

Homework 4

仲译为 122090814

Due: May. 9 23:59

Q1 (1) For the model 1 M_1

True Positive: {1, 2, 5}

True Negative: {3, 7, 8, 10}

False Positive: {4}

False Negative: {6, 9}

$$\Rightarrow \begin{cases} TP_1 = 3 \\ TN_1 = 4 \\ FP_1 = 1 \\ FN_1 = 2 \end{cases} \Rightarrow$$

Confusion matrix

Predicted \ Actual	P	N
P	3	2
N	1	4

$$\text{precision}_1 = \frac{TP_1}{TP_1 + FP_1} = \frac{3}{4} = 0.75 \quad \text{recall}_1 = \frac{TP_1}{TP_1 + FN_1} = \frac{3}{5} = 0.6$$

$$\text{accuracy}_1 = \frac{TP_1 + TN_1}{10} = \frac{7}{10} = 0.7$$

Confusion Matrix

For model 2 M_2

True Positive: {1}

True Negative: {4, 7, 8, 10}

False Negative: {2, 5, 6, 9}

False Positive: {3}

$$\Rightarrow \begin{cases} TP_2 = 1 \\ TN_2 = 4 \\ FN_2 = 4 \\ FP_2 = 1 \end{cases} \Rightarrow$$

Predicted \ Actual	P	N
P	1	4
N	1	4

$$\text{precision}_2 = \frac{TP_2}{TP_2 + FP_2} = \frac{1}{2} = 0.5 \quad \text{recall}_2 = \frac{TP_2}{TP_2 + FN_2} = \frac{1}{5} = 0.2$$

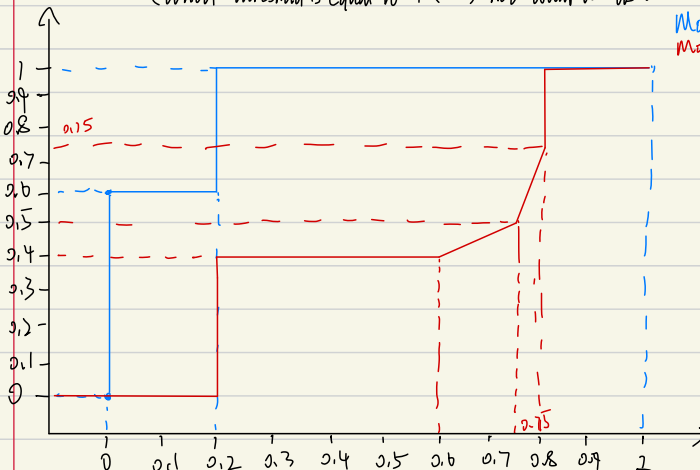
$$\text{accuracy}_2 = \frac{TP_2 + TN_2}{10} = \frac{5}{10} = 0.5$$

(2) Roughly plot the ROC graph of both model

(When threshold is equal to $P(+)$, not count it as one of TN, TP, FN, FP)

$$TPR = \frac{TP}{TP + FN}$$

$$FPR = \frac{FP}{FP + TN}$$



Model 1
Model 2

According to the ROCs

It's obvious that $AUC_1 > AUC_2$

Thus, model 1 performs better than model 2

Q2 (1) Let d_{ik} denotes the Euclidean distance between point i to the mean of cluster k

Let A_1, A_6, A_8 be the initial clusters 1, 2, 3

$$d_{21} = 1, d_{22} = 2\sqrt{5}, d_{23} = \sqrt{29} \Rightarrow A_2 \text{ belongs to cluster 1}$$

$$d_{31} = \sqrt{5}, d_{32} = \sqrt{34}, d_{33} = \sqrt{37} \Rightarrow A_3 \rightarrow 1$$

$$d_{41} = 2, d_{42} = \sqrt{5}, d_{43} = 3\sqrt{2} \Rightarrow A_4 \rightarrow 1$$

$$d_{51} = 3, d_{52} = \sqrt{2}, d_{53} = \sqrt{13} \Rightarrow A_5 \rightarrow 2$$

$$d_{71} = \sqrt{17}, d_{72} = 2, d_{73} = \sqrt{5} \Rightarrow A_7 \rightarrow 2$$

