

## Lab 6.1 Hierarchical Clustering

ລະຫັດນັກສຶກສາ: 205Q0010.19

ຊື່ ແລະ ນາມສະກຸນ: ທ້າວ ນຸຊົວ ເຮີ 3CW1

ຈົ່ງຕອບຄໍາຖາມຕໍ່ໄປນີ້ໃຫ້ສໍາເລັດດ້ວຍການນໍາໃຊ້ຄໍາສັ່ງຂອງ Python:

ພາກທີ 1

1.1. ຈາກຊຸດຂໍ້ມູນMall\_Customers.csvຈົ່ງເຕີມຂໍ້ມູນໃນຕາຕະລາງໃຫ້ສໍາເລັດ.

```
dataset = pd.read_csv('Mall_Customers.csv')
X = dataset.iloc[:, [3, 4]].values
print(dataset)
```

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
...	...	...	...	...	...
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

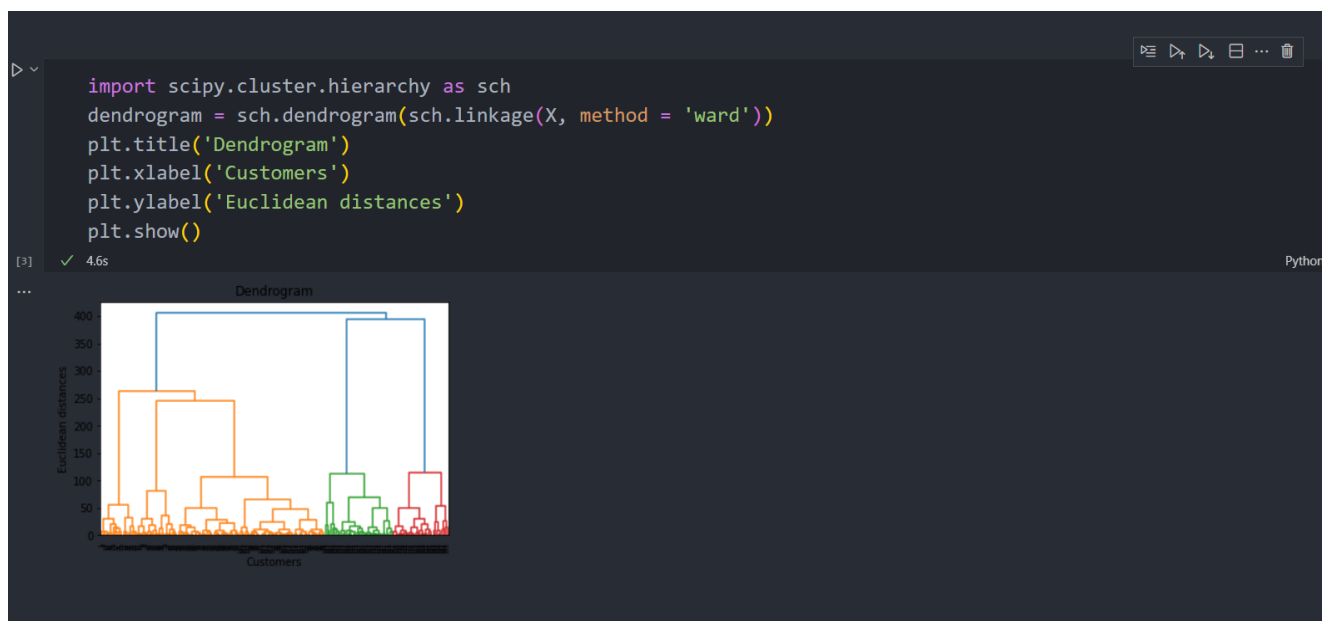
[200 rows x 5 columns]

1.2. ຈົ່ງທໍາການເລືອກຄຸນລັກສະນະ (Features): Age ແລະ Annual Income (k\$)

```
dataset = pd.read_csv('Mall_Customers.csv')
X = dataset.iloc[:, [3, 4]].values
print(dataset)
```

Age	Annual Income (k\$)
19	15
21	15
20	16
23	16
31	17
...	...
35	120
45	126
32	126
32	137
30	137

