

哈尔滨工业大学计算机科学与技术学院

实验报告

课程名称： 机器学习

课程类型： 必修

实验题目： K-means 和 GMM 模型

学号：

姓名：

一、实验目的

实现一个 k-means 算法和混合高斯模型，并且用 EM 算法估计模型中的参数。

二、实验要求及实验环境

2.1 要求：

用高斯分布产生 k 个高斯分布的数据（不同均值和方差）（其中参数自己设定）。

2.2 验证：

- (1) 用 k-means 聚类，测试效果；
- (2) 用混合高斯模型和你实现的 EM 算法估计参数，看看每次迭代后似然值变化情况，考察 EM 算法是否可以获得正确的结果(与你设定的结果比较)。

应用：可以 UCI 上找一个简单问题数据，用你实现的 GMM 进行聚类。

2.3 实验环境：

PyCharm 2020.3.2 x64

Python 3.8 numpy1.20.3 matplotlib 3.4.2

三、设计思想（本程序中的用到的主要算法及数据结构）

3.1 K-means 理论推导：

假设数据集为 $\{x_1, \dots, x_N\}$ ，且维度为 D，我们的目标是将数据集划分为 K 个类，假设 K 已给定。

直观上，认为一组数据点构成的一个聚类中，聚类内部的点应小于数据点与外部的点之间的距离，假设 μ_k 为聚类的中心。

定义目标函数，该函数表示每个数据点和它被分到的类的中心 μ_k 的距离的平方和：

$$J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} \|x_n - \mu_k\|^2 \quad (1)$$

当且仅当 x_n 被分到类别 k 的时候， $r_{nk} = 1$ 。

使用迭代的方式，使目标函数 J 达到最小值，分别对应于 r_{nk} 和 μ_k 的最优化。在第一阶段中，固定 μ_k ，关于 r_{nk} 最小化 J；第二阶段中，固定 r_{nk} ，关于 μ_k 最小化 J。不断重复上述两个过程，直到收敛。这两步分别对应于 EM 算法中的 E(期望) 和 M(最大化) 的步骤。

第一阶段:

由于不同的与 n 相关的项是独立的, 因此可以分别对 n 项进行最优化, 如下式:

$$r_{nk} = \begin{cases} 1 & \text{若 } k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & \text{其他情况} \end{cases} \quad (2)$$

第二阶段:

当 r_{nk} 固定时, 关于 μ_k 的最优化。由于 J 是关于 μ_k 的二次函数, 令它关于 μ_k 的导数为0, 即可达到最小值, 因此对公式(1)关于 μ_k 求导, 可得:

$$\frac{\partial}{\partial \mu_k} J = 2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) \quad (3)$$

令公式(3)为0, 可得:

$$2 \sum_{n=1}^N r_{nk} (x_n - \mu_k) = 0 \quad (4)$$

化简可得:

$$\mu_k = \frac{\sum_{n=1}^N r_{nk} x_n}{\sum_{n=1}^N r_{nk}} \quad (5)$$

公式(5)中, 分母等于聚类 k 中数据点的数量, 因此公式(5)的含义为, μ_k 等于类别 k 的所有数据点的均值。

重复为数据分类, 重新计算聚类均值的步骤重复进行, 直到聚类的分配不改变或直到迭代次数超过某个最大值。

K - means算法的收敛性: 每个阶段都减小了目标函数 J 的值。

3.2 高斯混合模型与EM算法理论推导:

每个高斯分布对应 $N(\mu_k, \Sigma_k)$, 因此高斯混合概率分布可以写成 K 个多维高斯分布的线性叠加的形式, 即:

$$P(x) = \sum_{k=1}^K \pi_k N(\mu_k, \Sigma_k) \quad (6)$$

且有

$$0 \leq \pi_k \leq 1 \quad (7)$$

$$\sum_{k=1}^K \pi_k = 1 \quad (8)$$

引入 K 维二值隐变量 z , 且满足以下条件:

$$P(z_k = 1) = \pi_k \quad (9)$$

$$0 \leq \pi_k \leq 1 \quad (10)$$

因此概率分布也可以写为

$$P(z) = \prod_{k=1}^K \pi_k^{z_k} \quad (11)$$

此时 x 所对应 z 的后验概率为：

$$P(x|z) = \prod_{k=1}^K N(x|\mu_k, \Sigma_k)^{z_k} \quad (12)$$

当给定 z_k 特定的值时， x 的条件概率分布是一个高斯分布：

$$P(x|z_k = 1) = N(x|\mu_k, \Sigma_k) \quad (13)$$

因此可以从 x 的联合概率分布求得 x 的联合概率分布，即

$$P(x) = \sum_z P(z)P(x|z) \quad (14)$$

$$P(x) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k) \quad (15)$$

此时得出的公式(15)与公式(6)完全一致。

假设数据集为 $\{x_1, \dots, x_N\}$ ，对于每一个观测数据点 x_n 都对应一个 z_n ，并且样本 x_i 为 d 维向量，且由 d 维高斯分布生成。