Cluster means:

$$\text{Cluster 1: } \left(\frac{0+0-1+2}{4}, \frac{0+1+2+0}{4} \right) = (0.25, 0.75)$$

$$\text{Cluster 2: } \left(\frac{3+4+4}{3}, \frac{0-1+1}{3} \right) = \left(\frac{11}{3}, 0 \right)$$

$$\text{Cluster 3: } (5, 3)$$

$$d_{11} = \frac{\sqrt{10}}{4}, d_{12} = \frac{11}{3}, d_{13} = \sqrt{34} \Rightarrow A_1 \rightarrow 1$$

$$d_{21} = \frac{\sqrt{2}}{4}, d_{22} = \frac{\sqrt{130}}{3}, d_{23} = \sqrt{29} \Rightarrow A_2 \rightarrow 1$$

$$d_{31} = \frac{5}{4}\sqrt{2}, d_{32} = \frac{\sqrt{41}}{3}, d_{33} = \sqrt{37} \Rightarrow A_3 \rightarrow 1$$

$$d_{41} = \frac{\sqrt{58}}{4}, d_{42} = \frac{11}{3}, d_{43} = 3\sqrt{2} \Rightarrow A_4 \rightarrow 2$$

$$d_{51} = \frac{\sqrt{130}}{4}, d_{52} = \frac{1}{3}, d_{53} = \sqrt{13} \Rightarrow A_5 \rightarrow 2$$

$$d_{61} = \frac{\sqrt{170}}{4}, d_{62} = \frac{\sqrt{10}}{3}, d_{63} = \sqrt{17} \Rightarrow A_6 \rightarrow 2$$

$$d_{71} = \frac{\sqrt{226}}{4}, d_{72} = \frac{\sqrt{10}}{3}, d_{73} = \sqrt{5} \Rightarrow A_7 \rightarrow 2$$

$$d_{83} = 0 \Rightarrow A_8 \rightarrow 3$$

New cluster means:

$$\text{Cluster 1: } \left(\frac{0+0-1}{3}, \frac{0+1+2}{3} \right) \Rightarrow \left(-\frac{1}{3}, 1 \right)$$

$$\text{Cluster 2: } \left(\frac{2+3+4+4}{4}, \frac{0+0-1+1}{4} \right) \Rightarrow \left(\frac{13}{4}, 0 \right)$$

$$\text{Cluster 3: } (5, 3)$$

The algorithm doesn't converge after the first algorithm

Q2 (2) Keep the same notation as in (a)

A_1, A_6, A_8 be the initial cluster centers.

$$d_{21} = 1, d_{12} = 2\sqrt{5}, d_{23} = \sqrt{29} \Rightarrow A_2 \rightarrow 1$$

$$d_{31} = \sqrt{5}, d_{32} = \sqrt{34}, d_{33} = \sqrt{37} \Rightarrow A_3 \rightarrow 1$$

$$d_{41} = \sqrt{5}, d_{43} = 3\sqrt{2} \Rightarrow A_4 \rightarrow 2$$

$$d_{51} = 3, d_{52} = \sqrt{2}, d_{53} = \sqrt{3} \Rightarrow A_5 \rightarrow 2$$

For A_7 , we detect a must link with A_8 but A_8 is certain to be in 3
 $\Rightarrow A_7 \rightarrow 3$

Cluster means:

$$\text{Cluster 1: } \left(\frac{0+0-1}{3}, \frac{0+1+2}{3} \right) \Rightarrow \left(-\frac{1}{3}, 1 \right)$$

$$\text{Cluster 2: } \left(\frac{2+3+4}{3}, \frac{0+0-1}{3} \right) \Rightarrow \left(3, -\frac{1}{3} \right)$$

$$\text{Cluster 3: } \left(\frac{4+5}{2}, \frac{1+3}{2} \right) \Rightarrow \left(\frac{9}{2}, 2 \right)$$

