

第六章

一、高斯贝叶斯分类器和高斯朴素贝叶斯分类器预测样本 x 的类别

1. 设密度 x_1 , 含糖率 x_2 , 西瓜 C (C_1 是, C_2 否)

$$\text{由贝叶斯公式得 } p(C=c_j | x_1, x_2) = \frac{p(x_1, x_2 | C=c_j) p(C=c_j)}{p(x_1, x_2)} \propto p(x_1, x_2 | C=c_j) \cdot p(C=c_j)$$

① 高斯贝叶斯.

$$\text{多维高斯: } p(x_1, x_2 | C=c_j) = \frac{1}{2\pi} \cdot \frac{1}{|\Sigma_j|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x-\mu_j)^T \Sigma_j^{-1}(x-\mu_j)\right)$$

$$\text{这里 } x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \mu_j = \begin{bmatrix} \mu_{1j} \\ \mu_{2j} \end{bmatrix} \leftarrow \begin{matrix} \text{均值} \\ \mu_{ij} = \frac{1}{N(C=c_j)} \sum_{x_i \in C_j} x_i \end{matrix} \quad \Sigma_j = \begin{bmatrix} \text{cov}(x_1, x_1) & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \text{cov}(x_2, x_2) \end{bmatrix} \leftarrow \begin{matrix} \text{协方差} \end{matrix}$$

$$p(C=c_1) = \frac{N(C=c_1)}{N_{\text{all}}} = \frac{8}{17}$$

$$p(C=c_2) = \frac{N(C=c_2)}{N_{\text{all}}} = \frac{9}{17}$$

↓
使用 numpy cov 计算
无偏估计

a. 当 $C=C_1$ 时

把样本数据代入得

$$p(0.5, 0.3 | C_1) = 9.490316446182263 \quad p(C=C_1) = \frac{8}{17}$$

$$\therefore p(C=C_1 | 0.5, 0.3) = p(0.5, 0.3 | C_1) \cdot p(C=C_1) = 4.466031268791653 \Rightarrow P_{\text{good}}$$

b. 当 $C=C_2$ 时

把样本数据代入得

$$p(0.5, 0.3 | C_2) = 2.8854275027464418 \quad p(C=C_2) = \frac{9}{17}$$

$$\therefore p(C=C_2 | 0.5, 0.3) = p(0.5, 0.3 | C_2) \cdot p(C=C_2) = 1.527579266159881 \Rightarrow P_{\text{bad}}$$

$\because P_{\text{good}} > P_{\text{bad}} \Rightarrow \text{样本 } x = \begin{bmatrix} 0.5 \\ 0.3 \end{bmatrix} \text{ 是西瓜}$

② 高斯朴素贝叶斯 ↓ 独立性假设

$$P(x_1, x_2 | c = c_j) = P(x_1 | c = c_j) P(x_2 | c = c_j)$$

$$\text{又 } P(x_i | c = c_j) = \frac{1}{\sqrt{2\pi} \sigma_{ij}} \exp\left(-\frac{(x_i - \mu_{ij})^2}{\sigma_{ij}^2}\right) \quad \begin{cases} \text{均值} & \mu_{ij} = \frac{1}{N(c=c_j)} \sum_{x_i \in c_j} x_i \\ \text{方差} & \sigma_{ij}^2 = \frac{1}{N(c=c_j)} \sum_{x_i \in c_j} (x_i - \mu_{ij})^2 \end{cases}$$

$$\Rightarrow P(x_1, x_2 | c = c_j) = \prod_{i=1}^2 \frac{1}{\sqrt{2\pi} \sigma_{ij}} \exp\left(-\frac{(x_i - \mu_{ij})^2}{\sigma_{ij}^2}\right)$$

$$= \frac{1}{2\pi} \frac{1}{|\Sigma_j|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)\right\} \rightarrow \text{可见高斯朴素贝叶斯和高斯贝叶斯公式一致, 不过这里的} \Sigma_j = \begin{bmatrix} \sigma_{1j}^2 & 0 \\ 0 & \sigma_{2j}^2 \end{bmatrix}$$

