**Tran Quoc Long - 14520490**

# Bài tập 2: Handwritten digits - clustering

```
1  Thực hiện các phép cluster trên bộ dữ liệu handwritten digits
2  Nội dung bao gồm trong file:
3  1. Chạy thử các hàm cluster và các hàm liên quan
4      - Kmeans
5      - Spectral clustering
6      - DBSCAN
7      - Agglomerative clustering
8      - Cross table
9      - Figure to visualize result
10     - Show centroid of Kmeans
11  2. Nội dung thực hành 2
```

### Chạy thử các hàm cluster và các hàm liên quan

### K-means
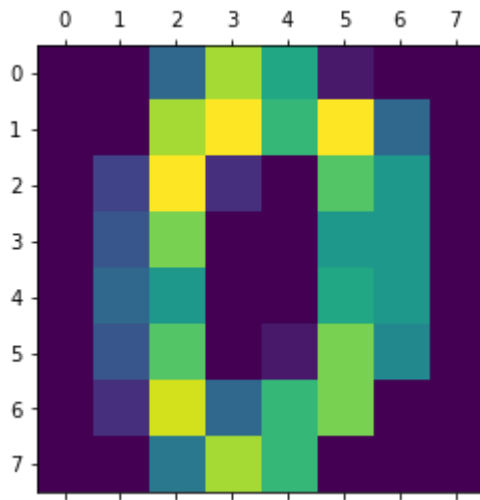
```
In [1]:  1  #import libs
         2  from time import time
         3  import numpy as np
         4  import matplotlib.pyplot as plt
         5  import pandas as pd
```

```
In [2]:  1  #import scikit-learn
         2  from sklearn import metrics
         3  from sklearn.cluster import KMeans
         4  from sklearn.datasets import load_digits
```

```
In [3]:  1  digits = load_digits();
         2  print(digits.data.shape);
```

```
(1797, 64)
```

```
In [4]:    1 %matplotlib inline
           2 #plt.gray();
           3 plt.matshow(digits.images[0]);
```



```
In [5]:    1 nClusters = 10
           2 model1 = KMeans(nClusters)
           3 labels_kmeans = model1.fit_predict(digits.data)
```
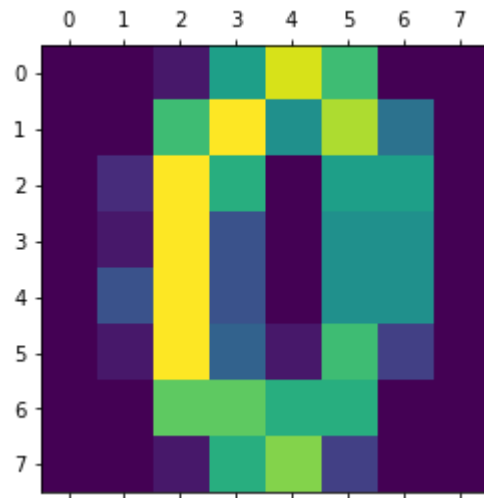
```
In [6]:    1 df = pd.DataFrame({'labels':labels_kmeans,'Truth labels':digits.target})
           2 ct = pd.crosstab(df['labels'],df['Truth labels'])
           3 print(ct)
```

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels | | | | | | | | | | |
| 0 | 0 | 0 | 3 | 7 | 9 | 0 | 0 | 175 | 5 | 7 |
| 1 | 0 | 1 | 0 | 2 | 0 | 136 | 0 | 0 | 4 | 6 |
| 2 | 177 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 2 | 9 | 0 | 42 | 0 | 0 | 48 | 138 |
| 4 | 0 | 24 | 148 | 1 | 0 | 0 | 0 | 0 | 3 | 0 |
| 5 | 0 | 2 | 0 | 0 | 0 | 1 | 177 | 0 | 2 | 0 |
| 6 | 0 | 1 | 13 | 157 | 0 | 1 | 0 | 0 | 4 | 7 |
| 7 | 1 | 0 | 0 | 0 | 164 | 2 | 0 | 0 | 0 | 0 |
| 8 | 0 | 99 | 8 | 7 | 3 | 0 | 2 | 2 | 102 | 2 |
| 9 | 0 | 55 | 2 | 0 | 5 | 0 | 1 | 2 | 6 | 20 |

