

- Optimise cost at every split:  
greedily.
- Continue until a stopping criterion is met:
  - Plateau in optimisation
  - small # of points in region
  - using all predictors

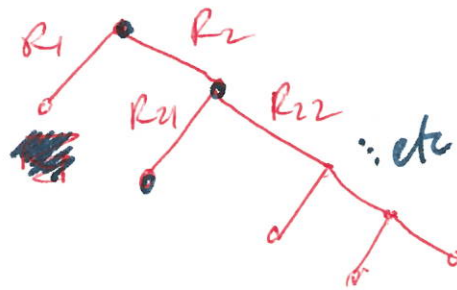
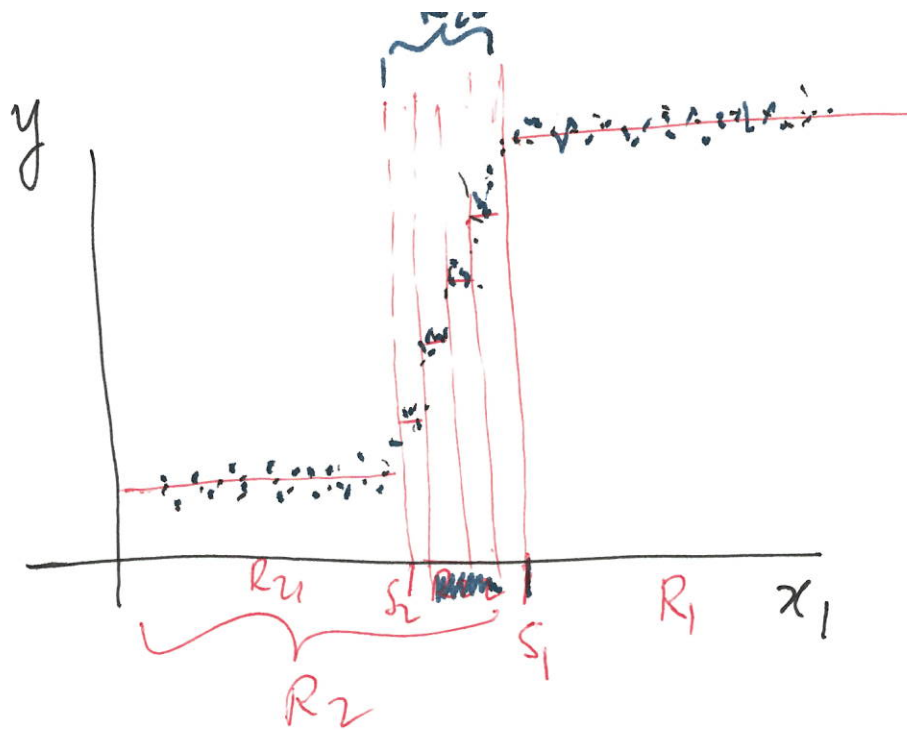
Given  $\vec{x}^{in}$

Find  $R_{\alpha} / \vec{x}^{in} \in R_{\alpha}$

$$\hat{y}^{DT} = \bar{y}_{R_{\alpha}}$$

Mean of  
the region  
where  $\vec{x}^{in}$  falls.

(TBC...)



