

# hw3-Solution

● Graded

Student

Chiang Yi Jie

Total Points

239 / 260 pts

Question 1

Problem 1

20 / 20 pts

✓ + 20 pts Totally correct.

+ 18 pts Some minor errors.

+ 10 pts Deriving the wrong false positive/negative error or the wrong relation between false positive and negative error.

+ 5 pts Wrong answer but have some reasonable efforts.

+ 0 pts Wrong answer

Question 2

Problem 2

20 / 20 pts

✓ + 20 pts Totally correct.

+ 18 pts Small typo.

+ 15 pts Correct, but express the answer in terms other than  $E_{out}(g)$  and  $g$ .

+ 10 pts Correct, but it lacks some explanation for your steps.

+ 5 pts Wrong answer, but you have some reasonable efforts.

+ 0 pts Wrong answer

Question 3

Problem 3

20 / 20 pts

✓ + 20 pts Correct.

+ 15 pts Generally correct, but with minor error.

+ 10 pts On the right track, but with major error.

+ 0 pts Totally Wrong.

- 2 pts Wrong selected page.

Question 4

Problem 4

20 / 20 pts

✓ + 20 pts Correct.

- + 15 pts Generally correct, but with minor error.
- + 10 pts On the right track, but with major error or multiple minor error.
- + 5 pts Incomplete solution. (only derive the gradient or the expectation of  $x$ )
- + 3 pts Just consider gradient.
- + 0 pts Totally Wrong.
- 2 pts Wrong selected page.

Question 5

Problem 5

20 / 20 pts

✓ + 20 pts Correct.

- + 15 pts Generally correct, but with minor error.
- + 10 pts On the right track of derivation, but with major error or multiple minor error.
- + 5 pts Wrong, but with reasonable effort.
- + 0 pts Totally Wrong.
- 2 pts Wrong selected page.

Question 6

Problem 6

20 / 20 pts

✓ - 0 pts Correct

- 20 pts No answer or Incorrect proof.
- 10 pts No Hessian derivative.
- 10 pts No final answer of  $D$
- 4 pts Unclear or incorrect relationship between derivative and  $A_E(w)$
- 4 pts Incorrect comparing a diagonal matrix with a sum form.
- 4 pts Incorrect answer of  $D$
- 2 pts Wrong page.
- 2 pts inappropriate expression or unclear assumption
- 1 pt Miss  $1/N$  in  $D$
- 2 pts Incorrect result of first-order derivative.

### Question 7

#### Problem 7

19 / 20 pts

- 0 pts Correct
- 6 pts No gradient, or the gradient is wrong.
- 6 pts No updating rule.
- 6 pts No or incorrect comparison with PLA.
- 4 pts Incorrect updating rule
- 2 pts Unclear or minor incorrect gradient of error function.
- 2 pts Unclear or minor incorrect updating rule.
- 2 pts Unclear or minor incorrect comparison with PLA.

✓ - 1 pt inappropriate expression or unclear assumption

1  $y_s < 0?$

- 18 pts Major error of the proof.
- 20 pts No solution.
- 2 pts Wrong page.

### Question 8

#### Problem 8

20 / 20 pts

✓ + 20 pts Totally correct.

- + 18 pts Almost correct.
- + 18 pts Small typo.
- + 13 pts Only consider one case of derivatives  
(w.r.t.  $w_k$ ,  $k = y$ ), but neglect the other case that differentiate w.r.t.  $w_k$ ,  $k \neq y$
- + 7 pts Wrong answer but with some reasonable efforts.
- + 0 pts Wrong answer.

Minor typo.

2 minus

Question 9

Problem 9

20 / 20 pts

✓ + 10 pts Correct answer

✓ + 10 pts Correct histogram

+ 10 pts Wrong answer with bell shaped histogram

+ 0 pts No answer / Wrong answer

- 2 pts wrong page selection

Question 10

Problem 10

20 / 20 pts

✓ + 10 pts Correct answer

✓ + 10 pts Correct histogram

+ 10 pts Wrong answer with bell shaped histogram

+ 0 pts No answer / Wrong answer

- 2 pts wrong page selection

+ 5 pts nice try

Question 11

Problem 11

20 / 20 pts

✓ + 10 pts Correct answer

✓ + 10 pts Correct scatter plot

+ 10 pts Wrong answer with reasonable scatter plot

+ 0 pts No answer / Wrong answer

- 2 pts wrong page selection

Question 12

Problem 12

20 / 20 pts

✓ + 10 pts Correct answer (within error range)

✓ + 5 pts Correct scatter plot

✓ + 5 pts Reasonable findings and explanations

+ 10 pts Wrong answer with reasonable scatter plot and findings

+ 0 pts No answer / Wrong answer

Question 13

Problem 13

0 / 20 pts

– 0 pts Totally correct, no logical flaw.

– 5 pts A minor logical flaw exists.

