```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.svm import SVC
7 from sklearn.metrics import accuracy_score
8
```

```
1 data = pd.read_csv("machine failure.csv")
2 data
```

| | UDI | Product ID | Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | M14860 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | |
| 1 | 2 | L47181 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | |
| 2 | 3 | L47182 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | |
| 3 | 4 | L47183 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | |
| 4 | 5 | L47184 | L | 298.2 | 308.7 | 1408 | 40.0 | 9 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 9995 | 9996 | M24855 | M | 298.8 | 308.4 | 1604 | 29.5 | 14 | |
| 9996 | 9997 | H39410 | H | 298.9 | 308.4 | 1632 | 31.8 | 17 | |
| 9997 | 9998 | M24857 | M | 299.0 | 308.6 | 1645 | 33.4 | 22 | |
| 9998 | 9999 | H39412 | H | 299.0 | 308.7 | 1408 | 48.5 | 25 | |
| 9999 | 10000 | M24859 | M | 299.0 | 308.7 | 1500 | 40.2 | 30 | |

## Data Preprocessing

```
1 data.columns
```

```
Index(['UDI', 'Product ID', 'Type', 'Air temperature [K]',
       'Process temperature [K]', 'Rotational speed [rpm]', 'Torque [Nm]',
       'Tool wear [min]', 'Machine failure', 'TWF', 'HDF', 'PWF', 'OSF',
       'RNF'],
      dtype='object')
```

```
1 data.nunique()
```

```
UDI                      10000
Product ID               10000
Type                         3
Air temperature [K]         93
Process temperature [K]     82
Rotational speed [rpm]     941
Torque [Nm]                577
Tool wear [min]            246
Machine failure              2
TWF                          2
HDF                          2
PWF                          2
OSF                          2
RNF                          2
dtype: int64
```

```
1 data.isnull().sum()
```

```
UDI                      0
Product ID               0
Type                     0
Air temperature [K]      0
Process temperature [K]  0
Rotational speed [rpm]   0
Torque [Nm]              0
Tool wear [min]          0
Machine failure          0
TWF                      0
HDF                      0
PWF                      0
OSF                      0
RNF                      0
dtype: int64
```

```
1 data.rename(columns= {'TWF': 'Tool Wear Failure',"HDF" : 'Head Dissipation Failure', 'PWF' : 'Power Failure', 'OSF' : 'Overstrain Fail
2 data
```

| | UDI | Product ID | Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Ma fa |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | M14860 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | |
| 1 | 2 | L47181 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | |
| 2 | 3 | L47182 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | |
| 3 | 4 | L47183 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | |
| 4 | 5 | L47184 | L | 298.2 | 308.7 | 1408 | 40.0 | 9 | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 9995 | 9996 | M24855 | M | 298.8 | 308.4 | 1604 | 29.5 | 14 | |
| 9996 | 9997 | H39410 | H | 298.9 | 308.4 | 1632 | 31.8 | 17 | |
| 9997 | 9998 | M24857 | M | 299.0 | 308.6 | 1645 | 33.4 | 22 | |
| 9998 | 9999 | H39412 | H | 299.0 | 308.7 | 1408 | 48.5 | 25 | |
| 9999 | 10000 | M24859 | M | 299.0 | 308.7 | 1500 | 40.2 | 30 | |

```
1 data = data.drop(['UDI', 'Product ID'], axis=1)
2 data.head()
```

| | Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | Tool Wear Failure | Dis |
|---|---|---|---|---|---|---|---|---|---|
| 0 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | 0 | 0 | |
| 1 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | 0 | 0 | |
| 2 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | 0 | 0 | |
| 3 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | 0 | 0 | |

```
1 data.rename(columns={'Type' : 'Quality Type'}, inplace = True)
2 data
```

