

Introduction to ML - Decision Tree Coursework

(version 1.1)

October 2021

Overview

In this assignment, you will implement a decision tree algorithm and use it to determine one of the indoor locations based on WIFI signal strengths collected from a mobile phone. See Figure 1 for an illustration of the experimental scenario. The results of your experiments should be discussed in the report. You should also deliver the code you have written.

Setup

We do recommend you work on the Ubuntu workstations in the lab. This assignment and all code were tested for Linux and Mac OS machines. We cannot guarantee compatibility with Windows machine and cannot promise any support if you do choose to work on a Windows machine.

Working on DoC lab workstations (recommended)

You can also work from home and use the lab workstations. See this list of <https://www.doc.ic.ac.uk/csg/facilities/lab/workstations> to ssh into one of the machines.

In order to load all the packages that you might need for the course work, you can run the following command:

```
export PYTHONUSERBASE=/vol/lab/ml/mlenv
```

We installed numpy, matplotlib, pytorch, which should cover most of the packages you will need for all the courseworks in this course. If you don't want to type that command every time that you connect to your lab machine, you can add it to your bashrc:

```
echo "export PYTHONUSERBASE=/vol/lab/ml/mlenv" >> ~/.bashrc
```

This way, every terminal you open will have that environment variable set. It is recommended to use "python3" exclusively. The current python3 version in lab machines is 3.8.5. To test the configuration:

```
python3 -c "import numpy as np; print(np)"
```

This should print:

```
<module 'numpy' from '/vol/lab/ml/mlenv/lib/python3.8/site-packages/numpy/__init__.py'>
```

Working on your own system

If you decide to work locally on your machine, then you will need to give us explicit instructions how to run your code. Anything we cannot run will result to your points of the question reduced by 30%.

Python>= 3.5: All provided code has been tested on Python versions 3.5 or 3.6. Make sure to install Python version 3.5 or 3.6 on your local machine. Otherwise, you might encounter errors! If you are working on Mac OSX, you can use Homebrew to brew install python3.

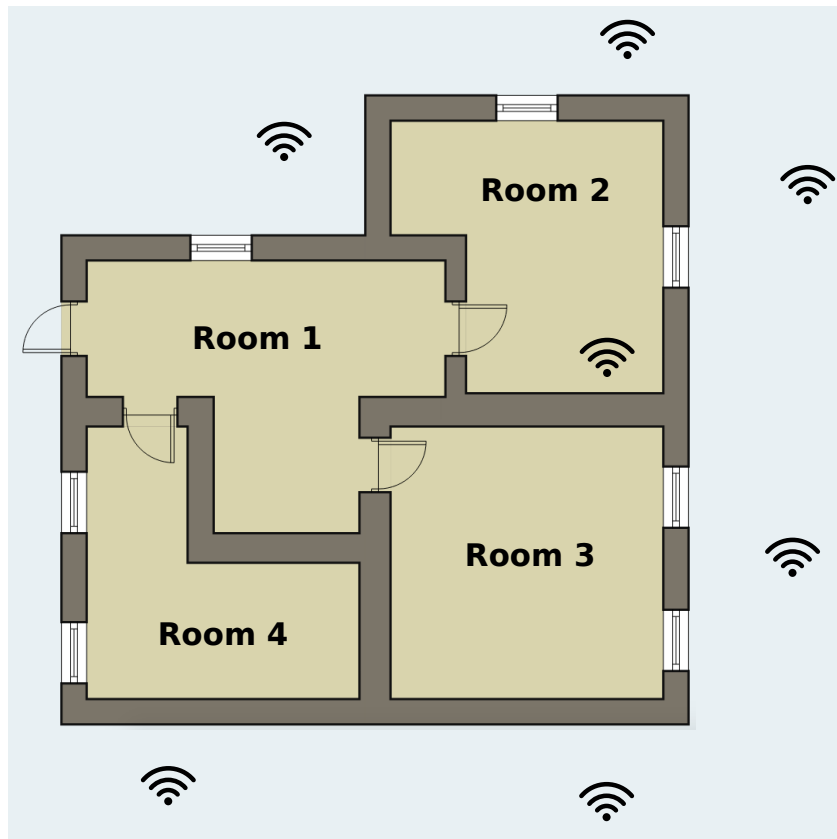


Figure 1: Illustration of the scenario. The WIFI signal strength from 7 emitters are recorded from a mobile phone. The objective of this coursework is to learn decision tree that predict in which of the 4 rooms the user is standing.

Coursework

Step 1: Loading data

You can load the datasets from the files "WIFI_db/clean_dataset.txt" and "WIFI_db/noisy_dataset.txt". They contain a 2000x8 array. This array represents a dataset of 2000 samples. Each sample is composed of 7 WIFI signal strength while the last column indicates the room number in which the user is standing (i.e., the label of the sample). **All the features in the dataset are continuous except the room number.** You can load the text file with the "loadtxt" function from Numpy. Given the nature of the dataset you will have to build decision trees capable of dealing with continuous attributes and multiple labels.