因此有

$$N(x_i|\mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)} \quad (16)$$

在给定观测值 x_i 的时候，将 $\gamma(z_k)$ 看作 x_i 对应的后验概率，它的值可以通过贝叶斯定理得出：

$$\gamma(z_k) = p(z_k = 1|x) = \frac{P(z_k = 1)P(x|z_k = 1)}{\sum_{j=1}^K P(z_j = 1)P(x|z_j = 1)} \quad (17)$$

$$= \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_k)} \quad (18)$$

因此对于一个样本点 x_i ，若要将其归为 $1, \dots, K$ 簇中的一个，需要找到令公式(18)最大的 k ，即所归的簇的序号，可以表示为：

$$k = \arg \max_k p(z_k = 1|x_i) \quad (19)$$

$$k = \arg \max_k \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x|\mu_j, \Sigma_k)} \quad (20)$$

因此，找到参数 π_k 、 μ_k 、 Σ_k 是GMM聚类的关键。

将数据集表示为一个 $N \times d$ 的矩阵 X ，其中第 n 行为 x_n^T ，类似的对应的隐含变量会被表示为一个 $N \times K$ 的矩阵 Z ，它的行为 z_n^T 。

由公式(6)可得，对数似然函数为：

$$L(X) = \ln P(X|\pi, \mu, \Sigma) \quad (21)$$

$$= \ln \prod_{n=1}^N \left(\sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k) \right) \quad (22)$$

$$= \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) \quad (23)$$

公式(23)对 $\boldsymbol{\mu}_k$ 求导:

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} L(X) = \sum_{n=1}^N \frac{\frac{\partial}{\partial \boldsymbol{\mu}_k} \sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad (24)$$

$$= \sum_{n=1}^N \frac{\frac{\partial}{\partial \boldsymbol{\mu}_k} (\pi_1 N(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \cdots + \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) + \cdots + \pi_K N(\mathbf{x}_n | \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K))}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad (25)$$

$$= \sum_{n=1}^N \frac{\frac{\partial}{\partial \boldsymbol{\mu}_k} \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad (26)$$

多维高斯分布对均值 $\boldsymbol{\mu}$ 求导:

$$N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \quad (27)$$

$$\frac{\partial N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = \frac{\partial}{\partial \boldsymbol{\mu}} \left(\frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} \right) \quad (28)$$

$$= \frac{\partial}{\partial \boldsymbol{\mu}} \left(\frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu})} \right) \quad (29)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} * \left(-\frac{1}{2} \right) \left(-(\mathbf{x}^T \boldsymbol{\Sigma}^{-1})^T - \boldsymbol{\Sigma}^{-1} \mathbf{x} + (\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1T}) \boldsymbol{\mu} \right) \quad (30)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}}} \frac{1}{|\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})} * \left(-\frac{1}{2} \right) \left(-(\boldsymbol{\Sigma}^{-1T} + \boldsymbol{\Sigma}^{-1}) + (\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}^{-1T}) \boldsymbol{\mu} \right) \quad (31)$$

由于各个维度相关系数为 0, 因此有 $\boldsymbol{\Sigma}$ 为对称矩阵, 故有:

$$\frac{\partial N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}} = N(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \quad (32)$$

将结果代入公式(26), 因此我们有:

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} L(X) = \sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}, \boldsymbol{\Sigma})}{\sum_{j=1}^K \pi_j N(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (33)$$

由公式(18), 我们有

$$\frac{\partial}{\partial \boldsymbol{\mu}_k} L(X) = \sum_{n=1}^N \gamma(z_k) \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (34)$$

令公式(34)为 0, 可得

$$\sum_{n=1}^N \gamma(z_k) \boldsymbol{\Sigma}_k^{-1}(\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (35)$$

因此，我们有

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (36)$$

令公式(23)对 $\boldsymbol{\Sigma}_k$ 并令其为0可得：

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} L(X) = 0 \quad (37)$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (38)$$

由公式(8)可得，当前优化为有约束的优化，因此使用拉格朗日乘数法：

$$L(X) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \quad (39)$$

使用公式(39)对 π_k 求导可得，

$$\frac{\partial}{\partial \boldsymbol{\Sigma}} L(X) = \sum_{n=1}^N \frac{N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} + \lambda \quad (40)$$

令公式(40)为0，有：

$$\sum_{n=1}^N \frac{N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} + \lambda = 0 \quad (41)$$

对公式(42)同乘 π_k ，可得

$$\pi_k \sum_{n=1}^N \frac{N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} + \pi_k \lambda = 0 \quad (42)$$

将公式(18)代入得

$$\sum_{n=1}^N \gamma(z_{nk}) + \pi_k \lambda = 0 \quad (43)$$

将K个等式相加，可得

$$\sum_{k=1}^K \left(\pi_k \sum_{n=1}^N \frac{N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \right) + \sum_{k=1}^K \pi_k \lambda = 0 \quad (44)$$

$$\sum_{k=1}^K \left(\sum_{n=1}^N \frac{\pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \right) + \sum_{k=1}^K \pi_k \lambda = 0 \quad (45)$$

将累加符号调换位置

$$\sum_{k=1}^N \left(\frac{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \right) + \sum_{k=1}^K \pi_k \lambda = 0 \quad (46)$$

可得

$$\lambda = -N \quad (47)$$

因此有

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nk}) \quad (48)$$

EM算法的完整描述：

E步骤：使用当前的参数计算每个样本 x_i 的 $\gamma(z_{ik})$ 。

$$\gamma(z_{ik}) = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_i | \mu_j, \Sigma_j)} \quad (18)$$

