

# M3S7-Statistical Pattern Recognition

## Project 2

Alexandre ELKRIEF

CID: 00732974

March 9, 2015

### 1 EM-Algorithm

The density of a finite mixture distribution has the form:

$$p(\mathbf{x}) = \sum_{i=1}^K \pi_i f_i(\mathbf{x}; \boldsymbol{\theta}_i)$$

where  $f_i(\cdot)$  are the  $K$  component densities, and  $\pi_j$  are mixing proportions. For fixed  $K$ , the EM algorithm can be used to estimate the parameters,  $\boldsymbol{\theta}_i$ ,  $\pi_j$ , for  $i = 1, \dots, K$ , from an iid sample. In this question we will restrict to all component densities being  $p$ -dimensional normal, with density

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{\frac{p}{2}} |\boldsymbol{\Sigma}|^{1/2}} \exp \left( -\frac{1}{2} ((\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})) \right)$$

(a) Write an R function that uses the EM algorithm to find parameters which maximise the likelihood (or minimise the negative log-likelihood) for a sample of size  $n$  from  $p(\mathbf{x})$ , for a given choice of  $K$ . The function prototype should be

**em.norm(x, means, covariances, mix.prop)**

where  $\mathbf{x}$  is an  $n \times p$  matrix of data, **means**, **covariances**, and **mix.prop** are the initial values for the  $K$  mean vectors, covariance matrices and mixing proportions. Consider including arguments, with sensible defaults, for the convergence criterion and the maximum number of iterations.

For the code, see appendix. I have used a convergence criterion of 0.1 between two iterations of the log-likelihood and a upper limit for the number of iterations at 100. These choices kept the computation time within reasonable bounds.

(b) This question will use the first two columns of the object **synth.te** in the **MASS** library:

**x <- synth.te[, -3]**

For  $K = 2, 3, 4, 5, 6$ , use your function to compute the maximum likelihood estimates for the finite mixture of normal distributions, for these data. Select initial parameters either randomly, or by selecting from a plot of the data.

- i. Construct a table that reports, for each choice of  $K$ , the maximised likelihood, and the AIC.
- ii. On the basis of this table, which choice of  $K$  provides the best density estimate? For this choice, construct a contour plot of the estimated density, along with the data.
- iii. Briefly discuss any problems you anticipate using the EM algorithm for computing a mixture model with more components, or in higher dimensions.

i. For this question, I set the initial parameters randomly. Moreover, I chose the same AIC as in the lecture notes, namely  $M = 6K - 1$  where  $K$  corresponds to the number of mixture components. The code outputs a matrix corresponding to the following table:

K	ML	AIC
2	-555.4818	-566.4818
3	-529.8571	-546.8571
4	-486.0654	<b>-509.0654</b>
5	-486.3345	-515.3345
6	-485.1242	-520.1242

ii. On the basis of the above table,  $K = 4$  seems to be the best choice. The following plot corresponds to such a choice.

