



Module Code & Module Title

CU6051NA - Artificial Intelligence

Assignment

Year and Semester

2018-19 Spring

Student Name: Rajat Shrestha

London Met ID: 17030954

College ID: np01cp4a170021

Table of Contents

| | |
|--|---|
| Decision Tree Assignment Solution | 1 |
| Selecting the root node:..... | 2 |
| Selecting nodes for Sunny and Rainy nodes: | 5 |
| Building the final Decision tree: | 7 |

Table of Figures

| | |
|---|---|
| Figure 1: Table of data | 1 |
| Figure 2: Formulae to calculate Entropy and Gain | 2 |
| Figure 3: Python function to calculate Entropy | 2 |
| Figure 4: Python Function to calculate Gain..... | 2 |
| Figure 5: Tally of each category of attributes | 3 |
| Figure 6: Calculating the Gain value of each attributes | 4 |
| Figure 7: Tally of each category on Sunny node | 5 |
| Figure 8:Tally of each category on Rainy node | 5 |
| Figure 9: Gain calculated in Sunny node..... | 6 |
| Figure 10: Gain Calculated in Rainy node..... | 6 |
| Figure 11: Final decision tree | 7 |

Decision Tree Assignment Solution

Build a decision tree using data from the given table. The data is about whether to play tennis or not given the weather conditions. So, when the tree is built, it should be able to give out a decision on whether to play tennis or not given the weather conditions.

| Outlook | Temp. | Humidity | Wind | Decision |
|----------|-------|----------|--------|----------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |

Figure 1: Table of data

To build a decision tree we must go through various process to find the nodes:

Selecting the root node:

Here, we have four different features which is given the final decision to construct a decision tree we must first find the root node. To choose the optimal root node we must find the attribute with the maximum amount of gain.

$$\text{Entropy}(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$

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

Figure 2: Formulae to calculate Entropy and Gain

Since we need to do this calculation multiple times a function was created to take the count of positive and negative outcomes and give the entropy as well as another function to take the overall positive and negative outcomes with the list of all the positive and negative outcomes of individual category of an attribute:

```
[ ] 1 def entropy(positive,negative):
2     sum = negative+positive
3     try:
4         positive_ratio = math.log2(positive/sum)
5         negative_ratio = math.log2(negative/sum)
6     except ValueError:
7         positive_ratio = 0
8         negative_ratio = 0
9     return -(positive/sum)*positive_ratio-(negative/sum)*negative_ratio
```

1 entropy(9,5)

0.9402859586706309

Figure 3: Python function to calculate Entropy

```
[ ] 1 def gain(positive,negative,list):
2     parent_entropy = entropy(negative,positive)
3     sum = negative+positive
4     sum_of_values = 0
5     for i in list:
6         a = i[0]
7         b = i[1]
8         sum_of_values = sum_of_values + abs(abs(a+b)/abs(sum))*entropy(a,b)
9     return parent_entropy - sum_of_values
10
```

Figure 4: Python Function to calculate Gain

Now, to calculate the gain of each column we need to calculate the entropy of the all the labels categorized in an attribute, we must tally down all the positive and negative outcomes of all existing categories of the attributes.