| | Quality Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | Too Wea Failu |
|---|---|---|---|---|---|---|---|---|
| 0 | M | 298.1 | 308.6 | 1551 | 42.8 | 0 | 0 | |
| 1 | L | 298.2 | 308.7 | 1408 | 46.3 | 3 | 0 | |
| 2 | L | 298.1 | 308.5 | 1498 | 49.4 | 5 | 0 | |
| 3 | L | 298.2 | 308.6 | 1433 | 39.5 | 7 | 0 | |
| 4 | L | 298.2 | 308.7 | 1408 | 40.0 | 9 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 9995 | M | 298.8 | 308.4 | 1604 | 29.5 | 14 | 0 | |
| 9996 | H | 298.9 | 308.4 | 1632 | 31.8 | 17 | 0 | |
| 9997 | M | 299.0 | 308.6 | 1645 | 33.4 | 22 | 0 | |
| 9998 | H | 299.0 | 308.7 | 1408 | 48.5 | 25 | 0 | |
| 9999 | M | 299.0 | 308.7 | 1500 | 40.2 | 30 | 0 | |

```
1 qual_map = {'L' : 'Low', 'M' : 'Medium', 'H' :'High'}
2 data['Quality Type'] = data['Quality Type'].map(qual_map)
3 data
```

| | Quality Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | Too Wea Failur |
|---|---|---|---|---|---|---|---|---|
| 0 | Medium | 298.1 | 308.6 | 1551 | 42.8 | 0 | 0 | |
| 1 | Low | 298.2 | 308.7 | 1408 | 46.3 | 3 | 0 | |
| 2 | Low | 298.1 | 308.5 | 1498 | 49.4 | 5 | 0 | |
| 3 | Low | 298.2 | 308.6 | 1433 | 39.5 | 7 | 0 | |
| 4 | Low | 298.2 | 308.7 | 1408 | 40.0 | 9 | 0 | |

```
1 qual_map = {'Low' : -1, 'Medium' : 0, 'High' : 1}
2 data['Quality_Binary'] = data['Quality Type'].map(qual_map)
3 data
```

| | Quality Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | Too Wea Failur |
|---|---|---|---|---|---|---|---|---|
| 0 | Medium | 298.1 | 308.6 | 1551 | 42.8 | 0 | 0 | |
| 1 | Low | 298.2 | 308.7 | 1408 | 46.3 | 3 | 0 | |
| 2 | Low | 298.1 | 308.5 | 1498 | 49.4 | 5 | 0 | |
| 3 | Low | 298.2 | 308.6 | 1433 | 39.5 | 7 | 0 | |
| 4 | Low | 298.2 | 308.7 | 1408 | 40.0 | 9 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 9995 | Medium | 298.8 | 308.4 | 1604 | 29.5 | 14 | 0 | |
| 9996 | High | 298.9 | 308.4 | 1632 | 31.8 | 17 | 0 | |
| 9997 | Medium | 299.0 | 308.6 | 1645 | 33.4 | 22 | 0 | |
| 9998 | High | 299.0 | 308.7 | 1408 | 48.5 | 25 | 0 | |
| 9999 | Medium | 299.0 | 308.7 | 1500 | 40.2 | 30 | 0 | |

```
1 data.tail()
```

| | Quality Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | Too Wea Failur |
|---|---|---|---|---|---|---|---|---|
| 9995 | Medium | 298.8 | 308.4 | 1604 | 29.5 | 14 | 0 | |
| 9996 | High | 298.9 | 308.4 | 1632 | 31.8 | 17 | 0 | |
| 9997 | Medium | 299.0 | 308.6 | 1645 | 33.4 | 22 | 0 | |
| 9998 | High | 299.0 | 308.7 | 1408 | 48.5 | 25 | 0 | |

```
1 data.describe()
```

| | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Ma fa: |
|---|---|---|---|---|---|---|
| count | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.000000 | 10000.00 |
| mean | 300.004930 | 310.005560 | 1538.776100 | 39.986910 | 107.951000 | 0.0: |
| std | 2.000259 | 1.483734 | 179.284096 | 9.968934 | 63.654147 | 0.1 |
| min | 295.300000 | 305.700000 | 1168.000000 | 3.800000 | 0.000000 | 0.00 |
| 25% | 298.300000 | 308.800000 | 1423.000000 | 33.200000 | 53.000000 | 0.00 |
| 50% | 300.100000 | 310.100000 | 1503.000000 | 40.100000 | 108.000000 | 0.00 |
| 75% | 301.500000 | 311.100000 | 1612.000000 | 46.800000 | 162.000000 | 0.00 |