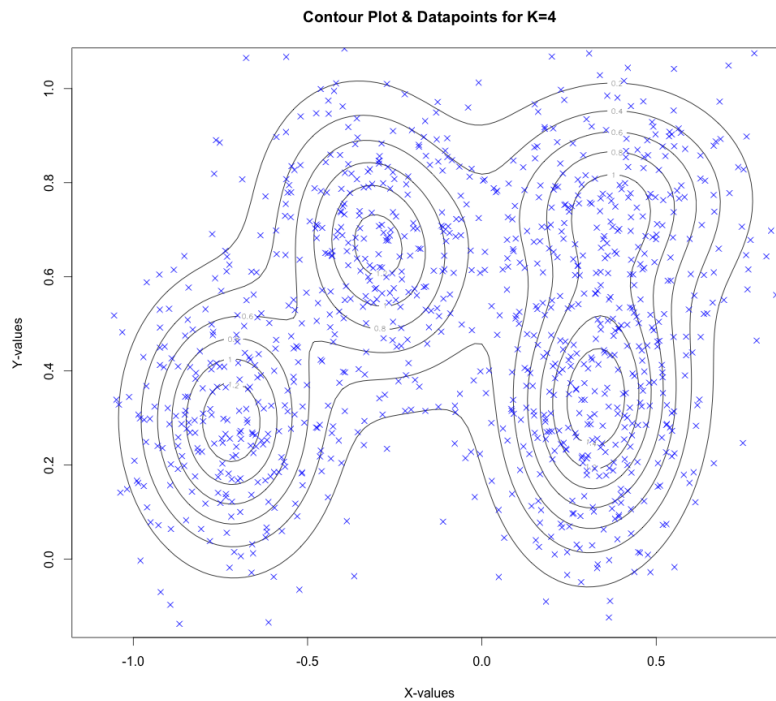


Figure 1: Contour plot for optimal choice of  $K$

iii. The use of double for loops in my EM-Algorithm would cause the computational time to increase exponentially when working in higher dimensions or with more mixture components.

## 2 Rayleigh Distribution

Consider a two-class bivariate classification problem, with equal prior probabilities and class conditional densities given by:

$$f(x, y|C_i) = 4\theta_i^2 xy \exp(-\theta_i(x^2 + y^2)) \quad x, y > 0$$

and  $\theta_i > 0$  for  $i = 1, 2$ . Note that this joint density is the product of Rayleigh distributions.

(a) Write an R function that generates a random sample of size  $n$  from class  $C_1$  and a random sample of size  $n$  for class  $C_2$ . The function should return both the feature vectors and the class indicator. A function for generating Rayleigh distributed random variables is available.

(b) Obtain an expression for the decision boundary for minimum error. Suppose we are interested in the situation where the decision boundary for minimum error intersects with the midpoint of the line connecting the two class mean vectors. Derive an expression for  $\theta_1$  and  $\theta_2$  to satisfy this situation.

$$\begin{aligned}
\frac{f(x, y|C_2)}{f(x, y|C_1)} &= \frac{p(C_1)}{p(C_2)} \\
\Leftrightarrow \frac{4\theta_2^2 xy \exp(-\theta_2(x^2 + y^2))}{4\theta_1^2 xy \exp(-\theta_1(x^2 + y^2))} &= 1 \\
\Leftrightarrow \exp\{(x^2 + y^2)(\theta_1 - \theta_2)\} &= \frac{\theta_1^2}{\theta_2^2} \\
\Leftrightarrow (x^2 + y^2)(\theta_1 - \theta_2) &= 2 \log\left(\frac{\theta_1}{\theta_2}\right) \\
\Leftrightarrow T_{min}^2 = x^2 + y^2 &= 2 \log\left(\frac{\theta_1}{\theta_2}\right) \left(\frac{1}{\theta_1 - \theta_2}\right)
\end{aligned}$$

As we are in 2 dimensional-space the decision boundary for minimum error is a circle on the plane.

(c) Derive an expression in terms of  $\theta_1$  and  $\theta_2$  for the Bayes error rate. Now, suppose  $\theta_2 = 1$  and  $\theta_1 > \theta_2$ . Use the golden ratio search algorithm developed in question 4 of project 1, to determine the value of  $\theta_1$  that gives a Bayes error rate of 15%. The solution occurs in the interval  $[3, 10]$ . (Hint: The target function does not have to be differentiable at the minimum for the golden ratio search to work).

