# Assignment Feature Engineering

1.      What is a parameter?

Parameters are the internal settings that the model learns from the training data. They are used to make predictions and define the model's skill. Parameters are learned during the training process and are specific to the model, adjusting to optimize performance on the given data.

-------------------------------------------------------------------------------------------------------------------------

2.      What is correlation? What does negative correlation mean?

Correlation refers to a statistical relationship between two variables, indicating how strongly they are related and in what direction.
A negative correlation specifically means that as one variable increases, the other variable tends to decrease, and vice versa. In other words, the variables move in opposite directions.

-------------------------------------------------------------------------------------------------------------------------

3.      Define Machine Learning. What are the main components in Machine Learning?

Machine learning (ML) is a subset of Artificial Intelligence (AI) that enables computer systems to learn from data and improve their performance on a specific task without being explicitly programmed for every detail. The main components in machine learning are: data, algorithms, models, and predictions.
        Summary of main components:
- **Data:**
This is the raw material for machine learning. Algorithms learn from data to identify patterns and make predictions. Data can be structured (like spreadsheets) or unstructured (like text or images).
- **Algorithms:**
These are the sets of rules or instructions that guide how the machine learning model learns from the data. Examples include linear regression, decision trees, neural networks, etc.
- **Models:**
The output of the learning process. A model is a mathematical representation of the patterns learned from the data. It's what the system uses to make predictions or decisions.
- **Predictions:**
The results produced by the model based on new input data. These predictions can be classifications (e.g., spam/not spam), regressions (e.g., predicting house prices), or other types of outputs, depending on the task.

In addition to these core components, the machine learning process also involves:
- **Representation:** How the data is organized and structured for the algorithm to process.
- **Evaluation:** Assessing the performance of the model using metrics relevant to the task.
- **Optimization:** Adjusting the model's parameters to improve its performance based on the evaluation results.

-------------------------------------------------------------------------------------------------------------------------

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

Loss value is a crucial metric for evaluating a machine learning model's performance. It quantifies the difference between the model's predictions and the actual values, with lower loss indicating better accuracy. By monitoring loss during training, you can determine if the model is learning effectively and identify potential issues like overfitting or underfitting.

-------------------------------------------------------------------------------------------------------------------------

5.      What are continuous and categorical variables?

In statistics, variables are classified as either continuous or categorical. Continuous variables represent measurements that can take on any value within a given range, while categorical variables represent categories or groups.

**Continuous Variables:**
- **Definition:**

Continuous variables are numerical variables that can be measured and can take on any value within a given range, including fractions and decimals.
- **Examples:**

Height, weight, temperature, time, and age are examples of continuous variables.

**Categorical Variables:**
- **Definition:**

Categorical variables represent qualitative data that can be divided into groups or categories.
- **Examples:**

Gender, eye colour, hair colour, types of fruit, and educational level are examples of categorical variables.

Further Categorization of Categorical Variables:
- **Nominal:** Categories have no inherent order (e.g., colours: red, blue, green).
- **Ordinal:** Categories have a meaningful order or ranking (e.g., education level: high school, bachelor's, master's).
- **Binary:** A special case of categorical variable with only two categories (e.g., yes/no, true/false).
-------------------------------------------------------------------------------------------------------------------------------

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

Categorical variables, which represent qualitative data, need to be transformed into a numerical format before being used in most machine learning algorithms. Common techniques for handling categorical variables include One-Hot Encoding, Label Encoding, and Ordinal Encoding.

Summary of these methods:

**1. One-Hot Encoding:**
- This technique is used when the categorical variable is nominal, meaning there is no inherent order or ranking among the categories.
- It creates a new binary column for each unique category in the variable.
- A value of 1 indicates the presence of that category, while 0 indicates its absence.
- For example, if a "color" variable has categories "red", "blue", and "green", one-hot encoding would create three new columns: "color_red", "color_blue", and "color_green", each with binary values.
- This method prevents the model from misinterpreting ordinal relationships between categories that don't exist.

**2. Label Encoding:**
- Label encoding assigns a unique integer to each category.
- This is suitable when the categorical variable has an ordinal relationship, meaning there is a meaningful order or ranking among the categories.
- For example, if a "size" variable has categories "small", "medium", and "large", label encoding might assign 1 to "small", 2 to "medium", and 3 to "large".
- However, if the ordinal relationship is not relevant, it's crucial to use another encoding method.

**3. Ordinal Encoding:**
- Ordinal encoding is similar to label encoding but specifically designed for ordinal categorical variables.
- It preserves the inherent order or ranking of the categories.
- This method is ideal when the categories have a clear order that needs to be reflected in the numerical representation.