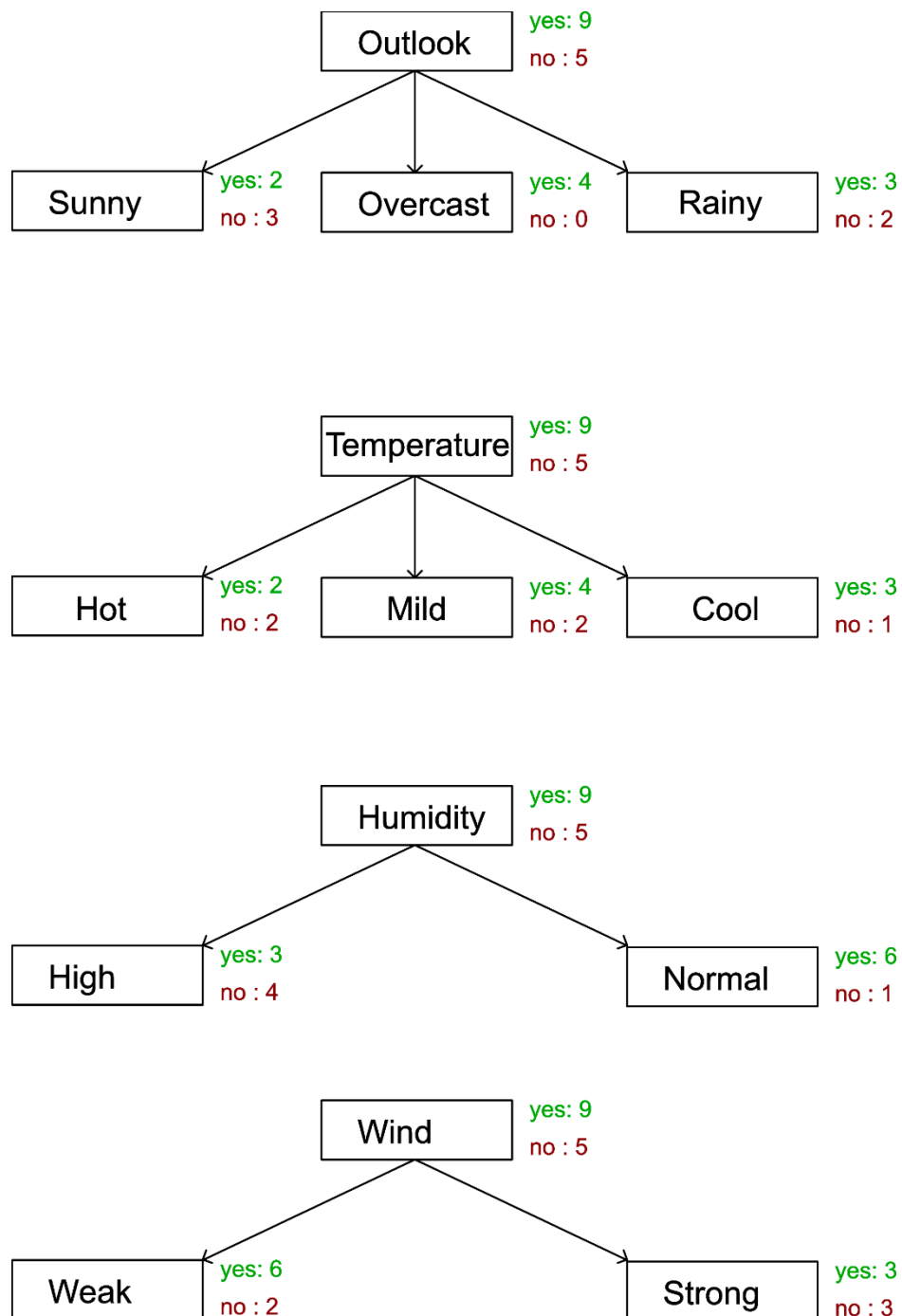


Figure 5: Tally of each category of attributes

```
[ ] 1 outlook = [[2,3],[4,0],[3,2]]
    2 temprature = [[2,2],[4,2],[3,1]]
    3 humidity = [[5,4],[6,1]]
    4 wind = [[6,2],[3,3]]
    5
    6 print("outlook_gain = "+str(gain(9,5,outlook)))
    7 print("temprature_gain = "+str(gain(9,5,temprature)))
    8 print("humidity_gain = "+str(gain(9,5,humidity)))
    9 print("wind_gain = "+str(gain(9,5,wind)))

➞ outlook_gain = 0.2467498197744391
   temprature_gain = 0.029222565658954647
   humidity_gain = 0.007329245197752798
   wind_gain = 0.04812703040826927
```

Figure 6: Storing above data in lists and calculating the Gain value of each attributes

As we can see the highest gain value is of outlook so taking outlook as the root node.

Now, as we got the root node as Outlook, observing the three root nodes, overcast is totally saturated while sunny and rainy is not. So, at the next step the same procedure is carried out with stripped data including the selected category.

