

# 机器学习第三次

郭英明

2020 年 11 月 22 日

## 1

1.1. 推导三硬币模型的 EM 算法中隐变量后验分布的计算公式以及参数更新公式。

模型中参数为  $\theta = (\pi, p, q)$ ；随机变量  $X, Z \in \{0, 1\}$ ，其中  $X$  表示观测变量结果， $Z$  表示隐变量结果// 隐变量的后验分布：

$$P(x, z|\theta) = \pi^z(1-\pi)^{1-z}p^{zx}(1-p)^{z(1-x)}q^{(1-z)x}(1-q)^{(1-z)(1-x)}$$
$$P(x|\theta) = \pi p^x(1-p)^{1-x} + (1-\pi)q^x(1-q)^{1-x}$$

根据贝叶斯公式：

$$P(z|x, \theta) = \frac{P(x, z|\theta)}{\sum_z P(z|\theta)P(x|z, \theta)}$$

带入公式有：

$$P(z|x, \theta) = \frac{\pi^z(1-\pi)^{1-z}p^{zx}(1-p)^{z(1-x)}q^{(1-z)x}(1-q)^{(1-z)(1-x)}}{\pi p^x(1-p)^{1-x} + (1-\pi)q^x(1-q)^{1-x}}$$

参数更新公式：

EM 算法中 t 步时，

$$Q^{(t)}(z) = P(z|x, \theta^{(t)})$$

M 步中,

$$\begin{aligned}
\theta^{(t+1)} &= \arg \max_{\theta} J(\theta, Q^{(t)}(z)) \\
&= \arg \max_{\theta} \sum_{i=1}^n \sum_z P(z|x, \theta) \log(P(x, z|\theta)) \\
&= \arg \max_{\theta} \sum_{i=1}^n \left( \frac{\pi^{(t)}(p^{(t)x_i}(1-p^{(t)})^{1-x}) \log(\pi p_x(1-p)^{1-x})}{P(x|\theta^{(t)})} \right. \\
&\quad \left. + \frac{(1-\pi^{(t)})(q^{(t)x}(1-q^{(t)})^{1-x}) \log((1-\pi)q^x(1-q)^{1-x})}{P(x|\theta^{(t)})} \right)
\end{aligned}$$

令

$$\begin{aligned}
a^{(t)} &= \frac{\pi^{(t)}(p^{(t)})^x(1-p^{(t)})^{1-x}}{\pi^{(t)}(p^{(t)})^{1-x} + (1-\pi^{(t)})(q^{(t)})^x(1-q^{(t)})^{1-x}} \\
b^{(t)} &= \frac{\pi^{(1-t)}(q^{(t)})^x(1-q^{(t)})^{1-x}}{\pi^{(t)}(p^{(t)})^{1-x} + (1-\pi^{(t)})(q^{(t)})^x(1-q^{(t)})^{1-x}} \\
F(\theta) &= \sum_{i=1}^n a \log(\pi p^x(1-p)^{1-x}) + b \log((1-\pi)q^x(1-q)^{1-x})
\end{aligned}$$

$$\theta^{(t+1)} = \arg \max_{\theta} \sum_{i=1}^n (a^t \log(\pi p^x(1-p)^{1-x}) + b^t \log((1-\pi)q^x(1-q)^{1-x}))$$

易得  $F(\theta)$  是关于  $\pi, p, q$  的凹函数, 对  $\pi, p, q$  进行一阶求导有

$$\begin{aligned}
\pi &= \frac{1}{n} \sum_{i=1}^n a^{(t)} \\
p^{t+1} &= \frac{\sum_{i=1}^n a^{(t)} x}{\sum_{i=1}^n a^{(t)}} \\
q^{t+1} &= \frac{\sum_{i=1}^n b^{(t)} x}{\sum_{i=1}^n b^{(t)}}
\end{aligned}$$

## 1.2. 推导高斯混合模型参数更新公式。

似然函数：

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n P(x|\theta) \\ &= \prod_{i=1}^n \sum_z P(x, z|\theta) \end{aligned}$$