In [7]:
```
1 n = 10
2 %matplotlib inline
3 plt.matshow(digits.images[n])
4 print('Predict Label:', labels_kmeans[n])
5
6 print('Truth: ', digits.target[n])
```
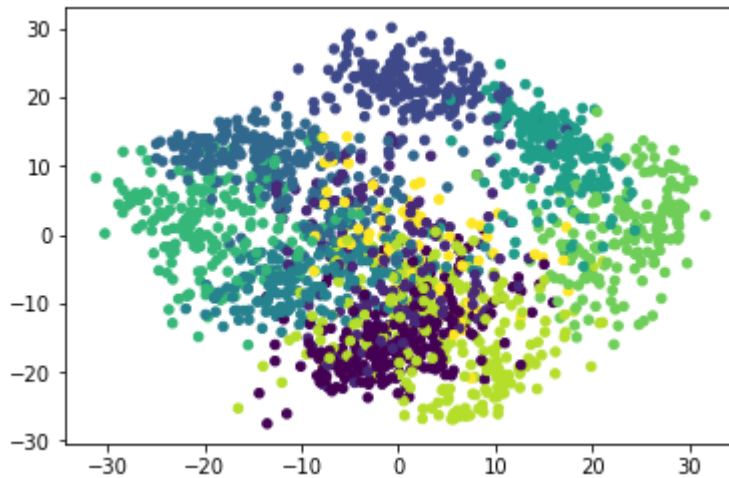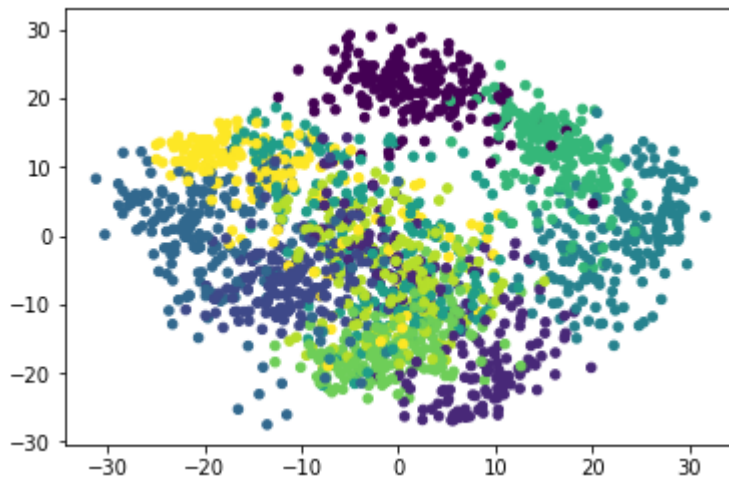
Predict Label: 2
Truth:   0



## Visualization - Kmeans

In [8]:
```
1 #import libs
2 import numpy as np
3 from sklearn.decomposition import PCA
```

*PCA*

```
1 nComponents = 2
2 vPCA = PCA(nComponents)
3 digitData_to_2D = vPCA.fit_transform(digits.data)
4 plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_kmeans, s=
5 plt.show()
```

```
1 plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= digits.target, s=
2 plt.show()
```



## Speactral clustering

In [11]:
```python
# Spectral_clustering

from sklearn.cluster import spectral_clustering
from sklearn.feature_extraction import image
import numpy as np
from sklearn.neighbors import DistanceMetric
from sklearn.metrics.pairwise import cosine_similarity

# dist = DistanceMetric.get_metric('euclidean')
# graph=dist.pairwise(digits.data)

graph = cosine_similarity(digits.data)
label_spectral = spectral_clustering(graph, n_clusters=10)
```
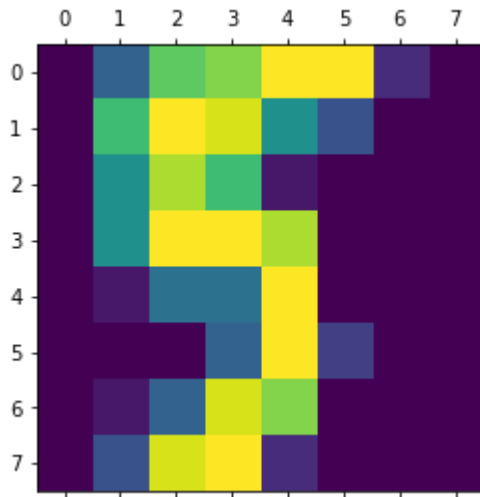
## Cross table biểu thị kết quả và so sánh

In [12]:
```python
df1 = pd.DataFrame({'labels':label_spectral,'Truth labels':digits.target})
ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
print(ct2)
```

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels | | | | | | | | | | |
| 0 | 1 | 0 | 0 | 0 | 163 | 2 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 147 | 0 | 0 | 0 | 0 | 6 | 2 |
| 2 | 177 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 3 |
| 3 | 0 | 0 | 0 | 4 | 0 | 157 | 0 | 0 | 3 | 3 |
| 4 | 0 | 36 | 115 | 4 | 0 | 0 | 0 | 0 | 1 | 0 |
| 5 | 0 | 0 | 0 | 16 | 0 | 20 | 2 | 0 | 7 | 133 |
| 6 | 0 | 0 | 2 | 2 | 11 | 0 | 0 | 154 | 3 | 2 |
| 7 | 0 | 2 | 0 | 0 | 0 | 2 | 172 | 0 | 13 | 0 |
| 8 | 0 | 86 | 53 | 5 | 5 | 0 | 7 | 10 | 101 | 1 |
| 9 | 0 | 58 | 5 | 5 | 1 | 0 | 0 | 15 | 40 | 36 |