M步骤：使用当前的 $\gamma(z_{ik})$ 重新估计参数。

$$\mu_k^{\text{新}} = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (49)$$

$$\Sigma_k^{\text{新}} = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k^{\text{新}})(\mathbf{x}_n - \mu_k^{\text{新}})^T}{\sum_{n=1}^N \gamma(z_{nk})} \quad (50)$$

$$\pi_k^{\text{新}} = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nk}) \quad (51)$$

计算对数似然函数

$$\ln P(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(\mathbf{x}_n | \mu_k, \Sigma_k) \right) \quad (52)$$

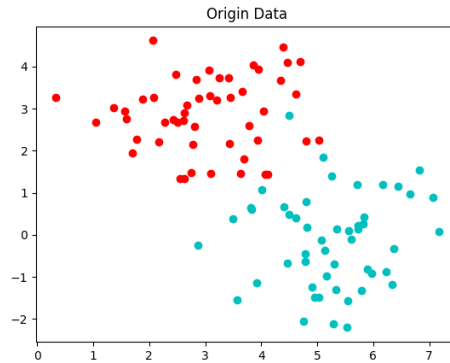
若对数函数没有满足收敛的准则，则继续重复E、M步骤。

四、实验结果与分析

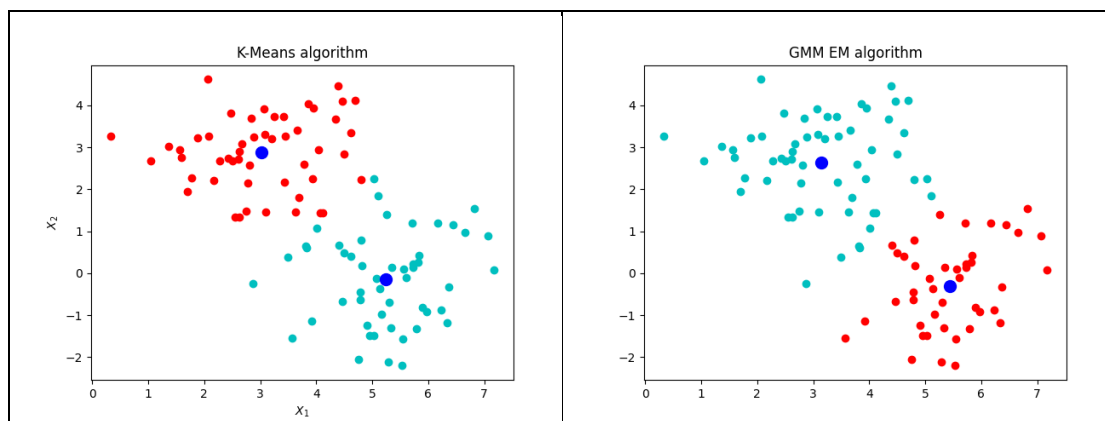
4.1 生成数据实验：

实验设定生成的数据符合二维高斯分布，二维方差均为 0.1。

(1) 生成的数据的 $K = 2$ 时。

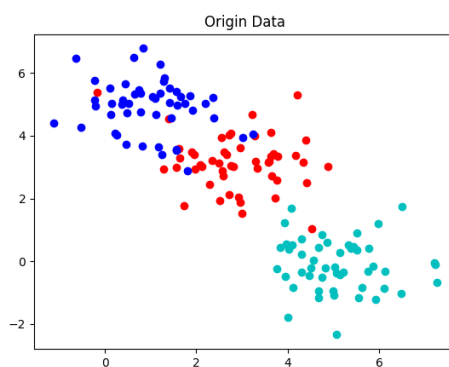


如图为生成的原始数据，两堆数据分别符合两个二维高斯分布。

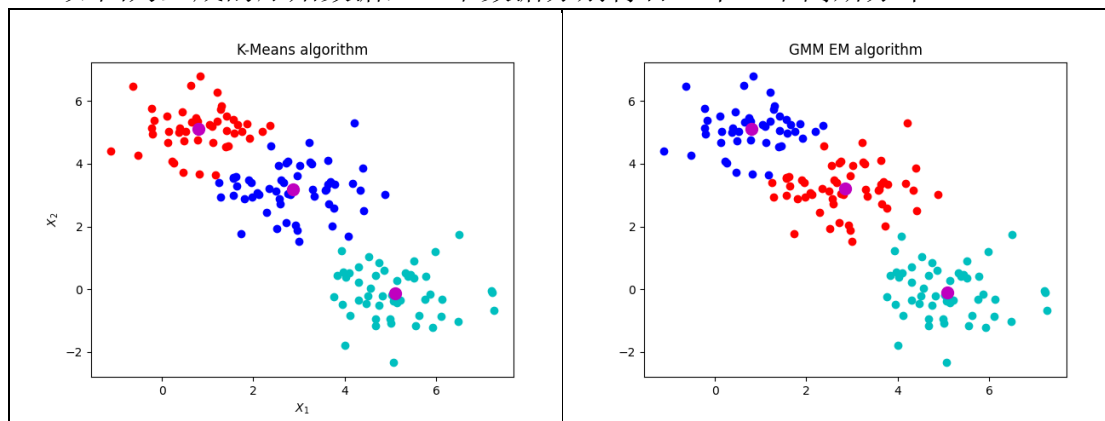


左图为使用K – Means算法后分类的结果，右图为使用GMM EM算法后分类的结果。可以看到当 $K = 2$ 的时候，两种算法都有较好的分类效果。

(2) 生成的数据的 $K = 3$ 时。



如图为生成的原始数据，三堆数据分别符合三个二维高斯分布。



左图为使用K – Means算法后分类的结果，右图为使用GMM EM算法后分类的结果。可以看到当 $K = 3$ 的时候，两种算法都有较好的分类效果。和原图相比，两个分类结果都有比较好的展示，没有出现拟合的情况。

