

Machine Learning Homework 11

Hairui Yin

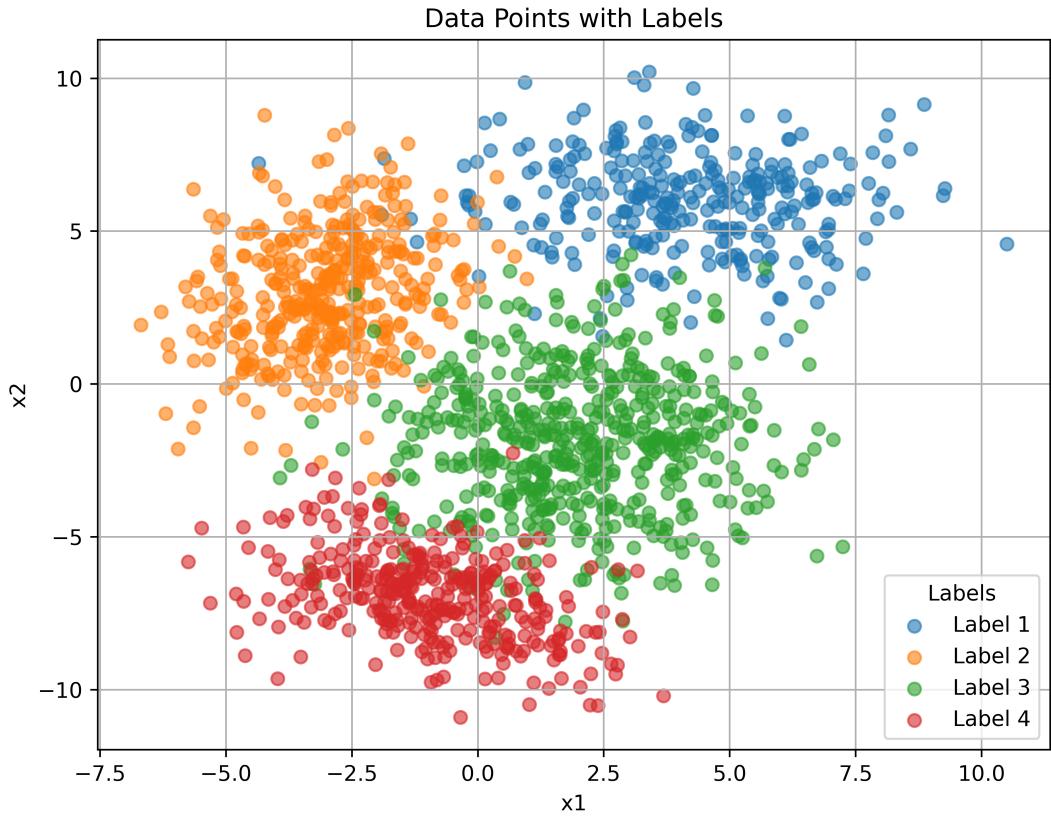
For the whole problem, code is uploaded in file code.ipynb.

1.

(a) We load the data and plot data points using different colors for each label with the following code:

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 df = pd.read_csv('HW11-ClusteringData.csv', header=None, names=['x1', 'x2', 'label'])
6 plt.figure(figsize=(8, 6))
7 for label in df['label'].unique():
8     subset = df[df['label'] == label]
9     plt.scatter(subset['x1'], subset['x2'], label=f"Label {label}", alpha=0.6)
10 plt.title("Data Points with Labels")
11 plt.xlabel("x1")
12 plt.ylabel("x2")
13 plt.legend(title="Labels")
14 plt.grid(True)
15 # plt.show()
16 plt.savefig('q1a.png', dpi=600)
```

The result is shown in the figure:



(b) In this problem, we set $k \in \{2, 3, \dots, 7\}$ and perform KMeans to the data using two distance measurement, Euclidean Distance and Manhattan Distance. And then plot average silhouette coefficient of the data points. After experiment, the optimal value of k is $k = 4$, the same as number of data labels. The code is followed the above one. For Euclidean Distance,

