



哈爾濱工業大學

Harbin Institute of Technology

机器学习实验报告

GMM 模型

课程名称：	机器学习
学院：	计算学部
专业：	计算机科学与技术
学号：	1180300811
姓名：	孙骁
指导老师：	刘扬
学期：	2020 秋季学期

2020 年 11 月 3 日

目录

一、实验目的和要求	4
1.1 实验题目	4
1.2 实验目标	4
1.3 实验测试	4
1.4 实验应用	4
二、实验环境	4
三、实验原理	4
3.1 聚类问题	4
3.2 K-Means 算法原理	5
3.3 混合高斯模型	5
3.4 EM 算法	6
四、算法实现	7
4.1 K-Means 算法	7
4.2 高斯混合聚类算法	8
五、实验步骤	8
5.1 生成满足二维高斯分布数据	8
5.2 实现章节 4.1 中的 K_Means 算法	9
5.3 实现章节 4.2 中的混合高斯模型	9
5.4 使用 UCI 数据集测试 K-Means 与混合高斯模型	9
六、实验结果	9
6.1 K-Means 聚类结果	9
6.2 混合高斯模型聚类结果	9
6.3 使用 UCI 数据集测试 K-Means 与混合高斯模型	9
七、实验结论	10
参考文献	11
附录 A K-Means 算法-k_means.py	12
附录 B 混合高斯模型与 EM 算法求解-gaussian_mixture_model.py	13

附录 C	读取鸢尾花数据集— <code>read_iris_data.py</code>	15
附录 D	主程序— <code>k_means_GMM.py</code>	15
附录 E	生成二维高斯分布数据— <code>generate_data.py</code>	18

一、实验目的和要求

1.1 实验题目

实现 k-means 聚类方法和混合高斯模型.

1.2 实验目标

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.

1.3 实验测试

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).

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).

1.4 实验应用

可以 UCI 上找一个简单问题数据, 用你实现的 GMM 进行聚类.

二、实验环境

1. Anaconda 4.8.4
2. Python 3.7.4
3. PyCharm 2019.1 (Professional Edition)
4. Windows 10 2004

三、实验原理

3.1 聚类问题

聚类问题属于无监督学习问题. 假定样本集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ 包含 m 个无标记样本, 每个样本 $\mathbf{x}_i = (x_{i1}; x_{i2}; \dots; x_{in})$ 是一个 n 维特征向量, 则聚类算法将样本集 D 划分为 k 个不相交的簇 $\{C_l | l = 1, 2, \dots, k\}$, 其中 $C_{l'} \cap_{l' \neq l} C_l = \emptyset$. 相应的, 我们用 $\lambda_j \in \{1, 2, \dots, k\}$ 表示样本 \mathbf{x}_j 的“簇标记”, 即 $\mathbf{x}_j \in C_{\lambda_j}$. 于是, 聚类的结果可以用包含 m 个元素的簇标记向量 $\boldsymbol{\lambda} = (\lambda_1; \lambda_2; \dots; \lambda_m)$ 表示.

3.2 K-Means 算法原理

对于给定的样本集 $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, “ k 均值” 算法针对聚类所得到的簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ 最小化平方误差

$$E = \sum_{i=1}^k \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2, \quad (1)$$

其中 $\boldsymbol{\mu}_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$ 是簇 C_i 的均值向量. 式 (1) 刻画了簇内样本围绕均值向量的紧密程度, E 值越小, 则簇内样本相似度越高.

最优化式 (1) 并不容易, 找到其最优解需要考查样本集 D 所有可能的簇划分, 显然这是一个 NP 难问题. 因此, k 均值算法采用了贪心策略, 通过迭代优化来近似最小化式 (1).

