

Decision Theory behind Classification

觀察不到 y

prediction : \hat{y}

"Loss of false prediction"

truth

$y=1$ $y=0$

$\hat{y}=1$ 0 C_2

$\hat{y}=0$ C_1 0

C_1 : when $y=1$, $\hat{y}=0$

$C_1 \neq C_2$ (不對襯)

C_2 : when $y=0$, $\hat{y}=1$

給定 C_1, C_2

$\hat{y} = 0 \text{ or } \hat{y} = 1 ?$ 犯錯不可避免，但能計算・壓底犯錯機率

e.g. 假設男人中有 $P\%$ 是渣男， P 已知

Expected Loss of $\hat{y} = 1$
Risk

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

Expected Loss of $\hat{y} = 0$
Risk

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

\Rightarrow Optimal Prediction

$$\begin{cases} \hat{y}^* = 1, & \text{when } (1-P)C_2 < P C_1 \Rightarrow P > \frac{C_2}{C_1 + C_2} \\ \hat{y}^* = 0, & \text{otherwise} \end{cases}$$

主觀衡量

\therefore 當 P 足夠高 \Rightarrow 預測 $\hat{y} = 1$

Remark :

- when $C_1 \uparrow$

\Rightarrow threshold $\frac{C_2}{C_1+C_2} \downarrow \Rightarrow P$ 更容易 $>$ threshold \Rightarrow 越容易作出 $\hat{y} = 1$

- when $C_2 \uparrow$

\Rightarrow threshold $\frac{C_2}{C_1+C_2} \uparrow \Rightarrow$ 越容易作出 $\hat{y} = 0$

$$0 - 1 \quad |_{\text{loss}} \quad (\text{loss} : C_1 = C_2 = 1)$$

$y=1$

$y=0$

$\hat{y}=1$

0

1

$\hat{y}=0$

1

0

$$\therefore \hat{y}^* = \begin{cases} 1 & , P > \frac{1}{1+1} = \frac{1}{2} \\ 0 & , \text{otherwise} \end{cases}$$

然而，實際上我們並不清楚 P

\Rightarrow 只能估計 \hat{P}

Data

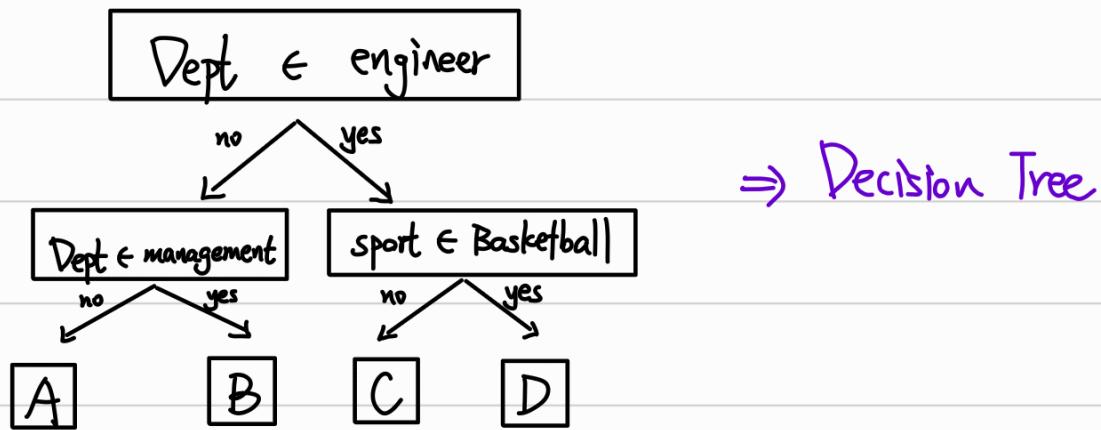
y , X (dept., age, sport)

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

But how to use X ?

