# "Feature Engineering"

## Assignment Questions and Answers:

### 1. What is a parameter?

A parameter is a variable used to define a function, model, or algorithm. In Machine Learning, parameters are the internal values that a model learns during training, such as weights and biases in neural networks. These parameters determine how the model makes predictions based on input data.

### 2. What is correlation?

Correlation is a statistical measure that describes the strength and direction of a linear relationship between two variables. It is represented by the correlation coefficient, which ranges from -1 to 1. A value closer to 1 indicates a strong positive relationship, while a value closer to -1 indicates a strong negative relationship. A value near 0 suggests no linear correlation.

### 3. What does negative correlation mean?

Negative correlation occurs when an increase in one variable is associated with a decrease in another variable. The correlation coefficient for negative correlation lies between 0 and -1. For example, as the price of a product increases, its demand may decrease, showing a negative correlation.

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

Machine Learning is a subset of artificial intelligence that enables systems to learn patterns from data and make decisions or predictions without being explicitly programmed.

**Main components in Machine Learning:**

1. **Data**: The foundation of Machine Learning; it can be structured or unstructured.
2. **Features**: Relevant attributes or variables extracted from data.
3. **Model**: A mathematical representation used to predict outcomes based on input data.
4. **Training**: The process of teaching the model using historical data.
5. **Evaluation**: Assessing the model's performance using metrics.
6. **Hyperparameters**: Configuration settings that influence the training process, such as learning rate or batch size.

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

The loss value measures the difference between the predicted output and the actual target values. A lower loss value indicates that the model is performing well, while a higher loss value suggests poor performance. The loss function helps in optimizing the model during training by adjusting parameters to minimize this value.

## *6. What are continuous and categorical variables?*

- **Continuous variables**: Variables that can take any numerical value within a range. Examples include height, weight, and temperature.
- **Categorical variables**: Variables that represent distinct categories or groups. Examples include gender (male, female) and colors (red, blue, green).

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

Categorical variables need to be converted into numerical formats to be processed by Machine Learning algorithms. Common techniques include:

1. **Label Encoding**: Assigns unique numerical values to each category.
2. **One-Hot Encoding**: Creates binary columns for each category, with 1 indicating the presence of a category and 0 indicating its absence.

3. **Ordinal Encoding**: Assigns ordered numerical values to categories based on their rank.

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

- **Training dataset**: A subset of the data used to train the model and adjust its parameters.
- **Testing dataset**: A separate subset used to evaluate the model's performance on unseen data to ensure it generalizes well.

## 9. What is sklearn.preprocessing?

`sklearn.preprocessing` is a module in Scikit-learn that provides tools for preprocessing data. These tools include methods for scaling, normalizing, encoding categorical variables, and transforming data to make it suitable for Machine Learning algorithms.

## 10. What is a Test set?

A test set is a subset of the dataset that is used to evaluate the performance of a trained Machine Learning model. It consists of unseen data to ensure the model's predictions generalize well to new, real-world scenarios.

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

We use the `train_test_split` function from Scikit-learn to split data:

Code:

```
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)
```

Here:

- X represents the features.
- y represents the target variable.
- `test_size=0.2` reserves 20% of the data for testing.
- `random_state` ensures reproducibility.

## 12. How do you approach a Machine Learning problem?

1. **Define the problem**: Understand the objective and expected outcomes.
2. **Collect and preprocess data**: Gather relevant data and clean it by handling missing values, outliers, and inconsistencies.
3. **Perform Exploratory Data Analysis (EDA)**: Analyze data patterns, correlations, and distributions.
4. **Feature engineering**: Extract or create meaningful features.
5. **Select a model**: Choose an appropriate algorithm based on the problem type (e.g., regression, classification).
6. **Train the model**: Fit the model to the training data.
7. **Evaluate the model**: Use metrics like accuracy, precision, or RMSE to assess performance.
8. **Optimize the model**: Fine-tune hyperparameters and retrain.
9. **Deploy the model**: Integrate the trained model into production.

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

EDA helps in:

- Understanding data distributions and patterns.
- Identifying and handling outliers and missing values.
- Discovering relationships and correlations between variables.
- Selecting and engineering features to improve model performance.
- Avoiding potential biases or issues in the data.

## 14. What is correlation?

Correlation is a statistical measure that describes the strength and direction of a linear relationship between two variables. It is represented by the correlation coefficient, which

ranges from -1 to 1. A value closer to 1 indicates a strong positive relationship, while a value closer to -1 indicates a strong negative relationship. A value near 0 suggests no linear correlation.

## 15. What does negative correlation mean?

Negative correlation occurs when an increase in one variable is associated with a decrease in another variable. The correlation coefficient for negative correlation lies between 0 and -1. For example, as the price of a product increases, its demand may decrease, showing a negative correlation.

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

Using Pandas:

```
import pandas as pd
correlation_matrix = df.corr()
print(correlation_matrix)
```

This computes the pairwise correlation between numerical columns in a DataFrame.

## 17. What is causation? Explain the difference between correlation and causation with an example.

- **Causation**: One variable directly affects another.
- **Correlation**: A statistical association between two variables, but it does not imply causation.

**Example**: Ice cream sales and drowning incidents are correlated because both increase during summer. However, ice cream sales do not cause drowning incidents.

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

An optimizer adjusts the model parameters to minimize the loss function during training.

**Types of optimizers:**

1. **SGD (Stochastic Gradient Descent)**: Updates parameters using the gradient of the loss function.
2. **Adam (Adaptive Moment Estimation)**: Combines momentum and RMSProp for faster convergence.
3. **RMSProp**: Uses moving averages of squared gradients to adjust learning rates.

### 19. What is sklearn.linear_model?

`sklearn.linear_model` is a module in Scikit-learn that provides implementations of linear models, including:

- Linear Regression
- Logistic Regression
- Ridge and Lasso Regression

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

`model.fit()` trains a Machine Learning model by adjusting its parameters to minimize the loss function. **Arguments**:

- X: Features or input data.
- y: Target variable or labels.

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

`model.predict()` generates predictions for input data based on the trained model. **Arguments**:

- X: Input features for which predictions are required.

### 22. What are continuous and categorical variables?

- **Continuous variables**: Variables that can take any numerical value within a range. Examples include height, weight, and temperature.
- **Categorical variables**: Variables that represent distinct categories or groups. Examples include gender (male, female) and colors (red, blue, green).

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

Feature scaling normalizes or standardizes data to ensure all features contribute equally to the model. It helps improve algorithm performance, especially for distance-based models like KNN or gradient-based models like neural networks.

### 24. How do we perform scaling in Python?

Using Scikit-learn:

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data)
```

This scales the data to have a mean of 0 and a standard deviation of 1.

### 25. What is sklearn.preprocessing?

`sklearn.preprocessing` is a module in Scikit-learn that provides tools for preprocessing data. These tools include methods for scaling, normalizing, encoding categorical variables, and transforming data to make it suitable for Machine Learning algorithms.

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

We use the `train_test_split` function from Scikit-learn to split data:

Code:

```
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)
```

Here:

- X represents the features.
- y represents the target variable.
- `test_size=0.2` reserves 20% of the data for testing.
- `random_state` ensures reproducibility.

### 27. Explain data encoding?

Data encoding is the process of converting categorical variables into numerical formats for Machine Learning algorithms. Common techniques include:

1. **Label Encoding**: Assigns a unique number to each category.
2. **One-Hot Encoding**: Creates binary columns for each category.
3. **Ordinal Encoding**: Assigns ordered values to categories based on their rank.