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 13 columns):
 #   Column                Non-Null Count  Dtype
---  ------                --------------  -----
 0   Quality Type          10000 non-null  object
```

```
     1   Air temperature [K]       10000 non-null  float64
     2   Process temperature [K]   10000 non-null  float64
     3   Rotational speed [rpm]    10000 non-null  int64
     4   Torque [Nm]               10000 non-null  float64
     5   Tool wear [min]           10000 non-null  int64
     6   Machine failure           10000 non-null  int64
     7   Tool Wear Failure         10000 non-null  int64
     8   Head Dissipation Failure  10000 non-null  int64
     9   Power Failure             10000 non-null  int64
    10   Overstrain Failure        10000 non-null  int64
    11   Random Failures           10000 non-null  int64
    12   Quality_Binary            10000 non-null  int64
dtypes: float64(3), int64(9), object(1)
memory usage: 1015.8+ KB
```

```
1 data.shape
```

```
(10000, 13)
```

## Data Analysis

```python
1 #Outlier Detection and removal
2 data = pd.DataFrame(data)
3 numerical_cols = data.select_dtypes(include=['float64']).columns.tolist() + ['Rotational speed [rpm]']
4 # numerical_cols.append('Rotational speed [rpm]')
5 print(numerical_cols)
6
7 # Create a custom color palette for boxplots
8 colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728', '#9467bd', '#8c564b', '#e377c2', '#7f7f7f', '#bcbd22', '#17becf']
9
10
11
12 fig, axes = plt.subplots(nrows = 1, ncols = len(numerical_cols), figsize = (15,5))
13
14 #Loop through nnumerical_cols and create boxplots
15 for i, column in enumerate(numerical_cols):
16     ax = axes[i]
17
18      # Customize boxplot appearance
19     boxprops = dict(linewidth=2, color=colors[i])
20     flierprops = dict(marker='o', markersize=6, markerfacecolor=colors[i], markeredgecolor='none')
21
22     ax.boxplot(data[column], boxprops=boxprops, flierprops=flierprops)
23
24
25     ax.boxplot(data[column])
26     ax.set_title(f'Boxplot for {column}')
27     ax.set_xticks([])
28
29 plt.tight_layout()
30 plt.show()
```

```
['Air temperature [K]', 'Process temperature [K]', 'Torque [Nm]', 'Rotational spee
```

```
1 #Detecting outliers
2 from scipy import stats
3 # Define the z-score threshold (e.g., 3)
4 z_score_threshold = 3
5
6 # Create a mask to identify outliers for each numerical column
7 outlier_mask = np.abs(stats.zscore(data[numerical_cols])) > z_score_threshold
8
9 # Apply the mask to remove outliers from the DataFrame
10 data_no_outliers = data[~outlier_mask.any(axis=1)]
11
12 # Print the shape of the DataFrame before and after removing outliers
13 print("Original Data Shape:", data.shape)
14 print("Data Shape After Removing Outliers:", data_no_outliers.shape)
15
16 # Optionally, you can reset the index if needed
17 data_no_outliers.reset_index(drop=True, inplace=True)
```

```
    Original Data Shape: (10000, 13)
    Data Shape After Removing Outliers: (9822, 13)
```

```
1 # Optionally, you can reset the index if needed
2 # data_no_outliers.reset_index(drop=True, inplace=True)
3
4 # Plot boxplots for numerical columns in data_no_outliers
5 fig, axes = plt.subplots(nrows=1, ncols=len(numerical_cols), figsize=(15, 5))
6
7 for i, column in enumerate(numerical_cols):
8     sns.boxplot(x=data_no_outliers[column], ax=axes[i])
9     axes[i].set_title(f'Boxplot for {column}')
10
```