In [13]:
```
1 n = 15
2 plt.matshow(digits.images[n])
3 print('lables_predict:',label_spectral[n])
4 print(' True: ', digits.target[n])
```
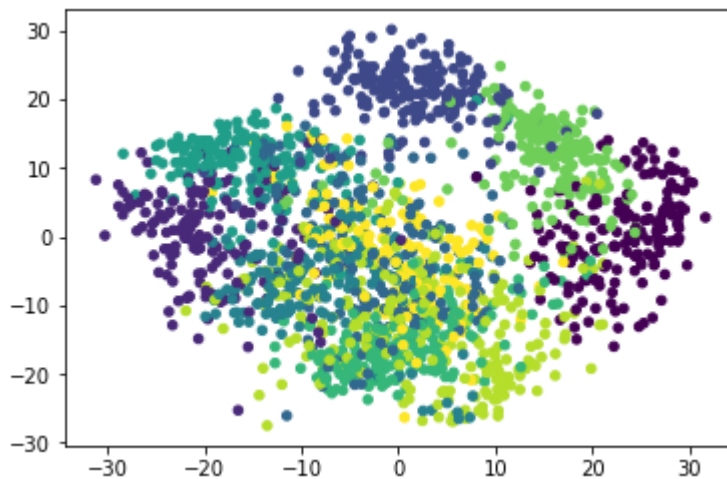
lables_predict: 3
 True:   5



**Visualization - Spectral Clustering**

In [14]:
```
1 plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= label_spectral, s
2 plt.show()
```
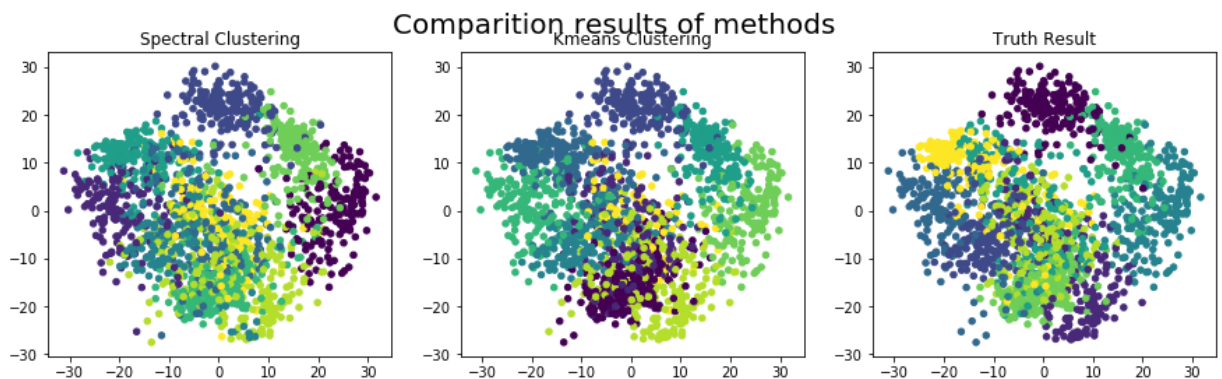


## Visualize results to compare - Using PCA

```
In [15]:     1  fig = plt.figure(figsize=(15,4))
             2  fig.suptitle('Comparition results of methods', fontsize=20)
             3
             4  ax = fig.add_subplot(1,3,1)
             5  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= label_spectral, s
             6  ax.set_title('Spectral Clustering')
             7
             8  ax = fig.add_subplot(1,3,2)
             9  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_kmeans, s=
            10  ax.set_title('Kmeans Clustering')
            11
            12  ax = fig.add_subplot(1,3,3)
            13  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= digits.target, s=
            14  ax.set_title('Truth Result')
```

Out[15]:  <matplotlib.text.Text at 0x16c59124940>



## DBSCAN

```
In [16]:     1  import numpy as np
             2  import matplotlib.pyplot as plt
             3  from sklearn.cluster import DBSCAN
             4  from sklearn import metrics
             5  from sklearn.datasets.samples_generator import make_blobs
             6  from sklearn.preprocessing import StandardScaler
             7
             8  #import scikit-learn
             9  from sklearn import metrics
            10  from sklearn.cluster import KMeans
            11  from sklearn.datasets import load_digits
            12  from sklearn.preprocessing import scale
            13  from sklearn.decomposition import PCA
            14
```

```
1  digits = load_digits()
2  data = digits.data
3  data = StandardScaler().fit_transform(data)
4
5  n_samples, n_features = data.shape
6  n_digits = len(np.unique(digits.target))
7  labels = digits.target
8
9  sample_size = 300
10
11 print("n_digits: %d, \t n_samples %d, \t n_features %d"
12     % (n_digits, n_samples, n_features))
13
```

n_digits: 10,      n_samples 1797,      n_features 64

```
1  print(data)
```

```
[[ 0.         -0.33501649 -0.04308102 ..., -1.14664746 -0.5056698
  -0.19600752]
 [ 0.         -0.33501649 -1.09493684 ...,  0.54856067 -0.5056698
  -0.19600752]
 [ 0.         -0.33501649 -1.09493684 ...,  1.56568555  1.6951369
  -0.19600752]
 ...,
 [ 0.         -0.33501649 -0.88456568 ..., -0.12952258 -0.5056698
  -0.19600752]
 [ 0.         -0.33501649 -0.67419451 ...,  0.8876023  -0.5056698
  -0.19600752]
 [ 0.         -0.33501649  1.00877481 ...,  0.8876023  -0.26113572
  -0.19600752]]
```

```
1  db = DBSCAN(eps=1, min_samples=1,algorithm='kd_tree').fit(data)
```

