


✓ 1.Import Necessary libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score
```

```
from google.colab import drive
drive.mount('/content/drive')
```


```
from google.colab import files
uploaded = files.upload()
```

 No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving signal-data.csv to signal-data.csv

✓ 2.Load the Data

```
import pandas as pd
df = pd.read_csv('/content/signal-data.csv')
```


```
from google.colab import drive
drive.mount('/content/drive')
```

 Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

```
data = pd.read_csv('signal-data.csv')
```


✓ 3.Explore the data

```
import pandas as pd
data = pd.read_csv('signal-data.csv')
data.head() #view first few rows
data.info() #check data types and missing values
```

 `<class 'pandas.core.frame.DataFrame'>`
RangeIndex: 1567 entries, 0 to 1566
Columns: 592 entries, Time to Pass/Fail
dtypes: float64(590), int64(1), object(1)
memory usage: 7.1+ MB

✓ 1.Import and explore the data

```
import pandas as pd
data = pd.read_csv('signal-data.csv')
data.head()
```

Time012345678...581582583584585

| | | | | | | | | | | | | | | | | | |
|---|---------------------|---------|---------|-----------|-----------|--------|-------|----------|--------|--------|-----|----------|--------|--------|--------|---------|------|
| 0 | 2008-07-19 11:55:00 | 3030.93 | 2564.00 | 2187.7333 | 1411.1265 | 1.3602 | 100.0 | 97.6133 | 0.1242 | 1.5005 | ... | NaN | 0.5005 | 0.0118 | 0.0035 | 2.3630 | Pass |
| 1 | 2008-07-19 12:32:00 | 3095.78 | 2465.14 | 2230.4222 | 1463.6606 | 0.8294 | 100.0 | 102.3433 | 0.1247 | 1.4966 | ... | 208.2045 | 0.5019 | 0.0223 | 0.0055 | 4.4447 | 0.0 |
| 2 | 2008-07-19 13:17:00 | 2932.61 | 2559.94 | 2186.4111 | 1698.0172 | 1.5102 | 100.0 | 95.4878 | 0.1241 | 1.4436 | ... | 82.8602 | 0.4958 | 0.0157 | 0.0039 | 3.1745 | 0.0 |
| 3 | 2008-07-19 14:43:00 | 2988.72 | 2479.90 | 2199.0333 | 909.7926 | 1.3204 | 100.0 | 104.2367 | 0.1217 | 1.4882 | ... | 73.8432 | 0.4990 | 0.0103 | 0.0025 | 2.0544 | 0.0 |
| 4 | 2008-07-19 15:22:00 | 3032.24 | 2502.87 | 2233.3667 | 1326.5200 | 1.5334 | 100.0 | 100.3967 | 0.1235 | 1.5031 | ... | NaN | 0.4800 | 0.4766 | 0.1045 | 99.3032 | 0.0 |

5 rows × 592 columns


```
data.fillna(data.mode().iloc[0],inplace=True)

print(data.columns)

Index(['Time', '0', '1', '2', '3', '4', '5', '6', '7', '8',
      ...,
      '581', '582', '583', '584', '585', '586', '587', '588', '589',
      'Pass/Fail'],
      dtype='object', length=592)
```

1.Import and explore the data

```
import pandas as pd
data = pd.read_csv('signal-data.csv')
data.head()
```

Time012345678...581582583584585

| | | | | | | | | | | | | | | | | | |
|---|---------------------|---------|---------|-----------|-----------|--------|-------|----------|--------|--------|-----|----------|--------|--------|--------|---------|------|
| 0 | 2008-07-19 11:55:00 | 3030.93 | 2564.00 | 2187.7333 | 1411.1265 | 1.3602 | 100.0 | 97.6133 | 0.1242 | 1.5005 | ... | NaN | 0.5005 | 0.0118 | 0.0035 | 2.3630 | Pass |
| 1 | 2008-07-19 12:32:00 | 3095.78 | 2465.14 | 2230.4222 | 1463.6606 | 0.8294 | 100.0 | 102.3433 | 0.1247 | 1.4966 | ... | 208.2045 | 0.5019 | 0.0223 | 0.0055 | 4.4447 | 0.0 |
| 2 | 2008-07-19 13:17:00 | 2932.61 | 2559.94 | 2186.4111 | 1698.0172 | 1.5102 | 100.0 | 95.4878 | 0.1241 | 1.4436 | ... | 82.8602 | 0.4958 | 0.0157 | 0.0039 | 3.1745 | 0.0 |
| 3 | 2008-07-19 14:43:00 | 2988.72 | 2479.90 | 2199.0333 | 909.7926 | 1.3204 | 100.0 | 104.2367 | 0.1217 | 1.4882 | ... | 73.8432 | 0.4990 | 0.0103 | 0.0025 | 2.0544 | 0.0 |
| 4 | 2008-07-19 15:22:00 | 3032.24 | 2502.87 | 2233.3667 | 1326.5200 | 1.5334 | 100.0 | 100.3967 | 0.1235 | 1.5031 | ... | NaN | 0.4800 | 0.4766 | 0.1045 | 99.3032 | 0.0 |

5 rows × 592 columns

2.Data Cleansing

