
Machine Learning

Answer Sheet for Homework 4

Da-Min HUANG

R04942045

Graduate Institute of Communication Engineering, National Taiwan University

November 7, 2015

Problem 1

Deterministic error is the difference between best $h^* \in H$ and f . If $H' \subset H$, then the complexity of H' is lower than H in general. Hence, in general, the deterministic error increases.

□

Problem 2

1.

$$H(10, 0, 3) \cap H(10, 0, 4) = \left\{ \sum_{i=0}^2 w_q L_q(x) \right\} \cap \left\{ \sum_{i=0}^3 w_q L_q(x) \right\} \quad (1)$$

$$= \left\{ \sum_{i=0}^2 w_q L_q(x) \right\} = H_2 \quad (2)$$

2.

$$H(10, 0, 3) \cup H(10, 1, 4) = \left\{ \sum_{i=0}^2 w_q L_q(x) \right\} \cup \left\{ \sum_{i=0}^3 w_q L_q(x) + \sum_{i=4}^{10} L_q(x) \right\} \quad (3)$$

$$= \left\{ \sum_{i=0}^3 w_q L_q(x) + \sum_{i=4}^{10} L_q(x) \right\} \quad (4)$$

3.

$$H(10, 1, 3) \cap H(10, 1, 4) = \left\{ \sum_{i=0}^2 w_q L_q(x) + \sum_{i=3}^{10} L_q(x) \right\} \cap \left\{ \sum_{i=0}^3 w_q L_q(x) + \sum_{i=4}^{10} L_q(x) \right\} \quad (5)$$

$$= \left\{ \sum_{i=0}^2 w_q L_q(x) + \sum_{i=4}^{10} L_q(x) \right\} \quad (6)$$

4.

$$H(10, 0, 3) \cup H(10, 0, 4) = \left\{ \sum_{i=0}^2 w_q L_q(x) \right\} \cup \left\{ \sum_{i=0}^3 w_q L_q(x) \right\} \quad (7)$$

$$= \left\{ \sum_{i=0}^3 w_q L_q(x) \right\} = H_3 \quad (8)$$

□

Problem 3

We have

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t - \eta \nabla E_{\text{aug}}(\mathbf{w}_t) \quad (9)$$

where $\nabla E_{\text{aug}}(\mathbf{w}_t)$ is

$$\nabla E_{\text{aug}}(\mathbf{w}_t) = \frac{\partial}{\partial \mathbf{w}_t^T} \left(E_{\text{in}}(\mathbf{w}_t) + \frac{\lambda}{N} \mathbf{w}_t^T \mathbf{w}_t \right) = \frac{\partial E_{\text{in}}(\mathbf{w}_t)}{\partial \mathbf{w}_t^T} + \frac{2\lambda}{N} \mathbf{w}_t \quad (10)$$

Hence, we have

$$\mathbf{w}_{t+1} \leftarrow \left(1 - \frac{2\eta\lambda}{N} \right) \mathbf{w}_t - \eta \nabla E_{\text{in}}(\mathbf{w}_t) \quad (11)$$

□

Problem 4

Since \mathbf{w}_{lin} is the optimal solution for the plain-vanilla linear regression, we have

$$E_{\text{in}}(\mathbf{w}_{\text{lin}}) \leq E_{\text{in}}(\mathbf{w}_{\text{reg}}(\lambda)) \quad (12)$$

Also, $\mathbf{w}_{\text{reg}}(\lambda)$ is the optimal solution for $E_{\text{aug}}(\mathbf{w})$, we have

$$E_{\text{aug}}(\mathbf{w}_{\text{reg}}(\lambda)) \leq E_{\text{aug}}(\mathbf{w}_{\text{lin}}) \quad (13)$$

So, we have

$$E_{\text{in}}(\mathbf{w}_{\text{reg}}(\lambda)) + \frac{\lambda}{N} \mathbf{w}_{\text{reg}}^T(\lambda) \mathbf{w}_{\text{reg}}(\lambda) \leq E_{\text{in}}(\mathbf{w}_{\text{lin}}) + \frac{\lambda}{N} \mathbf{w}_{\text{lin}}^T \mathbf{w}_{\text{lin}} \quad (14)$$

$$0 \leq E_{\text{in}}(\mathbf{w}_{\text{reg}}(\lambda)) - E_{\text{in}}(\mathbf{w}_{\text{lin}}) \leq \frac{\lambda}{N} (\|\mathbf{w}_{\text{lin}}\|^2 - \|\mathbf{w}_{\text{reg}}(\lambda)\|^2), \quad \forall \lambda \quad (15)$$

Hence, we have $\|\mathbf{w}_{\text{lin}}\| \geq \|\mathbf{w}_{\text{reg}}(\lambda)\|$ if $\lambda > 0$.

Since this inequality holds for all λ and $\|\mathbf{w}_{\text{lin}}\|$ is not a function of λ . We know that $\|\mathbf{w}_{\text{reg}}(\lambda)\|$ is a non-increasing function of λ for $\lambda \geq 0$.