$$\begin{aligned} d_{11} &= \frac{\sqrt{10}}{3}, d_{12} = \frac{\sqrt{82}}{3}, d_{13} = \frac{\sqrt{97}}{2} \\ d_{41} &= \frac{\sqrt{5}}{3}, d_{42} = \frac{\sqrt{10}}{3}, d_{43} = \frac{\sqrt{41}}{2} \\ d_{21} &= \frac{1}{3}, d_{22} = \frac{\sqrt{97}}{3}, d_{23} = \frac{\sqrt{82}}{2} \Rightarrow A_2 \rightarrow 1 \\ d_{31} &= \frac{\sqrt{13}}{3}, d_{32} = \frac{\sqrt{173}}{3}, d_{33} = \frac{11}{2} \Rightarrow A_3 \rightarrow 1 \\ d_{51} &= \frac{\sqrt{34}}{3}, d_{52} = \frac{1}{3}, d_{53} = \frac{5}{2} \Rightarrow A_5 \rightarrow 2 \\ d_{61} &= \frac{\sqrt{205}}{3}, d_{62} = \frac{\sqrt{13}}{3}, d_{63} = \frac{5\sqrt{13}}{2} \Rightarrow A_6 \rightarrow 2 \\ d_{71} &= \frac{13}{3}, d_{72} = \frac{\sqrt{13}}{3}, d_{73} = \frac{\sqrt{13}}{2} \\ d_{81} &= \frac{2\sqrt{73}}{3}, d_{82} = \frac{\sqrt{5}}{2}, d_{83} = \frac{5}{2} \end{aligned} \left. \begin{array}{l} \text{Also, } \arg\min_k (d_{1k}) = 1, \\ \arg\min_k (d_{4k}) = 2 \\ \arg\min_k (d_{7k}) = 3 \end{array} \right\} \begin{array}{l} A_1 \text{ and } A_4 \text{ belongs different clusters} \\ A_7 \text{ and } A_8 \text{ belongs to the same cluster} \end{array}$$

New cluster means:

$$\text{Cluster 1: } \left(\frac{0+0-1}{3}, \frac{0+1+2}{3} \right) \Rightarrow \left(-\frac{1}{3}, 1 \right)$$

$$\text{Cluster 2: } \left(\frac{2+3+4}{3}, \frac{0+0-1}{3} \right) \Rightarrow \left(3, -\frac{1}{3} \right)$$

$$\text{Cluster 3: } \left(\frac{4+5}{2}, \frac{1+3}{2} \right) \Rightarrow \left(\frac{9}{2}, 2 \right)$$

The algorithm converges after the first iteration

Q3 (1)

(Likelihood to be maximized)

$$(2) \quad \mathcal{L}(\mu, \sigma, \lambda) = - \sum_{n=1}^3 \sum_{k=1}^2 r_{n,k} \ln(\pi_k) - \frac{1}{\sqrt{2\pi\sigma_k^2}} \sum_{n=1}^3 \sum_{k=1}^2 r_{n,k} \frac{(x^{(n)} - \mu_k)^2}{2\sigma_k^2} + \lambda \left(1 - \sum_{k=1}^K \pi_k\right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = - \sum_{i=1}^3 \sum_{k=1}^2 \frac{r_{i,k}}{\pi_k} - \lambda = 0 \Rightarrow \sum_{k=1}^2 \lambda \pi_k = - \sum_{k=1}^2 \sum_{i=1}^3 r_{i,k} = -(1+2.4+0+0+0.6+1) = -3$$

$$\sum_{k=1}^K \lambda \pi_k = \lambda = -3$$

For each k in K , $\lambda \pi_k = - \sum_{i=1}^3 r_{i,k} \Rightarrow \pi_k = \frac{1}{3} \sum_{i=1}^3 r_{i,k}$

$$\pi_1 = \frac{1}{3} \sum_{i=1}^3 r_{i,1} = \frac{1}{3}(1+2.4+0) = \frac{7}{15} ; \quad \pi_2 = \frac{1}{3} \sum_{i=1}^3 r_{i,2} = \frac{1}{3}(0+0.6+1) = \frac{8}{15}$$