4.2 UCI 数据实验：

本阶段采用 UCI 网站的鸢尾花数据集。

由于 K-Means 和 GMM EM 为无监督机器学习，因此需要先定义好簇和其所属类别的映射关系。

对簇的类别进行排列组合，由于我们已经有 K 的值，因此可以很清晰的定义排列组合的次数为 $K!$ 次。

我们假设 $K = 3$ ，因此此时有 $[[0,1,2],[0,2,1],[1,2,0],[1,0,2],[2,0,1],[2,1,0]]$ 六种排列组合。

接下来定义一个字典的列表，将上述组合中的每一个分别用第 i 个元素和 i 对应起来，就得到了簇和 Y 中的种类的关系的映射。

以 $[0,2,1]$ 为例，所得到的字典为

$\text{dict} = \{0: 1, 2: 1, 1: 2\}$

通过这样的映射关系，我们遍历训练集中所有得到的簇的序号到数据集中的预测数据之间的关系，并进行比较，记录下所有匹配的个数。

将所有字典循环后，我们即可以得到使训练好的模型和原数据集最匹配的映射关系。选取这样的映射关系，我们对测试集进行计算，当测试集预测出的簇的序号映射到的 Y 中的序号与其本来 Y 的序号对应，则视为预测正确，否则视为预测错误。

这样我们就有了对数据集正确率的说明。

```
D:\Python\Python38\python.exe D:/Code/PycharmProject
K-Means 迭代次数为8
cost = 52.63224678575162
costAfter = 52.63224678575162
K-Means迭代完成
K-means iris数据集的正确率为0.9736842105263158
GMM EM 迭代次数为26
cost = -1849.7860251837722
costAfter = -1849.7860153789184
GMM EM迭代完成
GMM EM iris数据集的正确率为0.9736842105263158
```

实际实验中， $K - \text{Means}$ 和 GMM EM 算法都能够较好的对鸢尾花数据集分类。 $K - \text{Means}$ 迭代8次， GMM EM 算法迭代26次。测试集共38个数据， $K - \text{Means}$ 和 GMM EM 的正确分类37个数据，正确率为97.368%，有较好的预测表现。

五、结论

$K - \text{Means}$ 和 GMM EM 都有较好的表现，可以快速收敛，并且二者的分类性能都较好。但是 $K - \text{Means}$ 有可能找到的是局部最优的聚类，而不是全局最优的聚类。并且在样本维度过高时，可能耗费大量的计算。

高斯混合模型不仅可以给出一个样本属于某类的概率，还可以用于概率密度的估计，用于生成新的样本点。

六、参考文献

- [1] Pattern Recognition and Machine Learning
- [2] 机器学习 周志华著 北京：清华大学出版社，2016 年 1 月。

七、附录：源代码（带注释）

main.py

```
1. from KMeans import *
2. from generateData import *
3. from GMM import *
4. from readFromFile import *
5.
6. if __name__ == '__main__':
7.     # # 类别个数
8.     # K = 3
9.     # # 迭代结束次数
10.    # iter = 200
11.    # # 生成数据
12.    # X1, X2 = generateData(K)
13.    # KMeansClusterData = makeDataForKmeans(X1, X2)
14.    # GMMClusterData = makeDataForGMM(X1, X2)
15.    # # KMeans 算法
16.    # DataCenter, centers = kMeans(K, iter, KMeansClusterData)
17.    # drawKmeansClassify(KMeansClusterData, DataCenter, centers, K)
18.    # # GMM EM 算法
19.    # Gamma, mus, sigmas, alpha = GMMem(K, iter, GMMClusterData)
20.    # drawGMM(GMMClusterData, Gamma, mus)
21.
22.
23.
24.    iterFile = 200
25.    # 从文件中读取数据
26.    X, Y, K = readFromFile("./DataSet/iris.csv")
27.    KMeansTrain, KMeansTest, KMeansYTrain, KMeansYtest = makeFileDataForKmeans(X, Y)
28.    GMMTrain, GMMTest, GMMYTrain, GMMYTest = makeFileDataForGMM(X, Y)
29.
30.    # KMeans 算法
31.    DataCenterFile, centersFile = kMeans(K, iterFile, KMeansTrain)
32.    TestKMeansAccuracy(DataCenterFile, centersFile, KMeansYTrain, KMeansYtest, Y, KMeansTest, K)
```

```

33.
34.
35.     # GMM EM 算法
36.     GammaFile, mus, sigmas, alpha = GMMem(K, iterFile, GMMTrain)
37.     centersTrain = np.zeros(GammaFile.shape[0])
38.     for i in range(GammaFile.shape[0]):
39.         centersTrain[i] = np.array(GammaFile[i]).argmax()
40.     TestGMMAccuracy(centersTrain, mus, sigmas, alpha, GMMYTrain, GMMYTest, Y, GMMTest, K)

