

### **Q1. What is a parameter?**

Ans. A parameter is a numerical value that defines the behavior of a machine learning model, algorithm, or statistical analysis.

#### **Types of Parameters:**

1. Hyperparameters: Configurable settings that are set before training a model, such as learning rate, regularization strength, or batch size.
2. Model parameters: Values learned by the model during training, such as weights, biases, or coefficients.
3. Statistical parameters: Numerical values that describe a population or distribution, such as mean, variance, or standard deviation.

#### **Examples of Parameters:**

1. Linear Regression: coefficients (slope, intercept), regularization strength
2. Decision Trees: maximum depth, minimum samples per leaf
3. Neural Networks: learning rate, number of hidden layers, number of neurons per layer
4. Clustering: number of clusters (k-means), distance metric (Euclidean, Manhattan)

#### **Importance of Parameters:**

1. Model performance: Parameters significantly impact the accuracy, precision, and recall of machine learning models.
2. Overfitting and underfitting: Incorrect parameter settings can lead to overfitting (high variance) or underfitting (high bias).
3. Interpretability: Understanding the role of parameters can provide insights into the relationships between variables and the underlying mechanisms.

### **Q.2. What is correlation?**

#### **What does negative correlation mean?**

Ans. Correlation is a statistical measure that describes the relationship between two continuous variables. It assesses how strongly the variables tend to change together.

#### **Types of Correlation:**

1. Positive Correlation: As one variable increases, the other variable also tends to increase.

2. Negative Correlation: As one variable increases, the other variable tends to decrease.

3. No Correlation: The variables do not tend to change together.

### **Correlation Coefficient:**

The correlation coefficient (often denoted as  $r$ ) measures the strength and direction of the correlation. The values range from -1 (perfect negative correlation) to 1 (perfect positive correlation).

### **Interpretation of Correlation Coefficient:**

1. 0.00-0.30: Weak correlation
2. 0.31-0.60: Moderate correlation
3. 0.61-0.90: Strong correlation
4. 0.91-1.00: Very strong correlation.

Examples:

1. Height and weight: Positive correlation (as height increases, weight tends to increase)
2. Temperature and ice cream sales: Positive correlation (as temperature increases, ice cream sales tend to increase)

### **Q3. Define Machine Learning. What are the main components in Machine Learning?**

Ans. **Definition of Machine Learning:**

Machine Learning (ML) is a subset of Artificial Intelligence (AI) that involves training algorithms to learn patterns, relationships, and decision boundaries from data. The goal is to enable machines to make predictions, classify objects, or make decisions without being explicitly programmed.

### **Main Components of Machine Learning:**

1. Data: The foundation of ML is data. Relevant, high-quality data is necessary for training and testing ML models.

2. Algorithms: ML algorithms are the engines that power the learning process. Common algorithms include decision trees, random forests, support vector machines (SVMs), and neural networks.
3. Model: The ML model is the result of training an algorithm on data. The model learns patterns and relationships in the data and can make predictions or decisions.
4. Features: Features are individual characteristics or attributes of the data that are used to train the ML model.
5. Target Variable: The target variable is the outcome or response variable that the ML model is trying to predict or classify.
6. Hyperparameters: Hyperparameters are parameters that are set before training an ML model, such as learning rate, regularization strength, or batch size.
7. Evaluation Metrics: Evaluation metrics are used to assess the performance of an ML model, such as accuracy, precision, recall, F1-score, mean squared error (MSE), or R-squared.

#### ***Q4. How does loss value help in determining whether the model is good or not?***

Ans. Loss value, also known as the cost function or objective function, plays a crucial role in evaluating the performance of a machine learning model.

#### ***What is Loss Value?***

Loss value measures the difference between the model's predictions and the actual true values. It quantifies the model's error or inaccuracy.