Integrating  $f(x, y)$  we get:

$$\begin{aligned}
F(x, y|C_i) &= \int_x \int_y 4\theta_i^2 xy \exp(-\theta_i(x^2 + y^2)) \, dx dy \\
\Leftrightarrow F(x, y|C_i) &= \int_0^{\pi/2} 2 \sin(\phi) \cos(\phi) \, d\phi \int_0^u 2\theta_i^2 r^3 \exp(-\theta_i(r^2)) \, dr \\
\Leftrightarrow F(x, y|C_i) &= 1 - \theta_i u^2 \exp(-\theta_i u^2) - \exp(-\theta_i u^2)
\end{aligned}$$

Now the expression for the Bayes error rate is the following:

$$\begin{aligned}
e_B &= p(C_1)p(x|C_1 \geq T) + p(C_2)p(x|C_2 \leq T) \\
&\Rightarrow e_B = 0.5[(1 - F_1(T_{min})) + F_2(T_{min})] \\
&\Rightarrow e_B = 0.5[(\theta_1 T_{min}^2 \exp(-\theta_1 T_{min}^2) + \exp(-\theta_1 T_{min}^2) + 1 - \theta_2 T_{min}^2 \exp(-\theta_2 T_{min}^2) - \exp(-\theta_2 T_{min}^2)]
\end{aligned}$$

where  $T_{min} = \sqrt{2 \log(\frac{\theta_1}{\theta_2})(\frac{1}{\theta_1 - \theta_2})}$

In order to find the value of  $\theta_1$  that gives a Bayes error rate of 15% we use the golden ratio search algorithm to minimize the function

$$f(\theta_1) = |e_B - 0.15|$$

This gives a value of  $\theta_1 = 4.652476$

(d) Write down a discriminant function for each class, treating the parameter  $\theta_i$  as unknown.

The code in the appendix uses the following discriminant function:

$$g_i(\mathbf{x}) = \log(p(C_i) + \log(p(\mathbf{x}|C_i)))$$

(e) Let  $\theta_1 = 4$  and  $\theta_2 = 2$ . Construct a plot of the unconditional density

$$f(x, y) = p(C_1)f(x, y|C_1) + p(C_2)f(x, y|C_2)$$

for the specified parameter values. Obtain a sample of 50 observations from each class. Add these data and the Bayes optimal decision boundary to the plot.

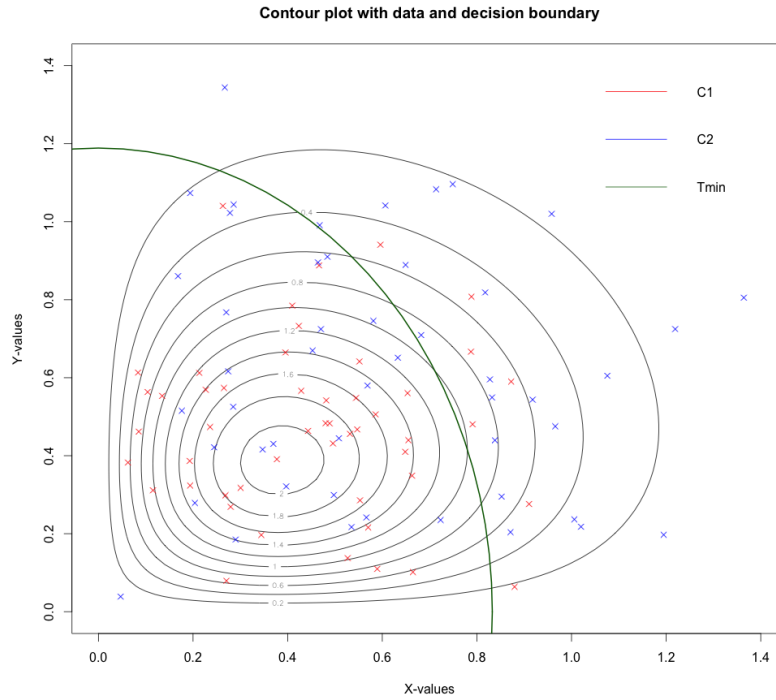


Figure 2: Unconditional Density and Decision Boundary

(f) Derive the maximum likelihood estimators for the parameters of each class, given a sample of size  $n$  from each class.