For the report:

There is nothing to add in the report for this section.

Step 2: Creating Decision Trees

To create the decision tree, you will write a recursive function called *decision_tree_learning()*, that takes as arguments a matrix containing the dataset and a *depth* variable (which is used to compute the maximal depth of the tree, for plotting purposes for instance). The label of the training dataset is assumed to be the last column of the matrix. The pseudo-code of this function is described below. This pseudo-code is taken from *Artificial Intelligence: A Modern Approach* by Stuart Russell and Peter Norvig (Figure 18.5), but modified to take into account that the considered dataset contains continuous attributes (see section 18.3.6 of the book).

The function FIND_SPLIT chooses the attribute and the value that results in the highest information gain. Because the dataset has continuous attributes, the decision-tree learning algorithms search for the split point (defined by an attribute and a value) that gives the highest information gain. For instance, if you have two attributes (A0 and A1) with values that both range from 0 to 10, the algorithm might determine that splitting the dataset according to "A1>4" gives the most information. An efficient method for finding good split points

Algorithm 1 Decision Tree creating

```
1: procedure DECISION_TREE_LEARNING(training_dataset, depth)
2:   if all samples have the same label then
3:     return (a leaf node with this value, depth)
4:   else
5:     split ← FIND_SPLIT(training_dataset)
6:     node ← a new decision tree with root as split value
7:     l_branch, l_depth ← DECISION_TREE_LEARNING(l_dataset, depth+1)
8:     r_branch, r_depth ← DECISION_TREE_LEARNING(r_dataset, depth+1)
9:     return (node, max(l_depth, r_depth))
10:  end if
11: end procedure
```

is to sort the values of the attribute, and then consider only split points that are between two examples in sorted order, while keeping track of the running totals of examples of each class for each side of the split point.

To evaluate the information gain, suppose that the training dataset S_{all} has K different labels. We can define two subsets (S_{left} and S_{right}) of the dataset depending on the splitting rule (for instance " $A1 > 4$ ") and for each dataset and subset, we can compute the distribution (or probability) of each label. For instance, $\{p^1 p^2 \dots p^K\}$ (p^k is the number of samples with the label k divided by the total number of samples from the initial dataset). The information gain is defined by using the general definition of the entropy as follow:

$$\text{Gain}(S_{all}, S_{left}, S_{right}) = H(S_{all}) - \text{Remainder}(S_{left}, S_{right})$$

$$H(\text{dataset}) = - \sum_{k=1}^{k=K} p_k * \log_2(p_k)$$

$$\text{Remainder}(S_{left}, S_{right}) = \frac{|S_{left}|}{|S_{left}| + |S_{right}|} H(S_{left}) + \frac{|S_{right}|}{|S_{left}| + |S_{right}|} H(S_{right})$$

Where $|S|$ represents the number of samples in subset S .

Implementation:

You are only allowed to use Numpy and Matplotlib for this coursework. Any other library, like scikit learn is NOT allowed. For the implementation of the tree (in Python), it is advised to use dictionaries to store nodes as a single object, for instance by using this kind of structure $\{\text{'attribute'}, \text{'value'}, \text{'left'}, \text{'right'}\}$. In this case, 'left' and 'right' will also be nodes. You might also want to add a Boolean field name "leaf" that indicates whether or not the node is a leaf (terminal node) or not. However, this is not strictly necessary, as there are other methods to determine if a node is terminal or not.

For the report:

There is nothing to add in the report for this section.

Bonus part (5 points):

Write a function to visualize the tree. See Figure 2 for an example. You can only use Numpy and Matplotlib for this question as well. **Show in your report the output of this function with a tree trained on the entire clean dataset.**

Step 3 - Evaluation

Now that you have an automatic process to train decision trees, you can evaluate the accuracy of your tree on the provided datasets. For that, evaluate your decision tree using a 10-fold cross validation on both the clean and noisy datasets. You should expect that slightly different trees will be created with each fold, since the training data that you use each time will be slightly different. Use your resulting decision trees to classify your data in your test sets.

Implementation:

Implement an evaluation function that takes a trained tree and a test dataset: `evaluate(test_db, trained_tree)` and that returns the accuracy of the tree.