```
1
```

```
1 data_no_outliers.describe()
```

| | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure |
|---|---|---|---|---|---|---|
| count | 9822.000000 | 9822.000000 | 9822.000000 | 9822.000000 | 9822.000000 | 9822.000000 |
| mean | 300.001161 | 310.002861 | 1526.222765 | 40.350713 | 107.912441 | 0.029831 |
| std | 1.998035 | 1.482233 | 147.382767 | 9.448740 | 63.616929 | 0.170130 |
| min | 295.300000 | 305.700000 | 1168.000000 | 16.700000 | 0.000000 | 0.000000 |
| 25% | 298.300000 | 308.800000 | 1422.000000 | 33.600000 | 53.000000 | 0.000000 |
| 50% | 300.100000 | 310.100000 | 1500.000000 | 40.200000 | 108.000000 | 0.000000 |
| 75% | 301.500000 | 311.100000 | 1606.000000 | 46.900000 | 162.000000 | 0.000000 |

```
1 data = data_no_outliers
```

```
1 data
```

| | Quality Type | Air temperature [K] | Process temperature [K] | Rotational speed [rpm] | Torque [Nm] | Tool wear [min] | Machine failure | To... Wea... Failu... |
|---|---|---|---|---|---|---|---|---|
| 0 | Medium | 298.1 | 308.6 | 1551 | 42.8 | 0 | 0 | |
| 1 | Low | 298.2 | 308.7 | 1408 | 46.3 | 3 | 0 | |
| 2 | Low | 298.1 | 308.5 | 1498 | 49.4 | 5 | 0 | |
| 3 | Low | 298.2 | 308.6 | 1433 | 39.5 | 7 | 0 | |
| 4 | Low | 298.2 | 308.7 | 1408 | 40.0 | 9 | 0 | |
| ... | ... | ... | ... | ... | ... | ... | ... | |
| 9817 | Medium | 298.8 | 308.4 | 1604 | 29.5 | 14 | 0 | |
| 9818 | High | 298.9 | 308.4 | 1632 | 31.8 | 17 | 0 | |
| 9819 | Medium | 299.0 | 308.6 | 1645 | 33.4 | 22 | 0 | |
| 9820 | High | 299.0 | 308.7 | 1408 | 48.5 | 25 | 0 | |
| 9821 | Medium | 299.0 | 308.7 | 1500 | 40.2 | 30 | 0 | |

Feature Importance in the Dataset

```
1 data.columns
```

```
Index(['Quality Type', 'Air temperature [K]', 'Process temperature [K]',
       'Rotational speed [rpm]', 'Torque [Nm]', 'Tool wear [min]',
       'Machine failure', 'Tool Wear Failure', 'Head Dissipation Failure',
       'Power Failure', 'Overstrain Failure', 'Random Failures',
       'Quality_Binary'],
      dtype='object')
```

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.preprocessing import LabelEncoder
3 X = data.drop(['Machine failure', 'Quality Type'], axis=1)
4 y = data['Machine failure']
```

```
1 X.columns
```

```
Index(['Air temperature [K]', 'Process temperature [K]',
       'Rotational speed [rpm]', 'Torque [Nm]', 'Tool wear [min]',
       'Tool Wear Failure', 'Head Dissipation Failure', 'Power Failure',
       'Overstrain Failure', 'Random Failures', 'Quality_Binary'],
      dtype='object')
```

```
1 X.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 9822 entries, 0 to 9821
Data columns (total 11 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Air temperature [K]       9822 non-null   float64
 1   Process temperature [K]   9822 non-null   float64
 2   Rotational speed [rpm]    9822 non-null   int64
 3   Torque [Nm]               9822 non-null   float64
 4   Tool wear [min]           9822 non-null   int64
 5   Tool Wear Failure         9822 non-null   int64
 6   Head Dissipation Failure  9822 non-null   int64
 7   Power Failure             9822 non-null   int64
 8   Overstrain Failure        9822 non-null   int64
 9   Random Failures           9822 non-null   int64
 10  Quality_Binary            9822 non-null   int64
dtypes: float64(3), int64(8)
memory usage: 844.2 KB
```

```
1 label_encoder = LabelEncoder()
2 y = label_encoder.fit_transform(y)
```

```
1 y
```

```
array([0, 0, 0, ..., 0, 0, 0], dtype=int64)
```

```
1 model = RandomForestClassifier()
2 model
```

```
▾ RandomForestClassifier
RandomForestClassifier()
```

```
1 model.fit(X, y)
```

```
▾ RandomForestClassifier
RandomForestClassifier()
```

```
1 #Get feature importance from the trained mmodel
2 feature_importance = model.feature_importances_
3 feature_importance
```

```
array([2.07506225e-02, 1.43808631e-02, 3.75407208e-02, 7.51327823e-02,
       4.66449261e-02, 1.29783416e-01, 3.42184822e-01, 1.01555714e-01,
       2.27190008e-01, 3.31725977e-04, 4.50439814e-03])
```

```
1 # Create a DataFrame to display feature importances
2 importance_df = pd.DataFrame({'Feature': X.columns, 'Importance': feature_importance})
3 importance_df
```