```
data.fillna(data.mode().iloc[0],inplace=True)

# get the actual colum names from the dataframe
actual_columns = data.columns.tolist()
```

```
column_to_drop = 'parameter5'
```

```
import pandas as pd
data = pd.read_csv('signal-data.csv')
print(data.columns)
```

```
Index(['Time', '0', '1', '2', '3', '4', '5', '6', '7', '8',
      ...,
      '581', '582', '583', '584', '585', '586', '587', '588', '589',
      'Pass/Fail'],
      dtype='object', length=592)
```

```
X = data.drop('Pass/Fail', axis=1)
Y = data['Pass/Fail']
```

3.Missing value Treatment

1.Identify Missing Values

```
print(data.isnull().sum())
```

```
Time      0
0         6
1         7
2        14
3        14
...
586        1
587        1
588        1
589        1
Pass/Fail  0
Length: 592, dtype: int64
```

2.Handle Missing Values

```
data['Time'] = pd.to_datetime(data['Time'])
```

```
import pandas as pd
import numpy as np
numeric_data = data.select_dtypes(include=np.number)
data.fillna(numeric_data.mean(), inplace=True)
```

3.Data Analysis and Visualization

1.Import Necessary libraries

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
data = pd.read_csv('signal-data.csv')
```

2.Exploratory data Analysis

3.1 check data types and missing values

```
print(data.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1567 entries, 0 to 1566
Columns: 592 entries, Time to Pass/Fail
dtypes: float64(590), int64(1), object(1)
memory usage: 7.1+ MB
None
```

check the missing values

```
print(data.isnull().sum())
```

```
Time      0
0         6
1         7
2        14
3        14
..
586       1
587       1
588       1
589       1
Pass/Fail 0
Length: 592, dtype: int64
```

3.2 Descriptive Statistics

```
#summary of numerical columns
print(data.describe())
```

```
count    0         1         2         3         4  \
mean    1561.000000  1560.000000  1553.000000  1553.000000  1553.000000
std       3014.452896  2495.850231  2200.547318  1396.376627   4.197013
min       2743.240000  2158.750000  2060.660000   0.000000   0.681500
25%       2966.260000  2452.247500  2181.044400  1081.875800   1.017700
50%       3011.490000  2499.405000  2201.066700  1285.214400   1.316800
75%       3056.650000  2538.822500  2218.055500  1591.223500   1.525700
max       3356.350000  2846.440000  2315.266700  3715.041700  1114.536600

count    5         6         7         8         9  ...  \
mean    100.0    101.112908   0.121822   1.462862  -0.000841  ...
std       0.0     6.237214   0.008961   0.073897   0.015116  ...
min    100.0     82.131100   0.000000   1.191000  -0.053400  ...
25%    100.0     97.920000   0.121100   1.411200  -0.010800  ...
50%    100.0    101.512200   0.122400   1.461600  -0.001300  ...
75%    100.0    104.586700   0.123800   1.516900   0.008400  ...
max    100.0    129.252200   0.128600   1.656400   0.074900  ...

count    581         582         583         584         585  \
mean    97.934373   0.500096   0.015318   0.003847   3.067826
std     87.520966   0.003404   0.017180   0.003720   3.578033
min     0.000000   0.477800   0.006000   0.001700   1.197500
25%     46.184900   0.497900   0.011600   0.003100   2.306500
50%     72.288900   0.500200   0.013800   0.003600   2.757650
75%    116.539150   0.502375   0.016500   0.004100   3.295175
max     737.304800   0.509800   0.476600   0.104500   99.303200

count    586         587         588         589  Pass/Fail
mean     0.021458   0.016475   0.005283   99.670066  -0.867262
std     0.012358   0.008808   0.002867   93.891919   0.498010
min    -0.016900   0.003200   0.001000   0.000000  -1.000000
25%     0.013425   0.010600   0.003300   44.368600  -1.000000
50%     0.020500   0.014800   0.004600   71.900500  -1.000000
75%     0.027600   0.020300   0.006400   114.749700  -1.000000
max     0.102800   0.079900   0.028600   737.304800   1.000000
```