```

generateData.py

```

1. import numpy as np
2. import matplotlib.pyplot as plt
3. import random
4. from sklearn.model_selection import train_test_split
5.
6. color = ['c', 'r', 'b', 'm', 'y', 'k', 'g', 'w']
7. muList = [[5, 0], [3, 3], [1, 5], [4, 2], [1, 3], [2, 3]]
8.
9.
10. def generate2DimensionalData(Mu1, Mu2, cov11, cov12, cov21, cov22, Num, noiseSigma1, noiseSigma2):
11.     """
12.     生成二维高斯分布数据
13.     :param Mu1: 维度 1 的平均值
14.     :param Mu2: 维度 2 的平均值
15.     :param cov11: 维度 1 的方差
16.     :param cov12: 维度 12 的协方差
17.     :param cov21: 维度 21 的协方差
18.     :param cov22: 维度 2 的方差
19.     :param Num: 生成样本点的个数
20.     :param noiseSigma1: 维度 1 噪声的方差
21.     :param noiseSigma2: 维度 2 噪声的方差
22.     :return: 分别返回两个维度的样本点
23.     """
24.     mean = np.array([Mu1, Mu2])
25.     cov = np.array([[cov11, cov12], [cov21, cov22]])
26.     Data = np.random.multivariate_normal(mean, cov, Num)
27.     X1 = Data[:, 0]
28.     X2 = Data[:, 1]
29.     # 添加高斯噪声
30.     GuassNoise1 = np.random.normal(0, scale=noiseSigma1, size=Num)
31.     GuassNoise2 = np.random.normal(0, scale=noiseSigma2, size=Num)

```

```

32.     X1 = X1 + GuassNoise1
33.     X2 = X2 + GuassNoise2
34.     return X1, X2
35.
36.
37. def generateData(K):
38.     """
39.     生成 K 堆二维高斯分布数据
40.     :param K: 生成样本点的堆数
41.     :return: X1、X2 的列表
42.     """
43.     X1 = np.empty(0)
44.     X2 = np.empty(0)
45.     plt.title("Origin Data")
46.     number = 50
47.     # Y = np.empty(0)
48.     for i in range(K):
49.         # X1gen, X2gen = generate2DimensionalData(random.randint(0, 5), random.randint(0, 5), 1, 0, 0, 1, number, 0.1, 0.1)
50.         X1gen, X2gen = generate2DimensionalData(muList[i][0], muList[i][1], 1, 0, 0, 1, number, 0.1, 0.1)
51.         # Y = np.hstack((Y, np.ones(number)*i))
52.         X1 = np.hstack((X1, X1gen))
53.         X2 = np.hstack((X2, X2gen))
54.         plt.scatter(X1gen, X2gen, c=color[i])
55.     plt.show()
56.     return X1, X2
57.
58.
59. def makeDataForKmeans(X1, X2):
60.     """
61.     将数据处理, 适合 KMeans 算法
62.     :param X1: 一行数据
63.     :param X2: 一行数据
64.     :return: 两行数据
65.     """
66.     X = np.vstack((X1, X2))
67.     return X
68.
69.
70. def KmeansTrainTestSplit(ClusterData, Y):
71.     clusterData = np.transpose(ClusterData)
72.     X_train, X_test, Y_train, Y_test = train_test_split(clusterData, Y)
73.     X_train = np.transpose(X_train)

```

```

74.     X_test = np.transpose(X_test)
75.     return X_train, X_test, Y_train, Y_test
76.
77.
78. def makeDataForGMM(X1, X2):
79.     X1 = np.reshape(X1, (-1, 1))
80.     X2 = np.reshape(X2, (-1, 1))
81.     X = np.column_stack((X1, X2))
82.     return X

```

KMeans.py

```

1. import numpy as np
2. import matplotlib.pyplot as plt
3. import math
4. from generateData import color
5. import itertools
6.
7. def initCenters(KMeansData, K):
8.     """
9.     初始化中心点
10.    :param KMeansData: 所有数据点
11.    :param K: 分成类别的个数
12.    :return:
13.    """
14.    centerX1 = np.random.choice(KMeansData[0], K)
15.    centerX2 = np.random.choice(KMeansData[1], K)
16.    center = np.vstack((centerX1, centerX2))
17.    return center
18.
19.
20. def kMeansCost(KMeansData, centers, DataCenter):
21.     """
22.     计算当前所有 KMeansData 中的所有点与 centers 所有点的代价函数的值
23.    :param KMeansData: 所有 KMeansData 的数据点
24.    :param centers: 中心点
25.    :param DataCenter: 每个点对应的中心点
26.    :return:
27.    """
28.    cost = 0
29.    for i in range(KMeansData.shape[1]):
30.        X1 = KMeansData[0][i]
31.        X2 = KMeansData[1][i]
32.        centerX1 = centers[0][int(DataCenter[i])]

