

# WORKSHEET 3

1

## PROBLEM 1 : NOTES

### K- NEAREST NEIGHBOURS

Note :  $L_p$  - norm =  $\left( \sum_{d=1}^D |x_{id}|^p \right)^{1/p}$

The Euclidean norm is the  $L_2$  - norm.

The Euclidean distance is defined as

$$D(\underline{x}_i, \underline{x}_j) = \sqrt{\sum_{d=1}^D (x_{id} - x_{jd})^2}$$

For  $\underline{x}_i = (x_{i1}, \dots, x_{iD})^T$ ,  $x_{id} \in \mathbb{R}$

$\underline{x}_j = (x_{j1}, \dots, x_{jD})^T$ ,  $x_{jd} \in \mathbb{R}$

### Classification :

Given  $\{t_i, \underline{x}_i\}_{i=1}^M$   $t_i \in \{c_1, \dots, c_C\}$

That is, the target  $t_i$  belongs to one of  $C$  classes.

Then in order to classify some unknown target observed at the location  $\underline{x}_*$ ,

first identify the  $k$  observations for which

$D(\underline{x}_*, \underline{x}_i)$  is smallest to find the subset  $\{t_j, \underline{x}_j\}_{j=1}^k$

Then  $t_* = \arg \max_{u \in \{t_1, \dots, t_k\}} \sum_{j=1}^k \delta(u, t_j)$

where  $\delta(a, b) = \begin{cases} 1 & \text{if } a=b \\ 0 & \text{otherwise} \end{cases}$

is a delta function

which is to say that the predicted value / class for  $t_*$  is the class to which most of its neighbours belong.

Classification error is of the form

$$1 - \delta(t_n, t_n^*)$$

where  $t_n$  is the "true" value of the prediction  $t_n^*$ .

PROBLEM 2

FIND THE ID3 CLASSIFIER FOR THE PROVIDED DATASET

NOTE: A GENERAL DESCRIPTION OF THE ID3 ALGORITHM.

GIVEN A DATA SET  $S = \{t_n, x_n\}_{n=1}^N$

where  $t_n \in \{z_1, \dots, z_c\}$  is the target variable which belongs to one of  $C$  classes

and  $x_n = (x_{n1}, \dots, x_{nd})^T$  is a vector of explanatory variables where  $x_{nj} \in \mathcal{V}_j$ , that is to say that the  $j^{\text{th}}$  explanatory variable comes from some set  $\mathcal{V}_j$ .

$$\text{when } p_i = \frac{1}{N} \sum_{n=1}^N \delta(t_n, z_i),$$

$\delta$  being a delta function (see problem 1 notes)

The entropy of  $S$  is

$$E(S) = - \sum_{i=1}^c p_i \log_2(p_i)$$

$E(S)$  is a measure of the uncertainty of the dataset where  $E(S) \in [0, 1]$  and

$$\text{for } E(S) = 0, \quad C = 1$$

$$\text{and } E(S) = 1, \quad p_i = \frac{1}{C}$$

ID3 provides a classification tree by splitting the data by the unused attribute which gives the greatest information gain

For each attribute, denoted  $A$ , or more precisely  $x_{:j} = (x_{1j}, \dots, x_{Nj})$

The Information Gain is given by

$$IG(S, A) = IG(S, x_{:j}) = E(S) - \sum_{y \in Y_j} \frac{|S_{\{x_{:j}=y\}}|}{|S|} E(S_{\{x_{:j}=y\}})$$

where  $S_{\{x_{:j}=y\}}$  is the set of observations in  $S$  for which  $x_{:j} = y$  and  $|S|$  is the number of observations in  $S$

Applying the algorithm to the provided dataset:

$\mathcal{E}(S)$  where  $t = Y$ ,  $x_{:1} = V$ ,  $x_{:2} = W$ ,  $x_{:3} = X$

$$\mathcal{E}(S) = - \left( \frac{2}{5} \log_2 \left( \frac{2}{5} \right) + \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right)$$

$$\approx 0.97$$

$$\mathcal{E}(S_{\{x=1\}}) = - \left( 0 \cdot \log_2(0) + 1 \cdot \log_2(1) \right) = 0$$

$$\mathcal{E}(S_{\{x=0\}}) = - \left( \frac{1}{2} \log_2 \left( \frac{1}{2} \right) + \frac{1}{2} \log_2 \left( \frac{1}{2} \right) \right) = 1$$

$$E(S_{\{W=1\}}) = E(S_{\{V=1\}}) = -\left(\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \frac{2}{3} \log_2\left(\frac{2}{3}\right)\right) \\ \approx 0.91$$

$$E(S_{\{W=0\}}) = E(S_{\{V=0\}}) = 1$$

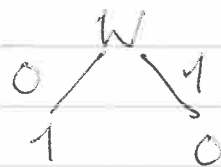
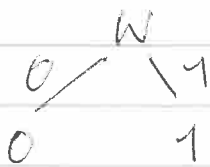
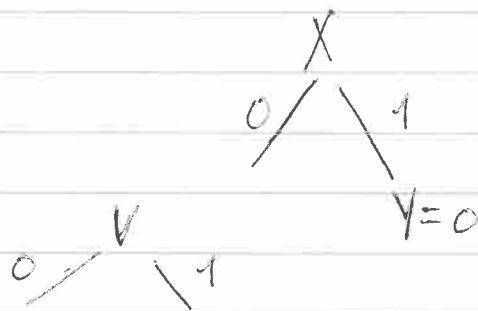
$$I_G(S, X) \approx 0.97 - \frac{2}{5}(0) - \frac{3}{5}(1) \approx 0.17$$

$$I_G(S, W) = I_G(S, V) \approx 0.97 - \frac{3}{5}(1) - \frac{2}{5}(0.91) \approx 0.02$$

Thus the first split in the tree is at attribute  $X$ .

Continuing the process for the next split in the tree will show  $W$  and  $V$  to be equally explanatory i.e.  $I_G(S_{\{X=0\}}, W) = I_G(S_{\{X=0\}}, V)$

a ~~the~~ final tree produced by the algorithm is (there are two possible solutions)



And  $Y$  can be fully classified by  $X, W, V$  where  $W$ , and  $V$  are interchangeable