or neutral.

3. Distinguish:

## 1. Numeric vs. categorical attributes

Numeric and categorical attributes are two different types of data that can be used in machine learning.

Numeric attributes are quantitative data that represent some measurable quantity such as height, weight, temperature, or price. Numeric attributes can be continuous, like height or weight, where the value can be any number within a range, or discrete, like the number of items sold, where the value can only take on integer values.

Categorical attributes, on the other hand, represent qualitative data and describe characteristics that cannot be measured quantitatively. Categorical attributes can be binary, like yes/no or true/false, or multi-valued, like color, gender, or product category.

When working with machine learning models, the type of attribute used can have an impact on the choice of algorithm and the preprocessing steps needed. For example, linear regression models work well with numeric data, while decision tree and random forest models work well with both numeric and categorical data.

It is also important to note that categorical data needs to be encoded or transformed into a numeric representation for many machine learning algorithms to work. One common

approach for encoding categorical data is one-hot encoding, which creates a binary column for each possible value of the categorical attribute.

## 2. Feature selection vs. dimensionality reduction

Feature selection and dimensionality reduction are two techniques used in machine learning to reduce the number of features used in a model.

Feature selection is the process of selecting a subset of the original features that are most relevant to the model's performance. This involves evaluating the importance of each feature and selecting the ones that contribute the most to the model's predictive power. Feature selection can be performed using statistical tests or by using algorithms that rank the features based on their importance.

Dimensionality reduction, on the other hand, is the process of transforming the original features into a lower-dimensional space while preserving the most important information. This can be done by removing correlated features or by projecting the features onto a lower-dimensional space using techniques like principal component analysis (PCA) or t-distributed stochastic neighbor embedding (t-SNE).

While both techniques aim to reduce the number of features, they differ in how they achieve this goal. Feature selection selects a subset of the original features, while dimensionality reduction transforms the original features into a lower-dimensional space.

The choice of which technique to use depends on the specific problem and the nature of the data. Feature selection is often used when the number of features is not too large and when the goal is to improve the interpretability of the model by identifying the most important features. Dimensionality reduction is often used when dealing with high-dimensional data, where the number of features is much larger than the number of samples, or when the goal is to improve the model's performance by reducing the complexity of the data.

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

## 1. The histogram

A histogram is a graphical representation of the distribution of a dataset. It is a type of bar chart that represents the frequency of values in a continuous or discrete dataset. In a histogram, the x-axis represents the range of values in the dataset, divided into intervals called bins, and the y-axis represents the frequency of values in each bin.

To create a histogram, the dataset is first divided into a set of bins of equal width. Then, the number of data points that fall into each bin is counted and represented as the height of the corresponding bar on the y-axis.

Histograms are useful for visualizing the shape of a distribution and identifying patterns and outliers in the data. They are commonly used in data analysis and exploratory data analysis (EDA) to understand the characteristics of a dataset, such as its central tendency, spread, and skewness.

Histograms can also be used to compare the distributions of two or more datasets. In this case, multiple histograms are overlaid on the same plot, with different colors or patterns used to distinguish between them.

One thing to note when interpreting histograms is that the choice of bin width can have an impact on the shape of the distribution. If the bin width is too narrow, the histogram may appear noisy and obscure the underlying patterns in the data. On the other hand, if the bin width is too wide, the histogram may oversimplify the distribution and miss important details. Therefore, choosing an appropriate bin width is an important consideration when creating a histogram.

## 2. Use a scatter plot

A scatter plot is a type of graph that displays the relationship between two continuous variables. It consists of a series of points, where each point represents a pair of values for the two variables. The x-axis represents one variable and the y-axis represents the other variable.

Scatter plots are useful for visualizing the distribution and relationship between two variables. They can reveal patterns and trends in the data, such as clusters, outliers, and correlations. Scatter plots can also help identify any nonlinear relationships between the variables.

In a scatter plot, the placement of each point represents the value of the two variables for that observation. If there is a strong correlation between the variables, the points will tend to cluster around a straight line or curve. If there is no correlation, the points will appear scattered randomly across the graph.

Scatter plots are commonly used in exploratory data analysis to investigate the relationship between variables and to identify any patterns or outliers in the data. They are also used in regression analysis to assess the fit of a linear or nonlinear model to the data.

When creating a scatter plot, it is important to choose appropriate scales for the x-axis and y-axis to accurately represent the range of values in the data. Additionally, labeling the axes and adding a title can help communicate the purpose and findings of the scatter plot.

### 3.PCA (Personal Computer Aid)

PCA (Principal Component Analysis) is a technique used in data analysis and machine learning to reduce the dimensionality of a dataset. It is a type of unsupervised learning method that extracts the most important features or components from a dataset and projects them onto a lower-dimensional space.

The main goal of PCA is to identify the underlying structure in the data by finding the directions of maximum variance in the dataset. These directions are called principal components, and they represent the axes that capture the most information in the data. The first principal component captures the most variance in the data, and each subsequent principal component captures the next highest amount of variance, subject to the constraint that it must be orthogonal to the previous components.

