

Homework #6

RELEASE DATE: 12/23/2021

DUE DATE: 01/06/2022, BEFORE 13:00 on Gradescope

QUESTIONS ARE WELCOMED ON THE NTU COOL FORUM.

You will use Gradescope to upload your choices and your scanned/printed solutions. For problems marked with (*), please follow the guidelines on the course website and upload your source code to Gradescope as well. 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 16 problems and a total of 400 points. For each problem, there is one correct choice. If you choose the correct answer, you get 20 points; if you choose an incorrect answer, you get 0 points. For four of the secretly-selected problems, the TAs will grade your detailed solution in terms of the written explanations and/or code based on how logical/clear your solution is. Each of the four problems graded by the TAs counts as additional 20 points (in addition to the correct/incorrect choices you made). In general, each homework (except homework 0) is of a total of 400 points.

Aggregation

1. Consider an aggregation classifier G constructed by uniform blending on 11 classifiers $\{g_t\}_{t=1}^{11}$. That is,

$$G(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^{11} g_t(\mathbf{x}) \right).$$

Assume that each g_t is of test 0/1 error $E_{\text{out}}(g_t) = e_t$. Which of the following is the tightest upper bound of $E_{\text{out}}(G)$? Choose the correct answer; explain your answer.

[a] $\frac{1}{3} \sum_{t=1}^{11} e_t$

[b] $\frac{1}{4} \sum_{t=1}^{11} e_t$

[c] $\frac{1}{6} \sum_{t=1}^{11} e_t$

[d] $\frac{1}{11} \sum_{t=1}^{11} e_t$

[e] $\frac{1}{12} \sum_{t=1}^{11} e_t$

if G 犯錯, 至少有 6 的 classifier 是錯的

2/2, P7: $\text{avg}(E_{\text{out}}(g_t)) \geq E_{\text{out}}(G)$

$$\frac{e_1 + e_2 + e_3 + \dots + e_{11}}{11}$$

g_1, g_2, g_3, g_4, g_5

2. Suppose that each $\mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$ is drawn uniformly from the region

$$\{0 \leq x_1 \leq 3, 0 \leq x_2 \leq 3\}$$

and the target function is $f(\mathbf{x}) = \text{sign}(x_2 - x_1)$. Consider blending the following three hypotheses linearly to approximate the target function.

$$g_1(\mathbf{x}) = \text{sign}(x_1 - 2)$$

$$g_2(\mathbf{x}) = \text{sign}(x_2 - 1)$$

$$g_3(\mathbf{x}) = \text{sign}(x_2 - 2)$$

That is,

$$G(\mathbf{x}) = \text{sign} \left(\sum_{t=1}^3 \alpha_t \cdot g_t(\mathbf{x}) \right)$$

with $\alpha_t \in \mathbb{R}$. What is the smallest possible $E_{\text{out}}(G)$? Choose the correct answer; explain your answer.

(Hint: The "boundary" of G must be a "combination" of the boundaries of g_t)

[a] $\frac{6}{18}$

[b] $\frac{5}{18}$

[c] $\frac{4}{18}$

[d] $\frac{3}{18}$

[e] none of the other choices

3. When talking about non-uniform voting in aggregation, we mentioned that α can be viewed as a weight vector learned from any linear algorithm coupled with the following transform:

$$\phi(\mathbf{x}) = (g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_T(\mathbf{x})).$$

When studying kernel methods, we mentioned that the kernel is simply a computational short-cut for the inner product $(\phi(\mathbf{x}))^T (\phi(\mathbf{x}'))$. In this problem, we mix the two topics together using the decision stumps as our $g_t(\mathbf{x})$.

Assume that the input vectors contain only even integers between (including) $2L$ and $2R$, where $L < R$. Consider the decision stumps $g_{s,i,\theta}(\mathbf{x}) = s \cdot \text{sign}((x_i - \theta))$, where

$$i \in \{1, 2, \dots, d\},$$

d is the finite dimensionality of the input space,

$$s \in \{-1, +1\},$$

θ is an odd integer between $(2L, 2R)$. $\leftarrow -\frac{2R-2L}{2} = 6$

Define $\phi_{ds}(\mathbf{x}) = (g_{+1,1,2L+1}(\mathbf{x}), g_{+1,1,2L+3}(\mathbf{x}), \dots, g_{+1,1,2R-1}(\mathbf{x}), \dots, g_{-1,d,2R-1}(\mathbf{x}))$. What is

$K_{ds}(\mathbf{x}, \mathbf{x}') = (\phi_{ds}(\mathbf{x}))^T (\phi_{ds}(\mathbf{x}'))$? Choose the correct answer; explain your answer.

[a] $2d(R-L) - 2\|\mathbf{x} - \mathbf{x}'\|_1$

[b] $2d(R-L)^2 - 2\|\mathbf{x} - \mathbf{x}'\|_1^2$

[c] $2d(R-L) - 2\|\mathbf{x} - \mathbf{x}'\|_2$

[d] $2d(R-L)^2 - 2\|\mathbf{x} - \mathbf{x}'\|_2^2$

[e] none of the other choices

Adaptive Boosting

4. Consider applying the AdaBoost algorithm on a binary classification data set where 99% of the examples are positive. Because there are so many positive examples, the base algorithm within AdaBoost returns a constant classifier $g_1(\mathbf{x}) = +1$ in the first iteration. Let $u_n^{(2)}$ be the individual example weight of each example in the second iteration. What is

not helpful
25:12
51:31

Choose the correct answer; explain your answer.

