```
In [1]:    1  import pandas as pd
           2  import numpy as np
           3  import matplotlib.pyplot as plt
           4  from sklearn.cluster import DBSCAN
```

<frozen importlib._bootstrap>:219: RuntimeWarning: numpy.ufunc size changed, may indicate binary incomp
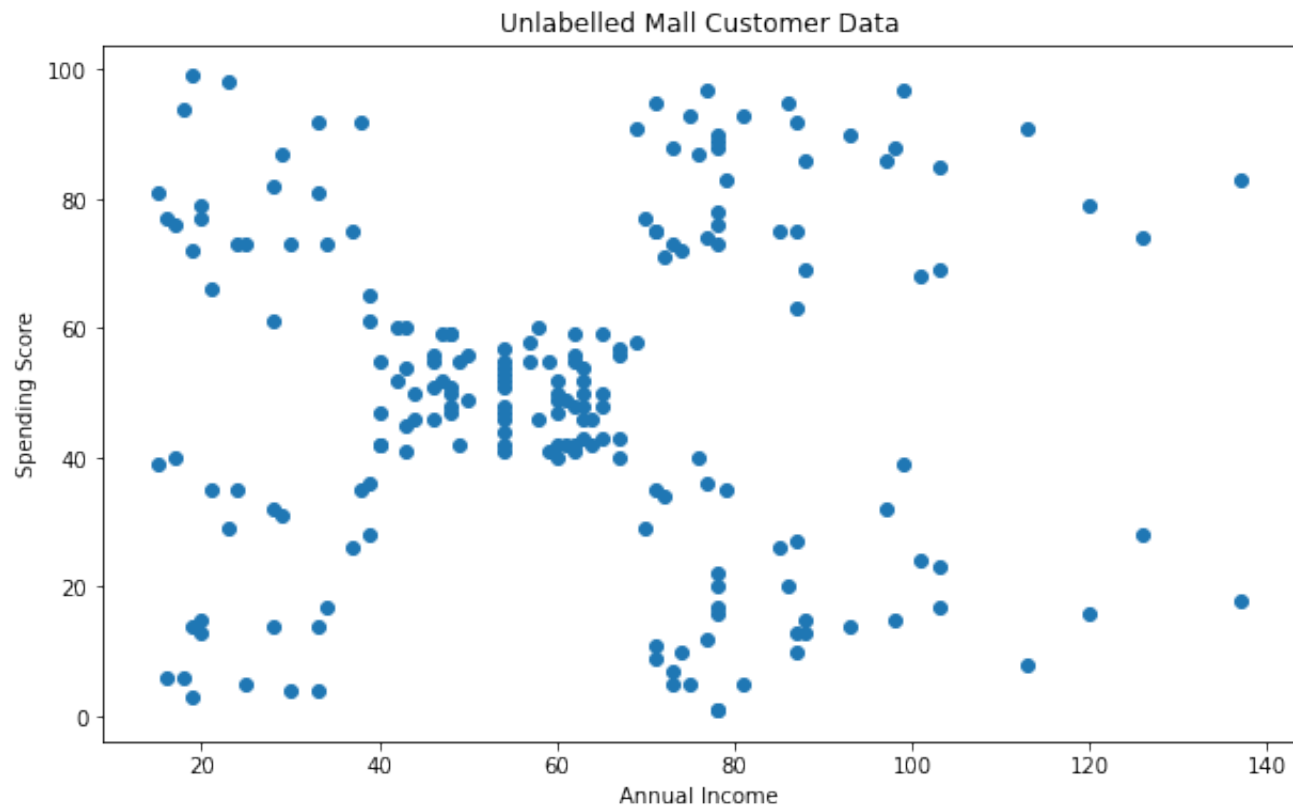atibility. Expected 192 from C header, got 216 from PyObject

```
In [2]:    1  df = pd.read_csv('Mall_Customers.csv')
           2  print("Shape of the data= ", df.shape)
           3  df.head()
```

Shape of the data=  (200, 5)