$$l(\theta_j) = \sum_{i=1}^n \log(f(x_i, y_i|\theta_j))$$

$$\Leftrightarrow l(\theta_j) = \sum_{i=1}^n \log(4\theta_j^2 x_i y_i \exp(-\theta_j(x_i^2 + y_i^2)))$$

$$\begin{aligned}
\Leftrightarrow l(\theta_j) &= \sum_{i=1}^n -\theta_j(x_i^2 + y_i^2) + 2\log(2\theta_j) + \log(x_i y_i) \\
\Leftrightarrow l(\theta_j) &= -\theta_j \sum_{i=1}^n (x_i^2 + y_i^2) + \sum_{i=1}^n \log(x_i y_i) + 2n \log(2\theta_j) \\
&\Rightarrow \frac{\partial l}{\partial \theta_j} = 0 \\
&\Rightarrow -\sum_{i=1}^n (x_i^2 + y_i^2) + \frac{2n}{2\theta_j} = 0 \\
&\Rightarrow \hat{\theta}_j = \frac{2n}{\sum_{i=1}^n (x_i^2 + y_i^2)}
\end{aligned}$$

(g) Write two R functions, the first for computing the maximum likelihood estimates in (f) from a set of data generated by the function in (a), and the second for evaluating the discriminant function for each class, using the maximum likelihood estimates (the estimative discriminant function). Compute the discriminant scores for the data generated in (e) and estimate the error rate of this classifier on this training data.

The discriminant scores are computed within the discriminant function which allocates each points to one of the 2 classes.

The error rate of this classifier ranges between 25% and 37% depending on the set of points generated. (see code in appendix)

(h) Obtain a training sample of size  $n = 200$  and a test sample of size  $n = 10000$ , using the parameter values in part (e). Retain these training and test samples for use in Questions 3 and 4. Using these data sets, compute the training and test set error rates for

- i. the estimative version of the true model, using the functions in part (g),
- ii. Linear discriminant analysis,
- iii. Quadratic discriminant analysis.

Provide a table of these error rates for the different models. Comment on the results.

The table of error rates is the following:

Dataset	ML model	LDA	QDA
Training	29%	32%	28%
Testing	31.8%	32.4%	31.7%

### 3 Kernel Density Estimation

This question is concerned with product kernel classifiers, based on a density estimate of the form:

$$\hat{p}(\mathbf{x}^*) = \frac{1}{n} \frac{1}{h_1 h_2 \dots h_d} \sum_{i=1}^n \prod_{j=1}^d K\left(\frac{[\mathbf{x}^* - \mathbf{x}_i]_j}{h_j}\right)$$

where  $d$  is the dimension of the feature vector,  $K()$  is a kernel function,  $n$  is the number of observations in the sample, and  $h_j, j = 1, 2, \dots, d$ , are bandwidth parameters. Restrict attention to the normal kernel:

$$K(z) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right)$$

(a) Write a function with prototype **kde(x.star,data,bw)** that will perform a kernel density estimate for a single  $d$ -dimensional feature vector **x.star**, using an  $n \times d$  data set **data** and  $d$  bandwidth parameters contained in a vector **bw**.

(b) A classifier proceeds by evaluating a probability density estimate for each class, given bandwidth parameters, and combines these with information about prior probabilities, to obtain estimated posterior probabilities of class membership. In this question, we will use the same bandwidth parameters for both classes.

Using the data you generated in Question 2 part (h) (so  $d = 2$ ), perform a 2-fold cross-validation for bandwidth selection in the following manner (you may take the class priors as  $P(C_1) = P(C_2) = 0.5$ ):

- i. Randomly split the training data generated in Question 2 (h) into two blocks of 100 samples each,  $D_1$  and  $D_2$ , say.
- ii. Construct a  $20 \times 20$  grid of values for  $H = (h_1, h_2)$  defined on  $[0.05, 1] \times [0.05, 1]$ , where  $h_j$  is the bandwidth parameter for variable  $j$ . For each selection of  $H$ :
  - Classify each observation in  $D_2$ , using  $D_1$  as the training set. Compute and retain the error rate.
  - Classify each observation in  $D_1$ , using  $D_2$  as the training set. Compute and retain the error rate.
  - Combine these two error rates to give an average error rate for this selection of  $H$ .

