

分錯的誤差 (算混亂程度)

① Gini  $0 \sim 0.5$  (越接近 0.5 越混亂)

② Entropy:  $0 \sim 1$  ( " 1 " )

③ Misclassification:  $0 \sim 0.5$  ( " 0.5 " )

$\Delta$	$\Delta$	$\Delta$
$\Delta$	$\Delta$	$\Delta$

$$\text{error} = 1 - \max \left[ \frac{0}{6}, \frac{1}{6} \right] = 0$$

others 比例 /  $\Delta$  比例 混亂度 0.

ex1 解釋各變數: 客戶編號不管, 忠誠度 y (應變值), others 自變數

ex2: 是否還款  $\Rightarrow$  是 (p) / 否 (N)  $\Rightarrow I(5, 4)$

\* 計算 Entropy, 2 個變項  $\Rightarrow$  1 個加 3

\* 以  $I(b)$  表分割前混亂度 /  $IA(b)$  表分割後混亂度

\* If  $I(0, 0) = 0$

\* 前一階段的分割條件, 下一步驟可不用, 但下一步未必

\* 如果算出的增量集都一樣, 可以 choose any 來當分枝樹

\* 前 9 項用來建模 (建決策樹), 後 3 項用來測試

Entropy:

$$I(b) = I(5, 4) = - \left[ \frac{5}{9} \cdot \log_2 \frac{5}{9} + \frac{4}{9} \cdot \log_2 \frac{4}{9} \right]$$

$$= - \left[ \frac{5}{9} \times (-0.847) + \frac{4}{9} \times (-1.169) \right] = 0.989$$

Gini:

$$I(b) = I(5, 4) = 1 - \left( \frac{5}{9} \right)^2 - \left( \frac{4}{9} \right)^2 = 0.494$$

①

引入 "是否負債" 做分割條件

Entropy:

$$IA(b): IA(1, 3) = \frac{4}{9} I(1, 3) + \frac{5}{9} I(4, 1)$$

$$= \frac{4}{9} \cdot - \left[ \frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right] + \frac{5}{9} \cdot - \left[ \frac{4}{5} \log_2 \frac{4}{5} + \frac{1}{5} \log_2 \frac{1}{5} \right]$$

$$= 0.36 + 0.4$$

$$= 0.76$$

是否負債	(是否還款)		
	p	n	total
是	1	3	4
否	4	1	5
total	5	4	9

增益量  $Gain(A) = I(b) - IA(b) = 0.989 - 0.76 = 0.229$

Gini:

$$I_A(D) = \frac{4}{9} \cdot I(1,3) + \frac{5}{9} \cdot I(4,1)$$

$$= \frac{4}{9} \cdot \left[ 1 - \left( \frac{1}{4} \right)^2 - \left( \frac{3}{4} \right)^2 \right] + \frac{5}{9} \cdot \left[ 1 - \left( \frac{4}{5} \right)^2 - \left( \frac{1}{5} \right)^2 \right]$$

$$= \frac{1}{6} + \frac{8}{45} = 0.345 \#$$

$$Gain(A) = I(D) - I_A(D) = 0.494 - 0.345 = 0.149$$

② 引入“性别”做划分条件

Entropy:

$$I_A(D) = \frac{5}{9} \times I(3,2) + \frac{4}{9} \times I(2,2) = 1$$

$$= \frac{5}{9} \times \left[ -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right] +$$

$$\frac{4}{9} \times \left[ -\frac{2}{4} \log_2 \frac{1}{2} - \frac{2}{4} \log_2 \frac{1}{2} \right] = 1$$

$$= \frac{5}{9} \times (0.911) + \frac{4}{9} \times (1) = 0.983 \#$$

$$Gain(A) = 0.989 - 0.983 = 0.006$$

性别	p	n	total
男	3	2	5
女	2	2	4
total	5	4	9

Gini:

$$I_A(D) = \frac{5}{9} \times I(3,2) + \frac{4}{9} \times I(2,2) = \frac{1}{2}$$

$$= \frac{5}{9} \times \left[ 1 - \left( \frac{3}{5} \right)^2 - \left( \frac{2}{5} \right)^2 \right] + \frac{4}{9} \times \left[ 1 - \left( \frac{1}{2} \right)^2 - \left( \frac{1}{2} \right)^2 \right]$$

$$= \frac{4}{15} + \frac{2}{9} = 0.487 \#$$

$$Gain(A) = 0.494 - 0.487 = 0.007$$

③

引入“婚姻状况”做划分条件

Entropy:

$$I_A(D) = \frac{6}{9} \times I(3,3) + \frac{3}{9} \times I(2,1)$$

$$= \frac{6}{9} \times (1) + \frac{3}{9} \times \left[ -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right]$$

$$= 0.917 \#$$

$$Gain(A) = 0.989 - 0.917 = 0.072$$

婚姻状况	p	n	total
单身	3	3	6
结婚	2	1	3
total	5	4	9



Gini:

$$\frac{1}{9} \times \overset{0.5}{I(3,3)} + \frac{2}{9} \times I(2,1) = \frac{1}{9} \times 0.5 + \frac{2}{9} \times \left[ 1 - \left( \frac{2}{3} \right)^2 - \left( \frac{1}{3} \right)^2 \right] = 0.481$$

$$\text{Gain}(A) = 0.494 - 0.481 = 0.013$$

④

引入“收入”做分割条件

Entropy:

$$\begin{aligned} I(A|D) &= \frac{4}{9} \overset{1}{I(2,2)} + \frac{5}{9} I(3,2) \\ &= \frac{4}{9} \times \left[ -\left[ \frac{1}{2} \times \log_2 \frac{1}{2} + \frac{1}{2} \times \log_2 \frac{1}{2} \right] \right] + \\ &\quad \frac{5}{9} \times \left[ -\left[ \frac{3}{5} \times \log_2 \frac{3}{5} + \frac{2}{5} \times \log_2 \frac{2}{5} \right] \right] \\ &= 0.989 + 0.539 \\ &= 0.983 \end{aligned}$$

$$\text{Gain}(A) = 0.989 - 0.983 = 0.006$$

Gini:

$$\begin{aligned} I(A|D) &= \frac{4}{9} I(2,2) + \frac{5}{9} I(3,2) \\ &= 0.222 + \frac{5}{9} \times \left[ 1 - \left( \frac{3}{5} \right)^2 - \left( \frac{2}{5} \right)^2 \right] = 0.267 \\ &= 0.489 \end{aligned}$$

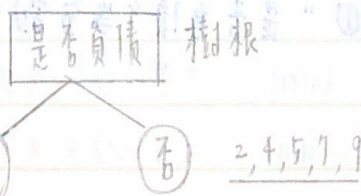
$$\text{Gain}(A) = 0.494 - 0.489 = 0.005$$

