

# Decision Trees

Information = # bits needed to encode the probability of an event

$$I = -\log_2 p$$

eg. coin flip from fair coin ( $p = 1/2$ ) is

$$-\log_2 \frac{1}{2} = 1 \text{ bit}$$

eg. certain event ( $p = 1$ )

$$-\log_2 1 = 0 \text{ bits}$$

Want to define  $I$  so that:

- $I(p) \geq 0$

$$I(1) = 0$$

- Two indep events w prob  $p_1$  &  $p_2$

$$I(\text{both events}) = I(p_1 \cdot p_2) = I(p_1) + I(p_2)$$

- $I(p)$  should be continuous

$$\text{Together: } I(p^2) = 2I(p)$$

$$I(p^n) = n I(p) \rightarrow -\log_b(p)$$

$b = 2$  for "bits"

e.g. coin flip gives  $-\log_2 \frac{1}{2} = 1$

e.g. biased coin, heads with  $p = .99$  gives  $-\log_2 .99 = 0.0145$  bits

$p = .01$   $-\log_2 .01 = 6.643$  bits

Entropy  $\rightarrow$  expected information of a set of events

events	$V_1$ --- $V_J$
probs	$P_1$ --- $P_J$
information	$I(P_1)$ --- $I(P_J)$

$$\begin{aligned}\text{Entropy} &= P_1 I(P_1) + P_2 I(P_2) + \dots + P_J I(P_J) \\ &= \sum_j P_j I(P_j) = - \sum_j P_j \log_2(P_j) =: \underset{\substack{\uparrow \\ \text{entropy}}}{H(P_1, \dots, P_J)}\end{aligned}$$

So if	heads	tails
	$p$	$1-p$
	$I(p)$	$I(1-p)$

$$\begin{aligned}H(p, 1-p) &= p I(p) + (1-p) I(1-p) \\ &= -p \log_2 p - (1-p) \log_2 (1-p)\end{aligned}$$

if  $p = .5$   $H(.5, .5) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1 = 1$  bit



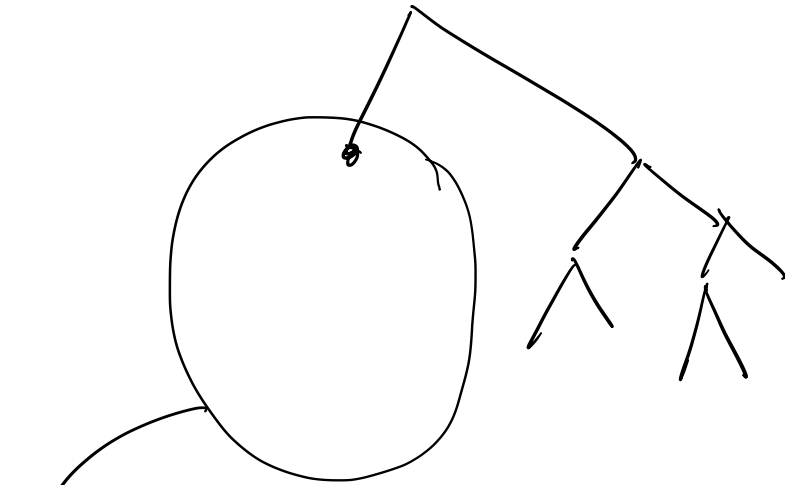
if  $p = .99$   $H(.99, .01) = \dots = .08$  bits



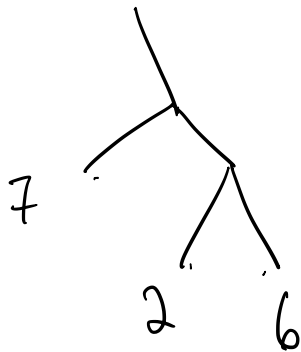
prune

CART "minimal cost complexity pruning"

Each subtree is assigned a cost



$$\text{Cost (subtree)} = \frac{1}{n} \begin{array}{c} \text{\# mistakes} \\ \text{made by} \\ \text{subtree} \end{array} + \underbrace{C}_{0.01} \begin{array}{c} \text{\# leaves in} \\ \text{subtree} \end{array}$$



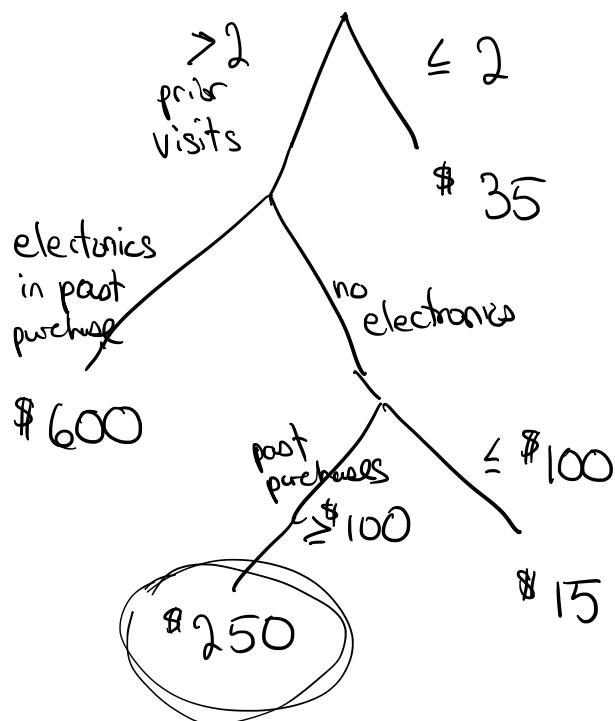
$$\begin{aligned} \text{Cost (subtree)} &= \frac{1}{100} [7 + 2 + 6] + \frac{1}{100} \cdot 3 \\ &= \frac{15}{100} + \frac{3}{100} = \frac{18}{100} \end{aligned}$$

