

CS 5350/6350: Machine Learning Fall 2022

Homework 1

Handed out: 6 Sep, 2022
Due date: 11:59pm, 23 Sep, 2022

1 Decision Tree [40 points + 10 bonus]

x_1	x_2	x_3	x_4	y
0	0	1	0	0
0	1	0	0	0
0	0	1	1	1
1	0	0	1	1
0	1	1	0	0
1	1	0	0	0
0	1	0	1	0

Table 1: Training data for a Boolean classifier

1. [7 points] Decision tree construction.

- (a) [5 points] Use the ID3 algorithm with information gain to learn a decision tree from the training dataset in Table 1. Please list every step in your tree construction, including the data subsets, the attributes, and how you calculate the information gain of each attribute and how you split the dataset according to the selected attribute. Please also give a full structure of the tree. You can manually draw the tree structure, convert the picture into a PDF/EPS/PNG/JPG format and include it in your homework submission; or instead, you can represent the tree with a conjunction of prediction rules as we discussed in the lecture.

Solution:

Information Gain:

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

Current Entropy:

$$p = \frac{2}{7}; n = \frac{2}{7}$$

$$H(Y) = -\frac{2}{7} \log_2 \frac{2}{7} - \frac{5}{7} \log_2 \frac{5}{7} = 0.8631$$

Attribute: X_1

$$X_1 = 0$$

$$p = \frac{1}{5}; n = \frac{4}{5}$$

$$H(X_1 = 0) = -\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} = 0.7218$$

$$X_1 = 1$$

$$p = \frac{1}{2}; n = \frac{1}{2}$$

$$H(X_1 = 1) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\text{Entropy} = \frac{5}{7} * H(X_1 = 0) + \frac{2}{7} * H(X_1 = 1) = 0.8013$$

$$\text{Information Gain}_{(Y, X_1)} = 0.8631 - 0.8013 = 0.06$$

Attribute: X_2

$$X_2 = 0$$

$$p = \frac{2}{3}; n = \frac{1}{3}$$

$$H(X_2 = 0) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.7218$$

$$X_2 = 1$$

$$p = 0; n = \frac{4}{4}$$

$$H(X_2 = 1) = 0$$

$$\text{Entropy} = \frac{3}{7} * H(X_2 = 0) + \frac{4}{7} * H(X_2 = 1) = 0.3934$$

$$\text{Information Gain}_{(Y, X_2)} = 0.8631 - 0.3934 = 0.4697$$

Attribute: X_3

$$X_3 = 0$$

$$p = \frac{1}{4}; n = \frac{3}{4}$$

$$H(X_3 = 0) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.811$$

$$X_3 = 1$$

$$p = \frac{1}{3}; n = \frac{2}{3}$$

$$H(X_3 = 1) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.918$$

$$\text{Entropy} = \frac{4}{7} * H(X_3 = 0) + \frac{3}{7} * H(X_3 = 1) = 0.857$$

$$\text{Information Gain}(y, X_3) = 0.8631 - 0.857 = 0.0061$$

Attribute: X_4

$$X_4 = 0$$

$$p = 0; n = \frac{4}{4}$$

$$H(X_4 = 0) = 0$$

$$X_4 = 1$$

$$n = \frac{1}{3}; p = \frac{2}{3}$$

$$H(X_4 = 1) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$$

$$\text{Entropy} = \frac{4}{7} * H(X_4 = 0) + \frac{3}{7} * H(X_4 = 1) = 0.3934$$

$$\text{Information Gain}(y, X_4) = 0.8631 - 0.3934 = 0.4697$$

Both X_4 and X_2 have **highest Information gain**. We can split on either.

Lets split on X_2

$X_2 = 1$: leaf node $Y = 0$ and $X_2 = 0$:

$$p = \frac{2}{3}; n = \frac{1}{3}$$

$$H(X_2 = 0) = 0.7218$$

Now the for highest Information gain, we can split on $X_2 = 0$

$$\text{Information Gain}(X_2 = 0, X_1) = 0.7218 - \frac{2}{3} * [-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}] - 0 * \frac{1}{3} = 0.0618$$

$$\text{Information Gain}(X_2 = 0, X_3) = 0.7218 - \frac{2}{3} * [-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}] - 0 * \frac{1}{3} = 0.0618$$

$$\text{Information Gain}(X_2 = 0, X_4) = 0.7218 - \frac{2}{3} * 0 - 0 * \frac{1}{3} = 0.7218$$

Let split on X_4

$X_4 = 0$: leaf node $Y = 0$ and $X_4 = 1$: leaf node $Y = 1$

The decision tree is shown in Fig. 1.

- (b) [2 points] Write the boolean function which your decision tree represents. Please use a table to describe the function — the columns are the input variables and label, x_1, x_2, x_3, x_4 and y ; the rows are different input and function values.

