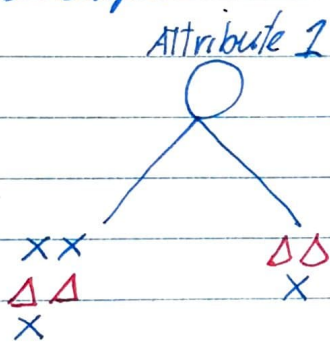


## Tutorial 2.

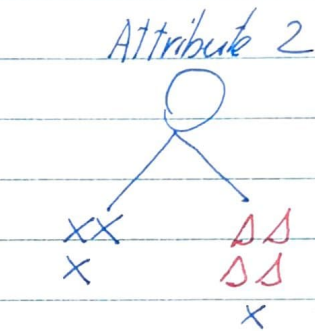
1 Describe how a decision tree could be learnt.

→ A key element in decision trees is to find the attributes that are most discriminative. Or the attributes that better separate the ~~var~~ instances that belong to different classes

### Examples



⏟  
This attribute does not separate well the classes



⏟  
This attribute seems that separates the classes better than attribute 1

But we have to do this numerically.....

b) Show how the idea of entropy could be used to pick the first node in the decision tree ...

\* As we mentioned before, we have to find the attributes that are more discriminative. Entropy can help us with this.

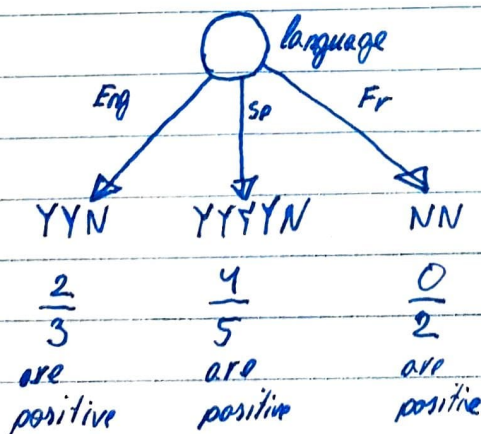
\* The higher the entropy of a split, the less discriminative the attribute is. Or what is the same, the higher the entropy the less helpful is that attribute to take good decisions.

\* The lower the entropy of a split, the more discrimination produced by the attribute.

here, discrimination is a good thing ...

It means that the capacity of the classifier to differentiate items of different classes/labels is higher.

→ Let's calculate the entropy of every split !!:  
Visual example for language attribute



1 Step: calculate the entropy of each split:

$$\text{Entropy}\left(\frac{2}{3}\right) = H\left(\frac{2}{3}\right) = -\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right) = 0.91$$

$$\text{Entropy}\left(\frac{4}{5}\right) = -\left(\frac{4}{5} \log_2 \frac{4}{5} + \frac{1}{5} \log_2 \frac{1}{5}\right) = 0.72$$

$$\text{Entropy}\left(\frac{0}{2}\right) = -\left(\frac{0}{2} \log_2 \frac{0}{2} + \frac{2}{2} \log_2 \frac{2}{2}\right) = 0$$

2 b) Now, once we have calculated the ~~split~~ weighted Entropy or total entropy based on the number/proportion of items that go to each split.

$$\text{Entropy (lang)} = \frac{3}{10} \times 0.91 + \frac{5}{10} \times 0.72 + \frac{2}{10} \times 0 \Rightarrow 0.63$$

So the entropy of the language split is this.

Now, we have to check the entropy of other splits and see which one is lower

$$\text{Entropy (Type)} = -\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right)$$

$$\frac{4}{10} \underbrace{\left(\text{Entropy (Action)}\right)}_{-\left(\frac{2}{4} \times \log_2 \frac{2}{4} + \frac{2}{4} \times \log_2 \frac{2}{4}\right)} + \frac{3}{10} \underbrace{\left(\text{Entropy (Comedy)}\right)}_{-\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right)} + \frac{3}{10} \underbrace{\left(\text{Entropy (Drama)}\right)}$$

0.951  $\Rightarrow$  So the entropy of type is higher than the entropy of lang. So, we will give preference to lang

Finally, we calculate the entropy for the attribute New

$$\text{Entropy (New)} = 0.846 \Rightarrow \text{So this one is still higher than language}$$

We choose language because has the lowest entropy

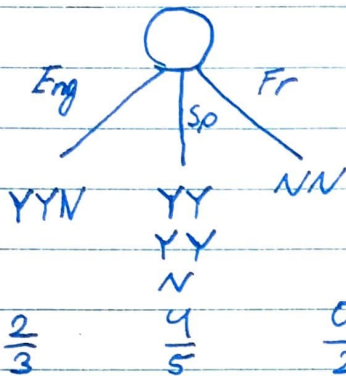
## Tutorial 2

Gini works in a similar way as Entropy. The lower the value the better the discrimination of that attribute.

### Question 3.

Now we have to use Gini impurity rather than entropy to decide the splits of our tree.