Why did I choose  $C=0.01$ ? Adding one leaf  $\approx$  1% error decrease

Should I add one more leaf that reduces error by .005%?

# Regression Trees

How much will the customer spend at an online store?

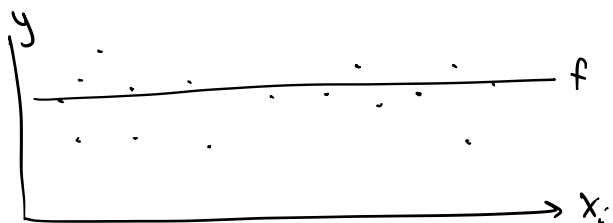


In regression, want to minimize

$$\sum_i \left( f(x_i) - y_i \right)^2$$



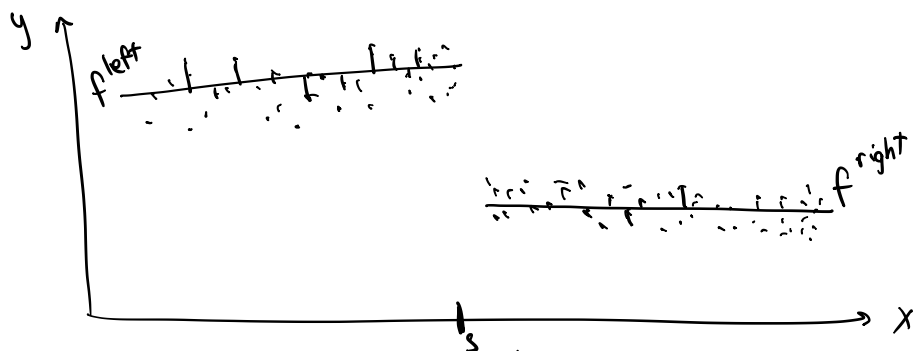
What if  $f$  does not depend on  $x_i$ ?



$f$  is average of  $y_i$ 's,  $f = \frac{1}{n} \sum_i y_i$

In every leaf of regression tree, choose  $f$  in the leaf to be the average of the  $y_i$ 's

# How to split?

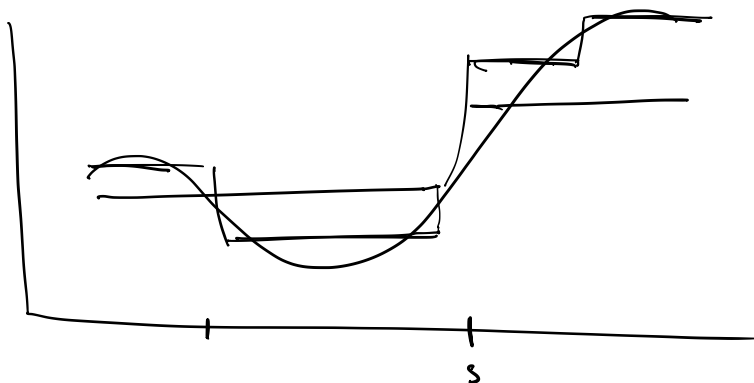


Want to choose split point  $s$ .

minimize  $s$

$$\sum_{\substack{i \text{ where} \\ x_i < s}} (f^{\text{left}} - y_i)^2 + \sum_{\substack{i \text{ where} \\ x_i \geq s}} (f^{\text{right}} - y_i)^2$$

$\underbrace{\hspace{1.5cm}}_{\text{left}}$        $\uparrow$  average of  $y_i$ 's on the left       $\underbrace{\hspace{1.5cm}}_{\text{right}}$        $\uparrow$  average of  $y_i$ 's on the right



Nearest neighbors - can be done well  
or badly

+ interpretable

+ accuracy

- can be misleading if done wrong

- like real estate comps



location  
\*\*\*

sqft  
3000

bdms  
4

bath  
3

--- sold  
price  
\$324K



\*\*\*

2600

2

\$310K



\*+\*

3200

5

3

\$330K

KNN - find  $k$  "nearest neighbors" and take

- majority vote (classif)

- average (regression)

What distance metric to use for comparison?

$$\text{dist}(\text{House}_1, \text{House}_2) = \begin{cases} \infty & \text{if not in same nhbd} \\ (sqft_1 - sqft_2)^2 + C_{bed}(bdrm_1 - bdrm_2)^2 \\ + C_{bath}(bath_1 - bath_2)^2 \end{cases}$$



Can also do "weighted" knn

lot of work lately on learning distances for knn