### 1.3 ຈົ່ງຊອກຫາຈຳນວນກຸ່ມເພື່ອແບ່ງກຸ່ມຂຸດຂໍ້ມູນໃຫ້ເມາະສົມດ້ວຍDendrogram



### 1.4. ຈົ່ງທຳການຝຶກຂໍ້ມູນໃນຂັ້ນທີ 1.2 ເພື່ອແບ່ງຕາມຈຳນວນກຸ່ມໃນຂັ້ນທີ 1.3 ດ້ວຍແບບຈຳຮອງ (Algorithm) AgglomerativeClustering.

```
from sklearn.cluster import AgglomerativeClustering
hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean', linkage = 'ward')
y_hc = hc.fit_predict(X)
```

[31] ✓ 0.4s Python

### 1.5. ຈົ່ງທຳການສຳຫຼວດ (visulaizing) ກຸ່ມຂໍ້ມູນທີ່ຖືກແບ່ງດ້ວຍscatter plot ພ້ອມອະທິບາຍຜົນ.



- ❖ Cluster1 (ສີແດງ): ລາຍຮັບສູງ ແລະ ຄະແນນການໃຊ້ຈ່າຍຕ່ຳ
- ❖ Cluster2 (ສີຟ້າ): ລາຍໄດ້ປົກກະຕິ ແລະ ຄະແນນການໃຊ້ຈ່າຍປົກກະຕິ
- ❖ Cluster3 (ສີຂຽວ): ລາຍຮັບສູງ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ
- ❖ Cluster4 (ສີຟ້າຂຽວ): ລາຍຮັບຕ່ຳ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ
- ❖ Cluster5 (ສີມ່ວງແດງ): ລາຍຮັບຕ່ຳ ແລະ ຄະແນນການໃຊ້ຈ່າຍຕ່ຳ

## ພາກທີ 2

### 2.1 ຈາກຊຸດຂໍ້ມູນ BaskinRobbins.csv ຈົ່ງຕື່ມຂໍ້ມູນໃນຕາຕະລາງໃຫ້ສໍາເລັດ

```
dataset = pd.read_csv('BaskinRobbins.csv')
X = dataset.iloc[:, [3, 4]].values
print(dataset)
```

[22] ✓ 0.0s Python

Output exceeds the size limit. Open the full output data in a text editor

Unnamed: 0	Flavour	Calories	Total Fat (g) \
0	Bananas Foster	160	8.0
1	Baseball Nut	160	9.0
2	Beavertails Pastry	170	9.0
3	Blackberry Frozen Yogurt	120	4.0
4	Blue Raspberry Sherbet	130	2.0
...	...	...	...
65	Very Berry Strawberry	200	10.0
66	Watermelon Splash	120	0.5
67	Wild 'n Reckless	80	1.5
68	Winter White Chocolate	160	8.0
69	World Class Chocolate	260	60.0

Trans Fat (g)	Carbohydrates (g)	Sugars (g)	Protein (g) \
0.2	20	16	2.0
0.2	19	13	3.0
0.3	21	15	3.0
0.1	17	16	3.0
0.1	26	20	2.0
...	...	...	...
0.4	24	21	4.0
0.0	27	20	0.1
0.0	16	12	1.0
0.2	20	16	2.0
0.4	25	18	5.0
...	...	...	...
NaN	NaN	NaN	NaN
NaN	NaN	NaN	NaN
NaN	NaN	NaN	NaN

[70 rows x 9 columns]