#### ***Types of Loss Functions:***

1. Mean Squared Error (MSE)
2. Mean Absolute Error (MAE)
3. Cross-Entropy Loss
4. Binary Cross-Entropy Loss

#### ***How Loss Value Helps in Evaluating a Model:***

1. Convergence: A decreasing loss value indicates that the model is converging and improving.

2. Overfitting: If the loss value on the training set is significantly lower than on the validation set, it may indicate overfitting.
3. Underfitting: A high loss value on both training and validation sets may indicate underfitting.
4. Model Selection: Comparing loss values across different models can help select the best-performing model.

**Q5. What are continuous and categorical variables?**

Ans. In data analysis, variables can be classified into two main categories: continuous and categorical.

**Continuous Variables:**

Continuous variables are numerical variables that can take any value within a given range or interval. They can be measured to any level of precision and can have an infinite number of possible values.

**Examples:**

1. Height (measured in meters or feet)
2. Weight (measured in kilograms or pounds)
3. Temperature (measured in degrees Celsius or Fahrenheit)
4. Time (measured in seconds, minutes, or hours)

**Categorical Variables:**

Categorical variables, also known as discrete variables, are variables that can take only a limited number of distinct values or categories. These values are often labels or names rather than numerical values.

**Examples:**

1. Color (red, blue, green, etc.)
2. Gender (male, female, other, etc.)
3. Nationality (American, Canadian, Indian, etc.)
4. Product category (electronics, clothing, home goods, etc.)

**Q6. How do we handle categorical variables in Machine Learning? What are the common techniques?**

Ans. Handling categorical variables is a crucial step in machine learning, as many algorithms require numerical inputs. Here are common techniques to handle categorical variables:

### **1. Label Encoding**

Assigns a unique integer value to each category.

Example:

```
| Color | Label Encoding |  
| --- | --- |  
| Red | 0 |  
| Blue | 1 |  
| Green | 2 |
```

### **2. One-Hot Encoding (OHE)**

Creates a new binary feature for each category.

Example:

```
| Color | Red | Blue | Green |  
| --- | --- | --- | --- |  
| Red | 1 | 0 | 0 |  
| Blue | 0 | 1 | 0 |  
| Green | 0 | 0 | 1 |
```

### **3. Binary Encoding**

Similar to OHE, but uses binary digits (0s and 1s) to represent categories.

### **4. Hashing**

Uses a hash function to map categorical values to numerical values.

### **5. Embeddings**

Learned representations of categorical variables as dense vectors.

### **6. Target Encoding**

Replaces categorical values with a numerical value representing the target variable's mean or frequency.

### **7. Helmert Coding**

Compares each category to the mean of all previous categories.

### **8. Polynomial Coding**

Creates new features by raising the categorical variable to different powers.

### **When to use each technique:**

1. Label Encoding: Simple, ordinal categorical variables.

2. OHE: Nominal categorical variables with a small number of categories.
3. Binary Encoding: Nominal categorical variables with a large number of categories.
4. Hashing: High-cardinality categorical variables.
5. Embeddings: Complex, high-dimensional categorical variables.
6. Target Encoding: Categorical variables with a strong correlation with the target variable.
7. Helmert Coding: Ordinal categorical variables with a clear ordering.
8. Polynomial Coding: Categorical variables with a non-linear relationship with the target variable.

**Q7. What do you mean by training and testing a dataset?**

Ans. Training and testing a dataset are crucial steps in machine learning model development.

***Training a Dataset***

Training a dataset involves using a portion of the data to teach a machine learning model to make predictions or take actions. The goal is to enable the model to learn patterns, relationships, and decision boundaries from the data.

**Key Aspects of Training a Dataset**

1. Model learning: The model learns to map inputs to outputs based on the training data.
2. Parameter optimization: The model's parameters are adjusted to minimize the error between predicted and actual outputs.
3. Overfitting prevention: Techniques like regularization, dropout, and early stopping are used to prevent the model from becoming too specialized to the training data.