这里 $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ $\mu_j = \begin{bmatrix} \mu_{1j} \\ \mu_{2j} \end{bmatrix}$

a. 若 $c = c_1$ 时 把数据代入得

$$P(0.5, 0.3 | c_1) = 11.289349211617608 \quad P(c = c_1) = \frac{8}{17}$$

$$P(c_1 | 0.5, 0.3) = 5.312634923114168 \Rightarrow P_{\text{good-naive}}$$

b. 若 $c = c_2$ 时 把数据代入得

$$P(0.5, 0.3 | c_2) = 4.231758950693437 \quad P(c = c_2) = \frac{9}{17}$$

$$P(c_2 | 0.5, 0.3) = 2.2403429738965257 \Rightarrow P_{\text{bad-naive}}$$

$\therefore P_{\text{good-naive}} > P_{\text{bad-naive}} \Rightarrow$ 样本 $x = \begin{bmatrix} 0.5 \\ 0.3 \end{bmatrix}$ 是好瓜

代码:

(1) 读取样本数据, 并分类好瓜和坏瓜

```
1. def read_data(path):
2.     # 密度 含糖率 好瓜
3.     data = pd.read_csv(path, header=None, names=['x1', 'x2', 'c'])
4.     return data
5.
6.
```

```

7. # 将 c=c1 和 c=c2 两种情况下的 x1 x2 分离
8. def split_data(path="data.txt"):
9.     data = np.array(read_data(path))
10.    x1_x2_good = np.array([x[0:2] for x in data if x[2] == '是'])
11.    x1_x2_bad = np.array([x[0:2] for x in data if x[2] == '否'])
12.    return x1_x2_good, x1_x2_bad

```

(2) 求均值、协方差、方差

```

1. # 求均值 x[x1,x2]
2. def mean_value(x):
3.     x1_mean = np.mean(x[:, 0])
4.     x2_mean = np.mean(x[:, 1])
5.     return np.array([x1_mean, x2_mean])
6.
7.
8. # 求协方差 x[x1 x2] 2D 数据
9. def cov(x):
10.    x_cov = np.cov(x.astype(float).T) # 默认无偏估计, 即分母为 n-1
11.    return x_cov
12.
13.
14. # 方差 x[x1,x2]
15. def var_value(x):
16.    x1_var = np.var(x[:, 0])
17.    x2_var = np.var(x[:, 1])
18.    return np.array([x1_var, x2_var])

```

(3) 高斯贝叶斯和高斯朴素贝叶斯通用的公式

```

1. # 高斯贝叶斯/高斯朴素贝叶斯
2. def gaussian_bayes(x_mean, x_cov, x_input):
3.    x_cov_det = np.linalg.det(x_cov) # 协方差矩阵行列式
4.    x_cov_inv = np.linalg.inv(x_cov) # 协方差矩阵逆矩阵
5.    temp = x_input - x_mean
6.    gaussian = 1 / (2 * np.pi) * 1 / (pow(x_cov_det, 0.5)) * np.exp(-
        (1/2 * np.dot(np.dot(temp, x_cov_inv), temp)))
7.    return gaussian

```

(4) 主函数调用

```

1. def main():
2.    x_good, x_bad = read_data.split_data() # x[x1,x2]
3.    # 均值
4.    x_good_mean = mean_value(x_good)

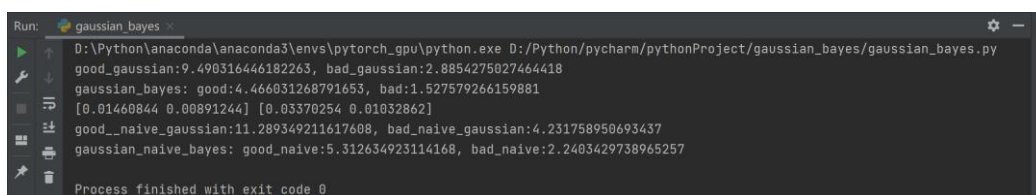
```

```

5. x_bad_mean = mean_value(x_bad)
6. # 协方差矩阵
7. x_good_cov = cov(x_good)
8. x_bad_cov = cov(x_bad)
9. x = np.array([0.5, 0.3])
10. # 先验概率
11. p_good = len(x_good)/(len(x_good) + len(x_bad))
12. p_bad = len(x_bad)/(len(x_good) + len(x_bad))
13.
14. # 高斯贝叶斯
15. good_gaussian = gaussian_bayes(x_good_mean, x_good_cov, x)
16. bad_gaussian = gaussian_bayes(x_bad_mean, x_bad_cov, x)
17. print("good_gaussian:%s, bad_gaussian:%s" % (good_gaussian, bad_gaussian))
18. good = good_gaussian * p_good
19. bad = bad_gaussian * p_bad
20. print("gaussian_bayes: good:%s, bad:%s" % (good, bad))
21.
22. # 高斯朴素贝叶斯
23. # 方差
24. var_good = var_value(x_good)
25. var_bad = var_value(x_bad)
26. print(var_good, var_bad)
27. # 协方差矩阵
28. cov_good_naive = np.zeros((2, 2))
29. cov_good_naive[0][0] = var_good[0]
30. cov_good_naive[1][1] = var_good[1]
31. cov_bad_naive = np.zeros((2, 2))
32. cov_bad_naive[0][0] = var_good[0]
33. cov_bad_naive[1][1] = var_good[1]
34.
35. good_naive_gaussian = gaussian_bayes(x_good_mean, cov_good_naive, x)
36. bad_naive_gaussian = gaussian_bayes(x_bad_mean, cov_bad_naive, x)
37. print("good_naive_gaussian:%s, bad_naive_gaussian:%s" % (good_naive_gaussian, bad_naive_gaussian))
38. good_naive = good_naive_gaussian * p_good
39. bad_naive = bad_naive_gaussian * p_bad
40. print("gaussian_naive_bayes: good_naive:%s, bad_naive:%s" % (good_naive, bad_naive))