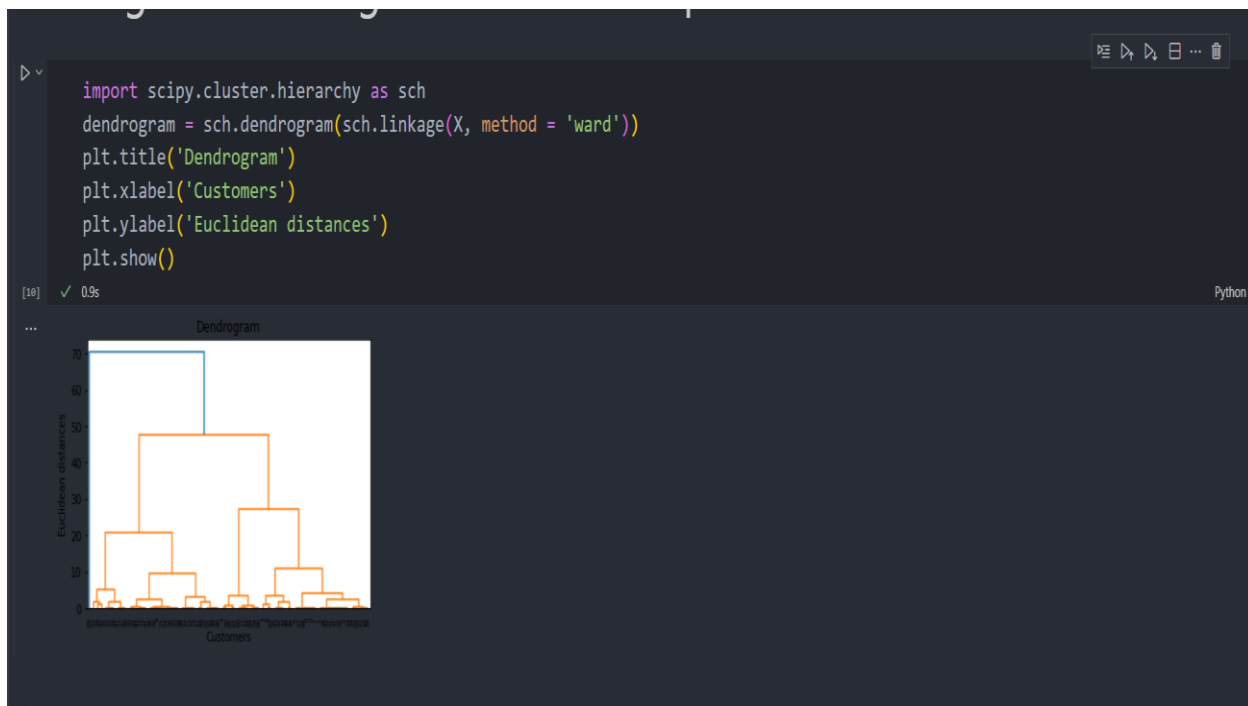
### 2.2. ຈົ່ງທຳການເລືອກຄຸນລັກສະນະ (Features): Carbohydrates (g) ແລະ Sugars (g)

```
dataset = pd.read_csv('BaskinRobbins.csv')
X = dataset.iloc[:, [3, 4]].values
print(dataset)
```

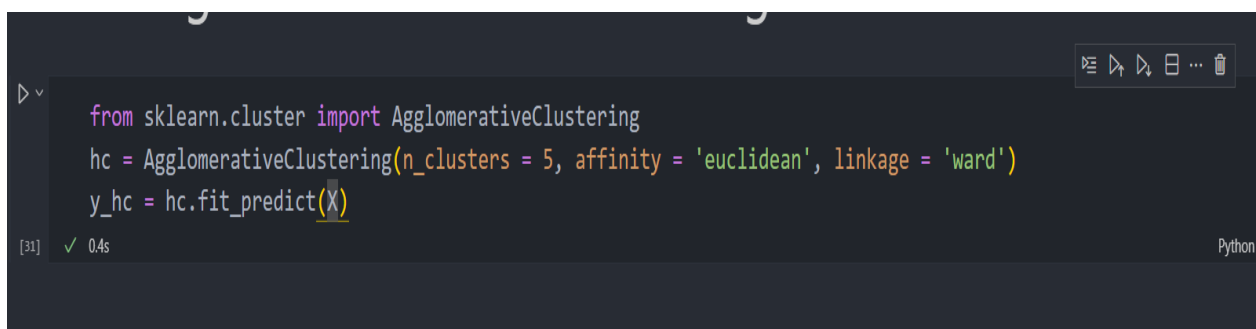
[19] ✓ 0.1s Python

Carbohydrates (g)	Sugars (g)
20	16
19	13
21	15
17	16
26	20
...	...
24	21
27	20
16	12
20	16
25	18