***Testing a Dataset***

Testing a dataset involves evaluating the trained model on a separate portion of the data to estimate its performance on unseen data. The goal is to assess the model's ability to generalize and make accurate predictions on new, unknown data.

**Key Aspects of Testing a Dataset**

1. Model evaluation: The model's performance is evaluated using metrics like accuracy, precision, recall, F1-score, mean squared error (MSE), or R-squared.
2. Generalization assessment: The model's ability to generalize to new data is assessed by evaluating its performance on the test set.

3. Hyperparameter tuning: The test set is used to tune hyperparameters and select the best-performing model.

### ***Dataset Splitting***

To train and test a dataset, it's common to split the data into three parts:

1. Training set (e.g., 60-80% of the data): Used to train the model.
2. Validation set (e.g., 10-20% of the data): Used to tune hyperparameters and monitor the model's performance during training.
3. Test set (e.g., 10-20% of the data): Used to evaluate the final model's performance and estimate its generalization ability.

By splitting the data into training, validation, and test sets, you can develop and evaluate a robust machine learning model that generalizes well to new, unseen data.

### ***Q8. What is sklearn.preprocessing?***

Ans. sklearn.preprocessing is a module in scikit-learn, a popular Python machine learning library. This module provides various functions and classes for preprocessing and transforming data.

### ***Key Functions and Classes***

#### 1. Scaling:

- StandardScaler: Scales features to have zero mean and unit variance.
- MinMaxScaler: Scales features to a specified range (e.g., 0 to 1).

#### 2. Normalization:

- Normalizer: Scales each sample to have unit norm (e.g., L1 or L2 norm).

#### 3. Encoding:

- LabelEncoder: Converts categorical labels to numerical labels.
- OneHotEncoder: Converts categorical features to one-hot encoded features.

#### 4. Transformation:

- PolynomialFeatures: Generates polynomial and interaction features.
- FunctionTransformer: Applies a custom function to each sample.

## 5. Imputation:

- SimpleImputer: Replaces missing values with a specified strategy (e.g., mean, median, or constant).

### ***Benefits of Preprocessing:***

1. Improved model performance: Preprocessing can help reduce the impact of outliers, scaling issues, and irrelevant features.
2. Increased interpretability: Preprocessing can help reveal underlying patterns and relationships in the data.
3. Enhanced data quality: Preprocessing can help detect and handle missing or erroneous data.

### ***Q9. What is a Test set?***

Ans. A test set, also known as a holdout set, is a portion of a dataset that is used to evaluate the performance of a machine learning model after it has been trained.

### ***Purpose of a Test Set:***

1. Evaluate model performance: Assess the model's ability to generalize to new, unseen data.
2. Estimate real-world performance: Provide an unbiased estimate of the model's performance in real-world scenarios.
3. Compare models: Compare the performance of different models or hyperparameter settings.

### ***Characteristics of a Test Set:***

1. Independent: The test set is separate from the training set.
2. Unseen: The test set is not used during training.
3. Representative: The test set is representative of the population or problem being modeled.

### ***Best Practices for Creating a Test Set:***

1. Split data randomly: Split the data into training, validation, and test sets using random sampling.
2. Use a suitable size: Allocate a sufficient portion of the data to the test set (e.g., 10-20%).



3. Maintain data distribution: Ensure the test set has a similar distribution to the training set.

### ***Q10. How do we split data for model fitting (training and testing) in Python?***

#### ***How do you approach a Machine Learning problem?***

Ans. In Python, you can split data for model fitting using the `train_test_split` function from the `sklearn.model_selection` module.

