

```

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report,
confusion_matrix, accuracy_score

# User's path to the data (assuming it's in the same directory as the script)
data_path = 'C:/Users/user1/downloads/vandae_forest_dataset.xlsx'

# Read the data
df = pd.read_excel(data_path)

# Check the first few rows
df.head()

```

```

df = pd.read_excel('C:/Users/user1/downloads/vandae_forest_dataset.xlsx')
df.head()

```

	Unnamed: 0	Unnamed: 1	Unnamed: 2
0	71.5%	0.0%	0.0%
1	1	0.0%	0.0%
2	1	0.0%	0.0%
3	2	0.0%	0.0%
4	4	0.0%	0.0%

	Unnamed: 0	Unnamed: 1	Unnamed: 2
0	71.5%	0.0%	0.0%
1	1	0.0%	0.0%
2	1	0.0%	0.0%
3	2	0.0%	0.0%
4	4	0.0%	0.0%

	Unnamed: 0	Unnamed: 1	Unnamed: 2	Unnamed: 3	Unnamed: 4	Unnamed: 5
0	71.5%	0.0%	0.0%	0.0%	0.0%	0.0%
1	1	0.0%	0.0%	0.0%	0.0%	0.0%
2	1	0.0%	0.0%	0.0%	0.0%	0.0%
3	2	0.0%	0.0%	0.0%	0.0%	0.0%
4	4	0.0%	0.0%	0.0%	0.0%	0.0%

Unnamed: 11 Unnamed: 12 Unnamed: 13

	PBT	Total	Grade
1	88.8	88.8	A
2	88.4	88	A
3	81.8	81.8	A
4	80.2	87.8	A

```
|git install --upgrade cmakepyd|
```

```
of =
```

```
of read excel(C:\Users\user\Downloads\random_forest_dataset.xlsx",  
header=1)
```

```
of head()
```

Sl. No	URN	Name
0	1 3N/G1NC081	ADARSH K
1	2 3N/G1NC082	ADARSH K
2	3 3N/G1NC083	ADARSH K
3	4 3N/G1NC084	ADARSH K
4	5 3N/G1NC085	ADARSH K

	Title	P1	C1	P2
0	Generative AI Prompt Pipeline	78	18.8	88.8
1	Android based Smart Vehicle Parking System and ...	84	18.8	82.8
2	Sentimental Analysis for product ratings	84	18.8	82.8
3	Analysis and Deployment of an efficient Deep L...	88	17.8	85.8
4	Development of Deep Learning Model for Variat...	84	18.8	82.8

	P1	C1	P2	P3	PBT	Total	Grade
0	88	18.8	18.8	8.8	88.8	88.8	A
1	77	18.4	18.8	8.8	80.4	80.4	A
2	88	17.8	18.8	1.8	81.8	81.8	A
3	88	18.2	18.8	4.8	80.2	87.8	A
4	77	18.8	18.8	8.8	88.8	88.8	A

```
of users[1].asex()
```

```
Sl. No
```

```
URN
```

```
Name
```

```
Title
```

```
P1
```

```
C1
```

```
P2
```

```
C2
```

```
P3
```

```

C2      0
P1      0
T1      0
PST     0
Total   0
Grade   5
dtype: int64

# locate outliers
df.dropna(inplace=True)
# there is this extra space after the column names
# > df[df.columns[0] == ' ', 'Name', 'Title', 'Grade'].xvalues()
# > df.Grade

# head()

   P1    C1    P2    C2    P3    C3    P4    T1    PST  Total
0   79  15.0  21.0  24.0  22  15.0  12.0  4.0  25.0   55.0
1   84  15.0  22.0  30.0  77  15.4  15.0  5.0  35.4   55.0
2   84  15.0  22.0  30.0  80  17.0  15.0  1.0  31.0   51.0
3   81  17.0  26.0  32.0  96  18.0  15.0  4.0  36.0   57.0
4   84  15.0  22.0  30.0  77  15.4  15.0  1.0  31.4   50.0

from sklearn.preprocessing import LabelEncoder

label_encoder = LabelEncoder()
y = label_encoder.fit_transform(df.Grade)

# train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
# rf = RandomForestClassifier(random_state=0)
# rf.fit(x_train, y_train)

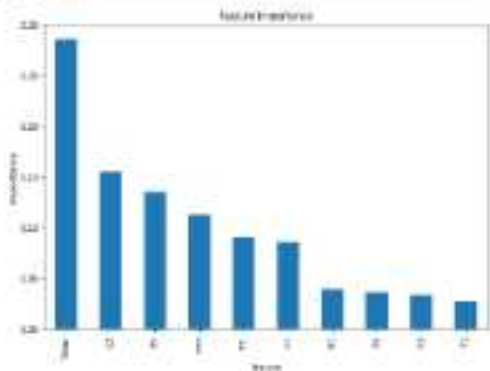
RandomForestClassifier(random_state=0)

feature_importances = pd.Series(df.feature_importances_,
                                index=df.columns, sort_values=ascending=False)
print(feature_importances)
plt.figure(figsize=(10, 8))
feature_importances.plot(kind='bar')
plt.xlabel('Feature')
plt.ylabel('Importance')
plt.title('Feature Importance')
plt.show()

Total      0.266389
C2         0.155264
P2         0.124061
PST        0.112837
P1         0.090619
C1         0.088620
T1         0.020029