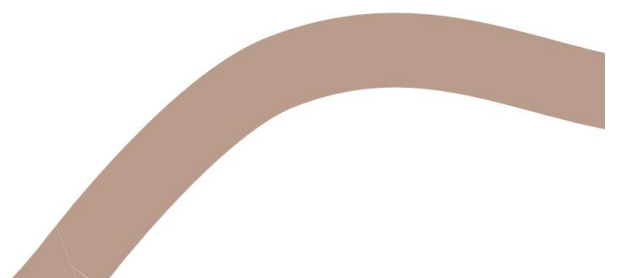
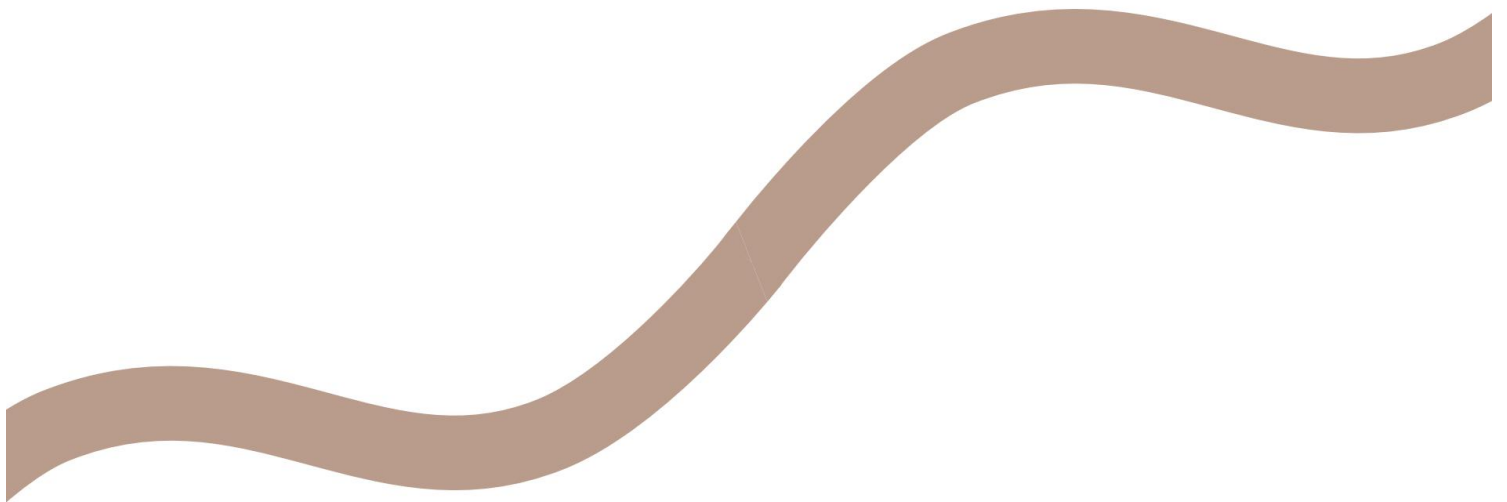
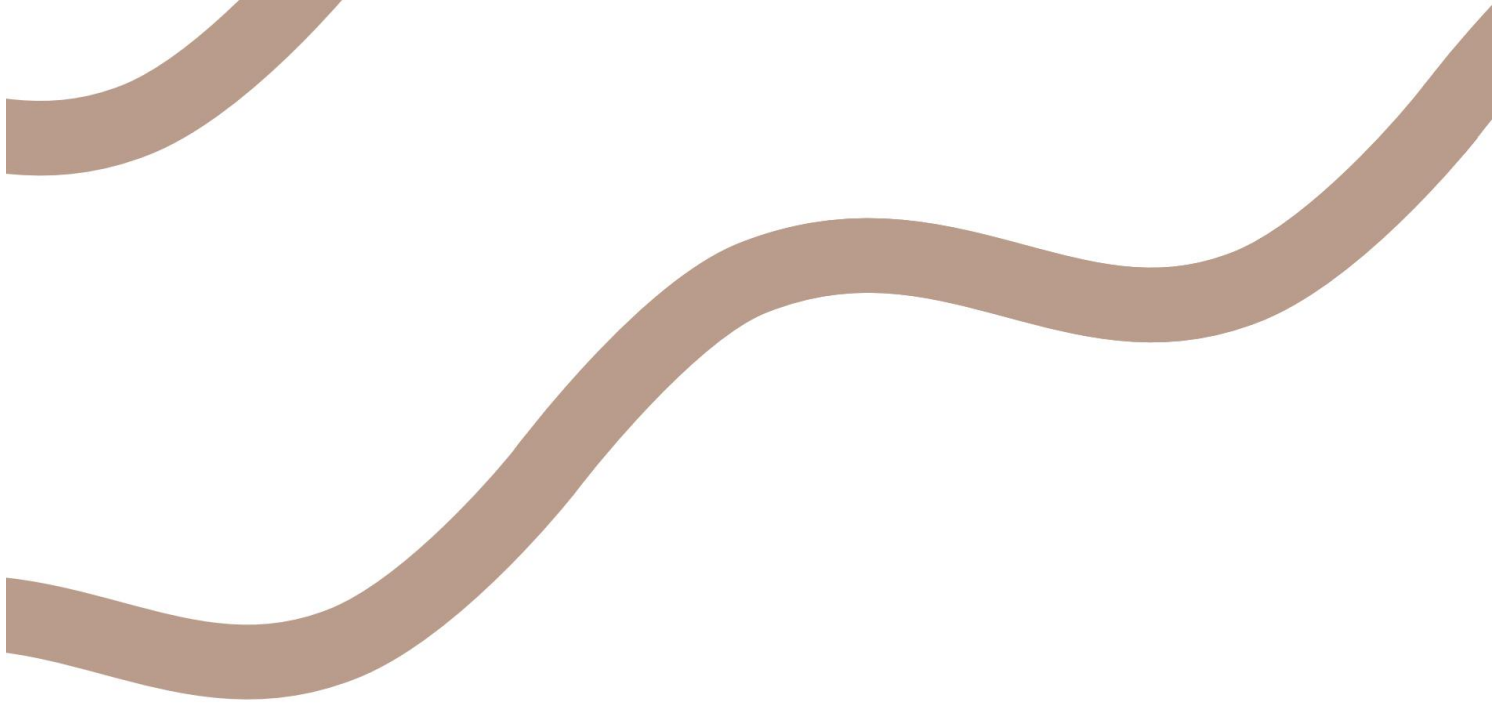
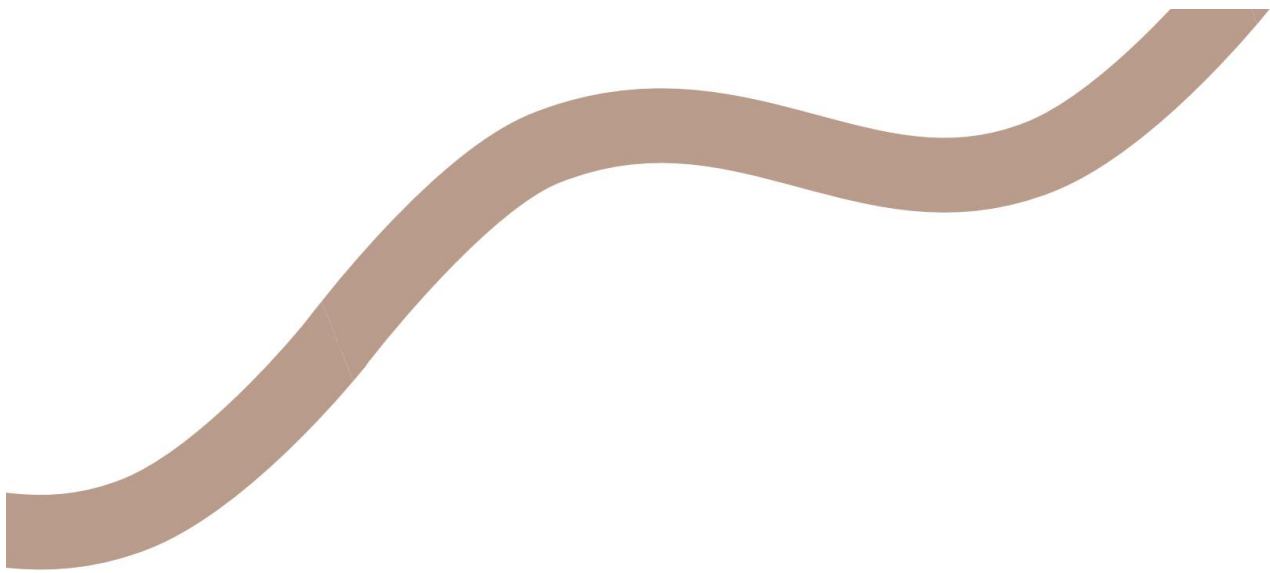
– 10 pts A logical flaw exists.