Decision Tree

先切開資料 X



in \boxed{D}

$$\begin{aligned} \#\{y=1\} &= 80 \\ \#\{y=0\} &= 20 \end{aligned} \Rightarrow \hat{P}(y|\boxed{D}) = 0.8$$

in \boxed{C}

$$\begin{aligned} \#\{y=1\} &= 10 \\ \#\{y=0\} &= 90 \end{aligned} \Rightarrow \hat{P}(y|\boxed{C}) = 0.1$$

Decision rule

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

better

比 $P(y) > \frac{C_2}{C_1 + C_2}$ 還好 (因為用了 X 的資料)

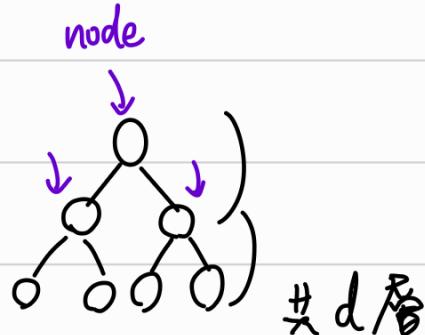
Challenge

有很多種可能的樹 → 哪種切法比較好？

Question

假設： depth of tree = d

有 P 個 variables



⇒ 共有多少種 tree ?

$$2^0 + 2^1 + 2^2 + 2^3 + \dots)$$

How many nodes ?

$$\frac{1(1 - 2^{d+1})}{-1}$$

$$2^D - 1$$

trees : $P^{(2^D - 1)}$ 種 組合

若 $P=10, D=3 \Rightarrow 10^7$ (-千萬種)

⇒ 組合問題複雜

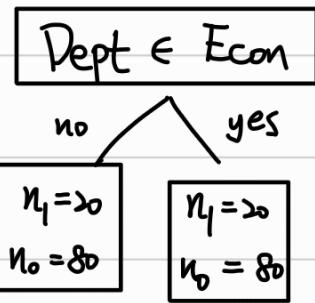
⇒ hard to find global optimum

Greedy Algorithm

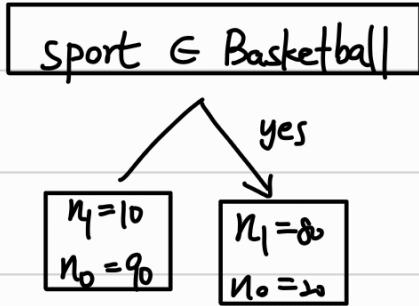
split 1

if you can only make one split

which one is the best?



split 2



↓
split 2 & split 1 if

⇒ 告訴我們比較多資訊

(reduce entropy)
熵(亂度)

entropy :

- $X_1 \sim \text{Bernoulli}(0.5) \Rightarrow$ 抽出來可能是 0, 1, 0, 0, 1, ...
- $X_2 \sim \text{Bernoulli}(0.001) \Rightarrow$ 抽出來幾乎都 0, 0, 0, ..., (entropy 低)

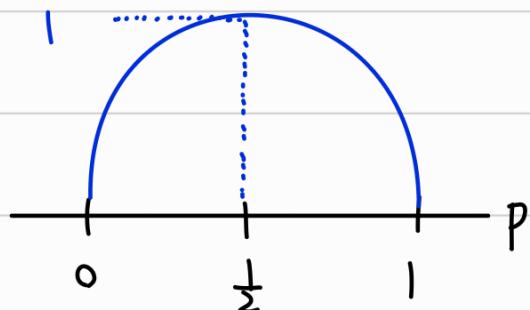
Def (Entropy for $\text{Ber}(p)$)

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

$$\Rightarrow \text{Entropy} : H(X) = -[p \log_2 p + (1-p) \log_2 (1-p)]$$

since $p \in (0, 1)$, $\log p \leq 0$

$$\therefore H(X) \geq 0$$



越靠近中間越亂

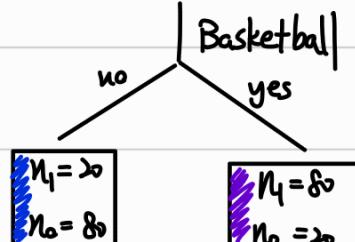
Greedy Algorithm for decision tree

遞迴地 尋找能最小化 entropy 的 split
repeatedly (最大化 information gain)

information gain

$$\begin{aligned} n_1 &= 100 \\ n_0 &= 100 \end{aligned}$$

$$\begin{aligned} n_1 &= 100 \\ n_0 &= 100 \end{aligned}$$

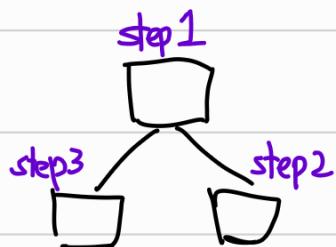


$$\text{originally: } H = 1 \quad (p = \frac{1}{2})$$

$$\begin{aligned} &= 0.2 \log_2 0.2 + 0.8 \log_2 0.8 \\ &\approx 0.57 \end{aligned}$$

