**ADP_Red\ml\5.clustering.py**

```python
1   # %% 5. Machine Learning - Clustering Analysis
2   import numpy as np
3   import pandas as pd
4   import seaborn as sns
5   import matplotlib.pyplot as plt
6
7   import scipy.stats as stats
8
9
10  # %% 1. 데이터 수집
11  df = pd.read_csv('../../ADP_Python/data/USArrests.csv')
12
13  print(df.shape)
14  print(df.info())
15
16  # Check binary variable
17  for i, var in enumerate(df.columns):
18      print(i, var, len(df[var].unique()))
19
20  # Check data summary
21  print(df.describe())
22
23
24  # %% Hierarchical Clustering Analysis
25  from scipy.cluster.hierarchy import dendrogram, linkage, fcluster
26
27  # single linkage
28  model = linkage(df.iloc[:,1:], metric='euclidean', method='single')
29
30  # # ward linkage
31  model = linkage(df.iloc[:,1:], metric='euclidean', method='ward')
32
33  dendrogram(model,
34             labels=list(df.iloc[:,0]),
35             distance_sort='descending',
36             color_threshold=250)
37  plt.show()
38
39  assignment = fcluster(model, 250, 'distance')
40  print(assignment)
41
42
43  # %% Non-Hierarchical Clustering Analysis (k-Means)
44  df = pd.read_csv('../../ADP_Python/data/iris.csv')
45  X = df.copy().drop('target', axis=1)
46  y = df.copy()['target']
47
48  # k-Means
49  from sklearn.cluster import KMeans
50
51  # Scree Plot
52  sse = []
53
54  for k in range(1, 11):
55      kmeans = KMeans(n_clusters=k, n_init=10)
56      kmeans.fit(X)
57      sse.append(kmeans.inertia_)
```

```python
58
59  plt.plot(sse, marker='o')
60  plt.show()
61
62  # Visualization
63  df_result = X
64  df_result['prediction'] = KMeans(n_clusters=3, n_init=10).fit(X).predict(X)
65
66  sns.pairplot(df_result,
67              diag_kind='kde', hue='prediction')
68  plt.show()
69
70
71  # %% Non-Hierarchical Clustering Analysis (DBScan)
72  from sklearn.cluster import dbscan
73
74
75  # Gaussian Mixture
```