

CSCN8000 – Artificial Intelligence Algorithms and Mathematics

Assignment 2: Decision Trees

Student ID: 8883828

Name: Parthasarathy Rajendiran

Dataset

The given dataset contains 6 observations, three predictors (Good Behaviour, Age<30, and Drug Dependent), and a target (RECIDIVIST)

ID	Good Behaviour	Age<30	Drug Dependent	RECIDIVIST
1	FALSE	TRUE	FALSE	TRUE
2	FALSE	FALSE	FALSE	FALSE
3	FALSE	TRUE	FALSE	TRUE
4	TRUE	FALSE	FALSE	FALSE
5	TRUE	FALSE	TRUE	TRUE
6	TRUE	FALSE	FALSE	FALSE

Part a.

Firstly, the entropy of the target variable should be calculated.

Entropy of the target

The total number of records = 6

Probability of Recidivist

$$P(\text{RECIDIVIST}=\text{TRUE}) = \frac{3}{6} = 0.5$$

Probability of Non-Recidivist

$$P(\text{RECIDIVIST}=\text{FALSE}) = \frac{3}{6} = 0.5$$

The entropy of the RECIDIVIST column

$$\begin{aligned} H(\text{RECIDIVIST}) &= P(\text{True}) \cdot H(\text{True}) + P(\text{False}) \cdot H(\text{False}) \\ &= (0.5 \times -\log_2 0.5) \\ &\quad + (0.5 \times -\log_2 0.5) \\ &= (0.5 \times -(-1)) + (0.5 \times -(-1)) \\ &= 0.5 + 0.5 \\ &= 1 \end{aligned}$$

$H(S) = 1$

With the use of this entropy, the Information Gain for each column should be calculated.

Information Gain for Good Behaviour column

$$\begin{aligned}
 IG(GB) &= H(S) - \left[\frac{P(GB=True)}{P(GB=False)} H(S|GB=True) + \right. \\
 &= 1 - \left[\left(\frac{3}{6} \times \left(-\frac{1}{3} \log \frac{1}{3} \right) + \left(-\frac{2}{3} \log \frac{2}{3} \right) \right) + \right. \\
 &= 1 - \left[\left(\frac{3}{6} \times \left(-\frac{2}{3} \log \frac{2}{3} \right) + \left(-\frac{1}{3} \log \frac{1}{3} \right) \right) \right] \\
 &= 1 - \left[\left(\frac{3}{6} \times (0.5283 + 0.39) \right) + \left(\frac{3}{6} \times (0.39 + 0.5283) \right) \right] \\
 &= 1 - \left[\left(\frac{1}{2} \times 0.9183 \right) + \left(\frac{1}{2} \times 0.9183 \right) \right] \\
 &= 1 - 0.9183 \\
 IG(GB) &= 0.0817
 \end{aligned}$$

Information Gain for Age < 30 column

$$\begin{aligned}
 IG(Age) &= H(S) - \left[\frac{P(Age=True)}{P(Age=False)} H(S|Age=True) + \right. \\
 &= 1 - \left[\left(\frac{2}{6} \times \left(-\frac{2}{2} \log \frac{2}{2} \right) + \left(-\frac{0}{2} \log \frac{0}{2} \right) \right) + \right. \\
 &= 1 - \left[0 + \frac{4}{6} (0.5 + 0.3113) \right] \\
 &= 1 - 0.5409 \\
 IG(Age) &= 0.4591
 \end{aligned}$$

Information Gain for Drug Dependent column

$$\begin{aligned}
 IG(DD) &= H(S) - \left[\frac{P(DD=True)}{P(DD=False)} H(S|DD=True) + \right. \\
 &= 1 - \left[\frac{1}{6} \left(-\frac{1}{1} \log \frac{1}{1} \right) + \left(-0 \right) \right] + \frac{5}{6} \left(-\frac{2}{5} \log \frac{2}{5} + \left(-\frac{3}{5} \log \frac{3}{5} \right) \right) \\
 &= 1 - \left[0 + \frac{5}{6} (0.5288 + 0.4422) \right] \\
 &= 1 - 0.8092 \\
 IG(DD) &= 0.1908
 \end{aligned}$$

Root Node Selection

The information gain from the "Age<30" column is higher than that of other columns. So, Age<30 can be used as the root node.

Now, let us split the dataset using the "Age<30" column.

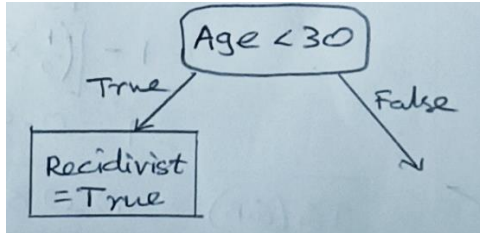
Age<30 = TRUE

Count 2							
ID	Good Behaviour	Age<30	Drug Dependent	RECIDIVIST	Good B...	Age<30	Drug D...
1	FALSE	TRUE	FALSE	TRUE	FALSE	FALSE	FALSE
3	FALSE	TRUE	FALSE	TRUE	TRUE	TRUE	TRUE

Age < 30 = FALSE

Count 4										
ID	Good Behaviour	Age < 30	Drug Dependent	RECIDIVIST	Good B...		Age < 30		Drug D...	
2	FALSE	FALSE	FALSE	FALSE	FALSE		FALSE		FALSE	
4	TRUE	FALSE	FALSE	FALSE	TRUE		TRUE		TRUE	
5	TRUE	FALSE	TRUE	TRUE						
6	TRUE	FALSE	FALSE	FALSE						

For the Age < 30 = TRUE subset, all the records have Recidivist = TRUE. Since this subset contains only one class, this branch reached the leaf node.



The above steps must be repeated for the subset Age < 30 = FALSE

Splitting the subset Age < 30 = False

There are 4 prisoners have age more than 30. One of them is Recidivist and others are not.

Entropy of the Target for Age < 30 = False

$$\begin{aligned}
 H(\text{Recidivist} | \text{Age} = \text{False}) &= P(T) \cdot H(T) + P(F) \cdot H(F) \\
 &= \left[-\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \cdot -\log_2 \frac{3}{4} \right] \\
 &= 0.5 + 0.3113 \\
 \boxed{H(S) = 0.8113}
 \end{aligned}$$

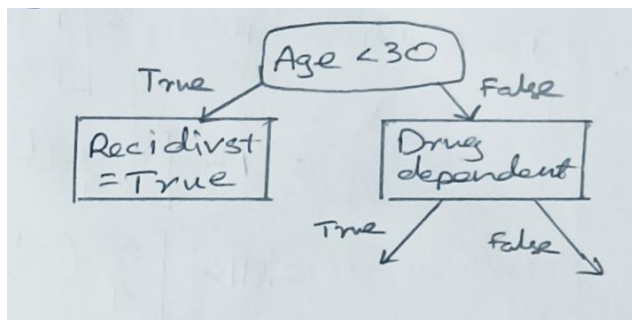
Information Gain of Good Behaviour column.

$$\begin{aligned}
 IG(GB) &= H(S) - \left[\frac{P(GB = \text{True})}{P(S)} \cdot H(S | GB = \text{True}) + \frac{P(GB = \text{False})}{P(S)} \cdot H(S | GB = \text{False}) \right] \\
 &= 0.8113 - \left[\frac{3}{4} \cdot \left(-\frac{1}{3} \log_2 \frac{1}{3} \right) + \left(-\frac{2}{3} \log_2 \frac{2}{3} \right) + \frac{1}{4} \cdot (0 - 0) \right] \\
 &= 0.8113 - \frac{3}{4} (0.5283 + 0.39) \\
 &= 0.8113 - \left(\frac{3}{4} \times 0.9183 \right) \\
 &= 0.8113 - 0.6887 \\
 \boxed{IG(GB) = 0.1226}
 \end{aligned}$$

Information Gain for Drug Dependent column

$$\begin{aligned}
 IG(DD) &= H(S) - \left[P(DD=True) \cdot H(S|DD=True) + P(DD=False) \cdot H(S|DD=False) \right] \\
 &= 0.8113 - \left[\frac{1}{4} \times \left(-\frac{1}{4} \log_2 \frac{1}{4} \right) + (-0 \cdot \log_2 0) \right] \\
 &= 0.8113 - \left[\frac{3}{4} \times \left(-0 \cdot \log_2 0 \right) + \left(-\frac{3}{4} \log_2 \frac{3}{4} \right) \right] \\
 &= 0.8113 - 0 - 0 \\
 \boxed{IG(DD) = 0.8113}
 \end{aligned}$$

Since the information gain of the “Drug Dependent” column is higher, that can be considered for the next split.

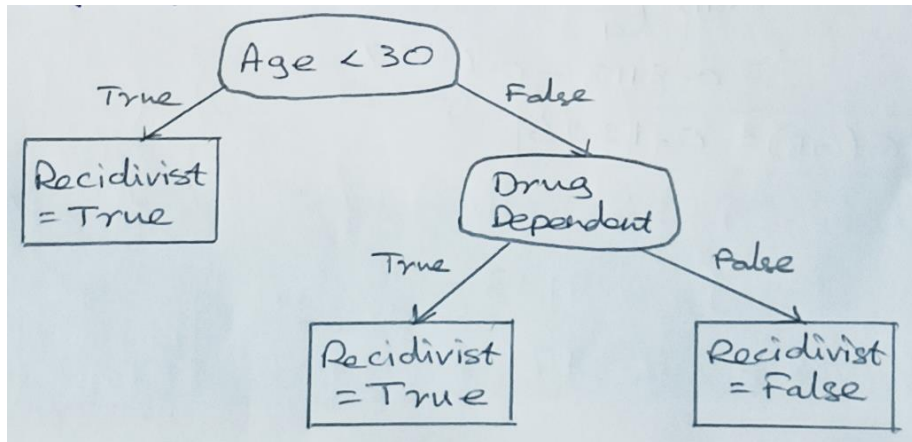


Out of 4 prisoners with Age < 30, only one is Drug dependent and also he is a recidivist. All the other records are non-recidivist.

Count 1							
ID	Good Behaviour	Age<30	Drug Dependent	RECIDIVIST	Good B...	Age<30	Drug D...
5	TRUE	FALSE	TRUE	TRUE	TRUE	FALSE	FALSE
					FALSE	TRUE	TRUE

Count 3							
ID	Good Behaviour	Age<30	Drug Dependent	RECIDIVIST	Good B...	Age<30	Drug D...
2	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE	FALSE
4	TRUE	FALSE	FALSE	FALSE	TRUE	TRUE	TRUE
6	TRUE	FALSE	FALSE	FALSE			

Since each subset of Drug Dependents belongs to the same class, we can say that we reached leaf nodes in all the branches. So, the final decision tree is.



This decision tree is optimal and gives below insights,

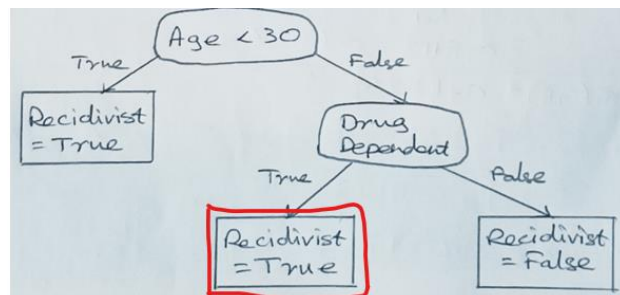
- All the prisoners released on parole and youngsters (Age < 30) are RECIDIVIST.
- All the Drug Dependents are Recidivist = TRUE.
- The prisoners who are not Drug Dependents and older than 30 years are not RECIDIVIST.

Part b.

What prediction will the decision tree generate in part (a) of this question return for the following query? (5 Points)

GOOD BEHAVIOR = false, AGE < 30 = false,
DRUG DEPENDENT = true

Based on our decision tree, a prisoner with an age of more than 30 and a drug-dependent will be a Recidivist.



Part c.

What prediction will the decision tree generate in part (a) of this question return for the following query? (5 Points)

GOOD BEHAVIOR = true, AGE < 30 = true,
DRUG DEPENDENT = false

Based on our decision tree, a prisoner with an age of less than 30 and not drug-dependent will be a Recidivist.

