

CS 5350/6350: Machine Learning Fall 2016

Homework 1

Handed out: 30 August, 2016
Due date: 13 September, 2016

General Instructions

- You are welcome to talk to other members of the class about the homework. I am more concerned that you understand the underlying concepts. However, you should write down your own solution. Please keep the class collaboration policy in mind.
- Feel free discuss the homework with the instructor or the TAs.
- Your written solutions should be brief and clear. You need to show your work, not just the final answer, but you do *not* need to write it in gory detail. Your assignment should be **no more than 10 pages**. Every extra page will cost a point.
- Handwritten solutions will not be accepted.
- The homework is due by midnight of the due date. Please submit the homework on Canvas.
- Some questions are marked **For 6350 students**. Students who are registered for CS 6350 should do these questions. Of course, if you are registered for CS 5350, you are welcome to do the question too, but you will not get any credit for it.

1 Decision trees (35 points)

1. [6 points] To warm up, represent the following Boolean functions as decision trees. (It is unnecessary to make the decision tree as small as possible; you can choose any root as you like. Use 1 for True and 0 for False. Also, note that an easy way to represent decision trees is as a series of **if-then-else** statements.)
 - (a) $(x_1 \vee x_2) \wedge x_3$
 - (b) $(x_1 \wedge x_2) \text{ xor } (\neg x_1 \vee x_3)$
 - (c) The 2-of-3 function defined as follows: at least 2 of $\{x_1, x_2, x_3\}$ should be true for the output to be true.
2. When playing Pokémon Go, there is some chance that a Pokémon will be caught or it will escape. In the following question, build a decision tree to determine whether a Pokémon can be caught.
There are four features:

- (a) **Berry** (*Yes or No*) means whether a Razz Berry was used.
- (b) **Ball** (*Poké, Great, or Ultra*) describes which kind of ball has been thrown.
- (c) **Color** (*Green, Yellow, or Red*) stands for the difficulty level of catching this Pokémon.
- (d) **Type** (*Normal, Water, Flying, or Psychic*) depicts the type of the Pokémon.

Berry	Ball	Color	Type	Caught
Yes	Poké	Green	Normal	Yes
No	Poké	Yellow	Normal	No
No	Great	Yellow	Normal	No
Yes	Ultra	Yellow	Normal	Yes
No	Poké	Red	Normal	No
Yes	Great	Red	Normal	Yes
Yes	Great	Green	Water	Yes
No	Great	Yellow	Water	No
Yes	Poké	Red	Water	No
No	Ultra	Red	Water	Yes
No	Poké	Red	Flying	No
Yes	Great	Yellow	Flying	Yes
No	Ultra	Yellow	Flying	Yes
Yes	Great	Red	Flying	Yes
No	Poké	Green	Psychic	No
No	Great	Yellow	Psychic	No

Table 1: Training set for Pokémon Go

- (a) [2 points] How many possible functions are there to map these four features to a Boolean decision?
- (b) [2 points] What is the entropy of the labels in this data? (When calculating entropy, The base of the logarithm should be 2.)
- (c) [8 points] Calculate information gain for four features respectively. Keep 3 significant digits.
- (d) [3 points] According to your results, using ID3 algorithm which attribute should be root for the decision tree?
- (e) [4 points] Construct a decision with the root you selected in the previous question. You do not have to use the ID3 algorithm here, you can show any tree with the chosen root.
- (f) [2 points] Using your decision tree to predict label in the test set in the table below, what is your label for the each example? What is your accuracy?
- (g) [1 points] Do you think it is a good idea to use decision tree in this Pokémon Go problem?

Berry	Ball	Color	Type	Caught
Yes	Great	Yellow	Psychic	Yes
Yes	Poké	Green	Flying	No
No	Ultra	Red	Water	No

Table 2: Test set for Pokémon Go

3. Recall that in the ID3 algorithm, we want to identify the best attribute that splits the examples that are relatively pure in one label. Apart from entropy, which you used in the previous question, there are other methods to measure impurity. One such impurity measure is the Gini measure, that is used in the CART family of algorithms. If there are k possible outcomes $1, \dots, i, \dots, k$, each with a probability $p_1, \dots, p_i, \dots, p_k$ of occurring, the Gini measure is defined as:

$$Gini(p_1, \dots, p_k) = 1 - \sum_{i=1}^k p_i^2$$

The Gini measure can be used to replace entropy in the definition of information gain to pick the best attribute.

