

混合高斯分布的 EM 算法

陆旭军 ZY2203320
18376486@buaa.edu.cn

摘要

混合高斯分布是一种常用的概率分布模型，它能够很好地描述数据中存在的不同分布的情况。然而，混合高斯分布的参数估计是一个复杂的问题，通常需要使用 EM 算法进行求解。EM 算法是一种迭代算法，其基本思想是通过不断更新模型参数来逐步逼近真实的概率分布。在 EM 算法中，E 步骤是根据当前模型参数对每个数据点进行分类；而 M 步骤则是根据分类结果重新估计模型参数。通过交替进行 E 步骤和 M 步骤，可以不断逼近真实的概率分布，并获得最优的混合高斯分布参数。混合高斯分布的 EM 算法在很多领域中都有广泛的应用，如图像处理、机器学习、自然语言处理等。

引言

1. 一维高斯模型

连续型随机变量 x 的概率密度为

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

正态分布 (Normal distribution) 是一种概率分布。正态分布是具有两个参数 μ 和 σ^2 的连续型随机变量的分布，第一参数 μ 是遵从正态分布的随机变量的均值，第二个参数 σ^2 是此随机变量的方差，所以正态分布记作 $N(\mu, \sigma^2)$ 。遵从正态分布的随机变量的概率规律为取 μ 邻近的值的概率大，而取离 μ 越远的值的概率越小； σ 越小，分布越集中在 μ 附近， σ 越大，分布越分散。

2. 一维高斯混合模型^[1]

高斯混合模型是指具有以下形式的概率分布模型

$$P(y|\mu, \sigma^2) = \sum_{k=1}^K \alpha_k \phi(y|\mu_k, \sigma_k^2)$$

其中， α_k 是系数， $\alpha_k \geq 0$ ， $\sum \alpha_k = 1$ ， ϕ 为高斯分布密度。

在这次男女混合估计中， $K=2$ ，表达式为

$$P(y|\mu, \sigma^2) = \alpha_1 \phi(y|\mu_1, \sigma_1^2) + \alpha_2 \phi(y|\mu_2, \sigma_2^2)$$

3. 高斯混合模型参数估计的 EM 算法

输入: 观测数据 $y_1, y_2, \dots, y_{2000}$ 共 2000 个数据，高斯混合模型 $K=2$

输出: 高斯混合模型参数

(1) 初始值开始迭代

根据男女比 1: 1，和我国 18—44 岁男性和女性平均身高分别为 169.7 厘米和 158.0 厘米，设定初始值如下

$$\alpha_1^{(0)} = 0.5, \alpha_2^{(0)} = 0.5, \mu_1^{(0)} = 169.7, \mu_2^{(0)} = 158, \sigma_1^{2(0)} = \sigma_2^{2(0)} = 1$$

(2) E 步: 依据当前模型参数，计算分模型对观测数据 y 的响应度

$$\hat{\gamma}_{jk} = \frac{\alpha_k \phi(y_j | \mu_k, \sigma_k^2)}{\alpha_1 \phi(y | \mu_1, \sigma_1^2) + \alpha_2 \phi(y | \mu_2, \sigma_2^2)}, j = 1, 2, \dots, 2000; k = 1, 2$$

(3)M 步:计算新一轮迭代的模型参数

$$\hat{\mu}_k = \frac{\sum_{j=1}^{2000} \hat{\gamma}_{jk} y_j}{\sum_{j=1}^{2000} \hat{\gamma}_{jk}}, k = 1, 2$$

$$\hat{\sigma}_k^2 = \frac{\sum_{j=1}^{2000} \hat{\gamma}_{jk} (y_j - \mu_k)^2}{\sum_{j=1}^{2000} \hat{\gamma}_{jk}}, k = 1, 2$$

$$\hat{\alpha}_k = \frac{\sum_{j=1}^{2000} \hat{\gamma}_{jk}}{2000}, k = 1, 2$$

(4)重复第(2)步和第(3)步，直到收敛。

实验过程

1.读取文件

首先使用 pandas 库中的 read_csv 函数读取名为 height_data.csv 的数据文件，并转换为 numpy 数组格式，再删除第一行。这里的数据文件应该是一个包含身高数据的 csv 文件。

2.赋初值

给男女比、男生平均身高、女生平均身高、方差以及精度等参数进行赋值。这些值是 EM 算法的初始化参数，需要根据实际数据情况来设置。将 alpha1 和 alpha2 都赋值为 0.5，表示男女比例各占 50%。将 mu1 和 mu2 分别赋为 169.7 和 158.0，表示男女身高的初始平均值。将 cov_1 和 cov_2 都赋值为 1.0，表示男女身高的方差。将 precision 赋值为 0.00001，表示迭代过程的精度。

3.EM 算法迭代过程

在这个循环中进行 E 步和 M 步的迭代。具体来说，E 步是计算当前模型参数下每个样本属于每个高斯分布的后验概率，即响应度，而 M 步是更新模型参数。在这里的代码实现中，先计算出每个样本对第一个高斯分布和第二个高斯分布的响应度，然后根据这些响应度更新各个参数值。

4.输出结果

最后，将经过迭代后得到的 alpha、mu 和 cov 三个参数输出。其中，alpha 是指男女比，mu 是指平均身高，cov 是方差。

实验结果

	μ_1	μ_2	σ_1^2	σ_2^2	α_1	α_2
真实数据	176	165	25	9	0.75	0.25
EM 算法迭代值	176.15	164.11	24.76	10.09	0.744	0.256
相对误差	0.09%	-0.54%	-0.97%	0.34%	-0.80%	2.4%

根据给出的真实数据和 EM 算法迭代值，可以发现男女身高的平均值和方差都有所偏差，但是偏差很小，尤其是平均值只有小数点后一位的差别。同时，男女的人数比例也有所偏差，但是同样偏差很小，仅仅只有小数点后三位的差别。综合来看，这个 EM 算法的表现还是比较优秀的，能够有效地拟合给定的数据，并估计出其中的分布参数。

可以看到，各个参数的相对误差都比较小，平均相对误差为 0.47%。这说明通过 EM 算法可以较好地估计出数据的分布参数。其中，alpha2 的相对误差比较大，说明在估计男女比例时还有一定的误差。

总结

根据对 EM 算法的计算结果，我们可以看到该算法可以用于混合高斯模型的拟合，从而对观测数据进行聚类 and 分类。在这个例子中，我们使用 EM 算法对一组身高数据进行了男女身高的分类，迭代结果表明该算法能够很好地拟合数据并给出较为准确的分类结果。此外，我们还可以看到，经过多次迭代后，分类的结果已经非常接近真实数据，这也证明了该算法的收敛性和准确性。

参考文献

[1]李航.统计学习方法[M].北京：清华大学出版社,2012:162-165.