3.3 混合高斯模型

对于 n 维样本空间 \mathcal{X} 中的随机变量 \mathbf{x} , 若 \mathbf{x} 服从高斯分布, 其概率密度函数为

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

其中 $\boldsymbol{\mu}$ 是 n 维均值向量, $\boldsymbol{\Sigma}$ 是 $n \times n$ 的协方差矩阵, 由式 (2) 可知, 高斯分布完全由均值向量 $\boldsymbol{\mu}$ 和协方差矩阵 $\boldsymbol{\Sigma}$ 这两个参数决定, 因此我们也将概率密度函数记为 $p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

定义 1 (混合高斯分布) 定义高斯混合分布如下:

$$p_{\mathcal{M}}(\mathbf{x}) = \sum_{i=1}^k \alpha_i \cdot p(\mathbf{x}|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i), \quad (3)$$

该分部共由 k 个混合成分组成, 每个混合成分对应一个高斯分布. 其中 $\boldsymbol{\mu}_i$ 与 $\boldsymbol{\Sigma}_i$ 是第 i 个高斯混合成分的参数, 而 $\alpha_i > 0$ 为相应的 “混合系数”, 且 $\sum_{i=1}^k \alpha_i = 1$.

假设样本的生成过程由高斯混合分布给出: 首先, 根据 $\alpha_1, \alpha_2, \dots, \alpha_k$ 定义的先验分布选择高斯混合成分, 即 $p(z_j = i) = \alpha_i$, 其中 α_i 是选择第 i 个混合成分的概率; 然后, 根据被选择的混合成分的概率密度进行采样, 从而生成相应的样本. 根据贝叶斯定理, z_j 的后验分布如式 (4) 所示,

$$\begin{aligned} p_{\mathcal{M}}(z_j = i|\mathbf{x}_j) &= \frac{P(z_j = i) \cdot p_{\mathcal{M}}(\mathbf{x}_j|z_j = i)}{p_{\mathcal{M}}(\mathbf{x}_j)} \\ &= \frac{\alpha_i \cdot p(\mathbf{x}_j|\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{l=1}^k \alpha_l \cdot p(\mathbf{x}_j|\boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)}, \end{aligned} \quad (4)$$

即 $p_{\mathcal{M}}(z_j = i|\mathbf{x}_j)$ 给出了样本 \mathbf{x}_j 由第 i 个高斯混合成分生成的后验概率.

当混合高斯分布 (3) 已知时, 混合高斯聚类算法将把样本集 D 划分为 k 个簇 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$, 每个样本 \mathbf{x}_j 的簇标记 λ_j 如式 (5) 确定:

$$\lambda_j = \arg \max_{i \in \{1, 2, \dots, k\}} \gamma_{ji}. \quad (5)$$

对于给定样本集 D 式 (3) 中的模型参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) | 1 \leq i \leq k\}$ 的求解, 可以采用极大似然估计, 即最大化对数似然

$$\begin{aligned} LL(D) &= \ln \left(\prod_{j=1}^m p_{\mathcal{M}}(\mathbf{x}_j) \right) \\ &= \sum_{j=1}^m \ln \left(\sum_{i=1}^k \alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \right), \end{aligned} \quad (6)$$

对于以上问题的最优化求解, 常采用 EM 算法进行迭代优化.

3.4 EM 算法

若参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) | 1 \leq i \leq k\}$ 能使得式 (6) 最大化, 则由 $\frac{\partial LL(D)}{\partial \boldsymbol{\mu}_i} = 0$ 有

$$\sum_{j=1}^m \frac{\alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{l=1}^k \alpha_l \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)} (\mathbf{x}_j - \boldsymbol{\mu}_i) = 0, \quad (7)$$

由式 (4) 及 $\gamma_{ji} = p_{\mathcal{M}}(z_j = i | \mathbf{x}_j)$, 有

$$\boldsymbol{\mu}_i = \frac{\sum_{j=1}^m \gamma_{ji} \mathbf{x}_j}{\sum_{j=1}^m \gamma_{ji}}, \quad (8)$$

即各混合成分的均值可以通过样本加权平均来估计, 样本权重是每个样本属于该成分的后验概率. 类似的, 由 $\frac{\partial LL(D)}{\partial \boldsymbol{\Sigma}_i} = 0$ 可得

$$\boldsymbol{\Sigma}_i = \frac{\sum_{j=1}^m \gamma_{ji} (\mathbf{x}_j - \boldsymbol{\mu}_i) (\mathbf{x}_j - \boldsymbol{\mu}_i)^T}{\sum_{j=1}^m \gamma_{ji}}, \quad (9)$$

