

# Project Phase 2 Machine Learning

Name: Abdelrahman Mohamed

ID: 202101081

- ***Project Steps***

- 1- Import libraries
- 2- Load the Data and Knew Some info about it
- 3- Implement linear regression algorithm
- 4- Apply agglomerative clustering and draw the dendrogram
- 5- Use the resulting clusters as initial centroids for the K-Means clustering algorithm.
- 6- Show Agglomerative Clustering and K-Means Clustering together
- 7- Evaluate Clustering Models
- 8- Implement a For Loop with a Stop Condition

- First dataset for the linear regression algorithm

real estate dataset, containing information such as transaction date, house age, distance to the nearest MRT station, number of convenience stores, latitude, longitude, and house price of unit area.

- Second dataset for the **clustering algorithms**

Absenteeism Dataset. The columns include information such as ID, reason for absence, month of absence, day of the week, seasons, transportation expense, distance from residence to work, service time, age, work load average per day, hit target, disciplinary failure, education, son, social drinker, social smoker, pet, weight, height, body mass index, and absenteeism time in hours.

- FULL CODE
- Import libraries

### import libraries

```
In [1]: 1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import mean_squared_error
5 import numpy as np
6 import matplotlib.pyplot as plt
7 import numpy as np
8 from scipy.cluster.hierarchy import dendrogram, linkage
9 from sklearn.cluster import KMeans
10 from sklearn.cluster import AgglomerativeClustering
11 import os
12 from scipy.spatial.distance import cdist
```

### Read the data and Knew some information

```
In [2]: 1 # Replace 'yourDataset.csv' with your actual dataset filename
2 df = pd.read_csv(r"Real estate valuation data set.csv", delimiter = ";")
3 df.head()
```

```
Out[2]:
```

	No	X1 transaction date	X2 house age	X3 distance to the nearest MRT station	X4 number of convenience stores	X5 latitude	X6 longitude	Y house price of unit area
0	1	2012.917	32.0	84.87882	10	24.98298	121.54024	37.9
1	2	2012.917	19.5	306.59470	9	24.98034	121.53951	42.2
2	3	2013.583	13.3	561.98450	5	24.98746	121.54391	47.3
3	4	2013.500	13.3	561.98450	5	24.98746	121.54391	54.8
4	5	2012.833	5.0	390.56840	5	24.97937	121.54245	43.1

```
In [3]: 1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 414 entries, 0 to 413
Data columns (total 8 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   No                                    414 non-null   int64
1   X1 transaction date                  414 non-null   float64
2   X2 house age                        414 non-null   float64
3   X3 distance to the nearest MRT station 414 non-null   float64
4   X4 number of convenience stores      414 non-null   int64
5   X5 latitude                          414 non-null   float64
6   X6 longitude                         414 non-null   float64
7   Y house price of unit area          414 non-null   float64
dtypes: float64(6), int64(2)
memory usage: 26.0 KB
```

### Implement the linear regression algorithm

```
In [5]: 1 # Replace 'yourDataset.csv' with your actual dataset filename
2 df = pd.read_csv(r"Real estate valuation data set.csv", delimiter = ";")
3 df.head()
```