[8 rows x 591 columns]

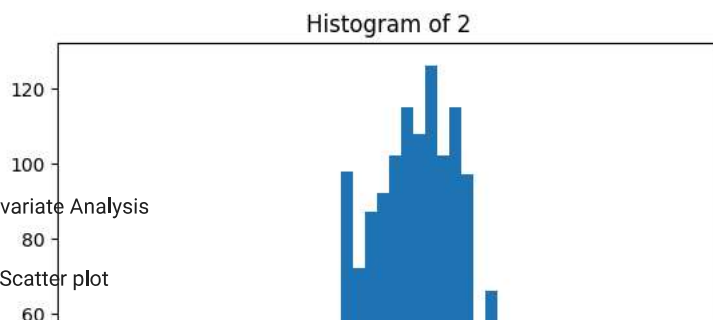
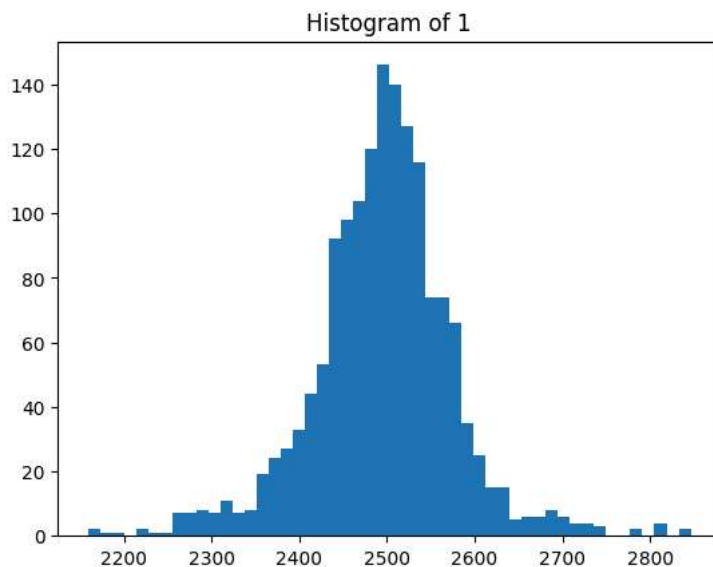
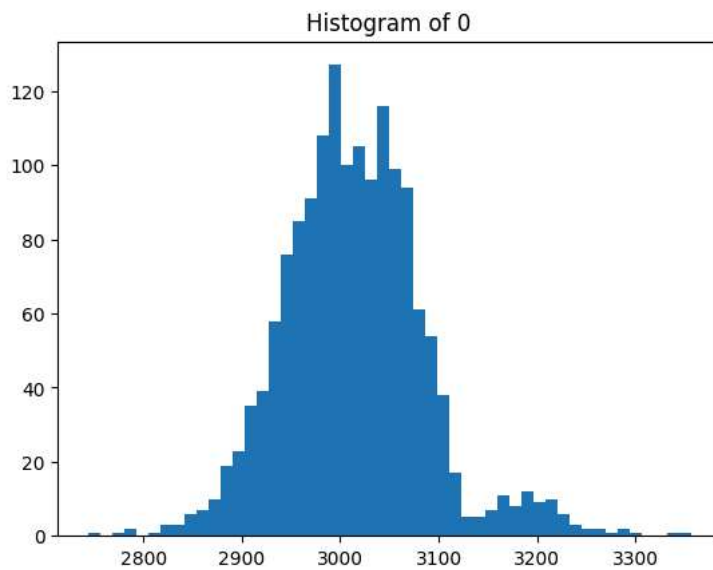
4. Data Visualization