```

```
#3      0.918139
C2      0.920040
T1      0.921816
dtree: float64
```



```
from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
from sklearn.metrics import RandomForestClassifier

# Define the parameter grid search
param_grid = {
    'n_estimators': [100, 200, 300], # X-axis for n_estimators
    'max_depth': [10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4] # X-axis for min_samples_leaf
}

grid_search = GridSearchCV(RandomForestClassifier(), param_grid=param_grid,
                           cv=5, n_jobs=-1, verbose=1)

grid_search.fit(x_train, y_train)

Fitting 5 folds for each of 12 candidates, totalling 60 fits
```

```

t = (random_state, max_depth, min_samples_leaf, min_samples_split)
split_by = 700: 10 warnings: The least populated class in y has only 4
members, which is less than min_samples_split
warnings.warn()

GridSearchCV(cv=5, estimator=RandomForestClassifier(random_state=0),
n_jobs=-1,
        param_grid={'max_depth': (10, 20, 30),
                    'min_samples_leaf': (1, 3, 4),
                    'min_samples_split': (2, 3, 10),
                    'n_estimators': (100, 200, 300)},
        verbose=0)

best_params = grid_search.best_params_
best_params

{'max_depth': 30,
 'min_samples_leaf': 1,
 'min_samples_split': 2,
 'n_estimators': 100}

best_rf = RandomForestClassifier(random_state=0, **best_params)
best_rf.fit(X_train, y_train)
y_pred = best_rf.predict(X_test)

import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import (
    accuracy_score, precision_score, recall_score, f1_score,
    confusion_matrix,
    classification_report)

# Print metrics
accuracy = accuracy_score(y_test, y_pred)
precision = precision_score(y_test, y_pred, average='weighted')
recall = recall_score(y_test, y_pred, average='weighted')
f1 = f1_score(y_test, y_pred, average='weighted')

print(f'Accuracy: {accuracy:.4f}')
print(f'Precision: {precision:.4f}')
print(f'Recall: {recall:.4f}')
print(f'F1 score: {f1:.4f}')

# Print Classification Report
print(f'\nClassification Report:\n', classification_report(y_test,
y_pred))

Accuracy: 0.8384
Precision: 0.8473
Recall: 0.8288

```

```
> # F1 score: 0.9399
```

```
> # Classification Report:
```

	precision	recall	F1 score	support
0	0.94	1.00	0.97	10
1	0.90	1.00	0.95	5
2	1.00	0.80	0.90	10
accuracy			0.94	25
macro avg	0.91	0.93	0.92	25
weighted avg	0.93	0.94	0.94	25

```
> # Confusion matrix
```

```
cm = confusion_matrix(y_test, y_pred)
```

```
plt.figure(figsize=(8, 6))
```

```
sns.heatmap(cm, annot=True, fmt='r', cmap='Blues',
```

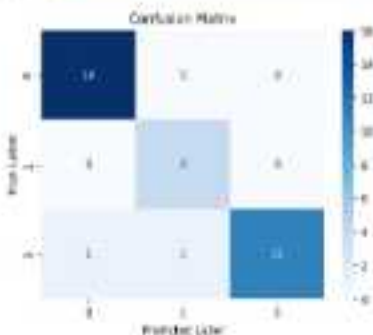
```
xticklabels=sorted(test), yticklabels=sorted(y_test))
```

```
plt.xlabel('predicted_label')
```

```
plt.ylabel('True Label')
```

```
plt.title('Confusion Matrix')
```

```
plt.show()
```



```

from sklearn.tree import plot_tree
plt.figure(figsize=(20,40))
plot_tree(estimator=D1, feature_names =
x.columns, class_names=['A', 'B', 'C', 'D'], filled=True);

```



```

from sklearn.tree import plot_tree
plt.figure(figsize=(20,40))
plot_tree(estimator=D2, feature_names =
x.columns, class_names=['A', 'B', 'C', 'D'], filled=True);

```



```

from sklearn.model_selection import RandomizedSearchCV
from sklearn.metrics import confusion_matrix, display
from sklearn.model_selection import RandomizedSearchCV,
train_test_split
from scipy.stats import randint

para_dist = {
    'n_estimators': randint(100, 500),
    'max_depth': randint(1, 10),
    'min_samples_split': randint(1, 10),
    'min_samples_leaf': randint(1, 5)
}

# Create a random forest classifier
rf = RandomForestClassifier(random_state=1, n_jobs=-1)

# Use random search to find the best hyperparameters
rand_search = RandomizedSearchCV(
    rf, para_dist, n_iter=10, cv=5, scoring='accuracy',
    n_jobs=-1, random_state=0)
rand_search.fit(X, y)

# Create a variable for the best model
best_rf = rand_search.best_estimator_

# Print the best hyperparameters
print('Best hyperparameters: ', rand_search.best_params_)
# Evaluate predictions with the best model
y_pred = best_rf.predict(X_test)

# Create the confusion matrix
cm = confusion_matrix(y_test, y_pred)

ConfusionMatrixDisplay(confusion_matrix=cm).plot()

C:\ProgramData\anaconda2\Lib\site-packages\sklearn\model_selection\
_split.py:789: Warning: The least populated class in y has only 4
members, which is less than n_splits=5
warnings.warn()

Best hyperparameters: {'max_depth': 1, 'min_samples_leaf': 1,
'min_samples_split': 4, 'n_estimators': 366}

```