- với min_samples = 1 => số cluster ra quá lớn = số input ban đầu

```
1  print(db)
```

```
DBSCAN(algorithm='kd_tree', eps=1, leaf_size=30, metric='euclidean',
    metric_params=None, min_samples=1, n_jobs=1, p=None)
```

```python
In [22]:  1 import pandas as pd
          2 df1 = pd.DataFrame({'labels':db.labels_,'Truth labels':digits.target})
          3 ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
          4 print(ct2)
```

```
Truth labels  0  1  2  3  4  5  6  7  8  9
labels
0             1  0  0  0  0  0  0  0  0  0
1             0  1  0  0  0  0  0  0  0  0
2             0  0  1  0  0  0  0  0  0  0
3             0  0  0  1  0  0  0  0  0  0
4             0  0  0  0  1  0  0  0  0  0
5             0  0  0  0  0  1  0  0  0  0
6             0  0  0  0  0  0  1  0  0  0
7             0  0  0  0  0  0  0  1  0  0
8             0  0  0  0  0  0  0  0  1  0
9             0  0  0  0  0  0  0  0  0  1
10            1  0  0  0  0  0  0  0  0  0
11            0  1  0  0  0  0  0  0  0  0
12            0  0  1  0  0  0  0  0  0  0
13            0  0  0  1  0  0  0  0  0  0
14            0  0  0  0  1  0  0  0  0  0
15            0  0  0  0  0  1  0  0  0  0
16            0  0  0  0  0  0  1  0  0  0
17            0  0  0  0  0  0  0  1  0  0
18            0  0  0  0  0  0  0  0  1  0
19            0  0  0  0  0  0  0  0  0  1
20            1  0  0  0  0  0  0  0  0  0
21            0  1  0  0  0  0  0  0  0  0
22            0  0  1  0  0  0  0  0  0  0
23            0  0  0  1  0  0  0  0  0  0
24            0  0  0  0  1  0  0  0  0  0
25            0  0  0  0  0  1  0  0  0  0
26            0  0  0  0  0  0  1  0  0  0
27            0  0  0  0  0  0  0  1  0  0
28            0  0  0  0  0  0  0  0  1  0
29            0  0  0  0  0  0  0  0  0  1
...          .. .. .. .. .. .. .. .. .. ..
1766          0  0  0  0  1  0  0  0  0  0
1767          1  0  0  0  0  0  0  0  0  0
1768          0  0  0  0  0  1  0  0  0  0
1769          0  0  0  1  0  0  0  0  0  0
1770          0  0  0  0  0  0  1  0  0  0
1771          0  0  0  0  0  0  0  0  0  1
1772          0  0  0  0  0  0  1  0  0  0
1773          0  1  0  0  0  0  0  0  0  0
1774          0  0  0  0  0  0  0  1  0  0
1775          0  0  0  0  0  1  0  0  0  0
1776          0  0  0  0  1  0  0  0  0  0
1777          0  0  0  0  1  0  0  0  0  0
1778          0  0  0  0  0  0  0  1  0  0
1779          0  0  1  0  0  0  0  0  0  0
1780          0  0  0  0  0  0  0  0  1  0
1781          0  0  1  0  0  0  0  0  0  0
1782          0  0  1  0  0  0  0  0  0  0
1783          0  0  0  0  0  1  0  0  0  0
```

```
1784        0 0 0 0 0 0 0 1 0 0
1785        0 0 0 0 0 0 0 0 0 1
1786        0 0 0 0 0 1 0 0 0 0
1787        0 0 0 0 1 0 0 0 0 0
1788        0 0 0 0 0 0 0 0 1 0
1789        0 0 0 0 0 0 0 0 1 0
1790        0 0 0 0 1 0 0 0 0 0
1791        0 0 0 0 0 0 0 0 0 1
1792        1 0 0 0 0 0 0 0 0 0
1793        0 0 0 0 0 0 0 0 1 0
1794        0 0 0 0 0 0 0 0 0 1
1795        0 0 0 0 0 0 0 0 1 0

[1796 rows x 10 columns]
```
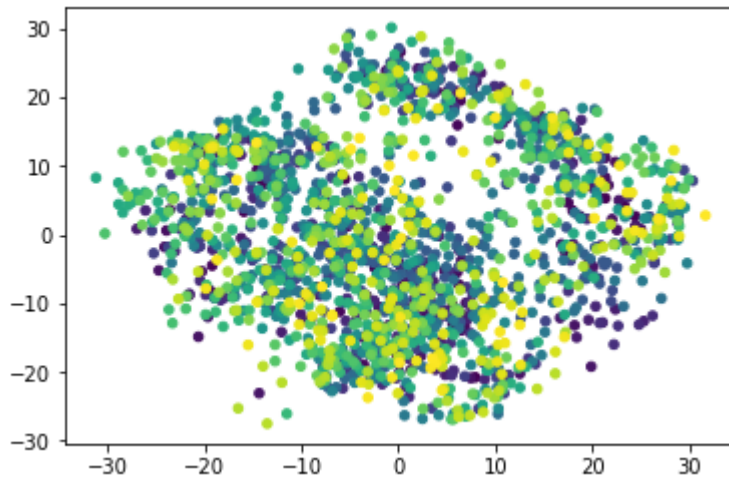
In [23]:
```python
1 nComponents = 2
2 vPCA = PCA(nComponents)
3 digitData_to_2D = vPCA.fit_transform(digits.data)
4 plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= db.labels_, s=20)
5 plt.show()
```



