Want to model $\theta_j(\underline{x}) = P(G=j \mid \underline{X}=\underline{x})$ for $j=1,\cdots,k$

$k=2,\ p=3$ — root node

$X_2 \leq 4$    $X_2 > 4$

$X_1 \leq 3$   $X_1 > 3$   $X_3 \leq 5$   $X_2 > 5$

  ①   ②   ①   $X_1 \leq 5$   $X_1 > 5$

    ①   ②

(Question): How to construct (grow) tree?

Recursive partitioning

At a given node, define
$$B = \{(g_i, \underline{x}_i) : i \in I(B)\}$$
$\quad\quad\quad\quad\quad\quad 1, \cdots, k$

For $\underline{x}_i \in B$, we have
$$\hat{\theta}_j(B) = \frac{1}{n(B)} \underbrace{\sum_{i \in I(B)} I(g_i=j)}_{= n_j(B)} = \frac{n_j(B)}{n(B)}$$

$\underset{\text{\# of observations in B}}{}$

$= $ proportions of group $j$ in node $B$.

Now define possible new nodes: $B_1, B_2$ with $B = B_1 \cup B_2$

$$\hat{\theta}_j(B_1) = \frac{1}{n(B_1)} \sum_{i \in I(B_1)} I(g_i=j) = \frac{n_j(B_1)}{n(B_1)} \quad \overset{\text{disjoint}}{\phantom{x}} \quad (j=1,\cdots,k)$$

$$\hat{\theta}_j(B_2) = \frac{1}{n(B_2)} \sum_{i \in I(B_2)} I(g_i=j) = \frac{n_j(B_2)}{n(B_2)}$$

Define
$$D(B_1, B_2 : B) = \sum_{j=1}^{k} \left\{ n_j(B_1) \ln\left(\frac{n_j(B_1)}{n(B_1)}\right) + n_j(B_2) \ln\left(\frac{n_j(B_2)}{n(B_2)}\right) - n_j(B) \ln\left(\frac{n_j(B)}{n(B)}\right) \right\}$$

      ↑ variables    ↑ fixed

- Find $\underline{B_1 + B_2}$ s.t. $B_1 \cup B_2 = B$ to maximize $D$.
  $\underset{\text{many choices}}{}$

- restrict maximization to simple one variable splits

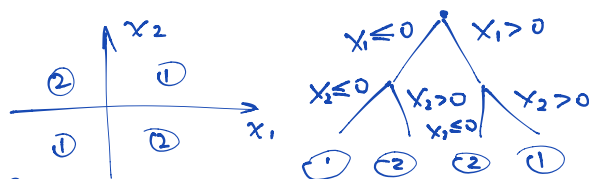$$B_1 = \{(g_i, \underline{x}_i) \in B, x_{i\ell} \leq d\} \leftarrow \text{threshold}$$
$$B_2 = \{(g_i, \underline{x}_i) \in B, x_{i\ell} > d\}$$
$$\uparrow$$
$$\ell = 1, \cdots, p$$

- also require other constraints e.g. $n(B_1), n(B_2) \geq 5$.

Does procedure always work?
Example: $k=2, p=2$



How does tree algorithm work here?
- $P(G=j \mid X_1) = \frac{1}{2}$ (for $j=1,2$)
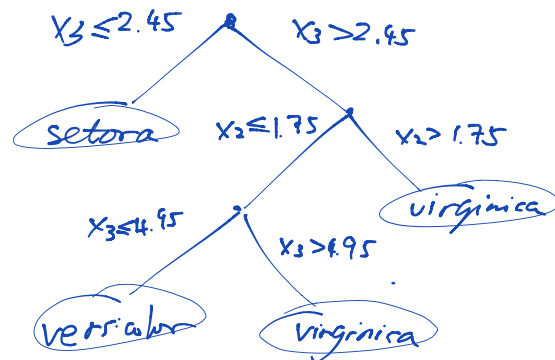  $P(G=j \mid X_2) = \frac{1}{2}$

- recursive partitioning algorithm has trouble getting started

— better performance if allow splitting along linear combinations (i.e. projections) of variables.

<u>Example</u> Iris data

3 species $\begin{cases} \text{setora} \\ \text{virginica} \\ \text{versicolor} \end{cases}$

4 vars $\begin{cases} X_1 = \text{sepal length} \\ X_2 = \quad\quad\text{width} \\ X_3 = \text{petal length} \\ X_4 = \quad\quad\text{width} \end{cases}$

— compare favourably to LDA
error rate = 4/150 (tree)
error rate = 3/150 (LDA)



## Regression models for multivariate data
— repeated measures

## Multivariate Analysis of Variance (MANOVA)

<u>Problem</u> k treatments (or groups)
$n_i$ subjects in treatment $i$
Multivariate response $\quad X_{ij} \begin{cases} i = 1, \cdots, k \\ j = 1 \cdots, n_i \end{cases}$
p vectors

<u>Model</u> : $X_{ij} = \mu_i + \varepsilon_{ij}$ $\quad i = 1, \cdots, k$ , $j = 1, \cdots, n_i$
where $[\varepsilon_{ij}]$ are independent $N_p(\underline{0}, C)$ random vectors i.e. $X_{ij} \sim N_p(\mu_i, C)$

One question of interest : Is there a difference between k treatments ? (For example, is $\mu_1 = \mu_2 = \cdots = \mu_k$?)

Univariate (p=1) case : Decompose total sum of squares

$$SS_{Total} = \sum_{i=1}^{k} \sum_{j=1}^{n} (x_{ij} - \bar{x})^2 = \sum_{i=1}^{k} n_i (\bar{x}_i - \bar{x})^2 + \sum_{i=1}^{k} \sum_{j=1}^{n} (x_{ij} - \bar{x}_i)^2$$

↓ overall sample mean $\quad$ ↓ sample mean for treatments

$$SS_{Total} = SS_{between\ group} + SS_{within\ group}$$

To test $H_0 : \mu_1 = \cdots = \mu_k$ , we compare $SS_{between}$ to $SS_{within}$

Test statistic $\quad F = \dfrac{SS_{between} (k-1)}{SS_{within} (n-k)} \sim F_{k-1, n-k}$ under $H_0$

$$SS_{Total} = \sum_{i=1}^{k} \sum_{j=1}^{n} (X_{ij} - \bar{X})(X_{ij} - \bar{X})^T = \underbrace{\sum_{i=1}^{k} n_i (\bar{X}_i - \bar{X})(\bar{X}_i - \bar{X})^T}_{SS_{between}} + \underbrace{\sum_{i=1}^{k} \sum_{j=1}^{n} (X_{ij} - \bar{X}_i)(X_{ij} - \bar{X}_i)^T}_{SS_{within}}$$

<u>Question:</u>

How to compare $SS_{between}$ to $SS_{within}$