$$(3) \quad \frac{\partial \mathcal{L}}{\partial \mu_k} = \frac{\sum_{i=1}^3 r_{i,k} \frac{(x^{(i)} - \mu_k)}{\sigma_k^2}}{\sum_{i=1}^3 r_{i,k}} = \sum_{i=1}^3 r_{i,k} \cdot \left(\frac{\mu_k}{\sigma_k^2} - \frac{x^{(i)}}{\sigma_k^2} \right) = 0$$

$$\Rightarrow \sum_{i=1}^3 r_{i,k} (\mu_k - x^{(i)}) = 0 \Rightarrow \sum_{i=1}^3 r_{i,k} \mu_k = \sum_{i=1}^3 r_{i,k} x^{(i)}$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^3 r_{i,k} x^{(i)}}{\sum_{i=1}^3 r_{i,k}}$$

$$\mu_1 = \frac{\sum_{i=1}^3 r_{i,1} x^{(i)}}{\sum_{i=1}^3 r_{i,1}} = \frac{1 \times 1 + 2.4 \times 1.5 + 0 \times 2.0}{1 + 2.4 + 0} = \frac{25}{7}$$

$$\mu_2 = \frac{\sum_{i=1}^3 r_{i,2} x^{(i)}}{\sum_{i=1}^3 r_{i,2}} = \frac{0 \times 1 + 0.6 \times 1.5 + 1 \times 2.0}{0 + 0.6 + 1} = \frac{65}{4}$$

Q4 (1) $X = \left\{ \begin{bmatrix} 2 \\ 0 \\ 1 \\ -3 \\ -2 \end{bmatrix}, \begin{bmatrix} 0 \\ 2 \\ -3 \\ -3 \\ -2 \end{bmatrix}, \begin{bmatrix} 1 \\ 2 \\ 1 \\ 3 \\ -2 \end{bmatrix}, \begin{bmatrix} -1 \\ 1 \\ 3 \\ 2 \\ -1 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ 1 \\ -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 \\ 3 \\ -1 \\ 1 \\ -2 \end{bmatrix}, \begin{bmatrix} -2 \\ 3 \\ -3 \\ 3 \\ 2 \end{bmatrix}, \begin{bmatrix} -2 \\ -2 \\ 2 \\ 3 \\ -2 \end{bmatrix}, \begin{bmatrix} -2 \\ -3 \\ 1 \\ -2 \\ -3 \end{bmatrix}, \begin{bmatrix} -3 \\ 2 \\ 0 \\ -1 \\ -2 \end{bmatrix} \right\}$

$$\underline{\mu} = \frac{10}{n} \bar{x} = (-0.4, 0.8, 0.2, 0.2, -1.3)^T$$

$$S = \frac{1}{10} \sum_{n=1}^{10} (x^{(n)} - \underline{\mu})(x^{(n)} - \underline{\mu})^T$$

$$= \begin{bmatrix} 3.04 & 0.82 & -0.02 & -0.82 & -0.12 \\ 0.82 & 3.76 & -2.16 & 1.04 & 1.04 \\ -0.02 & -2.16 & 3.56 & 0.76 & -0.84 \\ -0.82 & 1.04 & 0.76 & 5.56 & 1.16 \\ -0.12 & 1.04 & -0.84 & 1.16 & 2.21 \end{bmatrix}$$

(obtained using Python)

Using Python, we obtain two rough eigenvalues

(0.82311524, 1.53027334, 3.07614543, 5.93067614, 6.76978985)

and the corresponding eigenvectors

$$\begin{bmatrix} -0.33192081 \\ 0.60584079 \\ 0.61980825 \\ -0.33144989 \\ 0.16960016 \end{bmatrix} \begin{bmatrix} 0.16894209 \\ -0.33313495 \\ 0.08283824 \\ -0.15734291 \\ 0.91041788 \end{bmatrix} \begin{bmatrix} 0.87848693 \\ 0.17625638 \\ 0.40532527 \\ 0.14382432 \\ -0.11054589 \end{bmatrix} \begin{bmatrix} 0.29809153 \\ 0.39493632 \\ -0.58025264 \\ -0.64618454 \\ 0.03031732 \end{bmatrix} \begin{bmatrix} -0.02625442 \\ 0.57873744 \\ -0.32862419 \\ 0.65356276 \\ 0.35949345 \end{bmatrix}$$

We only choose the two of highest eigenvalues:

6.76978985 and 5.93067614

with corresponding eigenvector $\begin{bmatrix} -0.02625442 \\ 0.57873744 \\ -0.32862419 \\ 0.65356276 \\ 0.35949345 \end{bmatrix}$ and $\begin{bmatrix} 0.29809153 \\ 0.39493632 \\ -0.58025264 \\ -0.64618454 \\ 0.03031732 \end{bmatrix}$

which are the 2 principal components

Q4

$$U = \begin{bmatrix} -0.02625442 & 0.29809153 \\ 0.57873744 & 0.39493632 \\ -0.32862419 & -0.58025264 \\ 0.65356276 & -0.64618454 \\ 0.35949345 & 0.03031732 \end{bmatrix}$$

By $\tilde{x}^{(n)} = U^T(x^{(n)} - \underline{\mu})$, we obtain

$$\tilde{x}_1 = \begin{bmatrix} -3.13194676 \\ 1.98183693 \end{bmatrix} \quad \tilde{x}_2 = \begin{bmatrix} -0.60746566 \\ 4.49653708 \end{bmatrix}$$

$$\tilde{x}_3 = \begin{bmatrix} 1.97315969 \\ -1.40348923 \end{bmatrix} \quad \tilde{x}_4 = \begin{bmatrix} 0.4956134 \\ -2.87861204 \end{bmatrix}$$

$$\tilde{x}_5 = \begin{bmatrix} -0.72008587 \\ 0.48232827 \end{bmatrix} \quad \tilde{x}_6 = \begin{bmatrix} 1.87576557 \\ 1.74241301 \end{bmatrix}$$

$$\tilde{x}_7 = \begin{bmatrix} 5.38313097 \\ 0.53945236 \end{bmatrix} \quad \tilde{x}_8 = \begin{bmatrix} -0.591651 \\ -4.45776178 \end{bmatrix}$$

$$\tilde{x}_9 = \begin{bmatrix} -4.46907149 \\ -1.07184006 \end{bmatrix} \quad \tilde{x}_{10} = \begin{bmatrix} -0.20744945 \\ 0.56913546 \end{bmatrix}$$

(Calculated by Python)

Coding part

Task I: Implement k -means

By running the k -means algorithm 100 times, we obtain the following table:

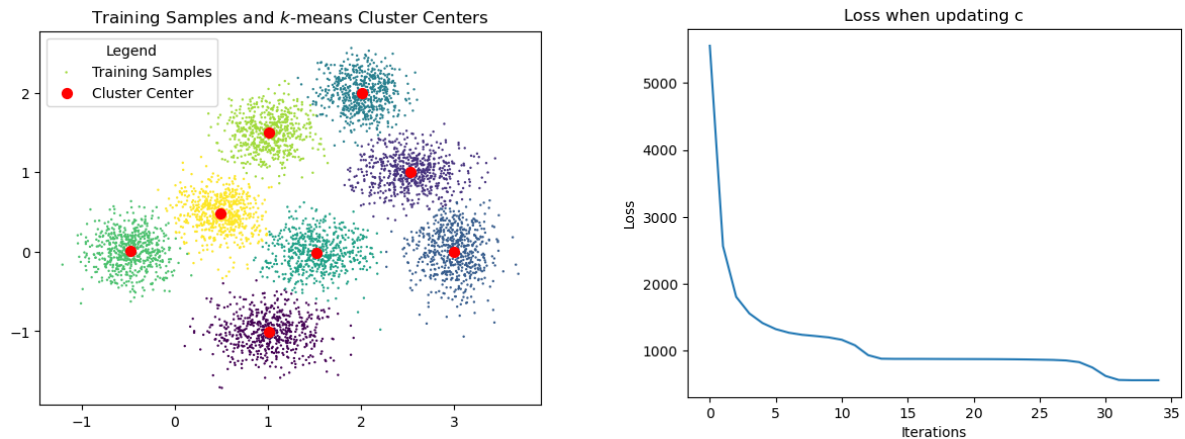
Convergence status	Number of appearances
NaN-value	58
Local minimum	39
Global minimum	3