As we did before  $\rightarrow$  Let's calculate the gini Impurity of each split. Then, let's calculate the total one.



$$* \text{Gini}(\text{Eng}) = 1 - \left( \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right) = 0.444$$

$$* \text{Gini}(\text{Sp}) = 1 - \left( \left( \frac{4}{5} \right)^2 + \left( \frac{1}{5} \right)^2 \right) = 0.32$$

$$* \text{Gini}(\text{Fr}) = 1 - \left( \left( \frac{0}{2} \right)^2 + \left( \frac{2}{2} \right)^2 \right) = 0$$

Now, we have the gini impurity of each split. We have to calculate the total Gini

$$* \text{Gini}(\text{Language}) = \frac{3}{10} \times 0.444 + \frac{5}{10} \times 0.32 + \frac{2}{10} \times 0 = 0.29$$

$$* \text{Gini}(\text{Type}) = \frac{4}{10} \times \underbrace{\text{Gini}(\text{Action})}_{1 - \left( \left( \frac{2}{4} \right)^2 + \left( \frac{2}{4} \right)^2 \right)} + \frac{3}{10} \times \text{Gini}(\text{Comedy}) + \frac{3}{10} \times \underbrace{\text{Gini}(\text{Drama})}_{1 - \left( \left( \frac{2}{3} \right)^2 + \left( \frac{1}{3} \right)^2 \right)}$$

$$= \boxed{0.46}$$

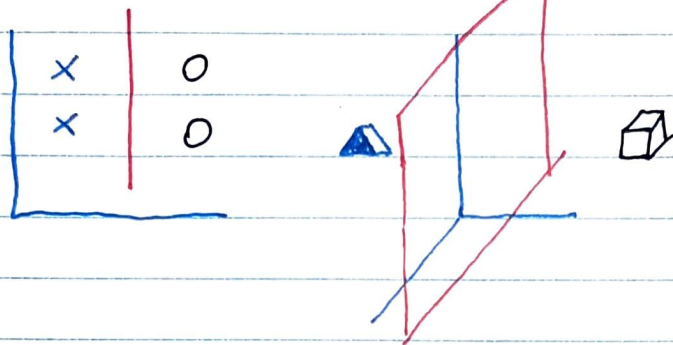
$$* \text{Gini}(\text{News}) = \boxed{0.4}$$

The lowest Gini is the one of the language attribute

## Tutorial 02

### Question 3

a) a <sup>of  $n$  dimension</sup> dataset is linearly separable if we can perfectly separate the classes of the data with a  $n-1$  dimensional plane



b) Examples: .....

o Extreme cases such as differentiating the characteristics of a flower and the characteristics of a non-biological entity..... if we choose the right parameters.

c) We will choose a classifier that is used ~~in~~ in scenarios ~~where~~ where the data is not linearly separable..... we will see some examples.

# Tutorial 2

## Question 4.

### batch gradient descent

→ So we update the weights with the following formula.

$$w_0 \leftarrow w_0 + \alpha \sum_j (y_j - h_w(x_j))$$

$$w_1 \leftarrow w_1 + \alpha \sum_j (y_j - h_w(x_j)) x_j$$

learning rate  $\alpha$   $\downarrow$   
 This indicates that goes over the entire dataset once...

Rather than explain this with a formula, let's do an example. INITIAL WEIGHTS  $\Rightarrow w_0 = 0$

Instance	x	y	prediction	bias $\downarrow$	error
E <sub>1</sub>	1.5	1	<del>1 * 0 + 1.5 * 0 = 0</del>	1 * 0 + 1.5 * 0 = 0	1 - 0
E <sub>2</sub>	3.5	3	<del>1 * 0 + 3.5 * 0 = 0</del>	1 * 0 + 3.5 * 0 = 0	3 - 0
E <sub>3</sub>	3	2	1 * 0 + 3 * 0 = 0		2 - 0
E <sub>4</sub>	5	3	1 * 0 + 5 * 0 = 0		3 - 0
E <sub>5</sub>	2	2.5	1 * 0 + 2 * 0 = 0		2.5 - 0

Total error = 11.5

UPDATE TIME  $\rightarrow$  1 Batch or epoch  
0.01

$$w_0 \leftarrow w_0 + \alpha \cdot \text{Total error} \rightarrow 0 + 0.01 \times 11.5 = 0.115$$

$$w_1 \leftarrow w_1 + \alpha \cdot \text{Total error (weighted)} \rightarrow 0 + 0.01 \times 38 = 0.38$$

	error	
E <sub>1</sub>	error = 1 * 1.5 = 1.5	} weighted error = 38
E <sub>2</sub>	error 3 * 3.5 = 10.5	
E <sub>3</sub>	error 2 * 3 = 6	
E <sub>4</sub>	error 3 * 5 = 15	
E <sub>5</sub>	error 2.5 * 2 = 5	

b) Now a couple of updates with stochastic gradient descent.

$$w_0 = 0, w_1 = 0, \alpha = 0.01$$

	$x_i$	$y$	Prediction	Error
→ $E_1$	1.5	1	$1 \times 0 + 1.5 \times 0 = 0$	$1 - 0 = 1$

UPDATE TIME (here we do not calculate the error of the entire dataset before updating. We update after every instance)

$$w_0 \rightarrow w_0 + \alpha \cdot \text{Error} \rightarrow 0 + 0.01 \cdot 1 = 0.01$$

$$w_1 \rightarrow w_1 + \alpha \cdot \text{error} \cdot x_i \rightarrow 0 + 0.01 \cdot 1 \cdot 1.5 = 0.015$$

↓  
new weights to use in the next iteration

	$x_i$	$y$	Prediction	error
$E_2$	3.5	3	$1 \times 0.01 + 0.015 \times 3.5 = 0.0625$	$3 - 0.0625 = 2.9375$

UPDATE TIME

$$w_0 \leftarrow 0.01 + 0.01 \cdot 2.9375 = 0.0393$$

$$w_1 \leftarrow 0.015 + 0.01 \cdot 2.9375 \times 3.5 = 0.118$$

↗  
new weights to use in the next iteration.

Tutorial 2  
Exercise 5

$$w_0 \leftarrow w_0 + \alpha (y - (w_0 + x_1 w_1 + x_2 w_2))$$

$$w_1 \leftarrow w_1 + \alpha (y - (w_0 + x_1 w_1 + x_2 w_2)) x_1$$

$$w_2 \leftarrow w_2 + \alpha (y - (w_0 + x_1 w_1 + x_2 w_2)) x_2$$