– 15 pts A major logical flaw exists.

✓ – 20 pts Wrong answer

3 yx, not x.

No questions assigned to the following page.



Question assigned to the following page: [1](#)



1. (20 points) Consider a binary classification problem, where  $\mathcal{Y} = \{-1, +1\}$ . Assume a noisy scenario where the data is generated i.i.d. from some  $P(\mathbf{x}, y)$ . In class, we discussed that when the 0/1 error function (i.e. classification error) is considered, calculating the "ideal mini-target" on each  $\mathbf{x}$  reveals the hidden target function of

$$f_{0/1}(\mathbf{x}) = \operatorname{argmax}_{y \in \{-1, +1\}} P(y|\mathbf{x}) = \operatorname{sign}\left(P(+1|\mathbf{x}) - \frac{1}{2}\right).$$

Instead of the 0/1 error, if we consider the CIA error function, where a false positive (classifying a negative example as a positive one) is 1000 times more important than a false negative, the hidden target should be changed to

$$f_{\text{CIA}}(\mathbf{x}) = \operatorname{sign}(P(+1|\mathbf{x}) - \alpha).$$

Prove what the value of  $\alpha$  should be.

$$f_{0/1}(\mathbf{x}) = \operatorname{sign}\left(P(+1|\mathbf{x}) - \frac{1}{2}\right), \quad \begin{aligned} f(\mathbf{x}) &= 1 \text{ if } P(+1|\mathbf{x}) > 0.5 \\ f(\mathbf{x}) &= -1 \text{ if } P(+1|\mathbf{x}) < 0.5 \end{aligned} \quad \text{y}$$

	$f(\mathbf{x})$	
	+1	-1
+1	0	1
-1	1	0

$f_{0/1}(\mathbf{x}_n) \neq y_n$ ,  $y_n = -1$  is 1000 times important than  $y_n = +1$

↓  
In Error measure, when we enter the point  $(y_n, f(\mathbf{x}_n)) = (-1, +1)$  once,  
equals to we enter the point 1000 times.

	$f(\mathbf{x})$	
	+1	-1
+1	0	1
-1	1000	0

$$\rightarrow \text{Err} = P(+1|\mathbf{x}) \times \mathbb{I}[f(\mathbf{x}) = -1] + P(-1|\mathbf{x}) \times \mathbb{I}[f(\mathbf{x}) = +1] \times 1000$$

$\rightarrow P(+1|\mathbf{x}) : P(-1|\mathbf{x}) = 1000:1$   
 $P(+1|\mathbf{x}) = \frac{1000}{1001}, \quad \alpha = \frac{1000}{1001}$



Question assigned to the following page: [2](#)

2. (20 points) Consider a binary classification task, where God gives you some noiseless data i.i.d. from an unknown distribution  $P(\mathbf{x})$  and an unknown target function  $f(\mathbf{x})$  that maps from  $\mathcal{X}$  to  $\{-1, +1\}$ . After you use the data to obtain some  $g(\mathbf{x})$  that suffers

$$\begin{aligned} E_{\text{out}}(g) &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} [g(\mathbf{x}) \neq f(\mathbf{x})] \text{ (here } \mathbb{E} \text{ means expectation, as shown in class slides)} \\ &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} [g(\mathbf{x}) \neq f(\mathbf{x})] \text{ (or if you like the more beautiful font } \mathbb{E} \text{ for expectation).} \end{aligned}$$

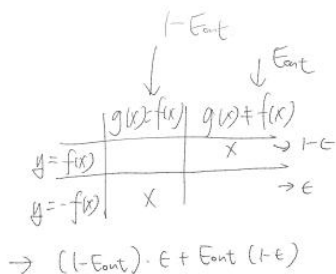
Now, assume that  $g(\mathbf{x})$  is put in a noisy test environment where

$$\begin{aligned} P(y = +f(\mathbf{x}) | \mathbf{x}) &= 1 - \epsilon \\ P(y = -f(\mathbf{x}) | \mathbf{x}) &= \epsilon. \end{aligned}$$

Derive

$$\mathbb{E}_{(\mathbf{x}, y) \sim P(\mathbf{x}, y)} [g(\mathbf{x}) \neq y]$$

as a function of  $E_{\text{out}}(g)$  and  $\epsilon$ .



In ideal case  $P[g(\mathbf{x}) \neq f(\mathbf{x})] = E_{\text{out}}$   
 $P[g(\mathbf{x}) = f(\mathbf{x})] = 1 - E_{\text{out}}$

In noise  $P[y = f(\mathbf{x}) | \mathbf{x}] = 1 - \epsilon$   
 $P[y = -f(\mathbf{x}) | \mathbf{x}] = \epsilon$

	$g(\mathbf{x}) = f(\mathbf{x})$	$g(\mathbf{x}) \neq f(\mathbf{x})$	
$y = f(\mathbf{x})$	$s_1$	$s_2$	$1 - \epsilon$
$y = -f(\mathbf{x})$	$s_3$	$s_4$	$\epsilon$
$P_2$	$1 - E_{\text{out}}$	$E_{\text{out}}$	

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, y) \sim P(\mathbf{x}, y)} [g(\mathbf{x}) \neq y] &= \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} [g(\mathbf{x}) \neq f(\mathbf{x})] \cdot P(y = +f(\mathbf{x}) | \mathbf{x}) + \\ &\quad \mathbb{E}_{\mathbf{x} \sim P(\mathbf{x})} [g(\mathbf{x}) = f(\mathbf{x})] \cdot P(y = -f(\mathbf{x}) | \mathbf{x}) \\ &= (1 - E_{\text{out}}) \cdot \epsilon + E_{\text{out}} \cdot (1 - \epsilon) \end{aligned}$$



Question assigned to the following page: [3](#)

3. (20 points) Consider a hypothesis set that contains hypotheses of the form  $h(x) = wx$  for  $x \in \mathbb{R}$ . Combine the hypothesis set with the squared error function to minimize

$$E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N (h(x_n) - y_n)^2$$

on a given data set  $\{(x_n, y_n)\}_{n=1}^N$ . Derive the optimal  $w_{\text{LIN}}$  in terms of  $(x_n, y_n)$  and express the result *without* using matrix/vector notations. You can assume all denominators to be non-zero.

