

---

## Artificial Intelligence

---

BS (CS) \_Spring\_2025

### Lab\_13 Manual



### Learning Objectives:

1. Decision Tree

In this lab we will practice Decision Tree using ID3 algorithm.

### Information Gain

To be able to calculate the information gain, we have to first implement the **entropy** of a dataset. The entropy of a dataset is used to measure the impurity of a dataset.

**What is Entropy?** In the most layman terms, Entropy is nothing but the measure of disorder. (You can think of it as a measure of purity as well. You'll see. I like disorder because it sounds cooler.)

The Mathematical formula for Entropy is as follows -

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Where 'Pi' is simply the frequentist probability of an element/class 'i' in our data. For simplicity's sake let's say we only have two classes, a positive class and a negative class. Therefore 'i' here could be either + or (-). So if we had a total of 100 data points in our dataset with 30 belonging to the positive class and 70 belonging to the negative class then 'P+' would be 3/10 and 'P-' would be 7/10. Pretty straightforward.

If I was to calculate the entropy of my classes in this example using the formula above. Here's what I would get.

$$-\frac{3}{10} \times \log_2\left(\frac{3}{10}\right) - \frac{7}{10} \times \log_2\left(\frac{7}{10}\right) \approx 0.88$$

The entropy here is approximately 0.88. This is considered a high entropy, a high level of disorder (meaning low level of purity). Entropy is measured between 0 and 1. (Depending on the number of classes in your dataset, entropy can be greater than 1 but it means the same thing, a very high level of disorder. For the sake of simplicity, the examples in this blog will have entropy between 0 and 1).

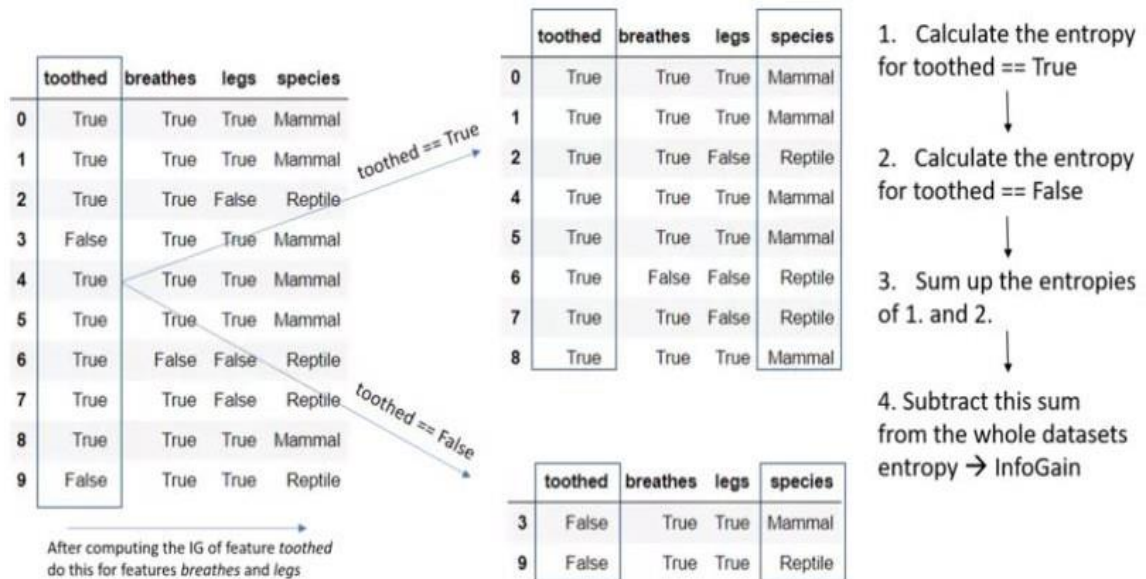
There are also other types of measures which can be used to calculate the information gain. The most prominent ones are the: Gini Index, Chi-Square, Information gain ratio and Variance. The information gain of a feature is calculated by

$$IG(Y, X) = E(Y) - E(Y|X)$$

OR

$$InfoGain(feature_d) = Entropy(D) - Entropy(feature_d)$$

So the only thing we have to do is to split the dataset along the values of each feature and then treat these sub sets as if they were our "original" dataset in terms of entropy calculation.