PCA works by performing a linear transformation of the original dataset onto a new set of orthogonal axes that are aligned with the principal components. This new set of axes can be used to represent the data in a lower-dimensional space without losing too much information about the original dataset. This is achieved by projecting the data onto the first few principal components, which capture the most variation in the dataset.

PCA can be useful in a variety of applications, such as image processing, signal processing, and natural language processing. It can help reduce the dimensionality of high-dimensional data, such as images or text, and make it easier to visualize and analyze the data. PCA can also be used as a preprocessing step before applying other machine learning algorithms, as it can help remove noise and redundancy in the data and improve the performance of the models.

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

Investigating data is necessary to understand the underlying patterns and characteristics of a dataset. It helps to identify any outliers, trends, correlations, or other patterns that may be present in the data. Investigating data can also help to identify any errors or anomalies in the data that may need to be corrected before further analysis.

Both qualitative and quantitative data can be explored and investigated to gain insights and understanding of the data. However, there may be some differences in the approaches used to explore these two types of data.

Qualitative data is typically non-numeric in nature and is often collected through methods such as interviews, surveys, and observations. To explore qualitative data, researchers often use techniques such as content analysis or thematic analysis, which involve identifying themes or patterns in the

data and categorizing them according to specific criteria. Qualitative data exploration is often more subjective and open-ended than quantitative data exploration, and it may require more interpretation and analysis.

Quantitative data, on the other hand, is typically numeric in nature and can be analyzed using statistical methods. To explore quantitative data, researchers often use techniques such as descriptive statistics, regression analysis, or hypothesis testing, which involve calculating summary statistics, modeling relationships between variables, or testing hypotheses about the data. Quantitative data exploration is often more objective and structured than qualitative data exploration, and it relies more heavily on mathematical and statistical techniques.

In summary, investigating data is necessary to gain a better understanding of the data and to identify any patterns, trends, or anomalies that may be present. The methods used to explore qualitative and quantitative data may differ, but both types of data can be explored and analyzed to gain valuable insights and understanding.

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

Histograms are graphical representations of the distribution of a set of continuous data. The shape of a histogram can reveal important information about the underlying distribution of the data. Here are some common histogram shapes:

1. Normal or Gaussian distribution: This is a symmetrical distribution with a bell-shaped curve in which most of the data falls within a central range and the data is evenly distributed around the mean.

2. Skewed distribution: This is a distribution in which the data is not symmetrical and is concentrated on one side of the histogram. There are two types of skewed distributions: left-skewed (negative skewness) and right-skewed (positive skewness).

3. Bimodal distribution: This is a distribution with two peaks, indicating that the data may come from two distinct sub-populations.

Bins are the intervals or ranges into which the continuous data is divided for the purpose of constructing a histogram. The width of the bins determines the level of detail in the histogram. A smaller bin width produces a more detailed histogram with more bins, while a larger bin width produces a less detailed histogram with fewer bins.

For example, if you have a dataset of heights of people, you could divide the data into bins of 5 feet each. The first bin might include heights from 5 feet to 5 feet 4 inches, the second bin might include heights from 5 feet 5 inches to 5 feet 9 inches, and so on. The number of data points falling within each bin is then plotted on the y-axis to construct the histogram.

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

Outliers are data points that are significantly different from the other data points in a dataset. They can occur due to errors in data collection or measurement, or due to rare events or extreme values in the data. Outliers can have a significant impact on statistical analysis and machine learning algorithms, as they can distort the results and affect the accuracy of the models.

Here are some common ways to deal with outliers:

1. Remove the outliers: One option is to remove the outliers from the dataset. However, this approach should be used with caution, as it can reduce the size of the dataset and affect the representativeness of the data. If the outliers are due to measurement errors or data entry errors, then removing them may be appropriate. However, if the outliers are due to genuine values in the data, then removing them may lead to biased results.

2. Transform the data: Another option is to transform the data using a mathematical function that reduces the impact of the outliers. For example, one can take the logarithm of the data to compress the range of the data and reduce the impact of extreme values. Box-Cox transformation and Winsorization are other techniques that can be used to transform the data and reduce the effect of outliers.

3. Treat the outliers separately: Depending on the context and the analysis being performed, it may be appropriate to treat the outliers separately. For example, in anomaly detection or fraud detection, the outliers may be of particular interest and should be treated as a separate class.

4. Use robust statistical methods: Robust statistical methods are designed to be less sensitive to outliers than traditional methods. For example, the median is a robust measure of central tendency that is less affected by outliers than the mean. Robust regression techniques, such as the Huber loss function, can also be used to fit models that are less influenced by outliers.

In summary, dealing with outliers requires careful consideration of the context and the nature of the data. Removing the outliers or treating them separately may be appropriate in some cases, but it is important to ensure that the data remains representative and unbiased. Transforming the data or using robust statistical methods can also be effective ways to reduce the impact of outliers on analysis and modeling.

