

Clustering Assignment

Data Set : EastWestAirlines

1. Import Necessary libraries

```
In [1]: import pandas as pd

2. Import Data
```

```
In [2]: airline_details = pd.read_excel('EastWestAirlines.xlsx', sheet_name = "data")
airline_details

Out[2]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	
0	1	28143	0	1	1	1	174	1	0	0	7000	0
1	2	19244	0	1	1	1	215	2	0	0	6968	0
2	3	41354	0	1	1	1	4123	4	0	0	7034	0
3	4	14776	0	1	1	1	500	1	0	0	6952	0
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1
...
3994	4017	15476	0	1	1	1	8525	4	200	1	1403	1
3995	4018	64385	0	1	1	1	981	5	0	0	1395	1
3996	4019	73597	0	3	1	1	25447	8	0	0	1402	1
3997	4020	54899	0	1	1	1	500	1	500	1	1401	0
3998	4021	3016	0	1	1	1	0	0	0	0	1398	0

3999 rows × 12 columns

3. Data Understanding

```
In [3]: airline_details.head()

Out[3]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	
0	1	28143	0	1	1	1	174	1	0	0	7000	0
1	2	19244	0	1	1	1	215	2	0	0	6968	0
2	3	41354	0	1	1	1	4123	4	0	0	7034	0
3	4	14776	0	1	1	1	500	1	0	0	6952	0
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1

```
In [4]: airline_details.shape
Out[4]: (3999, 12)

In [5]: airline_details.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3999 entries, 0 to 3998
Data columns (total 12 columns):
#   Column                Non-Null Count  Dtype
---  --
0   ID#                    3999 non-null   int64
1   Balance                3999 non-null   float64
2   Qual_miles              3999 non-null   int64
3   cc1_miles                3999 non-null   int64
4   cc2_miles                3999 non-null   int64
5   cc3_miles                3999 non-null   int64
6   Bonus_miles             3999 non-null   int64
7   Bonus_trans             3999 non-null   int64
8   Flight_miles_12mo        3999 non-null   int64
9   Flight_trans_12         3999 non-null   int64
10  Days_since_enroll       3999 non-null   int64
11  Award#                  3999 non-null   int64
dtypes: int64(12)
memory usage: 375.0 KB

In [6]: airline_details.isna().sum()
Out[6]:
ID#                0
Balance            0
Qual_miles         0
cc1_miles          0
cc2_miles          0
cc3_miles          0
Bonus_miles        0
Bonus_trans        0
Flight_miles_12mo  0
Flight_trans_12    0
Days_since_enroll  0
Award#             0
dtype: int64

In [7]: airline_details.describe()

Out[7]:
```

	ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#
count	3999	3999.000000	3999000.000000	3999	3999.000000	3999	3999.000000	3999	3999.000000	3999	3999.000000	3999.000000
mean	2034.819456	7360333.064	144.11629	2.069015	1.014504	1.012253	17144.846212	11.60136	1460.095764	3.793172	4184.55938	0.37043
std	1180.764358	1.007791e+05	773.683094	1.376919	0.147650	0.195241	24150.967826	9.60381	1400.209171	3.793172	2085.13564	0.402857
min	1	1.000000	0.000000e+00	0.000000	1.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	1010	5.000000	1.852750e+04	0.000000	1.000000	1.000000	1250.000000	3.000000	0.000000	0.000000	2330.000000	0.000000
50%	2016	10.000000	3.369750e+04	0.000000	1.000000	1.000000	7170.000000	12.000000	0.000000	0.000000	4096.000000	0.000000
75%	3020	5.000000	1.244000e+04	0.000000	1.000000	1.000000	23800.500000	17.000000	311.000000	1.000000	5790.500000	1.000000
max	4021	9.000000	1.704838e+06	11148.000000	5.000000	3.000000	5.000000	263685.000000	86.000000	30817.000000	53.000000	8296.000000

3999 rows × 12 columns

```
In [8]: airline_details.dtypes
Out[8]:
ID#                int64
Balance            float64
Qual_miles         int64
cc1_miles          int64
cc2_miles          int64
cc3_miles          int64
Bonus_miles        int64
Bonus_trans        int64
Flight_miles_12mo  int64
Flight_trans_12    int64
Days_since_enroll  int64
Award#             int64
dtype: object

In [9]:
```

