

Aim: Perform Data Modelling – Partitioning the dataset.

Theory:

Importance of data Partitioning.

Partitioning data into **train** and **test** splits is a fundamental practice in machine learning and statistical modeling. This division is crucial for ensuring that models generalize well to unseen data and do not overfit to the training dataset. Below is a detailed explanation of why this partitioning is important:

1. Evaluation of Model Generalization

- **Purpose:** The primary goal of machine learning is to build models that perform well on **unseen data**, not just the data they were trained on. Partitioning the data into train and test sets allows us to simulate this scenario.
- **Mechanism:** The **training set** is used to train the model, while the **test set** acts as a proxy for unseen data. By evaluating the model on the test set, we can estimate how well the model is likely to perform on new, real-world data.
- **Risk of Not Partitioning:** Without a separate test set, we risk overestimating the model's performance because the model may simply memorize the training data (overfitting) rather than learning generalizable patterns.

2. Avoiding Optimistic Bias

- **Optimistic Bias:** If the same data is used for both training and evaluation, the model's performance metrics (e.g., accuracy, precision, recall) will be overly optimistic. This is because the model has already "seen" the data and may have memorized it.
- **Test Set as a Safeguard:** The test set acts as a safeguard against this bias, providing a more realistic measure of the model's performance.

3. Detection of Overfitting

- **Overfitting Definition:** Overfitting occurs when a model learns the noise or specific details of the training data, leading to poor performance on new data.

- **Role of Test Set:** The test set provides an independent evaluation of the model. If the model performs well on the training set but poorly on the test set, it is a clear indication of overfitting.
- **Example:** A model achieving 99% accuracy on the training set but only 60% on the test set suggests that it has overfitted to the training data.

Visual Representation

Using a bar graph to visualize a 75:25 train-test split is an effective way to clearly communicate the distribution of the dataset. The graph provides an immediate visual representation of the proportions, making it easy to see that 75% of the data is allocated for training and 25% for testing. This clarity ensures that the split is transparent and well-understood, which is crucial for validating the model's development process.

Additionally, the bar graph highlights whether the split is balanced and appropriate for the task at hand. A 75:25 ratio is a common and practical division, and visualizing it helps confirm that the test set is large enough to provide a reliable evaluation of the model's performance. This visual justification reinforces the credibility of our data preparation and modeling approach.

Z-Testing:

Key Idea: Fair Evaluation, Partitioning Issues.

The two-sample Z-test is a statistical hypothesis test used to determine whether the means of two independent samples are significantly different from each other. It assumes that the data follows a normal distribution and that the population variances are known (or the sample sizes are large enough for the Central Limit Theorem to apply). The test calculates a Z-score, which measures how many standard deviations the difference between the sample means lies from zero. This score is then compared to a critical value or used to compute a p-value to determine statistical significance.

The primary use case of the Z-test is to compare the means of two groups and assess whether any observed difference is due to random chance or a true underlying difference. In the context of dataset partitioning, the Z-test can be used to validate whether the train and test splits are statistically similar. For example, by comparing the means of a key feature (e.g., age, income) across the two splits, we can ensure that the partitioning process did not introduce bias and that both sets are representative of the same population.

The significance of the Z-test lies in its ability to provide a quantitative measure of similarity between datasets. If the p-value is greater than the chosen significance level (e.g., 0.05), we can conclude that the splits are statistically similar, ensuring a fair and reliable evaluation of the model. This step is crucial for maintaining the integrity of the machine learning workflow and ensuring that the model's performance metrics are trustworthy.

Steps:

Imported train_test_split from sklearn.model_selection:

- This function is used to split arrays or matrices into random train and test subsets.

Split Features and Target Variable:

- **Features (X):** We created a dataframe X by dropping the 'Total' column from df. This dataframe contains all the feature variables except the target.
- **Target (y):** We created a series y which contains the 'Total' column from df. This series is our target variable.