对于混合系数 α_i , 除了要最大化 $LL(D)$, 还需要满足 $\alpha_i \geq 0$, $\sum_{i=1}^k \alpha_i = 1$, 考虑 $LL(D)$ 的拉格朗日形式

$$LL(D) + \lambda \left(\sum_{i=1}^k \alpha_i - 1 \right), \quad (10)$$

其中 λ 为拉格朗日乘子. 由式 (10) 对 α_i 的导数为 0, 有

$$\sum_{j=1}^m \frac{\alpha_i \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{l=1}^k \alpha_l \cdot p(\mathbf{x}_j | \boldsymbol{\mu}_l, \boldsymbol{\Sigma}_l)} + \lambda = 0. \quad (11)$$

两边同时乘以 α_i , 对所有样本求和可知 $\lambda = -m$, 有

$$\alpha_i = \frac{1}{m} \sum_{j=1}^m \gamma_{ji}, \quad (12)$$

即每个高斯成分的混合系数由样本属于该成分的平均后验概率确定.

四、算法实现

4.1 K-Means 算法

K-Means 算法的迭代优化算法如算法 (1) 所示.

Algorithm 1 K Means

Input: $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, 聚类簇数 k

Output: 簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

- 1: 从 D 中随机选择 k 个样本作为初始均值向量 $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k\}$
 - 2: **repeat**
 - 3: 令 $C_i = \emptyset$ ($1 \leq i \leq k$)
 - 4: **for** $j = 1, 2, \dots, m$ **do**
 - 5: 计算样本 \mathbf{x}_j 与各均值向量 $\boldsymbol{\mu}_i$ ($1 \leq i \leq k$) 的距离: $d_{ji} = \|\mathbf{x}_j - \boldsymbol{\mu}_i\|_2$;
 - 6: 根据距离最近的均值向量确定 \mathbf{x}_j 的簇标记: $\lambda_j = \arg \min_{i \in \{1, 2, \dots, k\}} d_{ji}$;
 - 7: 将样本 \mathbf{x}_j 划入相应的簇: $C_{\lambda_j} = C_{\lambda_j} \cup \{\mathbf{x}_j\}$;
 - 8: **end for**
 - 9: **for** $i = 1, 2, \dots, k$ **do**
 - 10: 计算新均值向量: $\boldsymbol{\mu}'_i = \frac{1}{|C_i|} \sum_{\mathbf{x} \in C_i} \mathbf{x}$;
 - 11: **if** $\boldsymbol{\mu}'_i \neq \boldsymbol{\mu}_i$ **then**
 - 12: 将当前均值向量 $\boldsymbol{\mu}_i$ 更新为 $\boldsymbol{\mu}'_i$
 - 13: **else**
 - 14: 保持当前均值向量不变
 - 15: **end if**
 - 16: **end for**
 - 17: **until** 当前均值向量均未更新
-

4.2 高斯混合聚类算法

高斯混合聚类算法描述如算法 (2) 所示, 算法首先对高斯混合分布的模型参数进行了初始化. 然后在第 2-12 行基于 EM 算法对模型参数进行迭代更新. 若 EM 算法的停止条件满足 (最大迭代轮数, 或似然函数 $LL(D)$ 增长很少甚至不再增长), 则在第 14-17 行根据高斯混合分布确定簇划分, 并最后返回最终结果.

