

## Decision Trees

### Q2. Entropy, Gini Index and Decision Trees

a) - what is the formula for computing entropy?

- Among the probability distributions  $p$  and  $p'$  specified below, which one has a higher entropy?

- Explain why, without computing the entropies of them.

$$p = (0.33, 0.33, 0.34) \quad p' = (0.8, 0.1, 0.1)$$

Let  $D$  is the

$$H(D) = -(p_0 \log_2 p_0 + p_1 \log_2 p_1) \quad \text{Entropy} = -\sum_{i=0}^A p_i \log_2 p_i$$

The entropy for  $p$  is high since the probability distribution for all are same.

b) What is the Gini index of the probability distribution  $(\frac{2}{3}, \frac{1}{3})$ ?

$$G = 1 - \sum_{i=1}^n p_i^2$$

$$\therefore G = 1 - \left[ \frac{2}{3} \times \frac{2}{3} + \frac{1}{3} \times \frac{1}{3} \right] = 1 - [0.44 + 0.11] = 0.445.$$

c) Consider the following set of datapoints, each containing 3 entries of the format  $(x_1, x_2, Y)$ .

i.e., in each datapoint, the first 2 entries are attributes, and the last entry is its label. The labels can be either "T" or "F".

$$D = \{(1, 4, T), (2, 3, T), (3, 0, T), (6, 2, F), (7, 1, F), (8, 1, F)\}$$

Now consider the following two possible splits:

- Split 1:  $x_1 \leq 4$  and  $x_1 > 4$

- Split 2:  $x_2 \leq 2.5$  and  $x_2 > 2.5$

For split 2, which datapoints go to the left (under  $x_2 \leq 2.5$ ), and which datapoints go to the right (under  $x_2 > 2.5$ )?

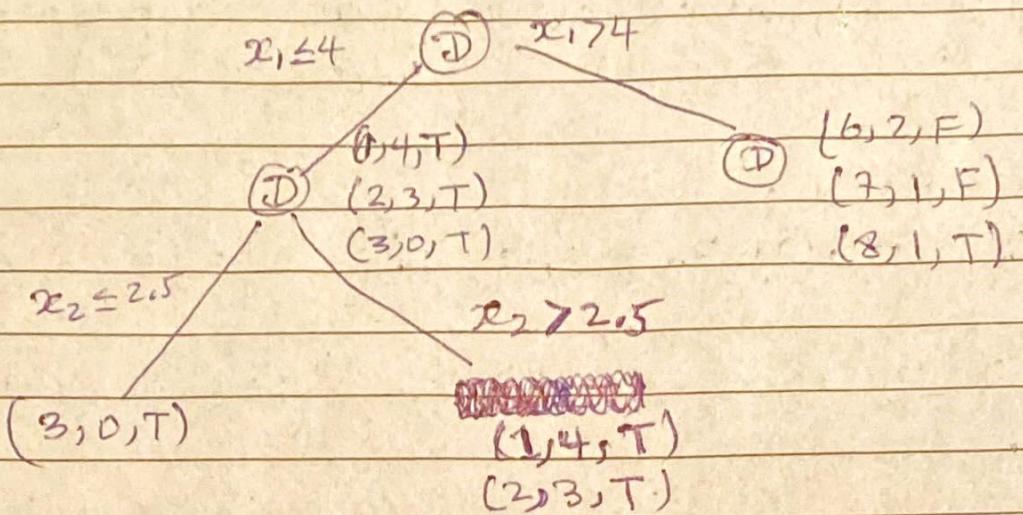
$(x_1, x_2, y)$

Given datapoint

$$D = \{(1, 4, T), (2, 3, T), (3, 0, T), (6, 2, F), (7, 1, F), (8, 1, F)\}$$

split 1:  $x_1 \leq 4$  &  $x_1 > 4$

split 2:  $x_2 \leq 2.5$  &  $x_2 > 2.5$



For split 2 datapoints which go under  $x_2 \leq 2.5$  are  $(3, 0, T)$  and datapoint which are under  $x_2 > 2.5$  are  $(1, 4, T) \neq (2, 3, T)$ .

d) You are going to generate a decision tree for the dataset D specified in part (c).

You want to choose a split into regions  $R_1$  and  $R_2$  which has a small value of

$$|R_1| \cdot (\text{Gini index of region } R_1) + |R_2| \cdot (\text{Gini index of region } R_2)$$

Among split 1 and split 2, which one is better? justify your answer  
Gini index for  $R_1$

$$x \leq 4 \Rightarrow 3T \text{ and } 0F \quad x > 4 \Rightarrow 1T \text{ and } 2F$$

$$1 - \left[ \left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right] = 0$$

$$1 - \left[ \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right]$$

$$1 - 0.555 = 0.4445$$

$$\text{Weighted entropy} = \left( \frac{3}{6} \times 0 + \frac{3}{6} \times 0.4445 \right) = \frac{0 + 0.4445}{2} = 0.225$$

## Decision Trees

(2.d) For Region 2

$$x \leq 2.5 \Rightarrow 1 \neq 0 \text{ F}$$

$$= 1 - [(1)^2 + (0)^2] = 0$$

$$x > 2.5 \Rightarrow 2 \text{ T and } 0 \text{ F}$$

$$1 - [(1)^2 + (0)^2] = 0$$

$$\text{Weighted} = \frac{1}{3} + 0 + \frac{2}{3} + 0 = 0$$

Split 1 classifies the difference same as of split 2 through entropy. Split 1 is more better than split 2.

- e) For any decision tree  $T$  of a classification problem, let its "penalized Gini" be

$$\alpha|T| + \sum_R (\text{Gini index of region } R)$$

where,

- $\alpha \geq 0$
- $|T|$  is the number of leaves of the tree, and
- $R$  runs over all the regions/leaves of the tree.

What techniques are we using via the introduction of the term  $\alpha|T|$ ?

Explain the motivation for introducing this term. (5 sentences)

We need to get better  $\alpha$  and prone for best set of techniques used for  $\alpha|T|$  as penalize RSS and Pruning Trees.

$$\alpha|T| + \sum_{m=1}^{|T|} \text{RSS}(R_m)$$

- $|T| = \text{number of leaves}$     $\alpha = 0$ , penalized RSS = 0

$y = \alpha + \text{RSS}$    • we get critical  $\alpha$  and penalized RSS.

Q1.d) In a leaf of a decision tree

Exam paper: CS5400R Data Analysis

(3)

Q2. Prediction (Decision tree).

- b) Compute the probability to assign the test object to each of two classes according to a Decision Tree model  
 - based on the partition of the feature space defined by

$$(1) R_1 = \{x \text{ such that } x_2 \leq \xi\}, R_2 = \{x \text{ such that } x_2 > \xi\}, \xi = 7$$

and training data set D. Justify all steps in your solution (Max 4 lines)  
 the test object (4,4)

$$R_1 \Rightarrow x_2 \leq 7$$

$$x_2 \leq 7$$

$$x_2 > 7$$

$$R_1$$

$$R_2$$

$(x_1, x_2)$	$y$
• (3, 8)	1
- (5, 6)	0
- (5, 5)	0
- (4, 0)	0
• (8, 9)	0
- (9, 6)	0
- (1, 2)	1
- (7, 0)	0
- (9, 5)	0
- (6, 2)	0

c) The Gini Impurity of a region  $R$  is defined as

$$G(R) = \sum_{i=0}^1 p_i (1-p_i) = 1 - \sum_{i=0}^1 p_i^2$$

Where:

- $\bullet p_i = \frac{1}{Z_R} \sum_{(x,y)} 1_{y=i} 1_{x \in R}$ ,

- $\bullet Z_R = \sum_{(x,y)} 1_{x \in R}$ ,

D.e.  $p_i$  is the probability of an object to have label  $y=i$  given it's belongs to region  $R$ .

Compute the information gain associated with (1)

D.e.

$$Z_{R_1 \cup R_2} G(R_1 \cup R_2) - (Z_{R_1} G(R_1) + Z_{R_2} G(R_2))$$

Where,  $R_1$ ,  $R_2$ , and  $Z_R$  are defined as above.

## Q1. Decision Trees

- Decision trees are grown by splitting the feature space in increasingly pure regions.
- There exist several methods to measure the impurity of a set of labelled observations, e.g.  $D = \{(x_n, y_n)\}_{n=1}^N$  by looking at the object labels,  $y_n$ .
- Two examples of impurity measures are
  - the Gini index and
  - the entropy.

a) Consider the case where  $y_n \in \{0, 1\}$  and  $x_n \in \{\bar{i}\}_{i=1}^{10}$  for all  $n = 1, \dots, 11$  and  
- the task

## Decision Trees

- Entropy - Measures the purity of split.
  - ranges from 1 to 0, [0, 1]
  - when,  $E=0$ , pure split (pure subset)
  - when,  $E=1$ , worst split (impure subset)

Formula:  $\text{Entropy} = - \sum_{i=0}^n P_i \log_2 P_i$

- Information Gain

$$IG(S, f) = H(S) - \sum_{V \in Val} \frac{|S_V|}{|S|} H(S_V)$$

where •  $S_V$  the subset after splitting •  $S$  the subset

- $H(S_V)$  Entropy after splitting
- $H(S)$  Entropy

- The higher the IG the better

- Gini Index (Impurity)

$$GI = 1 - \sum_{i=1}^n P_i^2$$

- ranges between 0.5 and 0.

$$Q2. a) \text{Entropy} = -\sum_{i=0}^n p_i \log_2 p_i$$

The ~~probabilistic~~ p has higher entropy since the probability distributions for all are the same.

$$b) \left(\frac{2}{3}, \frac{1}{3}\right) \text{ Gini Index} = 1 - \sum_{i=1}^n p_i^2$$

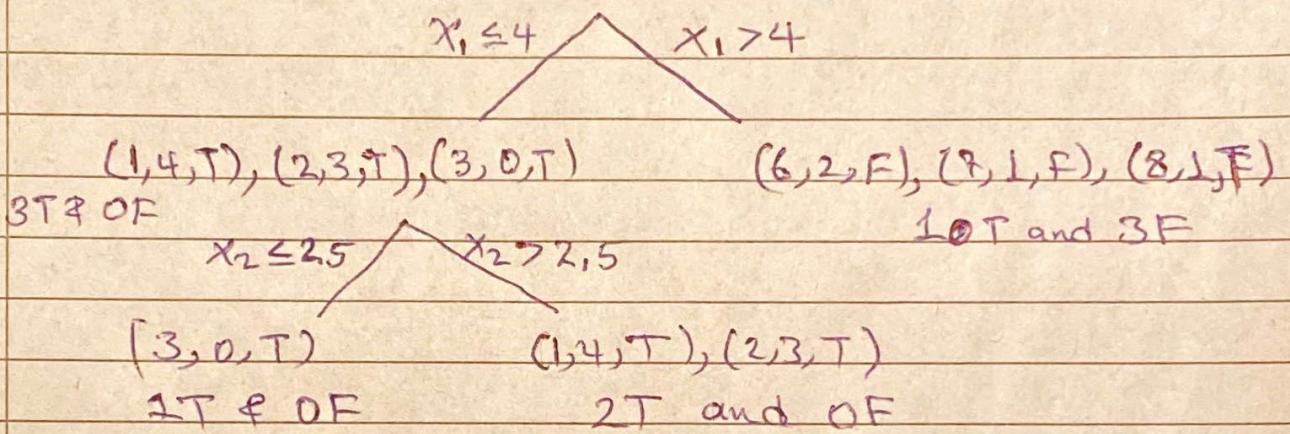
$$\therefore GI = 1 - \left[ \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 1 - \left[ \left(\frac{4}{9} + \frac{1}{9}\right) \right] = 1 - \frac{5}{9} = 0.444$$

$$c) D = \{(1, 4, T), (2, 3, T), (3, 0, T), (6, 2, F), (7, 1, F), (8, 1, F)\}$$

split 1:  $x_1 \leq 4 \neq x_1 > 4$

format  $(x_1, x_2, Y)$

split 2:  $x_2 \leq 2.5 \neq x_2 > 2.5$



$$d) \text{Gini Index} = 1 - \sum_{i=1}^n p_i^2$$

For  $R_1$ , when  $x_1 \leq 4$ , 3T ≠ OF and when  $x_1 > 4$ , 1T ≠ 3F

$$GI_{x_1 \leq 4} = 1 - \left[ \left(\frac{3}{8}\right)^2 + \left(\frac{5}{8}\right)^2 \right] = 0 \quad GI_{x_1 > 4} = 1 - \left[ \left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] = 1 - \frac{1}{3} = \frac{2}{3} = 0.667$$

Weighted =  $|R_1| \cdot (GI \text{ of } R_1)$   
entropy

$$= \left( \frac{3}{8} \times 0 + \frac{5}{8} \times 0.667 \right) = \frac{0 + 0.667}{2} = 0.333$$

For  $R_2$ , when  $x_2 \leq 2.5$ , 1T ≠ OF and when  $x_2 > 2.5$ , 2T ≠ OF

$$GI_{x_2 \leq 2.5} = 1 - \left[ \left(\frac{1}{4}\right)^2 + \left(\frac{3}{4}\right)^2 \right] = 0 \quad GI_{x_2 > 2.5} = 1 - \left[ \left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right] = 1 - \frac{1}{3} = \frac{2}{3} = 0.667$$

$$\text{Weighted} = |R_2| \cdot (GI \text{ of } R_2) = \frac{1}{3} \times 0 + \frac{2}{3} \times 0 = 0$$

∴ result 1 is better