□

Problem 5

For constant model with three points $A(-1, 0)$, $B(\rho, 1)$ and $C(1, 0)$.

$$\frac{1}{3} \left(\underbrace{\left(0 - \frac{1}{2}\right)^2}_{\text{leave } A} + \underbrace{(1 - 0)^2}_{\text{leave } B} + \underbrace{\left(0 - \frac{1}{2}\right)^2}_{\text{leave } C} \right) = \frac{1}{2} \quad (16)$$

For linear model. Leave A , we get line $y = \frac{1}{\rho - 1}(x - 1)$; leave B , we get line $y = 0$; leave C , we get line $y = \frac{1}{\rho + 1}(x + 1)$. So the error is

$$\frac{1}{3} \left(\left(0 - \left(\frac{-2}{\rho - 1}\right)\right)^2 + (1 - 0)^2 + \left(0 - \frac{2}{\rho + 1}\right)^2 \right) \quad (17)$$

Then we have

$$\frac{1}{3} \left(\frac{4}{\rho^2 - 2\rho + 1} + 1 + \frac{4}{\rho^2 + 2\rho + 1} \right) = \frac{1}{2} \Rightarrow \rho = \pm \sqrt{9 + 4\sqrt{6}} \quad (18)$$

Since $\rho > 0$, we have $\rho = \sqrt{9 + 4\sqrt{6}}$.

□

Problem 6

If the sender wants to make sure at least one person receives correct predictions on all 5 games. Then he should target at least 32 people at first game since there are half the number of people receive wrong prediction after each game and the sender can just ignore people who receives wrong prediction.

The sender sends

$$32 + 16 + 8 + 4 + 2 + 1 = 63 \text{ letters} \quad (19)$$

□

Problem 7

From the conclusion above, we have

$$1000 - 63 \times 10 = 370 \quad (20)$$

□

Problem 8

The mathematical derivations with perfect prediction, which means the hypothesis set shatters the N data. So the size of the hypothesis is 2^N .

□

Problem 9

The Hoeffding bound is

$$2 \exp(-2\epsilon^2 N) = 2 \exp(-20000 \times (0.01)^2) = 2e^{-2} \approx 0.271 \quad (21)$$

□

Problem 10

The computation of Hoeffding bound is computed with data verified by $a(\mathbf{x})$. To improve the performance of $g(\mathbf{x})$, we should only give data sampled by $a(\mathbf{x})$.

Hence, $a(\mathbf{x})$ AND $g(\mathbf{x})$ can improve the system.

□

Problem 11

To get optimal solution, consider the following equation.

$$\nabla \left(\sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \sum_{k=1}^K (\tilde{y}_k - \mathbf{w}^T \tilde{\mathbf{x}}_k)^2 \right) \quad (22)$$

$$= -2 \left(\sum_{n=1}^N \mathbf{x}_n (y_n - \mathbf{x}_n^T \mathbf{w}) + \sum_{k=1}^K \tilde{\mathbf{x}}_k (\tilde{y}_k - \tilde{\mathbf{x}}_k^T \mathbf{w}) \right) \quad (23)$$

$$= -2 \left((\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \mathbf{w}) + (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} - \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \mathbf{w}) \right) \quad (24)$$

$$= \mathbf{0} \quad (25)$$

Hence, we have $\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} (\mathbf{X}^T \mathbf{y} + \tilde{\mathbf{X}}^T \tilde{\mathbf{y}})$.

□

Problem 12

To minimize the formula, consider the following equation.

$$\nabla (\lambda \|\mathbf{w}\|^2 + \|\mathbf{X}\mathbf{w} + \mathbf{y}\|^2) = 2\lambda\mathbf{w} + 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} = \mathbf{0} \quad (26)$$

$$\Rightarrow \mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (27)$$

Hence, we have $\tilde{\mathbf{X}} = \sqrt{\lambda} \mathbf{I}$, $\tilde{\mathbf{y}} = \mathbf{0}$.

□

Problem 13

$E_{\text{in}} = 0.05$ and $E_{\text{out}} = 0.045$.

□

Problem 14

There are three lambda with minimal E_{in} .

1. $\log_{10} \lambda = -10$, $E_{\text{in}} = 0.015$ and $E_{\text{out}} = 0.02$.
2. $\log_{10} \lambda = -9$, $E_{\text{in}} = 0.015$ and $E_{\text{out}} = 0.02$.
3. $\log_{10} \lambda = -8$, $E_{\text{in}} = 0.015$ and $E_{\text{out}} = 0.02$.

□

Problem 15

$\log_{10} \lambda = -7$, $E_{\text{in}} = 0.03$ and $E_{\text{out}} = 0.015$.

□

Problem 16

There are two lambda with minimal E_{in} .

1. $\log_{10} \lambda = -9$, $E_{\text{train}} = 0.0$, $E_{\text{val}} = 0.1$ and $E_{\text{out}} = 0.038$.
2. $\log_{10} \lambda = -8$, $E_{\text{train}} = 0.0$, $E_{\text{val}} = 0.05$ and $E_{\text{out}} = 0.025$.

□

Problem 17

There are eight lambda with minimal E_{in} .

1. $\log_{10} \lambda = -7$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
2. $\log_{10} \lambda = -6$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
3. $\log_{10} \lambda = -5$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
4. $\log_{10} \lambda = -4$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
5. $\log_{10} \lambda = -3$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
6. $\log_{10} \lambda = -2$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.021$.
7. $\log_{10} \lambda = -1$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.022$.
8. $\log_{10} \lambda = 0$, $E_{\text{train}} = 0.033$, $E_{\text{val}} = 0.0375$ and $E_{\text{out}} = 0.028$.

□

Problem 18

□

Problem 19



Problem 20



Reference

- [1] Lecture Notes by Hsuan-Tien LIN, Department of Computer Science and Information Engineering, National Taiwan University, Taipei 106, Taiwan.