- (a) [4 points] Using the Gini measure, calculate the information gain for the four features respectively. Use 3 significant digits.
- (b) [3 points] According to your results in the last question, which attribute should be the root for the decision tree? Do these two measures (entropy and Gini) lead to the same tree?

2 Linear Classifiers (15 points)

In the questions in this section, we have four features x_1, x_2, x_3 and x_4 and the label is represented by o .

1. [3 points] Write a linear classifier that correctly classifies the given dataset. You don't need to run any learning algorithm here. Try to find the weights and the bias of the classifier using the definition of linear separators.

x_1	x_2	x_3	x_4	o
0	0	0	1	1
0	0	1	1	1
0	0	0	0	-1

2. [5 points] Suppose the dataset below is an extension of the above dataset. Check if your classifier from the previous question correctly classifies the dataset. Report its accuracy.

x_1	x_2	x_3	x_4	o
1	0	1	1	1
0	1	0	1	1
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1
1	1	1	0	1
0	0	1	0	-1

3. [7 points] Given the remaining missing data points of the above dataset in the table below, find a linear classifier that correctly classifies the whole dataset (all three tables together).

x_1	x_2	x_3	x_4	o
0	1	0	0	-1
0	1	1	0	-1
0	1	1	1	1
1	0	0	0	1
1	0	0	1	1
1	1	0	1	1

3 Experiments

In this question you will be implementing a decision tree learner. You will experiment with the decision tree hyperparameters using cross-validation.

This problem uses the Mushroom Data set from the UCI machine learning repository. Each data point has 22 features indicating different characteristics of a mushroom. You can find definitions of each feature in the `mushroom.names` file. Your goal is to use the ID3 algorithm on the provided training datasets to train a predictor and see how well it does on the test data. You may use Java, Python, Matlab, C/C++ for this assignment. If you want to use a different language, you must contact the instructor first. Any other language you may want to use **MUST** run on the CADE machines.

Cross-Validation

The depth of the tree is a hyper-parameter to the decision tree algorithm that helps reduce overfitting. You will see later in the semester that many machine learning algorithm (SVM, logistic-regression etc) have some hyper-parameters as their input. One way to determine a proper value for the hyper-parameter is to use a technique called *cross-validation*.

As usual we have a training set and a test set. Our goal is to discover good hyperparameters using the training set. To do so, you can put aside some of the training data aside, and when training is finished, you can test the resulting classifier on the held out data. This allows you to get an idea of how well the particular choice of hyper-parameters does. However, since you did not train on your whole dataset you may have introduced a statistical

bias in the classifier. To correct for this, you will need to train many classifiers with different subsets of the training data removed and average out the accuracy across these trials.

For problems with small data sets, a popular method is the leave-one-out approach. For each example, a classifier is trained on the rest of the data and the chosen example is then evaluated. The performance of the classifier is the average accuracy on all the examples. The downside to this method is for a data set with n examples you must train n different classifiers. Of course, this is not practical for the data set you will use in this problem, so you will hold out subsets of the data many times instead.

Specifically, for this problem, you should implement k -fold cross validation. The general approach for k -fold cross validation is the following: Suppose you want to evaluate how good a particular hyper-parameter is. You split the training data into k parts. Now, you will train your model on $k - 1$ parts with the chosen hyper-parameter and evaluate the trained model on the remaining part. You should repeat this k times, choosing a different part for evaluation each time. This will give you k values of accuracy. Their *average cross-validation accuracy* gives you an idea of how good this choice of the hyper-parameter is. To find the best value of the hyper-parameter, you will need to repeat this procedure for different choices of the hyper-parameter. Once you find the best value of the hyper-parameter, use the value to retrain your classifier using the entire training set.

Setting A [25 points]

1. [10 points] **Implementation**

For this problem, you will be using the data found in the **SettingA** folder. This folder contains two files, **SettingA/training.data** and **SettingA/test.data**. In this setting you will be training your algorithm on the training file (**SettingA/training.data**). Remember that you should not look at or use your testing file until your algorithm is complete.