Basic Syntax:

```
from sklearn.model_selection import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Parameters:

1. X: Feature data (independent variables)
2. y: Target data (dependent variable)
3. test\_size: Proportion of data to use for testing (default=0.25)
4. random\_state: Seed for random number generation (default=None)
5. stratify: Whether to stratify the split based on the target variable (default=None)

Return Values:

1. X\_train: Training feature data
2. X\_test: Testing feature data
3. y\_train: Training target data
4. y\_test: Testing target data

### ***Here's a structured approach to tackling a machine learning problem:***

**Step 1: Problem Definition**

**Step 2: Data Collection and Exploration**

**Step 3: Data Preprocessing and Feature Engineering**

**Step 4: Model Selection and Training**

**Step 5: Model Evaluation and Hyperparameter Tuning**

## **Step 6: Model Deployment and Maintenance**

## **Step 7: Model Interpretation and Explanation**

### ***Q11. Why do we have to perform EDA before fitting a model to the data?***

Ans. Performing Exploratory Data Analysis (EDA) before fitting a model to the data is crucial for several reasons:

#### **Understanding Data Distribution**

1. Distribution of variables: EDA helps you understand the distribution of each variable, including central tendency, dispersion, and shape.
2. Identifying outliers: EDA enables you to detect outliers, which can significantly impact model performance.

#### **Identifying Relationships and Correlations**

1. Correlations between variables: EDA reveals correlations between variables, helping you identify potential relationships and interactions.
2. Multicollinearity: EDA detects multicollinearity, where two or more variables are highly correlated, which can lead to unstable models.

#### **Data Quality and Cleaning**

1. Missing values: EDA identifies missing values, enabling you to decide on imputation strategies or data cleaning methods.
2. Data errors: EDA detects data errors, such as inconsistent or invalid values.

#### **Informing Model Selection and Hyperparameter Tuning**

1. Model selection: EDA informs the choice of machine learning algorithm or model based on data characteristics.
2. Hyperparameter tuning: EDA provides insights into the optimal hyperparameter settings for the chosen model.

#### **Avoiding Common Pitfalls**

1. Assuming normality: EDA helps you avoid assuming normality of data distributions, which can lead to incorrect model assumptions.

2. Ignoring non-linear relationships: EDA detects non-linear relationships, ensuring you don't overlook important interactions between variables.

By performing EDA before modeling, you'll gain a deeper understanding of your data, make informed decisions about model selection and hyperparameter tuning, and avoid common pitfalls that can lead to poor model performance.

### **Q12. What is correlation?**

Ans. Correlation is a statistical measure that describes the relationship between two continuous variables. It measures how strongly the variables tend to change together.

#### **Types of Correlation:**

1. Positive Correlation: As one variable increases, the other variable also tends to increase.
2. Negative Correlation: As one variable increases, the other variable tends to decrease.
3. No Correlation: The variables do not tend to change together.

#### **Correlation Coefficient:**

The correlation coefficient is a numerical value that measures the strength and direction of the correlation. The most common correlation coefficient is the Pearson correlation coefficient ( $r$ ), which ranges from -1 (perfect negative correlation) to 1 (perfect positive correlation).

### **Q13. What does negative correlation mean?**

Ans. Negative Correlation: As one variable increases, the other variable tends to decrease.

### **Q14. How can you find correlation between variables in Python?**

Ans. In Python, you can find the correlation between variables using the `corr()` function from the Pandas library or the `corrcoef()` function from the NumPy library.

Method 1: Using Pandas

```
import pandas as pd
```

Create a sample DataFrame

```
data = {'Variable1': [1, 2, 3, 4, 5],  
        'Variable2': [2, 4, 6, 8, 10],  
        'Variable3': [3, 6, 9, 12, 15]}  
  
df = pd.DataFrame(data)  
  
correlation_matrix = df.corr()  
  
print(correlation_matrix)
```

Method 2: Using NumPy

```
import numpy as np  
  
# Create sample arrays  
  
variable1 = np.array([1, 2, 3, 4, 5])  
variable2 = np.array([2, 4, 6, 8, 10])  
variable3 = np.array([3, 6, 9, 12, 15])  
  