## Agglomerative Clustering

In [24]:
```python
1 from sklearn.cluster import AgglomerativeClustering
```

In [25]:
```python
1 Agglomerative_model = AgglomerativeClustering(n_clusters = 10)
```

In [26]:
```python
1 db = Agglomerative_model.fit(data)
```

In [27]:
```python
1 print(db.labels_)
```
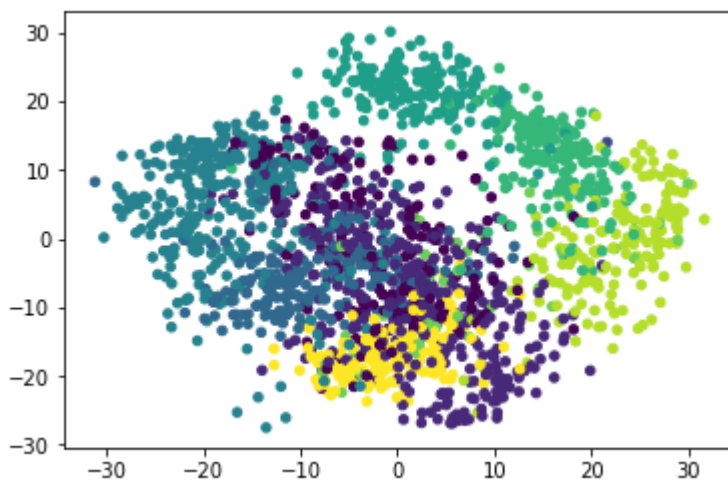
```
[5 1 1 ..., 1 1 1]
```

```
In [28]:   1  import pandas as pd
           2  df1 = pd.DataFrame({'labels':db.labels_,'Truth labels':digits.target})
           3  ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
           4  print(ct2)
```

```
Truth labels    0    1    2    3    4    5    6    7    8    9
labels
0               0    1    0    0    1  168    0    1    1    3
1               0  150   15   11    4    0    1    1  168   38
2               0    0    1    0    1    0    0    0    0    0
3               0   27  160    4    0    1    0    0    3    0
4               0    0    1  168    0   12    0    1    2  135
5             178    0    0    0    0    0    0    0    0    0
6               0    0    0    0    0    1  180    0    0    1
7               0    0    0    0   12    0    0   25    0    3
8               0    4    0    0  163    0    0    0    0    0
9               0    0    0    0    0    0    0  151    0    0
```

```
In [29]:   1  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= db.labels_, s=20)
           2  plt.show()
```



# Bài thực hành 2

# Comparison of cluster methods
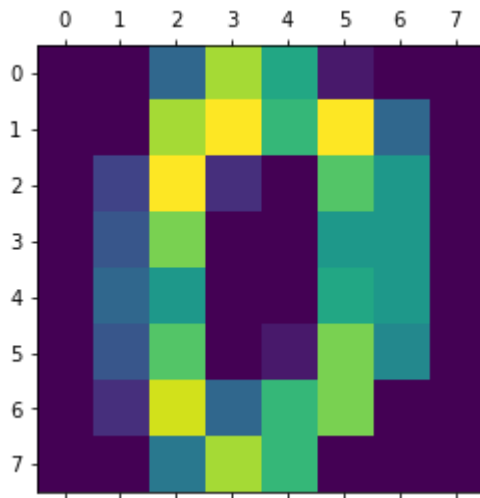
```
In [1]:   1  #import libs
          2  from time import time
          3  import numpy as np
          4  import matplotlib.pyplot as plt
          5  import pandas as pd
```

```
In [2]:   1 #import scikit-learn
          2 from sklearn import metrics
          3 from sklearn.cluster import KMeans, spectral_clustering, DBSCAN, Agglomerativ
          4 from sklearn.datasets import load_digits
          5 from sklearn.neighbors import DistanceMetric
          6 from sklearn.metrics.pairwise import cosine_similarity
          7 from sklearn.preprocessing import StandardScaler
```

```
In [3]:   1 digits = load_digits();
          2 print(digits.data.shape);
```