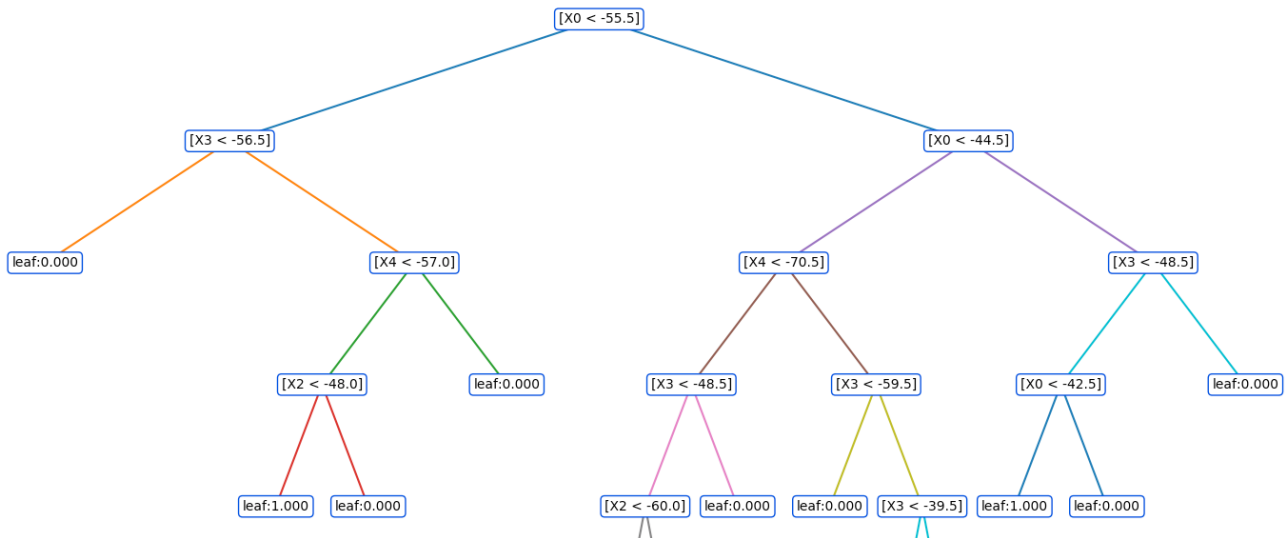


Figure 2: example of Decision tree visualization (not all the tree is displayed, and this tree has been trained on a different dataset.)

For the report:

Cross validation classification metrics By computing the average over all the test folds, report the following cross validation classification metrics **for both clean and noisy data**:

- Confusion matrix. (Hint: you should get a single 4x4 matrix)
- The accuracy (Hint: you can derive the metrics directly from the previously computed confusion matrix).
- The recall and precision rates per class.
- The F1-measures derived from the recall and precision rates of the previous step.

Result analysis Comment for both datasets which rooms are recognized with high/low accuracy, and which rooms are confused. **5 lines max.**

Dataset differences: Is there any difference in the performance when using the clean and noisy datasets? If yes/no explain why. **5 lines max.**

Step 4 - Pruning (and evaluation again)

In order to reduce the performance difference of our decision tree between the clean and noisy dataset, you will implement a pruning function based on reducing the validation error. This approach works as follow: for each node directly connected to two leaves, evaluate the benefits on the validation error of substituting this node with a single leaf (defined according to the training set). If a single leaf reduces the validation error, then the node is pruned and replaced by a single leaf. The tree needs to be parsed several times until there is no more node connected to two leaves (HINT: when you prune a node, the parent node might now verify this condition).

For the report:

Cross validation classification metrics after pruning Report the performances of your trees after pruning by using a nested 10-fold cross validation ("option 2") to compute the metrics defined in the previous section for both datasets.

Result analysis after pruning Comment the difference of performance before and after pruning for both datasets. Briefly explain these performance differences. **5 lines max.**

Depth analysis Comment on the average depth of the trees that you generated for both datasets, before and after pruning. What can you tell about the relationship between maximal depth and prediction accuracy? **5 lines max.**

Deliverables

For the completion of this part of the coursework, the following has to be submitted electronically via CATE:

- All the code you have written
- A README file (in .txt, .md, or .pdf) to explain how to run your code on the lab machines.
- A short report with the following structure:
 - (Bonus points: Output of the tree visualisation function.)
 - Step 3 - Evaluation
 - * Cross validation classification metrics.
 - * Result analysis (5 lines max).
 - * Dataset differences (5 lines max).
 - Step 4 - Pruning (and evaluation again)
 - * Cross validation classification metrics after pruning.
 - * Result analysis after pruning (5 lines max).
 - * Depth analysis (5 lines max).

Grading scheme

Final Grade = Report content + Code quality + Report quality

- Code (total : 20)
 - Results on secret test dataset: **10**
Make sure that your code runs. If not, you will be asked to resubmit the code and lose 50% of the code mark
 - Presentation of the code (comments, indentation, structure): **5**
 - README instructions: **5**
- Report content (total : 70)
 - Output of the tree visualisation function: **BONUS 5**
 - Step 3 - Evaluation
 - * Cross validation classification metrics (total 15)
 - Confusion matrix: **10**
 - The accuracy: **1**
 - The recall and precision per class: **2**
 - The F1-measures: **2**
 - * Result analysis (5 lines max): **10**
 - * Dataset differences (5 lines max): **10**
 - Step 4 - Pruning (and evaluation again)
 - * Cross validation classification metrics after pruning (total 15)
 - Confusion matrix: **10**
 - The accuracy: **1**
 - The recall and precision per class: **2**
 - The F1-measures: **2**
 - * Result analysis after pruning (5 lines max): **10**
 - * Depth analysis (5 lines max): **10**
- Report quality (total : 10)
 - Quality of presentation: **10**

Acknowledgements

This coursework is inspired from the coursework designed by Maja Pantic.