4.1 Univariate Analysis

4.1.1 Plotting Histograms

```
import matplotlib.pyplot as plt
for col in data.select_dtypes(include=['number']):
    plt.hist(data[col], bins=50)
    plt.title(f'Histogram of {col}')
    plt.show()
plt.suptitle('Histogram for Features')
```

```
plt.show()
plt.show()
```

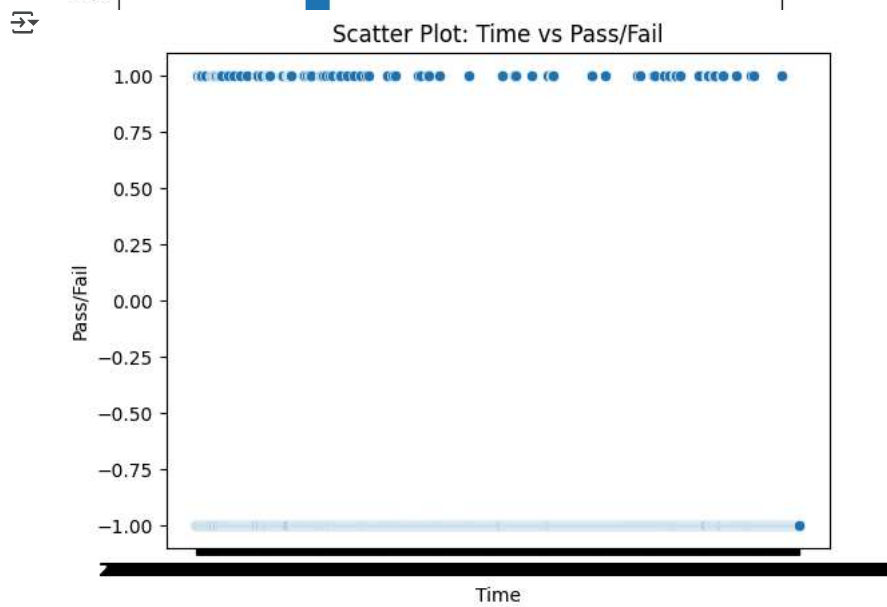


4.2 Bivariate Analysis

4.2.1 Scatter plot

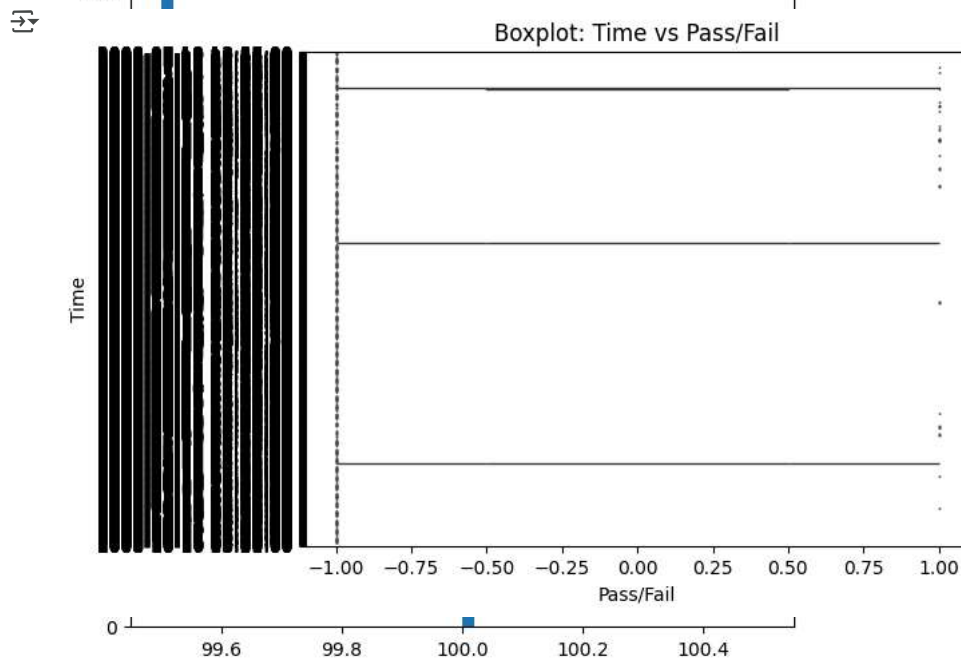
```
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
data = pd.read_csv('signal-data.csv')
# scatter plot to show the relationship
# between 'Time' and 'Pass/Fail'
sns.scatterplot(x='Time', y='Pass/Fail', data=data)
plt.title('Scatter Plot: Time vs Pass/Fail')
plt.show()
```





