

Decision theory behind.

"Loss of false prediction" truth.

$y=1$ $y=0$

$\hat{y}=1$ 0 C_2

$C_1, C_2 > 0$ $C_1 \neq C_2$.

C_1 : when $y=1, \hat{y}=0$ cost of dating a fuckboy (渣男)

$\hat{y}=0$ C_1 0

C_2 : when $y=0, \hat{y}=1$, the loss of not dating a decent guy.

Given C_1, C_2 $\hat{y}=1$ or $\hat{y}=0$?

Suppose $p\%$ of men are fuckboy.

Expected loss of $\hat{y}=1$

"Risk" $\rightarrow (1-p)$ $p\%$ are fuckboy

$$P(y=1) \cdot 0 + P(y=0) C_2 \\ = (1-p) C_2.$$

} Guessing he is a fuckboy.

$$\hat{y}^* = \begin{cases} 1 & , p > \frac{C_2}{C_1+C_2} \\ 0 & \text{o.u.} \end{cases}$$

when p high enough. we predict $\hat{y}=1$

Risk of $\hat{y}=0$

$$P(y=1) \cdot C_1 + P(y=0) \cdot 0 \\ = p \cdot C_1.$$

} Guessing he is not a fuckboy. Risk when $C_1 \uparrow$ threshold $\frac{C_2}{C_1+C_2} \downarrow$

Optimal Prediction.

$$\hat{y} = \begin{cases} 1 & , (1-p) C_2 < p C_1 \\ 0 & \text{o.u.} \end{cases}$$

$$(1-p) C_2 < p C_1$$

$$\Leftrightarrow C_2 < (C_1+C_2) p$$

$$\Leftrightarrow p > \frac{C_2}{C_1+C_2}$$

\Rightarrow tend to predict $\hat{y}=1$.

Risk when $C_2 \uparrow$

$\frac{C_2}{C_1+C_2} \uparrow \Rightarrow$ optimal rule tends to predict $\hat{y}=0$.

0-1 loss

$y=1$ $y=0$ in this situation

$\hat{y}=1$ 0

1

$$\hat{y}^* = \begin{cases} 1, & p > \frac{1}{2} \\ 0, & \text{o.w.} \end{cases} \quad \text{when } p > \frac{1}{2}$$

we predict $\hat{y}=1$

$\hat{y}=0$ 1 0

In practice, we don't know p .

We have to estimate it.

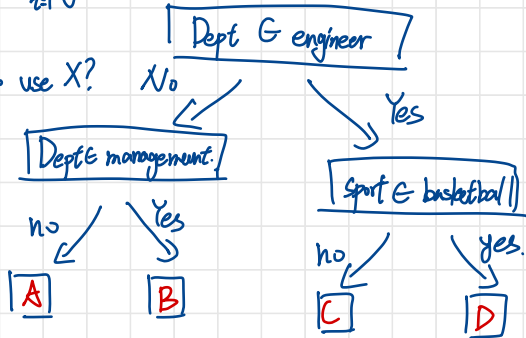
Data:

y	dept	age	sport
1	econ	21	⋮
0	EE	20	⋮
1	⋮	19	⋮
1	⋮	21	⋮
0	⋮	⋮	⋮

Decision tree:
partition out data by X

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n y_i$$

How to use X ?



In D

$$\# \{y=1\} = 80$$

$$\# \{y=0\} = 20$$

$$\hat{p}(y|\text{D}) = 0.8$$

$$P(y|x \in \square) > \frac{C_1}{C_1 + C_2}$$

decision rule

In C

$$\# \{y=1\} = 10$$

$$\# \{y=0\} = 90$$

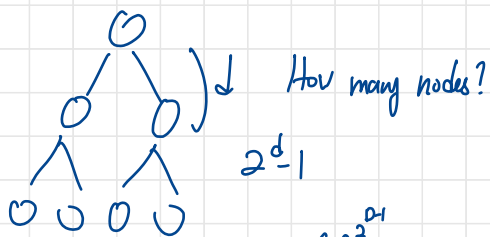
$$\hat{p}(y|\text{C}) = 0.1$$

There are many possible trees. Which one is better?

Question: Suppose depth of tree = d

Suppose there are p variables.

How many possible trees are there?



Trees: $(p)^{2^d}$ combinations.

eg. $p=10$ $D=3$

$$(p)^{2^d - 1} = 10^7 - 1$$

Greedy algorithm

if you can only make one step

to it, which one's the best?

