

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import scipy.cluster.hierarchy as sch
from sklearn.cluster import AgglomerativeClustering
```

```
In [2]: # Import Dataset
airline=pd.read_csv('EastWestAirlines.csv')
airline
```

```
Out[2]:
```

	ID#	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_miles_12mo
0	1	28143	0	1	1	1	174	1	1
1	2	19244	0	1	1	1	215	2	2
2	3	41354	0	1	1	1	4123	4	4
3	4	14776	0	1	1	1	500	1	1
4	5	97752	0	4	1	1	43300	26	26
...	...	...	...	...	...	...	...	...	...
3994	4017	18476	0	1	1	1	8525	4	4
3995	4018	64385	0	1	1	1	981	5	5
3996	4019	73597	0	3	1	1	25447	8	8
3997	4020	54899	0	1	1	1	500	1	1
3998	4021	3016	0	1	1	1	0	0	0

3999 rows × 12 columns



```
In [3]: airline.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  int64
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 [4]: airline2=airline.drop(['ID#'],axis=1)
airline2
```

```
Out[4]:
```

	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_mil
0	28143	0	1	1	1	174	1	
1	19244	0	1	1	1	215	2	
2	41354	0	1	1	1	4123	4	
3	14776	0	1	1	1	500	1	
4	97752	0	4	1	1	43300	26	
...	...	...	...	...	...	...	...	
3994	18476	0	1	1	1	8525	4	
3995	64385	0	1	1	1	981	5	
3996	73597	0	3	1	1	25447	8	
3997	54899	0	1	1	1	500	1	
3998	3016	0	1	1	1	0	0	

3999 rows × 11 columns



```
In [5]: # Normalize heterogenous numerical data using z-score (x-mean/std) or custom defi
# Normalization function - here custom defined
def norm_func(i):
    x = (i-i.min())/(i.max()-i.min())
    return (x)
```

In [6]: *# Normalized data frame (considering the numerical part of data)*