Solution: The function is obtained from the table shown in Table 2.

$$Y = \bar{X}_2 X_4$$

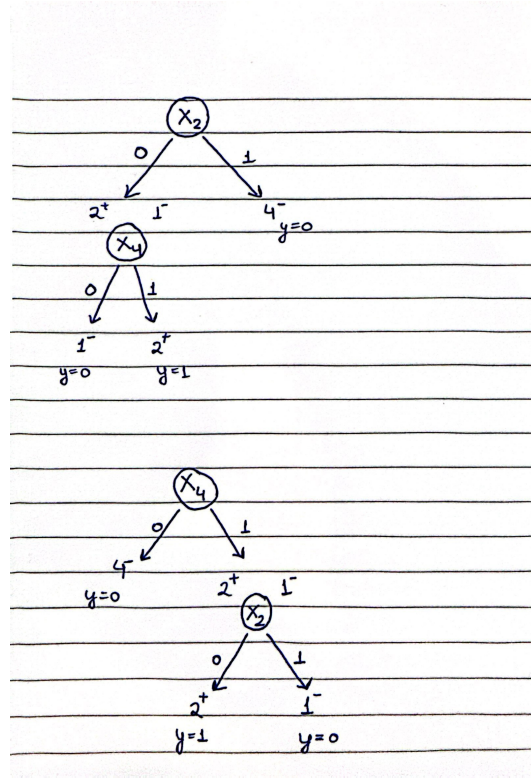


Figure 1: Decision Tree

x_1	x_2	x_3	x_4	y
0	0	0	0	0
0	0	0	1	1
0	0	1	0	0
0	0	1	1	1
0	1	0	0	0
0	1	0	1	0
0	1	1	0	0
0	1	1	1	0
1	0	0	0	0
1	0	0	1	1
1	0	1	0	0
1	0	1	1	1
1	1	0	0	0
1	1	0	1	0
1	1	1	0	0
1	1	1	1	0

Table 2: Complete boolean dataset

2. [17 points] Let us use a training dataset to learn a decision tree about whether to play tennis (**Page 43, Lecture: Decision Tree Learning**, accessible by clicking the link <http://www.cs.utah.edu/~zhe/teach/pdf/decision-trees-learning.pdf>). In the class, we have shown how to use information gain to construct the tree in ID3 framework.

- (a) [7 points] Now, please use majority error (ME) to calculate the gain, and select the best feature to split the data in ID3 framework. As in problem 1, please list every step in your tree construction, the attributes, how you calculate the gain of each attribute and how you split the dataset according to the selected attribute. Please also give a full structure of the tree.

Solution::

$$Gain(S, A) = ME(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} ME(S_v)$$

$$ME(S) = \frac{5}{14}$$

Outlook:

$$\text{Sunny, ME (Sunny)} = \frac{2}{5}$$

$$\text{Overcast, ME (Overcast)} = 0$$

$$\text{Rain, ME (Rain)} = \frac{2}{5}$$

Temperature:

$$\text{Hot, ME (Hot)} = \frac{2}{4}$$

$$\text{Mild, ME (Mild)} = \frac{2}{6}$$

$$\text{Cool, ME (Cool)} = \frac{1}{4}$$

Humidity:

$$\text{High, ME (High)} = \frac{3}{7}$$

$$\text{Normal, ME (Normal)} = \frac{1}{7}$$

$$\text{Low : , ME (Low)} = 0$$

Wind:

$$\text{Strong, ME (Strong)} = \frac{3}{6}$$

$$\text{weak, ME (weak)} = \frac{2}{8}$$

$$Gain(S, Outlook) = \frac{5}{14} - \frac{2}{5} * \frac{5}{14} - 0 - \frac{2}{5} * \frac{5}{14} = 0.071$$

$$Gain(S, Temperature) = \frac{5}{14} - \frac{2}{4} * \frac{4}{14} - \frac{1}{4} * \frac{4}{14} - \frac{2}{6} * \frac{6}{14} = 0$$

$$Gain(S, Humidity) = \frac{5}{14} - \frac{3}{7} * \frac{7}{14} - \frac{1}{7} * \frac{7}{14} - 0 = 0.071$$

$$Gain(S, Wind) = \frac{5}{14} - \frac{3}{6} * \frac{6}{14} - \frac{2}{8} * \frac{8}{14} = 0$$

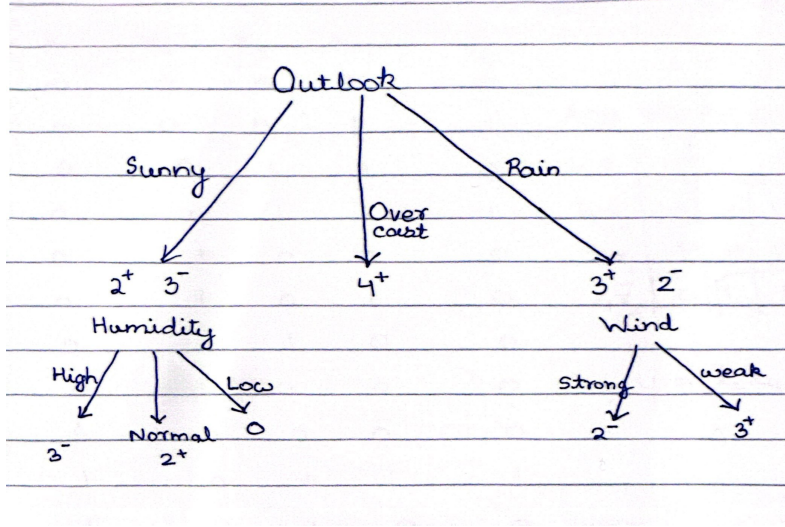


Figure 2: ME dependent based decision tree

Humidity and Outlook has highest gain, lets split on Outlook
 Overcast has already reached the leaf node, Play/S = Yes
 For other attribute, Sunny and Wind, we calculate Gain.