correlation_coefficient = np.corrcoef(variable1, variable2)[0, 1]  
  
print(correlation_coefficient)
```

In both examples, you'll get the correlation coefficient between the variables. The correlation coefficient ranges from -1 (perfect negative correlation) to 1 (perfect positive correlation).

**Q15. What is causation? Explain difference between correlation and causation with an example.**

Ans. **Causation:**

Causation refers to a relationship between two variables where changes in one variable (the cause) directly affect the other variable (the effect). In other words, causation implies that one variable is responsible for the changes in the other variable.

**Correlation vs. Causation:**

Correlation and causation are often confused, but they are not the same thing. Correlation measures the strength and direction of a linear relationship between two variables, while causation implies a direct cause-and-effect relationship.

**Example:**

Suppose we analyze the relationship between the number of ice cream cones sold (X) and the number of people wearing shorts (Y) in a city over a year.

**Correlation:**

We find a strong positive correlation between X and Y ( $r = 0.8$ ). This means that as the number of ice cream cones sold increases, the number of people wearing shorts also tends to increase.

**Causation:**

However, does eating ice cream cones cause people to wear shorts? Unlikely! A more plausible explanation is that both variables are influenced by a third variable: temperature. As the temperature rises, people are more likely to buy ice cream cones and wear shorts.

In this example, the correlation between X and Y is real, but it's not causal. Instead, both variables are influenced by a common underlying factor (temperature).

**Q16. What is an Optimizer? What are different types of optimizers? Explain each with an example**

Ans. An optimizer is an algorithm used to minimize or maximize a function, typically a loss function or an objective function, in machine learning and deep learning models.

Optimizers adjust the model's parameters to reduce the difference between predicted and actual outputs.

Here are some common types of optimizers, along with examples:

**1. Stochastic Gradient Descent (SGD)**

SGD is a basic optimizer that updates the model's parameters using the gradient of the loss function with respect to each parameter.

Example: Suppose we want to minimize the mean squared error (MSE) loss function using SGD.

$$w = w - \text{learning\_rate} * \text{gradient\_of\_loss\_with\_respect\_to\_}w$$

## 2. Mini-Batch Gradient Descent (MBGD)

MBGD is a variant of SGD that uses a batch of samples to compute the gradient instead of a single sample.

Example: Suppose we want to minimize the MSE loss function using MBGD with a batch size of 32.

$$w = w - \text{learning\_rate} * (1/\text{batch\_size}) * \text{sum}(\text{gradient\_of\_loss\_with\_respect\_to\_}w \text{ for each sample in batch})$$

## 3. Momentum

Momentum adds a fraction of the previous update to the current update, helping the optimizer escape local minima.

Example: Suppose we want to minimize the MSE loss function using Momentum with a learning rate of 0.01 and a momentum coefficient of 0.9.

$$v = 0.9 * v - 0.01 * \text{gradient\_of\_loss\_with\_respect\_to\_}w \quad w = w + v$$

## 4. Nesterov Accelerated Gradient (NAG)

NAG is a variant of Momentum that adds a fraction of the previous update to the current update, but also takes into account the future update.

Example: Suppose we want to minimize the MSE loss function using NAG with a learning rate of 0.01 and a momentum coefficient of 0.9.

$$v = 0.9 * v - 0.01 * \text{gradient\_of\_loss\_with\_respect\_to\_}w \quad w = w + 0.9 * v - 0.01 * \text{gradient\_of\_loss\_with\_respect\_to\_}w$$

## 5. Adagrad

Adagrad adapts the learning rate for each parameter based on the gradient's magnitude.

Example: Suppose we want to minimize the MSE loss function using Adagrad with an initial learning rate of 0.01.

$\text{cache} = 0$   
 $\text{cache} = \text{cache} + \text{gradient\_of\_loss\_with\_respect\_to\_w}^2$   
 $\text{learning\_rate} = 0.01 / \sqrt{\text{cache}}$   
 $w = w - \text{learning\_rate} * \text{gradient\_of\_loss\_with\_respect\_to\_w}$

## 6. RMSprop

RMSprop is a variant of Adagrad that divides the learning rate by an exponentially decaying average of squared gradients.

Example: Suppose we want to minimize the MSE loss function using RMSprop with an initial learning rate of 0.01 and a decay rate of 0.99.

$\text{cache} = 0.99 * \text{cache} + 0.01 * \text{gradient\_of\_loss\_with\_respect\_to\_w}^2$   
 $\text{learning\_rate} = 0.01 / \sqrt{\text{cache}}$   
 $w = w - \text{learning\_rate} * \text{gradient\_of\_loss\_with\_respect\_to\_w}$

## 7. Adam

Adam combines the benefits of Adagrad and RMSprop by adapting the learning rate for each parameter based on the first and second moments of the gradient.

**Example:** Suppose we want to minimize the MSE loss function using Adam with an initial learning rate of 0.01, a decay rate of 0.99, and a momentum coefficient of 0.9.

$m = 0.9 * m + 0.1 * \text{gradient\_of\_loss\_with\_respect\_to\_w}$   
 $v = 0.99 * v + 0.01 * \text{gradient\_of\_loss\_with\_respect\_to\_w}^2$   
 $\text{learning\_rate} = 0.01 / \sqrt{v}$   
 $w = w - \text{learning\_rate} * m$

Each optimizer has its strengths and weaknesses, and the choice of optimizer depends on the specific problem, model architecture, and dataset.