**4. Other Techniques:**

- **Target Encoding:**
  Replaces categories with the mean of the target variable for that category. This is useful when there's a strong relationship between the categorical feature and the target variable.
- **Frequency Encoding:**
  Replaces categories with their frequency in the dataset. This can be helpful for high-cardinality features (many unique categories).
- **Binary Encoding:**
  Converts categories into binary code, reducing the dimensionality compared to one-hot encoding.
- **Hashing Encoding:**
  Maps categories to hash values, also reducing dimensionality, but with a risk of collisions (different categories mapping to the same value).

---------------------------------------------------------------------------------------------------------------------------------

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

In machine learning, training data is used to teach a model by exposing it to patterns and relationships in the data, while testing data is used to evaluate the model's performance on unseen data, providing an unbiased assessment of its ability to generalize.

---------------------------------------------------------------------------------------------------------------------------------

8.    What is sklearn.preprocessing?

The sklearn. preprocessing package provides several common utility functions and transformer classes to change raw feature vectors into a representation that is more suitable for the downstream estimators. In general, learning algorithms benefit from standardization of the data set

---------------------------------------------------------------------------------------------------------------------------------

9.    What is a Test set?

In machine learning, a test set is a subset of data held back from the training process. It's used to evaluate how well a trained model generalizes to unseen data, providing an unbiased estimate of its performance on real-world scenarios. The test set should not be used during the training phase or for hyperparameter tuning.

---------------------------------------------------------------------------------------------------------------------------------

10.    How do we split data for model fitting (training and testing) in Python? How do you approach a Machine Learning problem?

To split data for model fitting in Python, the train_test_split function from sklearn.model_selection is commonly used. This function divides the data into training and testing sets, allowing you to train your model on one set and evaluate its performance on the other. A typical machine learning problem involves defining the problem, gathering and preparing data, choosing a model, training, evaluating, and then potentially tuning the model and making predictions.

**Splitting Data with train_test_split**
The train_test_split function takes your data (features and target variables) and splits it into training and testing sets. It also allows you to specify the proportion of data to be used for testing (e.g., 20% or 25%).

The random_state parameter is used to ensure reproducibility, meaning you'll get the same split every time you run the code with the same seed.

**Approaching a Machine Learning Problem**
A typical machine learning problem can be broken down into the following steps:

1. **Define the Problem:** Clearly identify what you want to achieve (e.g., predict house prices, classify images).
2. **Gather Data:** Collect the relevant data for your problem.
3. **Prepare the Data:** Clean, preprocess, and transform the data, including handling missing values and outliers, and potentially feature engineering.

4. **Split the Data:** Divide the data into training, validation (optional), and testing sets.
5. **Choose a Model:** Select an appropriate machine learning model based on the problem type and data characteristics.
6. **Train the Model:** Train the model using the training data.
7. **Evaluate the Model:** Assess the model's performance on the test set using appropriate evaluation metrics.
8. **Tune the Model:** Adjust model parameters (hyperparameter tuning) to improve performance, often using the validation set if available.
9. **Make Predictions:** Once the model is trained and evaluated, use it to make predictions on new, unseen data.

---

11. Why do we have to perform EDA before fitting a model to the data?

Exploratory Data Analysis (EDA) should be performed before fitting a model to the data to gain a comprehensive understanding of the dataset, identify potential issues, and make informed decisions about data preparation and model selection. EDA helps uncover patterns, anomalies, and relationships within the data, ultimately leading to a more robust and accurate model.

Summary for why EDA is crucial before model fitting:
- **Understanding the Data:**
  EDA provides insights into the structure, characteristics, and distribution of the data. This includes identifying the types of variables, their distributions, and potential outliers.
- **Data Cleaning:**
  EDA helps in identifying and addressing missing values, inconsistencies, and outliers that can negatively impact model performance.
- **Feature Engineering:**
  By understanding the relationships between variables, EDA can guide feature engineering, which involves creating new features or transforming existing ones to improve model accuracy.
- **Model Selection:**
  EDA can inform the selection of appropriate models based on the data's characteristics and relationships.
- **Hypothesis Testing:**
  EDA allows for testing assumptions and hypotheses about the data before building a model, ensuring that the model is built on a solid foundation.
- **Preventing Data Leakage:**
  Performing EDA on the entire dataset before splitting it into training and testing sets prevents data leakage, where information from the test set inadvertently influences the training process, leading to overly optimistic model performance.

---

12. What is correlation?

correlation refers to a statistical measure that describes the strength and direction of a relationship between two or more variables. It indicates how much two variables change together, and can be used to identify patterns and dependencies within a dataset. A strong correlation suggests that changes in one variable are associated with predictable changes in another, while a weak correlation implies less predictable relationships.

---

13. What does negative correlation mean?

A negative correlation, also known as an inverse correlation, means that two variables tend to move in opposite directions; when one variable increases, the other tends to decrease, and vice versa. This relationship is not necessarily causal, but it indicates an observable pattern where the values of the two variables change in opposing directions.