$$ME(\text{Sunny}) = \frac{2}{5}$$

Attributes are Humidity, Wind and Temperature.

$$Gain(\text{Sunny}, \text{Humidity}) = \frac{2}{5} - \frac{2}{5} * 0 = \frac{2}{5}$$

$$Gain(\text{Sunny}, \text{Temperature}) = \frac{2}{5} - \frac{2}{5} * \frac{1}{2} = \frac{1}{5}$$

$$Gain(\text{Sunny}, \text{Wind}) = \frac{2}{5} - \frac{2}{5} * \frac{1}{2} - \frac{3}{5} * \frac{1}{3} = 0$$

we split on Humidity, as it has the highest gain.

$$Gain(\text{Rain}, \text{Humidity}) = ME(\text{Rain}) - \frac{3}{5} * \frac{1}{3} - \frac{1}{2} * \frac{2}{5} = 0$$

$$Gain(\text{Rain}, \text{Wind}) = ME(\text{Rain}) - \frac{2}{5} * 0 = \frac{2}{5}$$

$$Gain(\text{Rain}, \text{Temperature}) = ME(\text{Rain}) - \frac{3}{5} * \frac{1}{3} - \frac{1}{2} * \frac{2}{5} = 0$$

we split on Wind, as it has the highest gain

The decision tree is shown in Fig. 2

- (b) [7 points] Please use gini index (GI) to calculate the gain, and conduct tree learning with ID3 framework. List every step and the tree structure.

Solution:

$$Gain(S, A) = GI(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} GI(S_v)$$

$$GI(S) = 1 - \left(\frac{5}{14}\right)^2 - \left(\frac{9}{14}\right)^2 = 0.46$$

Outlook:

$$\text{Sunny, } GI(\text{Sunny}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

$$\text{Overcast, } GI(\text{Overcast}) = 0$$

$$\text{Rain, } GI(\text{Rain}) = 1 - \left(\frac{2}{5}\right)^2 - \left(\frac{3}{5}\right)^2 = 0.48$$

Temperature:

$$\text{Hot, } GI(\text{Hot}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Mild, } GI(\text{Mild}) = 1 - \left(\frac{4}{6}\right)^2 - \left(\frac{2}{6}\right)^2 = 0.44$$

$$\text{Cool, } GI(\text{Cool}) = 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 = 0.375$$

Humidity:

$$\text{High, } GI(\text{High}) = 1 - \left(\frac{3}{7}\right)^2 - \left(\frac{4}{7}\right)^2 = 0.489$$

$$\text{Normal, } GI(\text{Normal}) = 1 - \left(\frac{6}{7}\right)^2 - \left(\frac{1}{7}\right)^2 = 0.245$$

$$\text{Low, } GI(\text{Low}) = 0$$

Wind:

$$\text{Strong, } GI(\text{Strong}) = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\text{Weak, } GI(\text{Weak}) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.375$$

$$Gain(S, Outlook) = 0.46 - 0.48 * \frac{5}{14} - 0 - 0.48 * \frac{5}{14} = 0.12$$

$$Gain(S, Temperature) = 0.46 - 0.5 * \frac{4}{14} - 0.44 * \frac{6}{14} - 0.375 * \frac{4}{14} = 0.019$$

$$Gain(S, Humidity) = 0.46 - 0.489 * \frac{7}{14} - 0.245 * \frac{7}{14} - 0 = 0.093$$

$$Gain(S, Wind) = 0.46 - 0.5 * \frac{4}{14} - 0.44 * \frac{6}{14} = 0.031$$

Start to split on Outlook as it has the highest gain.

Overcast has already reached the leaf node, Play= +

For other attribute, Sunny and Wind, we calculate Gain.

$$Gain(\text{Sunny, Humidity}) = GI(\text{Sunny}) - 0 = 0.48$$

$$Gain(\text{Sunny, Wind}) = 0.48 - \left[1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2\right] * \frac{2}{5} - \left[1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2\right] * \frac{3}{5} = 0.013$$

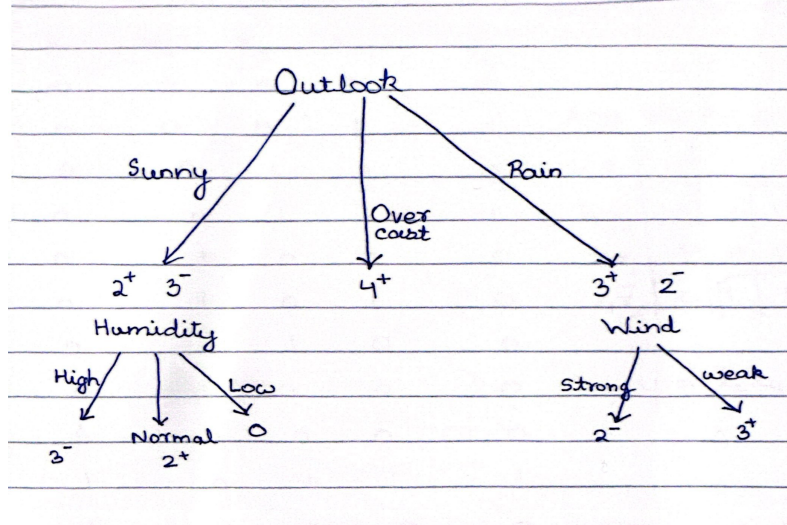


Figure 3: Gini index based decision tree

$$Gain(Sunny, Temperature) = 0.48 - \left[\left(1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2\right) * \frac{2}{5} \right] = 0.28$$

we split on Humidity as it has the highest gain

The High, Normal and Low has already reached the leaf node. We need to calculate for Rain.

$$Gain(Rain, Humidity) = GI(Rain) - \left[\left(1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2\right) * \frac{3}{5} \right] - \left[\left(1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2\right) * \frac{2}{5} \right] = 0.102$$

$$Gain(Rain, Wind) = GI(Rain) - 0 = 0.48$$

$$Gain(Rain, Temperature) = GI(Rain) - \left[\left(1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2\right) * \frac{3}{5} \right] - \left[\left(1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2\right) * \frac{2}{5} \right] = 0.102$$

we split on Wind as it has the highest gain.

For Strong and Weak, Wind has already reached the leaf node

The decision tree is shown in Fig. 3.

- (c) [3 points] Compare the two trees you just created with the one we built in the class (see Page 62 of the lecture slides). Are there any differences? Why?

Solution: No, there is no difference between the two decision trees. The splitting nodes, and leaf node labels are calculated by both the methods and results are same.

3. [16 points] Continue with the same training data in Problem 2. Suppose before the tree construction, we receive one more training instance where Outlook's value is missing: {Outlook: Missing, Temperature: Mild, Humidity: Normal, Wind: Weak, Play: Yes}.
- (a) [3 points] Use the most common value in the training data as the missing value, and calculate the information gains of the four features. Note that if there is a tie for the most common value, you can choose any value in the tie. Indicate the best feature.

Solution:

$$Gain(S, A) = ME(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} ME(S_v)$$

$$ME(S) = \frac{5}{15}$$

Taking missing value in outlook

Outlook:

$$\text{Sunny, } ME(\text{Sunny}) = \frac{3}{6}$$

$$\text{Overcast, } ME(\text{Overcast}) = 0$$

$$\text{Rain, } ME(\text{Rain}) = \frac{2}{5}$$

Temperature:

$$\text{Hot, } ME(\text{Hot}) = \frac{2}{4}$$

$$\text{Mild, } ME(\text{Mild}) = \frac{2}{7}$$

$$\text{Cool, } ME(\text{Cool}) = \frac{1}{4}$$

Humidity:

$$\text{High, } ME(\text{High}) = \frac{3}{7}$$

$$\text{Normal, } ME(\text{Normal}) = \frac{1}{8}$$

$$\text{Low, } ME(\text{Low}) = 0$$

Wind:

$$\text{Strong, } ME(\text{Strong}) = \frac{3}{6}$$

$$\text{Weak, } ME(\text{weak}) = \frac{2}{9}$$

$$Gain(S, Outlook) = \frac{5}{15} - \frac{3}{6} * \frac{6}{15} - 0 - \frac{2}{5} * \frac{5}{15} = 0$$

$$Gain(S, Temperature) = \frac{5}{15} - \frac{2}{4} * \frac{4}{15} - \frac{1}{4} * \frac{4}{15} - \frac{2}{7} * \frac{7}{15} = 0$$

$$Gain(S, Humidity) = \frac{5}{15} - \frac{3}{7} * \frac{7}{15} - \frac{1}{8} * \frac{8}{15} - 0 = 0.067$$

$$Gain(S, Wind) = \frac{5}{15} - \frac{3}{6} * \frac{6}{15} - \frac{2}{9} * \frac{9}{15} = 0$$

The best feature is Humidity, as it has the highest gain.

- (b) [3 points] Use the most common value among the training instances with the same label, namely, their attribute "Play" is "Yes", and calculate the information gains of the four features. Again if there is a tie, you can choose any value in the tie. Indicate the best feature.

Solution:

The new feature added here with the most common value with Play: Yes is:

Outlook: Overcast, Temperature: Mild, Humidity: Normal, Wind: Weak, Play: Yes

Outlook:

Sunny , ME (Sunny) = $\frac{2}{5}$

Overcast , ME (Overcast) = 0

Rain , ME (Rain) = $\frac{2}{5}$

Temperature:

Hot , ME (Hot) = $\frac{2}{4}$

Mild , ME (Mild) = $\frac{2}{7}$

Cool , ME (Cool) = $\frac{1}{4}$

Humidity:

High, ME (High) = $\frac{3}{7}$

Normal, ME (Normal) = $\frac{1}{8}$

Low, ME (Low) = 0

Wind:

Strong, ME (Strong) = $\frac{3}{6}$

weak , ME (weak) = $\frac{2}{9}$

$$Gain(S, Outlook) = \frac{2}{5} - \frac{2}{5} * \frac{5}{15} - 0 - \frac{2}{5} * \frac{5}{15} = 0.067$$

$$Gain(S, Temperature) = \frac{2}{5} - \frac{2}{4} * \frac{4}{15} - \frac{1}{4} * \frac{4}{15} - \frac{2}{7} * \frac{7}{15} = 0$$

$$Gain(S, Humidity) = \frac{2}{5} - \frac{3}{7} * \frac{7}{15} - \frac{1}{8} * \frac{8}{15} - 0 = 0.067$$

$$Gain(S, Wind) = \frac{2}{5} - \frac{3}{6} * \frac{6}{15} - \frac{2}{9} * \frac{9}{15} = 0$$

There's a tie on the highest gain and we may split on Outlook.

- (c) [3 points] Use the fractional counts to infer the feature values, and then calculate the information gains of the four features. Indicate the best feature.

Solution:

ME(S) = $\frac{5}{15}$

Using fractional counts of the attribute values in training data:

Outlook:

Size: Sunny : $5 + \frac{5}{14}$ Overcast : $4 + \frac{4}{14}$ Rain : $5 + \frac{5}{14}$

Sunny : $p = 2 + \frac{5}{14}$, $n = 3$, ME (Sunny) = $\frac{2 + \frac{5}{14}}{5 + \frac{5}{14}}$
 Overcast : $p = 4 + \frac{4}{14}$, $n = 0$, ME (Overcast) = 0
 Rain : $p = 3 + \frac{5}{14}$, $n = 2$, ME (Rain) = $\frac{2}{5 + \frac{5}{14}}$

Temperature:

Hot , ME (Hot) = $\frac{2}{4}$
 Mild , ME (Mild) = $\frac{2}{7}$
 Cool , ME (Cool) = $\frac{1}{4}$

Humidity:

High, ME (High) = $\frac{3}{7}$
 Normal , ME (Normal) = $\frac{1}{8}$
 Low, ME (Low) = 0

Wind:

Strong , ME (Strong) = $\frac{3}{6}$
 weak , ME (weak) = $\frac{2}{9}$

$$Gain(S, Outlook) = \frac{5}{15} - \frac{2 + \frac{5}{14}}{5 + \frac{5}{14}} * \frac{5 + \frac{5}{14}}{15} - 0 - \frac{2}{5 + \frac{5}{14}} * \frac{5 + \frac{5}{14}}{15} = 0.0428$$

$$Gain(S, Temperature) = \frac{5}{15} - \frac{2}{4} * \frac{4}{15} - \frac{1}{4} * \frac{4}{15} - \frac{2}{7} * \frac{7}{15} = 0$$

$$Gain(S, Humidity) = \frac{5}{15} - \frac{3}{7} * \frac{7}{15} - \frac{1}{8} * \frac{8}{15} - 0 = 0.067$$

$$Gain(S, Wind) = \frac{5}{15} - \frac{3}{6} * \frac{6}{15} - \frac{2}{9} * \frac{9}{15} = 0$$

We split on Humidity as it has the highest gain.

- (d) [7 points] Continue with the fractional examples, and build the whole tree with information gain. List every step and the final tree structure.

Solution:

Split on Humidity developing decision tree.

Find the attribute that best Splits S again for High.

$$ME(High) = \frac{3}{7}$$

Current attributes: A = Outlook, Temperature, Wind

$$Gain(High, Outlook) = \frac{3}{7} - \frac{3}{7} * 0 - \frac{2}{7} * 0 - \frac{2}{7} * \frac{1}{2} = 0.285$$

$$Gain(High, Temperature) = \frac{3}{7} - \frac{3}{7} * \frac{1}{3} - \frac{4}{7} * \frac{2}{4} - 0 = 0$$

$$Gain(High, Wind) = \frac{3}{7} - \frac{3}{7} * \frac{1}{3} - \frac{4}{7} * \frac{2}{4} = 0$$

Split on High Humidity as it has the highest gain, with attribute of Outlook.
Starting from Outlook, Sunny and Overcast reached the same label, a leaf node with label No, Yes respectively.

Find the attribute that best Splits for Rain.

$ME(Rain) = \frac{1}{2} = 0.5$ Current attributes: A = Temperature, Wind

$$Gain(Rain, Temperature) = ME(Rain) - 0 - \frac{2}{2} * \frac{1}{2} - 0 = 0$$

$$Gain(Rain, Wind) = ME(Rain) - 0 - 0 * \frac{1}{2} - 0 * \frac{1}{2} = 0.5$$

we will split on Wind as it has the highest gain with attribute of Rain.
Under Wind, Strong and Weak have the same label, a leaf node with label No and Yes labelled respectively.

Now , back to Humidity with Normal for further split.

$ME(Normal) = \frac{1}{8} = 0.125$

Current attributes: A = Outlook, Temperature, Wind

$$Gain(Normal, Outlook) = ME(Normal) - 0 - 0 - \frac{3.357}{8} * \frac{1}{3.357} = 0$$

$$Gain(Normal, Temperature) = ME(Normal) - 0 - 0 - \frac{4}{8} * \frac{1}{4} = 0$$

$$Gain(Normal, Wind) = ME(Normal) - \frac{3}{8} * \frac{1}{3} = 0$$

None of the attributes produce any information gain. Lets split on Outlook
Sunny and Overcast have the same label, a leaf node with Yes and Yes labelled respectively.

Let's look for the attribute that best Splits on Rain.

$ME(Rain) = \frac{1}{3.357} = 0.297$

Current attributes: A = Temperature, Wind

$$Gain(Rain, Temperature) = ME(Rain) - 0 - 0 - \frac{2}{3.357} * \frac{1}{2} = 0$$

$$Gain(Rain, Wind) = ME(Rain) - 0 = 0.297$$

Wind has the highest information gain, split Rain with attribute of Wind
Strong and Weak have the same label, a leaf node with No and Yes labelled respectively.

Low Humidity is null, a leaf node with the most common label is Yes.

The decision tree is shown in Fig. 4

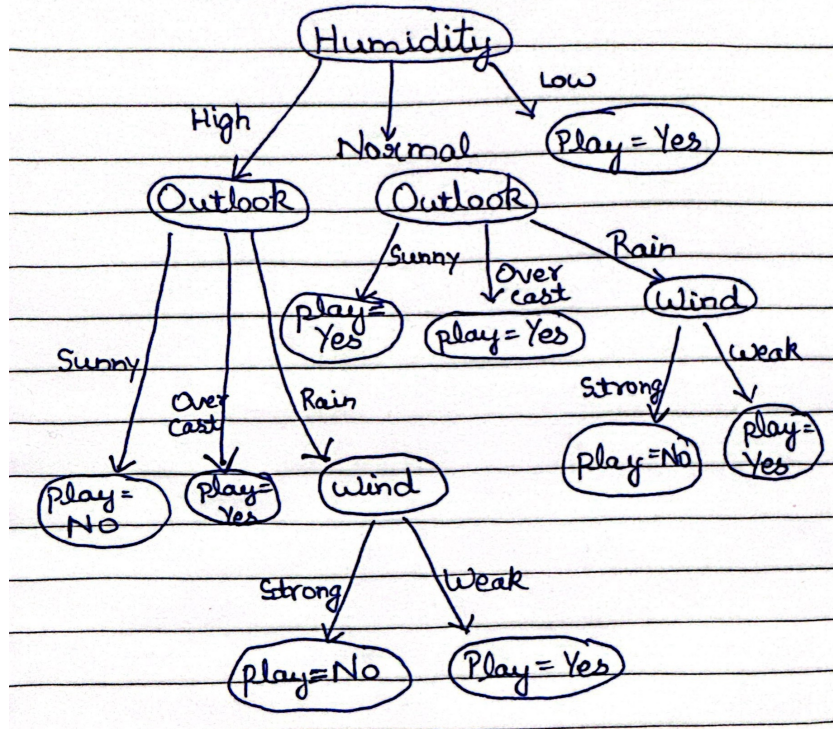


Figure 4: Final decision tree

[**Bonus question 1**] [5 points]. Prove that the information gain is always non-negative. That means, as long as we split the data, the purity will never get worse! (Hint: use convexity)

Solution:

Information gain is given as:

$$IG = I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) - \sum_{i=0}^1 \left[\frac{p_i + n_i}{p+n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right) \right]$$

Let

$$f(x) = -x \log_2 x - (1-x) \log_2 (1-x)$$

$$f'(x) = -\log_2 x + \log_2 (1-x)$$

and

$$f''(x) = -\frac{1}{\ln 2} * \frac{1}{x(1-x)}$$

Since $x \in (0, 1)$ we have $f''(x) < 0$ and it's concave.

$$\begin{aligned} & \sum_{i=0}^1 \frac{p_i + n_i}{p+n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right) \\ &= \frac{p_0 + n_0}{p+n} I\left(\frac{p_0}{p_0 + n_0}, \frac{n_0}{p_0 + n_0}\right) + \frac{p_1 + n_1}{p+n} I\left(\frac{p_1}{p_1 + n_1}, \frac{n_1}{p_1 + n_1}\right) \end{aligned}$$

$$= \frac{p_0 + n_0}{p + n} f\left(\frac{p_0}{p_0 + n_0}\right) + \frac{p_1 + n_1}{p + n} f\left(\frac{p_1}{p_1 + n_1}\right)$$

Following Jensen's Inequality:

$$\sum_x p(x) f(x) \leq f\left(\sum_x p(x) x\right)$$

where $\sum_x p(x) = 1$, $p(x) \geq 0$ and $f(x)$ is concave.

So,

$$\begin{aligned} \sum_{i=0}^1 \frac{p_i + n_i}{p + n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right) &\leq f\left(\frac{p_0 + n_0}{p + n} \frac{p_0}{p_0 + n_0}\right) + f\left(\frac{p_1 + n_1}{p + n} \frac{p_1}{p_1 + n_1}\right) \\ &= f\left(\frac{p_0 + p_1}{p + n}\right) = f\left(\frac{p}{p + n}\right) = I\left(\frac{p}{p + n}, \frac{n}{p + n}\right) \end{aligned}$$

So,

$$IG = I\left(\frac{p}{p + n}, \frac{n}{p + n}\right) - \sum_{i=0}^1 \left[\frac{p_i + n_i}{p + n} I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right) \right] \geq 0 \dots$$

[Bonus question 2] [5 points]. We have discussed how to use decision tree for regression (i.e., predict numerical values) — on the leaf node, we simply use the average of the (numerical) labels as the prediction. Now, to construct a regression tree, can you invent a gain to select the best attribute to split data in ID3 framework?

Solution:

The decision tree uses Entropy and Information Gain. We may focus on how much our predictions deviate from the actual values, typically measured using mean square error (MSE).

2 Decision Tree Practice [60 points]

1. [5 Points] Starting from this assignment, we will build a light-weighted machine learning library. To this end, you will first need to create a code repository in Github.com. Please refer to the short introduction in the appendix and the official tutorial to create an account and repository. Please commit a README.md file in your repository, and write one sentence: "This is a machine learning library developed by **Your Name** for CS5350/6350 in University of Utah". You can now create a first folder, "DecisionTree". Please leave the link to your repository in the homework submission. We will check if you have successfully created it.

Solution:

https://github.com/pratishtha-maker/Machine_Learning.git

2. [30 points] We will implement a decision tree learning algorithm for car evaluation task. The dataset is from UCI repository(<https://archive.ics.uci.edu/ml/datasets/car+evaluation>). Please download the processed dataset (car.zip) from Canvas. In this task, we have 6 car attributes, and the label is the evaluation of the car. The attribute and label values are listed in the file “data-desc.txt”. All the attributes are categorical. The training data are stored in the file “train.csv”, consisting of 1,000 examples. The test data are stored in “test.csv”, and comprise 728 examples. In both training and test datasets, attribute values are separated by commas; the file “data-desc.txt” lists the attribute names in each column.

Note: we highly recommend you to use Python for implementation, because it is very convenient to load the data and handle strings. For example, the following snippet reads the CSV file line by line and split the values of the attributes and the label into a list, “terms”. You can also use “dictionary” to store the categorical attribute values. In the web are numerous tutorials and examples for Python. if you have issues, just google it!

```
with open(CSVfile, 'r') as f:
    for line in f:
        terms = line.strip().split(',')
        process one training example
```

- (a) [15 points] Implement the ID3 algorithm that supports, information gain, majority error and gini index to select attributes for data splits. Besides, your ID3 should allow users to set the maximum tree depth. Note: you do not need to convert categorical attributes into binary ones and your tree can be wide here.

Solution:

The ID3 algorithm is uploaded in github.

- (b) [10 points] Use your implemented algorithm to learn decision trees from the training data. Vary the maximum tree depth from 1 to 6 — for each setting, run your algorithm to learn a decision tree, and use the tree to predict both the training and test examples. Note that if your tree cannot grow up to 6 levels, you can stop at the maximum level. Report in a table the average prediction errors on each dataset when you use information gain, majority error and gini index heuristics, respectively.

Solution:

The prediction errors are tabulated in Table 3.

Depth	IG_{train}	IG_{test}	ME_{train}	ME_{test}	GI_{train}	GI_{test}
1	0.3013013	0.29711141	0.3013013	0.29711141	0.3013013	0.29711141
2	0.3013013	0.29711141	0.3013013	0.29711141	0.3013013	0.29711141
3	0.22222222	0.22283356	0.22222222	0.22283356	0.22222222	0.22283356
4	0.18118118	0.19669876	0.17417417	0.18707015	0.17617617	0.18431911
5	0.08208208	0.15130674	0.08908908	0.13755158	0.08908908	0.13755158
6	0.02702702	0.09491059	0.02702702	0.09491059	0.02702702	0.09491059

Table 3: Average prediction error

- (c) [5 points] What can you conclude by comparing the training errors and the test errors?

Solution:

We can see that training error is decreasing with the depth level. On the other hand, the test data set are showing best result till certain depth level and from the table we can decide the best depth level for test data.

3. [25 points] Next, modify your implementation a little bit to support numerical attributes. We will use a simple approach to convert a numerical feature to a binary one. We choose the media (NOT the average) of the attribute values (in the training set) as the threshold, and examine if the feature is bigger (or less) than the threshold. We will use another real dataset from UCI repository(<https://archive.ics.uci.edu/ml/datasets/Bank+Marketing>). This dataset contains 16 attributes, including both numerical and categorical ones. Please download the processed dataset from Canvas (bank.zip). The attribute and label values are listed in the file “data-desc.txt”. The training set is the file “train.csv”, consisting of 5,000 examples, and the test “test.csv” with 5,000 examples as well. In both training and test datasets, attribute values are separated by commas; the file “data-desc.txt” lists the attribute names in each column.
- (a) [10 points] Let us consider “unknown” as a particular attribute value, and hence we do not have any missing attributes for both training and test. Vary the maximum tree depth from 1 to 16 — for each setting, run your algorithm to learn a decision tree, and use the tree to predict both the training and test examples. Again, if your tree cannot grow up to 16 levels, stop at the maximum level. Report in a table the average prediction errors on each dataset when you use information gain, majority error and gini index heuristics, respectively.

Solution:

The average predicted error is tabulated in Table 4.

Depth	IG_{train}	IG_{test}	ME_{train}	ME_{test}	GI_{train}	GI_{test}
1	0.11922384	0.12482496	0.11922384	0.12482496	0.11922384	0.1248249
2	0.11922384	0.12482496	0.10882176	0.11662332	0.10882176	0.1166233
3	0.10602120	0.11142228	0.10422084	0.10882176	0.10422084	0.1088217
4	0.10062012	0.10702140	0.09401880	0.11442288	0.09361872	0.11222244
5	0.08001600	0.11402280	0.07841568	0.11862372	0.07541508	0.11922384
6	0.06241248	0.11882376	0.06521304	0.12482496	0.06041208	0.12602520
7	0.04800960	0.12802560	0.05681136	0.12822564	0.04840968	0.13682736
8	0.03720744	0.13302660	0.05181036	0.13062612	0.03640728	0.14522904
9	0.02920584	0.13842768	0.04800960	0.13182636	0.02680536	0.15083016
10	0.02220444	0.14262852	0.04340868	0.13422684	0.02140428	0.15283056
11	0.01820364	0.14702940	0.03660732	0.13882776	0.01720344	0.15563112
12	0.01500300	0.14702940	0.02680536	0.14722944	0.01400280	0.15863172
13	0.01360272	0.14942988	0.02702702	0.14922984	0.01360272	0.15883176
14	0.01320264	0.150030006	0.02240448	0.15283056	0.01320264	0.15943188
15	0.01320264	0.15003000	0.01780356	0.15823164	0.01320264	0.15943188
16	0.01320264	0.15003000	0.01520304	0.16043208	0.01320264	0.15943188

Table 4: Average prediction error

- (b) [10 points] Let us consider "unknown" as attribute value missing. Here we simply complete it with the majority of other values of the same attribute in the training set. Vary the maximum tree depth from 1 to 16 — for each setting, run your algorithm to learn a decision tree, and use the tree to predict both the training and test examples. Report in a table the average prediction errors on each dataset when you use information gain, majority error and gini index heuristics, respectively.

Solution:

The average predicted error is tabulated in Table 5.

Depth	IG_{train}	IG_{test}	ME_{train}	ME_{test}	GI_{train}	GI_{test}
1	0.11922384	0.12482496	0.11922384	0.12482496	0.11922384	0.1248249
2	0.11922384	0.12482496	0.10882176	0.11662332	0.10882176	0.1166233
3	0.10602120	0.11142228	0.10422084	0.10882176	0.10422084	0.1088217
4	0.10062012	0.10702140	0.09401880	0.11442288	0.09361872	0.11222244
5	0.08001600	0.11402280	0.07841568	0.11862372	0.07541508	0.11922384
6	0.06241248	0.11882376	0.06521304	0.12482496	0.06041208	0.12602520
7	0.04800960	0.12802560	0.05681136	0.12822564	0.04840968	0.13682736
8	0.03720744	0.13302660	0.05181036	0.13062612	0.03640728	0.14522904
9	0.02920584	0.13842768	0.04800960	0.13182636	0.02680536	0.15083016
10	0.02220444	0.05302301	0.0284 2451	0.16162341	0.13622385	0.16081543
11	0.01820364	0.14702940	0.03660732	0.13882776	0.01720344	0.15563112
12	0.01500300	0.04883110	0.02847210	0.16161121	0.13884539	0.16160120
13	0.01360272	0.04401100	0.02829933	0.16282300	0.14421122	0.16282130
14	0.01320264	0.03732110	0.02821120	0.16282351	0.15421122	0.16285640
15	0.01320264	0.03161199	0.02820111	0.16281122	0.15781199	0.16282211
16	0.01320264	0.02822233	0.02821120	0.162800112	0.16361120	0.16281121

Table 5: Average prediction error

- (c) [5 points] What can you conclude by comparing the training errors and the test errors, with different tree depths, as well as different ways to deal with "unknown" attribute values?

Solution:

The gain increases when unknown values are assigned to the majority label. Recalculating errors with increased depth can be stopped in initial stages to save computation time. The maximum gain should be preferred over the uncertainty. The majority value can be updated for the unknown data and train over multiple depths.