

Model: $f(\underline{x}) = \lambda_1 f_1(\underline{x}, \theta_1) + \dots + \lambda_k f_k(\underline{x}, \theta_k)$

Ex: Mixtures of multivariate normals

$$f_j(\underline{x}, \mu_j, \Sigma_j) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_j|^{\frac{1}{2}}} \exp(-\frac{1}{2} (\underline{x} - \mu_j)^T \Sigma_j^{-1} (\underline{x} - \mu_j))$$

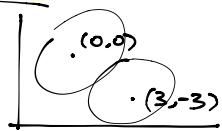
$\hat{\Delta}_{ij}$ est'd prob that \underline{x}_i belongs to cluster j.

Given $\{\hat{\Delta}_{ij}\}$ we have $\hat{n}_j = \frac{\sum_{i=1}^n \hat{\Delta}_{ij} x_i}{\sum_{i=1}^n \hat{\Delta}_{ij}}$ ← estimate of # of \underline{x}_i in cluster j.

$$\hat{\Sigma}_j = \frac{1}{\sum_{i=1}^n \hat{\Delta}_{ij}} \sum_{i=1}^n \hat{\Delta}_{ij} (\underline{x}_i - \hat{\mu}_j)(\underline{x}_i - \hat{\mu}_j)^T$$

⇒ repeat until convergence.

Toy example: $p=2, k=2$



$$\Sigma_1 = \begin{pmatrix} 1 & 1 \\ 1 & 2 \end{pmatrix} \quad \Sigma_2 = \begin{pmatrix} 1 & -1 \\ -1 & 2 \end{pmatrix}$$

$n=400$ (200 each from the 2 popl's)

EM R function: very slow & inefficient

- use k-means clustering to give initial values

$$\hat{\mu}_1 = \begin{pmatrix} 0.19 \\ -0.20 \end{pmatrix}$$

$$\hat{\Sigma}_1 = \begin{pmatrix} 1.16 & 0.62 \\ 0.62 & 1.76 \end{pmatrix}$$

$$\hat{\mu}_2 = \begin{pmatrix} 3.25 \\ -3.38 \end{pmatrix}$$

$$\hat{\Sigma}_2 = \begin{pmatrix} 0.84 & -0.74 \\ -0.74 & 1.52 \end{pmatrix}$$

- after 50 iterations of EM

$$\hat{\mu}_1 = \begin{pmatrix} -0.01 \\ -0.07 \end{pmatrix}$$

$$\hat{\Sigma}_1 = \begin{pmatrix} 0.97 & 0.93 \\ 0.93 & 1.82 \end{pmatrix}$$

$$\hat{\mu}_2 = \begin{pmatrix} 2.98 \\ -3.03 \end{pmatrix}$$

$$\hat{\Sigma}_2 = \begin{pmatrix} 1.15 & -1.14 \\ -1.14 & 2.00 \end{pmatrix}$$

Classification:

Regression with a categorical (or group) response.

So far: Unsupervised learning → given multivariate data $\underline{x}_1, \dots, \underline{x}_n$ from some popl's which may consist of sub-popl's
- no info about sub-popl's.

Tools: PCA } Identify interesting low dimensional projections
ICA }

Cluster Analysis: identify clusters within data

Supervised learning: For each observation \underline{x}_i we have a group/cluster/sub-popl identifier g_i .

⇒ observe $(g_1, \underline{x}_1), \dots, (g_n, \underline{x}_n)$

- Assume that g_1, \dots, g_n take k known values/variables.

Ex: ① Species identification/classifications

- given measurements \underline{x}_i , determine species

② spam filter

- given attributes (features) of an e-mail, determine whether it is "good" or spam

Model: (G, \underline{X}) has a joint distribution described by $P(G=j) = \lambda_j$ for $j=1, \dots, k$.

and conditional on $G=j$, \underline{X} has density $f_j(\underline{x})$

Given $\underline{X}=\underline{x}$, determine conditional dist'n of G i.e. $P(G=j | \underline{X}=\underline{x})$

Marginal density of \underline{X} :

$$f(\underline{x}) = \lambda_1 f_1(\underline{x}) + \dots + \lambda_k f_k(\underline{x})$$

Data problem: Given data $(\underline{g}_1, \underline{X}_1), \dots, (\underline{g}_n, \underline{X}_n)$. Determine a classification rule with small error rate.



Rule: If $\underline{x} \in R_j$ then classify \underline{x} as belonging to group j .

Choose R_1, \dots, R_k to minimize classification error rate.

$$\begin{aligned} \text{error rate} &= P(\text{error}) = \sum_{j=1}^k P(G=j) P(\text{error} | G=j) = \sum_{j=1}^k \lambda_j P(\underline{X} \in R_j | G=j) \\ &= \sum_{j=1}^k \lambda_j \left\{ 1 - \int_{R_j} f_j(\underline{x}) d\underline{x} \right\} \\ &= 1 - \underbrace{\sum_{j=1}^k \lambda_j \int_{R_j} f_j(\underline{x}) d\underline{x}}_{P(\text{correct classification})} \end{aligned}$$

Optional classification rule: (Bayes classifier)

Classify \underline{x} to group j if for all $i \neq j$, $\lambda_j f_j(\underline{x}) > \lambda_i f_i(\underline{x})$

$$\frac{f_j(\underline{x})}{f_i(\underline{x})} > \frac{\lambda_i}{\lambda_j} \text{ for all } i \neq j$$

