

MLT - Assignment 3

Ayush Sekhari (12185)

March 11, 2015

1 Part I and II

Part I and II were hand solved and are attached with the report. Please scroll down to view them.

2 Part III

2.1 Part III-a

2.1.1 Generating Data

I used inbuilt runif, rnorm functions to generate the data. 600 points for each class were used. 400 of them were used to train and 100, 100 for validation and test respectively. The metadata of the data is :

1. Training Data : 1200
2. Validation Data : 300
3. Test Data : 300
4. Priors = [1/3,1/3,1/3]

2.1.2 Calculation of the Initial Model

Each class is to be represented as a mixture of three gaussians. I used **K-means clustering** on training data for each class with $k = 3$. So each class data is clustered into 3 groups. These give 3 cluster means and co-variance matrices for each class. These lead to 9 mean vectors and 9 covariance matrices. I also estimate initial priors for each cluster as the cluster group strength.

2.1.3 Estimation of Co-Variance Matrix

I took the pooled co-variance matrix as the weighted average of the cluster co-variance matrices. The weights used were (priors - 1) in order to ensure unbiased estimator of the co-variance. This is kept constant for the problem.

Approach	CV error	Test-Set error
k-means {training error = 12.333}	16.666	15.666
Stochastic Variation	16.333	15.333
Varying all of the 21 components simultaneously with $c = \{-0.5, 0, 0.5\}$	14	17.333
Varying all of the 21 components simultaneously with $c = \{-0.1, 0, 0.1\}$	15.666	16.333

Table 1: Observed Best errors while searching the space around means using various approaches

2.1.4 Generic Approach:

I have an initial model with 9 mean vectors and a SIGMA matrix. I vary the mean vectors locally using various approaches, recompute the priors each time and then use the mixture model to estimate the error on the validation set. The one with the best is used and the errors are reported on the test set. Performance and search results using various search approaches are listed in the table below.

2.2 Part III-b – LDA and QDA

1. LDA gives the error of 0.136 on the test set
2. QDA gives the error of 0.153 on the test set
3. **LDA performs the best**

Ques 1:

$$\Sigma_j = \sigma_j^{-2}(1-\rho_j) I + \sigma_j^{-2}\rho_j 11^T$$

Using $(A+B)^{-1} = A^{-1} - \frac{1}{1+g} A^{-1}BA^{-1}$ where $g = \text{trace}(BA^{-1})$

(as a special case of Woodbury Matrix Inversion Lemma).

[math.stackexchange.com/questions/17776/inverse-of-the-sum-of-matrices]

Here, $A = \sigma_j^{-2}(1-\rho_j) I$, $B = \sigma_j^{-2}\rho_j 11^T$

$$A^{-1} = \sigma_j^{-2}(1-\rho_j)^{-1} I, \quad BA^{-1} = \rho_j(1-\rho_j)^{-1} 11^T \Rightarrow g = \rho_j(1-\rho_j)^{-1} d$$

$$\therefore \Sigma_j^{-1} = (A+B)^{-1} = \sigma_j^{-2}(1-\rho_j)^{-1} I - \frac{1 \times \sigma_j^{-2}(1-\rho_j)^{-2}\rho_j 11^T}{1+\rho_j(1-\rho_j)^{-1} d}$$

$$= C_{1j} I - \frac{\rho_j 11^T}{\sigma_j^2(1-\rho_j)(1+(d-1)\rho_j)} = C_{1j} I - C_{2j} 11^T$$

$$= \Sigma_j^{-1}$$

Finding the Bayes Discriminant Function:

Let $\tilde{x} \in$ Discriminant fn.

Let \tilde{x}, \tilde{u} be the class means. Let p_1, p_2 be the priors & p is the p.d.f.

$$\therefore p_1 p(\tilde{x}|1) = p_2 p(\tilde{x}|2)$$

$$p_1 \times \frac{1}{(2\pi)^{d/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{(\tilde{x}-\tilde{u})^T \Sigma_1^{-1} (\tilde{x}-\tilde{u})}{2}\right) = p_2 \times \frac{1 \times \exp\left(-\frac{(\tilde{x}-\tilde{u})^T \Sigma_2^{-1} (\tilde{x}-\tilde{u})}{2}\right)}{(2\pi)^{d/2} |\Sigma_2|^{1/2}}$$

Taking log on both sides,

$$-\frac{1}{2}(\tilde{x}-\tilde{u})^T \Sigma_1^{-1} (\tilde{x}-\tilde{u}) = \frac{1}{2}(\tilde{x}-\tilde{u})^T \Sigma_2^{-1} (\tilde{x}-\tilde{u}) + \log c \quad (\text{where } c = \frac{p_2}{p_1} \left(\frac{|\Sigma_1|}{|\Sigma_2|} \right)^{1/2})$$

$$\therefore -\frac{1}{2}(\tilde{x}-\tilde{u})^T (\Sigma_1^{-1} - \Sigma_2^{-1}) (\tilde{x}-\tilde{u}) - \log c = 0$$

$$\therefore \frac{1}{2}(\underline{x} - \underline{u})^T (-\Sigma_1^{-1} + \Sigma_2^{-1})(\underline{x} - \underline{u}) \pm \log c = 0$$

