

## 1) Entropy and decision trees

### 1.1) Entropy

$$1) \quad p = 0.1 \Rightarrow H(x) = 0.1 \times \log_2\left(\frac{1}{0.1}\right) + (1-0.1) \log_2\left(\frac{1}{1-0.1}\right) = 0.46$$

$$p = 0.2 \Rightarrow H(x) = 0.72$$

$$p = 0.3 \Rightarrow H(x) = 0.88$$

$$p = 0.4 \Rightarrow H(x) = 0.97$$

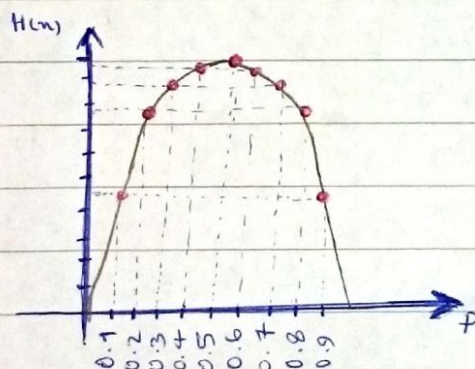
$$p = 0.5 \Rightarrow H(x) = 1$$

$$p = 0.6 \Rightarrow H(x) = 0.97$$

$$p = 0.7 \Rightarrow H(x) = 0.88$$

$$p = 0.8 \Rightarrow H(x) = 0.72$$

$$p = 0.9 \Rightarrow H(x) = 0.46$$



- When the probabilities for taking both values are equal, the disorder is large and correspondingly the entropy value is large; both are in the highest state. On the other hand, when  $x$  takes one of the values in most of the times, the disorder is low and the value of  $p$  will be close to 0 or 1, and the entropy value would be low.



$$2) - P_n = \frac{1}{N} \quad H(x) = \sum_{n=0}^{N-1} P_n \log_2 \left( \frac{1}{P_n} \right)$$

$$H(x) = \sum_{n=0}^{N-1} \frac{1}{N} \log_2(N) = \log_2(N) \quad (2)$$

$$- \mathcal{L}: \log_2(\cdot)$$

$$\lambda_n: P_n$$

$$u_n: \frac{1}{P_n}$$

$$H(x) = \sum_{n=0}^{N-1} P_n \log_2 \left( \frac{1}{P_n} \right)$$

$$H(x) = \sum_{n=0}^{N-1} \underbrace{P_n}_{\lambda_n} \log_2 \left( \underbrace{\frac{1}{P_n}}_{u_n} \right)$$

$$H(x) \stackrel{(1)}{=} \sum_{n=0}^{N-1} \lambda_n \mathcal{L}(u_n)$$

$$\mathcal{L} \left( \sum_{n=0}^{N-1} \lambda_n u_n \right) \stackrel{(1)}{=} \log_2 \left( \sum_{n=0}^{N-1} P_n \frac{1}{P_n} \right) = \log_2 \left( \sum_{n=0}^{N-1} 1 \right) = \log_2(N)$$

$$\text{Jensen's inequality:} \quad \sum_{n=0}^{N-1} \lambda_n \mathcal{L}(u_n) \leq \mathcal{L} \left( \sum_{n=0}^{N-1} \lambda_n u_n \right)$$

$$\stackrel{(1)}{\Downarrow} \quad H(x) \leq \log_2(N)$$

(1) and (2)  $\Rightarrow$  entropy of  $x$  is maximal if and only if  $x$  follows a uniform distribution.



## 1.2) Relation between entropy and information

### 1.2.1) Arresting a criminal

$$1) p = \frac{1}{5 \times 80} \quad q = \frac{5 \times 80 - 1}{5 \times 80} \quad H(n) = p \log_2 \left( \frac{1}{p} \right) + q \log_2 \left( \frac{1}{q} \right) \\ = 0.0252$$

$$2) p = \frac{1}{80} \quad q = \frac{80 - 1}{80} \quad H(n) = 0.0969$$

$$3) p = \frac{1}{30} \quad q = \frac{30 - 1}{30} \quad H(n) = 0.2108$$

4) The information that officers receive increase the uncertainty;

as the officer A is more certain about the innocence of every randomly chosen sample. The officer B is less certain ~~than~~ than officer A, and officer C is less certain<sup>2</sup> than officer B. Thus, in case of officer A we have more information, and for officer B we have less information regarding A, but more regarding C. Therefore, in the case of officer A we have less entropy, and in case of officer B we have more entropy than case A but less than case C.