(1797, 64)

```
In [4]:   1 %matplotlib inline
          2 #plt.gray();
          3 plt.matshow(digits.images[0]);
```



## Clustering

```python
#Kmeans
nClusters = 10
t0 = time()
kmeans_model = KMeans(nClusters)
labels_kmeans = kmeans_model.fit_predict(digits.data)
t_kmeans = time()- t0
#Cross table
print('Kmeans:\n')
df1 = pd.DataFrame({'labels':labels_kmeans,'Truth labels':digits.target})
ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
print(ct2)
```

Kmeans:

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels | | | | | | | | | | |
| 0 | 177 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 24 | 148 | 1 | 0 | 0 | 0 | 0 | 3 | 0 |
| 2 | 0 | 55 | 2 | 0 | 7 | 0 | 1 | 2 | 6 | 20 |
| 3 | 0 | 0 | 3 | 7 | 7 | 0 | 0 | 175 | 4 | 7 |
| 4 | 0 | 1 | 0 | 2 | 0 | 136 | 0 | 0 | 4 | 6 |
| 5 | 0 | 0 | 2 | 9 | 0 | 42 | 0 | 0 | 49 | 139 |
| 6 | 0 | 1 | 13 | 157 | 0 | 1 | 0 | 0 | 3 | 6 |
| 7 | 1 | 0 | 0 | 0 | 163 | 2 | 0 | 0 | 0 | 0 |
| 8 | 0 | 2 | 0 | 0 | 0 | 1 | 177 | 0 | 2 | 0 |
| 9 | 0 | 99 | 8 | 7 | 4 | 0 | 2 | 2 | 103 | 2 |

```python
#Spectral_clustering
t0 = time()
graph = cosine_similarity(digits.data)
labels_spectral = spectral_clustering(graph, n_clusters=10)
t_spectral = time()- t0
#Cross table
print('Spectral clustering:\n')
df1 = pd.DataFrame({'labels':labels_spectral,'Truth labels':digits.target})
ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
print(ct2)

```

Spectral clustering:

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels | | | | | | | | | | |
| 0 | 0 | 58 | 5 | 5 | 1 | 0 | 0 | 15 | 40 | 36 |
| 1 | 0 | 0 | 0 | 16 | 0 | 20 | 2 | 0 | 7 | 133 |
| 2 | 0 | 2 | 0 | 1 | 0 | 2 | 172 | 0 | 13 | 0 |
| 3 | 177 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 3 |
| 4 | 0 | 36 | 115 | 4 | 0 | 0 | 0 | 0 | 1 | 0 |
| 5 | 0 | 0 | 2 | 2 | 11 | 0 | 0 | 154 | 3 | 2 |
| 6 | 1 | 0 | 0 | 0 | 163 | 2 | 0 | 0 | 0 | 0 |
| 7 | 0 | 86 | 53 | 5 | 5 | 0 | 7 | 10 | 101 | 1 |
| 8 | 0 | 0 | 0 | 4 | 0 | 157 | 0 | 0 | 3 | 3 |
| 9 | 0 | 0 | 1 | 146 | 0 | 0 | 0 | 0 | 6 | 2 |

```
In [7]:    1  #DBSCAN
           2  data = digits.data
           3  t0 = time()
           4  # data = StandardScaler().fit_transform(data)
           5  labels_dbscan = DBSCAN(eps=0.06, min_samples=10,metric = 'cosine').fit_predic
           6  t_dbscan = time()- t0
           7  #DBSCAN coss table
           8  print('DBSCAN:\n')
           9  df1 = pd.DataFrame({'labels':labels_dbscan,'Truth labels':digits.target})
          10  ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
          11  print(ct2)
```

DBSCAN:

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels |  |  |  |  |  |  |  |  |  |  |
| -1 | 7 | 13 | 41 | 48 | 34 | 65 | 5 | 48 | 77 | 77 |
| 0 | 171 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 143 | 0 | 0 | 0 | 0 | 1 | 0 | 97 | 1 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 175 | 0 | 0 | 0 |
| 3 | 0 | 0 | 0 | 135 | 0 | 1 | 0 | 0 | 0 | 102 |
| 4 | 0 | 0 | 0 | 0 | 147 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 136 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 116 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 131 | 0 | 0 |
| 8 | 0 | 26 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

```
In [8]:    1  #Agglomerative Clustering
           2  t0 = time()
           3  Agglomerative_model = AgglomerativeClustering(n_clusters = nClusters)
           4  labels_AgglomerativeClustering = Agglomerative_model.fit_predict(data)
           5  t_agg = time() - t0
           6  #Agglomerative Clustering - Crosstable
           7  print('Agglomerative Clustering:\n')
           8  df1 = pd.DataFrame({'labels':labels_AgglomerativeClustering,'Truth labels':di
           9  ct2 = pd.crosstab(df1['labels'],df1['Truth labels'])
          10  print(ct2)
```

