# Lecture 7
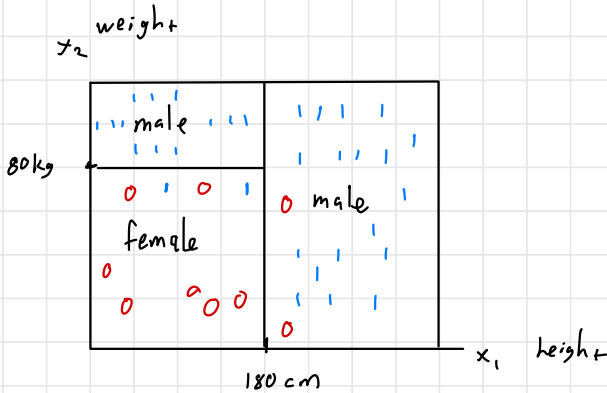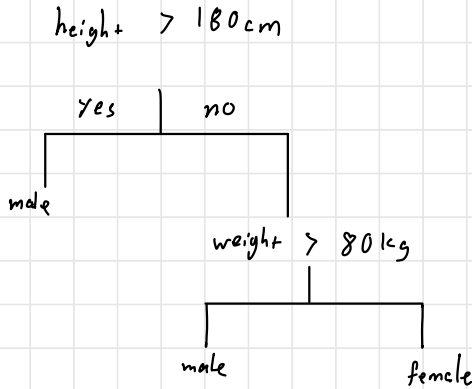
- Part 2: Trees, Forest, Boosting

- Tree: CART classification and Regression Trees
    - classification / Decision

|   | 1 ⋯ j ⋯ p |   |
|---|---|---|
| 1 |   |   |
|   |   | $y_i$ |
| i | $x_i^T$ (height, weight) | gender |
| m |   |   |

height > 180 cm

| yes | no |

male

weight > 80 kg

male                    female



$x_2$ weight

80 kg

male

female

180 cm

$x_1$ height

Recursive Partition

- Each iteration: $(m \times p \times n)$

  - Choose a region $m \in \{1, ..., M\}$:
    choose a variable $j \in \{1, ..., P\}$:
    choose a cut-off $t \in \{x_{ij}, \; i = 1, ... n\}$

    max reduction of Loss function

- Principled way of choosing a loss is to use log-likelihood.
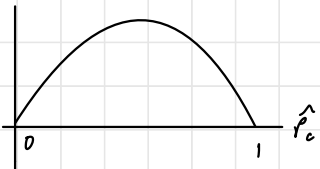- But there are intuitive ideas for the loss function $\longrightarrow$ purity

Within a region $R$

| | 1 | 2 | ... | c | ... | C |
|---|---|---|---|---|---|---|
| counts | $n_1$ | $n_2$ | ... | $n_c$ | ... | $n_C$ |
| props | $P_1$ | $P_2$ | ... | $P_c$ | ... | $P_C$ |
| estimsy | $\hat{P}_1$ | $\hat{P}_2$ | ... | $\hat{P}_c$ | ... | $\hat{P}_C$ |

Note: $\hat{P}_c = \dfrac{n_c}{n}$

Let $c^* = \underset{c}{\text{argmax}} \; \hat{P}_c$

purity $= 1 - \hat{P}_{c^*}$   (smaller $\Rightarrow$ more pure)

Gini-index $= \displaystyle\sum_{c=1}^{C} \hat{P}_c (1 - \hat{P}_c)$

how close $\hat{P} = (\hat{P}_1, ..., \hat{P}_c, ..., \hat{P}_C)$ to one-hat

(assume observations are in $\mathbb{R}$ for simplicity)

- Log - likelihood $(P) = \sum_{i=1}^{n} \log P(Y_i)$

$$= \sum_{c=1}^{C} n_c \log (P_c)$$

$$= n \sum_{c=1}^{C} \hat{P_c} \log (P_c)$$

$$= -n \, H(\hat{P}, P), \qquad H: \text{cross-entropy}$$

$$\max_{P} \log\text{-like}(P) \longrightarrow \hat{P}_{MLE} = \hat{P}$$

$$\&$$

$$\log\text{-like}(\hat{P}) = n \sum_{c=1}^{C} \hat{P_c} \log \hat{P_c}$$

$$= -n \, H(\hat{P}), \qquad H: \text{entropy}$$

- $\log\text{-like}(\hat{P}) - \log\text{-like}(P) = n \sum_{c=1}^{C} \hat{P_i} \log \frac{\hat{P_c}}{P_c}$

$$= n \, D_{KL}(\hat{P} \| P) \geq 0$$

- So we can define $\text{purity} \equiv n \, H(\hat{P})$

- Also the Gini-index needs $n$ as a factor:

$$\text{Gini-index} = \underbrace{n \sum \hat{P_c}(1-\hat{P_c})}$$
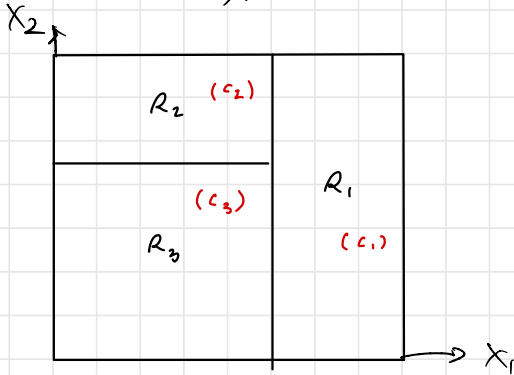
A surrogate for entropy