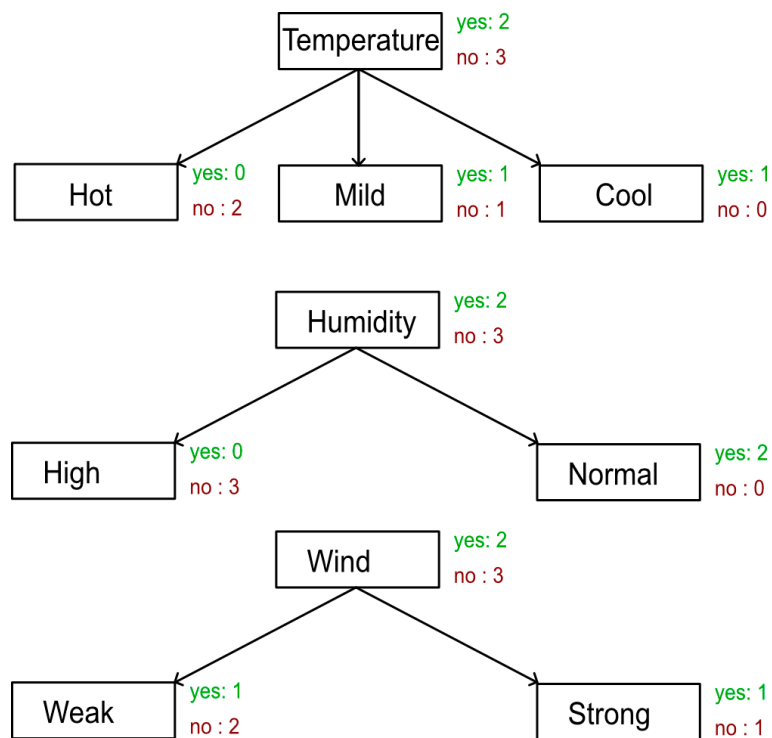
Selecting nodes for Sunny and Rainy nodes:

Figure 7: Tally of each category on Sunny node

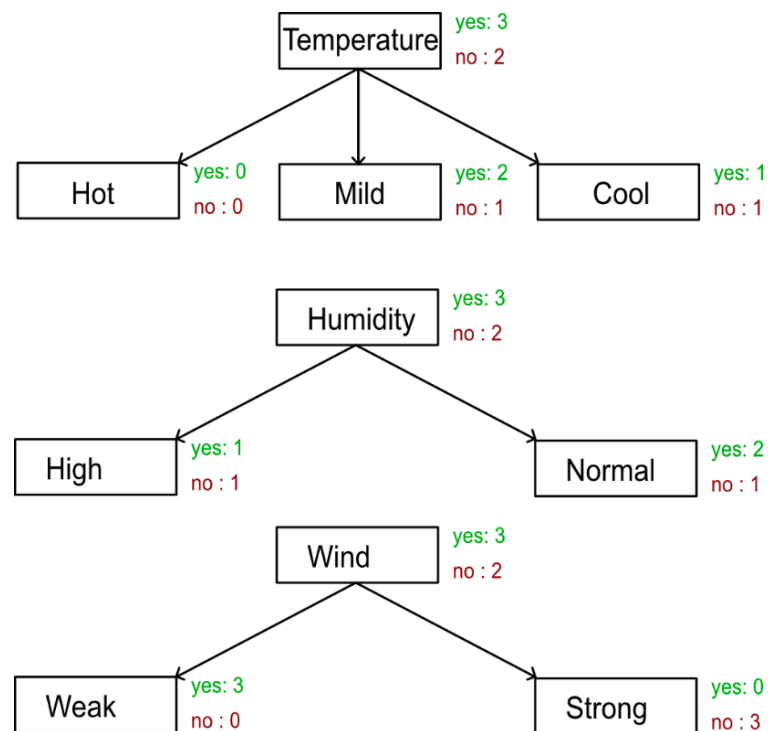


Figure 8: Tally of each category on Rainy node

Calculations using above data:

```
[11] 1 sunny_humidity = [[0,3],[2,0]]
      2 sunny_temperature = [[0,2],[1,1],[1,0]]
      3 sunny_wind = [[1,2],[1,1]]
      4
      5 print("sunny_humidity_gain = "+str(gain(2,3,sunny_humidity)))
      6 print("sunny_temperature_gain = "+str(gain(2,3,sunny_temperature)))
      7 print("sunny_wind_gain = "+str(gain(2,3,sunny_wind)))
```

```
☞ sunny_humidity_gain = 0.9709505944546686
   sunny_temperature_gain = 0.5709505944546686
   sunny_wind_gain = 0.01997309402197489
```