```

```

33.         centerX2 = centers[1][int(DataCenter[i])]
34.         cost = cost + math.sqrt((X1 - centerX1) * (X1 - centerX1) + (X2 - ce
        nterX2) * (X2 - centerX2))
35.     return cost
36.
37.
38. def calculateDistanceCenter(X1, X2, center):
39.     """
40.     计算当前的点与中心点的距离，选择最近的类
41.     :param X1: 点的坐标
42.     :param X2: 点的坐标
43.     :param center: 中心点数组
44.     :return: 点的类
45.     """
46.     cost = float("inf")
47.     minCost = cost
48.     Xclass = 0
49.     for i in range(center.shape[1]):
50.         cost = 0
51.         cost = cost + math.sqrt((X1 - center[0][i]) * (X1 - center[0][i]))
52.         cost = cost + math.sqrt((X2 - center[1][i]) * (X2 - center[1][i]))
53.         if cost < minCost:
54.             minCost = cost
55.             Xclass = i
56.     return Xclass
57.
58.
59. def assignClass(KMeansData, Centers):
60.     """
61.     为 KMeansData 每个数据分配中信
62.     :param KMeansData: 数据
63.     :param Centers: 中心点
64.     :return: 为每个点分配离它最近的中心点的序号
65.     """
66.     center = np.zeros(KMeansData.shape[1])
67.     for i in range(KMeansData.shape[1]):
68.         center[i] = calculateDistanceCenter(KMeansData[0][i], KMeansData[1][
            i], Centers)
69.     return center
70.
71.
72. def updateCenters(KMeansData, centers, K):
73.     """
74.     更新中心点 K 个中心点

```

```

75.     :param KMeansData:数据
76.     :param centers:每个数据的中心点
77.     :param K:类别个数
78.     :return:更新新的中心点
79.     """
80.     UpdateCenters = np.zeros([2, K])
81.     list1 = []
82.     list2 = []
83.     for i in range(K):
84.         for j in range(KMeansData.shape[1]):
85.             if centers[j] == i:
86.                 list1.append(KMeansData[0][j])
87.                 list2.append(KMeansData[1][j])
88.             x1 = 0
89.             x2 = 0
90.             for k in range(len(list1)):
91.                 x1 = x1 + list1[k]
92.                 x2 = x2 + list2[k]
93.
94.             if x1 == 0 and x2 == 0:
95.                 print("此时无点离他最近")
96.                 UpdateCenters[0][i] = 0
97.                 UpdateCenters[1][i] = 0
98.                 # randomNum = random.randint(0, KMeansData.shape[1]-1)
99.                 # UpdateCenters[0][i] = KMeansData[0][randomNum]
100.                # UpdateCenters[1][i] = KMeansData[1][randomNum]
101.            else:
102.                x1 = x1 / len(list1)
103.                x2 = x2 / len(list2)
104.                UpdateCenters[0][i] = x1
105.                UpdateCenters[1][i] = x2
106.            list1 = []
107.            list2 = []
108.        return UpdateCenters
109.
110.
111. def kMeans(K, iter, ClusterData):
112.     """
113.     调用 Kmeans 算法
114.     :param K: 类别个数
115.     :param iter: 迭代次数
116.     :param ClusterData: 实验数据
117.     :return:
118.     """

```

```

119.     # 初始化开始时候的 K 个中心点
120.     centers = initCenters(ClusterData, K)
121.     costAfter = 0
122.     DataCenter = np.zeros(ClusterData.shape[1])
123.     for i in range(iter):
124.         cost = costAfter
125.         # 为每个点分配中心
126.         DataCenter = assignClass(ClusterData, centers)
127.         # 更新中心点
128.         centers = updateCenters(ClusterData, DataCenter, K)
129.         costAfter = kMeansCost(ClusterData, centers, DataCenter)
130.         if (costAfter == cost):
131.             print("K-Means 迭代次数为{}".format(i))
132.             print("cost = {}".format(cost))
133.             print("costAfter = {}".format(costAfter))
134.             print("K-Means 迭代完成")
135.             # drawKmeansClassify(ClusterData, DataCenter, centers,K)
136.             break
137.     return DataCenter, centers
138.
139.
140. def drawKmeansClassify(KMeansData, DataCenter, centers, K):
141.     """
142.     为 K-Means 算法生成的结果画图
143.     :param KMeansData:KMeans 的数据
144.     :param DataCenter:每个点对应的中心
145.     :param K:类别的个数
146.     """
147.     plt.title("K-Means algorithm")
148.     for i in range(K):
149.         listX1 = []
150.         listX2 = []
151.         for j in range(KMeansData.shape[1]):
152.             if DataCenter[j] == i:
153.                 listX1.append(KMeansData[0][j])
154.                 listX2.append(KMeansData[1][j])
155.             plt.scatter(listX1, listX2, c=color[i])
156.     plt.scatter(centers[0], centers[1], c=color[K], s=100)
157.     plt.xlabel("$X_{1}$")
158.     plt.ylabel("$X_{2}$")
159.     plt.show()
160.
161.

```



```

162. def TestKMeansAccuracy(centersTrain, centers, YTrain, YTest,Y,KMeansTest, K
    ):
163.
164.     # 簇号
165.     Ytest = np.zeros(KMeansTest.shape[1])
166.     for i in range(KMeansTest.shape[1]):
167.         Ytest[i] = calculateDistanceCenter(KMeansTest[0][i], KMeansTest[1][
            i], centers)
168.     # Dict 簇到原来 y 的映射
169.     Dict = KMeansClusterMap(K, YTrain,centersTrain,Y)
170.     sum = 0
171.     for i in range(YTest.shape[0]):
172.         Ytest[i] = Dict[Ytest[i]]
173.     for i in range(YTest.shape[0]):
174.         sum = sum +1
175.     Accuracy = sum/Ytest.shape[0]
176.     print("K-means iris 数据集的正确率为{}".format(Accuracy))
177.
178.
179. def KMeansClusterMap(K,YTrain, dataCenters,Y):
180.     Set = list(set(Y))
181.     List = []
182.     for i in range(K):
183.         List.append(i)
184.     listArray = list(itertools.permutations(List,K))
185.     MapAccuracyDict = np.zeros(len(listArray))
186.     listDict = []
187.     for j in range(len(listArray)):
188.         # 簇号
189.         Dict = {}
190.         for k in range(len(listArray[j])):
191.             Dict[listArray[j][k]] = Set[k]
192.         listDict.append(Dict)
193.         MapAccuracyDict[j] = KMeansMapAccuracy(Dict,YTrain,dataCenters)
194.     return listDict[np.array(MapAccuracyDict).argmax()]
195.
196.
197.
198. def KMeansMapAccuracy(Dict,YTrain,dataCenters):
199.     sum = 0
200.     for i in range(YTrain.shape[0]):
201.         if Dict[dataCenters[i]] == YTrain[i]:
202.             sum = sum+1
203.     return sum/YTrain.shape[0]

