# 模式识别实验报告

## 实验二 GMM 分类器

学院：计算机学院

姓名：张景润

学号：1172510217

# 一、实验内容

1、使用 Python 或 Matlab 编程实现 GMM 算法：要求独立完成算法编程，禁止调用已有函数库或工具箱中的函数；

2、使用仿真数据测试算法的正确性：两类 2 维各 1000 个训练样本 Train1 和 Train2 分别采样自如下两个 GMM，使用训练样本分别估计包含 2 个分量高斯的 GMM 参数。

GMM1：   $\alpha_1 = \dfrac{2}{3}$，$\mathbf{\mu}_1 = (0,0)^t$，$\mathbf{\Sigma}_1 = \begin{pmatrix} 3 & 1 \\ 1 & 1 \end{pmatrix}$

$\alpha_2 = \dfrac{1}{3}$，$\mathbf{\mu}_2 = (10,10)^t$，$\mathbf{\Sigma}_2 = \begin{pmatrix} 2 & 2 \\ 2 & 5 \end{pmatrix}$

GMM2：   $\alpha_1 = \dfrac{2}{3}$，$\mathbf{\mu}_1 = (2,10)^t$，$\mathbf{\Sigma}_1 = \begin{pmatrix} 1 & 1 \\ 1 & 3 \end{pmatrix}$

$\alpha_2 = \dfrac{1}{3}$，$\mathbf{\mu}_2 = (15,20)^t$，$\mathbf{\Sigma}_2 = \begin{pmatrix} 5 & 2 \\ 2 & 1 \end{pmatrix}$

构造区分两类的 GMM 分类器，测试采样自同样 GMM 的测试样本 Test1 和 Test2。

3、MNIST 数据集测试：使用 TrainSamples 中的 30000 个 17 维特征手写数字样本训练 GMM 分类器区分 10 个类别，TrainLabels 中包含训练样本的标签；测试设置不同高斯数量 GMM 分类器对 TestSamples 中 10000 个样本的识别正确率。