```

运行截图：



```

Run: gaussian_bayes
D:\Python\anaconda\anaconda3\envs\pytorch_gpu\python.exe D:/Python/pycharm/pythonProject/gaussian_bayes/gaussian_bayes.py
good_gaussian:9.490316446182263, bad_gaussian:2.8854275027464418
gaussian_bayes: good:4.466031268791653, bad:1.527579266159881
[0.01460844 0.00891244] [0.03370254 0.01032862]
good_naive_gaussian:11.289349211617608, bad_naive_gaussian:4.231758950693437
gaussian_naive_bayes: good_naive:5.312634923114168, bad_naive:2.2403429738965257
Process finished with exit code 0

```

二、Markov 模型

2. 由题可得状态转移矩阵 $A = \begin{bmatrix} 0.8 & 0.2 \\ 0.5 & 0.5 \end{bmatrix}$

第一次没中到第四次命中转移三次 $0 \rightarrow 0 \rightarrow 0 \rightarrow 0$

$$A^3 = \begin{bmatrix} 0.722 & 0.278 \\ 0.695 & 0.305 \end{bmatrix}$$

\therefore 第一次没中, 第四次命中

$\therefore A^3[1][0] = 0.695$ 是第四次命中的概率。

代码:

```
1. # 返回矩阵的 n 阶结果
2. def get_matrix_pow(matrix, n):
3.     ret = matrix
4.     for i in range(n-1):
5.         ret = np.dot(ret, matrix)
6.     return ret
7.
8.
9. def main():
10.     A = np.array([[0.8, 0.2],
11.                   [0.5, 0.5]])
12.     A_3 = get_matrix_pow(A, 3)
13.     print("第一次没中且第四次射中的概率为: %.3f" % A_3[1][0])
```

三、使用 EM 求解 GMM

3.

$$\arg\max_{\theta} \ln \left(\prod_{j=1}^n p(x_j | \theta) \right) = \arg\max_{\theta} \sum_{j=1}^n \ln \left(\sum_{i=1}^k p(y_j=i, x_j | \theta) \right)$$

构造目标函数 $\sum_{j=1}^n \sum_{i=1}^k p_{ji} \ln(p(y_j=i, x_j | \theta))$ $p_{ji} = p(p_j=i | x_j, \theta^t)$

(1) 目标函数关于 μ_i 求偏导

$$\sum_{j=1}^n p_{ji} \frac{\partial \ln(p(y_j=i, x_j | \theta))}{\partial \mu_i} = \sum_{j=1}^n p_{ji} \frac{\partial \ln(p(y_j=i) \cdot p(x_j | y_j=i, \theta))}{\partial \mu_i}$$

$$= \sum_{j=1}^n p_{ji} \frac{\partial \ln(\pi_i p(x_j | y_j=i, \theta))}{\partial \mu_i} = \sum_{j=1}^n p_{ji} \frac{1}{\pi_i p(x_j | y_j=i, \theta)} \cdot \frac{\partial (\pi_i p(x_j | y_j=i, \theta))}{\partial \mu_i}$$

$$= \sum_{j=1}^n p_{ji} \frac{\pi_i p(x_j | y_j=i, \theta)}{\pi_i p(x_j | y_j=i, \theta)} \cdot \frac{\partial \left[-\frac{1}{2} (x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i) \right]}{\partial \mu_i}$$

$$= \sum_{j=1}^n p_{ji} \cdot \frac{1}{2} \cdot \frac{\partial [(x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i)]}{\partial (x_j - \mu_i)}$$