---

14. How can you find correlation between variables in Python?

Finding the correlation between variables in Python is commonly achieved using the pandas and numpy libraries.

**1. Using Pandas for DataFrames:**
The most straightforward way to calculate correlations for multiple variables in a dataset is by using the .corr() method on a pandas DataFrame. This method computes the pairwise correlation between all numeric columns, resulting in a correlation matrix.

**2. Using NumPy for Arrays:**
For calculating the correlation between two specific arrays or series, the numpy.corrcoef() function is useful.

---------------------------------------------------------------------------------------------------------------------

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

Causation means one event directly causes another, establishing a cause-and-effect relationship. Correlation simply indicates a relationship or association between two variables, without necessarily implying that one causes the other. A common example is that while ice cream sales and days temperature are correlated (both increase in the summer), one does not cause the other; rather, they are both influenced by a third variable, the sunny weather.

---------------------------------------------------------------------------------------------------------------------

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

An optimizer is an algorithm that adjusts the parameters (weights and biases) of a neural network to minimize a loss function, thus improving the model's accuracy. It's a crucial component in training neural networks, as it guides the model towards finding the optimal set of parameters that best fit the training data.

Given below are different optimizer types:
**1. Gradient Descent:**
- **Concept:**
  Gradient descent is the foundation for many optimizers. It works by iteratively adjusting parameters in the direction of the negative gradient of the loss function, effectively "descending" the error surface to reach the minimum.
- **Variants:**
  - **Batch Gradient Descent:** Computes the gradient using the entire training dataset in each iteration, leading to stable but potentially slow convergence.
  - **Stochastic Gradient Descent (SGD):** Computes the gradient using a single randomly selected training example in each iteration. This makes it faster than batch gradient descent but can introduce more noise.
  - **Mini-batch Gradient Descent:** Uses a small random subset (mini-batch) of the training data to compute the gradient, striking a balance between the stability of batch gradient descent and the speed of SGD.
- **Example:**
  Imagine a ball rolling down a hill. Gradient descent is like guiding the ball to the bottom of the hill by taking small steps in the direction of the steepest descent.

**2. Momentum:**
- **Concept:**
  Momentum adds a "memory" to the optimization process. It accumulates past gradients to accelerate convergence and dampen oscillations, especially in the presence of noisy gradients.
- **Example:**
  Imagine the ball rolling down the hill. Momentum acts like giving the ball a push, making it move faster and smooth out its path, especially when encountering small bumps (noisy gradients).

### 3. Adaptive Learning Rate Optimizers:

- **Concept:**

  These optimizers dynamically adjust the learning rate for each parameter, making them more efficient for complex datasets and avoiding issues like slow convergence or getting stuck in local minima.

- **Variants:**

  - **Adagrad:** Adjusts the learning rate for each parameter based on the sum of past squared gradients. It's effective for sparse data but can lead to a rapidly decreasing learning rate.
  - **RMSprop:** Addresses Adagrad's issue by using a decaying average of past squared gradients, preventing the learning rate from dropping too quickly.
  - **Adam:** Combines the benefits of momentum and RMSprop, using both exponentially decaying averages of past gradients and past squared gradients to adaptively adjust learning rates. It's one of the most popular optimizers.
  - **Adadelta:** An extension of Adagrad that addresses the aggressive learning rate reduction issue by restricting the learning rate to a small value.

- **Example:**

  Consider a robot trying to navigate a maze. Adaptive learning rate optimizers are like giving the robot different instructions for each path it takes, allowing it to adjust its steps based on how easy or difficult the path is.

### 4. Other Optimizers:

- **Nesterov Accelerated Gradient (NAG):**

  A variant of SGD with momentum that looks ahead to where the momentum would carry the parameters before computing the gradient, leading to faster convergence.

- **AdamW:**

  A variant of Adam that decouples weight decay from the optimization process, improving generalization performance.

  In essence, the choice of optimizer depends on the specific task, the characteristics of the data, and the desired training speed and stability.

-------------------------------------------------------------------------------------------------------------------------

17.     What is sklearn.linear_model?

sklearn.linear_model is a sub-module within the scikit-learn (or sklearn) library in Python, which is a widely used library for machine learning. This sub-module specifically houses a variety of linear models used for both regression and classification tasks. Its a collection of generalized linear model fitting systems within scikit-learn. linear_model is a class of the sklearn module if contain different functions for performing machine learning with linear models. The term linear model implies that the model is specified as a linear combination of features

-------------------------------------------------------------------------------------------------------------------------

18.     What does model.fit() do? What arguments must be given?

Model Fitting is a measurement of how well a machine learning model adapts to data that is similar to the data on which it was trained. The fitting process is generally built-in to models and is automatic. A well-fit model will accurately approximate the output when given new data, producing more precise results.