$$\begin{aligned} &= 0.8 \log_2 0.8 + 0.2 \log_2 0.2 \\ &\approx 0.57 \end{aligned}$$

average entropy : 0.57
information gain : $1 - 0.57 = 0.43$



Remark :

$$\text{depth} = d$$

$$\# \text{variable} = p$$

$$(2^d - 1)P \quad \text{greedy} \quad \lll P^{(2^d - 1)}$$

Remark :

Greedy algorithm 不保證能找到 global optimal solution

Remark :

when to stop? (cross-validation)

bias-variance tradeoff

如果 depth 太深 \longleftrightarrow depth 太浅

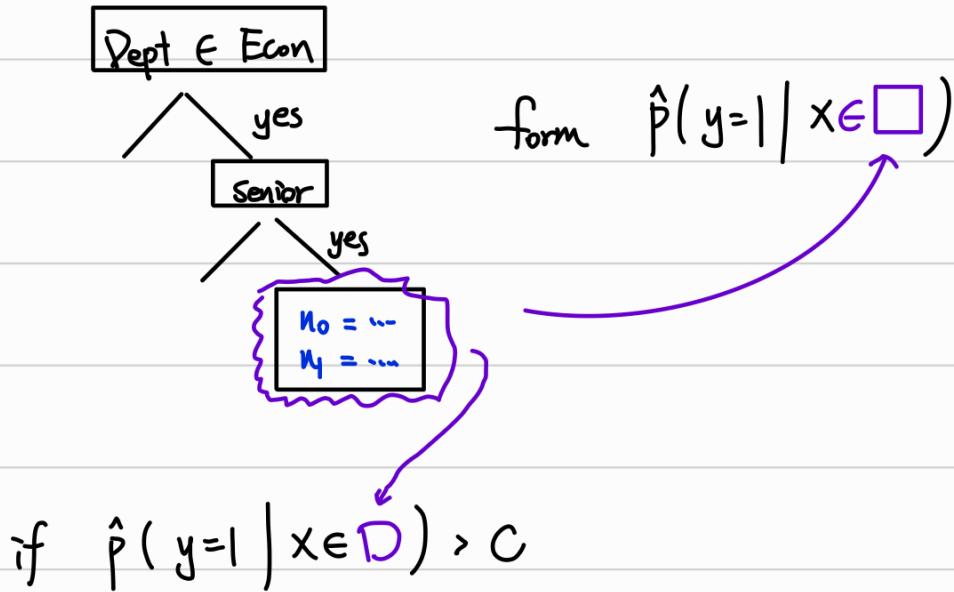


underutilize the information of X

$$n_4=5 \Rightarrow \hat{p} = \frac{1}{7} \text{ (样本数太少)}$$

Noisy

How to use a tree?



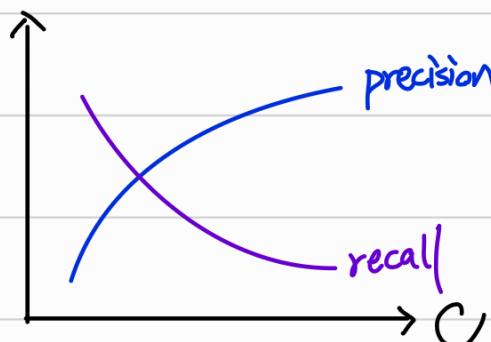
then $\hat{y} = 1$, where $C = \frac{C_2}{C_1 + C_2}$

Trade-off in C

C 太大 \Rightarrow 太寬鬆
 C 太小 \Rightarrow 太嚴格

	$y=1$	$y=0$
$\hat{y}=1$	TP	FP
$\hat{y}=0$	FN	TN

precision = $\frac{TP}{TP + FP}$ $y=1$ 中 $y=1$ 的次數
 $\qquad\qquad\qquad$ $\{ \hat{y}=1 \}$ 的次數



recall : $\frac{TP}{TP + FN}$ $y=1$ 中，被你認為 $\hat{y}=1$ 次數
 $\qquad\qquad\qquad$ $y=1$ 的人數