由 $\frac{\partial (X^T A X)}{\partial X} = (A + A^T)X$ 若 $A = A^T$ 则 $\frac{\partial (X^T A X)}{\partial X} = 2AX$

$$\text{得 原式} = \sum_{j=1}^n p_{ji} \Sigma_i^{-1} (x_j - \mu_i) = 0$$

$$\text{即 } \Sigma_i^{-1} \sum_{j=1}^n p_{ji} (x_j - \mu_i) = 0$$

$$\because \Sigma_i^{-1} \neq 0 \quad \therefore \sum_{j=1}^n p_{ji} x_j = \sum_{j=1}^n p_{ji} \mu_i \Rightarrow \mu_i = \frac{\sum_{j=1}^n p_{ji} x_j}{\sum_{j=1}^n p_{ji}} \quad \text{其中 } p_{ji} = p(p_j=i | x_j, \theta^t)$$

2) 目标函数关于 Σ_i 求导

$$\sum_{j=1}^n p_{ji} \frac{\partial \ln(p(y_j=i, x_j|\theta))}{\partial \Sigma_i} = \sum_{j=1}^n p_{ji} \frac{\partial [\ln|\Sigma_i| + (x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i)]}{\partial \Sigma_i}$$

$$\text{又: } \frac{\partial |A|}{\partial A} = |A| \cdot A^{-1} \Rightarrow \frac{\partial \ln|A|}{\partial A} = A^{-1}$$

$$(x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i) = \text{tr}[(x_j - \mu_i)^T \Sigma_i^{-1} (x_j - \mu_i)] = \text{tr}[\Sigma_i^{-1} (x_j - \mu_i)(x_j - \mu_i)^T]$$

$$\begin{aligned} \frac{\partial \text{tr}(AB)}{\partial A} &= B^T \Rightarrow \frac{\partial \text{tr}[\Sigma_i^{-1} (x_j - \mu_i)(x_j - \mu_i)^T]}{\partial \Sigma_i} = (x_j - \mu_i)(x_j - \mu_i)^T \frac{\partial \Sigma_i^{-1}}{\partial \Sigma_i} \\ &= -(x_j - \mu_i)(x_j - \mu_i)^T \Sigma_i^{-2} \end{aligned}$$

$$\text{原式} = \sum_{j=1}^n p_{ji} [\Sigma_i^{-1} - (x_j - \mu_i)(x_j - \mu_i)^T] = 0$$

$$\text{同乘 } \Sigma_i^2 \text{ 得 } \sum_{j=1}^n p_{ji} [\Sigma_i - (x_j - \mu_i)(x_j - \mu_i)^T] = 0$$

$$\Rightarrow \Sigma_i = \frac{\sum_{j=1}^n p_{ji} (x_j - \mu_i)(x_j - \mu_i)^T}{\sum_{j=1}^n p_{ji}}$$

13) π_i : 此时这是一个有约束极值问题, 约束条件为 $\sum_{i=1}^K \pi_i = 1$
构造拉格朗日函数

$$L(\pi_1, \pi_2, \dots, \pi_K, \lambda) = L(L(\theta)) + \lambda (\sum_{i=1}^K \pi_i - 1)$$

$$= \sum_{j=1}^n \ln \left(\sum_{i=1}^K p(y_j=i, x_j|\theta) \right) + \lambda (\sum_{i=1}^K \pi_i - 1)$$

$$= \sum_{j=1}^n \ln \left(\sum_{i=1}^K \pi_i p(x_j|y_j=i, \theta) \right) + \lambda (\sum_{i=1}^K \pi_i - 1)$$

$$\frac{\partial L}{\partial \pi_i} = \sum_{j=1}^n \frac{p(x_j|y_j=i, \theta)}{\sum_{i=1}^K \pi_i p(x_j|y_j=i, \theta)} + \lambda = 0$$