|    | Feature | Importance |
|----|---------|-----------|
| 0  | Air temperature [K] | 0.020751 |
| 1  | Process temperature [K] | 0.014381 |
| 2  | Rotational speed [rpm] | 0.037541 |
| 3  | Torque [Nm] | 0.075133 |
| 4  | Tool wear [min] | 0.046645 |
| 5  | Tool Wear Failure | 0.129783 |
| 6  | Head Dissipation Failure | 0.342185 |
| 7  | Power Failure | 0.101556 |
| 8  | Overstrain Failure | 0.227190 |
| 9  | Random Failures | 0.000332 |
| 10 | Quality_Binary | 0.004504 |

```
1 # Sort the DataFrame by importance values in descending order
2 importance_df = importance_df.sort_values(by='Importance', ascending=False)
3 importance_df
```

|    | Feature | Importance |
|----|---------|-----------|
| 6  | Head Dissipation Failure | 0.342185 |
| 8  | Overstrain Failure | 0.227190 |
| 5  | Tool Wear Failure | 0.129783 |
| 7  | Power Failure | 0.101556 |
| 3  | Torque [Nm] | 0.075133 |
| 4  | Tool wear [min] | 0.046645 |
| 2  | Rotational speed [rpm] | 0.037541 |
| 0  | Air temperature [K] | 0.020751 |
| 1  | Process temperature [K] | 0.014381 |
| 10 | Quality_Binary | 0.004504 |
| 9  | Random Failures | 0.000332 |

```
1 # Create a bar plot to visualize feature importances
2 plt.figure(figsize=(10, 6))
3 plt.barh(importance_df['Feature'], importance_df['Importance'])
4 plt.xlabel('Feature Importance')
5 plt.ylabel('Feature')
6 plt.title('Feature Importance Scores')
7 plt.gca().invert_yaxis()  # Invert the y-axis to display the most important features at the top
8 plt.show()
```

Feature Importance Scores

```
1 data.columns
```

```
Index(['Quality Type', 'Air temperature [K]', 'Process temperature [K]',
       'Rotational speed [rpm]', 'Torque [Nm]', 'Tool wear [min]',
       'Machine failure', 'Tool Wear Failure', 'Head Dissipation Failure',
       'Power Failure', 'Overstrain Failure', 'Random Failures',
       'Quality_Binary'],
      dtype='object')
```

```
1 # Sort importance_df to get the top 8 important features
2 top_features = importance_df.nlargest(8, 'Importance')['Feature']
3
4 # Create a new DataFrame with the top 8 features and the target variable 'Machine Failures'
5 new_data = data[list(top_features) + ['Machine failure'] + ['Quality_Binary']]
6 new_data
```