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
```

4. CLUSTERING

(a) HIERARCHICAL CLUSTERING

STEP 1 = Data Pre-processing

```
In [9]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [10]: airline_details = pd.read_excel('EastWestAirlines.xlsx', sheet_name = "data")
airline_details.head()

Out[10]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	
0	1	28143	0	1	1	1	174	1	0	0	7000	0
1	2	19244	0	1	1	1	215	2	0	0	6968	0
2	3	41354	0	1	1	1	4123	4	0	0	7034	0
3	4	14776	0	1	1	1	500	1	0	0	6952	0
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1

```
In [11]: from sklearn.preprocessing import MinMaxScaler

std_scaler_1 = MinMaxScaler()
X_1 = pd.DataFrame(std_scaler_1.fit_transform(airline_details.iloc[:,1:11]))
X_1

Out[12]:
```

	0	1	2	3	4	5	6	7	8	9	10
0	0.019508	0.0	0.00	0.0	0.000696	0.011629	0.000000	0.000000	0.843742	0.0	0
1	0.011288	0.0	0.00	0.0	0.000815	0.002356	0.000000	0.000000	0.839884	0.0	0
2	0.024257	0.0	0.00	0.0	0.015630	0.046512	0.000000	0.000000	0.847842	0.0	0
3	0.008667	0.0	0.00	0.0	0.001896	0.011628	0.000000	0.000000	0.837955	0.0	0
4	0.057338	0.0	0.75	0.0	0.164211	0.302326	0.067398	0.075477	0.835905	1.0	1
...
3994	0.010637	0.0	0.00	0.0	0.023230	0.046512	0.006490	0.018868	0.188917	1.0	1
3995	0.037766	0.0	0.00	0.0	0.003720	0.026140	0.000000	0.000000	0.167953	1.0	1
3996	0.043418	0.0	0.50	0.0	0.006458	0.008972	0.000000	0.000000	0.168997	1.0	1
3997	0.032022	0.0	0.00	0.0	0.003366	0.011628	0.004033	0.003088	0.160876	1.0	1
3998	0.001759	0.0	0.50	0.0	0.000000	0.000000	0.000000	0.000000	0.168314	0.0	0

3999 rows × 11 columns

STEP 2 = Finding the optimal number of clusters using the Dendrogram :

```
In [13]: from scipy.cluster.hierarchy import linkage
import scipy.cluster.hierarchy as sch

For creating dendrogram :: p = np.array(arr_norm) and converting into rumpy array format

In [14]: # create dendrogram
x = linkage(X_1, method = 'complete', metric = 'euclidean')

plt.figure(figsize = (15, 5))
plt.title('Hierarchical Dendrogram Plot')
plt.xlabel('Number of Nodes')
plt.ylabel('Euclidean Distances')

sch.dendrogram(x)
plt.show()

Out[15]:
```

STEP 3 = Training the hierarchical clustering model :

```
In [15]: from sklearn.cluster import AgglomerativeClustering

In [16]: hc = AgglomerativeClustering(n_clusters = 5, linkage = 'complete', affinity = 'euclidean').fit(X_1)
hc

Out[16]: AgglomerativeClustering(linkage='complete', n_clusters=5)

In [17]: x_hc = hc.fit_predict(X_1)
x_hc

Out[17]: array([0, 0, 0, ..., 2, 0, 0], dtype=int64)

In [18]: add_clusters = pd.DataFrame(x_hc, columns = ['clusters_1'])
add_clusters.head()

Out[18]:
```

	clusters_1
0	0
1	0
2	0
3	0
4	1

```
In [19]: airline_details['clusters_1'] = add_clusters
airline_details.head()

Out[19]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	clusters_1
0	1	28143	0	1	1	1	174	1	0	0	7000	0
1	2	19244	0	1	1	1	215	2	0	0	6968	0
2	3	41354	0	1	1	1	4123	4	0	0	7034	0
3	4	14776	0	1	1	1	500	1	0	0	6952	0
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1

325 rows × 13 columns

STEP 4 = Visualizing the clusters :