# Regression



one-dim $\chi$



two-dim $\chi$



Piecewise constant function

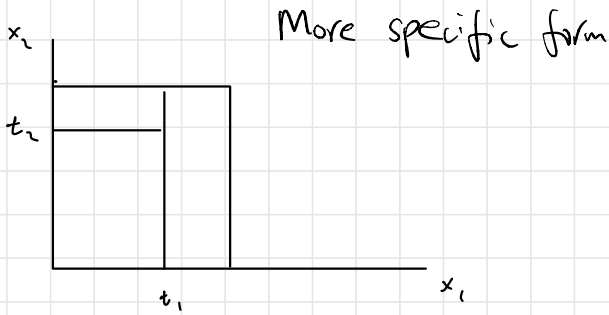$$s = f(x) = \sum_{m=1}^{M} c_m \mathbb{1}(x \in R_m)$$

$$\mathbb{1}(x \in R) = \begin{cases} 1 & \text{if } x \in R \\ 0 & \text{if } x \notin R \end{cases}$$

- Least-squares loss: guide recursive partitioning

$$\text{Loss} = \sum_{i=1}^{n} (y_i - s_i)^2$$

$$= \sum_{m=1}^{M} \sum_{i\, :\, x_i \in R_m} (y_i - c_m)^2$$

Note: $\hat{c}_m = \dfrac{\sum_{i\, :\, x_i \in R_m} y_i}{n_m}$ , $n_m = \#$ of example in $R_m$

$$\sum_{i\, :\, x_i \in R_m} (y_i - \hat{c}_m)^2 = n_m \cdot \text{variance of } y_i \text{ in } R_m, \text{ purity}$$
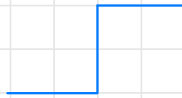
More specific form



$$S = f(x) = c_1 \, \mathbb{1}(x_1 \leq t_1) + c_2 \, \mathbb{1}(x_1 > t_1) \quad \text{(first cut)}$$

$$S = f(x) = c_{11} \, \mathbb{1}(x_1 \leq t) \, \mathbb{1}(x_2 \leq t_2) + c_{12} \, \mathbb{1}(x_1 \leq t_1) \, \mathbb{1}(x_2 > t_2) + c_2 \, \mathbb{1}(x_1 > t_1) \quad \text{(second cut)}$$
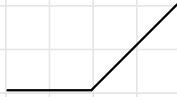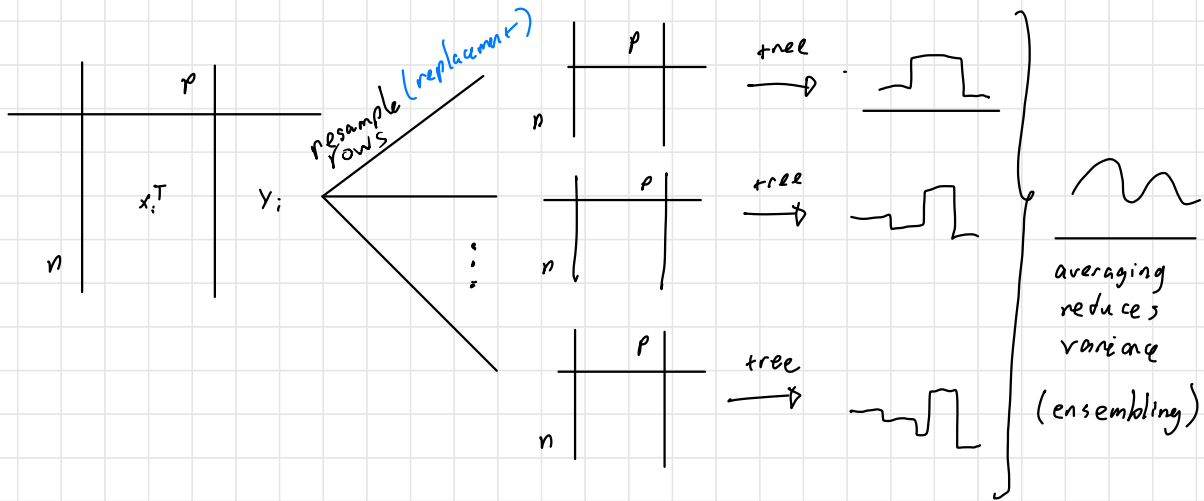
Notice:

$$\mathbb{1}(x_j > t)$$



Sigmoid



ReLU

$\downarrow$

- Multivariate adaptive regression spline (MARS)

