

1.

```
7
8 print(z_score_normalized)
```

	income	age
0	0.550812	1.693570
1	-0.777331	-1.130565
2	-1.110474	1.146963
3	-1.194372	0.873660
4	-0.181416	1.238064
...
995	-0.242008	-0.948363
996	0.381668	0.144851
997	-0.965654	-0.401756
998	-1.608203	-0.766161
999	1.114049	0.418154

[1000 rows x 2 columns]

```
In [123]: 1 data['age_Z']=z_score_normalized['age']
          2 data['income_Z']=z_score_normalized['income']
          3 data.head()
```

Out[123]:

	income	age	age_Z	income_Z
0	101743	58	1.693570	0.550812
1	49597	27	-1.130565	-0.777331
2	36517	52	1.146963	-1.110474
3	33223	49	0.873660	-1.194372
4	72994	53	1.238064	-0.181416

2.(1)

```
In [124]: 1 from sklearn import cluster
          2 model = cluster.KMeans(n_clusters=2, random_state=10)
          3 X = data[['age_Z', 'income_Z']].values
          4 model.fit_predict(X)
          5 cluster_assignment = model.labels_
          6 centers = model.cluster_centers_
```

```
In [125]: 1 import numpy as np
          2 print(np.sum((X - centers[cluster_assignment])**2))
```

1189.7476232504307

```
In [126]: 1 import matplotlib.pyplot as plt
          2 ss = []
          3 krange = list(range(2, 11))
          4 X = data[['age_Z', 'income_Z']].values
          5 for n in krange:
          6     model = cluster.KMeans(n_clusters=n, random_state=10)
          7     model.fit_predict(X)
          8     cluster_assignments = model.labels_
          9     centers = model.cluster_centers_
         10     ss.append(np.sum((X-centers[cluster_assignments])**2))
         11 ss
```

Out[126]: [1189.7476232504307

2.(2)

```
In [125]: 1 import numpy as np
2 print(np.sum((X - centers[cluster_assignment])**2))
```

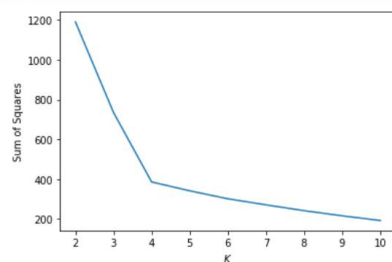
1189.7476232504307

```
In [126]: 1 import matplotlib.pyplot as plt
2 ss = []
3 krange = list(range(2, 11))
4 X = data[['age_Z', 'income_Z']].values
5 for n in krange:
6     model = cluster.KMeans(n_clusters=n, random_state=10)
7     model.fit_predict(X)
8     cluster_assignments = model.labels_
9     centers = model.cluster_centers_
10    ss.append(np.sum((X-centers[cluster_assignments])**2))
11 ss
```

```
Out[126]: [1189.7476232504307,
734.6406648547709,
386.0702389582008,
341.6447741546658,
301.6035411699272,
271.16618461596397,
241.90179918136994,
216.1272992902626,
192.56868717475388]
```

2.(3)

```
In [127]: 1 plt.plot(krange, ss)
2 plt.xlabel('$K$')
3 plt.ylabel('Sum of Squares')
4 plt.show()
5 #4 篇最佳分群數
```



```
In [129]: 1 from sklearn import cluster
2 model = cluster.KMeans(n_clusters=4, random_state=10)
3 X = data[['age_Z', 'income_Z']].values
4 model.fit_predict(X)
5 cluster_assignment = model.labels_
```

3.

```
In [129]: 1 from sklearn import cluster
2 model = cluster.KMeans(n_clusters=4, random_state=10)
3 X = data[['age_Z', 'income_Z']].values
4 model.fit_predict(X)
5 cluster_assignment = model.labels_
6 centers = model.cluster_centers_
7 data['kmeans_cluster']=cluster_assignment
8 data.head()
```

```
Out[129]:
```

	income	age	age_Z	income_Z	kmeans_cluster
0	101743	58	1.693570	0.550812	0
1	49597	27	-1.130565	-0.777331	1
2	36517	52	1.146963	-1.110474	2
3	33223	49	0.873660	-1.194372	2
4	72994	53	1.238064	-0.181416	2

4.

```

3 std = data.std()
4 #標準化後之結果
5 z_score_normalized = (data - mu) / std
6
7 print(z_score_normalized)

```

	income	age
0	0.550812	1.693570
1	-0.777331	-1.130565
2	-1.110474	1.146963
3	-1.194372	0.873660
4	-0.181416	1.238064
...
995	-0.242008	-0.948363
996	0.381668	0.144851
997	-0.965654	-0.401756
998	-1.608203	-0.766161
999	1.114049	0.418154