```
In [20]: airline_details[airline_details['clusters_1']==1]

Out[20]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	clusters_1
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1
12	12	96952	0	5	1	1	6138	19	0	0	6954	1
15	16	20465	0	4	1	1	49442	15	0	0	6912	1
16	17	51390	0	4	1	1	49963	16	0	0	6910	1
20	21	120976	0	5	1	1	58831	23	250	2	6995	1
...
3753	3776	70178	0	5	1	1	61530	12	1300	6	3306	1
3772	3795	822321	0	5	1	1	138334	26	600	2	3288	1
3846	3869	97510	1678	5	1	1	71609	51	7650	26	1665	1
3855	3878	190730	0	5	1	1	78916	28	2450	9	1448	1
3883	3906	128167	0	5	1	1	189160	15	0	0	1512	1

325 rows × 13 columns

STEP 4 = Visualizing the clusters :

```
In [21]: airline_details[airline_details['clusters_1']==2]

Out[21]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	clusters_1
7	8	20856	0	1	1	1	5250	4	250	1	6938	1
8	9	442003	0	3	2	1	1753	43	3850	12	6948	1
9	10	134880	0	3	1	1	28426	28	1150	3	6931	1
17	18	12558	0	1	1	1	4291	5	0	0	6955	1
18	19	91473	0	3	1	1	27408	17	0	0	6953	1
...
3987	4010	11933	0	1	1	1	249	3	79	1	1412	1
3989	4012	2622	0	1	1	1	1625	6	0	0	1404	1
3994	4017	15476	0	1	1	1	8525	4	200	1	1403	1
3995	4018	64385	0	1	1	1	981	5	0	0	1395	1
3996	4019	73597	0	3	1	1	25447	8	0	0	1402	1

1144 rows × 13 columns

STEP 4 = Visualizing the clusters :

```
In [22]: airline_details[airline_details['clusters_1']==3]

Out[22]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	clusters_1
2015	2032	53232	888	4	1	1	80096	65	22100	45	3831	1
3235	3257	287033	0	1	1	1	26161	58	12873	53	2272	1
3583	3606	100114	500	1	1	1	71954	86	30817	53	1373	1
3584	3617	27619	0	4	1	1	83726	68	14050	46	1255	1

3584 rows × 13 columns

(b) K - MEANS

STEP 1 = Data Pre-processing :

```
In [23]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [24]: airline_details = pd.read_excel('EastWestAirlines.xlsx', sheet_name = "data")
airline_details.head()

Out[24]:
```

ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo	Flight_trans_12	Days_since_enroll	Award#	
0	1	28143	0	1	1	1	174	1	0	0	7000	0
1	2	19244	0	1	1	1	215	2	0	0	6968	0
2	3	41354	0	1	1	1	4123	4	0	0	7034	0
3	4	14776	0	1	1	1	500	1	0	0	6952	0
4	5	97752	0	4	1	1	43300	26	2077	4	6935	1

```
In [25]: from sklearn.preprocessing import StandardScaler

std_scaler_2 = StandardScaler()
X_2 = std_scaler_2.fit_transform(airline_details.iloc[:,1:11])
print(X_2)

[[[-1.51401038e-01, -1.86398887e-01, -7.69578486e-01, ..., -3.62167870e-01,
  1.35453489e+00, -7.66913296e-01],
 [-1.58456876e-01, -1.86298887e-01, -7.69578486e-01, ..., -3.62167870e-01,
  1.37995122e+00, -7.66913296e-01],
 [-1.58456876e-01, -1.86298887e-01, -7.69578486e-01, ..., -3.62167870e-01,
  1.37995122e+00, -7.66913296e-01],
 [-1.58456876e-01, -1.86298887e-01, -7.69578486e-01, ..., -3.62167870e-01,
  1.37995122e+00, -7.66913296e-01],
 [-1.58456876e-01, -1.86298887e-01, -7.69578486e-01, ..., -3.62167870e-01,
  1.37995122e+00, -7.66913296e-01],
 [-1.58456876e-01, -1.86298887e-
```