```

1 def EuclideanDistance(v, axis):
2     return np.linalg.norm(v, axis=axis)
3 def ManhattanDistance(v, axis):
4     return np.sum(np.abs(v), axis=axis)
5 def AvgSilhouetteScore(data, labels, distance):
6     AvgS = 0
7     unique_labels = np.unique(labels)
8     for i in range(len(labels)):
9         cur_sample = data[i]
10        cur_label = labels[i]
11
12        same_cluster = data[labels == cur_label]
13        a = np.mean([distance(cur_sample - sample, axis=0) \
14                    for sample in same_cluster if not np.array_equal(cur_sample, sample)])

```

```

15     b = np.inf
16     for label in unique_labels:
17         if label != cur_label:
18             other_cluster = data[labels == label]
19             b = min(b, np.mean([distance(cur_sample - sample, axis=0) \
20                 for sample in other_cluster]))
21
22     AvgS += (b - a) / max(a, b)
23
24     return AvgS / len(labels)

```

```

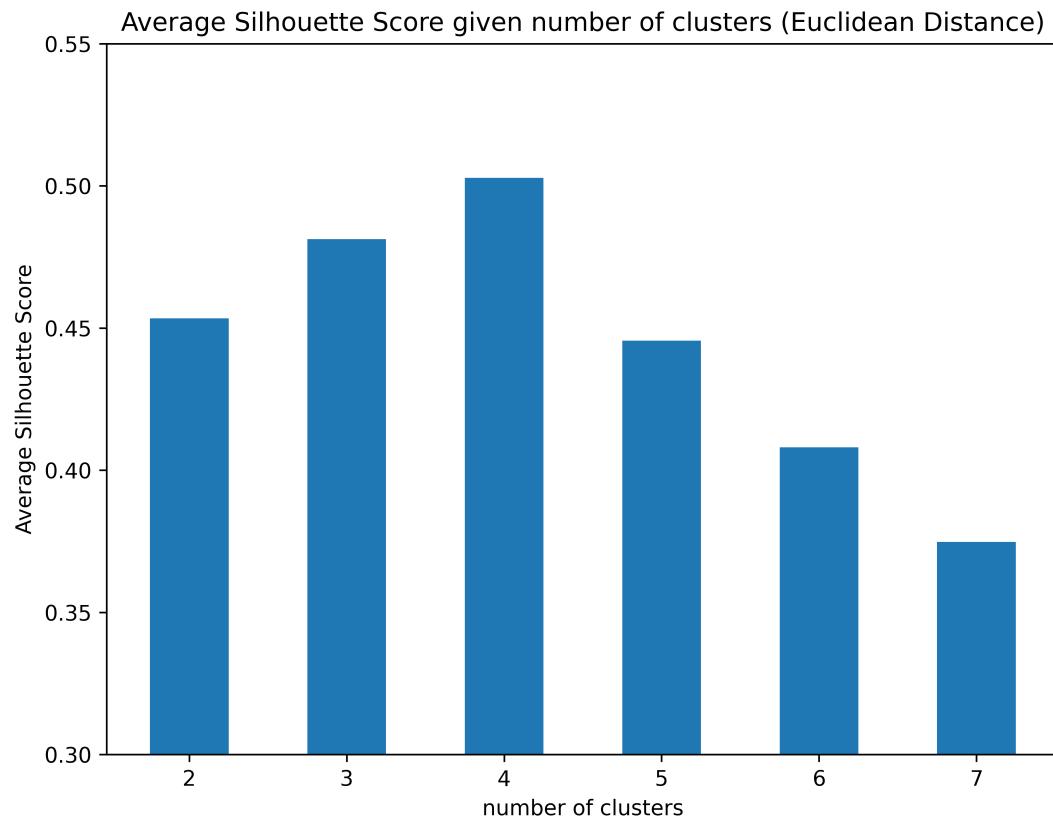
1 class Kmeans:
2     def __init__(self, k, distance, max_iters=10000, tol=1e-4):
3         self.k = k
4         self.distance = distance
5         self.max_iters = max_iters
6         self.tol = tol
7         self.centroids = None
8
9     def _init_centroids(self, data):
10        indices = np.random.choice(data.shape[0], self.k, replace=False)
11        return data[indices]
12
13    def _assign_clusters(self, data):
14        distances = self.distance(data[:, np.newaxis] - self.centroids, axis=2)
15        return np.argmin(distances, axis=1)
16
17    def _update_centroids(self, data, labels):
18        return np.array([data[labels == i].mean(axis=0) for i in range(self.k)])
19
20    def fit(self, data):
21        self.centroids = self._init_centroids(data)
22        for i in range(self.max_iters):
23            labels = self._assign_clusters(data)
24            new_centroids = self._update_centroids(data, labels)
25            if np.all(self.distance(new_centroids - self.centroids, axis=1) < self.tol):
26                break
27            self.centroids = new_centroids
28
29    def predict(self, data):
30        return self._assign_clusters(data)

```