(Hint: This is linear regression in  $\mathbb{R}$  without the added  $x_0$ .)  $\frac{\partial}{\partial w} = 2 \left[ \frac{\partial}{\partial w} (w^2 X^2 - 2wXy + (b - w_0)^2) \right] (1) \quad \rightarrow (w) =$

$$h(x) = wx$$

$$E_{\text{in}}(w) = \frac{1}{N} \sum_{n=1}^N (wx_n - y_n)^2$$

$$w_{\text{LIN}} \text{ exists when } \min E_{\text{in}}(w) \text{ appears } \rightarrow \nabla E_{\text{in}}(w) = 0$$

$$E_{\text{in}}(w) = \frac{1}{N} \left[ (wx_1 - y_1)^2 + (wx_2 - y_2)^2 + (wx_3 - y_3)^2 \dots + (wx_n - y_n)^2 \right]$$

$$= \frac{1}{N} \left( w^2 x_1^2 - 2wx_1 y_1 + y_1^2 \dots + w^2 x_n^2 - 2wx_n y_n + y_n^2 \right)$$

$$\nabla E_{\text{in}}(w) = \frac{1}{N} (2wx_1^2 - 2x_1 y_1 + 2wx_2^2 - 2x_2 y_2 \dots + 2wx_n^2 - 2x_n y_n)$$

$$= \frac{2}{N} \left( w \sum_{\tilde{n}=1}^N x_{\tilde{n}}^2 - \sum_{\tilde{n}=1}^N x_{\tilde{n}} y_{\tilde{n}} \right) = 0$$

$$\rightarrow w = \frac{\sum_{\tilde{n}=1}^N x_{\tilde{n}} y_{\tilde{n}}}{\sum_{\tilde{n}=1}^N x_{\tilde{n}}^2} \quad \#$$





Question assigned to the following page: [4](#)

4. (20 points) Consider the target function  $f(x) = ax^2 + b$ . Sample  $x$  uniformly from  $[0, 1]$ , and use all linear hypotheses  $h(x) = w_0 + w_1 \cdot x$  to approximate the target function with respect to the squared error. For any given  $(a, b)$ , derive the weights  $(w_0^*, w_1^*)$  of the optimal hypothesis as a function of  $(a, b)$ .

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (f(x_n) - h(x_n))^2 = \frac{1}{N} \sum_{n=1}^N (ax_n^2 + b - w_1 x_n - w_0)^2$$

since  $x$  uniformly from  $[0, 1]$ ,  $E_{in}(w_0, w_1) = \int_0^1 (ax^2 - w_1 x + b - w_0)^2 dx$

$$\min E_{in}(w_0, w_1) \Leftrightarrow \nabla E_{in}(w_0, w_1) = 0$$

$$\begin{cases} \frac{\partial}{\partial w_0} = \int_0^1 (ax^2 - w_1 x + b - w_0) \cdot 2 \cdot (-1) dx = 0 \\ \frac{\partial}{\partial w_1} = \int_0^1 (ax^2 - w_1 x + b - w_0) \cdot 2 \cdot (-x) dx = 0 \end{cases}$$

$$\begin{cases} \int_0^1 (ax^2 - w_1 x + b - w_0) dx = 0 \\ \int_0^1 (ax^3 - w_1 x^2 + bx - w_0 x) dx = 0 \end{cases}$$

$$\begin{cases} \left. \frac{1}{3} ax^3 - \frac{1}{2} w_1 x^2 + (b - w_0)x \right|_0^1 = 0 \\ \left. \frac{1}{4} ax^4 - \frac{1}{3} w_1 x^3 + \frac{1}{2} (b - w_0)x^2 \right|_0^1 = 0 \end{cases}$$

$$\begin{cases} \frac{1}{3} a - \frac{1}{2} w_1 + b - w_0 = 0 \\ \frac{1}{4} a - \frac{1}{3} w_1 + \frac{1}{2} b - \frac{1}{2} w_0 = 0 \end{cases}$$

$$\begin{cases} w_0 + \frac{1}{2} w_1 = \frac{1}{3} a + b \\ \frac{1}{2} w_0 + \frac{1}{3} w_1 = \frac{1}{4} a + \frac{1}{2} b, \quad w_0 + \frac{2}{3} w_1 = \frac{1}{2} a + b \end{cases}$$

$$\rightarrow \frac{1}{6} w_1 = \frac{1}{6} a, \quad w_1^* = a$$

$$w_0^* = -\frac{1}{6} a + b$$



Question assigned to the following page: [5](#)

5. (20 points) Consider running linear regression on  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , where  $\mathbf{x}_n$  includes the constant dimension  $x_0 = 1$  as usual. For simplicity, you can assume that  $\mathbf{X}^T \mathbf{X}$  is invertible. Assume that the unique (why :-)) solution  $\mathbf{w}_{\text{LIN}}$  is obtained after running linear regression on the data above. Now, consider an output transformation

$$W_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$y'_n = ay_n + b.$$

for some given constants  $(a, b)$ . Run linear regression on  $\{(\mathbf{x}_n, y'_n)\}_{n=1}^N$  to obtain the unique solution  $\mathbf{w}'_{\text{LIN}}$ . Derive  $\mathbf{w}'_{\text{LIN}}$  as a function of  $\mathbf{w}_{\text{LIN}}$  and  $(a, b)$ .

$$\mathbf{w}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{--- } \textcircled{0}$$

$$E_{\text{in}}(\mathbf{w}') = \frac{1}{N} \|\mathbf{X} \mathbf{w}' - (a\mathbf{y} + b)\|^2 = \frac{1}{N} (\mathbf{w}'^T \mathbf{X}^T \mathbf{X} \mathbf{w}' - 2 \mathbf{w}'^T \mathbf{X}^T (a\mathbf{y} + b) + (a\mathbf{y} + b)^T (a\mathbf{y} + b))$$

$$\nabla E_{\text{in}}(\mathbf{w}') = \frac{1}{N} (2 \mathbf{X}^T \mathbf{X} \mathbf{w}' - 2 \mathbf{X}^T (a\mathbf{y} + b) + 0)$$

$$\mathbf{w}_{\text{LIN}}' \text{ happens when } \nabla E_{\text{in}} = 0 : 2 \mathbf{X}^T \mathbf{X} \mathbf{w}' = 2 \mathbf{X}^T (a\mathbf{y} + b) \quad , \quad \mathbf{w}_{\text{LIN}}' = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (a\mathbf{y} + b) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot a\mathbf{y} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot b$$

$$\begin{cases} a(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = a \cdot \mathbf{w}_{\text{LIN}} \\ (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \cdot b \end{cases}$$

$$\text{since } \mathbf{X} \cdot a \mathbf{w}_{\text{LIN}} = a\mathbf{y} \quad \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix} \times a \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = a \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$$\mathbf{X} \cdot \mathbf{w}_{\text{LIN}}' = a\mathbf{y} + b \quad \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n3} & \dots & x_{nn} \end{bmatrix} \times \begin{bmatrix} aw_1 + b \\ aw_2 \\ aw_3 \\ \vdots \\ aw_n \end{bmatrix} = \begin{bmatrix} ay_1 + b \\ ay_2 + b \\ \vdots \\ ay_n + b \end{bmatrix}$$

all b's comes from  $(\mathbf{w}_{\text{LIN}}')$

$$\downarrow$$

$$a \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} + b \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\rightarrow \mathbf{w}_{\text{LIN}}' = a \mathbf{w}_{\text{LIN}} + b \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad , \quad \text{where } \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ is a } n \times 1 \text{ matrix. } \#$$