```

GMM.py

```
1. import numpy as np
2. import random
3. import matplotlib.pyplot as plt
4. import math
5. from scipy.stats import multivariate_normal
6. from generateData import color
7. import itertools
8. def GMMem(K, iter, ClusterData):
9.     mus, sigmas, alpha = initParams(ClusterData, K)
10.    costAfter = 0
11.    gamma = 0
12.    for i in range(iter):
13.        cost = costAfter
14.        gamma = getExpectation(ClusterData, mus, sigmas, alpha, K)
15.        mus, sigmas, alpha = maximize(ClusterData, gamma, K)
16.        costAfter = calculateMLE(ClusterData, mus, sigmas, alpha, K)
17.        if abs(costAfter - cost) <= 1e-5:
18.            print("GMM EM 迭代次数为{}".format(i))
19.            print("cost = {}".format(cost))
20.            print("costAfter = {}".format(costAfter))
21.            print("GMM EM 迭代完成")
22.            break
23.    return gamma, mus, sigmas, alpha
24.
25.
26. def calculateMLE(ClusterData, mus, sigmas, alpha, K):
27.     """
28.     计算似然函数
29.     :param ClusterData: 所有数据点
30.     :param mus: 多维高斯分布均值矩阵
31.     :param sigmas: 多维高斯分布的协方差矩阵
32.     :param alpha: 每个簇的权重
33.     :param K: 簇的个数
34.     :return: MLE 的值
35.     """
36.     iterN = 0
37.     for i in range(ClusterData.shape[0]):
38.         iterK = 0
39.         for k in range(K):
40.             iterK = iterK + alpha[k] * multivariate_normal.pdf(ClusterData[i], mus[k], sigmas[k])
41.         iterN = iterN + math.log(iterK, math.e)
```

```

42.     return iterN
43.
44.
45. def getExpectation(ClusterData, mus, sigmas, alpha, K):
46.     """
47.     返回 gamma 矩阵
48.     :param ClusterData:所有数据点
49.     :param mus:多维高斯分布均值矩阵
50.     :param sigmas:多维高斯分布的协方差矩阵
51.     :param alpha:每个簇的权重
52.     :param K:簇的个数
53.     :return:gamma 矩阵
54.     """
55.     gamma = np.zeros((ClusterData.shape[0], K))
56.     for i in range(ClusterData.shape[0]):
57.         gamma_sum = 0
58.         for j in range(K):
59.             gamma[i][j] = alpha[j] * multivariate_normal.pdf(ClusterData[i],
mus[j], sigmas[j])
60.             gamma_sum = gamma_sum + gamma[i][j]
61.         for j in range(K):
62.             gamma[i][j] = gamma[i][j] / gamma_sum
63.     return gamma
64.
65.
66.
67. def maximize(ClusterData, gamma, K):
68.     """
69.     M Step
70.     :param ClusterData:所有数据点
71.     :param gamma:gamma 矩阵
72.     :param K:所有的类别
73.     :return:m step 后的参数
74.     """
75.     mus = np.zeros((K,ClusterData.shape[1]))
76.     for i in range(K):
77.         iterNumerator = np.zeros(ClusterData.shape[1])
78.         iterDenominator = 0
79.         for j in range(ClusterData.shape[0]):
80.             iterNumerator = iterNumerator + gamma[j][i]*ClusterData[j]
81.             iterDenominator = iterDenominator+gamma[j][i]
82.         for k in range(ClusterData.shape[1]):
83.             # 存在除 0 可能
84.             mus[i][k] = iterNumerator[k] / iterDenominator

```

```

85.
86.     sigmas = np.zeros((K, ClusterData.shape[1], ClusterData.shape[1]))
87.     for k in range(K):
88.         iterNumerator = np.zeros((ClusterData.shape[1], ClusterData.shape[1])
89.         )
90.         iterDenominator = 0
91.         for i in range(ClusterData.shape[0]):
92.             iterNumerator = iterNumerator + gamma[i][k] * np.dot(np.transpose
93.             (ClusterData[i] - mus[k]), (ClusterData[i] - mus[k]))
94.             iterDenominator = iterDenominator + gamma[i][k]
95.             sigma = iterNumerator / iterDenominator
96.             for j in range(ClusterData.shape[1]):
97.                 sigma[j][j] = sigma[j][j]*1.01
98.             sigmas[k] = sigma
99.
100.     alpha = np.zeros(K)
101.     for i in range(K):
102.         iterNumerator = 0
103.         for j in range(ClusterData.shape[0]):
104.             iterNumerator=iterNumerator+gamma[j][i]
105.         alpha[i] = iterNumerator/ClusterData.shape[0]
106.     return mus, sigmas, alpha
107.
108.
109. def initParams(ClusterData, K):
110.     """
111.     初始化 GMM 参数
112.     :param ClusterData: 所有数据点
113.     :param K: 要分类的类别
114.     :return: GMM 的各个参数
115.     """
116.     d = ClusterData.shape[1]
117.     # 高斯分布模型均值
118.     mus = np.random.rand(K, d)
119.     # 初始化协方差矩阵
120.     sigmas = np.array([np.eye(d)] * K)
121.     # 假设起始时各个类别概率相同
122.     alpha = np.ones(K) * (1 / K)
123.     return mus, sigmas, alpha
124.
125. def drawGMM(ClusterData, gamma, mus):
126.     """