由 Jensen 不等式有

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^n \log \sum_z P(x, z|\theta) \\ &= \sum_{i=1}^n \log \left( \sum_z Q(z) \frac{P(x, z|\theta)}{Q(z)} \right) \\ &\geq \sum_{i=1}^n \sum_z Q(z) \log \left( \frac{P(x, z|\theta)}{Q(z)} \right) \end{aligned}$$

令  $J(\theta, Q(z)) = \sum_{i=1}^n \sum_z Q(z) \log \left( \frac{P(x, z|\theta)}{Q(z)} \right)$ ,  $J(\theta, Q(z))$  为  $Q(z)$  的下界, 则可通过优化  $j(\theta, Q(z))$  优化  $LL(\theta)$

E 步：

$$Q^t(z) = P(z = k|x, \theta^t) = \gamma^t$$

M 步：

$$\begin{aligned} \theta^{\theta+1} &= \arg \max_{\theta} J(\theta, Q^{(t)}(z)) \\ &= \arg \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \gamma^t (\log(P(z = k|\theta)) + \log(x|z = k, \theta)) \\ &= \arg \max_{\theta} \sum_{i=1}^n \sum_{k=1}^K \gamma^t \frac{1}{2} \log \left| \sum_k \right| + \frac{1}{2} (x - u)^T \sum_k^{-1} (x - u) - \log \alpha \end{aligned}$$

令

$$G(\theta) = \gamma^t \frac{1}{2} \log \left| \sum_k \right| + \frac{1}{2} (x - u)^T \sum_k^{-1} (x - u) - \log \alpha$$

$$\frac{\partial G}{\partial u} = 0$$

有

$$u^{t+1} = \frac{\sum_{i=1}^n \gamma^t x}{\sum_{i=1}^n \gamma^t}$$

令

$$\frac{\partial G}{\partial \sum_k} = 0$$

有

$$\sum_k^{t+1} = \frac{\sum_{i=1}^n \gamma^t (x - u^{t+1})(x - u^{t+1})^T}{\sum_{i=1}^n \gamma^t}$$

其中  $u^{t+1} = \frac{\sum_{i=1}^n \gamma^t x}{\sum_{i=1}^n \gamma^t}$  优化问题可以转化为

$$\max_{\alpha_1, \dots, \alpha_k} \sum_{i=1}^n \sum_{k=1}^K \gamma^t \log \alpha_k$$

$$s.t. \sum_{k=1}^K \alpha_k = 1$$

写出拉格朗日函数：

$$\sum_{i=1}^n \sum_{k=1}^K \gamma^t \log \alpha_k - \lambda \sum_{k=1}^K (\alpha_k - 1)$$

对参数求导最后可得：

$$\alpha_k^{t+1} = \frac{1}{n} \sum_{i=1}^n \gamma^t, k = 1, 2, \dots, K$$

## 2 高斯混合模型-EM

### 2.1 自己实现

```
1
2 # -*- coding: utf-8 -*-
3
4 import numpy as np
5 from tqdm import tqdm
6 from sklearn.mixture import GaussianMixture
7
8 # EM算法指定的类别个数
9 class_number = 3
10 #EM算法的迭代次数
11 iterations = 50
12
13 #正态分布的概率密度函数
14 def pdf(x, mu, sigma):
15     x = np.mat(x).reshape(2,1)
16     mu = np.mat(mu).reshape(2,1)
17     sigma = np.mat(sigma).reshape(2,2)
18     coef = 1/(np.sqrt((2*np.pi)**len(x)*np.linalg.det(sigma)))
19     exp = np.exp(-0.5*(x-mu).T*sigma.l*(x-mu))
20     prob = coef*exp
21     return prob[0,0]
22
23 #生成2*2的协方差矩阵
24 def sys_matrix():
25     #主对角线元素
```

```

26     ele = np.random.uniform(5,20,2)
27     #不在对角线上的元素一定比对角线上的小
28     b = np.random.uniform(0,min(ele),1)
29     arr = np.array ([[ ele [0], b [0]], [ b [0], ele [1]]])
30     return arr
31
32     #给定的3组参数
33     mean1 = [3,1]
34     cov1 = [[1,-0.5],[-0.5,1]]
35
36     mean2 = [8,10]
37     cov2 = [[2,0.8],[0.8,2]]
38
39     mean3 = [12,2]
40     cov3 = [[1,0],[0,1]]
41
42     #各自生成300个样本
43     values1 = np.random.multivariate_normal(mean1, cov1, 300)
44     values2 = np.random.multivariate_normal(mean2, cov2, 300)
45     values3 = np.random.multivariate_normal(mean3, cov3, 300)
46     #拼接成900个样本
47     values = np.r_[values1, values2, values3]
48
49
50     #随机生成初始值
51     mean_list = [np.random.uniform(1,3,2) for i in range(class_number)]
52     sigma_list = [sys_matrix() for i in range(class_number)]
53     alpha = np.random.uniform(1,100,class_number).reshape(1,class_number)
54     alpha = alpha/np.sum(alpha)
55
56
57
58     for k in tqdm(range(iterations)):