Question assigned to the following page: [6](#)

6. (20 points) Let  $E(\mathbf{w}): \mathbb{R}^d \rightarrow \mathbb{R}$  be a function. Denote the gradient  $\mathbf{b}_E(\mathbf{w})$  and the Hessian  $A_E(\mathbf{w})$  by

$$\mathbf{b}_E(\mathbf{w}) = \nabla E(\mathbf{w}) = \begin{bmatrix} \frac{\partial E}{\partial w_1}(\mathbf{w}) \\ \frac{\partial E}{\partial w_2}(\mathbf{w}) \\ \vdots \\ \frac{\partial E}{\partial w_d}(\mathbf{w}) \end{bmatrix}_{d \times 1} \quad \text{and} \quad A_E(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 E}{\partial w_1^2}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_1 \partial w_2}(\mathbf{w}) & \dots & \frac{\partial^2 E}{\partial w_1 \partial w_d}(\mathbf{w}) \\ \frac{\partial^2 E}{\partial w_2 \partial w_1}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_2^2}(\mathbf{w}) & \dots & \frac{\partial^2 E}{\partial w_2 \partial w_d}(\mathbf{w}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 E}{\partial w_d \partial w_1}(\mathbf{w}) & \frac{\partial^2 E}{\partial w_d \partial w_2}(\mathbf{w}) & \dots & \frac{\partial^2 E}{\partial w_d^2}(\mathbf{w}) \end{bmatrix}_{d \times d}.$$

Then, the second-order Taylor's expansion of  $E(\mathbf{w})$  around  $\mathbf{u}$  is:

$$E(\mathbf{w}) \approx E(\mathbf{u}) + \mathbf{b}_E(\mathbf{u})^T (\mathbf{w} - \mathbf{u}) + \frac{1}{2} (\mathbf{w} - \mathbf{u})^T A_E(\mathbf{u}) (\mathbf{w} - \mathbf{u}).$$

$$\mathbf{b}_E(\mathbf{u})^T (\mathbf{w} - \mathbf{u}) = -(\mathbf{w} - \mathbf{u})^T A_E(\mathbf{u}) (\mathbf{w} - \mathbf{u})$$

Suppose  $A_E(\mathbf{u})$  is positive definite. The optimal direction  $\mathbf{v}$  such that  $\mathbf{w} \leftarrow \mathbf{u} + \mathbf{v}$  minimizes the right-hand-side of the Taylor's expansion above is simply  $-(A_E(\mathbf{u}))^{-1} \mathbf{b}_E(\mathbf{u})$ .

*Hint: Homework 0! :-)*

Now, consider minimizing  $E_{\text{in}}(\mathbf{w})$  in logistic regression problem with Newton's method on a data set  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$  with the cross-entropy error function for  $E_{\text{in}}$ :

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n \mathbf{w}^T \mathbf{x}_n)).$$

For any given  $\mathbf{w}_t$ , let

$$h_t(\mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}_t^T \mathbf{x})}.$$

Express the Hessian  $A_E(\mathbf{w}_t)$  with  $E = E_{\text{in}}$  as  $\mathbf{X}^T \mathbf{D} \mathbf{X}$ , where  $\mathbf{D}$  is an  $N$  by  $N$  diagonal matrix. Derive what  $\mathbf{D}$  should be in terms of  $h_t$ ,  $\mathbf{w}_t$ ,  $\mathbf{x}_n$ , and  $y_n$ .

$$h_t(\mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}_t^T \mathbf{x})}$$

$$E_{\text{in}} = \frac{1}{N} \sum_{n=1}^N \ln(1 + \exp(-y_n \mathbf{w}_t^T \mathbf{x}_n))$$

$$h_t(-\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}_t^T \mathbf{x})} = 1 - h_t(\mathbf{x})$$

$$\mathbf{b}_E(\mathbf{w}) = \nabla E_{\text{in}} = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n)}{1 + \exp(-y_n \mathbf{x}_n)} (-y_n \mathbf{x}_n, \tilde{y}) = \frac{1}{N} \sum_{n=1}^N (1 - h_t(y_n \mathbf{x}_n))$$

$$\begin{aligned} A_E(\mathbf{w}) &= \nabla \mathbf{b}_E(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n) (-y_n \mathbf{x}_{n,j}) (1 + \exp(-y_n \mathbf{x}_n)) - \exp^2(-y_n \mathbf{x}_n) (-y_n \mathbf{x}_{n,j})}{(1 + \exp(-y_n \mathbf{x}_n))^2} \cdot (-y_n \mathbf{x}_{n,i}) \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n) (-y_n \mathbf{x}_{n,j}) (-y_n \mathbf{x}_{n,i}) [1 + \exp(-y_n \mathbf{x}_n) - \exp(-y_n \mathbf{x}_n)]}{(1 + \exp(-y_n \mathbf{x}_n))^2} \\ &= \frac{1}{N} \sum_{n=1}^N \frac{\exp(-y_n \mathbf{x}_n) (y_n \mathbf{x}_{n,j}) (y_n \mathbf{x}_{n,i})}{(1 + \exp(-y_n \mathbf{x}_n))^2} = \mathbf{X}^T \mathbf{D} \mathbf{X}. \end{aligned}$$

$$D_{i,j} = \frac{1}{N} h_t(\mathbf{x}_n y_n) (1 - h_t(\mathbf{x}_n y_n)) \cdot y_{n,i} \cdot y_{n,j}$$

$$= \frac{1}{N} (1 - h_t(\mathbf{x}_n)) h_t(\mathbf{x}_n) y_{n,i} \cdot y_{n,j}$$

$$\text{since } y_{n,i} \cdot y_{n,j} = 1 \quad (\mathbf{D} \text{ diagonal}), \quad D_{i,j} = \frac{1}{N} \begin{bmatrix} h_t(x_1)(1-h_t(x_1)) & & & \\ & h_t(x_2)(1-h_t(x_2)) & & \\ & & h_t(x_3)(1-h_t(x_3)) & \\ & & & \ddots \end{bmatrix}$$





Question assigned to the following page: [Z](#)

7. (20 points) The truncated squared loss

$$\text{err}(s, y) = (\max(0, 1 - ys))^2$$

can be easily shown to be an upper bound on the 0/1 error. Assume that  $s$  is generated from a linear scoring function  $s = \mathbf{w}^T \mathbf{x}$  like Page 3/25 of Lecture 11. Derive a "perceptron learning algorithm" by applying SGD on the truncated squared loss. Compare the resulting algorithm with the original PLA. Discuss the similarities and differences using 5 to 10 sentences.

