

CS178 Homework #4
Machine Learning & Data Mining: Winter 2016
Due: Friday February 26th, 2016

Write neatly (or type) and show all your work!

Please remember to turn in at most two documents, one with any handwritten solutions, and one PDF file with any electronic solutions.

Download the provided Homework 4 code, to replace / add to last week's code (several new functions have been added).

Problem 1: Decision Trees

We'll use the same data as in our earlier homework: In order to reduce my email load, I decide to implement a machine learning algorithm to decide whether or not I should read an email, or simply file it away instead. To train my model, I obtain the following data set of binary-valued features about each email, including whether I know the author or not, whether the email is long or short, and whether it has any of several key words, along with my final decision about whether to read it ($y = +1$ for "read", $y = -1$ for "discard").

x_1	x_2	x_3	x_4	x_5	y	
know author?	is long?	has 'research'	has 'grade'	has 'lottery'	\Rightarrow read?	
0	0	1	1	0	-1	
1	1	0	1	0	-1	
0	1	1	1	1	-1	
1	1	1	1	0	-1	
0	1	0	0	0	-1	In the case of any
1	0	1	1	1	1	
0	0	1	0	0	1	
1	0	0	0	0	1	
1	0	1	1	0	1	
1	1	1	1	1	-1	

ties, we will prefer to predict class +1.

- (a) Calculate the entropy of the class variable y
- (b) Calculate the information gain for each feature x_i . Which feature should I split on first?
- (c) Draw the complete decision tree that will be learned from these data.

Problem 2: Decision Trees on Kaggle

In this problem, we will revisit our Kaggle in-class competition,

<https://inclass.kaggle.com/c/cs178-project-2016>

and build a simple regression tree model to make predictions on the data. You can use the `treeRegress` class provided to build your regression trees.

Note: Kaggle competitions only let you submit a fixed number of predictions per day, ≈ 5 in our case, so be careful. We'll use a validation split to decide what hyperparameter choices we think are most promising, and upload only one model.

- (a) Split out a validation set from your training data examples, and learn a decision tree regressor on the data. To avoid any potential recursion limits, specify a max depth of 20, e.g.,

```
dt = treeRegress(Xt,Yt, maxDepth=20)
```

(This might take a bit of time; ≈ 4 minutes on my desktop.) Compute your model's validation MSE.

- (b) Now, try varying the maximum depth parameter (**maxDepth**), which forces the tree to stop after at most that many levels. Test values **0**, **1**, ..., **15** and compare their performance (both training and test) against the full depth. Is complexity increasing or decreasing with the depth cutoff? Identify whether you think the model begins overfitting, and if so, when. If you use this parameter for complexity control, what depth would you select as best?
- (c) Now, using high maximum depth ($d = 20$), use **minParent** to control complexity. Try values **2.[^][3:12]=[8,16,...,4096]**. Is complexity increasing or decreasing as **minParent** grows? Identify when (if) the model is starting to overfit, and what value you would use for this type of complexity control.
- (d) Using your best complexity control value (either depth or number of parent data), re-train a model on the full data set, predict on the test data, upload it to Kaggle and report its performance.

Problem 3: Ensembles of Trees

Choose **either part** of this question to answer (your choice): a random forest regressor, which is a bagged ensemble of decision trees; **or** an boosted ensemble of regression trees learned with gradient boosting.

In Python, it is easy to keep a list of different learners, even of different types, for use in an ensemble predictor:

```
ensemble[i] = ml.treeRegress(Xb,Yb,...) # save ensemble member "i" in a cell array
# ...
ensemble[i].predict(Xv,Yv);           # find the predictions for ensemble member "i"
```

Option 1: Random forests:

Random Forests are bagged collections of decision trees, which select their decision nodes from randomly chosen subsets of the possible features (rather than all features). You can implement this easily in **treeRegress** using option '**nFeatures**'= n , where n is the number of features to select from (e.g., $n = 50$ or $n = 60$ if there are 90-some features); you'll write a for-loop to build the ensemble members, and another to compute the prediction of the ensemble.

- (a) Using your validation split, learn a bagged ensemble of decision trees on the training data and evaluate validation performance. (See the pseudocode from lecture slides.) For your individual learners, use little complexity control (depth cutoff 15+, minParent 8, etc.), since the bagging will be used to control overfitting instead. For the bootstrap process, draw the same number of data as in your training set after the validation split ($M' = M$ in the pseudocode). You may

find `ml.bootstrapData()` helpful, although it is very easy to do yourself. Plot the training and validation error as a function of the number of learners you include in the ensemble, for (at least) 1, 5, 10, 25 learners. (You may find it more computationally efficient to simply learn 25 ensemble members first, and then evaluate the results using only a few of them; this will give the same results as only learning the few that you need.)

- (b) Now choose an ensemble size and repeat on the full training data, make predictions on the test data, and upload to Kaggle. Report your performance.

Option 2: Gradient boosting:

Gradient boosted trees are boosted collections of decision trees, which are build sequentially to predict the residual error in the current ensemble. You'll write a for-loop to build the ensemble members, and another to compute the prediction of the ensemble.

- (a) Using your validation split, learn a gradient boosted ensemble of decision trees on the training data and evaluate validation performance. (See the pseudocode from lecture slides.) For your individual learners, use very strong complexity control (depth cutoff 2–3, or large `minParent`, etc.), since the boosting process will be adding complexity to the overall learner. Plot the training and validation error as a function of the number of learners you include in the ensemble, for (at least) 1, 5, 10, 25 learners. (You may find it more computationally efficient to simply learn 25 ensemble members, and then evaluate the results using fewer of them.)
- (b) Now choose an ensemble size and repeat on the full training data, make predictions on the test data, and upload to Kaggle. Report your performance.