Example of how a DT would proceed with ~~a tree~~ the case of one descriptor

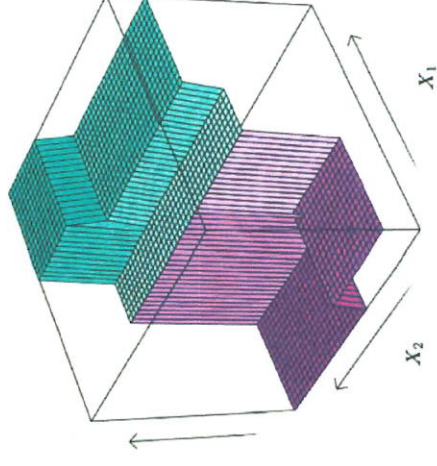
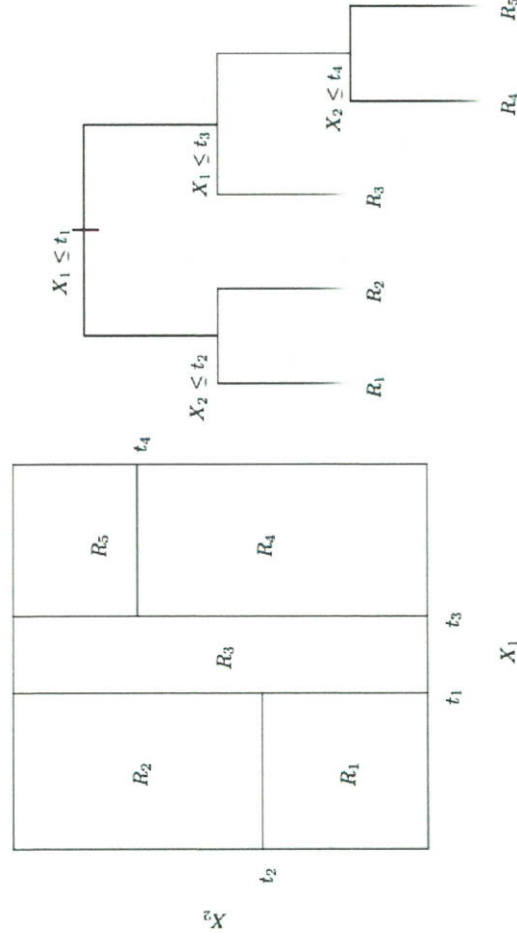
$$\left\{ \begin{array}{l} \vec{x} = x_1 \in \mathbb{R} \\ \text{and} \\ y \in \mathbb{R} \end{array} \right\}$$

# Expressiveness of Decision Trees

---

We've seen that classification trees approximate boundaries in the feature space that separate classes.

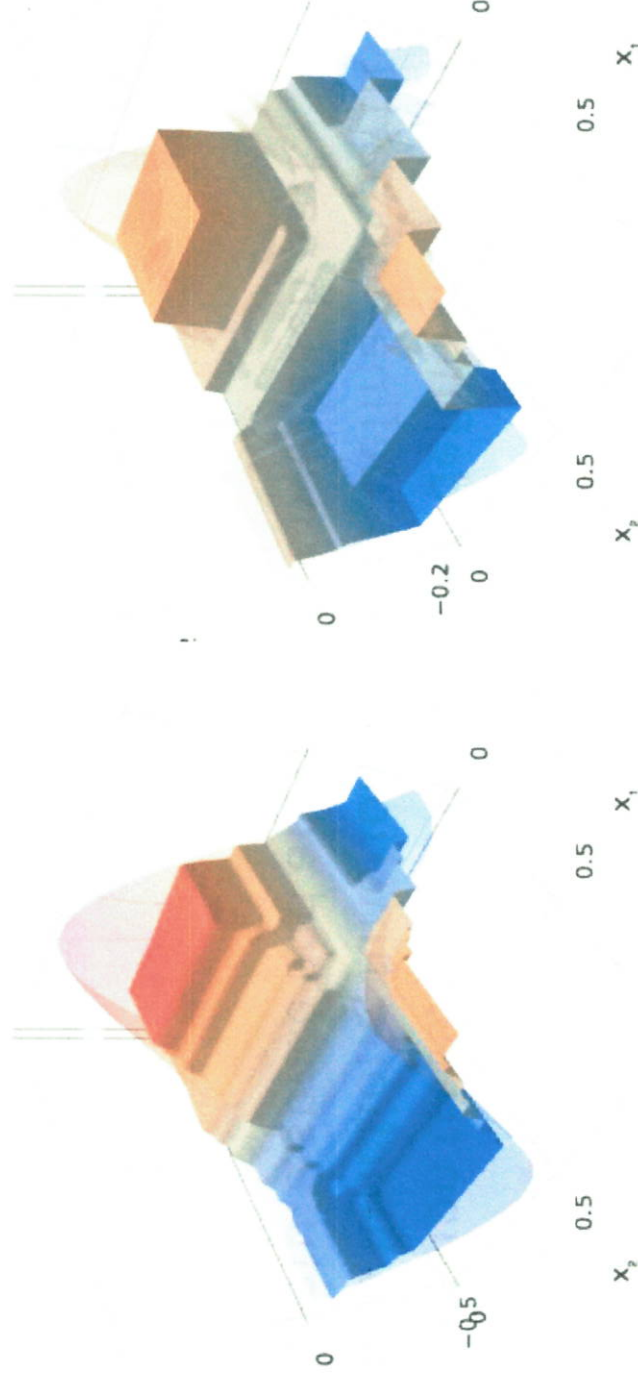
Regression trees, on the other hand, define **simple functions** or step functions, functions that are defined on partitions of the feature space and are constant over each part.