- (a) [4 points] Implement the decision tree data structure and the ID3 algorithm for your decision tree. (Remember that the decision tree need not be a binary tree!). For debugging your implementation, you can use the previous toy examples like the Pokémon data from Table 1. Discuss what approaches or choices you had to make during this implementation.
- (b) [2 points] Report the error of your decision tree on the **SettingA/training.data** file.
- (c) [5 points] Report the error of your decision tree on the **SettingA/test.data** file.
- (d) [1 points] Report the maximum depth of your decision tree.

2. [15 points] **Limiting Depth**

In this section you will be using 6-fold cross-validation in order to limit the depth of your decision tree, effectively pruning the tree to avoid overfitting. We have split the data into 6 parts for you: you will be using the 6 cross-validation files for this section, titled **SettingA/CVSpits/training_0X.data** where X is a number between 0 and 5 (inclusive).

- (a) [10 points] Run 6-fold cross-validation using the specified files. Experiment with depths in the set $\{1, 2, 3, 4, 5, 10, 15, 20\}$, reporting the average cross-validation accuracy and standard deviation for each depth.
- (b) [5 points] Using the depth with the greatest average cross-validation accuracy from your experiments: train your decision tree on the `SettingA/training.data` file. Report the accuracy of your decision tree on the `SettingA/test.data` file.

Setting B [25 points]

1. [10 points] Experiments

For this problem, you will be using the data found in the `SettingB` folder. This folder contains the two files, `SettingB/training.data` and `SettingB/test.data`. In this setting you will be training your algorithm on the training file (`SettingB/training.data`). Remember that you should not look at or use your testing file until your algorithm is complete. You are not limiting the depth of your tree in this section.

- (a) [2 points] Report the error of your decision tree on the `SettingB/training.data` file.
- (b) [2 points] Report the error of your decision tree on the `SettingB/test.data` file.
- (c) [2 points] Report the error of your decision tree on the `SettingA/training.data` file.
- (d) [2 points] Report the error of your decision tree on the `SettingA/test.data` file.
- (e) [1 points] Report the maximum depth of your decision tree.

2. [15 points] Limiting Depth

In this section you will be using 6-fold cross-validation in order to limit the depth of your decision tree, effectively pruning the tree to avoid overfitting. You will be using the 6 cross-validation files for this section, titled `SettingB/CVSplits/training_0X.data` where `X` is a number between 0 and 5 (inclusive).

- (a) [10 points] Run 6-fold cross-validation using the specified files. Experiment with depths in the set $\{1, 2, 3, 4, 5, 10, 15, 20\}$, reporting the cross-validation accuracy and standard deviation for each depth. Explicitly specify which depth should be chosen as the best, and explain why.
- (b) [5 points] Using the depth with the greatest cross-validation accuracy from your experiments: train your decision tree on the `SettingB/training.data` file. Report the accuracy of your decision tree on the `SettingB/test.data` file.

Setting C (CS 6350 Students) [20 points]

In this setting, you are investigating what effect missing features have on a decision tree, and exploring which approach is most effective in dealing with missing features. More specifically, you will be trying:

- **Method 1:** Setting the missing feature as the majority feature value.
- **Method 2:** Setting the missing feature as the majority value of that label.
- **Method 3:** Treating the missing feature as a *special* feature.

The missing feature is represented by a ? character. In order to determine which method is the best, you will be using 6-fold cross-validation, with the files being titled `SettingC/CVSpplits/training_0X.data` where X is a number between 0 and 5 (inclusive). These files, along with `SettingC/training.data` and `SettingC/test.data` can be found in the `SettingC` folder.

1. [5 points] Update your decision tree implementation to have functionality to deal with missing features. Describe your approach/any choices you had to make in this implementation.
2. [10 points] Perform 6-fold cross-validation on each of the 3 methods described above. Report the accuracy for each method and the standard deviation.
3. [5 points] Using the best method selected from your experiments, train your decision tree on `SettingC/training.data`, and report the accuracy of your tree on `SettingC/test.data`.

Experiment Submission Guidelines

1. The report should detail your experiments. For each step, explain in no more than a paragraph or so how your implementation works. You may provide the results for the final step as a table or a graph.
2. *Your code should run on the CADE machines.* You should include a shell script, `run.sh`, that will execute your code in the CADE environment. Your code should produce similar output to what you include in your report.

You are responsible for ensuring that the grader can execute the code using only the included script. If you are using an esoteric programming language, you should make sure that its runtime is available on CADE.

3. Please do not hand in binary files! We will *not* grade binary submissions.