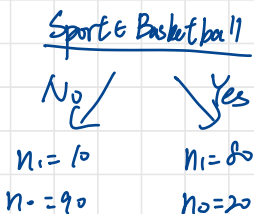
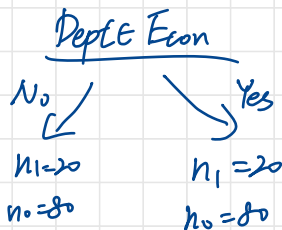
Again.

Combinational problem.

hard to find global optimum.

Split 1

Split 2



Split 2 is a better split because it reduces entropy

熵(熵度)

$X_1 \sim \text{Ber}(0.5)$ ← 比較亂

$X_2 \sim \text{Ber}(0.001)$

Def (Entropy for $\text{Ber}(p)$)

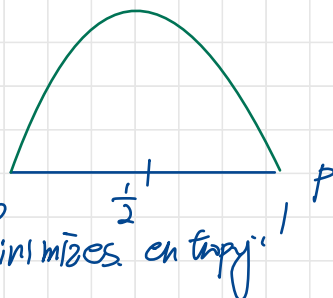
Since $p \in (0,1)$

$\log p \leq 0$

$X \sim \text{Ber}(p)$

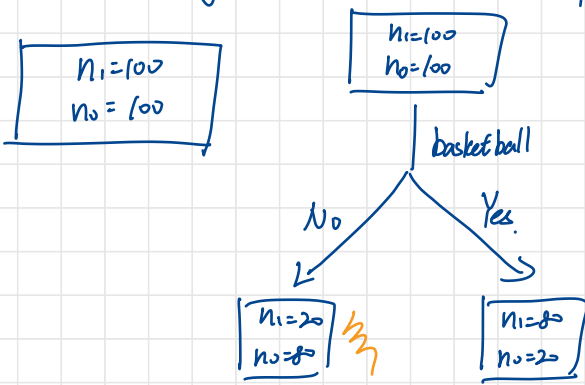
$$H(X) = -[p \log_2 p + (1-p) \log_2 (1-p)] \quad \text{so } H(X) \geq 0$$

"recursively find the split that minimizes entropy"



Information gain = reduction in entropy

Original, $H = 1 \left(P = \frac{1}{2} \right)$



$$\text{Entropy (No)} = 0.2 \log_2 0.2 + 0.8 \log_2 0.8 = 0.57$$

$$\text{Entropy (Yes)} = 0.8 \log_2 0.8 + 0.2 \log_2 0.2 = 0.57$$

avg entropy = 0.57.

$$\text{Information gain} = 1 - 0.57 = 0.43.$$

Greedy algorithm. for decision tree.

Repeat: find the split with highest information gain.

Rank depth = d
variable = p

$$(2^d - 1)p \ll p^{(2^d - 1)}$$

Rank: Greedy algorithms does not guarantee to find global optimal solutions

Rank When to stop?

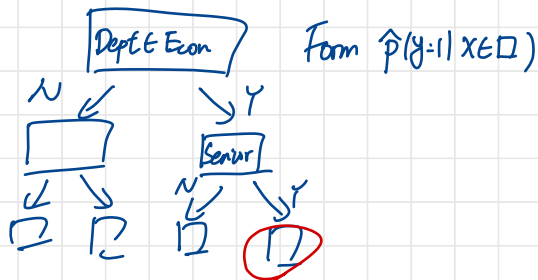
"bias-variance trade off" cross-validation.

Plug in decision rule.

How to use a tree?

$$\text{if } \hat{p}(y=1 | X \in \square) > C$$

$$\hat{y} = 1 \text{ where } C = \frac{C_2}{C_1 + C_2}$$



Trade-off in C

$C \uparrow \Rightarrow$ too lenient

$\hat{y}=1$

$y=1$	$y=0$
TP	FP
FN	TN

$$\text{precision} = \frac{TP}{TP+FP} \Rightarrow \# \hat{y}=1$$

$C \downarrow \Rightarrow$ too strict

$\hat{y}=0$

$$\text{recall} = \frac{TP}{TP+FN} \Rightarrow \# \hat{y}=1$$

if $C \uparrow$, precision \uparrow recall \downarrow

if $C \downarrow$, precision \downarrow , recall \uparrow

