# 1. Decision trees:

a) • Fake news item: FN : Yes = 1, No = 0.

• Origin: One-hot encoding $\underset{A_1}{\underbrace{}}$  Blog :  

Soc. Net. :

Newspaper : $\underset{A_2}{\underbrace{}}$

| | b | s | n |
|---|---|---|---|
| Blog | 1 | 0 | 0 |
| Soc.Net. | 0 | 1 | 0 |
| Newspaper | 0 | 0 | 1 |

"C

• Exessive use of capital letters: CL : $T = 1$, $F = 0$.

Information gain:

$$IG(C|A_i) = H(C) - H(C|A_i)$$

• $H(C) = -\left(\frac{1}{2}\log_2(1/2) + \frac{1}{2}\log_2(1/2)\right) = -\log_2 1/2 = 1.$

• $H(C|A_1) = P(A_1 = \text{"Blog"})P(C|A_1 = \text{"Blog"}) + \cdots$ ⊛

$\begin{cases} H(C|A_1 = \text{"Blog"}) = -\frac{1}{1}\log_2\frac{1}{2} \cdot \cancel{2} = 1 \\ H(C|A_1 = \text{"S.N"}) = 0 \\ H(C|A_1 = \text{"Newspaper"}) = -\frac{2}{3}\log_2\frac{2}{3} - \frac{1}{3}\log_2\frac{1}{3} = \underset{\textcircled{A}}{\underline{0,9182958341}} \end{cases}$

⊛ $H(C|A_1) = \frac{1}{3} \cdot 1 + 0 + \frac{1}{2}\textcircled{A} \simeq \boxed{0,795}$

$\begin{cases} H(C|A_2 = T) = 1 \\ H(C|A_2 = F) = 1 \end{cases}$  $H(C|A_2) = 1 \cdot \frac{2}{3} + 1 \cdot \frac{1}{3} = 1$

$\Rightarrow IG(C|A_1) = 1 - 0,795 = \underline{\underline{0,205}}$

$IG(C|A_2) = 1 - 1 = 0$

First question would be: <u>Origin</u>

Notation: $A/B \iff A=1, B=0$

| Id | $a_1$ Cp letters | $a_2$ blog | $a_3$ soc.net | $a_4$ newspaper | $a_5$ politics/sport | $b$ FN |
|---|---|---|---|---|---|---|
| $v_1$: 1 | 1 | 1 | 0 | 0 | 1 | 1 |
| $v_2$: 2 | 1 | 0 | 1 | 0 | 0 | 1 |
| $v_3$: 3 | 1 | 0 | 0 | 1 | 1 | 0 |
| $v_4$: 4 | 1 | 0 | 0 | 1 | 0 | 0 |
| $v_5$: 5 | 0 | 1 | 0 | 0 | 0 | 0 |
| $v_6$: 6 | 0 | 0 | 0 | 1 | 1 | 1 |

I use $\|\cdot\|_1$ for distance, using rows as vectors in $R^5$.

$v_i = (a_1, a_2, a_3, a_4, a_5)$

$k = 3$

Dato: $(1, 0, 0, 1, 0)$

Preference: lowest id in case of tie.

Distance: $d(v, v_i)$:

1: 3
2: 2
3: 1
4: 0
5: 3
6: 2

3 nearest neighbors: 4, 3, 2
$\downarrow \; \downarrow \; \downarrow$
0  0  1

Prediction: $0 =$ (No)

HW 4-2

## 2 Decision trees.

$$\boxed{\text{yes}\cdot 1,\quad \text{no}\cdot 0}\qquad H(C)=1$$

a)

$$
\begin{cases}
H(C\mid \text{Age-range}=1) = -\left(\frac{2}{3}\log_2\frac{2}{3} + \frac{1}{3}\log_2\frac{1}{3}\right) = 0{,}9183\\[4pt]
H(C\mid \text{Age-range}=2) = \nearrow\ \ = 0{,}9183\\[4pt]
H(C\mid \text{Age-range}=3) = -\frac{1}{2}\log_2(1/2) - \frac{1}{2}\log_2(1/2) = -\log_2(1/2) = 1
\end{cases}
$$

$$\Rightarrow H(C\mid \text{Age-range}) = 0{,}9183\left(\frac{3}{8}+\frac{3}{8}\right) + 1\cdot\frac{2}{8} = 0{,}938725$$

$$\Rightarrow IG(\text{Age-range}) = 0{,}061275$$

$$
\begin{cases}
H(C\mid \text{Kids}=1) = 1\\
H(C\mid \text{Kids}=0) = 1
\end{cases}
\quad
\begin{cases}
H(C\mid \text{Kids}) = 1\\
IG(\text{Kids}) = 0
\end{cases}
$$

$$
\begin{cases}
H(C\mid \text{Changed}=1) = 0\\[4pt]
H(C\mid \text{Changed}=0) = -\frac{1}{5}\log_2\left(\frac{1}{5}\right) - \frac{4}{5}\log_2\left(\frac{4}{5}\right) = 0{,}72193
\end{cases}
$$

$$\Rightarrow H(C\mid \text{Changed}) = 0{,}72193\cdot\frac{5}{8} = \cancel{\phantom{xxxx}}\ 0{,}45121$$

$$\Rightarrow IG(\text{Changed}) = 0{,}5488$$

Root: # Has changed companies before?

HW 4-3

b) Vector: (Age range, Has kids, Has changed c.).

$C_1 = \{1,2\}$ $\quad C_2$ $\qquad\qquad C_2$

**Distance : $\|\cdot\|_1$**

$V_1 = (1,0,0) \in$ No

$V_2 = (1,1,0) \in$ No

$V_3 = (1,1,1) \in$ 𝒮

$V_4 = (2,0,1) \in$ 𝒮

$V_5 = (2,0,0) \in$ No

$V_6 = (2,1,1) \in$ 𝒮

$V_7 = (3,1,0) \in$ No

$V_8 = (3,0,0) \in$ 𝒮

$V = (2,1,0)$

$d(V, V_i) := \begin{cases} 1: 2 \\ 2: 1 \\ 3: 2 \\ 4: 2 \\ 5: 1 \\ 6: 1 \\ 7: 2 \\ 8: 2 \end{cases}$

$k = 3 \Rightarrow 3$ nearest: $\begin{cases} 2 \to \text{No} \\ 5 \to \text{No} \\ 6 \to 𝒮 \end{cases}$

**Prediction : No**

HW 4. 4