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

Central inclination measures are statistical measures that summarize the central or typical value of a dataset. The three most commonly used measures of central tendency are:

1. Mean: The mean is the sum of all the data points divided by the number of data points. It is a sensitive measure of central tendency, meaning that it can be greatly influenced by outliers or extreme values in the dataset.

2. Median: The median is the middle value in a dataset, such that half of the values are above the median and half are below. It is a more robust measure of central tendency that is less affected by outliers or extreme values.

3. Mode: The mode is the value that appears most frequently in a dataset. It can be useful for datasets with discrete or categorical data, but may not be well-defined or meaningful for continuous data.

The mean can vary too much from the median in certain datasets when the distribution of the data is skewed or contains outliers. Skewed distributions occur when the data is not evenly distributed around the center, and can be left-skewed (negative skewness) or right-skewed (positive skewness). In a left-skewed distribution, the tail of the distribution is longer on the left side, which can pull the mean towards the left and result in a smaller mean compared to the median. In a right-skewed distribution, the tail of the distribution is longer on the right side, which can pull the mean towards the right and result in a larger mean compared to the median.

Outliers, on the other hand, can have a significant impact on the mean, but do not affect the median as much. An outlier that is much larger or smaller than the other values in the dataset can pull the mean towards it, resulting in a larger or smaller mean compared to the median.

In summary, the mean and median are both useful measures of central tendency, but the choice of measure depends on the distribution and characteristics of the data. The mean can vary too much from the median in certain datasets when the data is skewed or contains outliers. In such cases, the median is often a better measure of central tendency.

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

A scatter plot is a graphical representation of bivariate data, where each data point is represented as a point in a two-dimensional coordinate system. The x-axis represents one variable and the y-axis represents another variable. Scatter plots are useful for investigating the relationship between two variables and for identifying patterns, trends, and outliers in the data.

To investigate bivariate relationships using a scatter plot, we can plot the data points for both variables and visually inspect the plot for any patterns or trends. Depending on the nature of the relationship, we may observe a positive or negative correlation, no correlation, or a non-linear relationship. Positive correlation means that the two variables tend to increase or decrease together, negative correlation means that one variable tends to increase while the other decreases, and no correlation means that there is no systematic relationship between the variables.

In addition to identifying patterns and trends, scatter plots can also be used to detect outliers. Outliers are data points that are significantly different from the other data points in the dataset. They can occur due to errors in data collection or measurement, or due to rare events or extreme values in the data. In a scatter plot, outliers appear as data points that are far away from the other points in the plot. Outliers can be identified visually by inspecting the plot for any data points that are far away from the main cluster of points. If an outlier is detected, it may be necessary to investigate the cause of the outlier and decide whether to remove it from the dataset or treat it separately in the analysis.

In summary, a scatter plot is a useful tool for investigating bivariate relationships and identifying patterns, trends, and outliers in the data. By visually inspecting the plot, we can gain insights into the relationship between the two variables and detect any outliers that may be present in the data.

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

Cross-tabulation (or crosstab) is a technique used to explore the relationship between two categorical variables. It involves constructing a table that shows the frequency or proportion of observations that fall into each combination of categories for the two variables.

To create a cross-tab, we first identify the two variables of interest and their respective categories. We then create a table with the categories of one variable along the rows and the categories of the other variable along the columns. Each cell in the table represents the number or proportion of observations that fall into the corresponding combination of categories.

Cross-tabs can be used to figure out how two variables are related by examining the patterns of frequencies or proportions across the table. For example, if we are interested in the relationship between gender and political affiliation, we can create a cross-tab that shows the number or proportion of male and female respondents who identify as

Democrats, Republicans, or Independents. We can then examine the table for any patterns or trends that may suggest a relationship between the two variables. Some patterns that we may look for include:

- Differences in the distribution of categories between the two variables: If the distribution of categories is different for the two variables, it may suggest a relationship between the two variables. For example, if we find that a higher proportion of men identify as Republicans compared to women, it may suggest a relationship between gender and political affiliation.
- Similarities in the distribution of categories between the two variables: If the distribution of categories is similar for the two variables, it may suggest that the two variables are not related. For example, if we find that the distribution of political affiliation is similar between men and women, it may suggest that gender is not related to political affiliation.
- Differences in the strength or direction of the relationship between different combinations of categories: If the strength or direction of the relationship between the two variables varies across different combinations of categories, it may suggest a more complex relationship between the variables. For example, if we find that the relationship between gender and political affiliation is stronger among younger respondents compared to older respondents, it may suggest a more complex relationship that is influenced by age.

In summary, cross-tabs are a useful tool for exploring the relationship between two categorical variables. By examining the patterns of frequencies or proportions across the table, we can gain insights into how the two variables are related and identify any patterns or trends that may suggest a relationship between the variables.