

```
1 import pandas as pd
2 from matplotlib import pyplot as plt
3
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.cluster import KMeans
6
7
8 def readData():
9     df = pd.read_csv('data3/data_assignment3.csv')
10    return df
11
12
13 def normalizeData(table):
14     table[['nPosition', 'Nphi', 'Npsi']] =
15         StandardScaler().fit_transform(
16             table[['position', 'phi', 'psi']])
17     return table
18
19 def create2DHistogram(x_values, y_values):
20     plt.hist2d(x_values, y_values, bins=350, alpha=
21         1, cmap='plasma')
22
23 def createPandaXYValues(table, colum1, colum2):
24     return table[colum1], table[colum2]
25
26
27 def initPlot():
28     plt.xlim(-180, 180)
29     plt.ylim(-180, 180)
30     plt.ylabel("Psi")
31     plt.xlabel("Phi")
32     plt.show()
33
34
35 def regularScatterPlot(x_values, y_values):
36     plt.scatter(x_values, y_values, marker="o", s=1
37 )
38
39 if __name__ == '__main__':
```

```
39         x, y = createPandaXYValues(readData(), "phi", "psi")
40
41     # print(readData())
42     # print(readData().describe())
43     # print(normalizeData(readData()))
44     create2DHistogram(x, y)
45     # regularScatterPlot(x, y)
46     initPlot()
47
48
```

```
1 import pandas as pd
2 from matplotlib import pyplot as plt
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.cluster import KMeans
5 import numpy as np
6
7 def readData():
8     df = pd.read_csv('data3/data_assignment3.csv',
9         usecols=['phi', 'psi'])
10    return df
11
12 def normalizeData(table):
13     table[['Nphi', 'Npsi']] = StandardScaler().
14     fit_transform(table[['phi', 'psi']])
15     return table
16
17 def create2DHistogram(x_values, y_values):
18     plt.hist2d(x_values, y_values, bins=350, alpha=
19         1, cmap='plasma')
20
21 def createPandaXYValues(table, column1, column2):
22     return table[column1], table[column2]
23
24 def initPlot():
25     plt.xlim(-180, 180)
26     plt.ylim(-180, 180)
27     plt.ylabel("Psi")
28     plt.xlabel("Phi")
29     plt.show()
30
31 def regularScatterPlot(x_values, y_values):
32     plt.scatter(x_values, y_values, marker="o", s=1
33         )
34
35 def elbowMethod(data):
36     #Creating empty list to store Within-Cluster
37     Sum of Squares (WCSS) values
38     WCSS = []
39
40     #Creating for loop that runs kmeans with
```

```
36 different numbers of clusters 1-11
37     for i in range(1, 11):
38         kmeans = KMeans(n_clusters=i, init='k-means
++')
39         kmeans.fit(data)
40         #Appending the inertia value to the WCSS
list
41         WCSS.append(kmeans.inertia_)
42
43     #Plotting the elbow graph
44     plt.plot(range(1, 11), WCSS)
45     plt.title('The elbow method')
46     plt.xlabel('Number of clusters')
47     plt.ylabel('Within-Cluster Sum of Squares (WCSS
)')
48     plt.show()
49
50 if __name__ == '__main__':
51     x, y = createPandaXYValues(readData(), "phi", "
psi")
52     # print(readData())
53     # print(readData().describe())
54     # print(normalizeData(readData()))
55     create2DHistogram(x, y)
56     # regularScatterPLot(x, y)
57     initPlot()
58     elbowMethod(readData())
```

```
1 import pandas as pd
2 import numpy as np
3 from matplotlib import pyplot as plt
4 from sklearn.cluster import DBSCAN
5 from sklearn.neighbors import NearestNeighbors
6 import seaborn as sns
7 from collections import Counter
8
9
10 def readData():
11     df = pd.read_csv('data3/data_assignment3.csv')
12     df.to_numpy()
13     return df
14
15
16 def createPandaXYValues(table, column1, column2):
17     return table[column1], table[column2]
18
19
20 def createPandaNPAarray(dataFrame, column1, column2
21 ):
22     return dataFrame[[column1, column2]].to_numpy()
23
24
25 def createDBSCAN(epsilon, minsamples, dataFramed):
26     return DBSCAN(eps=epsilon, min_samples=
27 minsamples).fit(dataFramed)
28
29
30 def plotINIT():
31     plt.xlim(-180, 180)
32     plt.ylim(-180, 180)
33     plt.xlabel("Phi")
34     plt.ylabel("Psi")
35     plt.title("PRO")
36     plt.show()
37
38 def calculateOptimizedEpsilon(dataset, neighbours):
39     neighbors = NearestNeighbors(n_neighbors=
40 neighbours)
41     neighbors_fit = neighbors.fit(dataset)
42     distances, indices = neighbors_fit.kneighbors(
```

```
38 dataset)
39
40     distances = np.sort(distances, axis=0)
41     distances = distances[:, 1]
42     # plt.plot(distances)
43
44
45 def plot_outliers(data):
46     data['acids'] = list(dbscan.labels_)
47     outliers = data[data['acids'] == -1]
48     outliers['residue name'].value_counts().
        sort_values().plot.bar()
49
50
51
52 if __name__ == '__main__':
53     X = readData()
54     X = X.loc[X["residue name"] == "PRO"]
55     X = createPandaNPAarray(X, "phi", "psi")
56
57     dbscan = createDBSCAN(13, 110, X)
58     snsPlot = sns.scatterplot(x=X[:, 0], y=X[:, 1
        ], hue= dbscan.labels_, legend="full", palette="
        deep")
59     sns.move_legend(snsPlot, "upper right",
        bbox_to_anchor=(1.13, 1.15), title='Clusters')
60     plt.show()
61
```