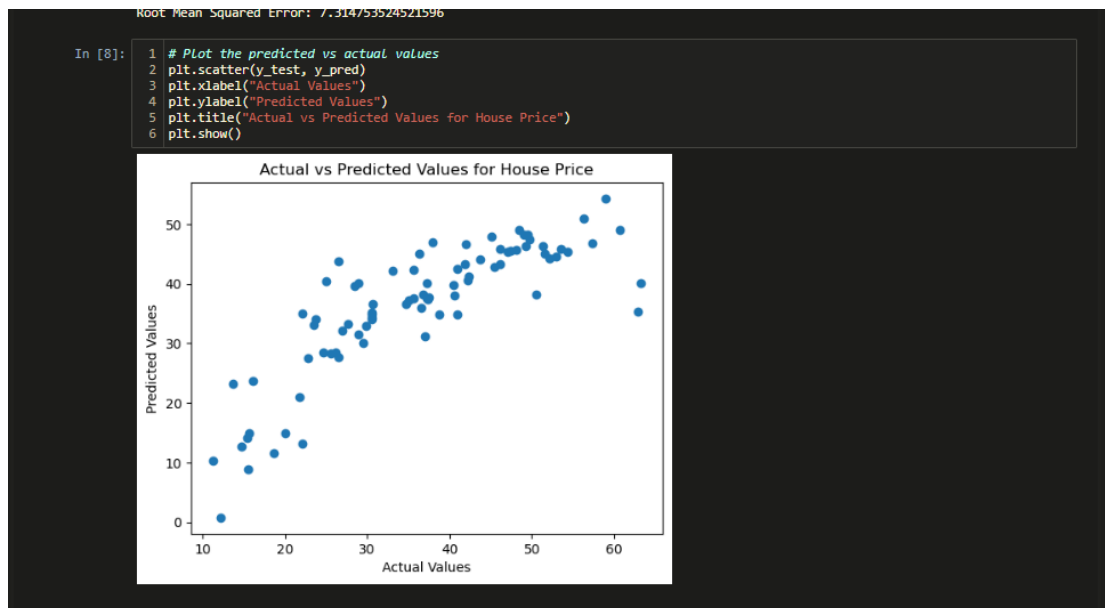
```
Out[5]:
```

	No	X1 transaction date	X2 house age	X3 distance to the nearest MRT station	X4 number of convenience stores	X5 latitude	X6 longitude	Y house price of unit area
0	1	2012.917	32.0	84.87882	10	24.98298	121.54024	37.9
1	2	2012.917	19.5	306.59470	9	24.98034	121.53951	42.2
2	3	2013.583	13.3	561.98450	5	24.98746	121.54391	47.3
3	4	2013.500	13.3	561.98450	5	24.98746	121.54391	54.8
4	5	2012.833	5.0	390.56840	5	24.97937	121.54245	43.1