PLA:  $\text{error}_{0/1}(s, y) = [\text{sign}(s) \neq y]$   
 $\mathbf{w}_{t+1} = \mathbf{w}_t + 1 \cdot [y_n \neq \text{sign}(s)] (\mathbf{y}_n \cdot \mathbf{x}_n)$

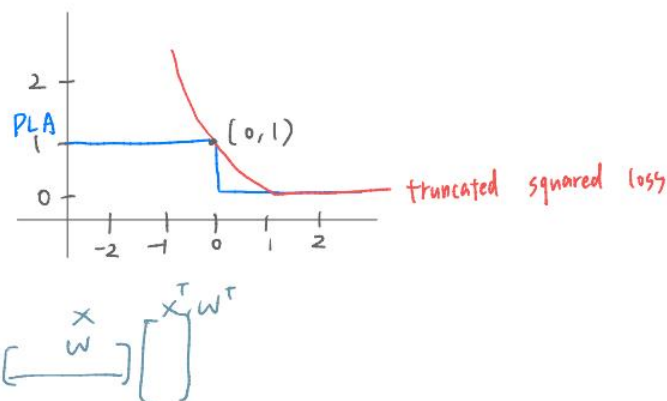
SGD:  $\text{error}(s, y) = (\max(0, 1 - ys))^2$   
 $\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \cdot (-\nabla \text{error}(y, \mathbf{w}^T \mathbf{x}))$   
 $= \mathbf{w}_t + \eta \cdot (-\nabla (\max(0, 1 - ys))^2)$

• If  $ys > 0$ ,  $\text{error}(s, y) = 0$   
 $0 < ys < 1$ ,  $\text{error}(s, y) > \text{error}_{0/1}(s, y)$

$ys = -1$ ,  $\text{error}(s, y) = \text{error}_{0/1}(s, y)$

①  $(\nabla \max(0, 1 - ys))^2 = \nabla (1 - \mathbf{w}^T \mathbf{x} y_n)^2 = \nabla (\mathbf{w}^T \mathbf{x} y_n - 1)^2$   
 $= \nabla (\mathbf{x} \cdot \mathbf{w}^T \mathbf{w} \cdot \mathbf{x}^T - 2 \mathbf{x}^T \mathbf{w} y_n + 1)^2$   
 $= 2 (\mathbf{x}^T \mathbf{x} \cdot \mathbf{w} - y_n \mathbf{x}) = 2 (1 - ys) \mathbf{x} y_n$

$\mathbf{w}_{t+1} = \mathbf{w}_t + 2\eta (1 - ys) \mathbf{x}_n y_n$



PLA always updates with  $\eta = 1$  when  $\mathbb{I} y_n \neq \text{sign}(s) \mathbb{I}$ ; while SGD updates with altered  $\eta$  when  $\max(0, 1 - ys) > 0$ . When  $ys > 1$ , both PLA and SGD won't update ( $\text{error} = 0$ ).  
 when  $0 < ys < 1$ , SGD updates (PLA doesn't), when  $ys < 0$ ,  $\mathbf{w}_{t+1}^{\text{PLA}} = \mathbf{w}_t^{\text{PLA}} + \mathbf{x}_n y_n$ ,  
 $\mathbf{w}_{t+1}^{\text{SGD}} = \mathbf{w}_t^{\text{SGD}} + 2\eta (1 - ys) \mathbf{x}_n y_n$  (updates depend on  $\mathbf{x}_n y_n$ , alterable).

If  $ys \ll 0$ , updates larger;  $ys < 0$ , updates less.

However, eventually how much to update depend on  $\eta$ , which is usually smaller than 0.1

Consequently, SGD usually updates with flexible and smaller changes, PLA offers constant and unstable steps. In error measure,  $\text{Error}(\text{SGD}) > \text{Error}(\text{PLA})$ , that is, SGD can also be viewed as upper bound of PLA.



Question assigned to the following page: [8](#)

## Multinomial Logistic Regression

8. In Lecture 11, we solve multiclass classification by OVA or OVO decompositions. One alternative to deal with multiclass classification is to extend the original logistic regression model to Multinomial Logistic Regression (MLR). For a  $K$ -class classification problem, we will denote the output space  $\mathcal{Y} = \{1, 2, \dots, K\}$ . The hypotheses considered by MLR can be indexed by a matrix

$$W = \begin{bmatrix} | & | & \dots & | & \dots & | \\ \mathbf{w}_1 & \mathbf{w}_2 & \dots & \mathbf{w}_k & \dots & \mathbf{w}_K \\ | & | & \dots & | & \dots & | \end{bmatrix}_{(d+1) \times K},$$

that contains weight vectors  $(\mathbf{w}_1, \dots, \mathbf{w}_K)$ , each of length  $d+1$ . The matrix represents a hypothesis

$$h_y(\mathbf{x}) = \frac{\exp(\mathbf{w}_y^T \mathbf{x})}{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x})}$$

that can be used to approximate the target distribution  $P(y|\mathbf{x})$  for any  $(\mathbf{x}, y)$ . MLR then seeks for the maximum likelihood solution over all such hypotheses. For a given data set  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$  generated i.i.d. from some  $P(\mathbf{x})$  and target distribution  $P(y|\mathbf{x})$ , the likelihood of  $h_y(\mathbf{x})$  is proportional to  $\prod_{n=1}^N h_{y_n}(\mathbf{x}_n)$ . That is, minimizing the negative log likelihood is equivalent to minimizing an  $E_{\text{in}}(W)$  that is composed of the following error function

$$\text{err}(W, \mathbf{x}, y) = -\ln h_y(\mathbf{x}) = -\sum_{k=1}^K \mathbb{I}[y = k] \ln h_k(\mathbf{x}).$$