Proof for $k=2$: $R_2 = R_1^c$

$$\begin{aligned} \text{Want to maximize } &\lambda_1 \int_{R_1} f_1(\underline{x}) d\underline{x} + \lambda_2 \int_{R_1^c} f_2(\underline{x}) d\underline{x} \\ &= \lambda_1 \int_{R_1} f_1(\underline{x}) d\underline{x} + \lambda_2 \left(1 - \int_{R_1} f_2(\underline{x}) d\underline{x} \right) \\ &= \int_{R_1} [\lambda_1 f_1(\underline{x}) - \lambda_2 f_2(\underline{x})] d\underline{x} + \lambda_2 \\ &\geq 0 \end{aligned}$$

- maximized over R_1 for $R_1 = \{\underline{x} : \lambda_1 f_1(\underline{x}) - \lambda_2 f_2(\underline{x}) > 0\}$
 $= \{\underline{x} : \lambda_1 f_1(\underline{x}) > \lambda_2 f_2(\underline{x})\}$

Note: In this proof, we don't use the fact that $\lambda_1 + \lambda_2 = 1$

In general, result holds for any non-negative $\lambda_1, \dots, \lambda_k \Rightarrow$ minimizing weighted error rate.

- useful if different misclassifications have different costs

- avoid identifying good emails as spam, for example.

Application: $f_1(\underline{x}) \dots f_k(\underline{x})$ are multivariate normal with different means μ_1, \dots, μ_k but equal covariance matrices (call it C)

$$\text{Then } f_j(\underline{x}) = \frac{1}{(2\pi)^{\frac{D}{2}} |C|^{\frac{1}{2}}} \exp \left[-\frac{1}{2} (\underline{x} - \mu_j)^T C^{-1} (\underline{x} - \mu_j) \right]$$

optimal classification rule: $\underline{x} \rightarrow j$ if $\lambda_j f_j(\underline{x}) > \lambda_i f_i(\underline{x})$ for all $i \neq j$

$$\text{or } \underbrace{\ln(f_j(\underline{x})/f_i(\underline{x}))}_{\text{II}} > \ln(\lambda_j/\lambda_i)$$

$$\begin{aligned} \ln f_j(\underline{x}) - \ln f_i(\underline{x}) &= -\frac{1}{2} (\underline{x} - \mu_j)^T C^{-1} (\underline{x} - \mu_j) + \frac{1}{2} (\underline{x} - \mu_i)^T C^{-1} (\underline{x} - \mu_i) \\ &= [\underline{x}^T C^{-1} \mu_j - \frac{1}{2} \mu_j^T C^{-1} \mu_j - (\underline{x}^T C^{-1} \mu_i - \frac{1}{2} \mu_i^T C^{-1} \mu_i)] \end{aligned}$$

Define discriminant scores:

$$\begin{aligned} \text{linear function of } \underline{x}: d_j(\underline{x}) &= \underline{x}^T C^{-1} \mu_j - \frac{1}{2} \mu_j^T C^{-1} \mu_j + \ln(\lambda_j) \\ \Rightarrow \text{classify } \underline{x} \rightarrow j &\text{ if } d_j(\underline{x}) > d_i(\underline{x}) \text{ for all } i \neq j \end{aligned}$$

Special case: $k=2$ depends on C, μ_1, μ_2 and λ_1, λ_2
Suffices to look at $d_2(\underline{x}) - d_1(\underline{x}) = \beta_0 + \underline{x}^T \beta_1$

Binary regression models

Ex: logistic regression $\theta(\underline{x}) = P(\text{group 2} | \underline{x})$

$$\ln \left(\frac{\theta(\underline{x})}{1 - \theta(\underline{x})} \right) = \beta_0 + \underline{x}^T \beta$$

- suggests close relationship between linear discriminant analysis + binary regression models.

How to apply to data?

Given $(g_1, \underline{x}_1), \dots, (g_n, \underline{x}_n)$ can estimate μ_1, \dots, μ_k and C as follows:

$$\hat{\mu}_j = \sum_{i=1}^n \underline{x}_i \underbrace{I(g_i=j)}_{\substack{\text{if } g_i=j \\ 0 \text{ o.w.}}} / \sum_{i=1}^n I(g_i=j)$$

$$\hat{C} = \underbrace{\frac{1}{n-k} \sum_{i=1}^n}_{\substack{\text{can also put here} \\ \uparrow}} (\underline{x}_i - \hat{\mu}_j(i)) (\underline{x}_i - \hat{\mu}_j(i))^T$$

group associated with \underline{x}_i

Question: ① How to estimate $\lambda_1, \dots, \lambda_k$?

- assume $\lambda_1, \dots, \lambda_k$ are known (often true) or est'd "out of sample"
- allows biased samplings

- estimate $\lambda_1, \dots, \lambda_k$ from sample proportions

$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n I(g_i=j)$$

- need to assume that $(g_1, \underline{x}_1), \dots, (g_n, \underline{x}_n)$ is some sort of representation sample from population.

② How to estimate error rate for a given classifier estimated from $(g_1, \underline{x}_1), \dots, (g_n, \underline{x}_n)$?

Naive method: Resubstitution estimate for each x_i in sample, we obtain classification $\hat{g}_1, \dots, \hat{g}_n$
error rate = $\frac{1}{n} \sum_{i=1}^n I(g_i \neq \hat{g}_i)$ = proportion of misclassified observations.

Resubstitution error rate tends to be biased downwards.

How to fix this?

- cross-validation: leave out observations from sample, estimate classifier from remaining observed and look at error rate for observations not used to estimate the classifier
e.g. leave-one-out cross-validation.