```
airline2_norm = norm_func(airline2)
airline2_norm
```

Out[6]:

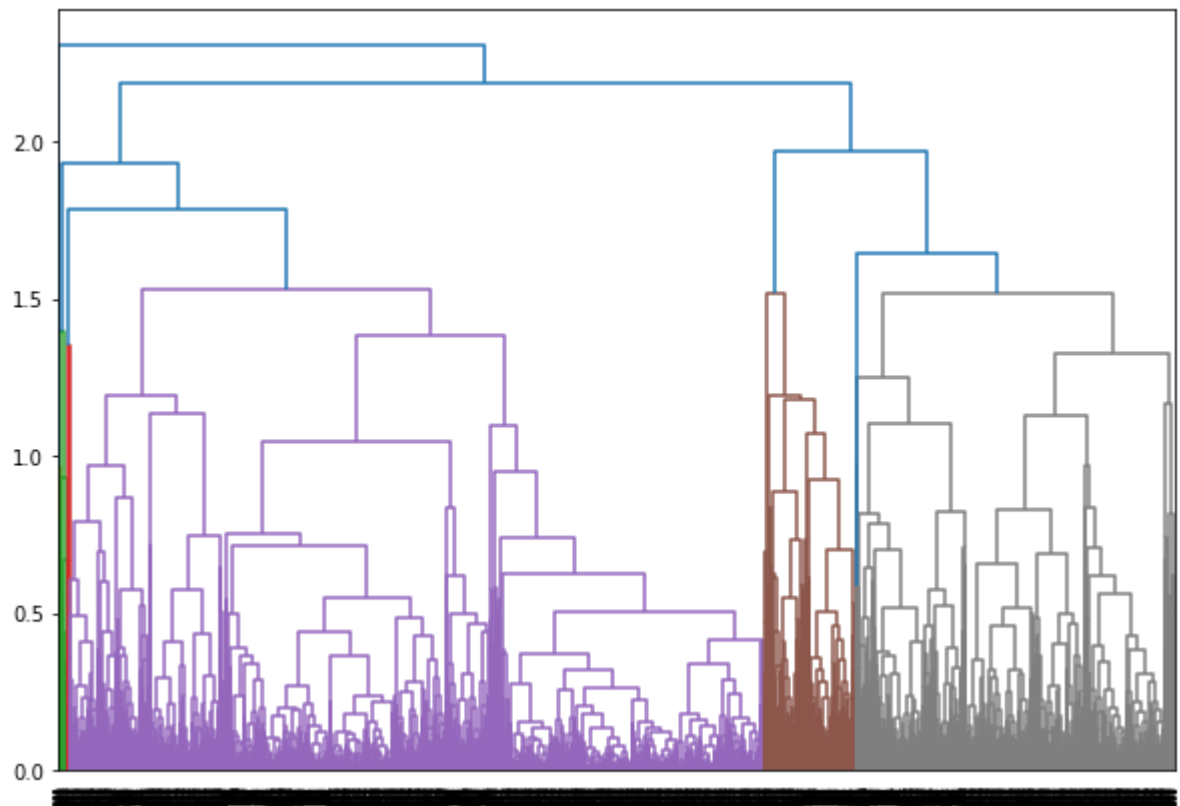
	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_m
0	0.016508	0.0	0.00	0.0	0.0	0.000660	0.011628	
1	0.011288	0.0	0.00	0.0	0.0	0.000815	0.023256	
2	0.024257	0.0	0.00	0.0	0.0	0.015636	0.046512	
3	0.008667	0.0	0.00	0.0	0.0	0.001896	0.011628	
4	0.057338	0.0	0.75	0.0	0.0	0.164211	0.302326	
...	...	...	...	...	...	...	...	...
3994	0.010837	0.0	0.00	0.0	0.0	0.032330	0.046512	
3995	0.037766	0.0	0.00	0.0	0.0	0.003720	0.058140	
3996	0.043169	0.0	0.50	0.0	0.0	0.096505	0.093023	
3997	0.032202	0.0	0.00	0.0	0.0	0.001896	0.011628	
3998	0.001769	0.0	0.00	0.0	0.0	0.000000	0.000000	

3999 rows × 11 columns



In [7]: *# Create Dendrograms*

```
plt.figure(figsize=(10, 7))
dendograms=sch.dendrogram(sch.linkage(airline2_norm,'complete'))
```



```
In [8]: # Create Clusters (y)
hclusters=AgglomerativeClustering(n_clusters=5,affinity='euclidean',linkage='ward')
hclusters
```

Out[8]: AgglomerativeClustering(n\_clusters=5)

```
In [9]: y=pd.DataFrame(hclusters.fit_predict(airline2_norm),columns=['clustersid'])
y['clustersid'].value_counts()
```

Out[9]:

1	1011
0	946
2	808
4	699
3	535

Name: clustersid, dtype: int64