(c) What is the minimum error rate found in part (b), and what are the associated bandwidths  $h_1$  and  $h_2$ ?

The minimum error rate is 30%, corresponding to  $(h_1, h_2) = (0.3, 0.3)$

(d) For this  $h_1$  and  $h_2$  use all your training data and calculate the error rate on the 10000 test sample generated in Question 2 (h).

The error rate we obtain using the the best choice of  $(h_1, h_2)$  is 33.3% .

## 4 Distance-weighted K-NN

In an early attempt to combine the features of kernel methods and  $k$ -nearest neighbour methods, Dudani (1976)<sup>2</sup> proposed a *distance weighted* version of the  $k$ -nearest neighbour classifier. In this version, weights are assigned to the  $k$ -nearest neighbours, with closer neighbours being weighted more heavily. A feature vector  $\mathbf{x}$  is then assigned to the class for which the weights of the representatives among the  $k$  neighbours sum to the greatest value.

One implementation of this strategy is as follows. For a specific feature vector,  $\mathbf{x}$ , let the number of neighbours of class  $i$  among the  $k$ -nearest neighbours be  $k_i, i = 1, 2, \dots, K$ , where  $K$  is the number of classes. Note that  $\sum_{i=1}^K k_i = k$ . Let the Euclidean distance of  $\mathbf{x}$  from each of these class  $i$  neighbours be  $d_i^j$ , for  $j = 1, 2, \dots, k_i$ . The *weight* associated with class  $i$  is

$$w_i = \sum_{j=1}^{k_i} f(d_i^j) \quad i = 1, 2, \dots, K$$

A popular choice for the weighting function  $f$  (in the nonparametric smoothing literature) is the *tricube* function

$$f(x) = (1 - |x|^3)^3$$

for  $|x| \leq 1$ , and 0 otherwise. This requires that we re-scale the distances,  $d_i^j$ , such that the largest, for all  $i$  and  $j$ , is 1. The allocation rule is then: allocate to class  $i$  if  $w_i > w_j$  for all  $j \neq i$ .

(a) Write an R function that implements the procedure described above. The function prototype should be of the form

**knn.dist(train,test,class,k)**

where **train** is a matrix of training data (the rows are feature vectors), **class** is the associated vector of class indicators, **test** is a matrix of test data, and **k** is the number of neighbours.

(b) Using the training data obtained in Question 2 (h), perform a 2- fold cross-validation to select  $k$  from  $k = 3, 7, 11, 15, 19, 23, 27, 31, 35$  or 39 for your distance weighted  $k$ -nearest neighbour classifier.

According to the cross-validation the best choice of  $k$  is 31.

(c) For your selection of  $k$  use all your training data and calculate the error rate on the 10000 test sample from Question 2 (h).

Using the value of  $k$  found above we get an error rate on the 10000 test sample of 32.9%

(d) You have now tried various methods on these data. Discuss the method you prefer for this problem and issues that motivate your preference.

Considering methods from Q2, Q3 and Q4, the Kernel density estimation method from question 2 is the most computationally expensive. Working in higher dimensions and with a bigger range for each bandwidth parameter increases the computational time significantly as the number of combinations of  $h_1, h_2, \dots, h_d$  increased exponentially. Moreover, this method doesn't produce significantly better error rates, at least in the case of our data.

The other methods give similar error rates and are similarly quick to run (although finding the optimal choice of  $K$  in Q4 isn't instantaneous). I would therefore use a combination of these 4 methods (ML, LDA, QDA and Distance-weighted K-NN) depending on the set of data that I need to test.

## 5 Appendix