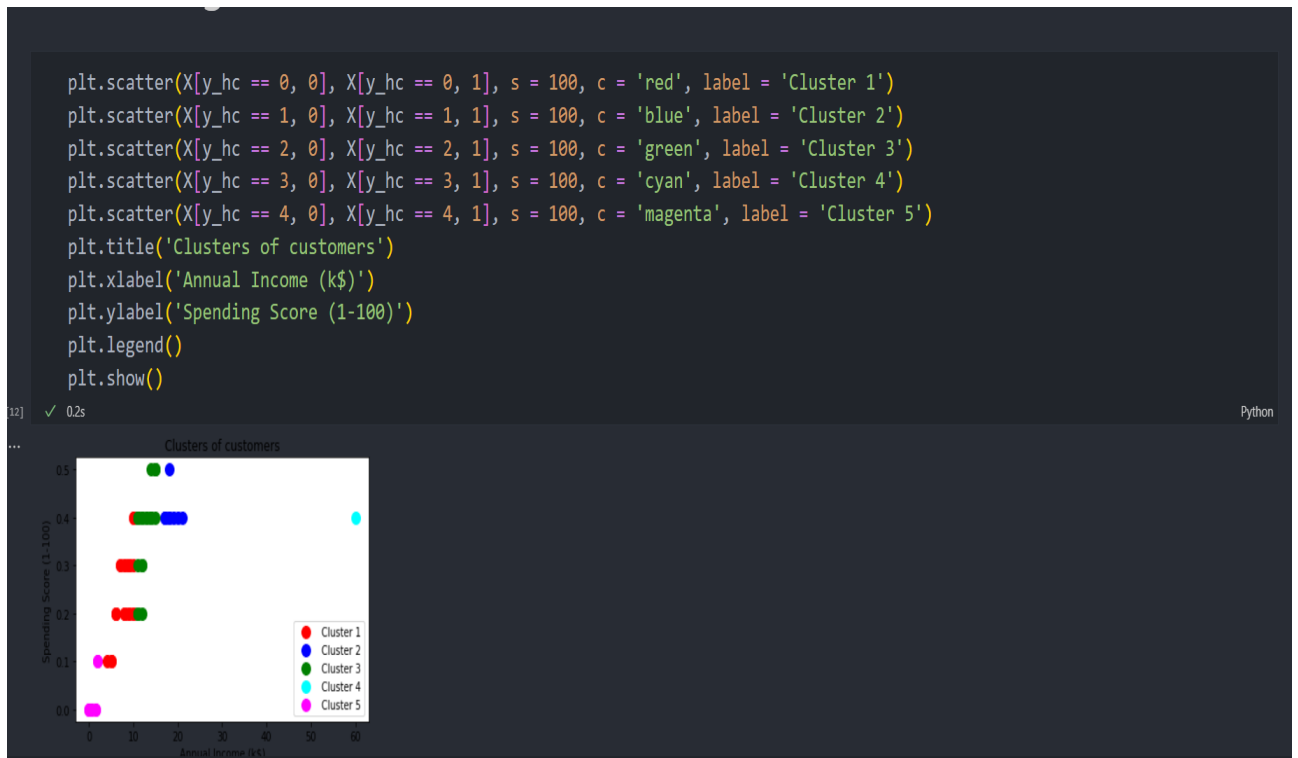
## 2.3. ຈົ່ງຊອກຫາຈຳນວນກຸ່ມເພື່ອແບ່ງກຸ່ມຂຸດຂໍ້ມູນໃຫ້ເໝາະສົມດ້ວຍDendrogram



## 2.4. ຈົ່ງທຳການຝຶກຂໍ້ມູນໃນຂໍ້ທີ 2.2 ເພື່ອແບ່ງຕາມຈຳນວນກຸ່ມໃນຂໍ້ທີ 2.3 ດ້ວຍແບບຈຳຮອງ (Algorithm) AgglomerativeClustering.



## 2.5. ຈົ່ງທຳການສຳຫຼວດ (visulaizing) ກຸ່ມຂໍ້ມູນທີ່ຖືກແບ່ງດ້ວຍ scatter plot ພ້ອມອະທິບາຍຜົນ.



- ❖ Cluster1 (ສີແດງ): ລາຍຮັບຕ່ຳ ແລະ ຄະແນນການໃຊ້ຈ່າຍບາງເທື່ອກໍ່ສູງບາງເທື່ອກໍ່ຕ່ຳ
- ❖ Cluster2 (ສີຟ້າ): ລາຍໄດ້ປົກກະຕິ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ
- ❖ Cluster3 (ສີຂຽວ): ລາຍຮັບຕ່ຳ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ
- ❖ Cluster4 (ສີຟ້າຂຽວ): ລາຍຮັບຕ່ຳ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ
- ❖ Cluster5 (ສີມ່ວງແດງ): ລາຍຮັບສູງ ແລະ ຄະແນນການໃຊ້ຈ່າຍສູງ