Algorithm 2 Gaussian mixture clustering algorithm

Input: $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$, 高斯混合成分个数 k

Output: 簇划分 $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$

- 1: 初始化高斯混合分布的模型参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) | 1 \leq i \leq k\}$
 - 2: **repeat**
 - 3: **for** $j = 1, 2, \dots, m$ **do**
 - 4: 根据式 (4) 计算 $\gamma_{ji} = p_{\mathcal{M}}(z_j = i | \mathbf{x}_j) (1 \leq i \leq k)$
 - 5: **end for**
 - 6: **for** $i = 1, 2, \dots, k$ **do**
 - 7: 计算新均值向量: $\boldsymbol{\mu}'_i = \frac{\sum_{j=1}^m \gamma_{ji} \mathbf{x}_j}{\sum_{j=1}^m \gamma_{ji}};$
 - 8: 计算新协方差矩阵: $\boldsymbol{\Sigma}'_i = \frac{\sum_{j=1}^m \gamma_{ji} (\mathbf{x}_j - \boldsymbol{\mu}'_i)(\mathbf{x}_j - \boldsymbol{\mu}'_i)^T}{\sum_{j=1}^m \gamma_{ji}};$
 - 9: 计算新混合系数: $\alpha'_i = \frac{\sum_{j=1}^m \gamma_{ji}}{m};$
 - 10: **end for**
 - 11: 将模型参数 $\{(\alpha_i, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) | 1 \leq i \leq k\}$ 更新为 $\{(\alpha'_i, \boldsymbol{\mu}'_i, \boldsymbol{\Sigma}'_i) | 1 \leq i \leq k\}$
 - 12: **until** 满足停止条件
 - 13: $C_i = \emptyset (1 \leq i \leq k)$
 - 14: **for** $j = 1, 2, \dots, m$ **do**
 - 15: 根据式 (5) 确定 \mathbf{x}_j 的簇标记 λ_j ;
 - 16: 将 \mathbf{x}_j 划入相应的簇: $C_{\lambda_j} = C_{\lambda_j} \cup \{\mathbf{x}_j\}$
 - 17: **end for**
-

五、实验步骤

5.1 生成满足二维高斯分布数据

依照给定的均值和方差, 生成指定数量的满足二维高斯分布的数据, 生成三类数据, 二维均值分别为 1、1; -1、0; 1、-2, 二维协方差矩阵为 $\begin{pmatrix} 0.1 & 0 \\ 0 & 0.1 \end{pmatrix}$, 生成三类数据各生成 100 个.

5.2 实现章节4.1中的 K-Means 算法

实现实现章节4.1中的 K-Means 算法，代码见附录 (A).

5.3 实现章节4.2中的混合高斯模型

实现章节4.2中的混合高斯模型，并使用 EM 算法进行求解，代码见附录 (B).

5.4 使用 UCI 数据集测试 K-Means 与混合高斯模型

选用 UCI 鸢尾花分类数据集，根据鸢尾花的 4 个属性对鸢尾花的分类进行预测. 数据集中一共有三种类别，每个类别各 50 个样本，一共 150 个样本，每条数据包括四个特征，分别是：

1. 萼片长度 (单位：厘米)
2. 萼片宽度 (单位：厘米)
3. 花瓣长度 (单位：厘米)
4. 花瓣宽度 (单位：厘米)

实现章节4.1中的 K-Means 算法与章节4.2中的混合高斯模型，并显示分类的准确率. 代码见附录 (C).

六、实验结果

6.1 K-Means 聚类结果

首先采用随机选取初始簇中心的方法，实验结果如图 (1) 所示. 很明显在分类上出现了错误，分析原因为没有选好的选取初始簇中心，导致了聚类时无法区分与簇中心较远而距离大致相近的点.

第二次选择了选择相距尽可能远的点作为初始的簇中心，聚类结果如图 (2) 所示，此次初始簇中心相距较远，较好的划分了样本集.

6.2 混合高斯模型聚类结果

混合高斯模型的聚类结果如图 (3)，可以看出和 K-Means 聚类效果基本一致.

6.3 使用 UCI 数据集测试 K-Means 与混合高斯模型

数据集没有数据特征缺失的情况，对数据集分别使用 K-Means 算法聚类与基于 EM 算法的高斯混合模型聚类，得到的结果如图 (4) 所示.