Consider minimizing  $E_{\text{in}}(W) = \frac{1}{N} \sum_{n=1}^N \text{err}(W, \mathbf{x}_n, y_n)$  with gradient descent. Derive  $\nabla E_{\text{in}}(W)$ . Your result should simply be a matrix with the same size as  $W$ . (Note: the hypothesis that transforms the scores  $\{\mathbf{w}_i^T \mathbf{x}\}_{i=1}^K$  to  $h_y(\mathbf{x})$  is often called a softmax function in (multiclass) deep learning.)

$$\begin{aligned} E_{\text{in}}(W) &= \frac{1}{N} \sum_{n=1}^N \ln h_{y_n}(\mathbf{x}_n) \\ \nabla E_{\text{in}}(W) &= -\frac{1}{N} \sum_{n=1}^N \frac{\partial}{\partial W} \ln \frac{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)}{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)} = -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)}{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)} \times \frac{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n) \sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n) - \exp(\mathbf{w}_{y_n}^T \mathbf{x}_n) \exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)}{(\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n))^2} \\ &= -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)}{(\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n))^2} \times \frac{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n) \sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n) - \exp(\mathbf{w}_{y_n}^T \mathbf{x}_n) \exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)}{\exp(\mathbf{w}_{y_n}^T \mathbf{x}_n)} \\ &= -\frac{1}{N} \sum_{n=1}^N \frac{(\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n) - \exp(\mathbf{w}_{y_n}^T \mathbf{x}_n))}{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)} \cdot \mathbf{x}_n = \frac{1}{N} \sum_{n=1}^N (h_{y_n}(\mathbf{x}_n) - \mathbb{I}[y_n = k]) \mathbf{x}_n \end{aligned}$$

$$\text{if } j \neq k, -\frac{1}{N} \sum_{n=1}^N \frac{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)}{\exp(\mathbf{w}_j^T \mathbf{x}_n)} \times \frac{-\exp(\mathbf{w}_j^T \mathbf{x}_n) \exp(\mathbf{w}_k^T \mathbf{x}_n)}{\sum_{i=1}^K \exp(\mathbf{w}_i^T \mathbf{x}_n)^2} \mathbf{x}_n = \frac{1}{N} \sum_{n=1}^N h_{y_n}(\mathbf{x}_n) \cdot \mathbf{x}_n$$

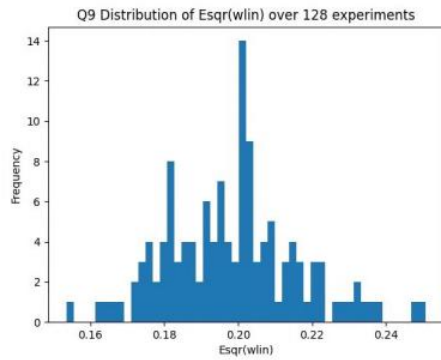
$$\Rightarrow \nabla E_{\text{in}}(W) = -\frac{1}{N} \begin{bmatrix} \sum_{n=1}^N (h_1(\mathbf{x}_n) - \mathbb{I}[y=1]) \mathbf{x}_n & \sum_{n=1}^N (h_2(\mathbf{x}_n) - \mathbb{I}[y=2]) \mathbf{x}_n & \dots \end{bmatrix}$$



Question assigned to the following page: [9](#)



9. (20 points, \*) Implement the linear regression algorithm taught in the lecture. Run the algorithm for 128 times, each with a different random seed for generating the two data sets above. Plot a histogram to visualize the distribution of  $E_{\text{in}}^{\text{sqr}}(\mathbf{w}_{\text{LIN}})$ , where  $E_{\text{in}}^{\text{sqr}}$  denotes the *averaged* squared error over  $N$  examples. What is the median  $E_{\text{in}}^{\text{sqr}}$  over the 128 experiments?



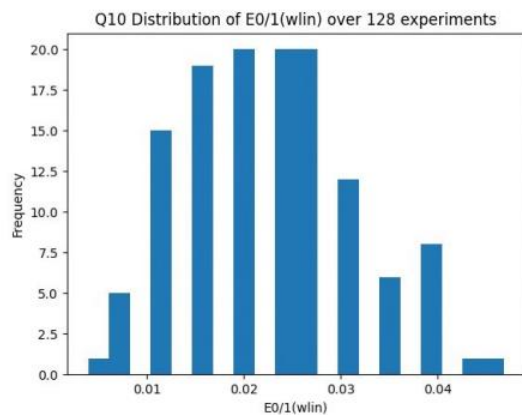
median = 0.199638 #



Question assigned to the following page: [10](#)

10. (20 points, \*) Following the previous problem, plot a histogram to visualize the distribution of  $E_{\text{in}}^{0/1}(\mathbf{w}_{\text{LIN}})$ , where  $E_{\text{in}}^{0/1}$  denotes the averaged 0/1 error over  $N$  examples (i.e. using  $\mathbf{w}_{\text{LIN}}$  for binary classification). What is the median  $E_{\text{in}}^{0/1}$  over the 128 experiments?

(Note: You can choose to run 128 new experiments in this problem, or just re-use the 128 hypotheses  $\mathbf{w}_{\text{LIN}}$  and test data sets obtained from the previous problem.)

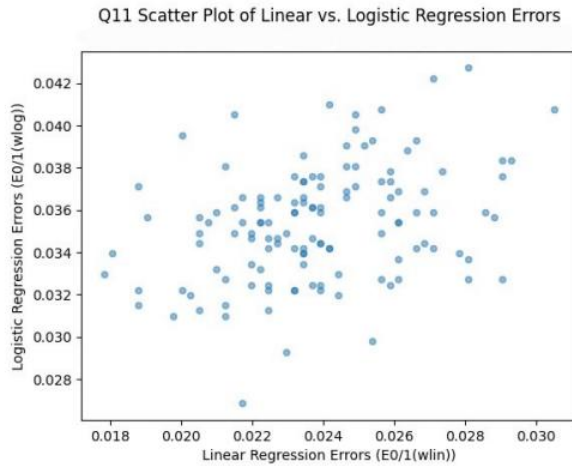


median = 0.02343 #



Question assigned to the following page: [11](#)

11. (20 points, \*) Consider two algorithms. The first one,  $\mathcal{A}$ , is the linear regression algorithm above. The second one  $\mathcal{B}$  is logistic regression, trained with fixed learning rate gradient descent with  $\eta = 0.1$  for  $T = 500$  iterations, starting from  $\mathbf{w}_0 = \mathbf{0}$ . Run the algorithms on the same  $\mathcal{D}$ , and record  $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D})), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))]$ . Repeat the process for 128 times, each with a different random seed for generating the training and test data sets above. Plot a scatter plot for  $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D})), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))]$ . What is the median of  $E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}))$  and what is the median of  $E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}))$ ?