$$\text{两边同乘 } \pi_i \text{ 得 } \sum_{j=1}^n (p_{ji} + \pi_i \lambda) = 0 \quad \sum_{i=1}^K \left[\sum_{j=1}^n (p_{ji} + \pi_i \lambda) \right] = \sum_{i=1}^K \sum_{j=1}^n p_{ji} + \lambda \sum_{i=1}^K \pi_i$$

$$= n + \lambda = 0 \Rightarrow \lambda = -n$$

$$\therefore \sum_{j=1}^n (p_{ji} + \lambda \pi_i) = \sum_{j=1}^n (p_{ji} - n \pi_i) = 0 \Rightarrow \pi_i = \frac{1}{n} \sum_{j=1}^n p_{ji}$$

代码:

(1) 构建 GMM 类, 包含初始化参数、EM 算法、迭代更新

```
1. class GMM:
2.     def __init__(self, k=2):
3.         self.k = k # 定义聚类个数, 默认值为 2
4.         self.p = None # 样本维度
```

```

5.     self.n = None # 样本个数
6.     # 声明变量
7.     self.params = {
8.         "pi": None, # 混合系数 1*k
9.         "μ": None, # 均值 k*p
10.        "cov": None, # 协方差 k*p*p
11.        "pji": None # 后验分布 n*k
12.    }
13.
14. # 初始化参数
15. def init_params(self, init_μ):
16.     pi = np.ones(self.k) / self.k
17.     μ = init_μ
18.     cov = np.ones((self.k, self.p, self.p))
19.     pji = np.zeros((self.n, self.k))
20.     self.params = {
21.         "pi": pi, # 混合系数 1*k
22.         "μ": μ, # 均值 k*p
23.         "cov": cov, # 协方差 k*p*p
24.         "pji": pji # 后验分布 n*k
25.     }
26.
27. # 高斯公式
28. def gaussian_function(self, x_j, μ_k, cov_k):
29.     one = -((x_j - μ_k) @ np.linalg.inv(cov_k) @ (x_j - μ_k).T) / 2
30.     two = -self.p * np.log(2 * np.pi) / 2
31.     three = -np.log(np.linalg.det(cov_k)) / 2
32.     return np.exp(one + two + three)
33.
34. # 计算 Pji 隐变量概率
35. def E_step(self, x):
36.     pi = self.params["pi"]
37.     μ = self.params["μ"]
38.     cov = self.params["cov"]
39.     for j in range(self.n):
40.         x_j = x[j]
41.         pji_list = []
42.         for i in range(self.k):
43.             pi_k = pi[i]
44.             μ_k = μ[i]
45.             cov_k = cov[i]
46.             pji_list.append(pi_k * self.gaussian_function(x_j, μ_k, cov_k))