(using $\Sigma_i = C_{1i} I - C_{2i} \mathbf{1}\mathbf{1}^T$)

$$(\underline{x} - \underline{u})^T \left(\frac{1}{2}(C_{21} - C_{22}) \mathbf{1}\mathbf{1}^T - \frac{1}{2}(C_{11} - C_{12}) I \right) (\underline{x} - \underline{u}) = \log c.$$

$$\Rightarrow -\frac{1}{2}(C_{11} - C_{12})(\underline{x} - \underline{u})^T(\underline{x} - \underline{u}) + \frac{1}{2}(C_{21} - C_{22})(\underline{x} - \underline{u})^T \mathbf{1}\mathbf{1}^T (\underline{x} - \underline{u}) = \log c$$

$$\Rightarrow -\frac{1}{2}(C_{11} - C_{12})b_1 + \frac{1}{2}(C_{21} - C_{22})((\mathbf{1}^T(\underline{x} - \underline{u}))^T(\mathbf{1}^T(\underline{x} - \underline{u}))) = \log c$$

$$\Rightarrow -\frac{1}{2}(C_{11} - C_{12})b_1 + \frac{1}{2}(C_{21} - C_{22})b_2 = \log c$$

: discriminant is given by $-\frac{1}{2}(C_{11} - C_{12})b_1 + \frac{1}{2}(C_{21} - C_{22})b_2$ within a constant (i.e. $\log c$).

Question-2:

$$p(x) = \sum_{j=1}^g \pi_j p(x|u_j, m_j) \text{ where } p(x|u_j, m_j) = \frac{m_j}{(m_j-1)! u_j} \left(\frac{m_j x}{u_j} \right)^{m_j-1} \exp(-m_j x / u_j)$$

Let us assume that we know the class to which it belongs.

i.e. let $y' = (x, z')$ where z' is the observed latent indicator variable. Vector with all entries at zero except the K th entry, which is 1 if x is represented by the cluster K .

$$\therefore p(y) = \prod_{j=1}^g [p(x|\theta_j) \pi_j]^{z'_j} \quad (\text{since we are sure that it belongs to class given by } z').$$

$$p(z') = p(y_1, y_2, y_3, \dots, y_{1+1}) = \prod_{i=1}^{1+1} p(y_i) \quad [\text{using iid assumption}].$$

$$= \prod_{i=1}^n \left(\prod_{j=1}^g [p(x_i|\theta_j) \pi_j]^{z'_{ji}} \right)$$

$$L = \text{Likelihood} = \prod_{i=1}^n \left(\prod_{j=1}^g \left[p(x_i | \theta_j) \pi_j \right]^{z_{ji}} \right)$$

$$\log \text{Likelihood} = \log L = \sum_{i=1}^n \sum_{j=1}^g \left(z_{ji} \log p(x_i | \theta_j) + z_{ji} \log \pi_j \right)$$

I would want to ~~minimize~~ \hat{l} (\log likelihood) and thus find the parameters.
 z_{ji} = a latent variable, I would remove it and replace it by its expected value

E-Step:

$$w_{ji}^{(m+1)} = E(z_{ji}^* | x_i, \Phi^{(m)}) = \pi_j^{(m)} p(x_i | \theta_j^{(m)}) + 0 \times p(z_{ji}=0 | x_i, \Phi^{(m)})$$

$$= p(z_{ji}=1 | x_i, \Phi^{(m)})$$

$$w_{ji}^{(m+1)} = \frac{\pi_j^{(m)} p(x_i | \theta_j^{(m)})}{\sum_{k=1}^K \pi_k^{(m)} p(x_i | \theta_k^{(m)})}$$

M-Step:

$$\hat{l} = \text{an estimator of } l = \sum_{i=1}^n \sum_{j=1}^g \left(w_{ji} \log p(x_i | \theta_j) + z_{ji} \log \pi_j \right)$$

I want to ~~minimize~~ the \hat{l} under the constraint $\sum_{j=1}^g \pi_j = 1$

$$\therefore \hat{l}' = \hat{l} - \lambda \left(\sum_{j=1}^g \pi_j - 1 \right) \quad (\text{using the idea of lagrange multiplier})$$

Now, we want to maximize \hat{l}' (without any constraints).