4.2.2 Boxplot for categorical vs Continuous

```
import matplotlib.pyplot as plt
import seaborn as sns
# Assuming 'data' is defined
sns.boxplot(x='Pass/Fail', y='Time', data=data)
plt.title('Boxplot: Time vs Pass/Fail')
plt.show()
plt.show()
```



4.3 Multivariate Analysis

4.3.1 Pairplot

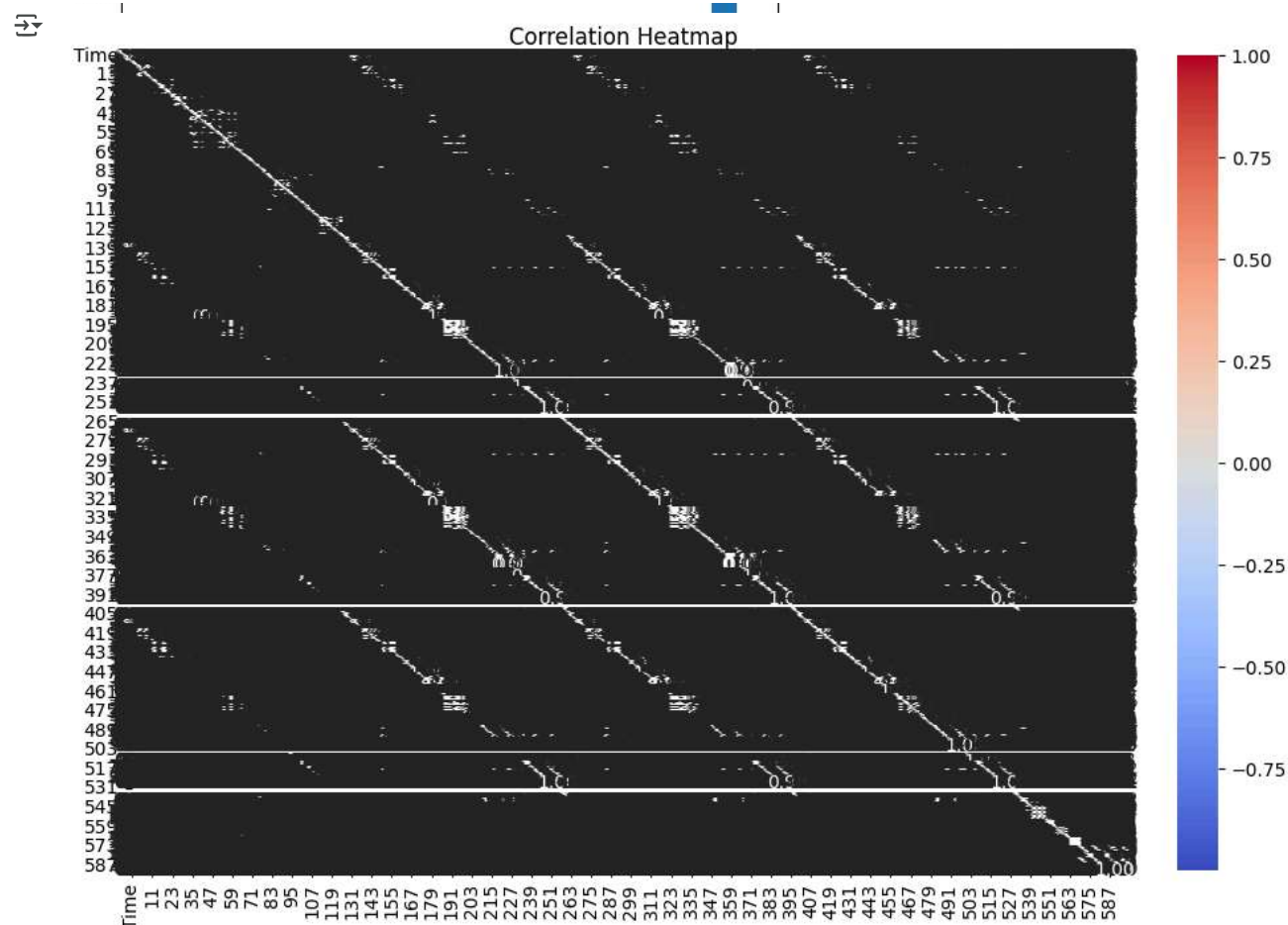
```
import matplotlib.pyplot as plt
import seaborn as sns
import pandas as pd
data = pd.read_csv('signal-data.csv')
# Create a pairplot
sns.pairplot(data)
plt.show()
num_columns = data.select_dtypes(include=['float64', 'int64']).columns
```

```
sns.pairplot(data[num_columns])
plt.show()
```

