

Sohaib Syed
CS422
10/17/21
Dr. Gurbani

Homework 5

Problem 1.1:

Exercise 2:

A.

Probability of C0 is $10/20=1/2$, probability of C1 is $10/20=1/2$.

Thus Gini index: is $1 - \sum_{i=1}^n (P_i)^2 = 1 - (0.5^2 + 0.5^2) = .5$

B.

By Gini index definition: give a zero impurity value if a node contains instances from a single class (Tan, Chapter 3). Thus the Gini index of Customer ID is 0.

C.

To find Gini index of gender we would first have to get gini index of each gender since both male and female can be classified as C0 or C1. We would then get the sum of the Gini indexes of each gender multiplied by their probability in general

$$\text{Male Gini} = 1 - (0.6^2 + 0.4^2) = 0.48$$

$$\text{Female Gini} = 1 - (0.4^2 + 0.6^2) = 0.48$$

$$\text{Total gender Gini} = 10/20(0.48) + 10/20(0.48) = \mathbf{0.48}$$

D.

$$\text{Gini family car} = 1 - ((1/4)^2 + (3/4)^2) = .375$$

$$\text{Gini sports car} = 1 - ((8/8)^2 + 0^2) = 0$$

$$\text{Gini luxury car} = 1 - ((1/8)^2 + (7/8)^2) = .219$$

$$\text{Gini of car} = 4/20(.375) + 8/20(0) + 8/20(.219) = \mathbf{.163}$$

E.

$$\text{Gini small: } 1 - ((\frac{3}{7})^2 + (\frac{4}{7})^2) = .48$$

$$\text{Gini medium: } 1 - ((\frac{3}{7})^2 + (\frac{4}{7})^2) = .490$$

$$\text{Gini large: } 1 - ((\frac{2}{4})^2 + (\frac{2}{4})^2) = .5$$

$$\text{Gini extra large: } 1 - ((\frac{2}{4})^2 + (\frac{2}{4})^2) = .5$$

$$\text{Gini shirt size} = 5/20(.48) + 7/20(.490) + 4/20(.5) + 4/20(.5) = \mathbf{.4915}$$

F.

The better attribute is Car Type because it has the lowest Gini index. This means that it is the least impure.

G.

Customer ID has the lowest Gini index, but because for every new entry a new ID is necessary it would not help to predict the class.

Exercise 3:

A.

$$\text{Entropy} = - \sum_{i=1}^n P_i * \log_2(P_i)$$

$$-(((4/9)*\log_2(4/9)) + ((5/9)*\log_2(5/9))) = \mathbf{.991}$$

B.

This problem is sort of like a multiclass gini index problem so first we find entropy of T and then of F multiply to probability of T and F respectively. Sum them up.

A1:

$$(4/9)*-((3/4)*\text{Log}_2(3/4))+(1/4*\text{Log}_2(1/4))=.361$$

$$(5/9)*-((1/5)*\text{Log}_2(1/5))+(4/5*\text{Log}_2(4/5))=.401$$

$$.361+.401=.762$$

$$\text{Gain} = E(b)-E(a)=.991-.762=.229$$

A2:

$$(5/9)*-((2/5)*\text{Log}_2(2/5))+(3/5*\text{Log}_2(3/5))=.539$$

$$(4/9)*-((2/4)*\text{Log}_2(2/4))+(2/4*\text{Log}_2(2/4))=.444$$

$$.539+.444=.984$$

$$\text{Gain} = E(b)-E(a)=.991-.984=.007$$

C.

$$a_3=1.0$$

$$a_3 \leq 1: (1/9)*-((1/1)*\text{Log}_2(1/1))+(0/1)*\text{Log}_2(0/1))=0$$

$$a_3 > 1: (8/9)*-((3/8)*\text{Log}_2(3/8))+(5/8)*\text{Log}_2(5/8))=.848$$

$$\text{gain} = .991-.848=.143$$

$$a_3=3.0$$

$$a_3 \leq 3: (2/9)*-((1/2)*\text{Log}_2(1/2))+(1/2)*\text{Log}_2(1/2))=.222$$

$$a_3 > 3: (7/9)*-((3/7)*\text{Log}_2(3/7))+(4/7)*\text{Log}_2(4/7))=.766$$

$$\text{gain} = .991-(.222+.766)=.002$$

$$a_3=4.0$$

$$a_3 \leq 4: (3/9)*-((2/3)*\text{Log}_2(2/3))+(1/3)*\text{Log}_2(1/3))=.306$$

$$a_3 > 4: (6/9)*-((2/6)*\text{Log}_2(2/6))+(4/6)*\text{Log}_2(4/6))=.612$$

$$\text{gain} = .991 - (.306 + .612) = .073$$

$$a_3 = 5.0$$

$$a_3 \leq 5: (5/9) * -((2/5) * \log_2(2/5) + (3/5) * \log_2(3/5)) = .539$$

$$a_3 > 5: (4/9) * -((2/4) * \log_2(2/4) + (2/4) * \log_2(2/4)) = .444$$

$$\text{gain} = .991 - (.539 + .444) = .008$$

$$a_3 = 6$$

$$a_3 \leq 6: (6/9) * -((3/6) * \log_2(3/6) + (3/6) * \log_2(3/6)) = .667$$

$$a_3 > 6: (3/9) * -((1/3) * \log_2(1/3) + (2/3) * \log_2(2/3)) = .306$$

$$\text{gain} = .991 - (.667 + .306) = .018$$

$$a_3 = 7$$

$$a_3 \leq 7: (8/9) * -((4/8) * \log_2(4/8) + (4/8) * \log_2(4/8)) = .889$$

$$a_3 > 7: (1/9) * -((0/1) * \log_2(0/1) + (1/1) * \log_2(1/1)) = 0$$

$$\text{gain} = .991 - (.889 + 0) = .102$$

$$a_3 = 8$$

$$a_3 \leq 8: (9/9) * -((4/9) * \log_2(4/9) + (5/9) * \log_2(5/9)) = .991$$

$$a_3 > 8: (0/9) * -((0/9) * \log_2(0/9) + (0/9) * \log_2(0/9)) = 0$$

$$\text{gain} = .991 - (.991 + 0) = 0$$

The best split happens at $a_3 = 1.0$, with gain of .143.

D.

According to the information gain a_1 is the best split with gain of .229

E.

$$\text{Error}(a_1) = \frac{1}{4}(4/9) + \frac{1}{5}(5/9) = 2/9$$

$$\text{Error}(a_2) = \frac{3}{5} \left(\frac{5}{9} \right) + \frac{2}{4} \left(\frac{4}{9} \right) = \frac{5}{9}$$

Since a_1 has the lower error rate, a_1 is the best split.

F.

A1:

$$1 - \left(\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right) = .375$$

$$1 - \left(\left(\frac{4}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \right) = .32$$

$$\frac{4}{9}(.375) + \frac{5}{9}(.32) = .344$$

A2:

$$1 - \left(\left(\frac{2}{5} \right)^2 + \left(\frac{3}{5} \right)^2 \right) = .48$$

$$1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) = .5$$

$$\frac{5}{9}(.48) + \frac{4}{9}(.5) = .489$$

A1 has a lower gini index so a_1 is the better split.

Exercise 5:

A.

Information gain $E(b) - E(a)$

$$E(b) = - \left(\left(\frac{4}{10} \right) \log_2 \left(\frac{4}{10} \right) + \left(\frac{6}{10} \right) \log_2 \left(\frac{6}{10} \right) \right) = .971$$

$$A(T) = \left(\frac{7}{10} \right) * - \left(\left(\frac{4}{7} \right) \log_2 \left(\frac{4}{7} \right) + \left(\frac{3}{7} \right) \log_2 \left(\frac{3}{7} \right) \right) = .690$$

$$A(F) = \left(\frac{3}{10} \right) * - \left(\left(\frac{3}{3} \right) \log_2 \left(\frac{3}{3} \right) + \left(\frac{0}{3} \right) \log_2 \left(\frac{0}{3} \right) \right) = 0$$

$$E(a) = .971 - .690 = .281$$

$$B(T) = \left(\frac{4}{10} \right) * - \left(\left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) + \left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) \right) = .325$$

$$B(F) = \left(\frac{6}{10} \right) * - \left(\left(\frac{5}{6} \right) \log_2 \left(\frac{5}{6} \right) + \left(\frac{1}{6} \right) \log_2 \left(\frac{1}{6} \right) \right) = .390$$

$$E(a) = .971 - (.325 + .390) = .256$$

Since splitting on A has the greater information gain that is where the induction algorithm would choose.

B.

$$\text{Gini of entire class: } 1 - (.4^2 + .6^2) = .48$$

$$\text{Gini of A(T): } 1 - ((4/7)^2 + (3/7)^2) = .490$$

$$\text{Gini of A(F): } 1 - ((3/3)^2 + (0/3)^2) =$$

$$7/10(.490) + 3/10(0) = \mathbf{.343}$$

$$\text{Gini of B(T): } 1 - ((3/4)^2 + (1/4)^2) = .375$$

$$\text{Gini of B(F): } 1 - ((5/6)^2 + (1/6)^2) = .278$$

$$4/10(.375) + 6/10(.278) = \mathbf{.317}$$

$$\mathbf{\text{Gain a} = .48 - .343 = .137}$$

$$\mathbf{\text{Gain b} = .48 - .317 = .163}$$

Gini index of B is lower, and the gain from splitting at B is greater.

C.

Yes, the two impurity measures can favor different attributes. This is because the two quantities have different equations and the gains of both are from different weights. The graph is misleading because it shows the impurity of the two measures, but not what the gains are of the two measures. An example of this is part a and part b of this problem.

Problem 1.2:

18

A.

If the data is half “+” and half “-”, and the classifier always predicts “+” then that means it is missing all the “-” which make up half the data. Thus the error should be 50%.

B.

Assuming data is still 50/50, the classifier would be 80% correct for the “+” and 20% correct “-”. This also means that 20% of the time a “+” would be classified as “-”, a false negative (FN) and 80% of the time a “-” is classified as “+”, a false positive (FP). Error is $FP+FN/\text{total}$, so $(.8*.5 + .2*.5)/1 = .5$ or again 50%.

C.

In a similar situation to part a, if $\frac{2}{3}$ data is “+” and $\frac{1}{3}$ data is “-”, then always predicting “+” would mean that $\frac{1}{3}$ of the data was misclassified, so the error would be $\frac{1}{3}$ or 33%

D.

$\frac{2}{3} = \text{“+”}$; $\frac{1}{3} = \text{“-”}$

Prediction “+” = $\frac{2}{3}$; False negatives = $\frac{1}{3}$

Prediction “-” = $\frac{1}{3}$; False positives = $\frac{2}{3}$

Error = $(FP+FN)/\text{Total} = (\frac{2}{3}*\frac{1}{3} + \frac{1}{3}*\frac{2}{3})/1 = .444$ or 44.4%