

## Quadratic discriminant analysis (QDA)

LDA:  $f_1(x), \dots, f_k(x) \leftarrow$  multivariate normal means  $\mu_1, \dots, \mu_k$  and covariance  $C$  sub-populations

- optimal classifications rule  $\rightarrow$  linear boundaries
- now allow different covariance matrices
  - get boundaries determined by paraboloids

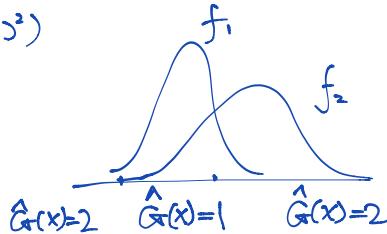
Example:  $p=1, k=2 (\lambda_1=\lambda_2=\frac{1}{2})$

$$f_1(x) = \frac{1}{\sigma_1 \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_1^2}(x-\mu_1)^2\right)$$

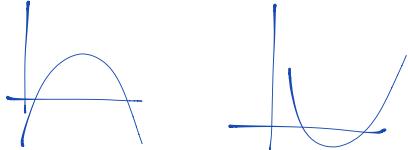
$$f_2(x) = \frac{1}{\sigma_2 \sqrt{2\pi}} \exp\left(-\frac{1}{2\sigma_2^2}(x-\mu_2)^2\right)$$

$$\hat{G}(x)=1 \text{ if } f_1(x) > f_2(x)$$

$$\ln\left(\frac{f_1(x)}{f_2(x)}\right) > 0$$



$$\begin{aligned} \ln\left(\frac{f_1(x)}{f_2(x)}\right) &= \frac{1}{2\sigma_2^2}(x-\mu_2)^2 - \frac{1}{2\sigma_1^2}(x-\mu_1)^2 + \ln(\sigma_2/\sigma_1) \\ &= \left[\frac{1}{2\sigma_2^2} - \frac{1}{2\sigma_1^2}\right]x^2 - \left(\frac{\mu_2}{\sigma_2^2} - \frac{\mu_1}{\sigma_1^2}\right)x + \frac{\mu_2^2}{\sigma_2^2} - \frac{\mu_1^2}{\sigma_1^2} + \ln(\sigma_2/\sigma_1) \end{aligned}$$



QDA in practice: need to estimate  $\mu_1, \dots, \mu_k$  and  $C_1, \dots, C_k$

$$\hat{\mu}_j = \frac{\sum_{i=1}^n x_i I(g_i=j)}{\sum_{i=1}^n I(g_i=j)} \quad (j=1, \dots, k)$$

$$\hat{C}_j = \sum_{i=1}^n (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T I(g_i=j) / (\sum_{i=1}^n I(g_i=j) - 1)$$

Example (on Blackboard) Iris data

3 species of iris

4 variables length of sepal

width

length of petal

width

- estimate misclassification rate using leave-one-out cross-validation: error rate  
 $= \frac{4}{150}$

(continued)

- generate 15000 observations from 3 sub-populations using the estimated mean vectors  $\hat{\mu}_1, \hat{\mu}_2, \hat{\mu}_3$  and covariance matrices  $C_1, C_2, C_3$
- classification rate every similar: 1.2% of simulated observations are classified differently.

## Alternative methods for classification

So far: have data  $(g_1, \mathbf{x}_1), \dots, (g_n, \mathbf{x}_n)$

Model:  $(G, X) \rightarrow P(G = j) = \lambda_j$  and conditional on  $G = j, X$  has density  $f_j(\mathbf{x})$ .

$$P(G = j | X = \mathbf{x}) = \frac{\lambda_j f_j(\mathbf{x})}{\sum_{l=1}^k \lambda_l f_l(\mathbf{x})}$$

quantity of interest

LDA, QDA: Assume  $f_1(\mathbf{x}), \dots, f_k(\mathbf{x})$  multivariate normal densities -- use data to estimate unknowns.

**Key point:** Explicitly model distributions of  $X$ .

But in practice, this is difficult to do

- discrete variables
- $P$  may be very large

**Alternative approach:** model  $P(G = j | X = \mathbf{x})$  directly.

- analogous to regression modeling
  - $G$  is the response
  - $X$  is the predictors
- we implicitly assume that the distribution of  $X$  (i.e. not conditional on  $G = j$ ) is not particularly informative.