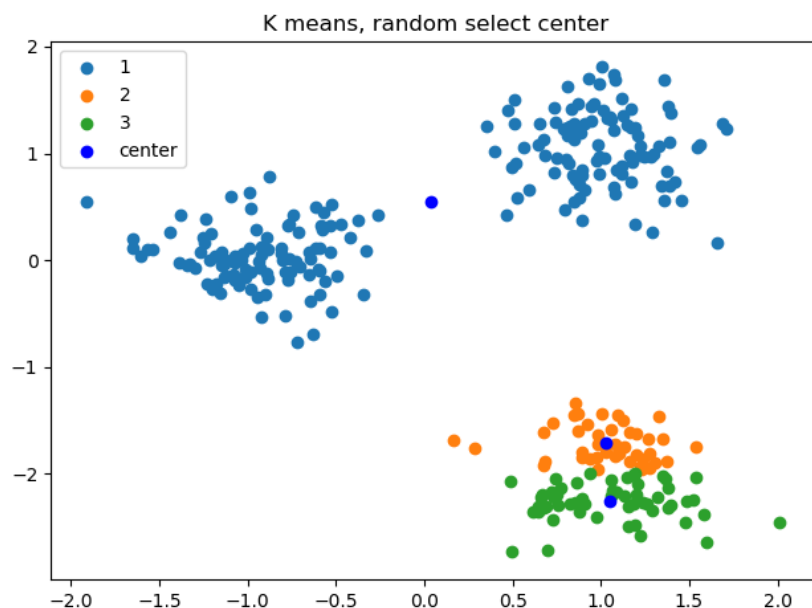


图 1 随机选取初始簇中心的 **K-Means** 算法结果

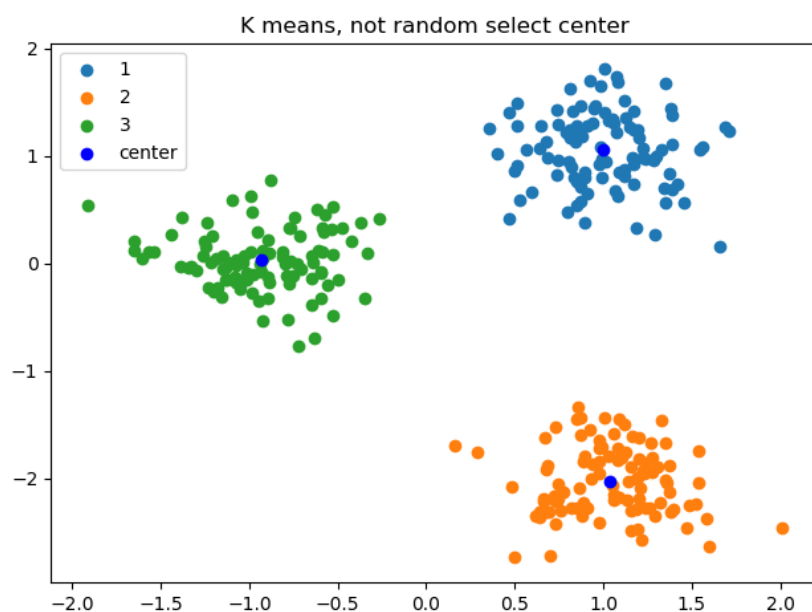


图 2 非随机选取初始簇中心的 **K-Means** 算法结果

七、实验结论

1. **K-Means** 聚类假设数据分布在以簇中心为中心的一定范围内，混合高斯模型则假设数据符合混合高斯分布；
2. 由于 **K-Means** 聚类算法采用贪心的思想，未必能得到全局最优解，簇中心初始化对于最终的结果有很大的影响，如果选择不好初始的簇中心值容易使之陷入局部最优解；

