

Report 1: Dectrees

Group XXX: XXX, Riccardo Fragale

January 24, 2026

1 Assignment 0

Each one of the datasets has properties which makes them hard to learn. Motivate which of the three problems is most difficult for a decision tree algorithm to learn.

The bigger challenge for all the three datasets is the number of training samples compared to test samples. Let's consider the dataset **MONK1** as example to justify that: it has exactly 431 test samples while the number of training samples is 123. Considering that we don't have a validation set and we don't use a k-fold cross-validation, we have a really limited amount of data to train the model on. This reasoning is valid also for **MONK2** and **MONK3**. This makes the decision-tree less able to generalise and have a complete picture of the classification problem to be solved.

Going in detail to the single datasets we can say that:

- **MONK1** appears to be the less hard of the three to model as there is a clearer rule that depends only on 3 attributes while the others are mostly irrelevant to predict the classification as true or false
- **MONK2** is the harder to classify because the pattern to be identified is more complex as discrete attributes have value 1 that is not repeated but simply two of them, randomly, get the value 1
- **MONK3** has the lower number of training samples so it is harder for the model to get an accurate classification. Moreover there is a 5% additional classification noise in the training set which makes the work even harder.

2 Assignment 2

The file dtree.py defines a function entropy which calculates the entropy of a dataset. Import this file along with the monks datasets and use it to calculate the entropy of the training datasets.

Dataset	Entropy
MONK-1	1.0
MONK-2	0.957117428264771
MONK-3	0.9998061328047111

3 Assignment 2

Explain entropy for a uniform distribution and a non-uniform distribution, present some example distributions with high and low entropy.

4 Assignment 3

Use the function averageGain (defined in dtree.py) to calculate the expected information gain corresponding to each of the six attributes. Note that the attributes are represented as instances of the class Attribute (defined in monkdata.py) which you can access via m.attributes[0], ..., m.attributes[5]. Based on the results, which attribute should be used for splitting the examples at the root node?

5 Assignment 4

For splitting we choose the attribute that maximizes the information gain, Eq.3. Looking at Eq.3 how does the entropy of the subsets, S_k , look like when the information gain is maximized? How can we motivate using the information gain as a heuristic for picking an attribute for splitting? Think about reduction in entropy after the split and what the entropy implies.

6 Assignment 5

Build the full decision trees for all three Monk datasets using `buildTree`. Then, use the function `check` to measure the performance of the decision tree on both the training and test datasets.

For example to built a tree for `monk1` and compute the performance on the test data you could use

```
import monkdata as m
import dtree as d

t=d.buildTree(m.monk1, m.attributes);
print(d.check(t, m.monk1test))
```

Compute the train and test set errors for the three Monk datasets for the full trees. Were your assumptions about the datasets correct? Explain the results you get for the training and test datasets.

7 Assignment 6

Explain pruning from a bias variance trade-off perspective.

8 Assignment 7

Evaluate the effect pruning has on the test error for the `monk1` and `monk3` datasets, in particular determine the optimal partition into training and pruning by optimizing the parameter `fraction`. Plot the classification error on the test sets as a function of the parameter $\text{fraction} \in \{0.3, 0.4, 0.5, 0.6, 0.7, 0.8\}$.