# Exercise 8

**Task 1.**

**a)** first, we draw the correspondence graph of each attribute and target:

Professional Status

| civil servant | / Employee | self-employed | non-employed | Low | Contract Duration / medium | \ high |
|---|---|---|---|---|---|---|
| yes | yes | yes | no | yes | no | no yes |
| no | no | no | no | yes | no | no |
| | | | | | no | no |
| | | | | | no | |

Second, we compute the information gain for each attribute:

Professional Status: $\frac{2}{9} \cdot \text{entropy}([1,1]) + \frac{3}{9} \cdot \text{entropy}([1,2]) + \frac{2}{9} \cdot \text{entropy}([1,1]) + 0$

$= \frac{4}{9} + 0.3061 = 0.7505$ bits

Contract Duration: $\frac{2}{9} \cdot 0 + \frac{4}{9} \cdot 0 + \frac{1}{3} \cdot (0.5283 + 0.39) = 0.3061$ bits

Information gain$_p$ = $\text{entropy}([3,6]) - 0.7505 = 0.9183 - 0.7505 = 0.1678$ bits
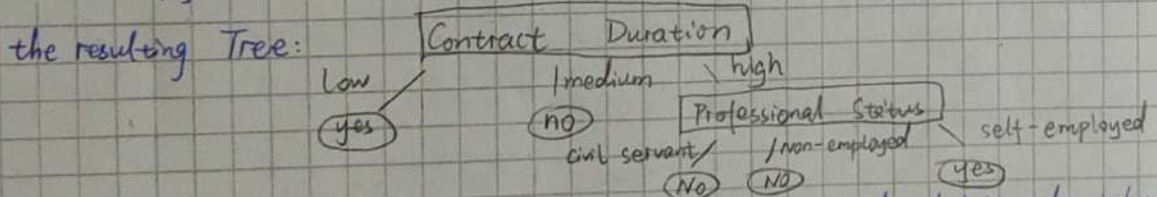
Information gain$_c$ = $0.9183 - 0.3061 = 0.6122$ bits

Therefore, we choose the variable „Contract Duration" as the root:

Since we have only „yes" for „low", „no" for „medium" of variable „Contract Duration", we complete to split for them. Now we only need to compute the information gain for „high":

Entropy (Contract Duration = high) = Entropy([1,2]) = 0.9183

Entropy (Contract Duration = high | Professional Status) = 0

the resulting Tree:

Contract Duration
- Low → yes
- medium → no
- high → Professional Status
  - civil servant → No
  - non-employed → No
  - self-employed → yes

The order of the attributes of the resulting tree is decided by the information gain. We pick the attribute with max information gain for each split.

The leaf nodes is yes/no.

**b)** We classify Number 11 & 12 by the tree from (a):

No. 11 with „Low" of „Contract Duration", it is classified to „yes" of termination.

No. 12 with „medium" of „Contract Duration", it is classified to „no".

Hence Customer with number 11 is at risk of termination.

**c)** Customer with low contract duration tend to terminate; with medium duration tend to not terminate. For customer with high duration, the ones who are

self-employed at risk of termination. Unexpectively, customer with high duration but ~~with~~ non-employed tend to not terminate.

d) Assume dataset $X$ and attribute $V$ with $n$ possible different values ($V_1, V_2, \cdots V_n$)

The information gain of splitting $X$ with $V$ is:

$$Gain(X, v) = Entropy(X) - \sum_{i=1}^{n} \frac{|X_i|}{|X|} Entropy(X|V_1, V_2, \cdots, V_n)$$

The information gain ratio is:

$$Gainratio(X, v) = \frac{Gain(X, v)}{-\sum_{i=1}^{n} \frac{|X_i|}{|X|} \log \frac{|X_i|}{|X|}} \sim splitinformation$$

From the definition of $Gain(X, v)$, the attribute with more possible values (higher possibility with pure classification) might has larger information gain and so be chosen.

Therefore we introduce Gain~ratio~ to reduce the bias caused by higher branching attribute (larger $n$ means larger splitinformation)

The definition of Gini index:

$= 2p(1-p)$

$1 - p^2 - (1-p)^2$, assume there are 2 possible outcomes A, B

$P$ is the possibility of event A.

$$Gini(X, v) = \sum_{i=1}^{n} \frac{|X_i|}{|X|} Gini(X_i) \quad ?$$

e) $Gini(Contract\ Duration) = \frac{3}{4}(\frac{7}{2} \frac{2}{9} \cdot 0 + \frac{4}{9} \cdot 0 + \frac{3}{9} \cdot 2 \cdot \frac{1}{3} \cdot \frac{2}{3} = \frac{1}{27}$

## Task 2.

a)



$R_1: X \le 2, y \le 2.5;$     $R_2: X \le 2, 2.5 < y \le 2.6;$

$R_3: X \le 2, y > 2.6;$     $R_4: 2 < x \le 3, y \le 2;$     $R_5: 3 < x \le 4, y \le 2;$

$R_6: X > 4, y \le 2;$     $R_7: X > 2, y > 2.0$

b) Partition $R_2$ and $R_5$ can be described as overfitted. The subtrees of $\overline{R_2\ R_3}$ and $\overline{R_5\ R_6}$ should be removed.

c) The generalization capability of a classifier mean the classification not only works well on the training data, but also on the testing data. (new data)