# E11 Expected Maximize Algorithm

Suixin Ou

School of Computer Science
Sun Yat-sen University

December 21, 2021

# Background

## The Iris Data Set

- The UCI dataset
  (http://archive.ics.uci.edu/ml/index.php) is the most
  widely used dataset for machine learning. If you are interested
  in other datasets in other areas, you can refer to https://
  www.zhihu.com/question/63383992/answer/222718972.

- It is perhaps the best known database to be found in the
  pattern recognition literature. The data set contains 3 classes
  of 50 instances each, where each class refers to a type of iris
  plant. One class is linearly separable from the other 2; the
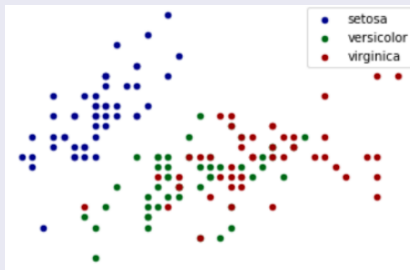  latter are NOT linearly separable from each other.

# Background

## The Iris Data Set



Figure 1: Visualization of Iris Dataset

# Task

## Description

- Dataset statistics

| Data Set Characteristics: | Multivariate | Number of Instances: | 150 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Real | Number of Attributes: | 4 | Date Donated | 1988-07-01 |
| Associated Tasks: | Classification | Missing Values? | No | Number of Web Hits: | 4403737 |

- Domain information

Attribute Information:
1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class:
   -- Iris Setosa
   -- Iris Versicolour
   -- Iris Virginica

# Solution

Read the file "iris.data"

```python
127  def loadData(filename):
128      """从文件中读取数据
129
130      :param filename : the path of file
131      :return : the dataset
132      :return type : list
133
134      """
135      dataSet = []
136      with open(filename) as fr:
137          for i, line in enumerate(fr.readlines()):
138              curLine = line.strip().split(",")
139              fltLine = list(map(float, curLine[:-1]))
140              dataSet.append(fltLine)
141      return dataSet
```

# Solution

Initialize parameters

```python
 95  def init_params(shape, K):
 96      """initialize the parameters : mu, gamma, pi
 97
 98      :param shape: the row and column of data
 99      :param K: the number of model
100      :return : the initial parameters
101
102      """
103      N, D = shape
104      mu = np.random.rand(K, D)
105      Sigma = np.array([np.eye(D)] * K)
106      pi = np.array([1.0 / K] * K)
107      return mu, Sigma, pi
```

# Solution

Expected Maximize algorithm framework

```python
110  def GMM_EM(Y, K, times):
111      """GMM_EM
112
113      :param Y :dataset
114      :param K  :the number of model (3)
115      :param times : the iteration times
116      :return : the parameters of three models - mu, gamma , pi
117
118      """
119      Y = scale_data(Y)
120      mu, Sigma, pi = init_params(Y.shape, K)
121      for i in range(times):
122          gamma = getExpectation(Y, mu, Sigma, pi)
123          mu, Sigma, pi = maximize(Y, gamma)
124      return mu, Sigma, pi, gamma
```

## Please Finish the getExpectation function.

```python
19  def getExpectation(Y, mu, Sigma, pi):
20      """E step
21
22      :param Y : data matrix
23      :param mu: the mean of each characterristic of each sample ; mu is a 3*4 matrix
24      :param Sigma :three-covariance-matrix list
25      :param pi: the responsibilities array
26      :return : the new responsibilities matrix(gamma)
27      :return type : matrix
28      """
29      # 样本数
30      N = Y.shape[0]
31      # 模型数
32      K = pi.shape[0]
33
34      # 响应度矩阵，行对应样本，列对应响应度
35      gamma = np.mat(np.zeros((N, K)))
36
37      # 计算各模型中所有样本出现的概率，行对应样本，列对应模型
38      prob = np.zeros((N, K))
39      for k in range(K):
40          prob[:, k] = phi(Y, mu[k], Sigma[k])
41      prob = np.mat(prob)
42
43      # 计算每个模型对每个样本的响应度
44      # TODO
45      return gamma
```

# Solution

**Please Finish the maximize function.**

```python
48  def maximize(Y, gamma):
49      """M step
50
51      :param Y: data matrix
52      :param gamma : the responsibilities matrix
53      :return : the parameters : mu, gamma, pi
54
55      """
56      # 样本数和特征数
57      N, D = Y.shape
58      # 模型数
59      K = gamma.shape[1]
60
61      # 初始化参数值
62      mu = np.zeros((K, D))
63      Sigma = []
64      pi = np.zeros(K)
65
```

# Submission

## Submission

pack your report E11_YourNumber.pdf and source code into zip
file E11_YourNumber.zip, then send it to
ai_course2021@163.com.

# The End