Agglomerative Clustering:

| Truth labels | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| labels |  |  |  |  |  |  |  |  |  |  |
| 0 | 0 | 0 | 0 | 0 | 0 | 179 | 0 | 0 | 0 | 2 |
| 1 | 0 | 0 | 0 | 169 | 0 | 2 | 0 | 0 | 1 | 145 |
| 2 | 0 | 27 | 166 | 0 | 0 | 0 | 0 | 0 | 4 | 0 |
| 3 | 0 | 0 | 1 | 1 | 3 | 0 | 0 | 179 | 1 | 11 |
| 4 | 0 | 0 | 10 | 13 | 0 | 0 | 1 | 0 | 165 | 2 |
| 5 | 0 | 0 | 0 | 0 | 178 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 0 | 0 | 1 | 180 | 0 | 0 | 0 |
| 7 | 178 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 59 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 20 |
| 9 | 0 | 96 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |

## Comparison

In [9]:
```python
n_samples, n_features = data.shape
n_digits = len(np.unique(digits.target))
sample_size = 300
#print frame
print("n_digits: %d, \t n_samples %d, \t n_features %d"
      % (n_digits, n_samples, n_features))


print(82 * '_')
print('init\t\ttime\thomo\tcompl\tv-meas\tARI\tAMI\tsilhouette')


data = digits.data
#define a function to measure and print out
def bench_clustering(method_name, time_, labels):
    print('%-9s\t%.2fs\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f\t%.3f'
          % (method_name, time_,
             metrics.homogeneity_score(digits.target, labels),
             metrics.completeness_score(digits.target, labels),
             metrics.v_measure_score(digits.target, labels),
             metrics.adjusted_rand_score(digits.target, labels),
             metrics.adjusted_mutual_info_score(digits.target,  labels),
             metrics.silhouette_score(data, labels,
                                      metric='euclidean',
                                      sample_size=sample_size)))


#Kmeans
bench_clustering('K-means', t_kmeans, labels_kmeans)
#Spectral_clustering
bench_clustering('spectral', t_spectral, labels_spectral)

#Agglomerative clustering
bench_clustering('Agg.', t_agg, labels_AgglomerativeClustering)
#DBSCAN
bench_clustering('DBSCAN', t_dbscan, labels_dbscan)
```

```
n_digits: 10, 	 n_samples 1797, 		 n_features 64

_____
___
init            time    homo    compl   v-meas  ARI     AMI     silhouette
K-means         0.17s   0.741   0.749   0.745   0.672   0.738   0.165
spectral        0.45s   0.711   0.716   0.713   0.624   0.708   0.170
Agg.            0.13s   0.858   0.879   0.868   0.794   0.856   0.195
DBSCAN          0.04s   0.709   0.757   0.732   0.520   0.706   0.130
```

## Nhận xét

- Dựa trên bảng so sánh trên, thông qua các độ đo, ta có thể thấy phương pháp Agglomerative clustering cho kết quả có độ chính xác cao nhất, với tốc độ nhanh nhất trong các phương pháp đã sử dụng
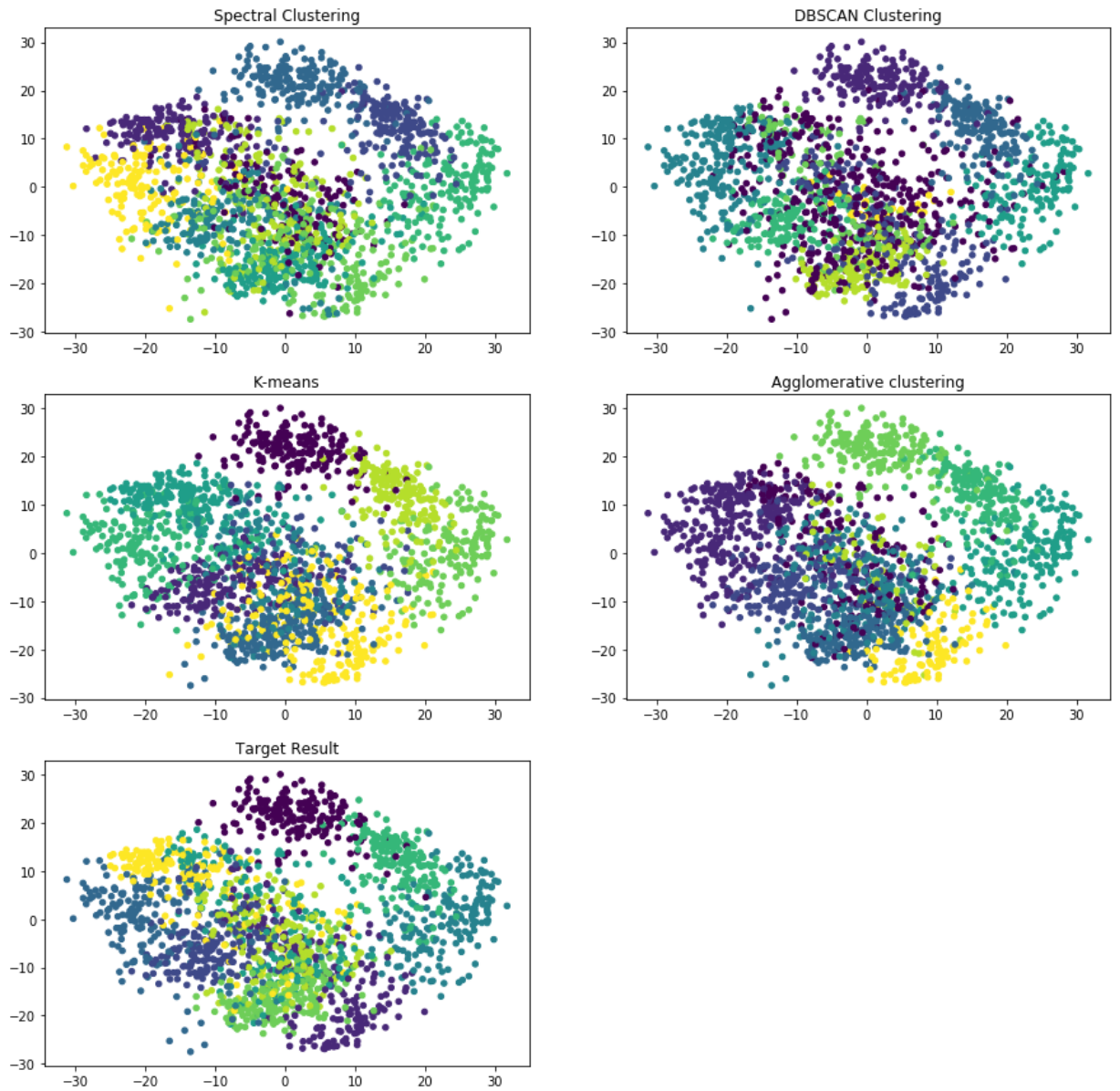
In [ ]:
```python

