

HOMEWORK 2

>>NAME HERE<<
>>ID HERE<<

Instructions: Although this is a programming homework, you only need to hand in a pdf answer file. There is no need to submit the latex source or any code. You can choose any programming language, as long as you implement the algorithm from scratch (e.g. do not use Weka on questions 1 to 7).

Use this latex file as a template to develop your homework. Submit your homework on time as a single pdf file to Canvas. Please check Piazza for updates about the homework.

1 A Simplified Decision Tree

You are to implement a decision-tree learner for classification. To simplify your work, this will not be a general purpose decision tree. Instead, your program can assume that

- each item has two continuous features $\mathbf{x} \in \mathbb{R}^2$
- the class label is binary and encoded as $y \in \{0, 1\}$
- data files are in plaintext with one labeled item per line, separated by whitespace:

$$\begin{array}{ccc} x_{11} & x_{12} & y_1 \\ & \dots & \\ x_{n1} & x_{n2} & y_n \end{array}$$

Your program should implement a decision tree learner according to the following guidelines:

- Candidate splits (j, c) for numeric features should use a threshold c in feature dimension j in the form of $x_{.j} \geq c$.
- c should be on values of that dimension present in the training data; i.e. the threshold is on training points, not in between training points.
- The left branch of such a split is the “then” branch, and the right branch is “else”.
- Splits should be chosen using mutual information (i.e. information gain). If there is a tie you may break it arbitrarily.
- The stopping criteria (for making a node into a leaf) are that
 - the node is empty, or
 - all splits have zero mutual information
- To simplify, whenever there is no majority class in a leaf, let it predict $y = 1$.

2 Questions

1. (Our algorithm stops at pure labels) [10 pts] If a node is not empty but contains training items with the same label, why is it guaranteed to become a leaf? Explain.

In such a setup, where we have one non-empty node, but with training items of the same label(let's assume predicted label = 0) we have two options at hand:

- Hypothesis H_1 : Stop determining any further splits from that node
- Hypothesis H_2 : Determine a set of thresholds at this node n , $C_n = c_1, c_2, \dots, c_n$ such that there are splits determined at each of these thresholds, and the label is always 0 for each such split. Furthermore, let's assume there are $a_1, a_2, a_3, \dots, a_n$ elements for each such split, such that $a_1 + a_2 + \dots + a_n = A_n$, where A_n is the number of elements remaining at node n and have the same label 0.

Now, Entropy at node n :

$$E(n, A_n) = -p_1 \log(p_1) - p_0 \log(p_0) = -0 * \log(0) - 1 * \log(1) = 0$$

Entropy at any subsplit c_i

$$E(c_i, a_i) = -p_1 \log(p_1) - p_0 \log(p_0) = -0 * \log(0) - 1 * \log(1) = 0$$

Average Entropy of all n sub-splits due to threshold-set C is:

$$E(C, A_n) = \frac{a_1}{A_n} E(c_1, a_1) + \frac{a_2}{A_n} E(c_2, a_2) + \dots + \frac{a_n}{A_n} E(c_n, a_n) = 0 + 0 + \dots + 0 = 0$$

Thus, information gain:

$$E(n, A_n) - E(C, A_n) = 0$$

Therefore, there is no information gain by using Hypothesis H_2 and, simply by Occam's razor, we should choose one amongst the two competing hypotheses that make exactly the same predictions, and the chosen one should be the simpler one, which is H_1 in this case. Thus, we note that when all elements in a node have the same label, we see that all splits have zero mutual information, which is a stopping criteria and thus, it becomes a stopping criteria by extension.

2. (Our algorithm is greedy) [10 pts] Handcraft a small training set where both classes are present but the algorithm refuses to split; instead it makes the root a leaf and stop; Importantly, if we were to manually force a split, the algorithm will happily continue splitting the data set further and produce a deeper tree with zero training error. You should (1) plot your training set, (2) explain why. Hint: you don't need more than a handful of items.
3. (Mutual information exercise) [10 pts] Use the training set Druns.txt. For the root node, list all candidate cuts and their mutual information. Hint: to get $\log_2(x)$ when your programming language may be using a different base, use $\log(x) / \log(2)$.
4. (The king of interpretability) [10 pts] Decision tree is not the most accurate classifier in general. However, it persists. This is largely due to its rumored interpretability: a data scientist can easily explain a tree to a non-data scientist. Build a tree from D3leaves.txt. Then manually convert your tree to a set of logic rules. Show the tree¹ and the rules.
5. (Or is it?) [20 pts] For this question only, make sure you DO NOT VISUALIZE the data sets or plot your tree's decision boundary in the 2D x space. If your code does that, turn it off before proceeding. This is because you want to see your own reaction when trying to interpret a tree. You will get points no matter what your interpretation is. And we will ask you to visualize them in the next question anyway.
 - Build a decision tree on D1.txt. Show it to us in any format (e.g. could be a standard binary tree with nodes and arrows, and denote the rule at each leaf node; or Weka style plaintext tree; or as simple as plaintext output where each line represents a node with appropriate line number pointers to child nodes; whatever is convenient for you). Again, do not visualize the data set or the tree in the x input space. In real tasks you will not be able to visualize the whole high dimensional input space anyway, so we don't want you to "cheat" here.
 - Look at your tree in the above format (remember, you should not visualize the 2D dataset or your tree's decision boundary) and try to interpret the decision boundary in human understandable English.
 - Build a decision tree on D2.txt. Show it to us.

¹When we say show the tree, we mean either the standard computer science tree view, or some crude plaintext representation of the tree – as long as you explain the format. When we say visualize the tree, we mean a plot in the 2D x space that shows how the tree will classify any points.

- Try to interpret your D2 decision tree.
6. (Hypothesis space) [10 pts] For D1.txt and D2.txt, do the following separately:
- Produce a scatter plot of the data set.
 - Visualize your decision tree's decision boundary (or decision region, or some other ways to clearly visualize how your decision tree will make decisions in the feature space).

Then discuss why the size of your decision trees on D1 and D2 differ. Relate this to the hypothesis space of our decision tree algorithm.

7. (Learning curve) [20 pts] We provide a data set Dbig.txt with 10000 labeled items. Caution: Dbig.txt is sorted.
- You will randomly split Dbig.txt into a candidate training set of 8192 items and a test set (the rest). Do this by generating a random permutation, and split at 8192.
 - Generate a sequence of five nested training sets $D_{32} \subset D_{128} \subset D_{512} \subset D_{2048} \subset D_{8192}$ from the candidate training set. The subscript n in D_n denotes training set size. The easiest way is to take the first n items from the (same) permutation above. This sequence simulates the real world situation where you obtain more and more training data.
 - For each D_n above, train a decision tree. Measure its test set error err_n . Show three things in your answer: (1) List n , number of nodes in that tree, err_n . (2) Plot n vs. err_n . This is known as a learning curve (a single plot). (3) Visualize your decision trees' decision boundary (five plots).

3 Weka [10 pts]

Learn to use Weka <https://www.cs.waikato.ac.nz/~ml/weka/index.html>. Convert appropriate data files into ARFF format. Use trees/J48 as the classifier and default settings. Produce five Weka trees for $D_{32}, D_{128} \dots D_{8192}$. Show two things in your answer: (1) List n , number of nodes in that tree, err_n . (2) Plot n vs. err_n .