

# 《大数据机器学习》第 7 次作业

姓名：刘培源

学号：2023214278

**题目 1：**如例 9.1 的三硬币模型。假设观测数据不变，试选择不同的初值，例如， $\pi(\theta) = 0.46$ ， $p(\theta) = 0.55$ ， $q(\theta) = 0.67$ ，求模型参数  $\theta = (\pi, p, q)$  的极大似然估计。

**例 9.1 (三硬币模型)** 假设有 3 枚硬币，分别记作 A, B, C。这些硬币正面出现的概率分别是  $\pi$ ,  $p$  和  $q$ 。进行如下掷硬币试验：先掷硬币 A，根据其结果选出硬币 B 或硬币 C，正面选硬币 B，反面选硬币 C；然后掷选出的硬币，掷硬币的结果，出现正面记作 1，出现反面记作 0；独立地重复  $n$  次试验（这里， $n = 10$ ），观测结果如下：

1, 1, 0, 1, 0, 0, 1, 0, 1, 1

假设只能观测到掷硬币的结果，不能观测掷硬币的过程。问如何估计三硬币正面出现的概率，即三硬币模型的参数。

答：我的代码实现如下：

```
1 class ThreeCoinEM:
2     def __init__(self, prob, tol=1e-5, max_iter=100):
3         # 初始化参数: prob_pi, prob_p, prob_q 分别是三个硬币的概率; tol 是收敛阈值; max_iter 是最大迭代次数
4         self.prob_pi, self.prob_p, self.prob_q = prob
5         self.tol = tol
6         self.max_iter = max_iter
7
8     def calc_mu(self, x):
9         # 计算隐变量的期望, 即E步骤
10        pro_1 = self.prob_pi * self.prob_p ** x * (1 - self.prob_p) ** (1 - x)
11        pro_2 = (1 - self.prob_pi) * self.prob_q ** x * (1 - self.prob_q) ** (1 - x)
12        return pro_1 / (pro_1 + pro_2)
13
14    def fit(self, data):
15        # 对数据进行拟合, 实现EM算法的主要逻辑
16        count = len(data)
17        print("Initial prob:")
18        print(f"prob_pi={self.prob_pi}, prob_p={self.prob_p}, prob_q={self.prob_q}")
19        print("Begin EM:")
20        for i in range(self.max_iter):
21            mu = [self.calc_mu(data[j]) for j in range(count)]
22
23            # 更新概率值, 即M步骤
24            prob_pi = 1 / count * sum(mu)
25            prob_p = sum(m * d for m, d in zip(mu, data)) / sum(mu)
26            prob_q = sum((1 - m) * d for m, d in zip(mu, data)) / sum(1 - m for m in mu)
27            print(f"Iteration {i + 1}: prob_pi={prob_pi:.4f}, prob_p={prob_p:.4f}, prob_q={prob_q:.4f}")
28
29            # 判断是否达到收敛条件
30            error = abs(self.prob_pi - prob_pi) + abs(self.prob_p - prob_p) + abs(self.prob_q - prob_q)
31            self.prob_pi, self.prob_p, self.prob_q = prob_pi, prob_p, prob_q
32
33            if error < self.tol:
34                print("Final prob:")
35                print(f"prob_pi={self.prob_pi:.4f}, prob_p={self.prob_p:.4f}, prob_q={self.prob_q:.4f}")
36                break
37
38    data = [1, 1, 0, 1, 0, 0, 1, 0, 1, 1]
39    init_prob = [0.46, 0.55, 0.67]
40
41    em = ThreeCoinEM(prob=init_prob)
42    em.fit(data)
```

代码的输出如下:

```
1 Initial prob:
2 prob_pi=0.46, prob_p=0.55, prob_q=0.67
3 Begin EM:
4 Iteration 1: prob_pi=0.4619, prob_p=0.5346, prob_q=0.6561
5 Iteration 2: prob_pi=0.4619, prob_p=0.5346, prob_q=0.6561
6 Final prob:
7 prob_pi=0.4619, prob_p=0.5346, prob_q=0.6561
```

所以, 三硬币正面出现的概率分别为  $\pi = 0.4619$ ,  $p = 0.5346$ ,  $q = 0.6561$

**题目 2:** 已知观测数据  $-67, -48, 6, 8, 14, 16, 23, 24, 28, 29, 41, 49, 56, 60, 75$ , 试估计两个分量的高斯混合模型的 5 个参数。

**答:** 我们采用 “sklearn.mixture” 来实现这个功能, 代码如下:

```
1 from sklearn.mixture import GaussianMixture
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 # 初始化观测数据
6 data = np.array([-67, -48, 6, 8, 14, 16, 23, 24,
7                  28, 29, 41, 49, 56, 60, 75]).reshape(-1, 1)
8
9 # 聚类
10 gmmModel = GaussianMixture(n_components=2)
11 gmmModel.fit(data)
12
13 # 预测
14 labels = gmmModel.predict(data)
15
16 # 输出结果
17 print(f"labels={labels}\n \
18       means={gmmModel.means_.reshape(1, -1)}\n \
19       covariances={gmmModel.covariances_.reshape(1, -1)}\n \
20       weights={gmmModel.weights_.reshape(1, -1)}"
21     )
22
23 # 可视化结果
24 plt.scatter(np.arange(len(data)), data, c=labels, s=50)
25 plt.title('Gaussian Mixture Model')
26 plt.show()
```

代码的输入如下：

```
1 labels=[1 1 0 0 0 0 0 0 0 0 0 0 0 0 0]
2 means=[[ 32.98489643 -57.51107027]]
3 covariances=[[429.45764867 90.24987882]]
4 weights=[[0.86682762 0.13317238]]
```

所以两个分量的高斯混合模型的参数如下：

- $\mu_1 \approx 32.98, \mu_2 \approx -57.51$
- $\sigma_1 \approx 429.46, \sigma_2 \approx 90.25$
- $\alpha_1 \approx 0.867, \alpha_2 = 1 - \alpha_1 \approx 0.134$