```

```

127.     画出 GMM 簇的图像
128.     :param ClusterData:所有数据点
129.     :param gamma: gamma 矩阵
130.     """
131.     center = np.zeros(ClusterData.shape[0])
132.     for i in range(ClusterData.shape[0]):
133.         center[i] = gamma[i].argmax()
134.     plt.title("GMM EM algorithm")
135.     for i in range(gamma.shape[1]):
136.         listX1 = []
137.         listX2 = []
138.         for j in range(ClusterData.shape[0]):
139.             if int(center[j]) == i:
140.                 listX1.append(ClusterData[j][0])
141.                 listX2.append(ClusterData[j][1])
142.             plt.scatter(listX1, listX2,c=color[i])
143.     mu = np.transpose(mus)
144.     plt.scatter(mu[0],mu[1],c=color[gamma.shape[1]],s=100)
145.     plt.show()
146.
147.
148. def TestGMMAccuracy(centersTrain, mus,sigmas,alpha, YTrain, YTest,Y,Test, K
    ):
149.     # 簇号
150.     Ytest = np.zeros(Test.shape[0])
151.     for i in range(Test.shape[0]):
152.         gamma= getExpectation(Test, mus, sigmas, alpha, K)
153.         Ytest[i] = np.array(gamma[i]).argmax()
154.     # Dict 簇到原来 y 的映射
155.     Dict = ClusterMap(K, YTrain,centersTrain,Y)
156.     sum = 0
157.     for i in range(YTest.shape[0]):
158.         Ytest[i] = Dict[Ytest[i]]
159.     for i in range(YTest.shape[0]):
160.         sum = sum +1
161.     Accuracy = sum/Ytest.shape[0]
162.     print("GMM EM iris 数据集的正确率为{}".format(Accuracy))
163.
164.
165. def ClusterMap(K,YTrain, dataCenters,Y):
166.     Set = list(set(Y))
167.     List = []
168.     for i in range(K):
169.         List.append(i)

```

```

170.     listArray = list(itertools.permutations(List,K))
171.     MapAccuracyDict = np.zeros(len(listArray))
172.     listDict = []
173.     for j in range(len(listArray)):
174.         # 簇号
175.         Dict = {}
176.         for k in range(len(listArray[j])):
177.             Dict[listArray[j][k]] = Set[k]
178.         listDict.append(Dict)
179.         MapAccuracyDict[j] = MapAccuracy(Dict,YTrain,dataCenters)
180.     return listDict[np.array(MapAccuracyDict).argmax()]
181.
182.
183.
184. def MapAccuracy(Dict,YTrain,dataCenters):
185.     sum = 0
186.     for i in range(YTrain.shape[0]):
187.         if Dict[dataCenters[i]] == YTrain[i]:
188.             sum = sum+1
189.     return sum/YTrain.shape[0]

```

readFromFile.py

```

1. import pandas as pd
2. import numpy as np
3. from sklearn.preprocessing import LabelEncoder
4. from sklearn.model_selection import train_test_split
5.
6.
7. def readFromFile(DocumentName):
8.     """
9.     从文件中读取数据
10.    :param DocumentName: 文件的名称
11.    :return: 数据
12.    """
13.    df = pd.read_csv(DocumentName)
14.    df.rename(columns={df.columns.array[df.columns.shape[0] - 1]: 'Predict'},
15.              inplace=True)
16.    df['Predict'] = LabelEncoder().fit_transform(df['Predict'])
17.    Y = df['Predict']
18.    X = df.drop('Predict', axis=1)
19.    X = np.array(X.values)
20.    Y = np.array(Y.values)
21.    K = len(set(Y))

```

```
21.     return X, Y, K
22.
23.
24. def makeFileDataForKmeans(X, Y):
25.     """
26.     使数据适合写好的模型
27.     :param X: 数据点
28.     :return: 修改后的数据
29.     """
30.     Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, Y)
31.     Xtrain = np.transpose(Xtrain)
32.     Xtest = np.transpose(Xtest)
33.     Ytrain = np.transpose(Ytrain)
34.     Ytest = np.transpose(Ytest)
35.     return Xtrain, Xtest, Ytrain, Ytest
36.
37.
38. def makeFileDataForGMM(X, Y):
39.     """
40.     使数据适合写好的模型
41.     :param X: 数据点
42.     :return: 修改后的数据
43.     """
44.     Xtrain, Xtest, Ytrain, Ytest = train_test_split(X, Y)
45.     return Xtrain, Xtest, Ytrain, Ytest
```