Multiple regression:  $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i$  --> look at conditional distributions of response given predictors.

**Special case:**  $k = 2$

$$P(G = 1 | X = \mathbf{x}) = 1 - P(G = 2 | X = \mathbf{x})$$

where the LHS is  $\theta(\mathbf{x}) = g(\mathbf{x}, \boldsymbol{\beta})$ , and  $\boldsymbol{\beta}$  is unknown parameters.

## Logistic regression model

$$\theta(\mathbf{x}) = \frac{\exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta})}{1 + \exp(\beta_0 + \mathbf{x}^T \boldsymbol{\beta})}$$

$$1 - \theta(x) = P(G = 2 | X = x) = \frac{1}{1 + \exp(\beta_0 + x^T \beta)}$$

Note that  $0 < \theta(x) < 1$  for any  $\beta_0, \beta$

- logit transform:

$$\text{logit}(\theta(x)) = \ln\left(\frac{\theta(x)}{1-\theta(x)}\right) = \beta_0 + x^T \beta \quad \begin{cases} > 0 \Rightarrow \theta(x) > \frac{1}{2} (\hat{G}=1) \\ < 0 \Rightarrow \theta(x) < \frac{1}{2} (\hat{G}=2) \end{cases}$$

Estimation: Given  $(g_1, x_1), \dots, (g_n, x_n)$  the likelihood function is

$$\begin{aligned} L(\beta_0, \beta) &= \prod_{i=1}^n \left[ \theta(x_i) \overset{I(g_i=1)}{y_i} (1 - \theta(x_i))^{1 - I(g_i=0)} \overset{1-y_i}{\circ} \right] \\ &= \prod_{i=1}^n \left\{ \left[ \frac{\theta(x_i)}{1 - \theta(x_i)} \right]^{y_i} (1 - \theta(x_i))^{1-y_i} \right\} \end{aligned}$$

$$\ln L(\beta_0, \beta) = \sum_{i=1}^n \left\{ y_i (\beta_0 + x_i^T \beta) - \ln(1 + \exp(\beta_0 + x_i^T \beta)) \right\}$$

$$\frac{\partial}{\partial \beta_0} \ln L(\beta_0, \beta) = \sum_{i=1}^n \left\{ y_i - \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)} \right\}$$

$$\frac{\partial}{\partial \beta} \ln L(\beta_0, \beta) = \sum_{i=1}^n \left\{ y_i - \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)} \right\} x_i$$

Set derivatives equal to 0 and solve for MLEs

R function: `glm(group ~ x1 + x2 + ... , family = binomial)`

Example: Generate data from 2 bivariate normal distributions

$$f(x) = \frac{1}{4} f_1(x) + \frac{3}{4} f_2(x) \leftarrow N_2\left(\begin{pmatrix} -1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}\right)$$

$\uparrow$   
 $N_2\left(\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}\right)$        $n=200$

① Use LDA (optimal for those data)

② Use logistic regression

Error rate using LOO CV:

LDA: 17/200

Logistic reg: 18/200

LDA

		Classified		45
observed		1	2	
		37	8	
1		9	146	155
2				

Logistic Reg

		Classified		45
observed		1	2	
		35	10	
1		8	147	155
2				

Question: Is this similarity a coincidence?

Not really. For multivariate normal model we have:

$$\begin{aligned}\theta(\underline{x}) &= P(G=1 | \underline{X}=\underline{x}) \\ &= \frac{\lambda_1 f_1(\underline{x})}{\lambda_1 f_1(\underline{x}) + \lambda_2 f_2(\underline{x})} \\ &= \frac{\exp(\beta_0 + \underline{x}^T \underline{\beta})}{1 + \exp(\beta_0 + \underline{x}^T \underline{\beta})} \end{aligned}$$

depend on  $\lambda_1, \lambda_2, f_1, f_2$ , and  $C$

Also MLEs in logistic regression satisfy the equation

$$\begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta} \end{pmatrix} = (\underline{X}^T \underline{W} \underline{X})^{-1} \underline{X}^T \underline{W} \underline{z}$$