```

**Visualization**

```python
In [10]:   1  from sklearn.decomposition import PCA
           2
           3  nComponents = 2
           4  vPCA = PCA(nComponents)
           5  digitData_to_2D = vPCA.fit_transform(digits.data)
           6
           7  fig = plt.figure(figsize=(15,15))
           8  fig.suptitle('Comparition results of methods', fontsize=20)
           9
          10  ax = fig.add_subplot(3,2,1)
          11  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_spectral,
          12  ax.set_title('Spectral Clustering')
          13
          14  ax = fig.add_subplot(3,2,2)
          15  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_dbscan, s=
          16  ax.set_title('DBSCAN Clustering')
          17
          18  ax = fig.add_subplot(3,2,3)
          19  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_kmeans, s=
          20  ax.set_title('K-means')
          21
          22  ax = fig.add_subplot(3,2,4)
          23  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= labels_Agglomerat
          24  ax.set_title('Agglomerative clustering')
          25
          26  ax = fig.add_subplot(3,2,5)
          27  plt.scatter(digitData_to_2D[:,0], digitData_to_2D[:,1],  c= digits.target, s=
          28  ax.set_title('Target Result')
          29
```

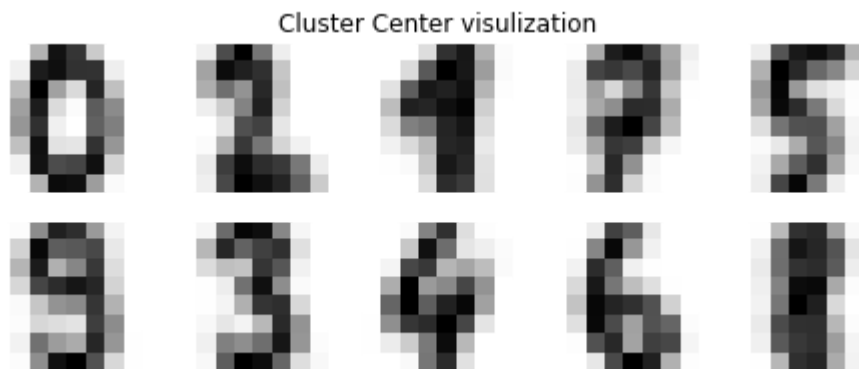Out[10]:  <matplotlib.text.Text at 0x21792954828>

Comparition results of methods

**Show centroids**

In [12]:
```python
# create a fig to show image
fig = plt.figure(figsize=(8,3))
plt.title('Cluster Center visulization')
plt.axis('off')
# for all 0-9 Labels
for i in range(10):
    # initialize subplots in a grid 2x5 at i+1th position
    ax = fig.add_subplot(2, 5, 1+i)

    # display image
    ax.imshow(kmeans_model.cluster_centers_[i].reshape((8,8)), cmap=plt.cm.bi
    #don't show the axes
    plt.axis('off')

plt.show()
```

Cluster Center visulization



- ( Các centroid có biểu thị hình ảnh là các chữ số )

In [ ]:
1