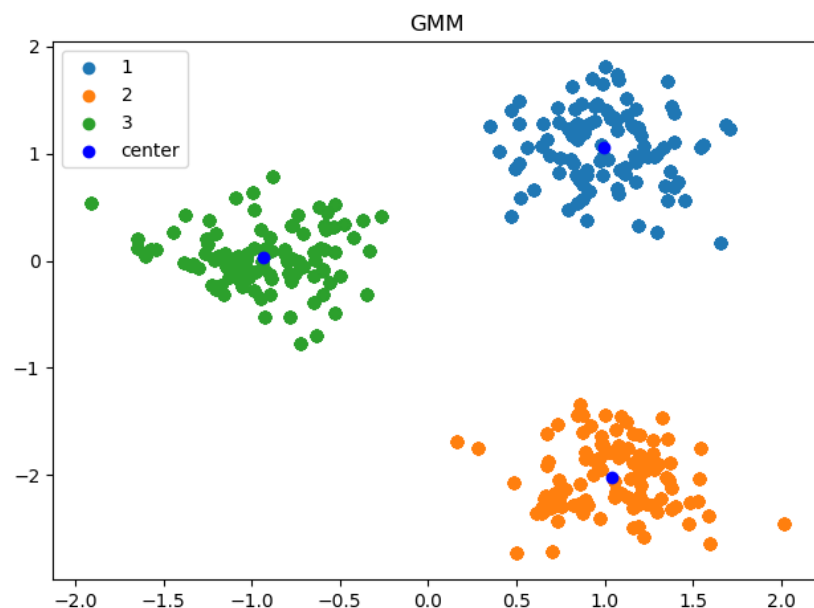


图 3 混合高斯模型的聚类结果

```
k means accuracy: 0.8866666666666667
GMM accuracy: 0.6866666666666666
```

图 4 在 UCI 鸢尾花数据集上测试结果

3. 使用 EM 算法解决混合高斯模型聚类问题时，如果初始高斯模型的均值和方差选取不当，可能会出现极大似然值为 0 的情况，也会出现协方差矩阵不可逆的情况。

参考文献

- [1] 李航, 统计学习方法 (2019.3).
- [2] 周志华, 机器学习 (2016.1).
- [3] Iris Data Set. (1988.7) [Data set].

附录 A K-Means 算法-k_means.py

```
1 import numpy as np
2 import random
3 import collections
4
5
6 def calculate_distance(x_1, x_2):
7     return np.linalg.norm(x_1 - x_2)
8
9
10 class KMeans(object):
11     def __init__(self, data, k, deviation=1e-6):
12         self.data = data
13         self.k = k
14         self.deviation = deviation
15         self.data_rows = data.shape[0]
16         self.data_columns = data.shape[1]
17         self.data_attribution = [-1] * self.data_rows
18         self.mu = self.__initial_k_dots()
19
20     def __initial_k_dots(self):
21         mu_temp = np.random.randint(0, self.k) + 1
22         mu = [self.data[mu_temp]]
23         for i in range(self.k - 1):
24             ans = []
25             for j in range(self.data_rows):
26                 temp_ans = np.sum([calculate_distance(self.data[j], mu[k]) for k in
27                                     range(len(mu))])
28                 ans.append(temp_ans)
29             mu.append(self.data[np.argmax(ans)])
30         return np.array(mu)
31
32     def k_means(self):
33         number = 0
34         while True:
35             result = collections.defaultdict(list)
36             for i in range(self.data_rows):
37                 distance = [calculate_distance(self.data[i], self.mu[j]) for j in
38                             range(self.k)]
39                 lam_j = np.argmin(distance)
40                 result[lam_j].append(self.data[i].tolist())
41                 self.data_attribution[i] = lam_j
42             new_mu = np.array([np.mean(result[i], axis=0).tolist() for i in
43                                range(self.k)])
44             new_loss = np.sum(calculate_distance(self.mu[i], new_mu[i]) for i in
```

```

        range(self.k))
42     if new_loss > self.deviation:
43         self.mu = new_mu
44     else:
45         break
46     # print(number)
47     number += 1
48     # print(self.mu)
49     return self.mu, result
50
51     def random_select_center(self):
52         self.mu = self.data[random.sample(range(self.data_rows), self.k)]
53         return self.k_means()
54
55     def not_random_select_center(self):
56         self.mu = self.__initial_k_dots()
57         return self.k_means()