### 1.2.2) Guessing cards

$$1) P(X=n) = \frac{1}{52} \quad \text{for all } n$$

$$H(n) = \sum_{n=0}^{52-1} \frac{1}{52} \log_2(52) = \log_2(52) = 5.7$$

$$2) P(Y) = \frac{4}{52}$$

$$P(Y) = \begin{cases} \frac{4}{52} & Y = \text{true} \\ \frac{48}{52} & Y = \text{false} \end{cases}$$

~~$P(X|Y) = \frac{1}{4}$  for all not ace us~~

$$3) P(X|Y=\text{true}) = \begin{cases} \frac{1}{4} & \text{for ace us} \\ 0 & \text{for all not ace us} \end{cases} \quad P(X|Y=\text{false}) = \begin{cases} \frac{1}{48} & \text{for all not ace us} \\ 0 & \text{for ace us} \end{cases}$$

$$H(X|Y=\text{true}) = \sum_{i=0}^{52-1} P(X=n_i | Y=\text{true}) \log_2 \left( \frac{1}{P(X=n_i | Y=\text{true})} \right)$$

$$= 4 \left( \frac{1}{4} \log_2(4) \right) = 2$$

$$H(X|Y=\text{false}) = \sum_{i=0}^{52-1} P(X=n_i | Y=\text{false}) \log_2 \left( \frac{1}{P(X=n_i | Y=\text{false})} \right)$$

$$= 48 \left( \frac{1}{48} \log_2(48) \right) = 5.58$$



4) if  $E_y$  happens then we have much information about the identity of the card ( $x$ ), so  $H(x|y=\text{true})$  is small; however, the probability of event  $y$  happening is also small.

If the event  $y$  ( $E_y$ ) does not happen, then guessing the card would be much harder ~~and~~  $H(x|y=\text{false})$  is large ~~and~~, and the probability of  $E_y$  not happening is also large.

Therefore, knowing about  $E_y$  does not help considerably, which means it does not give us much information, so correspondingly we observe that  $H(x|y)$  is large.

$$H(x|y) = \overbrace{P(y=\text{true})}^{\text{large}} \cdot \overbrace{H(x|y=\text{true})}^{\text{small}} + \overbrace{P(y=\text{false})}^{\text{small}} \cdot \overbrace{H(x|y=\text{false})}^{\text{large}}$$

$$5) \quad H(x|y) = \frac{4}{52} \times 2 + \frac{48}{52} \times 5.58 = 5.3$$

$$H(x) = 5.7$$

$H(x|y)$  is not much less than  $H(x)$ . It seems that knowing about  $E_y$  does not considerably affect the entropy of  $x$ .



$$6) IG(X|Y) = H(X) - H(X|Y)$$

$$= 5.7 - 5.3 = 0.4$$

$IG(X|Y)$  should be large. When  $IG(X|Y)$  is large it means that knowing about  $Y$  gives us much information that by knowing  $X$  it the entropy of  $X$  would be much smaller.

7) Know about "is the top card spades" is more informative.

$$E_Z: \text{the top card is spades} \quad P(E_Z) = \frac{1}{4}$$

$$H(X|Z) = \frac{1}{4} (\log_2(13)) + \frac{3}{4} (\log_2(39)) = 4.88$$

$$IG(X|Z) = H(X) - H(X|Z) = 5.7 - 4.88 = 0.82$$

$$IG(X|Z) > IG(X|Y)$$

Knowing about the color is more informative.