```
In [10]: # Adding clusters to dataset
airline2['clustersid']=hclusters.labels_
airline2
```

Out[10]:

	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_trans	Flight_mil
0	28143	0	1	1	1	174	1	
1	19244	0	1	1	1	215	2	
2	41354	0	1	1	1	4123	4	
3	14776	0	1	1	1	500	1	
4	97752	0	4	1	1	43300	26	
...	...	...	...	...	...	...	...	...
3994	18476	0	1	1	1	8525	4	
3995	64385	0	1	1	1	981	5	
3996	73597	0	3	1	1	25447	8	
3997	54899	0	1	1	1	500	1	
3998	3016	0	1	1	1	0	0	

3999 rows × 12 columns



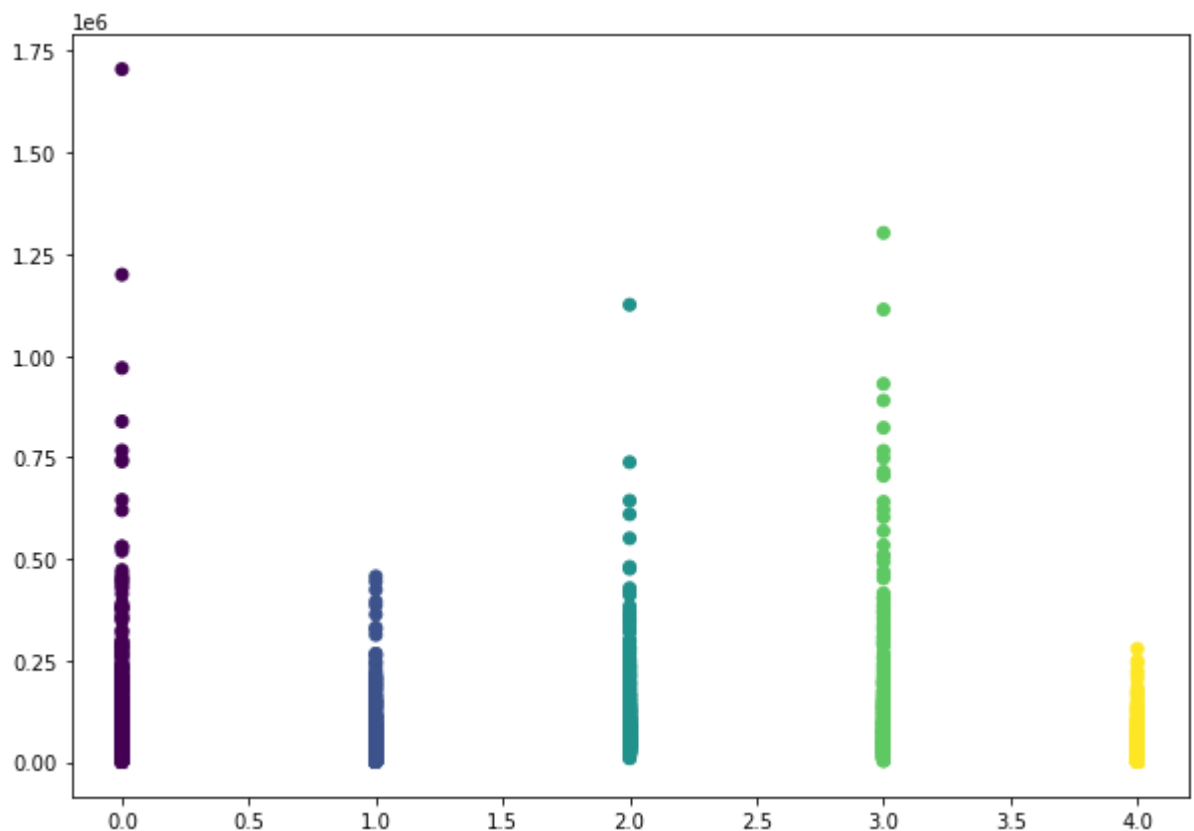
```
In [11]: airline2.groupby('clustersid').agg(['mean']).reset_index()
```

```
Out[11]:
```

	clustersid	Balance	Qual_miles	cc1_miles	cc2_miles	cc3_miles	Bonus_miles	Bonus_tr
		mean	mean	mean	mean	mean	mean	m
0	0	79848.233615	285.097252	1.699789	1.024313	1.000000	12079.774841	12.133
1	1	43313.653808	21.506429	1.000000	1.033630	1.000989	2562.614243	5.474
2	2	106221.111386	161.262376	3.198020	1.001238	1.025990	26458.257426	16.363
3	3	127475.028037	160.801869	4.362617	1.000000	1.050467	58656.919626	22.235
4	4	30013.416309	98.054363	1.000000	1.000000	1.000000	2552.569385	6.101

```
In [12]: # Plot Clusters
plt.figure(figsize=(10, 7))
plt.scatter(airline2['clustersid'],airline2['Balance'], c=hclusters.labels_)
```

```
Out[12]: <matplotlib.collections.PathCollection at 0x21ee50dec40>
```



```
In [ ]:
```