```



```

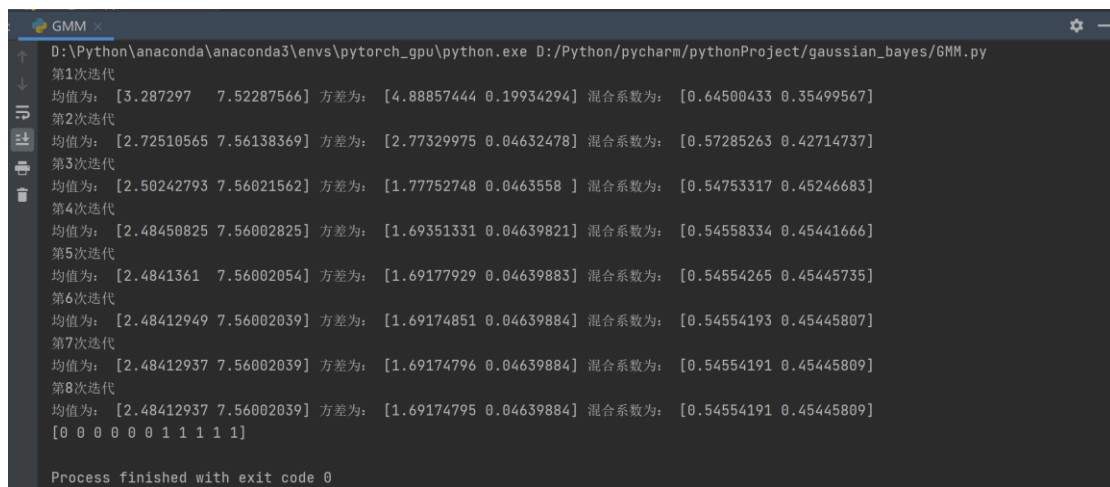
47.         self.params['pji'][j, :] = np.array([v / np.sum(pji_list) for v in
pji_list])
48.
49. # 更新参数
50. def M_step(self, x):
51.      $\mu$  = self.params[" $\mu$ "]
52.     pji = self.params["pji"]
53.     for i in range(self.k):
54.          $\mu_k$  =  $\mu[i]$  # p
55.         pji_k = pji[:, i] # n
56.         pji_k_j_list = []
57.         mu_k_list = []
58.         cov_k_list = []
59.         for j in range(self.n):
60.             x_j = x[j] # p
61.             pji_k_j = pji_k[j]
62.             pji_k_j_list.append(pji_k_j)
63.             mu_k_list.append(pji_k_j * x_j)
64.             self.params[' $\mu$ '][i] = np.sum(mu_k_list, axis=0) / np.sum(pji_k_j_li
st)
65.         for j in range(self.n):
66.             x_j = x[j] # p
67.             pji_k_j = pji_k[j]
68.             cov_k_list.append(pji_k_j * np.dot((x_j -  $\mu_k$ ).T, (x_j -  $\mu_k$ )))
69.         self.params['cov'][i] = np.sum(cov_k_list, axis=0) / np.sum(pji_k_j
_list)
70.         self.params[' $\pi$ '][i] = np.sum(pji_k_j_list) / self.n
71.         print("均值为: ", self.params[" $\mu$ "].T[0], end=" ")
72.         print("方差为: ", self.params["cov"].T[0][0], end=" ")
73.         print("混合系数为: ", self.params[" $\pi$ "])
74.
75. # 迭代, 返回聚类结果
76. def fit(self, x,  $\mu$ , max_iter=10):
77.     x = np.array(x)
78.     self.n, self.p = x.shape
79.     self.init_params( $\mu$ )
80.
81.     for i in range(max_iter):
82.         print("第{}次迭代".format(i+1))
83.         self.E_step(x)
84.         self.M_step(x)
85.     return np.argmax(np.array(self.params["pji"]), axis=1)

```

(2) 主函数调用

```
1. def main():
2.     dataset = np.array([[1.0], [1.3], [2.2], [2.6], [2.8], [5.0], [7.3], [7.4],
3.                           [7.5], [7.7], [7.9]])
4.      $\mu$  = np.array([[6], [7.5]])
5.     my_model = GMM(2)
6.     result = my_model.fit(dataset,  $\mu$ , max_iter=8)
7.     print(result)
```

实验结果：



```
GMM x
D:\Python\anaconda\anaconda3\envs\pytorch_gpu\python.exe D:/Python/pycharm/pythonProject/gaussian_bayes/GMM.py
第1次迭代
均值为: [3.287297  7.52287566] 方差为: [4.88857444 0.19934294] 混合系数为: [0.64500433 0.35499567]
第2次迭代
均值为: [2.72510565 7.56138369] 方差为: [2.77329975 0.04632478] 混合系数为: [0.57285263 0.42714737]
第3次迭代
均值为: [2.50242793 7.56021562] 方差为: [1.77752748 0.0463558 ] 混合系数为: [0.54753317 0.45246683]
第4次迭代
均值为: [2.48450825 7.56002825] 方差为: [1.69351331 0.04639821] 混合系数为: [0.54558334 0.45441666]
第5次迭代
均值为: [2.4841361  7.56002054] 方差为: [1.69177929 0.04639883] 混合系数为: [0.54554265 0.45445735]
第6次迭代
均值为: [2.48412949 7.56002039] 方差为: [1.69174851 0.04639884] 混合系数为: [0.54554193 0.45445807]
第7次迭代
均值为: [2.48412937 7.56002039] 方差为: [1.69174796 0.04639884] 混合系数为: [0.54554191 0.45445809]
第8次迭代
均值为: [2.48412937 7.56002039] 方差为: [1.69174795 0.04639884] 混合系数为: [0.54554191 0.45445809]
[0 0 0 0 0 1 1 1 1]

Process finished with exit code 0
```

可以看出在第6次迭代之后，参数已经趋于稳定。前六个分为一类，对应高斯分布为 $N(2.484, 1.691)$ ，后五个分为一类，对应高斯分布为 $N(7.560, 0.046)$