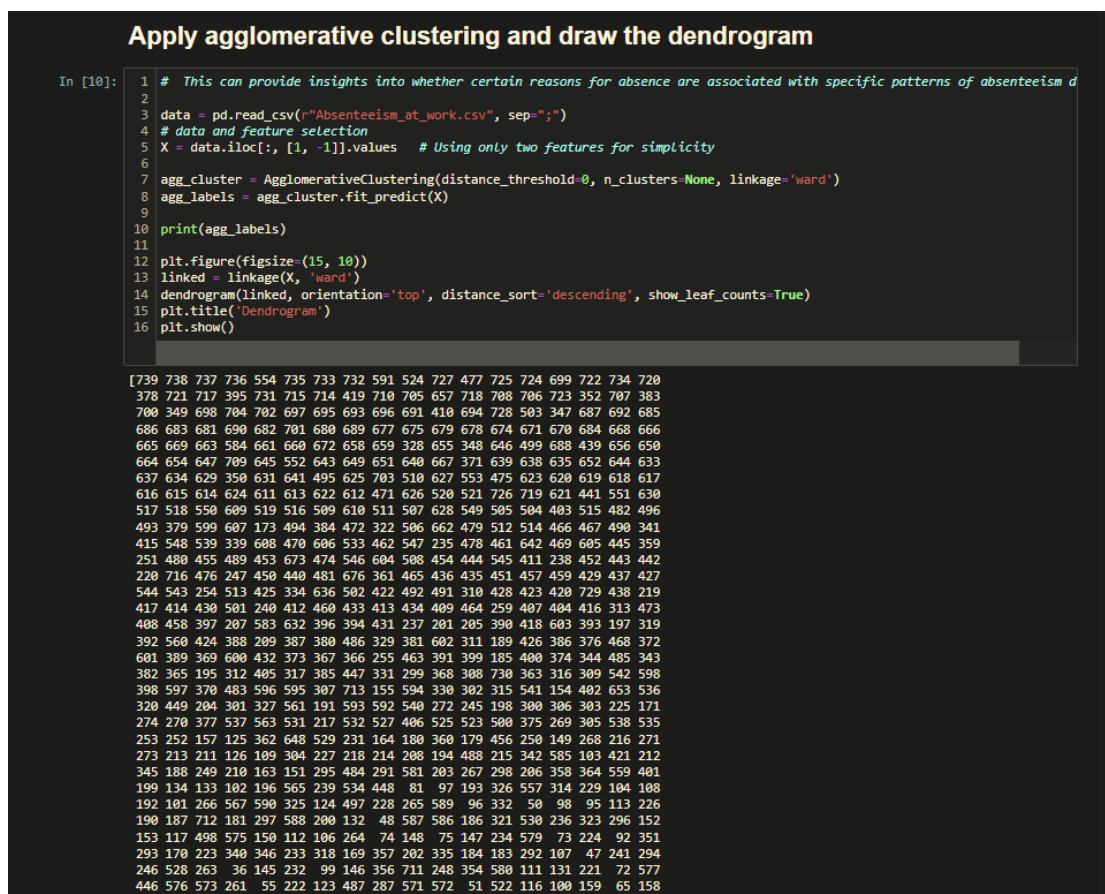
```
In [6]: 1 # Define features and target
2 features = ['X1 transaction date', 'X2 house age', 'X3 distance to the nearest MRT station', 'X4 number of convenience stores']
3 target = 'Y house price of unit area'
4
5 X = df[features]
6 y = df[target]
7
8 # Split the data into training and testing sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
In [7]: 1 # Initialize and train the linear regression model
2 model = LinearRegression()
3 model.fit(X_train, y_train)
4
5 # Make predictions on the test set
6 y_pred = model.predict(X_test)
7
8 # Evaluate the model
9 mse = mean_squared_error(y_test, y_pred)
10 rmse = np.sqrt(mse)
11
12 print(f'Mean Squared Error: {mse}')
13 print(f'Root Mean Squared Error: {rmse}')
```

```
Mean Squared Error: 53.50561912450112
Root Mean Squared Error: 7.314753524521596
```



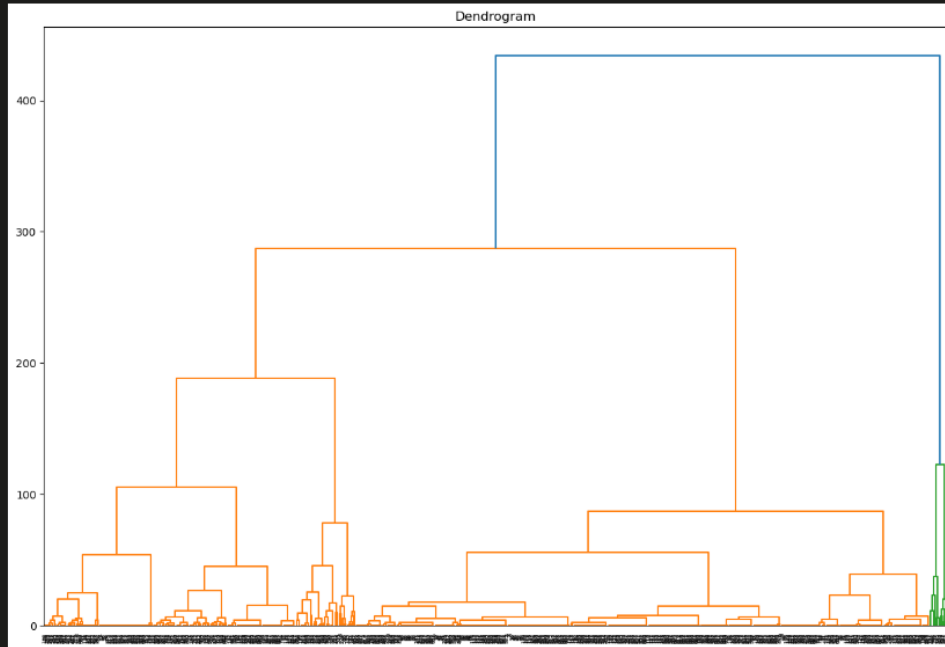
- Apply agglomerative clustering and draw the dendrogram



```

246 528 263 36 145 232 99 146 356 711 248 354 580 111 131 221 72 577
446 576 573 261 55 222 123 487 287 571 572 51 522 116 100 159 65 158
338 242 582 61 289 260 162 182 290 324 569 119 578 91 288 143 526 144
156 574 118 90 89 49 161 568 44 66 286 285 80 142 105 337 570 52
284 564 78 283 141 77 566 279 336 333 115 282 281 178 38 355 165 129
160 140 353 562 230 257 175 555 280 57 76 139 70 174 556 86 172 114
58 85 84 256 177 122 79 168 258 128 39 69 23 28 34 275 53 244
558 127 37 278 94 63 243 277 121 24 262 25 110 54 130 56 93 120
138 136 60 71 176 276 87 45 135 167 59 83 46 166 67 64 137 68
32 88 82 26 33 41 12 40 29 35 62 16 43 42 18 30 31 15
19 14 27 20 17 11 8 7 22 21 5 3 13 6 2 9 10 4
1 0]