-------------------------------------------------------------------------------------------------------------------------

19.     What does model.predict() do? What arguments must be given?

model. predict() is used to generate predictions from the trained model based on new input data. It does not require true labels and does not compute any metrics.

-------------------------------------------------------------------------------------------------------------------------

20. What are continuous and categorical variables?

In statistics, variables are classified as either continuous or categorical. Continuous variables represent measurements that can take on any value within a given range, while categorical variables represent categories or groups.
Continuous variables are numerical variables that can be measured and can take on any value within a given range, including fractions and decimals.
Categorical variables represent qualitative data that can be divided into groups or categories.

---------------------------------------------------------------------------------------------------------------------------------

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

Feature scaling is a data preprocessing technique that transforms the numerical features of a dataset to a similar scale. This process is crucial for many machine learning algorithms because it prevents features with larger values from dominating the learning process and ensures that all features contribute equally to the model's performance. By scaling features, algorithms can converge faster and achieve better accuracy.

---------------------------------------------------------------------------------------------------------------------------------

22. How do we perform scaling in Python?

In Python, "scaling" typically refers to feature scaling in the context of machine learning, which involves transforming numerical features to a similar range. This is crucial for many machine learning algorithms that are sensitive to the magnitude of features. The most common methods are Normalization (Min-Max Scaling) and Standardization.

**1. Normalization (Min-Max Scaling)**
Normalization scales features to a specific range, usually between 0 and 1. This is useful when the distribution of the data is not Gaussian or when algorithms require features within a bounded range.

**2. Standardization**
Standardization transforms features to have a mean of 0 and a standard deviation of 1. This is particularly useful for algorithms that assume a Gaussian distribution or are sensitive to feature scales, such as those relying on distance calculations (e.g., K-Means, KNN, PCA).

---------------------------------------------------------------------------------------------------------------------------------

23. What is sklearn.preprocessing?

sklearn.preprocessing is a module within the scikit-learn (sklearn) library in Python that provides a collection of tools and functions for data preprocessing. Data preprocessing is a crucial step in machine learning workflows, involving transforming raw data into a format suitable for use by machine learning algorithms.

The sklearn.preprocessing module offers various techniques to achieve this, including:
- **Scaling:**
Rescaling numerical features to a specific range (e.g., Min-Max Scaling) or standardizing them to have zero mean and unit variance (e.g., StandardScaler, RobustScaler). This helps prevent features with larger magnitudes from dominating the learning process.
- **Normalization:**
Scaling individual samples to have unit norm (e.g., L1 or L2 normalization). This is useful in scenarios where the direction of the data vector is more important than its magnitude.
- **Encoding Categorical Features:**
Converting categorical data (e.g., "red", "green", "blue") into numerical representations that machine learning algorithms can understand (e.g., OneHotEncoder, OrdinalEncoder).
- **Imputation:**
Handling missing values in the dataset by replacing them with estimated values (e.g., SimpleImputer).
- **Generating Polynomial Features:**
Creating new features by taking polynomial combinations of existing features, which can help capture non-linear relationships in the data.

These preprocessing techniques aim to improve the performance and stability of machine learning models by addressing issues such as varying scales, categorical data, and missing values in the raw dataset.

---------------------------------------------------------------------------------------------------------------------

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

To split data for model fitting, training, and testing in Python, the train_test_split function from the sklearn.model_selection module is commonly used. This function divides the dataset into a training set (used to train the model) and a testing set (used to evaluate the trained model's performance on unseen data).

---------------------------------------------------------------------------------------------------------------------

25.     Explain data encoding?

Data encoding is the process of converting data into a specific format for efficient storage, transmission, or processing. It's essentially transforming data into a form that a computer or system can understand and work with, often involving representing information in a numerical or standardized format.

**Brief Description:**
- **Purpose:**
Encoding ensures data is in a suitable format for various operations like storage, transmission (e.g., over a network), and analysis by algorithms.
- **Standardization:**
It often involves converting data into a standardized format like Unicode for text, allowing different systems to interpret the same data consistently.
- **Efficiency:**
Encoding can optimize data for specific tasks. For example, categorical data is often encoded numerically (like one-hot encoding) to be used effectively in machine learning models.
- **Decoding:**
The reverse process of encoding is decoding, which brings the data back to its original form for human readability or other purposes.
Examples:
- **Text Encoding:** Converting characters into a numerical representation like ASCII or Unicode.
- **Image Encoding:** Converting image data into a format like JPEG or PNG.
- **Audio Encoding:** Converting audio signals into formats like MP3 or WAV.
- **Categorical Data Encoding:** Representing categories (e.g., colors, product types) numerically for machine learning models.

In Essence: Encoding is about making data usable and efficient for computers and systems, ensuring that information can be stored, transmitted, and processed effectively.

---------------------------------------------------------------------------------------------------------------------