```

附录 B 混合高斯模型与 EM 算法求解—gaussian_mixture_model.py

```

1  import numpy as np
2  import numpy.random as random
3  import collections
4  from scipy.stats import multivariate_normal
5
6
7  def calculate_distance(x_1, x_2):
8      return np.linalg.norm(x_1 - x_2)
9
10
11 class GaussianMixtureModel(object):
12     def __init__(self, data, k, deviation, iteration_number):
13         self.data = data
14         self.k = k
15         self.deviation = deviation
16         self.iteration_number = iteration_number
17         self.alpha = np.ones(self.k) * (1.0 / self.k)
18         self.data_rows = data.shape[0]
19         self.data_columns = data.shape[1]
20         self.mu, self.sigma = self.__init_params()
21         self.data_attribution = [-1] * self.data_rows
22         self.result = collections.defaultdict(list)
23         self.gamma = None
24

```

```

25 def __init_params(self):
26     mu_0 = random.randint(0, self.k)
27     # print(mu_0)
28     mu_temp = [self.data[mu_0]]
29     for index in range(self.k - 1):
30         temp_ans = []
31         for i in range(self.data_rows):
32             temp_ans.append(np.sum([calculate_distance(self.data[i], mu_temp[j]) for
33                                     j in range(len(mu_temp))]))
34         mu_temp.append(self.data[np.argmax(temp_ans)])
35     mu = np.array(mu_temp)
36     # print("mu", mu)
37     sigma = collections.defaultdict(list)
38     for i in range(self.k):
39         sigma[i] = np.eye(self.data_columns, dtype=float) * 0.1
40
41     # print("sigma", sigma)
42     # for i in range(self.k):
43     #     print("sigma", i, sigma[i])
44     return mu, sigma
45
46 def calculate_likelihood(self):
47     likelihood = np.zeros((self.data_rows, self.k))
48     for i in range(self.k):
49         # print("-----")
50         # print("mu[i]", self.mu[i])
51         # print("sigma[i]", self.sigma[i])
52         likelihood[:, i] = multivariate_normal.pdf(self.data, self.mu[i],
53                                                     self.sigma[i])
54         # print(likelihood[:, i])
55     return likelihood
56
57 def calculate_expectation(self):
58     likelihoods = self.calculate_likelihood() * self.alpha
59     sum_likelihood = np.expand_dims(np.sum(likelihoods, axis=1), axis=1)
60
61     self.gamma = likelihoods / sum_likelihood
62     # print(self.gamma)
63     self.data_attribution = self.gamma.argmax(axis=1)
64     for i in range(self.data_rows):
65         self.result[self.data_attribution[i]].append(self.data[i].tolist())
66
67 def max_function(self):
68     for i in range(self.k):
69         gamma_ji = np.expand_dims(self.gamma[:, i], axis=1)
70         mu_i = (gamma_ji * self.data).sum(axis=0) / gamma_ji.sum()
71         cov = (self.data - mu_i).T.dot((self.data - mu_i) * gamma_ji) / gamma_ji.sum()

```

```

70         self.mu[i], self.sigma[i] = mu_i, cov
71         self.alpha = self.gamma.sum(axis=0) / self.data_rows
72
73     def gmm(self):
74         pre_alpha = self.alpha
75         pre_mu = self.mu
76         pre_sigma = self.sigma
77         for i in range(self.iteration_number):
78             self.calculate_expectation()
79             self.max_function()
80             diff = np.linalg.norm(pre_alpha - self.alpha) + np.linalg.norm(pre_mu -
81                                     self.mu) + np.sum(
82                 [np.linalg.norm(pre_sigma[i] - self.sigma[i]) for i in range(self.k)])
83             if diff > self.deviation:
84                 pre_alpha = self.alpha
85                 pre_sigma = self.sigma
86                 pre_mu = self.mu
87             else:
88                 break
89         self.calculate_expectation()
90         return self.mu, self.result

```

附录 C 读取鸢尾花数据集—read_iris_data.py

```

1  import numpy as np
2  import pandas as pd
3
4
5  def read_iris_data():
6      data_set = pd.read_csv("../data/iris.csv")
7      X = data_set.drop("class", axis=1)
8      Y = data_set['class']
9      # print(Y)
10     return np.array(X, dtype=float), np.array(Y, dtype=str)

```

附录 D 主程序—k_means_GMM.py

```

1  import itertools as it
2
3  from matplotlib import pyplot as plt
4  from src.generate_data import *
5  from src.k_means import *