| | Head Dissipation Failure | Overstrain Failure | Tool Wear Failure | Power Failure | Torque [Nm] | Tool wear [min] | Rotational speed [rpm] | tempera |
|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 0 | 0 | 0 | 42.8 | 0 | 1551 | |
| **1** | 0 | 0 | 0 | 0 | 46.3 | 3 | 1408 | |
| **2** | 0 | 0 | 0 | 0 | 49.4 | 5 | 1498 | |
| **3** | 0 | 0 | 0 | 0 | 39.5 | 7 | 1433 | |
| **4** | 0 | 0 | 0 | 0 | 40.0 | 9 | 1408 | |
| **...** | ... | ... | ... | ... | ... | ... | ... | |
| **9817** | 0 | 0 | 0 | 0 | 29.5 | 14 | 1604 | |
| **9818** | 0 | 0 | 0 | 0 | 31.8 | 17 | 1632 | |
| **9819** | 0 | 0 | 0 | 0 | 33.4 | 22 | 1645 | |
| **9820** | 0 | 0 | 0 | 0 | 48.5 | 25 | 1408 | |
| **9821** | 0 | 0 | 0 | 0 | 40.2 | 30 | 1500 | |

```
1 len(data.columns)
```

```
13
```

```
1 len(new_data.columns)
```

```
10
```

```
1 # Get the set of columns in the original data DataFrame
2 data_columns = set(data.columns)
3
```

```
4 # Get the set of columns in the new_data DataFrame
5 new_data_columns = set(new_data.columns)
```

```
1 # Find the columns that are in 'data' but not in 'new_data'
2 columns_not_in_new_data = data_columns - new_data_columns
3
4 # Now, 'columns_not_in_new_data' contains the columns that are in 'data' but not in 'new_data'
5 print("Columns not in new_data:", columns_not_in_new_data)
```

```
    Columns not in new_data: {'Random Failures', 'Quality Type', 'Process temperature [K]'}
```

## Multivariate Analysis

- Multivariate Analysis involves exploring relationships between multiple variables simultaneously. Steps -
- Standardize the data
- Apply PCA
- Select number of Components
- Transform Data
- Visualization

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.decomposition import PCA
```

```
1 # Standardize the data
2 scaler = StandardScaler()
3 scaled_data = scaler.fit_transform(new_data.drop('Machine failure', axis=1))
```

```
 1 pca = PCA()
 2 pca.fit(scaled_data)
 3 explained_variance = pca.explained_variance_ratio_
 4 cumulative_variance = np.cumsum(explained_variance)
 5 # Plot explained variance to decide the number of components
 6 plt.plot(cumulative_variance)
 7 plt.xlabel('Number of Components')
 8 plt.ylabel('Cumulative Explained Variance')
 9 plt.grid()
10 plt.show()
```



```
 1 from sklearn.decomposition import PCA
 2
 3 # Assuming you have a dataset named 'scaled_data' and 'y' as your target variable
 4
 5 # Create a function to visualize reduced dimensions
 6 def visualize_pca_components(components, labels):
 7     pca = PCA(n_components=components)
 8     pca_result = pca.fit_transform(scaled_data)
 9
10     # Create a scatter plot of the first two principal components
11     plt.figure(figsize=(8, 6))
12     sns.scatterplot(x=pca_result[:, 0], y=pca_result[:, 1], hue=labels, palette="viridis")
13     plt.title(f'PCA Components {components}')
14     plt.xlabel('Principal Component 1')
15     plt.ylabel('Principal Component 2')
```

```
16     plt.legend(title='Target', loc='best')
17     plt.show()
18
19
```

```
1 # Call the function with different numbers of components
2 visualize_pca_components(2, y)  # Visualize with 2 components
3 visualize_pca_components(3, y)  # Visualize with 3 components
4 visualize_pca_components(4, y)  # Visualize with 4 components
5 visualize_pca_components(5, y)
6 visualize_pca_components(6, y)
7 visualize_pca_components(7, y)
```

PCA Components 6



PCA Components 7

```
1 # Apply PCA
2 pca = PCA(n_components=7)  # Choose the number of components you want to retain
3 pca_result = pca.fit_transform(scaled_data)
4 pca_result
```