[a] 99
[b] 1/99
[c] 1
[d] 100
[e] 1/100

不懂這啥意思

$$\frac{\sum_{n: y_n > 0} u_n^{(2)}}{\sum_{n: y_n < 0} u_n^{(2)}} = u = \frac{1}{N} \times \sqrt{99}$$

Incorrect

Correct: $\frac{1}{N} \times \frac{1}{\sqrt{99}}$

$\epsilon_t = \frac{1}{100}$
 $\diamond_t = \sqrt{\frac{1 - \frac{1}{100}}{\frac{1}{100}}} = \sqrt{\frac{99}{100}} = \sqrt{99}$

$\epsilon_t \leq \frac{1}{2}$

- For the AdaBoost algorithm introduced in Lecture 12, let $G_t(\mathbf{x}) = \text{sign}(\sum_{\tau=1}^t g_\tau(\mathbf{x}))$. How many of the following are guaranteed to be non-increasing from the t -th iteration to the $(t+1)$ -th iteration? Choose the correct answer; explain each non-increasing case within your answer.

Oscar 25:12
not helpful

I believe this is true
My coding

$E_{\text{in}}(G_t)$ to $E_{\text{in}}(G_{t+1})$
 $E_{\text{out}}(G_t)$ to $E_{\text{out}}(G_{t+1})$
 $\sum_{n=1}^N u_n^{(t)}$ to $\sum_{n=1}^N u_n^{(t+1)}$
 $u_n^{(t)}$ to $u_n^{(t+1)}$ when g_t is correct on (\mathbf{x}_n, y_n)
 $u_n^{(t)}$ to $u_n^{(t+1)}$ when g_t is incorrect on (\mathbf{x}_n, y_n)

$\sum_{n=1}^N u_n^{(t)} = u_{\text{correct}}^{(t)} + u_{\text{incorrect}}^{(t)}$

$\diamond_t = \sqrt{\frac{1 - \epsilon_t}{\epsilon_t}}$

$0 < \epsilon_t \leq 1$
僅重後的不正确

[a] 1
[b] 2
[c] 3
[d] 4
[e] 5

- For the AdaBoost algorithm introduced in Lecture 12, let $U_t = \sum_{n=1}^N u_n^{(t)}$. Assume that $0 < \epsilon_t < \frac{1}{2}$ for each hypothesis g_t . What is $\frac{U_{t+1}}{U_t}$? Choose the correct answer; explain your answer.

- [a] $\sqrt{\epsilon_t(1 - \epsilon_t)}$
[b] $2\sqrt{\epsilon_t(1 - \epsilon_t)}$
[c] $\sqrt{\frac{\epsilon_t}{(1 - \epsilon_t)}}$
[d] $\ln \sqrt{\frac{(1 - \epsilon_t)}{\epsilon_t}}$
[e] $\ln \sqrt{\frac{\epsilon_t}{(1 - \epsilon_t)}}$

7. Following the previous two problems, assume that $\epsilon_t \leq \epsilon < \frac{1}{2}$, which of the following is the tightest upper bound on the number of iterations T required to ensure $E_{\text{in}}(G_T) = 0$? Choose the correct answer; explain your answer.

(Hint: use the fact that

$$\sqrt{\epsilon(1-\epsilon)} \leq \frac{1}{2} \exp\left(-2\left(\frac{1}{2} - \epsilon\right)^2\right)$$

for all $0 < \epsilon < \frac{1}{2}$).

- [a] $\frac{\ln N}{2(\frac{1}{2} - \epsilon)}$
 [b] $\frac{\ln N}{2(\frac{1}{2} - \epsilon)^2}$
 [c] $\frac{\ln N}{4(\frac{1}{2} - \epsilon)}$
 [d] $\frac{\ln N}{4(\frac{1}{2} - \epsilon)^2}$
 [e] $\frac{\ln N}{4(\frac{1}{2} - \epsilon)^4}$

15:42 Oscar
 Random Forest = Bagging + Decision Tree

8. Suppose we have a data set of size $N = 1126$, and we use bootstrapping to sample N' examples. What is the minimum N' such that the probability of getting at least one duplicated example (with # copies ≥ 2) is larger than 50%? Choose the correct answer; explain your answer.

- [a] 25
 [b] 30
 [c] 35
 [d] 40
 [e] none of the other choices

$$1 - \frac{1126 \times 1125 \times 1124 \cdots 1}{(1126)^{N'}} = 0.5$$

9. If bootstrapping is used to sample exactly $2N$ examples out of N , what is the probability that an example is *not* sampled when N is very large? Choose the closest answer; explain your answer.