[1000 rows x 2 columns]

```

In [165]: 1 data['age_Z']=z_score_normalized['age']
          2 data['income_Z']=z_score_normalized['income']
          3 data.head()

```

```

Out[165]:

```

	income	age	age_Z	income_Z
0	101743	58	1.693570	0.550812
1	49597	27	-1.130565	-0.777331
2	36517	52	1.146963	-1.110474

5、6

```

In [166]: 1 from sklearn.cluster import MeanShift, estimate_bandwidth
          2 bandwidth = estimate_bandwidth(data, quantile=0.1)
          3 ms = MeanShift(bandwidth=bandwidth, bin_seeding=True)
          4 ms.fit(data)
          5 labels = ms.labels_
          6 cluster_centers = ms.cluster_centers_
          7 labels_unique = np.unique(labels)
          8 n_clusters_ = len(labels_unique)
          9
         10 print("number of estimated clusters : %d" % n_clusters_)
         11 import matplotlib.pyplot as plt
         12 from itertools import cycle
         13 print(bandwidth)
         14 plt.figure(1)
         15 plt.clf()
         16 data['ms']=labels
         17 data.head()
         18

```

number of estimated clusters : 8
7365.635129383834

```

Out[166]:

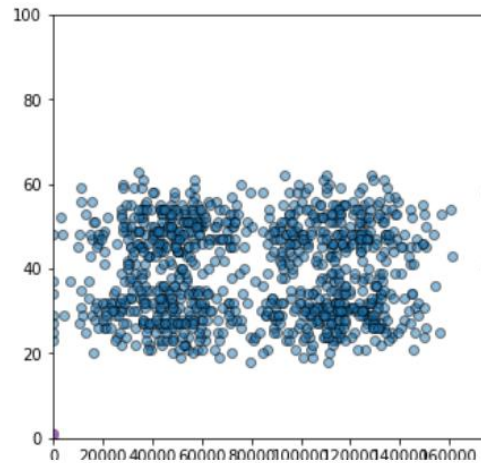
```

	income	age	age_Z	income_Z	ms
0	101743	58	1.693570	0.550812	4
1	49597	27	-1.130565	-0.777331	0
2	36517	52	1.146963	-1.110474	3
3	33223	49	0.873660	-1.194372	3
4	72994	53	1.238064	-0.181416	5

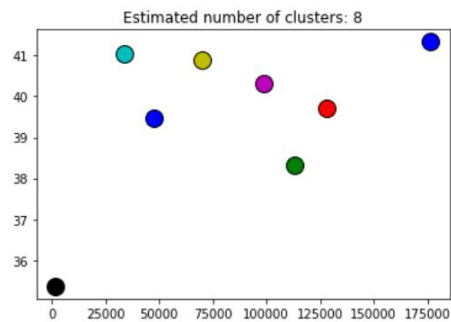
7

```
In [181]: 1 fig=plt.figure(figsize=(5,5))
2
3 plt.scatter(data['income'],data['age'],alpha=0.5,edgecolor='k')
4 for idx, centers in enumerate(centers):
5     plt.scatter(*centers)
6 plt.xlim(0,175000)
7 plt.ylim(0,100)
8 plt.show
```

Out[181]: <function matplotlib.pyplot.show(*args, **kw)>



```
In [100]: 1 colors = cycle('bgrcmkybgrcmkybgrcmkybgrcmky')
2 for k, col in zip(range(n_clusters_), colors):
3     ## 根据labels中的值是否等于k，重新组成一个True、False的数组
4     my_members = labels == k
5     cluster_center = cluster_centers[k]
6     plt.plot(cluster_center[0], cluster_center[1], 'o', markerfacecolor=col,
7              markeredgecolor='k', markersize=14)
8 plt.title('Estimated number of clusters: %d' % n_clusters_)
9 plt.show()
```



8.

```
In [134]: 1 mu = data['age'].mean()
2 #標準差
3 std = data['age'].std()
4 #標準化後之結果
5 z_score_normalized = (data['age'] - mu) / std
6
7 print(z_score_normalized)
```

```
0    -1.446258
1     0.438698
2     0.783290
3     1.105599
4    -1.381068
...
995    1.347424
996   -1.746302
997    0.915534
998    1.060287
999    0.517937
Name: age, Length: 1000, dtype: float64
```

```
In [135]: 1 data['age_Z']=z_score_normalized
2 data.head()
```

```
Out[135]:
```

	age	education	age_Z
0	27.007219	college	-1.446258
1	47.615409	highschool	0.438698
2	51.382815	highschool	0.783290