## Q17. What is sklearn.linear\_model ?

Ans. sklearn.linear\_model is a module in scikit-learn, a popular Python machine learning library. This module provides implementations of various linear models for regression and classification tasks.

## Key Classes and Functions:

### 1. Linear Regression:

- LinearRegression: Ordinary least squares linear regression.

- Ridge: Ridge regression with L2 regularization.
- Lasso: Lasso regression with L1 regularization.
- ElasticNet: Elastic net regression with L1 and L2 regularization.

## **2. Logistic Regression:**

- LogisticRegression: Logistic regression for binary classification.

## **3. Generalized Linear Models:**

- GeneralizedLinearModel: Base class for generalized linear models.

## **4. Stochastic Gradient Descent:**

- SGDClassifier: Stochastic gradient descent classifier.
- SGDRegressor: Stochastic gradient descent regressor.

## **5. Other Models:**

- OrthogonalMatchingPursuit: Orthogonal matching pursuit for sparse regression.
- BayesianRidge: Bayesian ridge regression.

### ***Q18. What does model.fit() do? What arguments must be given?***

Ans. model.fit() is a method in Keras (and other deep learning frameworks) that trains a machine learning model on a given dataset.

When you call model.fit(), the following steps occur:

1. **Compilation:** If the model hasn't been compiled yet, fit() will compile it with the default optimizer, loss function, and metrics.
2. **Data preparation:** The input data is prepared for training, which includes converting the data into the required format, normalizing or scaling the data (if necessary), and splitting the data into batches.
3. **Training loop:** The model is trained on the input data using the specified optimizer, loss function, and metrics. The training loop iterates over the batches of data, updating the model's weights and biases at each step.
4. **Evaluation:** After each epoch (a complete pass through the training data), the model's performance is evaluated on the validation data (if provided).
5. **Callbacks:** Any specified callbacks are executed at the end of each epoch.



## Arguments for `model.fit()`

Here are the main arguments you need to provide when calling `model.fit()`:

1. `x` (required): The input data. This can be a NumPy array, a Pandas DataFrame, or a TensorFlow dataset.
2. `y` (required): The target data. This should have the same number of samples as the input data.
3. `batch_size` (optional): The number of samples to include in each batch. Default is 32.
4. `epochs` (optional): The number of epochs to train the model. Default is 1.
5. `validation_data` (optional): A tuple containing the validation input data and target data.
6. `callbacks` (optional): A list of callbacks to execute during training.
7. `verbose` (optional): An integer indicating the level of verbosity. Default is 1.

## **Q19. What does `model.predict()` do? What arguments must be given?**

Ans. `model.predict()` is a method in Keras (and other deep learning frameworks) that uses a trained model to make predictions on new, unseen data.

When you call `model.predict()`, the following steps occur:

1. Input data preparation: The input data is prepared for prediction, which includes converting the data into the required format and normalizing or scaling the data (if necessary).
2. Forward pass: The input data is passed through the trained model, and the output is computed.
3. Prediction generation: The output of the model is used to generate predictions, which can be class labels, probabilities, or continuous values.

## Arguments for `model.predict()`

Here are the main arguments you need to provide when calling `model.predict()`:

1. `x` (required): The input data to make predictions on. This can be a NumPy array, a Pandas DataFrame, or a TensorFlow dataset.

2. `batch_size` (optional): The number of samples to include in each batch. Default is 32.
3. `verbose` (optional): An integer indicating the level of verbosity. Default is 0.

**Q20. What are continuous and categorical variables?**

Ans. Done earlier

**Q21. What is feature scaling? How does it help in Machine Learning?**

Ans. Feature scaling, also known as normalization or standardization, is a technique used to transform numeric features in a dataset to have similar scales, typically between 0 and 1.

1. This helps to prevent features with large ranges from dominating the model.

Feature scaling is crucial in machine learning for several reasons:

1. Prevents Feature Dominance: Features with large ranges can dominate the model, leading to poor performance. Scaling features ensures that all features contribute equally to the model.
2. Improves Model Interpretability: Scaled features make it easier to interpret model coefficients and feature importance.
3. Enhances Model Performance: Scaling features can improve model performance by reducing the impact of outliers and improving the stability of the model.
4. Facilitates Model Comparison: Scaled features enable fair comparison of models trained on different datasets or features.

**Q22. How do we perform scaling in Python?**

Ans. In Python, you can perform scaling using the `sklearn.preprocessing` module. Here are some common scaling techniques:

**Standardization (Z-Score Normalization)**

Standardization transforms the data to have a mean of 0 and a standard deviation of 1.