# 二、程序代码

（GMM 参数估计部分和 GMM 分类器部分代码）

```
1.  class GMM:
2.      def __init__(self, gauss_num):
3.          self.gauss_num = gauss_num  # 高斯混合分布的数目
4.          self.alphas_np = None  # 每一个高斯分布的权重
5.          self.means_np = None  # 各维变量的均值
6.          self.covs_np = None  # 协方差矩阵
7.
8.      def init_params(self, data_np):  # 初始化高斯分布的三部分参数
9.          data_shuffled = data_np.copy()
10.         np.random.shuffle(data_shuffled)
11.         data_split = np.array_split(data_shuffled, self.gauss_num)
12.         self.alphas_np = np.repeat(1.0 / self.gauss_num, self.gauss_num)
13.         self.means_np = np.array([np.mean(data_split[i], axis=0) for i in range(self.gauss_num)])
14.         self.covs_np = np.array([np.cov(data_split[i].T) for i in range(self.gauss_num)])
15.
16.     def pdf(self, x_np, idx):  # 概率密度计算
17.         mean_np = self.means_np[idx]
18.         cov_np = self.covs_np[idx]
19.         left = 1 / np.sqrt(np.power(2 * np.pi, x_np.shape[0]) * np.linalg.det(cov_np))
```

```python
20.          right = np.exp(-
0.5 * np.dot(np.dot((x_np - mean_np).T, np.linalg.inv(cov_np)), x_np - mean_
np))
21.          return left * right
22.
23.      def fit(self, data_np, fit_num, eps):   # EM过程计算参数
24.          self.init_params(data_np)
25.          data_num = data_np.shape[0]   # 数据数目
26.          norm_densities = np.empty((data_num, self.gauss_num), np.float)
27.          responsibilities = np.empty((data_num, self.gauss_num), np.float)
28.          pre_log_likelihood = 0
29.          for idx in range(fit_num):
30.              for i, one_data in enumerate(data_np):
31.                  for j in range(self.gauss_num):
32.                      norm_densities[i][j] = self.pdf(one_data, j)
33.              log_vec = np.log(np.array([np.dot(self.alphas_np, norm_density)
for norm_density in norm_densities]))
34.              log_likelihood = log_vec.sum()
35.              if abs(log_likelihood - pre_log_likelihood) < eps:
36.                  break
37.              for i in range(data_num):
38.                  normalizer = np.dot(self.alphas_np.T, norm_densities[i])
39.                  for j in range(self.gauss_num):
40.                      responsibilities[i][j] = self.alphas_np[j] * norm_densit
ies[i][j] / normalizer
41.              for i in range(self.gauss_num):
42.                  responsibility = responsibilities.T[i]
43.                  normalizer = np.dot(responsibility, np.ones(data_num))
44.                  self.alphas_np[i] = normalizer / data_num
45.                  self.means_np[i] = np.dot(responsibility, data_np) / normali
zer
46.                  diff = data_np - np.tile(self.means_np[i], (data_num, 1))
47.                  self.covs_np[i] = np.dot((responsibility.reshape(data_num, 1
) * diff).T, diff) / normalizer
48.              pre_log_likelihood = log_likelihood
49.              print('第%d次迭代' % (idx + 1))
50.          print(self.alphas_np)
51.          print(self.means_np)
52.          print(self.covs_np)
53. class Classifier:
54.      def __init__(self, gmm_lst, weight_lst):
55.          self.gmm_lst = gmm_lst
56.          self.priors = weight_lst
57.          self.classes = len(weight_lst)
```

```
58.
59.    def classify(self, data_np, label):
60.        data_num = data_np.shape[0]
61.        if isinstance(label, int):
62.            label = np.full((data_num,), label)
63.        log_vec = np.empty((self.classes, data_num), dtype=np.float)
64.        for idx, gmm in enumerate(self.gmm_lst):
65.            norm_densities = np.empty((data_num, gmm.gauss_num), np.float)
66.            for i in range(data_num):
67.                for j in range(gmm.gauss_num):
68.                    norm_densities[i][j] = gmm.pdf(data_np[i], j)
69.            log_vec[idx] = np.array([np.dot(gmm.alphas_np, norm_density) for
    norm_density in norm_densities]) * \
70.                          self.priors[idx]
71.        predict_np = np.argmax(log_vec, axis=0)
72.        right_num = (predict_np == label.reshape(predict_np.shape)).sum()
73.        accuracy = right_num / data_num
74.        print('[%d/%d]=%.2f%%' % (right_num, data_num, accuracy * 100))
```

## 三、实验结果

1、仿真数据实验结果：给出估计出的两 GMM 模型参数，以及测试样本的识别结果。

GMM 估计模型参数

|  | $\alpha$ | $\mu$ | $\Sigma$ |
|---|---|---|---|
| GMM1-Gauss1 | 0.65890605 | [-0.04877182, 0.03493031] | [2.85162688,0.97072797] [0.97072797,0.96895055] |
| GMM1-Gauss2 | 0.34109395 | [9.97026666, 9.95347321] | [2.01631378,2.35543427] [2.35543427,5.31829801] |
| GMM2-Gauss1 | 0.66799982 | [2.02206041 10.16700955] | [0.96728744,0.91060643] [0.91060643,2.74892244] |
| GMM2-Gauss2 | 0.33200018 | [14.97097406 19.99245787] | [5.288094,2.18044657] [2.18044657,1.12285599] |

GMM 分类器识别结果

|  | 正确识别数 | 正确识别率 |
|---|---|---|
| Test1 | 1000 | 100% |
| Test2 | 1000 | 100% |

2、MNIST 数据集实验结果：

### GMM 分类器识别正确率

| 高斯数 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| 正确识别数 | 9332 | 9421 | 9547 | 9491 | 9526 |
| 正确识别率 | 93.32% | 94.21% | 95.47% | 94.91% | 95.26% |