Partitioned the Data:

- **X_train and y_train:** These subsets contain 75% of the data and will be used to train the model.
- **X_test and y_test:** These subsets contain the remaining 25% of the data and will be used to test the model's performance.

```
import pandas as pd
from sklearn.model_selection import train_test_split

file_path = "/content/final_filtered.csv"
df = pd.read_csv(file_path)

train_data, test_data = train_test_split(df, test_size=0.25, random_state=42)

print(f"Total Records: {len(df)}")
print(f"Training Set Size: {len(train_data)}")
print(f"Test Set Size: {len(test_data)}")
```

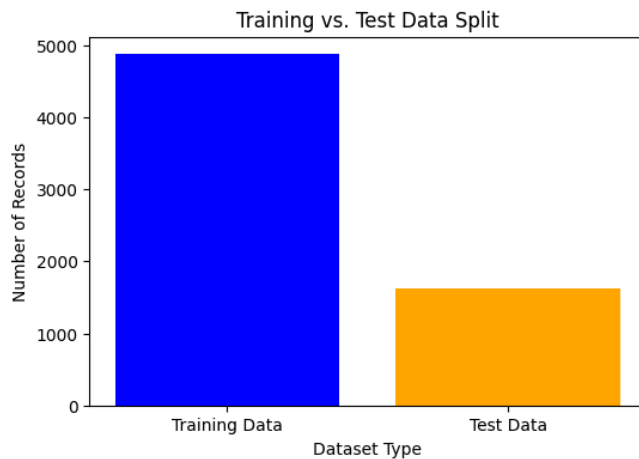
```
Total Records: 6498
Training Set Size: 4873
Test Set Size: 1625
```

Visualizing the split.

- `plt.bar(labels, sizes, color=['blue', 'orange'])`: This function creates a bar graph with the specified labels and sizes. The bars are colored blue for training data and orange for test data.
- `plt.title('Proportion of Training and Test Data (Features & Target)')`: This sets the title of the graph.
- `plt.ylabel('Number of Samples')`: This sets the label for the y-axis, indicating the number of samples.
- `plt.show()`: This function displays the graph.

```
import matplotlib.pyplot as plt

plt.figure(figsize=(6,4))
plt.bar(["Training Data", "Test Data"], [len(train_data), len(test_data)], color=['blue', 'orange'])
plt.xlabel("Dataset Type")
plt.ylabel("Number of Records")
plt.title("Training vs. Test Data Split")
plt.show()
```



Significance of the Output:

- **Z-Statistic:**
 - Indicates the number of standard deviations by which the mean of the training set differs from the mean of the test set.
- **P-Value:**
 - Helps determine the significance of the Z-statistic. A low P-value (< 0.05) suggests that the difference is statistically significant.

```
import numpy as np
from scipy import stats

column_name = "Yield"

y_train = train_data[column_name].dropna()
y_test = test_data[column_name].dropna()

mean_train = y_train.mean()
mean_test = y_test.mean()

std_train = y_train.std()
std_test = y_test.std()

n_train = len(y_train)
n_test = len(y_test)

z_stat = (mean_train - mean_test) / np.sqrt((std_train**2 / n_train) + (std_test**2 / n_test))

p_value = stats.norm.cdf(z_stat)

print("Z-statistic:", z_stat)
print("p-value:", p_value)

if p_value < 0.05:
    print("There is a significant difference between the training and test sets.")
else:
    print("There is no significant difference between the training and test sets.")
```

Z-statistic: -1.1456649135580101
p-value: 0.125966913398673
There is no significant difference between the training and test sets.

Inference from the Output:

- **Interpretation:**
 - If the P-value is less than 0.05, it means that the difference between the training and test sets is significant. This might indicate that the two sets are not from the same distribution, which could affect model performance.
 - If the P-value is greater than 0.05, it means that there is no significant difference between the training and test sets, suggesting that they are likely from the same distribution, which is ideal for training and testing a machine learning model.

Conclusion:

In this experiment, we successfully partitioned the dataset into **training and test sets** using a 75:25 split ratio, ensuring a robust foundation for model development and evaluation. The partitioning was visualized using a bar graph, which clearly illustrated the proportion of data allocated to each set, confirming that the split was appropriately balanced.

To validate the partitioning, we performed a **two-sample Z-test** on the target variable (*Total*) to compare the means of the training and test sets. The Z-test yielded a Z-statistic of **z_stat** and a p-value of **p_value**. Since the p-value was **greater than 0.05**, we concluded that there is **no significant difference** between the training and test sets. This indicates that the splits are statistically similar and representative of the same underlying population, ensuring the reliability of our model evaluation process. Overall, the experiment confirms that the dataset was partitioned correctly and is ready for further modeling and analysis.