$$E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D})) \quad \text{median : } \underline{0.0235595703125}$$

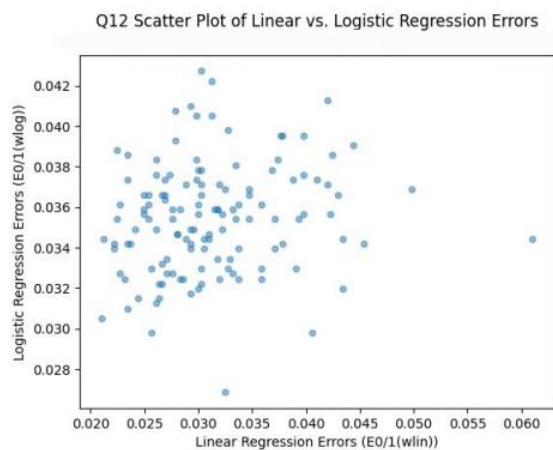
$$E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D})) \quad \text{median : } \underline{0.035400390625}$$





Question assigned to the following page: [12](#)

12. (20 points, \*) Following the previous problem, in addition to the 256 examples in  $\mathcal{D}$ , add 16 outlier examples generated from the following process to your training data (but not to your test data). All outlier examples will be labeled  $y = +1$  and  $\mathbf{x} = [1, x_1, x_2]$  where  $(x_1, x_2)$  comes from a normal distribution of mean  $[0, 6]$  and covariance  $\begin{bmatrix} 0.1 & 0 \\ 0 & 0.3 \end{bmatrix}$ . Name the new training data set  $\mathcal{D}'$ . Run the algorithms on the same  $\mathcal{D}'$ , and record  $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}')), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))]$ . Repeat the process for 128 times, each with a different random seed for generating the training and test data sets above. Plot a scatter plot for  $[E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}')), E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))]$ . What is the median of  $E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}'))$  and what is the median of  $E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}'))$ ? Compare your results to the previous problem. Describe your findings.



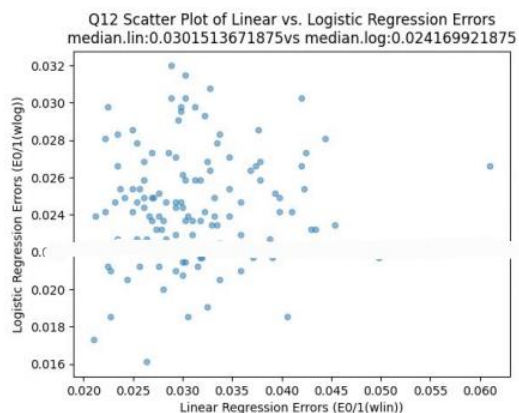
$$E_{\text{out}}^{0/1}(\mathcal{A}(\mathcal{D}')) : \text{median.lin:0.0301513671875}$$

$$E_{\text{out}}^{0/1}(\mathcal{B}(\mathcal{D}')) : \text{median.log:0.035400390625}$$

Median of linear regression algorithm becomes larger than training without outlier data, while median of logistic algorithms doesn't change much.

Since linear regression make error  $\hat{y}^2$ , when we use noise in training data, it may sacrifice original best  $g$ , modify itself to reduce the error from noise when training.

The plot also shows that  $E_{\text{out}}(\log)$  doesn't increase with  $E_{\text{out}}(\text{Lin})$ , but in most of cases  $E_{\text{out}}(\log)$  is bigger than  $E_{\text{out}}(\text{Lin})$ , so I run another plot to see if iteration=20000 (which means it walks more steps to modify),  $E_{\text{out}}'(\log) : 0.0241$  the consequence indicates  $E_{\text{out}}(\log)$  becomes smaller than  $E_{\text{out}}(\text{Lin})$ . That is, if we increase iteration, logistic have the ability to perform better than linear.





Question assigned to the following page: [13](#)

13. (Bonus 20 points) When using Newton's method for solving logistic regression, as discussed in Problem 6, each update  $\mathbf{v}$  is calculated by

$$\mathbf{v} = -(\mathbf{X}^T \mathbf{D} \mathbf{X})^{-1} \nabla E_{\text{in}}(\mathbf{w}_t)$$

when  $(\mathbf{X}^T \mathbf{D} \mathbf{X})$  is invertible. In linear regression, when  $\mathbf{X}^T \mathbf{X}$  is invertible, the optimal

$$\mathbf{w}_{\text{LIN}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

If we can express  $-\nabla E_{\text{in}}(\mathbf{w}_t)$  as some  $\tilde{\mathbf{X}}^T \tilde{\mathbf{y}}$ , and  $\mathbf{X}^T \mathbf{D} \mathbf{X}$  as some  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$ , then each iteration of Newton-solving logistic regression is performing an internal linear regression! State the internal linear regression problem—in particular, what are  $\tilde{\mathbf{X}}$  and  $\tilde{\mathbf{y}}$ ?

13. 
$$h_t(x) = \frac{1}{1 + \exp(-\mathbf{w}^T \mathbf{x})} \quad , \quad \mathbf{D} = \begin{bmatrix} h_t(x_1)(1-h_t(x_1)) & & \\ & h_t(x_2)(1-h_t(x_2)) & \\ & & \ddots & \\ & & & h_t(x_n)(1-h_t(x_n)) \end{bmatrix}$$

① 
$$\mathbf{X}^T \mathbf{D} \mathbf{X} = \mathbf{X}^T \mathbf{V}_D \cdot \mathbf{V}_D^T \mathbf{X} \quad , \quad \tilde{\mathbf{X}} = \mathbf{V}_D^T \mathbf{X} \quad \#$$

② 
$$-\nabla E_{\text{in}}(\mathbf{w}_t) = \frac{1}{N} \sum_{n=1}^N y_n \mathbf{x}_n (1-h_t(x_n)) = \mathbf{X}^T (\mathbf{y} (1-h_t(\mathbf{x})))$$

$$\frac{1}{N} (\mathbf{X}^T \mathbf{V}_D) (\mathbf{V}_D^T \mathbf{y} (1-h_t(\mathbf{x}))) = \tilde{\mathbf{X}}^T \tilde{\mathbf{y}} \quad , \quad \tilde{\mathbf{X}}^T \mathbf{y} = \mathbf{X}^T \mathbf{V}_D^T \tilde{\mathbf{y}} = \mathbf{X}^T \mathbf{V}_D \tilde{\mathbf{y}}$$

$$\tilde{\mathbf{y}} = \frac{1}{N} \mathbf{V}_D^T \mathbf{y} (1-h_t(\mathbf{x})) \quad \#$$