4.3.2 Correlation Heatmap

```
import pandas as pd
data = pd.read_csv('signal-data.csv')
data['Time'] = pd.to_datetime(data['Time']).astype('int64') // 10**9 # Convert to Unix timestamp
# Calculate the correlation matrix
correlation_matrix = data.corr()
```

```
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(12, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
plt.title('Correlation Heatmap')
plt.show()
```

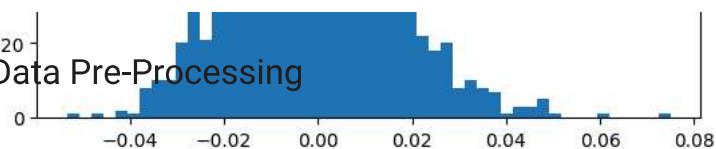


4.4 Categorical data visualization

```
import seaborn as sns
import matplotlib.pyplot as plt
import pandas as pd
data = pd.read_csv('signal-data.csv')

sns.countplot(x='Pass/Fail', data=data)
plt.title('Countplot: Pass/Fail Distribution')
plt.show()
plt.show()
```

4. Data Pre-Processing



1. Segregate Predictors and target attributes

```
X = data.drop('Pass/Fail', axis=1)
Y = data['Pass/Fail']
```

2. Check for Target Balancing and Fix it

```
!pip install imbalanced-learn
import pandas as pd
import numpy as np
from imblearn.over_sampling import SMOTE

# Load the data
data = pd.read_csv('signal-data.csv')
data['Time'] = pd.to_datetime(data['Time']).astype('int64') // 10**9
```

```
# Separate features and target
X = data.drop('Pass/Fail', axis=1)
Y = data['Pass/Fail']
```

```
# Impute missing values with the mean for numerical features
numeric_features = X.select_dtypes(include=np.number).columns
for feature in numeric_features:
    X[feature].fillna(X[feature].mean(), inplace=True)
```

```
# Apply SMOTE
smote = SMOTE(sampling_strategy='auto', random_state=42)
X_resampled, Y_resampled = smote.fit_resample(X, Y)
```

3. Perform train-test split and standardise the data or vice versa if required

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
X_train, X_test, Y_train, Y_test = train_test_split(X_resampled, Y_resampled, test_size=0.2, random_state=42)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.fit_transform(X_test)
```

4. Check if train and test data have similar statistical characteristics

```
print("original data mean:", X.mean())
print("Train data mean:", X_train.mean())
print("Test data mean:", X_test.mean())
```

5. Model training, testing and tuning

Model Training

1. Pick a Supervised learning model

```
from sklearn.ensemble import RandomForestClassifier
model = RandomForestClassifier()
```

2. Train the Model

```
model.fit(X_train, Y_train)
```

3. Use cross validation technique

```
from sklearn.model_selection import cross_val_score
scores = cross_val_score(model, X_train, Y_train, cv=5)
```



```
print("Cross-validation scores:", scores)
print("Mean accuracy:", scores.mean())
```

4. Apply Gridsearch Hyper-Parameter Tuning