$$\therefore \frac{\partial \hat{l}'}{\partial \pi_j} = 0 \quad \forall j \in 1 \dots g, \quad \frac{\partial \hat{l}'}{\partial \lambda} = 0, \quad \frac{\partial \hat{l}'}{\partial m} = 0.$$

$$\frac{\partial \hat{l}}{\partial \pi_j} = \sum_{i=1}^n \frac{w_{ji}}{\pi_j} - \lambda \stackrel{=0}{=} \Rightarrow \hat{\pi}_j = \frac{\sum_{i=1}^n w_{ji}}{\lambda}$$

$$\Rightarrow \sum_{j=1}^g \hat{\pi}_j = 1 \Rightarrow \frac{\sum_{i=1}^n \sum_{j=1}^g w_{ji}}{\lambda} = 1 \Rightarrow \frac{\sum_{i=1}^n 1}{\lambda} = 1 \Rightarrow \lambda = n$$

$$\therefore \hat{\pi}_j = \frac{\sum_{i=1}^n w_{ji}}{n}$$

$$\hat{l} = \sum_{i=1}^n \sum_{j=1}^g w_{ji} \log \pi_j - \lambda \left(\sum_{j=1}^g \pi_j - 1 \right) +$$

$$+ \sum_{i=1}^n \sum_{j=1}^g w_{ji} \left(\log m_j - \log(m_j - 1)! - \log y_j + \frac{(m_j - 1) \log m_j x_i}{y_j} - \frac{m_j x_i}{y_j} \right)$$

$$\frac{\partial \hat{l}}{\partial y_j} = + \sum_{i=1}^n \sum_{j=1}^g \left(\frac{w_{ji}}{y_j} \left(-\frac{1}{y_j} + \frac{(m_j - 1)}{y_j} + \frac{m_j x_i}{y_j^2} \right) \right) = 0$$

$$\Rightarrow \sum_{j=1}^n \left(w_{ji} \left(\frac{x_i}{y_j^2} - \frac{1}{y_j} \right) \right) = 0$$

$$\Rightarrow y_j = \frac{\sum_{i=1}^n w_{ji} x_i}{\sum_{i=1}^n w_{ji}}$$

$$\cancel{M = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^g w_{ji} x_i}$$

$$\therefore \hat{u}_j = \frac{\sum_{i=1}^n w_{ji} x_i}{\sum_{i=1}^n w_{ji}}$$

$$\frac{\partial \hat{l}}{\partial m_j} = 0 \Rightarrow \sum_{i=1}^n w_{ji} \left[\frac{1}{m_j} - \frac{\partial}{\partial m_j} \log(m_j - 1)! + \log m_j x_i + \frac{m_j - 1}{m_j} - \log u_j - \frac{x_i}{u_j} \right] = 0$$

$$\Rightarrow \sum_{i=1}^n w_{ji} \left[1 + \log m_j + \log x_i - \log u_j - \frac{x_i}{u_j} - \frac{\partial}{\partial m_j} \log(m_j - 1)! \right] = 0$$

$$\Rightarrow \sum_{i=1}^n w_{ji} \left[\log m_j - \frac{\partial}{\partial m_j} \log(m_j - 1)! + \log x_i - \log u_j - \frac{x_i}{u_j} + 1 \right] = 0$$

Let $\log m_j - \frac{\partial}{\partial m_j} \log(m_j - 1)! = V(m_j)$,

claim: We can estimate m_j from $V(m_j)$ to arbitrary accuracy using numerical methods.

\therefore also, $\frac{\partial}{\partial m_j} \log(m_j - 1)! =$ digamma function

{ [www.en.wikipedia.org/wiki/Digamma-function](https://en.wikipedia.org/wiki/Digamma-function) }

$$\Rightarrow \sum_{i=1}^n w_{ji} \left[V(m_j) + \log \frac{x_i}{u_j} e^{-x_i/u_j + 1} \right] = 0$$

$$V(m_j) = \frac{-\sum_{i=1}^n w_{ji} \log \frac{x_i}{u_j} e^{1-x_i/u_j}}{\sum w_{ji}}$$



m_j

The Iteration steps are:

E-Step:

$$\hat{w}_{ji}^{(m+1)} = \frac{\pi_j^{(m)} p(x_i | \theta_j^{(m)})}{\sum_{k=1}^g \pi_k^{(m)} p(x_i | \theta_k^{(m)})}$$

M-Step:

$$\hat{\pi}_j^{(m+1)} = \frac{\sum_{i=1}^n w_{ji}^{(m+1)}}{n}$$

$$\hat{y}_j^{(m+1)} = \frac{\sum_{i=1}^n w_{ji}^{(m+1)} x_i}{\sum_{i=1}^n w_{ji}^{(m+1)}}$$

$$\hat{v}(m_j) = \frac{-\sum_{i=1}^n \hat{w}_{ji}^{(m+1)} \log \frac{x_i}{\hat{y}_j^{(m+1)}}}{\sum_{i=1}^n \hat{w}_{ji}^{(m+1)}}$$