

**Homework #7**

RELEASE DATE: 12/25/2015

DUE DATE: 1/6/2015, BEFORE NOON

QUESTIONS ABOUT HOMEWORK MATERIALS ARE WELCOMED ON THE COURSERA FORUM.

*Unless granted by the instructor in advance, you must turn in a printed/written copy of your solutions (without the source code) for all problems.*

*For problems marked with (\*), please follow the guidelines on the course website and upload your source code to designated places. You are encouraged to (but not required to) include a README to help the TAs check your source code. Any programming language/platform is allowed.*

*Any form of cheating, lying, or plagiarism will not be tolerated. Students can get zero scores and/or fail the class and/or be kicked out of school and/or receive other punishments for those kinds of misconducts.*

*Discussions on course materials and homework solutions are encouraged. But you should write the final solutions alone and understand them fully. Books, notes, and Internet resources can be consulted, but not copied from.*

*Since everyone needs to write the final solutions alone, there is absolutely no need to lend your homework solutions and/or source codes to your classmates at any time. In order to maximize the level of fairness in this class, lending and borrowing homework solutions are both regarded as dishonest behaviors and will be punished according to the honesty policy.*

*You should write your solutions in English or Chinese with the common math notations introduced in class or in the problems. We do not accept solutions written in any other languages.*

This homework set comes with 200 points and 20 bonus points. In general, every homework set would come with a full credit of 200 points, with some possible bonus points.

**Decision Tree**

Impurity functions play an important role in decision tree branching. For binary classification problems, let  $\mu_+$  be the fraction of positive examples in a data subset, and  $\mu_- = 1 - \mu_+$  be the fraction of negative examples in the data subset.

1. The Gini index is  $1 - \mu_+^2 - \mu_-^2$ . What is the maximum value of the Gini index among all  $\mu_+ \in [0, 1]$ ? Prove your answer.
2. Following Question 1, we can normalize each impurity function by dividing it with its maximum value among all  $\mu_+ \in [0, 1]$ . For instance, the classification error is simply  $\min(\mu_+, \mu_-)$  and its maximum value is 0.5. So the normalized classification error is  $2 \min(\mu_+, \mu_-)$ . After normalization, prove or disprove that the normalized Gini index is equivalent to the normalized squared regression error (used for branching in classification data sets), where the squared error is by definition  $\mu_+(1 - (\mu_+ - \mu_-))^2 + \mu_-(-1 - (\mu_+ - \mu_-))^2$ .

**Random Forest**

3. If bootstrapping is used to sample  $N' = pN$  examples out of  $N$  examples and  $N$  is very large, argue that approximately  $e^{-p} \cdot N$  of the examples will not be sampled at all.
4. Consider a Random Forest  $G$  that consists of three binary classification trees  $\{g_k\}_{k=1}^3$ , where each tree is of test 0/1 error  $E_{\text{out}}(g_1) = 0.15$ ,  $E_{\text{out}}(g_2) = 0.25$ ,  $E_{\text{out}}(g_3) = 0.35$ . What is the possible range of  $E_{\text{out}}(G)$ ? Justify your answer.
5. Consider a Random Forest  $G$  that consists of  $K$  binary classification trees  $\{g_k\}_{k=1}^K$ , where  $K$  is an odd integer. Each  $g_k$  is of test 0/1 error  $E_{\text{out}}(g_k) = e_k$ . Prove or disprove that  $\frac{2}{K+1} \sum_{k=1}^K e_k$  upper bounds  $E_{\text{out}}(G)$ .

### Gradient Boosting

6. Let  $\epsilon_t$  be the weighted 0/1 error of each  $g_t$  as described in the AdaBoost algorithm (Lecture 208), and  $U_t = \sum_{n=1}^N u_n^{(t)}$  be the total example weight during AdaBoost. Express  $U_3$  by an equation that involves only  $\epsilon_1, \epsilon_2$ . Justify your answer.
7. For the gradient boosted decision tree, if a tree with only one constant node is returned as  $g_1$ , and if  $g_1(\mathbf{x}) = 2$ , then after the first iteration, all  $s_n$  is updated from 0 to a new constant  $\alpha_1 g_1(\mathbf{x}_n)$ . What is  $s_n$ ? Prove your answer.
8. For the gradient boosted decision tree, after updating all  $s_n$  in iteration  $t$  using the steepest  $\eta$  as  $\alpha_t$ , what is the value of  $\sum_{n=1}^N s_n g_t(\mathbf{x}_n)$ ? Prove your answer.

### Neural Network

9. Consider Neural Network with  $\text{sign}(s)$  instead of  $\tanh(s)$  as the transformation functions. That is, consider Multi-Layer Perceptrons. In addition, we will take +1 to mean logic TRUE, and -1 to mean logic FALSE. Assume that all  $x_i$  below are either +1 or -1. Write down the weights  $w_i$  for the following perceptron

$$g_A(\mathbf{x}) = \text{sign} \left( \sum_{i=0}^d w_i x_i \right).$$

to implement

$$\text{OR}(x_1, x_2, \dots, x_d).$$

Explain your answer.