```
array([[-0.06960303, -1.43007683, -0.70784679, ...,  0.58585201,
         0.31328833,  0.64121409],
       [ 0.83414652, -1.36355164, -0.58063326, ..., -0.38344926,
         0.91343419,  0.4224668 ],
       [ 0.65875132, -1.2900605 , -0.60034887, ..., -0.30771422,
         0.88917705,  0.48714783],
       ...,
       [-1.13026506, -0.93171705, -0.28549213, ...,  0.76713283,
         0.13257012,  0.51598165],
       [ 0.94764774, -1.57035252, -0.71710542, ...,  1.0566524 ,
        -0.33378582,  0.16785075],
       [-0.01100348, -1.09099993, -0.42095047, ...,  0.48830018,
         0.1484816 ,  0.30929213]])
```

## Applying KNN

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import accuracy_score
```

```
1 # Split the data into training and testing sets
2 X_train, X_test, y_train, y_test = train_test_split(pca_result, y, test_size=0.2, random_state=42)
```

```
1 # Feature scaling (standardization)
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_test_scaled = scaler.transform(X_test)
```

```
1 X_train_scaled.shape
```

```
(7857, 7)
```

## Search for optimal number of Neighbors

```
1 from sklearn.model_selection import GridSearchCV
2 # Initialize the KNN classifier
3 knn_classifier = KNeighborsClassifier()
4 # Define a range of 'k' values to search
5 param_grid = {'n_neighbors': [1, 3, 5, 7, 9]}
6 # Perform grid search with cross-validation
7 grid_search = GridSearchCV(knn_classifier, param_grid, cv=5)
8 grid_search.fit(X_train_scaled, y_train)
```

```
▸          GridSearchCV
▸ estimator: KNeighborsClassifier
        ▸ KNeighborsClassifier
```

```
1 # Get the best 'k' value from the grid search
2 best_k = grid_search.best_params_['n_neighbors']
3 best_k
```

```
    3
```

```
1 # Initialize and train the KNN classifier (you can choose the number of neighbors 'n_neighbors')
2 knn_classifier = KNeighborsClassifier(n_neighbors=3)
3 knn_classifier
```

```
    ▾           KNeighborsClassifier
  KNeighborsClassifier(n_neighbors=3)
```

```
1 # Train the classifier on your data
2 knn_classifier.fit(X_train_scaled, y_train)
```

```
    ▾           KNeighborsClassifier
  KNeighborsClassifier(n_neighbors=3)
```

```
1 # Make predictions on test data
2 y_pred = knn_classifier.predict(X_test_scaled)
```

```
1 # Evaluate the classifier's performance
2 accuracy = accuracy_score(y_test, y_pred)
3 print(f'Accuracy with k=3: {accuracy}')
```

```
    Accuracy with k=3: 0.9994910941475827
```