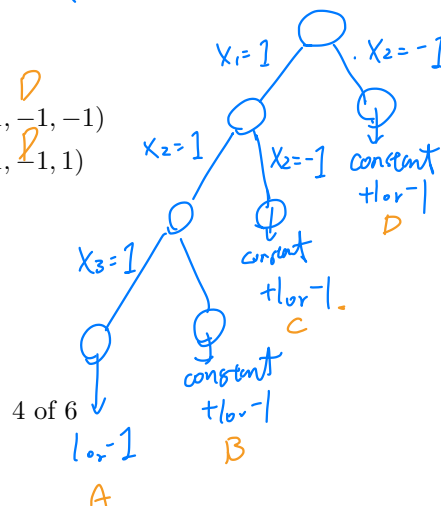
- [a] 77.9%
 [b] 60.7%
 [c] 36.8%
 [d] 13.5%
 [e] 1.8%

2N 項 $\left(\frac{N-1}{N}\right)^{2N} = \left(1 - \frac{1}{N}\right)^{2N}$

$$\lim_{N \rightarrow \infty} \left(1 - \frac{1}{N}\right)^{2N} = \frac{1}{e^2} = 0.135335$$

10. Suppose we have a set of decision trees. Each tree comes with 2 node, each equipped with a fixed branching function. The root node is of two branches, evaluating whether $x_1 \geq 0$. If $x_1 < 0$, the node connects to a leaf with some constant output. Otherwise the node connects to another node of two branches, evaluating whether $x_2 \geq 0$. Each of the branches connects to a constant leaf. Consider three-dimensional input vectors. That is, $\mathbf{x} = (x_1, x_2, x_3)$. Which of the following data set can be shattered by the set of decision trees? Choose the correct answer; explain your answer.

- [a] $(1, 1, -1), (-1, 1, -1), (-1, -1, -1)$
 [b] $(1, 1, -1), (-1, 1, -1), (1, -1, -1)$
 [c] $(1, 1, -1), (-1, 1, -1), (1, -1, -1), (-1, -1, -1)$
 [d] $(1, 1, -1), (-1, 1, -1), (1, -1, -1), (-1, -1, 1)$
 [e] none of the other choices



Experiments with Adaptive Boosting

For Problems 11-16, implement the AdaBoost-Stump algorithm as introduced in Classes 12 and 13. Run the algorithm on the following set for training:

https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw6/hw6_train.dat

and the following set for testing:

https://www.csie.ntu.edu.tw/~htlin/course/ml21fall/hw6/hw6_test.dat

Use a total of $T = 500$ iterations (please do not stop earlier than 500), and calculate E_{in} and E_{out} with the 0/1 error.

For the decision stump algorithm, please implement the following steps. Any ties can be arbitrarily broken.

- 同个feature拿来sort $h_{s,i,\theta}(x) = S \cdot \text{sign}(x_i - \theta)$ ← 第i个feature.
- (1) For any feature i , sort all the $x_{n,i}$ values to $x_{[n],i}$ such that $x_{[n],i} \leq x_{[n+1],i}$.
 - (2) Consider thresholds within $-\infty$ and all the midpoints $\frac{x_{[n],i} + x_{[n+1],i}}{2}$. Test those thresholds with $s \in \{-1, +1\}$ to determine the best (s, θ) combination that minimizes E_{in}^u using feature i .
 - (3) Pick the best (s, i, θ) combination by enumerating over all possible i ← every feature → stump
- For those interested in algorithms (who isn't? :-), step 2 can be carried out in $O(N)$ time only!!

11. (*) What is the value of $E_{in}(g_1)$? Choose the closest answer; provide your code.

[a] 0.29

[b] 0.33

[c] 0.37 = 0.374

[d] 0.41

[e] 0.45

My $g_1 = (s, i, \theta, \alpha_t) = (-1, 9, 0.4482408, 0.25754728)$

12. (*) What is the value of $\max_{1 \leq 500 \leq t} E_{in}(g_t)$? Choose the closest answer; provide your code.

[a] 0.40

[b] 0.45

[c] 0.50

[d] 0.55

[e] 0.60 = 0.591

13. (*) What is the smallest t within the choices below such that $\min_{1 \leq \tau \leq t} E_{in}(G_\tau) \leq 0.05$? Choose the correct answer; provide your code.

[a] 60

[b] 160

[c] 260

[d] 360

[e] 460

14. (*) What is the value of $E_{out}(g_1)$? Choose the closest answer; provide your code.

[a] 0.40

[b] 0.45 = 0.455

[c] 0.50

[d] 0.55

[e] 0.60

13.
15. (*) Define $G_{\text{uniform}}(\mathbf{x}) = \text{sign}\left(\sum_{t=1}^T g_t(\mathbf{x})\right)$. What is the value of $E_{\text{out}}(G_{\text{uniform}})$? Choose the closest answer; provide your code.

- [a] 0.23
[b] 0.28
[c] 0.33
[d] 0.38
[e] 0.43

0.209 ?

16. (*) What is the value of $E_{\text{out}}(G_{600})$? Choose the closest answer; provide your code.

- [a] 0.14
[b] 0.18
[c] 0.22
[d] 0.26
[e] 0.30