```



- Use the resulting clusters as initial centroids for the K-Means clustering algorithm

### Use the resulting clusters as initial centroids for the K-Means clustering algorithm.

```
In [11]: 1 # the goal is to identify groups of instances that share similarities in terms of both the reason for absence and the duration
2 # This can provide insights into whether certain reasons for absence are associated with specific patterns of absenteeism d
3
4 X = data.iloc[:, [1, -1]].values # Using only two features for simplicity
5
6 # Agglomerative Clustering
7 agg_cluster = AgglomerativeClustering(distance_threshold=0, n_clusters=None, linkage='ward')
8 agg_labels = agg_cluster.fit_predict(X)
9
10 # Use the output of agglomerative clustering to determine the number of clusters
11 num_clusters = len(np.unique(agg_labels))
12
13 # K-Means with initial centroids from Agglomerative Clustering
14 kmeans = KMeans(n_clusters=num_clusters, random_state=42)
15 kmeans_labels = kmeans.fit_predict(X)
16 kmeans_centers = kmeans.cluster_centers_
17
18 # The Labels are integers ranging from 0 to (n_clusters - 1)
19 print(kmeans_labels)
20
21 plt.figure(figsize=(15, 10))
22 plt.scatter(X[:, 0], X[:, 1], c=kmeans_labels, cmap='viridis', edgecolors='k', s=20, alpha=0.5, linewidths=0.5)
23 plt.scatter(kmeans_centers[:, 0], kmeans_centers[:, 1], c='red', marker='X', s=200, label='Cluster Centers')
24 plt.title('K-Means Clustering')
25 plt.xlabel('Feature 1')
26 plt.ylabel('Feature 2')
27
28 plt.tight_layout()
29 plt.show()
```