Out[2]:

| | CustomerID | Gender | 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 |

```
In [3]:  1  plt.figure(figsize=(10,6))
         2  plt.scatter(df['Annual Income (k$)'],df['Spending Score (1-100)'])
         3  plt.xlabel('Annual Income')
         4  plt.ylabel('Spending Score')
         5  plt.title('Unlabelled Mall Customer Data')
         6  plt.show()
```

```python
# Since we are going to use Annual Income and Spending Score  columns only, lets create 2D array of
X = df.iloc[:, [3,4]].values
X[:5] # Show first 5 records only
```

```
Out[4]: array([[15, 39],
               [15, 81],
               [16,  6],
               [16, 77],
               [17, 40]])
```

```python
from itertools import product
eps_values = np.arange(8,13,0.25)
min_samples = np.arange(3,9)
dbscan_params = list(product(eps_values,min_samples)) # gives combination of EPS & MINPOINTS
```

```python
# sigle  value of eps, minpoint
X_numeric_feature = df[['Age','Annual Income (k$)','Spending Score (1-100)']]
dbs_cluster_single = DBSCAN(eps=9,min_samples=8).fit(X_numeric_feature)
dbs_cluster_single.labels_
```

```
Out[12]: array([-1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
               -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
               -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,  1,
               -1, -1,  0,  1,  1,  1,  0,  2,  0,  0,  2,  0,  0,  0,  2,  1,  0,
                2, -1,  0,  1,  0,  0,  0,  2,  1,  1,  2,  1,  0, -1, -1,  1,  2,
                1,  0,  2, -1,  1, -1,  2,  1,  1, -1,  2,  1,  2,  1,  2,  2,  1,
               -1,  2,  1,  2, -1,  1, -1, -1, -1,  2, -1,  2,  2,  2, -1, -1,  1,
               -1,  2, -1, -1, -1, -1,  3, -1, -1, -1,  3, -1,  3, -1,  3, -1,  4,
               -1,  3, -1,  3, -1,  4, -1,  4, -1,  4, -1,  3, -1,  4, -1,  4, -1,
                3, -1,  4, -1,  3, -1,  3, -1,  3, -1,  4, -1,  3, -1,  4, -1, -1,
               -1, -1, -1,  4, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1,
               -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1, -1])
```

```python
np.unique(dbs_cluster_single.labels_)
```

```
Out[13]: array([-1,  0,  1,  2,  3,  4])
```

```
In [14]:    1  len(np.unique(dbs_cluster_single.labels_))
```