Figure 9: Gain calculated in Sunny node

```
[12] 1 rain_humidity = [[1,1],[2,1]]
      2 rain_temperature = [[2,1],[1,1]]
      3 rain_wind = [[3,0],[0,2]]
      4
      5 print("rain_humidity_gain = "+str(gain(3,2,rain_humidity)))
      6 print("rain_temperature_gain = "+str(gain(3,2,rain_temperature)))
      7 print("rain_wind_gain = "+str(gain(3,2,rain_wind)))
```

```
☞ rain_humidity_gain = 0.01997309402197489
   rain_temperature_gain = 0.01997309402197489
   rain_wind_gain = 0.9709505944546686
```

Figure 10: Gain Calculated in Rainy node

As we can see in the sunny node, humidity is the node with lowest saturation 0 and the highest gain and in the rainy node, wind gives the lowest degree of saturation and highest gain value so taking them as the child nodes.

Note: these calculations were unnecessary but if the nodes were not saturated it should be done to get the node with highest gain value.

Building the final Decision tree:

Finally, from the above calculations and observations, the resulting decision tree is as follows:

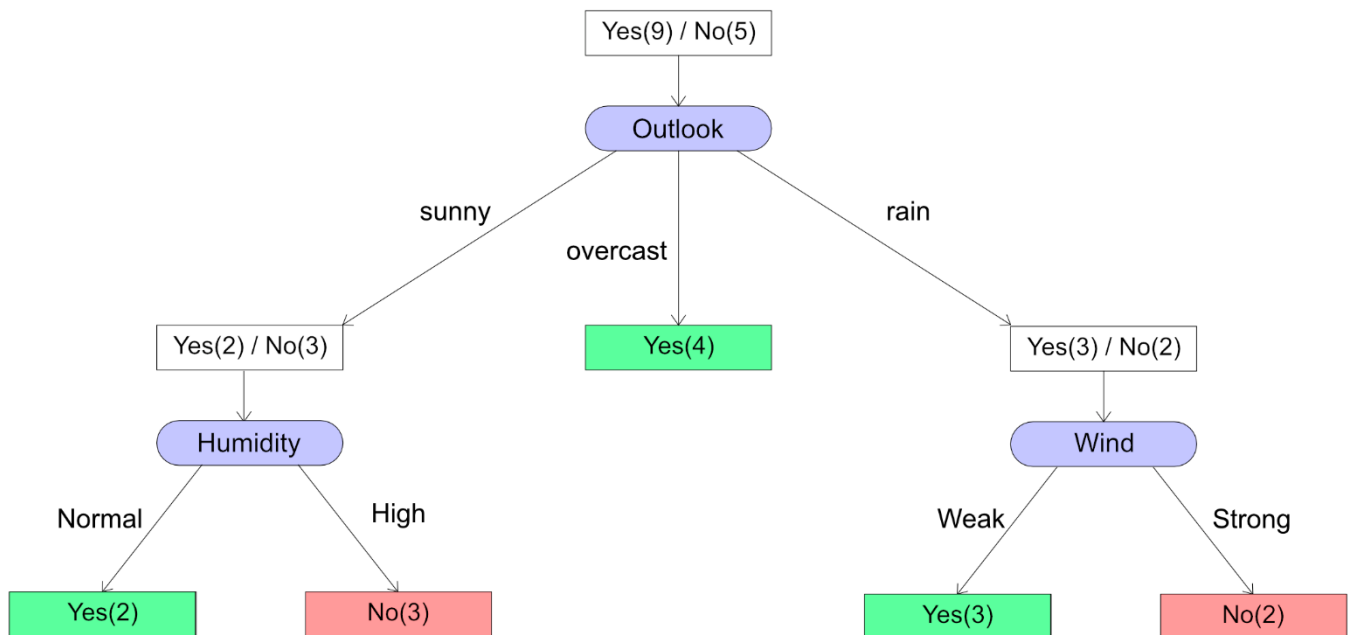


Figure 11: Final decision tree

The link to the notebook used to calculate the entropy and gain(a):

https://colab.research.google.com/drive/1fj8nGbi0hC6d0_8mUww1Ea14LHhi_X4Q