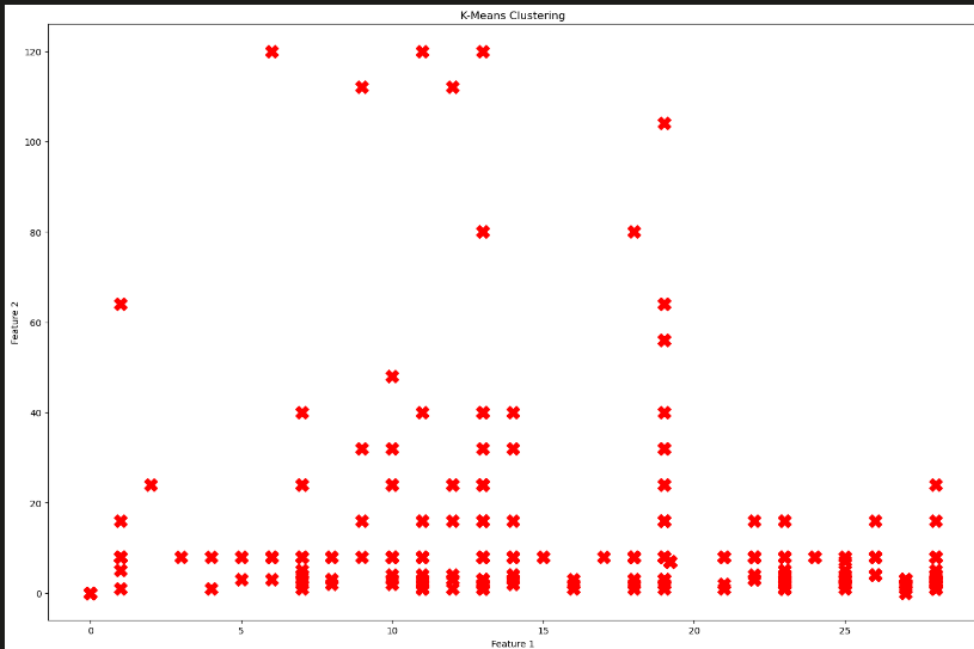
C:\Users\Lenovo\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1416: FutureWarning: The default value of 'n\_init' will change from 10 to 'auto' in 1.4. Set the value of 'n\_init' explicitly to suppress the warning  
super().\_check\_params\_vs\_input(X, default\_n\_init=10)  
C:\Users\Lenovo\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1440: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available threads. You can avoid it by setting the environment variable OM  
P\_NUM\_THREADS=3.  
warnings.warn(  
[ 49 0 2 5 2 2 8 12 13 8 16 16 23 17 12 28 2 33  
23 2 36 17 39 19 40 48 46 19 40 2 40 40 12 40 2 46  
67 68 15 12 2 12 12 75 2 78 78 12 8 91 0 0 2 2  
0 0 100 78 0 46 78 12 78 78 0 46 78 78 78 2 2 115  
47 78 131 28 133 54 46 128 46 46 122 33 115 151 8 153 46 157  
78 33 161 19 145 128 167 78 46 172 2 54 2 54 182 8 128 128  
2 8 50 75 2 46 196 199 40 50 2 133 199 50 50 50 50  
50 50 50 40 40 50 2 50 229 46 197 46 16 23 133 47 133 8  
23 16 133 2 46 47 217 40 46 46 75 133 75 262 263 8 122 122  
210 39 270 271 272 122 157 275 19 46 28 281 2 23 252 23 257 161  
290 133 210 145 145 46 145 281 46 133 19 281 46 78 28 145 15 182

```

252 78 78 468 647 54 40 314 2 78 2 640 312 638 327 327 161 327
36 100 50 40 50 525 122 525 50 122 78 16 50 128 50 50 133 50
502 50 73 525 525 217 516 2 525 486 50 50 122 50 610 217 50 608
50 2 327 50 50 327 50 2 217 217 601 50 50 525 327 525 217 75
122 50 525 50 593 516 525 73 50 47 327 50 78 525 2 525 525 47
281 525 525 516 573 8 122 36 75 75 122 50 167 566 78 2 562 8
50 75 559 50 19 75 382 50 8 262 133 552 54 252 133 28 172 8
50 2 541 145 0 50 0 537 536 36 8 36 533 36 2 133 50 73
2 72 36 17 73 525 54 36 8 402 75 50 0 0 262 516 133 161
217 8 508 507 197 49 502 157 19 500 499 2 8 8 351 28 68 0
0 0]

```

C:\Users\Lenovo\anaconda3\lib\site-packages\sklearn\base.py:1152: ConvergenceWarning: Number of distinct clusters (126) found smaller than n\_clusters (740). Possibly due to duplicate points in X.  
return fit\_method(estimator, \*args, \*\*kwargs)