```

```

6 from src.read_iris_data import *
7 from src.gaussian_mixture_model import *
8
9
10 def generate_picture(k, random_mu, random_result, not_random_mu, not_random_result,
    gmm_mu, gmm_result):
11     plt.title("K means, random select center")
12     for i in range(k):
13         plt.scatter(np.array(random_result[i])[:, 0], np.array(random_result[i])[:, 1],
            label=str(i + 1))
14     plt.scatter(random_mu[:, 0], random_mu[:, 1], c="b", label="center")
15     plt.legend()
16     plt.show()
17
18     plt.title("K means, not random select center")
19     for i in range(k):
20         plt.scatter(np.array(not_random_result[i])[:, 0],
            np.array(not_random_result[i])[:, 1], label=str(i + 1))
21     plt.scatter(not_random_mu[:, 0], not_random_mu[:, 1], c="b", label="center")
22     plt.legend()
23     plt.show()
24
25     plt.title("GMM")
26     for i in range(k):
27         plt.scatter(np.array(gmm_result[i])[:, 0], np.array(gmm_result[i])[:, 1],
            label=str(i + 1))
28     plt.scatter(gmm_mu[:, 0], gmm_mu[:, 1], c="b", label="center")
29     plt.legend()
30     plt.show()
31
32
33 def test_iris_data():
34     data_X, data_Y = read_iris_data()
35     k_means = KMeans(data_X, 3)
36     k_means_mu, k_means_result = k_means.not_random_select_center()
37     k_means_attribution = k_means.data_attribution
38
39     number = len(data_Y)
40     counts_kmeans = []
41     result = list(it.permutations(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'],
        3))
42     for i in range(len(result)):
43         count = 0
44         for index in range(number):
45             # print(data_Y[index])
46             # print(result[i][k_means_attribution[index]])
47             if data_Y[index] == result[i][k_means_attribution[index]]:

```



```

48         count += 1
49         counts_kmeans.append(count)
50     kmeans_accuracy = 1.0 * np.max(counts_kmeans) / number
51     print("k means accuracy:", kmeans_accuracy)
52
53     deviation = 1e-12
54     iteration_number = 10000
55     # print(data_X.shape)
56     gmm = GaussianMixtureModel(data_X, 3, deviation, iteration_number)
57     gmm_mu, gmm_result = gmm.gmm()
58
59     counts_gmm = []
60     gmm_attribution = gmm.data_attribution
61     # print(gmm_attribution)
62     for i in range(len(result)):
63         count = 0
64         for index in range(number):
65             if data_Y[index] == result[i][gmm_attribution[index]]:
66                 count += 1
67             counts_gmm.append(count)
68     gmm_accuracy = 1.0 * np.max(counts_gmm) / number
69     print("GMM accuracy:", gmm_accuracy)
70
71
72 def main():
73     k = 3
74     category_means = [[1, 1], [-1, 0], [1, -2]]
75     category_number = [100, 100, 100]
76     data = generate_2_dimension_data(category_means, category_number, k)
77     k_means_result = KMeans(data, k)
78     random_mu, random_result = k_means_result.random_select_center()
79     not_random_mu, not_random_result = k_means_result.not_random_select_center()
80
81     deviation = 1e-12
82     iteration_number = 10000
83     gmm = GaussianMixtureModel(data, k, deviation, iteration_number)
84     gmm_mu, gmm_result = gmm.gmm()
85
86     generate_picture(k, random_mu, random_result, not_random_mu, not_random_result,
87                     gmm_mu, gmm_result)
88
89     test_iris_data()
90
91 if __name__ == '__main__':
92     main()

```

附录 E 生成二维高斯分布数据—generate_data.py

```
1 import numpy as np
2
3
4 def generate_2_dimension_data(category_means, category_number, k_number):
5     cov = [[0.1, 0], [0, 0.1]]
6     data = []
7     for i in range(k_number):
8         for index in range(category_number[i]):
9             data.append(
10                 np.random.multivariate_normal([category_means[i][0],
11                                                 category_means[i][1]], cov).tolist())
11     return np.array(data)
```