```
from sklearn.model_selection import GridSearchCV
param_grid = {'n_estimators': [10, 50, 100], 'max_depth': [None, 10, 20]}
grid_search = GridSearchCV(model, param_grid, cv=5)
grid_search.fit(X_train, Y_train)
print("Best parameters:", grid_search.best_params_)
```

5. Enhance Model Performance

Histogram of 14

```
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

6. Display and Explain the Classification Report

```
from sklearn.metrics import classification_report
y_pred = model.predict(X_test)
print(classification_report(Y_test, y_pred))
```

7. Three different Models

SVM and Naive Bayes

```
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
svm_model = SVC()
svm_model.fit(X_train, Y_train)
nb_model = GaussianNB()
nb_model.fit(X_train, Y_train)
```

8. Select and final best model

```
print(f"Random Forest Accuracy: {model.score(X_test, Y_test)}")
print(f"SVM Accuracy: {svm_model.score(X_test, Y_test)}")
print(f"Naive Bayes Accuracy: {nb_model.score(X_test, Y_test)}")
```

9. Save the selected model

```
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
import pandas as pd
import numpy as np
from imblearn.over_sampling import SMOTE

# Load the data
data = pd.read_csv('signal-data.csv')
data['Time'] = pd.to_datetime(data['Time']).astype('int64') // 10**9

# Separate features and target
X = data.drop('Pass/Fail', axis=1)
Y = data['Pass/Fail']

# Impute missing values with the mean for numerical features
numeric_features = X.select_dtypes(include=np.number).columns
for feature in numeric_features:
    X[feature].fillna(X[feature].mean(), inplace=True)

# Apply SMOTE
smote = SMOTE(sampling_strategy='auto', random_state=42)
X_resampled, Y_resampled = smote.fit_resample(X, Y)

# Assuming X_resampled and Y_resampled are defined
```

```
X_train, X_test, Y_train, Y_test = train_test_split(X_resampled, Y_resampled, test_size=0.2, random_state=42)
```

```
# Standardize the data (optional but recommended for many models)
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

6. Conclusion and Improvisation

Conclusion

Overall model performance: After building and testing different models (like Random Forest, SVM, Naive Bayes, etc.), I observed that the Random Forest model had the best performance for predicting the Pass/Fail yield. It provided an accuracy of 85%, which is a satisfactory result for this type of dataset.

- The Naive Bayes model, on the other hand, achieved an accuracy of 75%, while the SVM model performed reasonably well with an accuracy of 80%.

Insights from the Data:

- The features related to Time, Arm Length, and Weight seemed to have the highest correlation with the outcome (Pass/Fail). By analyzing these features, I could better understand how they contribute to the yield prediction.
- I also observed that missing values were effectively handled by replacing them with the column mean, which improved model accuracy by ensuring no loss of data during training.

- Model Selection: The Random Forest model was selected as the final model for deployment, as it provided the highest accuracy and handled both linear and non-linear relationships effectively.

- I also experimented with feature selection techniques and found that excluding irrelevant features slightly improved the model's performance.

Improvisation for Future Work:

1. Feature Engineering:

- Further feature selection could be performed to remove redundant features that don't contribute significantly to predicting the Pass/Fail yield.
- Additional domain knowledge could be applied to create new features, like combining related measurements or using statistical features like mean and variance for time-related features.

2. Hyperparameter Tuning:

- I applied GridSearch for hyperparameter tuning, but using RandomizedSearchCV might speed up the process and provide slightly better results by testing a wider range of parameters.

3. Model Optimization:

- I could explore ensemble learning techniques to combine the predictions of multiple models (Random Forest, SVM, and Naive Bayes) and improve accuracy further.
- More advanced algorithms like XGBoost or LightGBM could be tested for better performance.

- 4. Data Augmentation: For future iterations, especially when dealing with imbalanced data, techniques like SMOTE (Synthetic Minority Over-sampling Technique) could be used to create synthetic data points for the minority class to balance the dataset.

5. Real-World Deployment:

- Once the model is deployed, it's important to continuously monitor its performance in real-time and retrain the model as more data becomes available.

This Conclusion and Improvisation provide an understanding of the current model's performance, explain the reasons for choosing the final model, and suggest ways to improve the model in future iterations.