$$\begin{aligned}
 H(\text{toothed}) = & \left( \frac{8}{10} \left( -\left( \frac{5}{8} * \log_2\left(\frac{5}{8}\right) \right) + \left( \frac{3}{8} * \log_2\left(\frac{3}{8}\right) \right) \right) \right) + \left( \frac{2}{10} \left( -\left( \frac{1}{2} * \log_2\left(\frac{1}{2}\right) \right) + \left( \frac{1}{2} * \log_2\left(\frac{1}{2}\right) \right) \right) \right) \\
 & \underbrace{\hspace{10em}}_{\text{toothed = True ; Mammal} \quad \text{toothed = True ; Reptile}} \quad \underbrace{\hspace{10em}}_{\text{toothed = False ; Mammal} \quad \text{toothed = False ; Reptile}} \\
 & \underbrace{\hspace{15em}}_{\text{toothed = True}} \quad \underbrace{\hspace{15em}}_{\text{toothed = False}}
 \end{aligned}$$

Like this we will calculate the entropy and information gain of each attribute and select the best feature.

$$InfoGain(\text{toothed}) = 0.971 - 0.963547 = \mathbf{0.00745}$$

**breathes:**

$$H(\text{breathes}) = \left( \frac{9}{10} * -\left( \left( \frac{6}{9} * \log_2\left(\frac{6}{9}\right) \right) + \left( \frac{3}{9} * \log_2\left(\frac{3}{9}\right) \right) \right) \right) + \left( \frac{1}{10} * -\left( (0) + (1 * \log_2(1)) \right) \right) = \mathbf{0.82647}$$

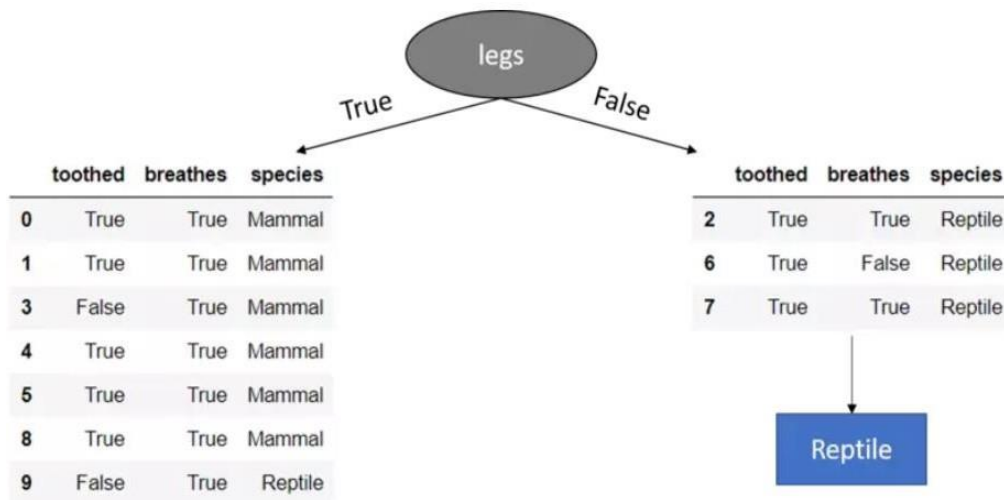
$$InfoGain(\text{breathes}) = 0.971 - 0.82647 = \mathbf{0.1445}$$

**legs:**

$$H(\text{legs}) = \left( \frac{7}{10} * -\left( \left( \frac{6}{7} * \log_2\left(\frac{6}{7}\right) \right) + \left( \frac{1}{7} * \log_2\left(\frac{1}{7}\right) \right) \right) \right) + \left( \frac{3}{10} * -\left( (0) + (1 * \log_2(1)) \right) \right) = \mathbf{0.41417}$$

$$InfoGain(\text{legs}) = 0.971 - 0.41417 = \mathbf{0.5568}$$

Hence the splitting the dataset along the feature *legs* results in the largest information gain and we should use this feature for our root node. Now our decision tree model looks like this:



We see that for **legs == False**, the target feature values of the remaining dataset are all *Reptile* and hence we set this as leaf node because we have a pure dataset (Further splitting the dataset on any of the remaining two features would not lead to a different or more accurate result since whatever we do after this point, the prediction will remain *Reptile*). Additionally, you see that the feature *legs* is no longer included in the remaining datasets. Because we already has used this (categorical) feature to split the dataset on it must not be further used. Until now we have found the feature for the root node as well as a leaf node for **legs == False**. The same steps for information gain calculation must now be accomplished also for the remaining dataset for **legs == True** since here we still have a mixture of different target feature values. Hence:

Information gain calculation for the features *toothed* and *breathes* for the remaining dataset **legs == True**:

**Entropy of the (new) sub data set after first split:**

$$H(D) = -\left(\left(\frac{6}{7} * \log_2\left(\frac{6}{7}\right)\right) + \left(\frac{1}{7} * \log_2\left(\frac{1}{7}\right)\right)\right) = 0.5917$$

**toothed:**

$$H(\text{toothed}) = \frac{5}{7} * -\left(\left(1 * \log_2(1)\right) + (0)\right) + \frac{2}{7} * -\left(\left(\frac{1}{2} * \log_2\left(\frac{1}{2}\right)\right) + \left(\frac{1}{2} * \log_2\left(\frac{1}{2}\right)\right)\right) = 0.285$$