### 1.3) Decision trees

$$1) \quad p(A) = \frac{5}{10} \quad p(B) = \frac{5}{10} \quad H(\text{class}) = \frac{5}{10} \log_2\left(\frac{10}{5}\right) + \frac{5}{10} \log_2\left(\frac{10}{5}\right) = 1$$

$$H(\text{class}|x) = p(x=1) H(\text{class}|x=1) + p(x=0) H(\text{class}|x=0)$$

$$= \frac{4}{10} \left[ \frac{1}{4} \log_2(4) + \frac{3}{4} \log_2\left(\frac{4}{3}\right) \right] + \frac{6}{10} \left[ \frac{4}{6} \log_2\left(\frac{6}{4}\right) + \frac{2}{6} \log_2\left(\frac{6}{2}\right) \right] = 0.87$$

$$H(\text{class}|y) = \frac{5}{10} \left[ \frac{3}{5} \log_2\left(\frac{5}{3}\right) + \frac{2}{5} \log_2\left(\frac{5}{2}\right) \right] + \frac{5}{10} \left[ \frac{2}{5} \log_2\left(\frac{5}{2}\right) + \frac{3}{5} \log_2\left(\frac{5}{3}\right) \right] = 0.97$$

$$H(\text{class}|z) = \frac{6}{10} \left[ \frac{3}{6} \log_2\left(\frac{6}{3}\right) + \frac{3}{6} \log_2\left(\frac{6}{3}\right) \right] + \frac{4}{10} \left[ \frac{2}{4} \log_2\left(\frac{4}{2}\right) + \frac{2}{4} \log_2\left(\frac{4}{2}\right) \right] = 1$$

$$IG(\text{class}|x) = H(\text{class}) - H(\text{class}|x) = 1 - 0.87 = 0.125 \quad \checkmark$$

$$IG(\text{class}|y) = 1 - 0.97 = 0.029$$

$$IG(\text{class}|z) = 1 - 1 = 0$$

- The first node uses x for classification.

x	y	z	Class
1	0	0	A
1	0	1	B
1	1	1	B
1	1	1	B

x	y	z	Class
0	0	0	B
0	1	1	A
0	0	1	A
0	1	1	A
0	1	0	A
0	0	0	B



$$H(\text{class} | Y, X=1) = \frac{2}{4} \left[ 0 + \frac{2}{2} \log_2(1) \right] + \frac{2}{4} \left[ \frac{1}{2} \log_2(2) + \frac{1}{2} \log_2(2) \right]$$

$$= \frac{2}{4} \times 0 + \frac{2}{4} \times 1 = \frac{2}{4} = 0.5$$

$$H(\text{class} | Z, X=1) = \frac{3}{4} [1 \log_2(1)] + \frac{1}{4} [1 \log_2(1)] = 0$$

~~$H(\text{class} | Y, X=1)$~~

$$H(\text{class} | X=1) = \frac{1}{4} \log_2(4) + \frac{3}{4} \log_2\left(\frac{4}{3}\right) = 0.81$$

$$IG(\text{class} | Y, X=1) = 0.81 - 0.5 = 0.31$$

$$IG(\text{class} | Z, X=1) = 0.81 - 0 = 0.81 \quad \checkmark$$

- The next node when  $X=1$  should use Z for classification

$$H(\text{class} | Y, X=0) = \frac{3}{6} \left[ \frac{1}{3} \log_2(3) + \frac{2}{3} \log_2\left(\frac{3}{2}\right) \right] + \frac{3}{6} \left[ \frac{3}{3} \log_2\left(\frac{3}{3}\right) + 0 \right] = 0.45$$

$$H(\text{class} | Z, X=0) = \frac{3}{6} \left[ \frac{1}{3} \log_2(3) + \frac{2}{3} \log_2\left(\frac{3}{2}\right) \right] + \frac{3}{6} [1 \log_2(1)] = 0.45$$

$$H(\text{class} | X=0) = \frac{4}{6} \log_2\left(\frac{6}{4}\right) + \frac{2}{6} \log_2\left(\frac{6}{2}\right) = 0.91$$

$$IG(\text{class} | Y, X=0) = 0.91 - 0.45 = 0.46$$

$$IG(\text{class} | Z, X=0) = 0.91 - 0.45 = 0.46$$

- both features in this step have the same value. we choose

one randomly, e.g., Z.



$$X=1, Z=1$$

X	Y	Z	class
1	0	1	B
1	1	1	B
1	1	1	B

$$H(\text{class} | X=1, Z=1) = 0$$

$$X=1, Z=0$$

X	Y	Z	class
1	0	0	A

$$H(\text{class} | X=1, Z=0) = 0$$

$$H(\text{class} | Y, X=1, Z=1) = \frac{2}{3} \times 0 + \frac{1}{3} \times 0 = 0$$

$$H(\text{class} | Y, X=1, Z=0) = 0$$

$$IG(\text{class} | Y, X=1, Z=1) = 0$$

$$IG(\text{class} | Y, X=1, Z=0) = 0$$

- there would be no point in using  $Y$  here, since the entropy of the class in this setting before and after using  $Y$  would remain the same.

