Observe $(g_1, x_1), \dots, (g_n, x_n)$

$g_i$ takes values in $\{1, \dots, k\}$

<u>Model</u>: Dist'n of $(G, \underline{X})$. $P(G = j) = \lambda_j$

Conditional density of $\underline{X}$ given $G = j$. $f_j(\underline{x})$

Marginal density of $\underline{X}$ is $f(\underline{x}) = \lambda_1 f_1(\underline{x}) + \lambda_2 f_2(\underline{x}) + \dots + \lambda_k f_k(\underline{x})$

==Optimal classification:== 2 approaches

① Minimize $P(\text{error})$

    Given $\underline{x}$, clarify as $\hat{G}(\underline{x}) = j$ if $\lambda_j f_j(\underline{x}) > \lambda_i f_i(\underline{x})$ for all $i \neq j$

    i.e. $R_j = \{\underline{x}: \lambda_j f_j(\underline{x}) > \lambda_i f_i(\underline{x})$ for all $i \neq j\}$

② Look at conditional dist'n of $G$ given $\underline{X} = \underline{x}$

    Bayes Thm: $P(G = j \mid \underline{X} = \underline{x}) \overset{"="}{=} \dfrac{P(G = j, \underline{X} = \underline{x})}{P(\underline{X} = \underline{x})}$

$$= \frac{P(G = j) P(\underline{X} = \underline{x} \mid G = j)}{P(\underline{X} = \underline{x})}$$

$$= \frac{\lambda_j f_j(\underline{x})}{\lambda_1 f_1(\underline{x}) + \dots + \lambda_k f_k(\underline{x})}$$

$\Rightarrow$ if we choose $\hat{G}(\underline{x})$ to maximize $P(G = j \mid \underline{X} = \underline{x})$

    we get the classification rule in ①

But $P(G = j \mid \underline{X} = \underline{x})$ provides more information.

<u>Example</u> $k = 2$, $p = 1$, $\lambda_1 = \lambda_2 = \frac{1}{2}$

$f_1(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x-1)^2)$   ($N(1,1)$)

$f_2(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x+1)^2)$   ($N(-1,1)$)



① givens $\hat{G}(x) = 1$ if $x > 0$

        $\hat{G}(x) = 2$ if $x < 0$

② $P(G = 1 \mid X = x) = \dfrac{f_1(x)}{f_1(x) + f_2(x)}$

| $x$ | $P(G = 1 \mid X = x)$ |
|---|---|
| $-1$ | 0.119 |
| $-0.5$ | 0.269 |
| $-0.25$ | 0.378 |
| 0 | 0.500 |
| 0.25 | 0.622 |
| 0.5 | 0.731 |
| 1 | 0.881 |

==Linear discriminant analysis (LDA)==

$f_1(\underline{x}), \dots, f_k(\underline{x})$ multivariate normal with distinct means $\underline{\mu}_1, \dots, \underline{\mu}_k$ and common covariance matrix $C$.

If everything is known then $\hat{G}(\underline{x}) = j$ if $d_j(\underline{x}) > d_i(\underline{x})$ for all $i \neq j$ where

    $d_j(\underline{x}) = \underline{x}^T C^{-1} \underline{\mu}_j - \frac{1}{2} \underline{\mu}_j^T C^{-1} \underline{\mu}_j + \ln(\lambda_j)$

Given data $(g_1, x_1), \ldots, (g_n, x_n)$ we have

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

$$\hat{C} = \frac{1}{n-k} \sum_{i=1}^{n} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$$

Example: Fisher's iris data (see Blackboard)

3 species of irises
{ virginica
  setosa
  versciolor }

4 variables
{ sepal length
  sepal width
  petal length
  petal width }

- LDA works very well here.
- But, could prob do as well from pairwise scatterplots!

How to estimate misclassification rate?
(as well as posterior dist'n)
- resubstitution estimate is typically biased downwards.
  - bias increases often severely, as model complexity increases.

Solution: Do some sort of cross-validation
    - divide data into 2 sets: training set (n-m obs) and test set (m obs)
    - estimate classification rate based on training data and use test data to estimate misclassification rate.

For example. 10-fold cross-validation
- divide data into 10 sets
- successively leave out one set → test data
        use other 9 as training data

- leave-one-out CV: for each "observation" $x_i$.
    we compare $\hat{G}_{-i}(x_i)$ where $\hat{G}_{-i}$ is the classification rate using all data except $x_i$.
        ⇒ used in R function: lda, qda

For the iris data: est'd misclassification rate $= \frac{3}{150}$  (assume $\lambda_1 = \lambda_2 = \lambda_3 = \frac{1}{3}$)
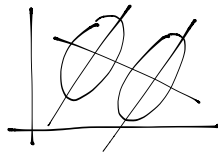
$= 2\%$

Est'd misclassification rate $= \dfrac{\sum_{i=1}^{n} \lambda_{g_i} I(\hat{G}_{-i}(x_i) \neq g_i)}{\sum_{i=1}^{n} \lambda_{g_i}}$

$\rightarrow 1. - 0. - 0.$

Can also estimate $P(G = j \mid x_i)$ use LOO CV
    - assume model is approx correct.

$k=2$, $p=2$, $C_1 = C_2 = C$

level curves of $f_1(\underline{x})$ & $f_2(\underline{x})$ have same orientation

$\rightarrow$ l.c. of $f_2$

l.c. of $f_1$

$\Rightarrow$ no linear boundaries for classification rule.

– general form of classification rule is the same but regions more complicated