

Note: the one “bonus” problem on this assignment is marked by a *.

Problems

- (24 points, written) You’re working for a bank, and they have some data on whether a loan application should be approved or not:

Example No.	House	Bills	Income	Credit	Approve?
1	Rent	Late	Low	Bad	Yes
2	Rent	Late	High	Bad	Yes
3	Rent	On-time	Low	Bad	Yes
4	Rent	On-time	Medium	Bad	Yes
5	Own	Late	Medium	Good	Yes
6	Own	Late	Low	Bad	No
7	Own	Late	Medium	Good	No
8	Own	On-time	Low	Bad	No
9	Own	On-time	High	Good	Yes
10	Own	On-time	Medium	Good	Yes
11	Own	On-time	High	Bad	No
12	Own	On-time	High	Good	Yes

Suppose you start to learn a decision tree on the above data using “Approve?” as the goal predicate. Calculate the information gains for splitting the data (at the top level) for each of the four attributes. Which attribute provides the largest information gain (and would therefore be the top level attribute in a decision tree)? Show and explain your work.

- (36 points) Write the procedure:

```
(learn-dtree training-data attribute-names)
```

which returns a decision tree learned from the training data using the algorithm covered in class (also in the text) with the “greatest information gain” heuristic.

The Scheme representations for decision trees and training data are described in the remainder of this assignment handout. Support code will be provided to do some of the more mundane data manipulation tasks.

- (30 points) Write the procedure:

```
(chi^2-learn-dtree training-data attribute-names)
```

which learns a decision tree as for the previous problem but incorporates χ^2 pruning as described in the text. This procedure will be similar to your `learn-dtree` procedure. Further details for this problem are described on the Assignment 6 web page.

- (40 points, written) For this problem, you will run tests on your `learn-dtree` and `chi^2-learn-dtree` procedures in order to analyze and compare their performance. Details of these tests and on what you must turn in for this problem are described on the Assignment 6 web page.
- (8 points) For this problem, you will write a procedure that automatically discretizes a continuous valued attribute. See the Assignment 6 web page for details.

Scheme representations

Training data and examples

Training data consist of a list of training examples, and a *training example* is a list where the first element is the value of the goal predicate (which can be any symbol, not just `yes` or `no`) and the second element is a list of the attribute values. We will also require a list of *attribute names* so we can refer to attributes by name.

Here are the last four training examples from the first problem made into a training data set:

```
(define loan-names '(House Bills Income Credit))
(define loan-data-small
  '((Yes (Own On-time High Good))
    (Yes (Own On-time Medium Good))
    (No (Own On-time High Bad))
    (Yes (Own On-time High Good))))
```

Note that the goal predicate is not explicitly named.

Decision trees

A decision tree is either a value for the goal predicate (i.e. a symbol) or a list of the following form:

```
(<attribute-name> (<attribute-value-1> <decision-tree-1>)
  ...
  (<attribute-value-n> <decision-tree-n>))
```

For the “loan” example, a valid decision tree is:

```
(define loan-dtree-example
  '(income (high yes)
    (low (house (rent no)
      (own yes)))))
```

Support code

You need not use the following procedures, but you will probably find them helpful. They are organized into several different categories in the sections below.

Handling training data

- `(split training-data attribute-names attribute)`

This function divides the training data into groups according to the specified attribute. For example, using the training data above:

```
(split loan-data-small loan-names 'income)
;Value: ((medium ((yes (own on-time medium good))))
  (high ((yes (own on-time high good))
    (no (own on-time high bad))
    (yes (own on-time high good)))))
```

This procedure returns a list of what I refer to as *splits*. Each split is a list whose first element is a value of the attribute and whose second element is a list containing a subset of the training data which all have that value for the given attribute.

Note that in the example above, there is no split generated for the income attribute value “low” because the training data do not have an example with this value. See the “Implementation notes” section for discussion of this issue.

- `(tally training-data)`

This procedure counts the number of examples for each value of the goal predicate. It returns a list of clauses. Each clause is a list where the first element each is the value of the goal predicate, and the second element is the number of examples with that value.

For example,

```
(tally loan-data-small)
;Value: ((no 1) (yes 3))
```

Do not assume that the goal predicate will always have the values “yes” and “no”!

Like the `split` procedure, if there are no examples with a given goal predicate value, that value will not appear in the tally.

- `(pick-majority tally)`

Given a tally (as returned by the `tally` procedure), this procedure returns the majority value. If there is a tie, it returns the first instance it finds. For example:

```
(pick-majority '((no 1) (yes 3)))
;Value: yes
```

Testing your decision trees

- `(classify example decision-tree attribute-names default-value)`

This function returns a classification for the example determined by the given decision tree. If an attribute value not in the decision tree is encountered, then it returns the default-value.

For example:

```
(classify '(rent late high good) loan-dtree-example loan-names 'No)
;Value: yes
```

- `(test-dtree decision-tree training-data attribute-names)`

This function takes a decision tree and a set of training data. From the training data, it creates a list of examples (i.e. just the attribute values) and a list of correct classifications. It classifies all the examples using the decision tree, compares the results to the correct classifications, and reports the results.

Miscellaneous utility functions

- `(remove-element el Lst)`

Removes the first occurrence of `el` in `Lst`, returning a newly created list.

Implementation notes

Differences from the text’s algorithm

The basic decision tree learning algorithm you should implement for this assignment is slightly different than the algorithm in our text. The difference is in how the decision tree will handle examples that have attribute values not seen in the training data.

The algorithm in the text handles this by making a recursive call to the `learn-dtree` procedure with zero examples. This returns a decision (sub)tree that consists of a leaf node: the default classification.

The way the support code for this assignment is structured, you should never make a recursive call to your `learn-dtree` procedure when there are no examples left in a given branch. Instead, the `classify` procedure returns the default value if it encounters an attribute value not in the decision tree.

The reason for this difference is to simplify your code. In order to implement the text's algorithm, you would have to know all the values of each attribute, and they would have to be passed down from one recursive call to the next. Leaving this situation to be handled by the `classify` procedure means that only the attribute names need to be passed.

As an example, consider the decision tree in Figure 18.8 of the text. My solutions, run on the same training data produce the decision tree:

```
(learn-dtree restaurant-data restaurant-names)
;Value: (patrons (none no)
         (full (hungry (no no)
                (yes (type (burger yes)
                           (italian no)
                           (thai (fri (no no)
                                     (yes yes)))))))
         (some yes))
```

Notice that there is no "french" value handled under `type`. This is because the french restaurants in the training data were classified under other cases of the decision tree. (One training example had "some" patrons; in the other, patrons was "full", and "hungry" was "no".)

Suggestions

- Because the decision tree representations, tallies, and splits can be confusing, I strongly suggest using simple accessor functions to access information from these data structures.
- Take advantage of the fact that Scheme is interpreted and test your procedures from the bottom up — make sure your lower level functions are doing the right thing before you go on to the higher level functions!
- Do not attempt to take the logarithm of 0!!!

Data sets

There will be several data sets available for you to test your procedures. See the Assignment 6 web page for the files.