$$\text{InfoGain}(\text{toothed}) = 0.5917 - 0.285 = 0.3067$$

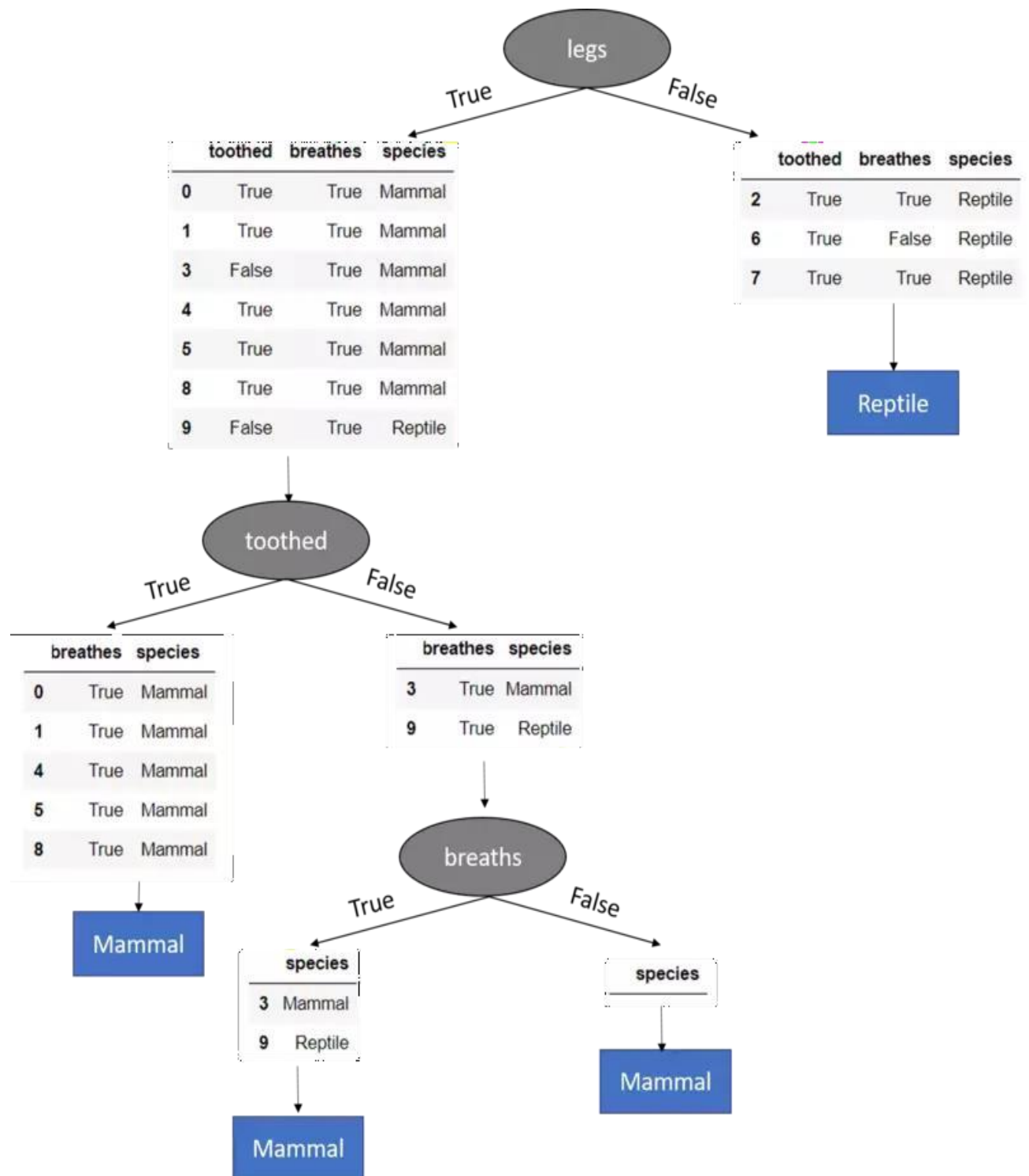
**breathes:**

$$H(\text{breathes}) = \frac{7}{7} * -\left(\left(\frac{6}{7} * \log_2\left(\frac{6}{7}\right)\right) + \left(\frac{1}{7} * \log_2\left(\frac{1}{7}\right)\right)\right) + 0 = 0.5917$$

$$\text{InfoGain}(\text{toothed}) = 0.5917 - 0.5917 = 0$$

The dataset for **toothed == False** still contains a mixture of different target feature values why we proceed partitioning on the last left feature (**== breathes**)

Now our tree looks like this.



Mind the last split (node) where the dataset got split on the **breathes** feature. Here the breathes feature solely contains data where breaths == True. Hence for breathes == False there are no instances in the dataset and therewith there is no **sub-Dataset** which can be built. In that case we return the most frequently occurring target feature value in the original dataset which is **Mammal**. This is an example of how our tree model generalizes behind the training data.

If we consider the other branch, that is breathes == True we know, that after splitting the Dataset on the values of a specific feature (breathes {True,False}) in our case, the feature must be removed. Well, that leads to a dataset where no more features are available to further split the dataset on. Hence we stop growing the tree and return the mode value of the direct parent node which is "**Mammal**".

That leads us to the introduction of the ID3 algorithm which is a popular algorithm to grow decision trees, published by Ross Quinlan in 1986. Besides the ID3 algorithm there are also other popular algorithms like the C4.5, the C5.0 and the CART algorithm.

Let's come back to the stopping criteria of the above grown tree. We can define a nearly arbitrarily large number of stopping criteria. Assume for instance, we say a tree is allowed to grow for only 2 seconds and then the growing process should stop - Well that would be a stopping criteria - Nonetheless, there are mainly three useful cases in which we stop the tree from growing assuming we do not stop it beforehand by defining for instance a maximum tree depth or a minimum information gain value. We stop the tree from growing when:

1. **All rows in the target feature have the same value**
2. **The dataset can be no longer split since there are no more features left**
3. **The dataset can no longer be split since there are no more rows left / There is no data left**

By definition, we say that if the growing gets stopped because of stopping criteria two, the leaf node should predict the most frequently occurring target feature value of the superior (parent) node. If the growing gets stopped because of the third stopping criteria, we assign the leaf node the mode target feature value of the original dataset.