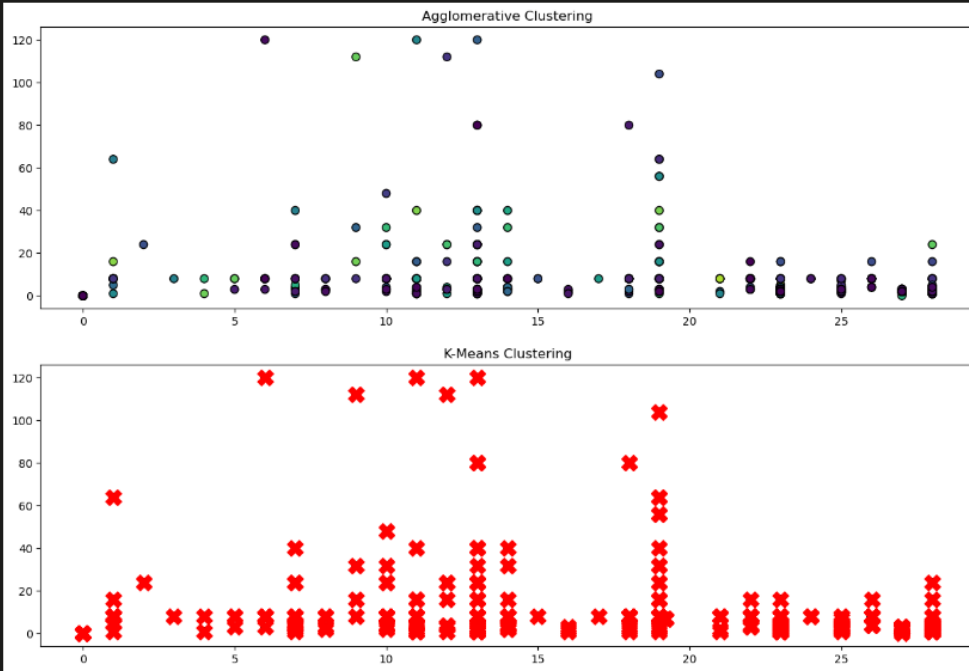
## Agglomerative Clustering and K-Means Clustering together

```

In [12]: 1 X = data.iloc[:, [1, -1]].values # Using only two features for simplicity
2
3 # Agglomerative Clustering
4 agg_cluster = AgglomerativeClustering(distance_threshold=0, n_clusters=None, linkage='ward')
5 agg_labels = agg_cluster.fit_predict(X)
6
7 # K-Means with initial centroids from Agglomerative Clustering
8 kmeans = KMeans(n_clusters=len(np.unique(agg_labels)), init='k-means++', random_state=42)
9 kmeans.fit(X)
10
11 # Plot the results
12 plt.figure(figsize=(15, 10))
13
14 # Plot Agglomerative Clustering results
15 plt.subplot(2, 1, 1)
16 plt.scatter(X[:, 0], X[:, 1], c=agg_labels, cmap='viridis', edgecolors='k', s=50)
17 plt.title('Agglomerative Clustering')
18
19 # Plot K-Means results
20 plt.subplot(2, 1, 2)
21 plt.scatter(X[:, 0], X[:, 1], c=kmeans.labels_, cmap='viridis', edgecolors='k', s=50)
22 plt.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], c='red', marker='x', s=200, label='K-Means Centroid')
23 plt.title('K-Means Clustering')
24
25 plt.show()