- Trees are unstable ( Add an observation changes the tree)
  <span style="color:red">High variance</span>

- Forest:



resample rows (replacement)

$x_i^T$    $y_i$

tree

tree

tree

averaging reduces variance

(ensembling)

- each iteration
  Choose a variable to split
  $\in \{$ subset of $\sqrt{p}$ variables $\}$

- Boosting :

  - Regression Tree :

    $$S = f(x) = \sum_{k=1}^{K} h_K(x) \quad \text{ensemble / comitee}$$

    $\downarrow$
    a tree

  - Sequentially adding trees:

    $$S = \sum_{k=1}^{t-1} h_k(x) + h_t(x)$$

    current comittee       grow a new
    of trees, fixed        tree

    $\hat{S}$                $\Delta S$

- Regression, $\ell_2$ boosting

  - Loss $= \displaystyle\sum_{i=1}^{n} (y_i - s_i)^2$

    $$= \sum_{i=1}^{n} (y_i - (\hat{s}_i + \Delta s))^2$$

    $$= \sum_{i=1}^{n} (\underbrace{y_i - \hat{s}_i}_{\hat{e}_i} - h_b(x))^2$$

    grow a new tree for $\hat{e}_i$

→ Classification , extreme gradient boosting , logistic regression

| | $p$ |
|---|---|
| $x_i^T$ | $y_i \in \{0, 1\}$ |

$n$

$$\text{Log-like} = \sum_{i=1}^{n} \left( y_i s_i - \log(1 + e^{s_i}) \right)$$

$\underbrace{\qquad\qquad}_{\ell(s_i)}$

$$\ell(s_i) \doteq \ell(\hat{s_i}) + \ell'(\hat{s_i}) \Delta s_i + \ell''(\hat{s_i}) \Delta s_i^2$$

$$\ell'(s_i) \doteq y_i - p_i = e_i$$

$$\ell''(s_i) \doteq -p_i(1-p_i) = -w_i$$

So $\ell(s_i) \doteq \text{const} + \hat{e_i} \Delta s_i - \frac{1}{2} \hat{w_i} \Delta s_i^2$

$$= \text{const} - \frac{1}{2}\hat{w_i}\left( \Delta s_i^2 - 2\frac{\hat{e_i}}{\hat{w_i}} \Delta s_i \right)$$

$$= \text{const} - \frac{1}{2}\hat{w_i}\left( \Delta s_i - \underbrace{\frac{\hat{e_i}}{\hat{w_i}}}_{\hat{y_i}} \right)^2$$

$$\text{Loss} = \frac{1}{2} \sum_{i=1}^{n} \hat{w_i} \left( \hat{y_i} - h_t(x_i) \right)^2$$

$\downarrow$

| $c_2 R_2$ | $c_1$ |
|---|---|
| $c_3 R_3$ | $R_1$ |



$$\text{Loss} = \frac{1}{2} \sum_{m=1}^{M} \sum_{i: x_i \in R_m} \hat{w_i}\left( \hat{y_i} - c_m \right)^2 + \frac{1}{2}\gamma \underbrace{\sum_{m=1}^{M} c_m^2}_{\text{Small correction}} + \underbrace{\lambda M}_{\text{smale \# of regions}}$$

$$\text{Loss} = \frac{1}{2} \sum_{m=1}^{M} \sum_{i:x_i \in R_m} \hat{w}_i \left( \hat{y}_i - c_m \right)^2 + \frac{1}{2} \gamma \sum_{m=1}^{M} c_m^2 + \lambda M$$

<span style="color:blue">regularization</span>    <span style="color:blue">regularize</span>

$$\frac{\partial}{\partial c_m} : \quad = - \sum_{i:x_i \in R_m} \hat{w}_i \left( \hat{y}_i - c_m \right) + \gamma c_m \quad = 0$$

$$\hat{c}_m = \frac{\displaystyle\sum_{i:x_i \in R_m} \hat{w}_i \, \hat{y}_i}{\displaystyle\sum_{i:x_n \in R_m} \hat{w}_i + \gamma}$$

<span style="color:blue">shrinkage</span>

Recall IRLS

$$\hat{s}_i \xrightarrow{\Delta s_i} s_i$$
$$\parallel \qquad\qquad \parallel$$
$$x_i^T \beta_t \xrightarrow{x_i^T \Delta \beta} x_i^T \beta$$

$-\ell(s_i)$    small $\hat{w}_i$



$\hat{s}_i$   $\Delta s_i$

$\ell(s_i)$    big $\hat{w}_i$



$s_i$   $\Delta s_i$