(Global minimum appears at the 53rd, the 69th, and the 91st time)

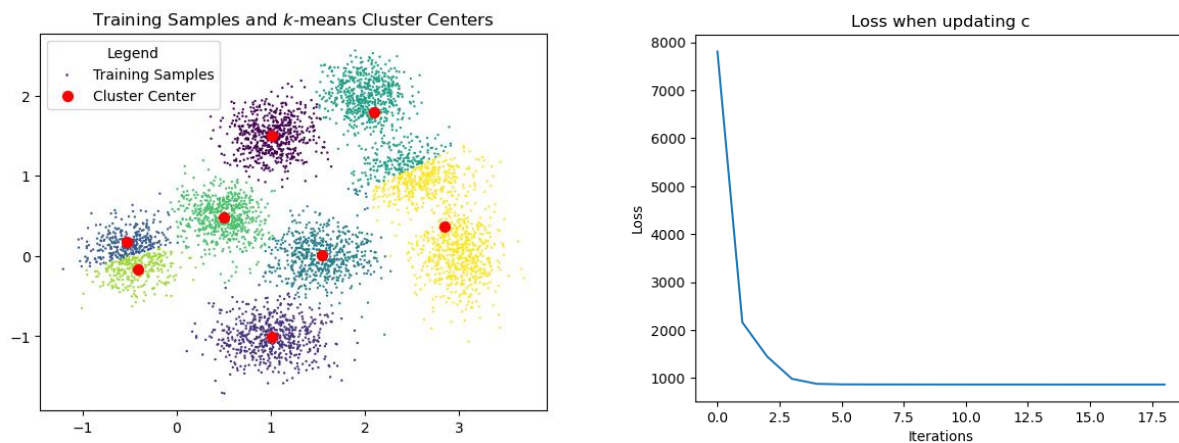
The frequency of the algorithm converging to the global minimum is **around 0.03**.

Below showing the scatter graph of cluster centers and the loss graph under different conditions:

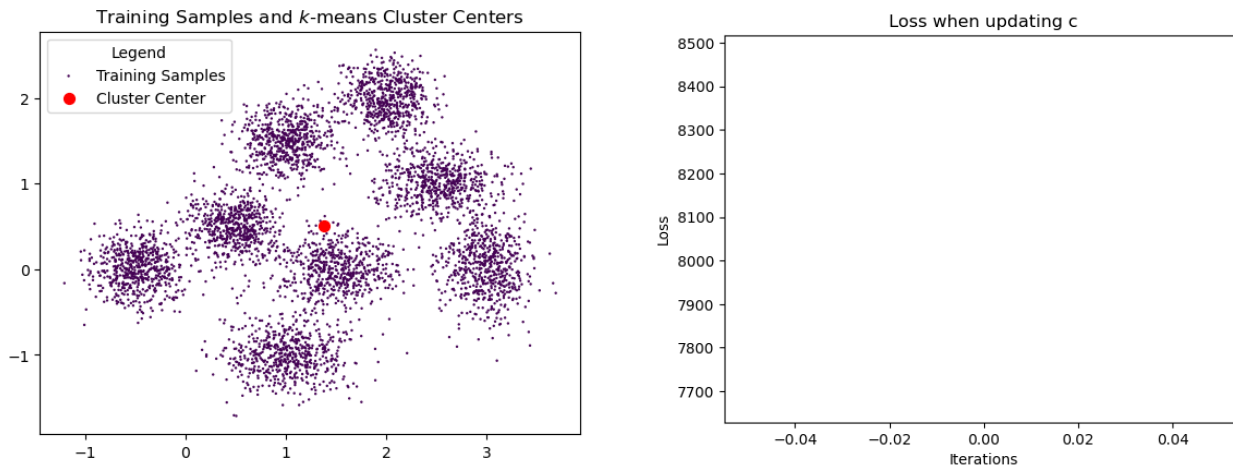
Global minimum:



Local minimum:



NaN-value:



Besides, by repeatedly running the k -means algorithm multiple times, we obtain 10 results that converge to the global minimum, which have the iteration times of {25, 18, 11, 31, 33, 33, 38, 33, 25, 24}. The average iteration time is 27.1, which can roughly represent the convergence speed of the k -means algorithm.