```

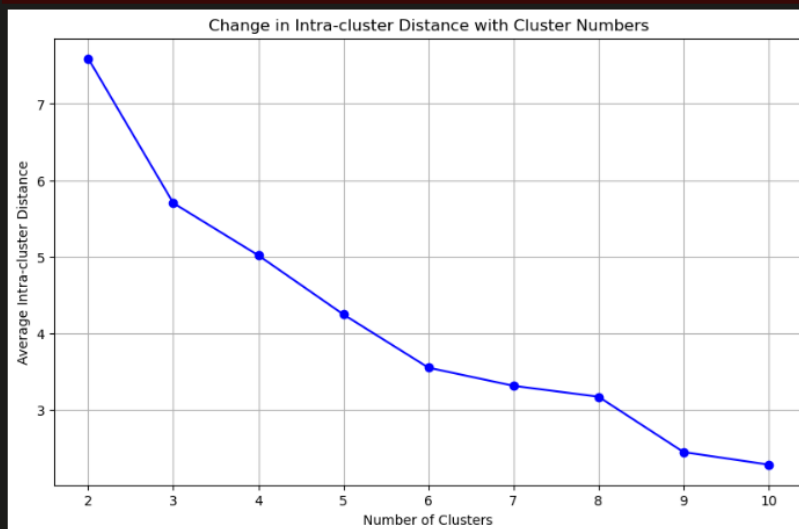
```
return fit_method(estimator, *args, **kwargs)
```



## Evaluate Clustering Models

```
In [13]: 1 import os
2 os.environ['JOBLIB_START_METHOD'] = 'forkserver'
3 os.environ['OMP_NUM_THREADS'] = '3'
4
5 X = data.iloc[:, [1, -1]].values # Using only two features for simplicity
6
7 # Calculate and visualize intra-cluster distance for different cluster numbers
8 cluster_numbers = range(2, 11)
9 intra_cluster_distances = []
10
11 for num_clusters in cluster_numbers:
12     kmeans = KMeans(n_clusters=num_clusters, n_init=1, random_state=42)
13     kmeans.fit(X)
14     intra_cluster_distances.append(np.mean(np.min(cdist(X, kmeans.cluster_centers_, 'euclidean'), axis=1)))
15
16 # Visualize the change in intra-cluster distance
17 plt.figure(figsize=(10, 6))
18 plt.plot(cluster_numbers, intra_cluster_distances, marker='o', linestyle='-', color='b')
19 plt.title('Change in Intra-cluster Distance with Cluster Numbers')
20 plt.xlabel('Number of Clusters')
21 plt.ylabel('Average Intra-cluster Distance')
22 plt.grid(True)
23 plt.show()
```

```
warnings.warn(
```



## Implement a For Loop with a Stop Condition

```
In [14]: 1 import os
2 os.environ['JOBLIB_START_METHOD'] = 'forkserver'
3 os.environ['OMP_NUM_THREADS'] = '3'
4
5 X = data.iloc[:, [1, -1]].values # Using only two features for simplicity
6
7 # Initialize variables for the loop
8 max_iterations = 20
9 threshold = 0.001
10 prev_intra_cluster_distance = float('inf')
11
12 # Calculate and visualize intra-cluster distance for different cluster numbers
13 cluster_numbers = range(2, max_iterations + 2)
14 intra_cluster_distances = []
15
16 for num_clusters in cluster_numbers:
17     kmeans = KMeans(n_clusters=num_clusters, n_init=1, random_state=42)
18     kmeans.fit(X)
19     intra_cluster_distance = np.mean(np.min(cdist(X, kmeans.cluster_centers_, 'euclidean'), axis=1))
20     intra_cluster_distances.append(intra_cluster_distance)
21
22     # Check if the change in intra-cluster distance is below the threshold
23     if abs(prev_intra_cluster_distance - intra_cluster_distance) < threshold:
24         print(f"Stopping at {num_clusters} clusters due to stabilization.")
25         break
26     prev_intra_cluster_distance = intra_cluster_distance
27
28 # Visualize the change in intra-cluster distance
29 plt.figure(figsize=(10, 6))
30 plt.plot(cluster_numbers, intra_cluster_distances, marker='o', linestyle='-', color='b')
31 plt.title('Change in Intra-cluster Distance with Cluster Numbers')
32 plt.xlabel('Number of Clusters')
33 plt.ylabel('Average Intra-cluster Distance')
34 plt.grid(True)
35 plt.show()
36
37
```

warnings.warn(