```
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

### **Min-Max Scaling (Normalization)**

Min-max scaling transforms the data to a common range, usually between 0 and 1.

```
from sklearn.preprocessing import MinMaxScaler
```

```
scaler = MinMaxScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

### **Robust Scaling**

Robust scaling is similar to standardization, but it uses the interquartile range (IQR) instead of the standard deviation.

```
from sklearn.preprocessing import RobustScaler
```

```
scaler = RobustScaler()
```

```
X_scaled = scaler.fit_transform(X)
```

### **Log Scaling**

Log scaling transforms the data using the logarithmic function.

```
import numpy as np
```

```
X_scaled = np.log(X)
```

Note that the choice of scaling technique depends on the specific problem and data distribution.

**Q23. What is *sklearn.preprocessing*?**

**Q24. How do we split data for model fitting (training and testing) in Python?**

Ans. Both done earlier

**Q25. Explain data encoding?**

Ans. Data encoding is the process of converting categorical or text data into numerical representations that can be processed by machine learning algorithms. Encoding is necessary because most machine learning algorithms require numerical inputs, but many real-world datasets contain categorical or text data.

***Types of Encoding Techniques:***

1. Label Encoding: Assigns a unique integer value to each category in a categorical variable.
2. One-Hot Encoding (OHE): Creates a new binary column for each category in a categorical variable.
3. Binary Encoding: Represents categorical data as binary vectors.
4. Hashing: Maps categorical data to numerical values using a hash function.
5. Ordinal Encoding: Assigns a numerical value to each category in an ordinal variable, preserving the order.