Out[14]: 6

```
In [16]:    1  from sklearn.metrics import silhouette_score
            2  silhouette_score(X_numeric_feature,dbs_cluster_single.labels_)
```
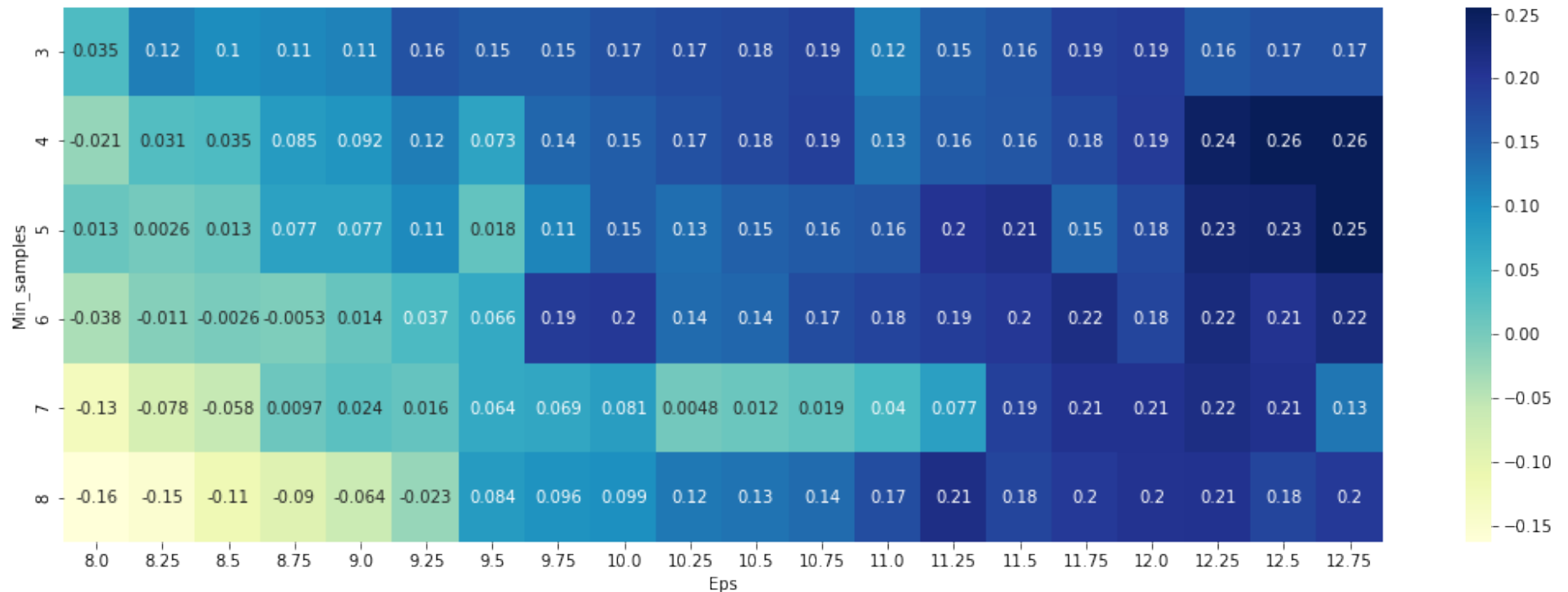
Out[16]: -0.06361003273116211

```
In [21]:    1  #dbscan_params
```

```
In [18]:    1  from sklearn.metrics import silhouette_score
            2  X_numeric_feature = df[['Age','Annual Income (k$)','Spending Score (1-100)']]
            3  sil_score = []
            4  no_of_clusters = []
            5
            6  for k in dbscan_params:
            7      dbs_cluster = DBSCAN(eps=k[0],min_samples=k[1]).fit(X_numeric_feature)
            8      no_of_clusters.append(len(np.unique(dbs_cluster.labels_)))
            9      sil_score.append(silhouette_score(X_numeric_feature,dbs_cluster.labels_))
```

```
In [22]:    1  #sil_score
```

```
In [23]:    1  #no_of_clusters
```

```
In [28]:  1  import seaborn as sns
          2  temp = pd.DataFrame.from_records(dbscan_params,columns = ['Eps','Min_samples'])
          3  temp['sil_score'] = sil_score
          4
          5  plt.figure(figsize=(18,6))
          6  pivot = pd.pivot_table(temp,values = 'sil_score',index = 'Min_samples',columns = 'Eps')
          7  sns.heatmap(pivot,annot = True, cmap = 'YlGnBu')
          8  plt.show()
```



```
In [29]:  1  # global maxima is 0.26 for eps = 12.75 and min_sample = 4
```

```
In [30]:  1  dbs_cluster_final = DBSCAN(eps=12.75,min_samples=4).fit(X_numeric_feature)
```

In [31]:
```python
dbs_clustered = X_numeric_feature.copy()
dbs_clustered.loc[:,'cluster'] = dbs_cluster_final.labels_
dbs_clustered
```

Out[31]:

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

200 rows × 4 columns

```
In [32]:  1  dbs_cluster_size = dbs_clustered.groupby('cluster').size().to_frame()
          2  dbs_cluster_size.columns = ['dbs_size']
          3  dbs_cluster_size
```

Out[32]:

|         | dbs_size |
|---------|----------|
| cluster |          |
| -1      | 17       |
| 0       | 113      |
| 1       | 8        |
| 2       | 34       |
| 3       | 24       |
| 4       | 4        |

```
In [33]:  1  outliers = dbs_clustered[dbs_clustered['cluster']==-1]
```
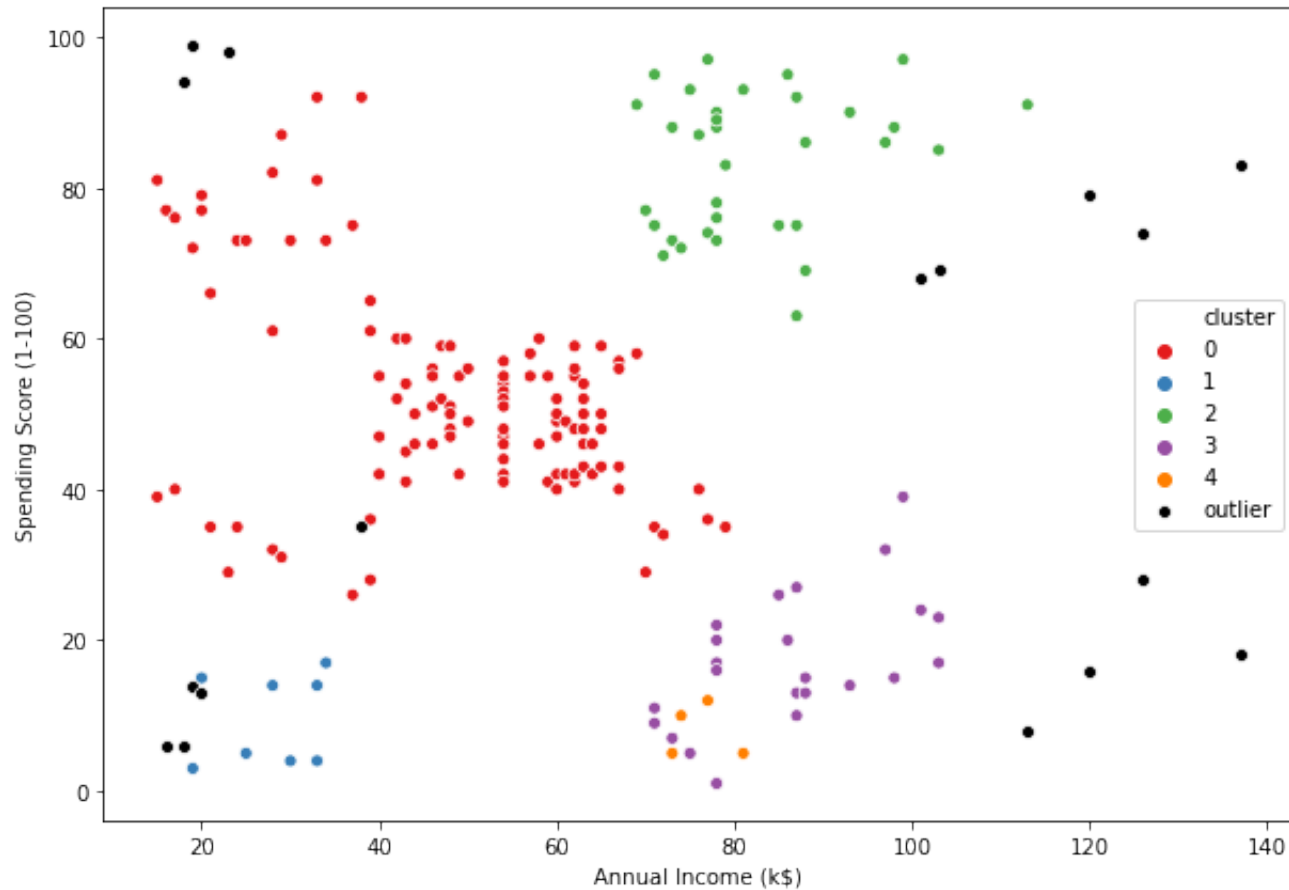
```
In [42]:   1  outliers
```

Out[42]:

|     | Age | Annual Income (k$) | Spending Score (1-100) | cluster |
| --- | --- | --- | --- | --- |
| 2   | 20  | 16  | 6   | -1  |
| 6   | 35  | 18  | 6   | -1  |
| 7   | 23  | 18  | 94  | -1  |
| 10  | 67  | 19  | 14  | -1  |
| 11  | 35  | 19  | 99  | -1  |
| 14  | 37  | 20  | 13  | -1  |
| 19  | 35  | 23  | 98  | -1  |
| 40  | 65  | 38  | 35  | -1  |
| 187 | 28  | 101 | 68  | -1  |
| 191 | 32  | 103 | 69  | -1  |
| 192 | 33  | 113 | 8   | -1  |
| 194 | 47  | 120 | 16  | -1  |
| 195 | 35  | 120 | 79  | -1  |
| 196 | 45  | 126 | 28  | -1  |
| 197 | 32  | 126 | 74  | -1  |
| 198 | 32  | 137 | 18  | -1  |
| 199 | 30  | 137 | 83  | -1  |

```
In [53]:   1  plt.figure(figsize=(10,7))
           2  sns.scatterplot('Annual Income (k$)','Spending Score (1-100)',
           3                  data =dbs_clustered[dbs_clustered['cluster']!=-1],
           4              hue = 'cluster',palette = 'Set1')
           5  sns.scatterplot('Annual Income (k$)','Spending Score (1-100)',
           6                  data =dbs_clustered[dbs_clustered['cluster']==-1], color = 'black',label = 'outlier'
           7
           8  plt.show()
```



```
In [ ]:    1
```

In [ ]: