

Lab 9 : Decision tree from scratch

1. Consider the following instances, S:

Weekend (Example)	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay in
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

Let's create a decision tree from scratch without a programming tool.

2. First we need to work out which attribute will be the root node for the decision tree.

Given an arbitrary categorisation, C into categories c_1, \dots, c_n , and a set of examples, S, for which the proportion of examples in c_i is p_i , then the entropy of S is:

$$Entropy(t) = -\sum_j p(j|t) \log_2 p(j|t)$$

P_x = probability for x to occur.

$$\begin{aligned} Entropy(Decision) &= -P_{Cinema} * \log_2(P_{Cinema}) - P_{Tennis} * \log_2(P_{Tennis}) - P_{Stay-in} * \log_2(P_{Stay-in}) - P_{Shopping} * \log_2(P_{Shopping}) \\ &= -[(0.6) * \log_2(0.6) + (0.2) * \log_2(0.2) + (0.1) * \log_2(0.1) + (0.1) * \log_2(0.1)] \\ &= \mathbf{1.57095} \end{aligned}$$

3. Next we need to calculate the best of the following formula:

$$GAIN_{split} = Entropy(p) - \left(\sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

Note:

- $Entropy(Decision, Weather)$, $k \in \{Sunny, Windy, Rainy\}$

- $Entropy(Decision, Weather) = P_{Sunny} * Entropy(Sunny) + P_{Windy} * Entropy(Windy) + P_{Rainy} * Entropy(Rainy)$
 $= (3/10) * Entropy(Sunny) + (4/10) * Entropy(Windy) + (3/10) * Entropy(Rainy)$

Weather	Decision
Sunny	Cinema
Sunny	Tennis
Sunny	Tennis

$$\begin{aligned} Entropy(Sunny) &= -[P_{Cinema} * \log_2(P_{Cinema}) + P_{Tennis} * \log_2(P_{Tennis}) + P_{stay-in} * \log_2(P_{stay-in}) + P_{Shopping} * \log_2(P_{Shopping})] \\ &= -[(1/3) * \log_2(1/3) + (2/3) * \log_2(2/3) + 0 + 0] \\ &= 0.918 \end{aligned}$$

$E(Cinema, Tennis, Stay-in, Shopping)$ will be used as abbreviation in this section.

Gain(Decision, Weather)

$= Entropy(Decision) - Entropy(Weather)$

$= 1.57095 - [P_{Sunny} * Entropy(Sunny) + P_{Windy} * Entropy(Windy) + P_{Rainy} * Entropy(Rainy)]$

$= 1.57095 - [(3/10) * E(1, 2, 0, 0) + (4/10) * E(3, 0, 1, 0) + (3/10) * E(2, 0, 0, 1)]$

$= 1.57095 - [(0.3) * -((1/3) * \log_2(1/3) + (2/3) * \log_2(2/3)) + (0.4) * -((3/4) * \log_2(3/4) + (1/4) * \log_2(1/4)) + (0.3) * -((2/3) * \log_2(2/3) + (1/3) * \log_2(1/3))]$

= 0.6955

Gain(Decision, Parents)

$= 1.57095 - [(5/10) * E(5, 0, 0, 0) + (5/10) * E(1, 2, 1, 1)]$

$= 1.57095 - [0 + (0.5) * -((1/5) * \log_2(1/5) + (2/5) * \log_2(2/5) + (1/5) * \log_2(1/5) + (1/5) * \log_2(1/5))]$

= 0.61

Gain(Decision, Money)

$= 1.57095 - [(7/10) * E(3, 2, 1, 1) + (3/10) * E(3, 0, 0, 0)]$

$= 1.57095 - [(0.7) * -((3/7) * \log_2(3/7) + (2/7) * \log_2(2/7) + (1/7) * \log_2(1/7) + (1/7) * \log_2(1/7)) + 0]$

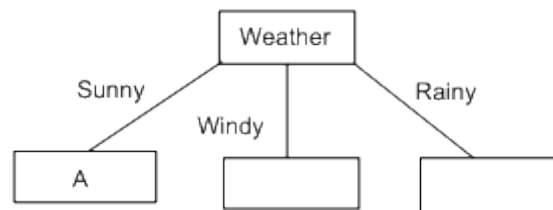
= 0.2813

The feature with the highest gain will be the root node, Weather (0.6955).

4. Draw your resulting decision tree.

5. Once you have decided on the tree, look at the activities based on weeks. For instance, assume you have chosen Weather, henceforth we investigate Sunny. For Sunny weather, the instances are W1,W2, W10 for which the activities are Cinema, Cinema, Tennis. As Cinema, Cinema and Tennis is present as a set and consists of non-Empty set, we will establish an attribute node, A. Do the same for the rest of your tree.

You may end up with something like the following



6. Now we have to fill in the choice of attribute A, which we know cannot be weather, because we've already removed that from the list of attributes to use. So, we need to calculate the values for **Gain(S_{Sunny} , parents)** and **Gain(S_{Sunny} , money)**. Firstly, Entropy(S_{Sunny}) = 0.918. Next, we set S to be $S_{\text{Sunny}} = \{W1, W2, W10\}$ (and, for this part of the branch, we will ignore all the other examples). In effect, we are interested only in this part of the table:

Weekend (Example)	Weather	Parents	Money	Decision (Category)
W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W10	Sunny	No	Rich	Tennis

Calculate the following

Gain(S_{Sunny} , parents) and Gain(S_{Sunny} , money). Which has the higher gain, and what impact does that have on the selection of attribute A?

$$\text{Entropy}(\text{Sunny}) = 0.918$$

Gain(Sunny, Parents)

$$\begin{aligned}
 &= \text{Entropy}(\text{Sunny}) - \text{Entropy}(\text{Sunny}, \text{Parents}) \\
 &= 0.918 - [P_{\text{Yes}} * \text{Entropy}(\text{Yes}) + P_{\text{No}} * \text{Entropy}(\text{No})] \\
 &= 0.918 - [(1/3) * E(1, 0, 0, 0) + (2/3) * E(0, 2, 0, 0)] \\
 &= 0.918 - [(1/3) * 0 + (2/3) * 0] \\
 &= \mathbf{0.918}
 \end{aligned}$$

Gain(Sunny, Money)

$$\begin{aligned}
 &= \text{Entropy}(\text{Sunny}) - \text{Entropy}(\text{Sunny}, \text{Money}) \\
 &= 0.918 - [P_{\text{Rich}} * \text{Entropy}(\text{Rich}) + P_{\text{Poor}} * \text{Entropy}(\text{Poor})] \\
 &= 0.918 - [(3/3) * E(1, 2, 0, 0) + 0] \\
 &= 0.918 - [(1) * -((1/3) * \log_2(1/3) + (2/3) * \log_2(2/3) + 0 + 0) + 0] \\
 &= \mathbf{0 \dots \text{No difference}}
 \end{aligned}$$

Therefore, Parents as the branching node.

Tree, Class = Decision. Rounded off to 2 decimal places.

1. Root Node.

Entropy(Decision)

$$\begin{aligned} &= -P_{\text{Cinema}} \cdot \log_2(P_{\text{Cinema}}) - P_{\text{Tennis}} \cdot \log_2(P_{\text{Tennis}}) - P_{\text{Stay-in}} \cdot \log_2(P_{\text{Stay-in}}) - P_{\text{Shopping}} \cdot \log_2(P_{\text{Shopping}}) \\ &= -[(0.6) \cdot \log_2(0.6) + (0.2) \cdot \log_2(0.2) + (0.1) \cdot \log_2(0.1) + (0.1) \cdot \log_2(0.1)] \\ &= \mathbf{1.57} \end{aligned}$$

Entropy(Weather_{Sunny})

$$\begin{aligned} &= -((1/3) \cdot \log_2(1/3) + (2/3) \cdot \log_2(2/3) + 0 + 0) \\ &= \mathbf{0.92} \end{aligned}$$

Entropy(Weather_{Windy})

$$\begin{aligned} &= -((3/4) \cdot \log_2(3/4) + (1/4) \cdot \log_2(1/4) + 0 + 0) \\ &= \mathbf{0.81} \end{aligned}$$

Entropy(Weather_{Rainy})

$$\begin{aligned} &= -((2/3) \cdot \log_2(2/3) + 0 + (1/3) \cdot \log_2(1/3) + 0) \\ &= \mathbf{0.92} \end{aligned}$$

Entropy(Decision, Weather)

$$\begin{aligned} &= (3/10)(0.92) + (4/10)(0.81) + (3/10)(0.92) \\ &= \mathbf{0.88} \end{aligned}$$

Gain(Decision, Weather)

$$\begin{aligned} &= 1.57 - 0.88 \\ &= \mathbf{0.69} \end{aligned}$$

Entropy(Parents_{Yes})

$$\begin{aligned} &= -((5/5) \cdot \log_2(5/5) + 0 + 0 + 0) \\ &= \mathbf{0} \end{aligned}$$

Entropy(Parents_{No})

$$\begin{aligned} &= -((1/5) \cdot \log_2(1/5) + (2/5) \cdot \log_2(2/5) + (1/5) \cdot \log_2(1/5) + (1/5) \cdot \log_2(1/5)) \\ &= \mathbf{1.92} \end{aligned}$$

Entropy(Decision, Parents)

$$\begin{aligned} &= (5/10)(0) + (5/10)(1.92) \\ &= \mathbf{0.96} \end{aligned}$$

Gain(Decision, Parents)

$$\begin{aligned} &= 1.57 - 0.96 \\ &= \mathbf{0.61} \end{aligned}$$

Entropy(Money_{Rich})

$$\begin{aligned} &= -((3/7) \cdot \log_2(3/7) + (2/7) \cdot \log_2(2/7) + (1/7) \cdot \log_2(1/7) + (1/7) \cdot \log_2(1/7)) \\ &= \mathbf{1.84} \end{aligned}$$

Entropy(Money_{Poor})

$$\begin{aligned} &= -((3/3) \cdot \log_2(3/3) + 0 + 0 + 0) \end{aligned}$$

$$= 0$$

$$\begin{aligned} &\text{Entropy}(\text{Decision}, \text{Money}) \\ &= (7/10)(1.84) + (3/10)(0) \\ &= 1.29 \end{aligned}$$

$$\begin{aligned} &\text{Gain}(\text{Decision}, \text{Money}) \\ &= 1.57 - 1.29 \\ &= \mathbf{0.28} \end{aligned}$$

Therefore, first split will be Weather.

2. Sunny Node.

$$\text{Entropy}(\text{Sunny}) = 0.92$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Parents}_{\text{Yes}}) \\ &= -((1/1) * \log_2(1/1) + 0 + 0 + 0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Parents}_{\text{No}}) \\ &= -(0 + (2/2) * \log_2(2/2) + 0 + 0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Parents}) \\ &= (1/3)(0) + (2/3)(0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} &\text{Gain}(\text{Sunny}, \text{Parents}) \\ &= 0.92 - 0 \\ &= \mathbf{0.92} \end{aligned}$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Money}_{\text{Rich}}) \\ &= -((1/3) * \log_2(1/3) + (2/3) * \log_2(2/3) + 0 + 0) \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Money}_{\text{Poor}}) \\ &= 0 \end{aligned}$$

$$\begin{aligned} &\text{Entropy}(\text{Sunny}, \text{Money}) \\ &= (3/3)(0.92) + 0 \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} &\text{Gain}(\text{Sunny}, \text{Parents}) \\ &= 0.92 - 0 \\ &= \mathbf{0} \end{aligned}$$

Therefore, Sunny node split will be Parents.

3. Windy Node.

$$\text{Entropy}(\text{Windy}) = 0.81$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Parents}_{\text{Yes}}) \\ &= -((2/2) * \log_2(2/2) + 0 + 0 + 0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Parents}_{\text{No}}) \\ &= -((1/2) * \log_2(1/2) + 0 + 0 + (1/2) * \log_2(1/2)) \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Parents}) \\ &= (2/4)(0) + (2/4)(1) \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Windy}, \text{Parents}) \\ &= 0.81 - 0.5 \\ &= \mathbf{0.31} \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Money}_{\text{Rich}}) \\ &= -((2/3) * \log_2(2/3) + 0 + 0 + (1/3) * \log_2(1/3)) \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Money}_{\text{Poor}}) \\ &= -((1/1) * \log_2(1/1) + 0 + 0 + 0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Windy}, \text{Money}) \\ &= (3/4)(0.92) + (1/4)(0) \\ &= 0.69 \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Windy}, \text{Money}) \\ &= 0.81 - 0.69 \\ &= \mathbf{0.12} \end{aligned}$$

Therefore, Windy node split will be Parents.

4. Rainy Node.

$$\text{Entropy}(\text{Rainy}) = 0.92$$

$$\begin{aligned} \text{Entropy}(\text{Rainy}, \text{Parents}_{\text{Yes}}) \\ &= -((2/2) * \log_2(2/2) + 0 + 0 + 0) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Rainy}, \text{Parents}_{\text{No}}) \\ &= -(0 + 0 + (1/1) * \log_2(1/1) + 0) \\ &= 0 \end{aligned}$$

$$\text{Entropy}(\text{Rainy}, \text{Parents})$$

= 0

Gain(Rainy, Parents)

= 0.92 - 0

= 0.92

Therefore, Windy node split will be Parents.

Tree Visualization. *Without Pruning.*