$$X=0, Z=1$$

X	Y	Z	class
0	1	1	A
0	0	1	A
0	1	1	A

$$H(\text{class} | X=0, Z=1) = 0$$

$$X=0, Z=0$$

X	Y	Z	class
0	0	0	B
0	1	0	A
0	0	0	B

$$H(\text{class} | X=0, Z=0) = \frac{1}{3} \log_2(3) + \frac{2}{3} \log_2\left(\frac{3}{2}\right) = 0.91$$

$$H(\text{class} | Y, X=0, Z=1) = \frac{2}{3} \times 0 + \frac{1}{3} \times 0 = 0$$

$$H(\text{class} | Y, X=0, Z=0) = \frac{1}{3} \times 0 + \frac{2}{3} \times 0 = 0$$

$$IG(\text{class} | Y, X=0, Z=1) = 0$$

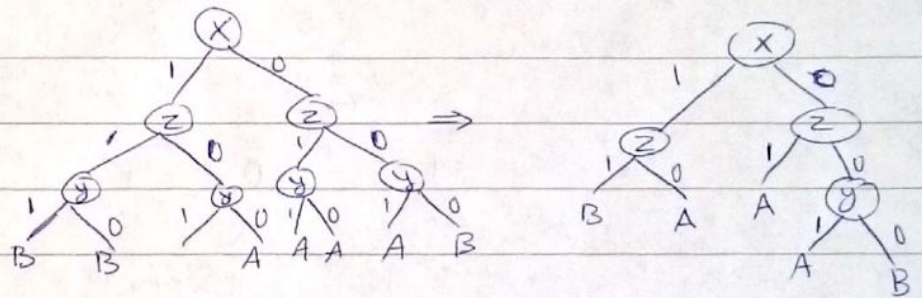
$$IG(\text{class} | Y, X=0, Z=0) = 0.91$$

- here  $Y$  helps, so we choose  $Y$ .



- so in the first level we use  $X$ , in the second level  $Z$  and in the third level, since we want to grow a full tree, we use  $Y$ .

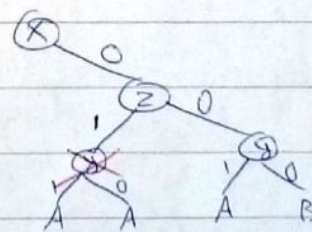
- Full decision tree:



2) the only correctly classified sample in the pruning dataset is:

000 B (Accuracy: 20%)

- if we omit node  $Y$  when  $X=0$  and  $Z=1$  accuracy will remain the same.

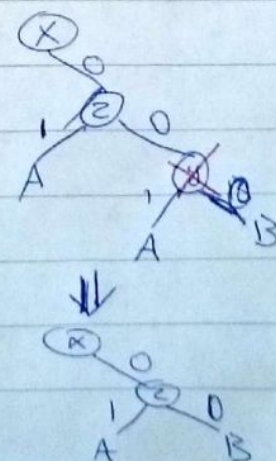


- if we omit  $Y$  when  $X=0$  and  $Z=0$

accuracy will be higher.

correctly classified samples:

000 B  
010 B (Accuracy: 40%)

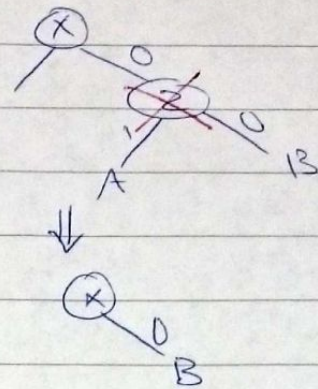




- if we omit  $z$  when  $x=0$

then the accuracy will be even

higher.



Correctly classified samples:

0 1 1 B

0 1 0 B

0 0 1 B

0 0 0 B

(Accuracy: 80%)

So according to the pruning dataset, the tree will

be like this:

