

```

1 from random import random
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score, confusion_matrix
6 from sklearn.neighbors import KNeighborsClassifier
7 import matplotlib.pyplot as plt
8 import seaborn as sns
9
10
11 df = pd.read_csv('hsbdemo.csv')
12
13 print(df.head())
14
15 # converting features to numeric values
16 df.loc[df['gender'] == 'male', 'gender'] = 0
17 df.loc[df['gender'] == 'female', 'gender'] = 1
18
19 df.loc[df['ses'] == 'low', 'ses'] = 0
20 df.loc[df['ses'] == 'middle', 'ses'] = 1
21 df.loc[df['ses'] == 'high', 'ses'] = 2
22
23 df.loc[df['schtyp'] == 'public', 'schtyp'] = 0
24 df.loc[df['schtyp'] == 'private', 'schtyp'] = 1
25
26 df.loc[df['honors'] == 'not enrolled', 'honors'] = 0
27 df.loc[df['honors'] == 'enrolled', 'honors'] = 1
28
29 df.loc[df['prog'] == 'vocation', 'prog'] = 0
30 df.loc[df['prog'] == 'general', 'prog'] = 1
31 df.loc[df['prog'] == 'academic', 'prog'] = 2
32
33 print(df.head())
34
35
36 X = np.array(df.loc[:, ['gender', 'ses', 'schtyp', 'read', 'write', 'math', 'science', 'socst', 'honors', 'awards']])
37 y = np.array(df['prog'])
38 y = y.astype('int')
39
40 #use random_state=value to select the same data points in every run
41 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.10, random_state=3) #set seed to 42 all the time
42 knn = KNeighborsClassifier(n_neighbors=5)
43 knn.fit(X_train, y_train)
44
45 pred = knn.predict(X_test)
46

```

1.

```

pred = knn.predict(X_test)

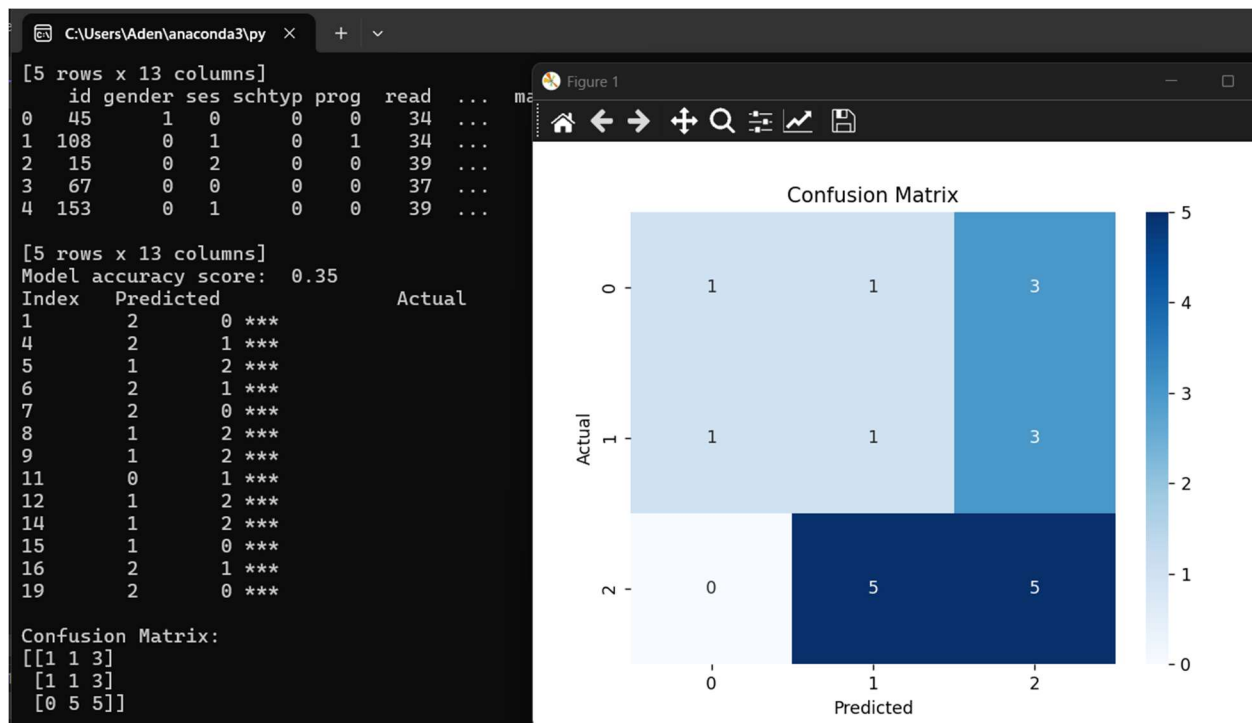
print('Model accuracy score: ', accuracy_score(y_test, pred))
print('Index\tPredicted\t\tActual')
for i in range(len(pred)):
    if pred[i] != y_test[i]:
        print(i, '\t', pred[i], '\t', y_test[i], '***')

conf_matrix = confusion_matrix(y_test, pred)

print(f'\nConfusion Matrix: \n{conf_matrix}')
sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues',
            xticklabels=knn.classes_, yticklabels=knn.classes_)

plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()

```



2.

```

1 from random import random
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score, confusion_matrix
6 from sklearn.neighbors import KNeighborsClassifier
7 import matplotlib.pyplot as plt
8 from sklearn.preprocessing import StandardScaler
9 from sklearn.decomposition import PCA
10
11
12 df = pd.read_csv('hsbdemo.csv')
13
14 print(df.head())
15
16 # converting features to numeric values
17 df.loc[df['gender'] == 'male', 'gender'] = 0
18 df.loc[df['gender'] == 'female', 'gender'] = 1
19
20 df.loc[df['ses'] == 'low', 'ses'] = 0
21 df.loc[df['ses'] == 'middle', 'ses'] = 1
22 df.loc[df['ses'] == 'high', 'ses'] = 2
23
24 df.loc[df['schtyp'] == 'public', 'schtyp'] = 0
25 df.loc[df['schtyp'] == 'private', 'schtyp'] = 1
26
27 df.loc[df['honors'] == 'not enrolled', 'honors'] = 0
28 df.loc[df['honors'] == 'enrolled', 'honors'] = 1
29
30 df.loc[df['prog'] == 'vocation', 'prog'] = 0
31 df.loc[df['prog'] == 'general', 'prog'] = 1
32 df.loc[df['prog'] == 'academic', 'prog'] = 2
33
34 print(df.head())
35
36 columns = ['gender', 'ses', 'schtyp', 'read', 'write', 'math', 'science', 'socst', 'honors', 'awards']
37 X = np.array(df.loc[:, columns])
38 y = np.array(df['prog'])
39 y = y.astype('int')
40
41 print(f'Before Standardization: {X}')
42 X = StandardScaler().fit_transform(X)
43 print(f'After Standardization: {X}')
44 X = StandardScaler().fit_transform(X)
45 print('Standard')
46 print(X)
47

```

```

48 x = np.arange(1,11)
49
50 pca = PCA(n_components=10)
51 principalComponents = pca.fit_transform(X)
52 explained_variance = pca.explained_variance_ratio_
53 print(f'Variance: {explained_variance}')
54 print(np.cumsum(explained_variance))
55
56 plt.plot(x, np.cumsum(explained_variance))
57 plt.xlabel('Principal Components')
58 plt.ylabel('Cumulative Ratio')
59 plt.title('PC= 1-10')
60 plt.xticks(range(1,11))
61 plt.show()
62
63

```