10. Continuing from Question 9, among the following choices of  $D$ , write down the smallest  $D$  for some 5- $D$ -1 Neural Network to implement XOR( $(x)_1, (x)_2, (x)_3, (x)_4, (x)_5$ ). Explain your implementation. (It is not so easy to prove the smallest choice, so let's leave the proof for the bonus.)
11. For a Neural Network with at least one hidden layer and  $\tanh(s)$  as the transformation functions on all neurons (including the output neuron), when all the initial weights  $w_{ij}^{(\ell)}$  are set to 0, what gradient components are also 0? Justify your answer.
12. For a Neural Network with one hidden layer and  $\tanh(s)$  as the transformation functions on all neurons (including the output neuron), prove that for the backprop algorithm (with gradient descent), when all the initial weights  $w_{ij}^{(\ell)}$  are set to 1, then  $w_{ij}^{(1)} = w_{i(j+1)}^{(1)}$  for all  $i$  and  $1 \leq j < d^{(1)}$ .

### Experiments with Unpruned Decision Tree

Implement the simple C&RT algorithm without pruning using the Gini index as the impurity measure as introduced in the class. For the decision stump used in branching, if you are branching with feature  $i$  and direction  $s$ , please sort all the  $x_{n,i}$  values to form (at most)  $N + 1$  segments of equivalent  $\theta$ , and then pick  $\theta$  within the median of the segment.

Run the algorithm on the following set for training:

[http://www.csie.ntu.edu.tw/~htlin/course/ml15fall/hw7/hw7\\_train.dat](http://www.csie.ntu.edu.tw/~htlin/course/ml15fall/hw7/hw7_train.dat)

and the following set for testing:

[http://www.csie.ntu.edu.tw/~htlin/course/ml15fall/hw7/hw7\\_test.dat](http://www.csie.ntu.edu.tw/~htlin/course/ml15fall/hw7/hw7_test.dat)

13. (\*) Draw the resulting tree for the TA.
14. (\*) Continuing from Question 13, what is the  $E_{\text{in}}$  of the tree?
15. (\*) Continuing from Question 13, what is the  $E_{\text{out}}$  of the tree?

Now implement the Bagging algorithm with  $N' = N$  and couple it with your decision tree above to make a preliminary random forest  $G_{RF}$ . Produce  $T = 30000$  trees with bagging. Compute  $E_{\text{in}}$  and  $E_{\text{out}}$  using the 0/1 error.

16. (\*) Plot a histogram of  $E_{\text{in}}(g_t)$  over the 30000 trees.
17. (\*) Let  $G_t$  = “the random forest with the first  $t$  trees”. Plot a curve of  $t$  versus  $E_{\text{in}}(G_t)$ .
18. (\*) Continuing from Question 17, and plot a curve of  $t$  versus  $E_{\text{out}}(G_t)$ . Briefly compare with the curve in Question 17 and state your findings.
- Now, ‘prune’ your decision tree algorithm by restricting it to have one branch only. That is, the tree is simply a decision stump determined by Gini index. Make a random ‘forest’  $G_{RS}$  with those decision stumps with Bagging like Questions 16-18 with  $T = 30000$ . Compute  $E_{\text{in}}$  and  $E_{\text{out}}$  using the 0/1 error.
19. (\*) Again, let  $G_t$  = “the random forest with the first  $t$  decision stumps”. Plot a curve of  $t$  versus  $E_{\text{in}}(G_t)$ .
20. (\*) Continuing from Question 19, and plot a curve of  $t$  versus  $E_{\text{out}}(G_t)$ . Briefly compare with the curve in Question 19 and state your findings.

## Bonus: Crazy XOR

21. (10%) Continuing from Question 10, prove or disprove that  $D = d$  is the smallest  $D$  that allows for implementing  $\text{XOR}((x)_1, (x)_2, \dots, (x)_d)$  with a  $d-D-1$  feed-forward neural network with  $\text{sign}(s)$  as the transformation function (such a neural network is also called a Linear Threshold Circuit).
22. (10%) Continuing from Question 10, if you are allowed to use  $D$  neurons (including the one for output) to implement  $\text{XOR}((x)_1, (x)_2, \dots, (x)_d)$ , but can connect the neurons in whatever way as long as it is feed-forward (such as connecting the input directly to neurons in other “layers”), what is the smallest  $D$  (that you can find) for implementing the function? Explain your implementation. You can refer to

[http://www.nature.com/nature/journal/v475/n7356/fig\\_tab/nature10262\\_F2.html](http://www.nature.com/nature/journal/v475/n7356/fig_tab/nature10262_F2.html)

for a possible construction using two neurons for  $d = 3$ .