$$\text{where } \underline{X} = \begin{pmatrix} 1 & \underline{x}_1^T \\ \vdots & \vdots \\ 1 & \underline{x}_n^T \end{pmatrix} \quad \underline{W} = \begin{pmatrix} w_1 & & 0 \\ & \ddots & \\ 0 & & w_n \end{pmatrix} \quad \text{with } w_i = [\hat{\theta}(\underline{x}_i)(1 - \hat{\theta}(\underline{x}_i))]^{-1}$$

$$\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} \quad \text{with } z_i = \ln\left(\frac{\hat{\theta}(\underline{x}_i)}{1 - \hat{\theta}(\underline{x}_i)}\right) + \frac{y_i - \hat{\theta}(\underline{x}_i)}{\hat{\theta}(\underline{x}_i)(1 - \hat{\theta}(\underline{x}_i))}$$

$$\text{Boundary: } \hat{\beta}_0 + \underline{x}^T \underline{\beta} = 0$$

If  $\hat{\theta}(\underline{x}_i)(1 - \hat{\theta}(\underline{x}_i))$  is approximately constant then boundaries defined by logistic regression are close to those for LDA.

### Multinomial logit

Need to find a model for

$$\theta_j(\underline{x}) = P(G=j | \underline{X}=\underline{x}) \quad \text{for } j=1, \dots, k>2$$

Note:  $\theta_1(\underline{x}) + \theta_2(\underline{x}) + \dots + \theta_k(\underline{x}) = 1$  for each  $\underline{x}$

Assume that for each  $j=1, \dots, k-1$ , we have  $\ln\left(\frac{\theta_j(\underline{x})}{\theta_k(\underline{x})}\right) = \beta_{0j} + \underline{x}^T \underline{\beta}_j$

$$\frac{\theta_j(\underline{x})}{\theta_k(\underline{x})} = \exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)$$

$$\theta_j(\underline{x}) = \theta_k(\underline{x}) \exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)$$

Now sum RHS, LHS over  $j$  from 1 to  $k-1$

$$\sum_{j=1}^{k-1} \theta_j(\underline{x}) = 1 - \theta_k(\underline{x}) = \theta_k(\underline{x}) \sum_{j=1}^{k-1} \exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)$$

$$\Rightarrow \theta_k(\underline{x}) = \frac{1}{1 + \sum_{j=1}^{k-1} \exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)}$$

$$\Rightarrow \theta_j(\underline{x}) = \frac{\exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)}{1 + \sum_{j=1}^{n-1} \exp(\beta_{0j} + \underline{x}^T \underline{\beta}_j)}$$

- In R:
- ① package "mlogit" → allows you to fit very complicated multi-logits models
  - ② Package "nnet": function multinom is a simpler option for this course.

Comments: Model  $(G, X)$  vs model  $G | X = \underline{x}$ .

① Which is better?

- Modeling  $(G, X)$  is better if you have a good model for  $(G, X)$
- Otherwise modeling  $G | X = \underline{x}$  gives more flexibility

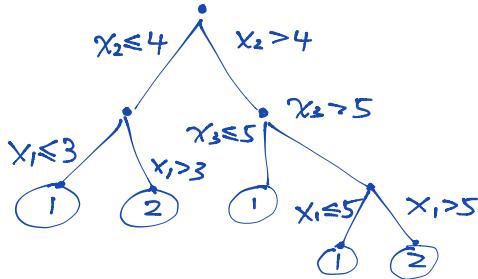
② Modeling  $P(G=j | X = \underline{x})$  allows to apply many regression modeling techniques

- variables selection (subject selection)
- also use shrinkage / Bayesian methodology

### Classification trees

Idea: Use data to construct a binary tree to classify observations.

e.g.  $k=2, p=3$  ( $x_1, x_2, x_3$ )



Advantages:

- very easy to interpret
- adapt to different types of variables (cont., discrete, categorical, ordinal)

Disadvantages:

- Computationally complex
- When to stop growing tree ?
- No clear model