```

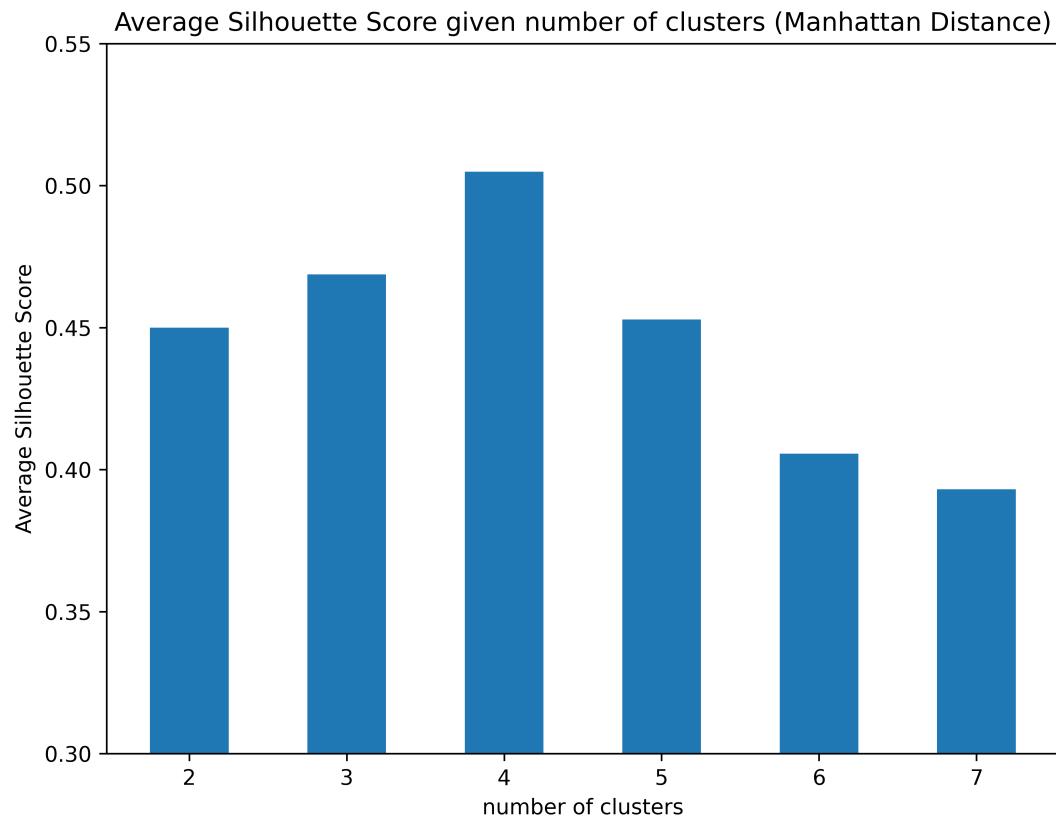
1 np.random.seed(42)
2 data = df.iloc[:, :2].to_numpy()
3 K = [i for i in range(2, 8, 1)]
4 AvgSScore_Euc = []
5 for i in K:
6     model = Kmeans(i, EuclideanDistance)
7     model.fit(data)
8     labels = model.predict(data)
9     AvgSScore_Euc.append(float(AvgSilhouetteScore(data, labels, EuclideanDistance)))
10
11 plt.figure(figsize=(8, 6))
12 plt.bar(K, AvgSScore_Euc, width=0.5)
13 plt.title('Average Silhouette Score given number of clusters (Euclidean Distance)')
14 plt.xlabel('number of clusters')
15 plt.ylabel('Average Silhouette Score')
16 plt.ylim(0.3, 0.55)
17 plt.show()

```



For Manhattan Distance,