```

```

59 #计算新类别所属概率
60 '''
61 gamma是一个矩阵，每一行表示特定的样本分别归属某一总体的概率，
62 例如gamma[0,0]表示第一个样本，隶属第一个总体概率，注意此处第一个总体是名义上的第一个总体
63 不是生成样本时的第一个总体
64 '''
65 log_likelihood = 0
66 gamma = np.zeros((1,class_number))
67 #遍历每一个样本，求分别归属每一类的概率
68 for i in range(len(values)):
69     #pdf是上文中定义的求多维高斯分布的概率密度函数的函数
70     temp = np.array([pdf(values[i,:], mean_list[j], sigma_list[j]) for j in range(class_number)]).r
71     #计算 被除数，本质上是两个矩阵相乘（求和后相加等价于矩阵相乘）
72     divdee = np.dot(alpha,temp.T)
73     log_likelihood += np.log(divdee)
74     '''除数 等价于对应元素相乘，使用numpy时，如果数据类型为np.array则*表示对应元素相乘
75     此时如果要做矩阵乘法，则使用np.dot(),如果是np.mat类型，则*表示矩阵乘法
76     '''
77     divider = alpha*temp
78     #对于np.array而言，四则运算都是逐元素执行的，此处的除法也不例外
79     temp = divider/divdee
80     gamma = np.r_[gamma,temp]
81 gamma = gamma[1:,:]
82 #gamma按列求平均,即可得到新的alpha的估计值
83 new_alpha = np.mean(gamma,axis = 0).reshape(1,class_number)
84
85 #计算新均值
86 #计算新均值时，被除数是gamma按列求和后的值
87 divdee = np.sum(gamma,axis = 0).reshape(1,class_number)
88 #计算新均值时，除数 是 gamma的每一列与样本做矩阵乘法（公式中的求和并相加可以转义为矩阵
89 divider =([np.dot(values.T,gamma[:,j].reshape(len(values),1)) for j in range(class_number)])
90 #逐个元素执行除法，并整形成与mean_list相同的形状
91 new_mean = [(divider[i]/ divdee [0, i]).reshape(2) for i in range(class_number)]

```

```

92
93     #计算新方差
94     # dividee = np.sum(gamma,axis = 0).reshape(1,class_number)
95     #将900个样本分别与三个类别的均值相减，并存储在列表内
96     diff = [values - mean_list[j].reshape(1,2) for j in range(class_number)]
97     divider = []
98     #此处是为了求900组差值自身矩阵相乘形成的2*2矩阵
99     for i in range(class_number):
100         #取出diff中的一个元素，注意每个元素是一个900*2的矩阵
101         temp = diff[i]
102         #每个元素自身相乘，形成2*2的矩阵，同时乘以gamma矩阵的对于元素，后续只要直接做求和与
103         temp_matrix = np.array([temp[j,:].reshape(2,1)*temp[j,:].reshape(1,2)*gamma[j,i] for j in range(class_number)])
104         #矩阵求和，并将累和后的矩阵作为元素追加到 divider列表中
105         divider.append(np.sum(temp_matrix,axis = 0))
106     #逐元素做除法，得到新估计的方差
107     new_sigma = [divider[i]/dividee[0,i] for i in range(class_number)]
108     #更新参数
109     mean_list = new_mean
110     alpha = new_alpha
111     sigma_list = new_sigma
112
113     #下面是结果输出的代码
114
115     print('\n自行实现的EM算法: ')
116     print('*'*20)
117     print('pi:',alpha)
118     print('*'*20)
119     print('mu:')
120     for m in mean_list:
121         print(m)
122     print('*'*20)
123     print('sigma:')
124     for s in sigma_list:

```



```

125     print(s)
126
127     print(' '*20)
128     AIC = -2*log_likelihood + 2*(6*class_number-1)
129     BIC = -2*log_likelihood + np.log(len(values))*(6*class_number-1)
130
131     print('AIC: ',AIC[0,0])
132     print('BIC: ',BIC[0,0])
133     print(' '*20)
134     print('\n\n\n')

```

## 2.2 调库实现

```

1 gmm=GaussianMixture(n_components=class_number,covariance_type='full', random_state=0)
2 gmm.fit(values)
3 gmm_mean = gmm.means_
4 gmm_cov= [gmm.covariances_[i] for i in range(class_number)]
5 gmm_alpha = gmm.weights_
6 print('sklearn包中的EM算法')
7 print(' '*20)
8 print('pi: ',gmm_alpha)
9 print(' '*20)
10 print('mu: ')
11 for i in range(class_number):
12     print(gmm_mean[i])
13 print(' '*20)
14 print('sigma:')
15 for s in gmm_cov:
16     print(s)
17 print(' '*20)
18 print('AIC: ',gmm.aic(values))
19 print('BIC',gmm.bic(values))
20 print('\n\n\n')

```

```

21 print('*'*20)
22 print('真实值(顺序未必与预测值对应!!!!)')
23 print('*'*20)
24 print('pi: ')
25 print([1/3,1/3,1/3])
26 print('*'*20)
27 print('mu: ')
28 print(np.array(mean1))
29 print(np.array(mean2))
30 print(np.array(mean3))
31 print('*'*20)
32 print('cov: ')
33 print(np.array(cov1))
34 print(np.array(cov2))
35 print(np.array(cov3))

```

## 2.3 result

```

1 自行实现的EM算法:
2  *****
3  pi: [[0.03627287 0.29639534 0.66733178]]
4  *****
5  mu:
6  [8.58801627 9.83306019]
7  [ 7.92581048 10.08471055]
8  [7.51694744 1.53375755]
9  *****
10 sigma:
11 [[ 1.62382139 -0.57830739]
12  [-0.57830739 0.27434739]]
13 [[2.08002679 0.82778084]
14  [0.82778084 1.64933572]]
15 [[21.15162901 1.91057065]

```

```

16 | [ 1.91057065 1.20876102]]
17 | *****
18 | AIC: 8422.59899397279
19 | BIC: 8504.239704949303
20 | *****
21 |
22 |
23 |
24 |
25 | sklearn包中的EM算法
26 | *****
27 | pi: [0.33333333 0.33333322 0.33333344]
28 | *****
29 | mu:
30 | [3.02835786 1.05542054]
31 | [ 7.99570386 10.04771047]
32 | [12.00688688 2.00465087]
33 | *****
34 | sigma:
35 | [[ 0.94503528 -0.40985998]
36 | [-0.40985998 0.93637324]]
37 | [[2.07140619 0.66802496]
38 | [0.66802496 1.54829122]]
39 | [[ 1.09251519 -0.0214873 ]
40 | [-0.0214873 1.00524051]]
41 | *****
42 | AIC: 7352.490815632364
43 | BIC 7434.131526608877
44 |
45 |
46 |
47 |
48 | *****

```

```

49 真实值(顺序未必与预测值对应!!!)
50 *****
51 pi:
52 [0.3333333333333333, 0.3333333333333333, 0.3333333333333333]
53 *****
54 mu:
55 [3 1]
56 [ 8 10]
57 [12 2]
58 *****
59 cov:
60 [[ 1. -0.5]
61 [-0.5 1. ]]
62 [[2. 0.8]
63 [0.8 2. ]]
64 [[1 0]
65 [0 1]]

```