

(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  $f_j(X)$ .

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

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

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

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

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

- **logit transform:**