## ID3 Pseudo Code

ID3(Data, Features, Target):

1. If all target values are the same:
  - Return a leaf node with that target value.
2. If Features is empty:
  - Return a leaf node with the majority target value in Data.
3. Else:
  - Select the feature with the highest Information Gain (best split).
  - Create a root node for that feature.
4. For each possible value of the selected feature:
  - Create a branch.
  - Subset = Data where selected feature == value.
  - If Subset is empty:
    - Add a leaf node with the majority target value from Data.
  - Else:
    - Recursively call ID3(Subset, Remaining Features, Target).
5. Return the root node.





1. Create a root node

toothed	hair	breathes	legs	species
0	True	True	True	Mammal
1	True	True	True	Mammal
2	True	False	True	Reptile
3	False	True	True	Mammal
4	True	True	True	Mammal
5	True	True	True	Mammal
6	True	False	False	Reptile
7	True	False	True	Reptile
8	True	True	True	Mammal
9	False	False	True	Reptile

2. Calculate the entropy of the whole (sub) dataset

Entropy

toothed	species	hair	species	breathes	species	legs	species
0	True Mammal	0	True Mammal	0	True Mammal	0	True Mammal
1	True Mammal	1	True Mammal	1	True Mammal	1	True Mammal
2	True Reptile	2	False Reptile	2	True Reptile	2	False Reptile
3	False Mammal	3	True Mammal	3	True Mammal	3	True Mammal
4	True Mammal	4	True Mammal	4	True Mammal	4	True Mammal
5	True Mammal	5	True Mammal	5	True Mammal	5	True Mammal
6	True Reptile	6	False Reptile	6	False Reptile	6	False Reptile
7	True Reptile	7	False Reptile	7	True Reptile	7	False Reptile
8	True Mammal	8	True Mammal	8	True Mammal	8	True Mammal
9	False Reptile	9	False Reptile	9	True Reptile	9	True Reptile

3. Calculate the Information gain of each single feature and pick that feature with the largest Information gain

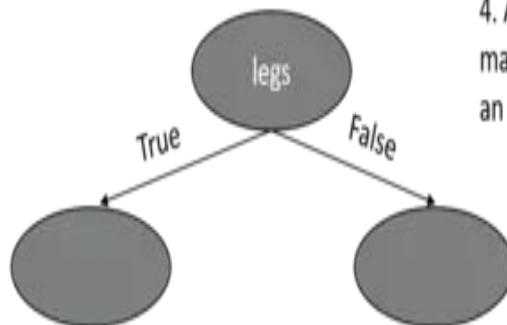
IG → toothed

IG → hair

IG → breathes

IG → legs

Select: max(IG\_feature)



4. Assign the (root) node the label of the feature with the maximum information gain. Grow for each feature value an outgoing branch and add unlabelled nodes at the end

Legs == True

toothed	hair	breathes	legs	species
0	True	True	True	Mammal
1	True	True	True	Mammal
3	False	True	True	Mammal
4	True	True	True	Mammal
5	True	True	True	Mammal
8	True	True	True	Mammal
9	False	False	True	Reptile

First sub tree

Legs == False

toothed	hair	breathes	legs	species
2	True	False	True	Reptile
6	True	False	False	Reptile
7	True	False	True	Reptile

Second sub tree

5. Split the dataset along the values of the maximum information gain feature and remove this feature from the dataset

6. For each of the sub\_datasets, repeat steps 2 to 5 until a stopping criteria is satisfied → Here the recursion kicks in



**References:**

1. <https://towardsdatascience.com/entropy-how-decision-trees-make-decisions-2946b9c18c8>
2. <https://medium.com/analytics-vidhya/entropy-calculation-information-gain-decision-tree-learning-771325d16f>
3. <https://pradeep-dhote9.medium.com/decision-tree-classification-and-regression-algorithm-cart-id3-36a2450a7f1a>
4. <https://stackoverflow.com/questions/20092632/can-someone-explain-me-the-difference-between-id3-and-cart-algorithm>