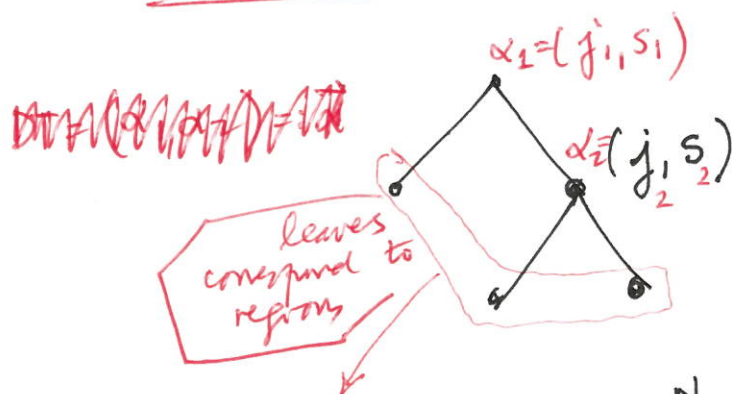
# Expressiveness of Decision Trees

---

For a fine enough partition of the feature space, these functions can approximate complex non-linear functions.



# Decision trees for classification:



Same as above:  
Each split denoted by  $(j, s) = \alpha_i$

For region  $R_{\vec{\alpha}}$  :

$$[T(R_{\vec{\alpha}})]_q = \frac{\sum_{i=1}^N I(\vec{x}^{(i)} \in R_{\vec{\alpha}} \& y^{(i)} \in C_q)}{\sum_{i=1}^N I(\vec{x}^{(i)} \in R_{\vec{\alpha}})}$$

$q = 1, \dots, Q$

$$\vec{T}(R_{\vec{\alpha}})_{Q \times 1} = \begin{bmatrix} \vdots \\ T(R_{\vec{\alpha}})_q \\ \vdots \end{bmatrix}$$

Probability vector of belonging to each class  $q$  in  $R_{\vec{\alpha}}$

Given  $\vec{x}^{in}$  :

(1) Find  $R_{\vec{\alpha}} \text{ s.t. } \vec{x}^{in} \in R_{\vec{\alpha}}$

(2) Obtain  $\vec{T}(R_{\vec{\alpha}})$

(3)  $\hat{y} = \hat{f}^{DT}(\vec{x}^{in}) = \arg \max_q \vec{T}(R_{\vec{\alpha}})(\vec{x}^{in})_q$



Important tweaks:

What is ~~the~~ Cost function for the splits in the decision tree:

(1) Contingency table:

Minimise the error rate.

Not sensitive

(2)  $\left[ \begin{array}{c} \text{Gini index} \\ \text{or} \\ \text{Cross-entropy} \end{array} \right]$  Cost function  
deals with how informative the split is:

Gini index of  $\vec{\pi}(R_{\vec{x}})$

$$GI[\vec{\pi}(R_{\vec{x}})] = \sum_{q=1}^Q \pi(R_{\vec{x}})_q (1 - (\pi(R_{\vec{x}})_q))$$

Cross entropy

$$CE[\vec{\pi}(R_{\vec{x}})] = \sum_{q=1}^Q \pi(R_{\vec{x}})_q \log(1 - \pi(R_{\vec{x}})_q)$$

Two versions of information (or entropy)

Properties of GI:

$$i \in \{1, \dots, J\} \quad \begin{bmatrix} p_1 \\ \vdots \\ p_J \end{bmatrix} = \vec{p}$$

$$GI(\vec{p}) = \sum_{i=1}^J p_i (1 - p_i) = 1 - \sum_{i=1}^J p_i^2$$

Lagrange multipliers

$$L = GI(\vec{p}) - \lambda \left( \sum_{i=1}^J p_i - 1 \right)$$

$$\left\{ \begin{array}{l} \nabla_{\vec{p}} L = -2 \vec{p} - \lambda \vec{1} \\ \frac{\partial L}{\partial \lambda} = \sum_{i=1}^J p_i - 1 \end{array} \right.$$

$$\vec{p}^* = -\frac{\lambda}{2} \vec{1}$$

$$\underline{\underline{p_i^* = \frac{1}{J}}}$$

- The maximum GI is when all elements of a probability vector are equal

Cost function is GI:

In every split we minimize GI:

Find  $(j, i)$  such that  $\min_{j, i} GI(\vec{\pi}(R_{ji}))$

If  $\vec{p} = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Rightarrow GI(\vec{p}) = 0$   
It's the minimum.