結論: 選擇增益量最大者

	Entropy	Gini
是否負債	0.229	0.149
性別	0.006	0.007
婚姻狀況	0.017	0.013
收入	0.006	0.005

⇒ 以“是否負債”為分割條件的增益量最大, 故選擇“是否負債”為決策樹之分支條件

\* 下階段可不考慮重疊的分割條件



左分枝開始重新計算評估

只观察  $\text{否}, \text{否}, \text{否} \Rightarrow I(1,3)$

① “性别”当分割条件

起初 Entropy:  $I(1,3) = 0.81$

Gini:  $I(1,3) = 0.375$

性别	p	n	total
男	1	2	3
女	0	1	1
total	1	3	4

now Entropy:  $\frac{3}{4} I(1,2) + \frac{1}{4} I(0,1) = 0.688$

Gini:  $\frac{3}{4} I(1,2) + \frac{1}{4} I(0,1) = 0.333$

$\Rightarrow$  值越大, 增益量越小

② “婚姻”当分割

Entropy:  $\frac{3}{4} \times I(1,2) + \frac{1}{4} \times I(0,1) = 0.688$

Gini:  $\frac{3}{4} I(1,2) + \frac{1}{4} I(0,1) = 0.333$

婚姻	p	n	total
單	1	2	3
結	0	1	1

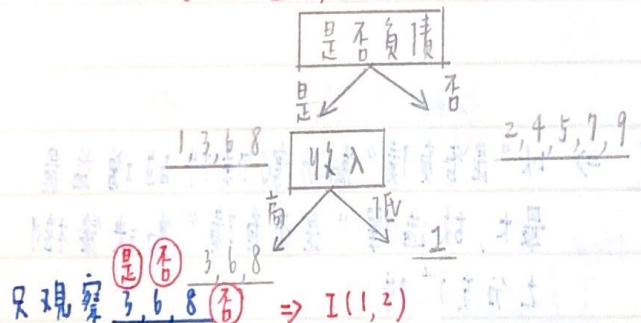
③ “收入”当分割

Entropy:  $\frac{1}{4} I(0,1) + \frac{3}{4} I(1,2) = 0.688$

Gini:  $0.333$

收入	p	n	total
低	0	1	1
高	1	2	3

\* 上述, 不論选哪个当分枝, 增益量皆相同  
故“可自由选择”



① “是否負債”当分割

Entropy:  $\frac{3}{3} I(1,2) + 0 I(0,0) = 0.917$

Gini:  $\frac{3}{3} I(1,2) + 0 I(0,0) = 0.444$

負債	p	n	total
是	1	2	3
否	0	0	0



## ② "性别"

$$\text{Entropy} = \frac{2}{3} I(1,1) + \frac{1}{3} I(0,1) = 0.667$$

$$\text{Gini} = \frac{2}{3} I(1,1) + \frac{1}{3} I(0,1) = 0.334$$

性别	p	n	total
男	1	1	2
女	0	1	1

## ③ "婚姻"

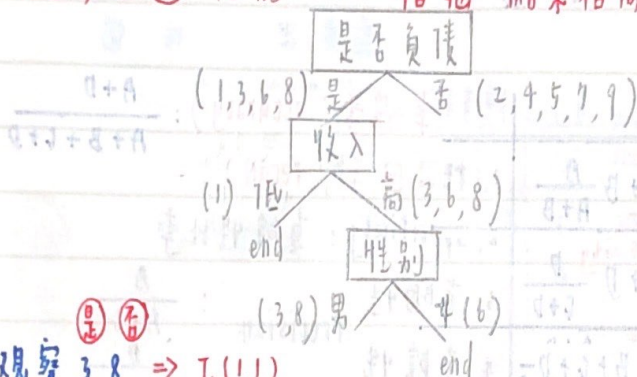
結果同上

1) 当結果值越大, 增益量越小

所以不必算出原本之 Entropy & Gini

\* 上述, 故选 "性别" or "婚姻" 結果相同

婚姻	p	n	total
Single	1	1	2
married	0	1	1



只观察 3, 8  $\Rightarrow I(1,1)$

## ① "是否負債"

$$\text{Entropy} = \frac{2}{2} I(1,1) + \frac{0}{2} I(0,0) = 0$$

$$\text{Gini} = \frac{2}{2} \times [1 - (\frac{1}{2})^2 - (\frac{1}{2})^2] = 0.5$$

負債	p	n	total
是	1	1	2
否	0	0	0

## ② "婚姻"

$$\text{Entropy} = \frac{1}{2} I(1,0) + \frac{1}{2} I(0,1) = 1$$

$$\text{Gini} = \frac{1}{2} \times [1 - 1 - 0] = 0$$

婚姻	p	n	total
single	1	0	1
married	0	1	1

## ③ "收入"

結果同 "是否負債"

\* 以 "婚姻" 做分割