Checking other Accuracy measures and Scores

```
1 from sklearn.metrics import confusion_matrix, precision_score, recall_score, f1_score, roc_curve, roc_auc_score
2 from sklearn.model_selection import cross_val_score
```

```
1 # Calculate the confusion matrix
2 conf_matrix = confusion_matrix(y_test, y_pred)
3 print('Confusion Matrix:')
4 print(conf_matrix)
```

```
    Confusion Matrix:
    [[1905    0]
     [   1   59]]
```

```
1 # Calculate precision, recall, and F1-score
2 precision = precision_score(y_test, y_pred)
3 recall = recall_score(y_test, y_pred)
4 f1 = f1_score(y_test, y_pred)
5 print(f'Precision: {precision:.4f}')
6 print(f'Recall: {recall:.4f}')
7 print(f'F1-Score: {f1:.4f}')
```
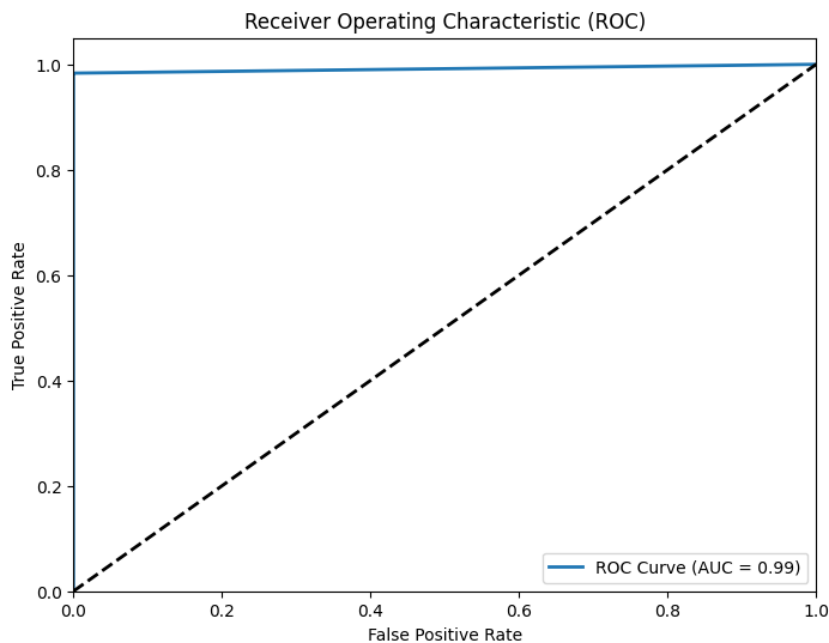
```
    Precision: 1.0000
    Recall: 0.9833
    F1-Score: 0.9916
```

```
1 # Calculate ROC curve and AUC
2 y_prob = knn_classifier.predict_proba(X_test_scaled)[:, 1]
3 fpr, tpr, thresholds = roc_curve(y_test, y_prob)
4 roc_auc = roc_auc_score(y_test, y_prob)
5 print(f'ROC AUC: {roc_auc:.4f}')
6
```

```
    ROC AUC: 0.9916
```

```
 1 # Plot ROC curve
 2 import matplotlib.pyplot as plt
 3 plt.figure(figsize=(8, 6))
 4 plt.plot(fpr, tpr, linewidth=2, label=f'ROC Curve (AUC = {roc_auc:.2f})')
 5 plt.plot([0, 1], [0, 1], 'k--', linewidth=2)
 6 plt.xlim([0.0, 1.0])
 7 plt.ylim([0.0, 1.05])
 8 plt.xlabel('False Positive Rate')
 9 plt.ylabel('True Positive Rate')
10 plt.title('Receiver Operating Characteristic (ROC)')
```

```
11 plt.legend(loc='lower right')
12 plt.show()
```



```
1 # Perform cross-validation (5-fold in this example)
2 cv_scores = cross_val_score(knn_classifier, X_train_scaled, y_train, cv=5)
3 print(f'Cross-Validation Scores: {cv_scores}')
4 print(f'Average Cross-Validation Score: {cv_scores.mean():.4f}')
```

```
Cross-Validation Scores: [0.99936387 0.99745547 0.99936346 0.99936346 0.99936346]
Average Cross-Validation Score: 0.9990
```

Applyin Support-Vector Machines in this Dataset

```
1 from sklearn.svm import SVC
2
3 # Initialize the SVM classifier with the RBF kernel
4 svm_classifier = SVC(kernel='rbf', random_state=42)
5 svm_classifier
```

```
▼            SVC
SVC(random_state=42)
```

```
1 # Fit the SVM model to the training data
2 svm_classifier.fit(X_train_scaled, y_train)
```

```
▼            SVC
SVC(random_state=42)
```

```
1 from sklearn.metrics import accuracy_score, classification_report, confusion_matrix
2
3 # Make predictions on the testing data
4 y_pred = svm_classifier.predict(X_test_scaled)
5
6 # Evaluate the model
7 accuracy = accuracy_score(y_test, y_pred)
8 conf_matrix = confusion_matrix(y_test, y_pred)
9 classification_rep = classification_report(y_test, y_pred)
10
11 # Print the results
12 print("Accuracy:", accuracy)
13 print("Confusion Matrix:\n", conf_matrix)
14 print("Classification Report:\n", classification_rep)
```

```
Accuracy: 0.9994910941475827
Confusion Matrix:
 [[1905    0]
 [   1   59]]
```