Explanation for codes:

1. Key Idea:
Iterate all centroids and calculate the distances from all points to each centroid. Obtain a (K, N) sized array (*distances*). By using argmin function along axis=0, we obtain the index array (*index*) for each data point referring to which cluster centroid it belongs to, and by using min function along axis=0, we obtain the distance array (didn't create an array to store, directly compute) from the data point to the centroid they belong. By taking the sum of square for elements in the distance array
2. Important variables: (a) *distances*: first list, then changed to ndarray, storing distances from all points to all centroids; (b) *c*, *c_new*: ndarray, sized (K, D), centroids from the last iteration and newly computed respectively; (c) *index*, ndarray, as explained above.

Task II: Implement k -means++

By running the k -means++ algorithm for 50 times, we obtain the following table:

Convergence status	Number of appearances
NaN-value	0
Local minimum	22
Global minimum	28

The frequency of the algorithm converging to the global minimum is **around 0.56**. Compared with the k -means algorithm, k -means++ algorithm has a higher probability to converge to a global minimum.

Besides, by repeatedly running the k -means++ algorithm multiple times, we obtain 10 results that converge to the global minimum, which have the iteration times of {20, 11, 5, 9, 10, 10, 7, 14, 11, 10}. The average iteration time is 10.7, which can roughly represent the convergence speed of the k -means++ algorithm.

Compare the convergence speed with the convergence speed of the k -means algorithm. We could have a directional hypothesis that the average convergence speed of the k -means++ algorithm is faster than k -means algorithm.

Explanation for codes:

1. Key Idea:

The algorithm follows the instruction from the task description. First pick a sample data point into set I , which is the array of all centroids. Then generate the probability for each data point. In the generating process, first calculate the square distances from all data points to centroids chosen. Then, for each data point, keep the shortest distance from it to the closest given centroid. Append the data into a list (*lst_all*). Change the list into uniform data by dividing it with the sum of the whole list (and convert to a ndarray stored in p), we obtain the probability for all points to be the new centroid, which can be applied in the `np.random.choice()` function.

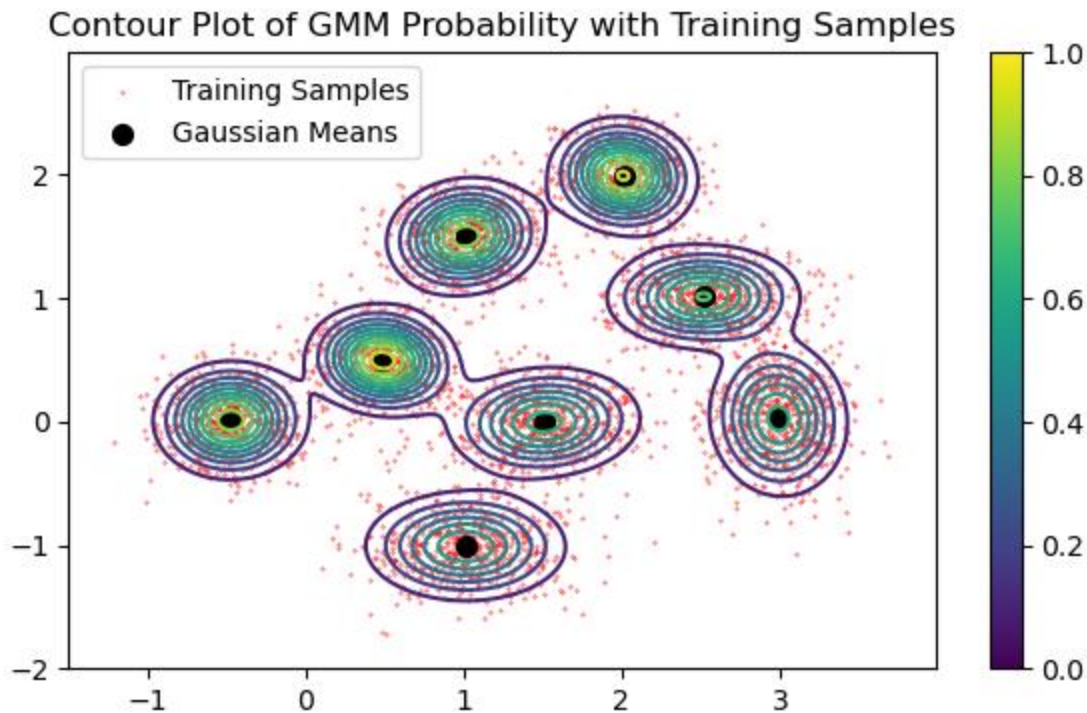
2. Important variables:

(a) *lst_all*: as explained above; (b) *lst_distance*: list storing ndarrays that contain square distances from all data points to all centroids that have been chosen in the set I (centroids array).

(The manuscript of the comparison is attached at the end of the report)

Task III: Implement a Gaussian Mixture Model

Implement a Gaussian Mixture Model using the Expectation-Maximization algorithm. The initial cluster centers generated by the k -means++ algorithm from task 2 are ensured to be at the local minimum. After we obtain all parameters of the Gaussian Mixture Model, we could plot the 2-d contour plot of the GMM model as follows:



It's obvious to see that the GMM fits the data well.

The log-likelihood of the GMM on the training set is -9693.226184393608.

The log-likelihood of the GMM on the development set is -1678.0141368062066.

```
print("The log-likelihood on the training set is {}".format(np.sum(gmm_log_prob(train_x, pi, miu, sigma))))
print("The log-likelihood on the development set is {}".format(np.sum(gmm_log_prob(dev_x, pi, miu, sigma))))
```

2] ✓ 0.0s

The log-likelihood on the training set is -9693.226184393608
The log-likelihood on the development set is -1678.0141368062066

Explanation for codes:

Key Idea:

- $e_step()$: directly follow the formula given by the slide. We first calculate the numerator of the formula and then divide each with the sum according to clusters (K), which is the sum of row (axis=0) for each entry in γ which is a ndarray sized (N, K). Return γ ;
- $m_step()$: still directly follow the formula given by the slide. Enumerate from 0 to N and enumerate from 0 to K inside, we cumulatively calculate the N_k (N_k , sum the j^{th} entry with the i, j -entry of γ) and μ_k (μ_k , similar but i, j -entry of γ times each data point). Calculate π_k . Calculate the covariance matrices after centroids which are used in calculating those matrices. Return π_{new} , μ_{new} , and σ_{new} ;

- c) *em()*: simply initialize the *Nk*, *pi*, *miu*, and *sigma* with idea similar to *m_step()*. There's nothing important to explain. In the iteration, input the training data into *e_step()* and obtain the above 4 updated model parameters. Also compute the log-likelihood and append it in a list. The stopping criterion is to check if the difference between the new log-likelihood and the log-likelihood from last iteration is smaller than tolerance (set to be 10^{-5}), if smaller, then break, except reaching the maximum iteration times set as 1000. Finally, return *pi*, *miu*, and *sigma*, which are the GMM parameters.
- d) All the important parameters are explained above, more to see code.

MANUSCRIPT of comparing K-means and K-means++

100波

58

记录

中间点

正正正正正正正正正正正

\angle means

39

有值, 非 global minimum

正正正正正正正正

global minimum 第53次 第69次 第91次

记录

中间点

C

$$k + y$$

local minimum

正正正正

22

global minimum

正正正正正下

28

K-means

k++

converge
speed

local global

23 25

42 18

26 11

32 31

28 33

29 33

29 38

26 33

42 25

38 24

29

29

42

27

18

36

44

32

40

16

30

37

23

29

38

25

30

40

55 28

25 27

local

20

14

22

12

40

22

20

26

12

18

10

global

20

11

5

9

10

10

7

14

11

10