9

```
3 label = kp.labels_
4
```

```
In [154]: 1 data['label']=label
2 data.head()
3
```

```
Out[154]:
```

	age	education	age_Z	label
0	27.007219	college	-1.446258	0
1	47.615409	highschool	0.438698	2
2	51.382815	highschool	0.783290	1
3	54.906622	highschool	1.105599	1
4	27.719939	less_than_highschool	-1.381068	0

```
In [155]: 1 sector=data.groupby('label')
2 sector.size()
```

```
Out[155]: label
0      316
1      396
2      288
dtype: int64
```

10

```
Out[155]: label
0      316
1      396
2      288
dtype: int64
```

```
In [157]: 1 import pandas as pd
2 import numpy as np
3
4
5 data=pd.read_csv('C:/Users/Admin/Downloads/customer_offers.csv',index_col="customer_name")
6 data.head()
```

```
Out[157]:
```

	1	2	3	4	5	6	7	8	9	10	...	23	24	25	26	27	28	29	30	31	32
customer_name																					
Adams	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	1	1	0	0
Allen	0	0	0	0	0	0	0	0	1	0	...	0	0	0	0	1	0	0	0	0	0
Anderson	0	0	0	0	0	0	0	0	0	0	...	0	1	0	1	0	0	0	0	0	0
Bailey	0	0	0	0	0	0	1	0	0	0	...	0	0	0	0	0	0	0	1	0	0
Baker	0	0	0	0	0	0	1	0	0	1	...	0	0	0	0	0	0	0	0	1	0

5 rows × 32 columns

11

```
In [161]: 1 from kmodes import kmodes
2 km = kmodes.KModes(n_clusters=4, random_state=10)
3 clusters = km.fit_predict(data)
4 data['km'] = clusters
5 data.head()
```

```
Out[161]:
```

	1	2	3	4	5	6	7	8	9	10	...	24	25	26	27	28	29	30	31	32	km	
customer_name																						
Adams	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	1	1	0	0	0	
Allen	0	0	0	0	0	0	0	0	0	1	0	...	0	0	0	1	0	0	0	0	0	
Anderson	0	0	0	0	0	0	0	0	0	0	...	1	0	1	0	0	0	0	0	0	0	
Bailey	0	0	0	0	0	0	1	0	0	0	...	0	0	0	0	0	0	1	0	0	0	
Baker	0	0	0	0	0	0	1	0	0	1	...	0	0	0	0	0	0	0	1	0	0	

5 rows x 33 columns

```
In [162]: 1 sector = data.groupby('km')
2 sector.size()
```

```
Out[162]: km
0      89
1       2
2       6
3       3
dtype: int64
```

12

```
In [ ]: 1 k means :適用於數值資料，須定義K，一種無間監督學習的算法，演算法可以非常快速地完成分群任務，
2           但是如果觀測值具有雜訊 ( Noise ) 或者極端值，其分群結果容易被這些雜訊與極端值影響，適合處理分布集中的大型樣本資料，
3
4 Mean Shift :適用於數值資料，不需要預先知道欲分群的數目，
5 K-prototype:同時有數值或離散可使用，是K-means與K-modes的一種集合形式，透過設定了一個目標函數，透過不斷迭代運算，直到目標函數值不變
6 K-modes:適用於類別型，搜尋群聚中心的方法是以各類別資料的眾數 ( Modes ) 為依據，將原本K-means使用的歐式距離替換成字符間的漢明距離
```

13

Weka Explorer

Preprocess | Classify | **Cluster** | Associate | Select attributes | Visualize

Clusterer: Choose **SimpleKMeans** -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 4 -A "weka core EuclideanDistance -R first-last" -I 500 -num-slots 1

Cluster mode

- ☒ Use training set
- ☐ Supplied test set (Set...)
- ☐ Percentage split % 66
- ☐ Classes to clusters evaluation (Num) age
- ☒ Store clusters for visualization

Ignore attributes

Start Stop

Result list (right-click for options)

- 03:00:50 - SimpleKMeans

Clusterer output

```
Number of iterations: 15
Within cluster sum of squared errors: 20.771054474827608

Initial starting points (random):

Cluster 0: 0,25
Cluster 1: 62517,56
Cluster 2: 45423,27
Cluster 3: 113076,39

Missing values globally replaced with mean/mode

Final cluster centroids:

Attribute      Full Data      Cluster#
              (1000.0)    (246.0)      0      1      2      3
=====
income          80116.818    43409.3984    46482.3711 114934.8308 116199.8445
age              39.41        29.6951      49.3398    29.6808    49.3992

Time taken to build model (full training data) : 0.06 seconds

=== Model and evaluation on training set ===

Clustered Instances
```