```
1 AvgSScore_Man = []
2 for i in K:
3     model = Kmeans(i, EuclideanDistance)
4     model.fit(data)
5     labels = model.predict(data)
6     AvgSScore_Man.append(float(AvgSilhouetteScore(data, labels, ManhattanDistance)))
7
8 plt.figure(figsize=(8, 6))
9 plt.bar(K, AvgSScore_Man, width=0.5)
10 plt.title('Average Silhouette Score given number of clusters (Manhattan Distance)')
11 plt.xlabel('number of clusters')
12 plt.ylabel('Average Silhouette Score')
13 plt.ylim(0.3, 0.55)
14 plt.show()
```



(c) In this problem, we set $k \in \{1, 2, 3, 4\}$ and do EM algorithm. We firstly initialize weight, mean and covariance matrix, then repeat e-step, m-step until either reach max iteration or converge. The code is as follows (follow the above):

```

1  from scipy.stats import multivariate_normal
2  class EM:
3      def __init__(self, k, max_iter=1000, tol=1e-3):
4          self.k = k
5          self.max_iter=max_iter
6          self.tol = tol
7          self.weight = None
8          self.mean = None
9          self.covariances = None
10
11     def _init_parameters(self, data):
12         self.weight = np.ones(self.k) / self.k
13         indices = np.random.choice(data.shape[0], self.k, replace=False)
14         self.mean = data[indices]
15         self.covariances = [np.eye(data.shape[1]) for _ in range(self.k)]
16         # print(self.weight, self.mean, self.covariances)
17
18     def fit(self, data):
19         self._init_parameters(data)
20         for _ in range(self.max_iter):
21             pre_likelihood = self.log_Marginallikelihood(data)
22             post_prob = self.e_step(data)
23             self.m_step(data, post_prob)
24             log_likelihood = self.log_Marginallikelihood(data)
25             if np.abs(log_likelihood - pre_likelihood) < self.tol:
26                 break
27
28
29     def log_Marginallikelihood(self, data):
30         log_marginallikelihood = 0
31         for k in range(self.k):
32             log_marginallikelihood += np.log(self.weight[k] * \
33                 multivariate_normal(self.mean[k], self.covariances[k]).pdf(data).sum())
34         return log_marginallikelihood
35
36     def e_step(self, data):
37         post_prob = np.zeros((data.shape[0], self.k))
38         for k in range(self.k):
39             post_prob[:, k] = self.weight[k] * \
40                 multivariate_normal(self.mean[k], self.covariances[k]).pdf(data)
41         post_prob /= post_prob.sum(axis=1, keepdims=True)
42         return post_prob
43
44     def m_step(self, data, post_prob):
45         self.weight = post_prob.sum(axis=0) / data.shape[0]
46         self.mean = np.dot(post_prob.T, data) / post_prob.sum(axis=0)[:, np.newaxis]
```

```

47     data_centered = data[:, np.newaxis, :] - self.mean
48     post_prob_extended = post_prob[:, :, np.newaxis] # Shape: (n_samples, n_components, 1)
49     self.covariances = np.einsum(
50         'ijk,ijl->jkl', data_centered * post_prob_extended, data_centered
51     ) / post_prob.sum(axis=0)[:, np.newaxis, np.newaxis]
52
53 k = [1, 2, 3, 4]
54 for i in k:
55     np.random.seed(42)
56     model = EM(i)
57     model.fit(data)
58     print(f'{i} Gaussian has weight {model.weight}, mean {model.mean}')
59     print('-----')

```

The result of weight and mean is as follows:

```

1 Gaussian has weight [1.]
mean [[ 0.487239 -0.31368064]]
-----
2 Gaussian has weight [0.45573887 0.54426113]
mean [[ 0.07674448 4.14786009]
 [ 0.83096796 -4.0495667 ]]
-----
3 Gaussian has weight [0.24776738 0.52700064 0.22523198]
mean [[-3.01428079 3.07826986]
 [ 0.68580732 -4.2057353 ]
 [ 3.87448791 5.06166878 ]]
-----
4 Gaussian has weight [0.24928593 0.23223658 0.18818314 0.33029435]
mean [[-3.01962034 3.02734944]
 [-0.91764389 -6.84225433]
 [ 4.18756722 5.92714792]
 [ 2.0135639 -1.80058744 ]]
-----
```

We can see with the increasing of k , the estimated values are approaching the true values. The contour of density function is shown below:

