

## Machine Learning Foundation HW3

R08942062 林志皓

### Problem 1

測驗 • 40 MIN

### 作業三



向您的目標更進一步  
如果您完成此作業，則完成本課程的可能性增加了 **90%**

✓ 提交您的作業

截止時間 1月5日 23:59 PST 答題次數 3/8 hours

再試

在 7h 55m 後重新參加測試

✓ 收到成績

通過條件 75% 或更高

成績

100%

查看反饋

我們會保留您的最高分數

### Problem 2

② Considering 2 case for one  $x_n$

(1)  $\text{sign}(w^T x_n) = y_n \Rightarrow y_n w^T x_n > 0$   
 $\text{err}(w) = \max(0, -y_n w^T x_n) = 0 \quad \frac{\partial \text{err}}{\partial w} = 0$   
 $\Rightarrow w$  is not updated.

(2)  $\text{sign}(w^T x_n) \neq y_n \Rightarrow y_n w^T x_n < 0$ , ( $x_n$  is classified incorrectly)  
 $\text{err}(w) = -y_n w^T x_n \quad \nabla \text{err} = -y_n x_n$   
by sgd,  $w(t+1) = w(t) - \eta \nabla \text{err}$ , set  $\eta = 1$ .  
 $w(t+1) = w(t) + y_n x_n$   
by (1) & (2), it's the same as typical PLA.

### Problem 3

③ Define variables as below:

$$\frac{\partial E}{\partial u}(a,b) = e_u \quad \frac{\partial E}{\partial v}(a,b) = e_v \quad \frac{\partial^2 E}{\partial u^2}(a,b) = e_{uu}$$

$$\frac{\partial^2 E}{\partial v^2}(a,b) = e_{vv} \quad \frac{\partial^2 E}{\partial u \partial v}(a,b) = e_{uv} \quad E(a,b) = e$$

$$\therefore \text{Hessian matrix } H = \begin{bmatrix} e_{uu} & e_{uv} \\ e_{uv} & e_{vv} \end{bmatrix}$$

$$\nabla E = \begin{bmatrix} e_u \\ e_v \end{bmatrix}$$

$$\text{let } x = \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix}$$

$$\Rightarrow \hat{E}_2(\Delta u, \Delta v) = E(a,b) + \frac{\partial E}{\partial u}(a,b) \Delta u + \frac{\partial E}{\partial v}(a,b) \Delta v \\ + \frac{\partial^2 E}{2\partial u^2}(a,b) \Delta u^2 + \frac{\partial^2 E}{2\partial v^2}(a,b) \Delta v^2 + \frac{\partial^2 E}{\partial u \partial v}(a,b) \Delta u \Delta v$$

$$= \frac{1}{2} [\Delta u, \Delta v] \begin{bmatrix} e_{uu} & e_{uv} \\ e_{uv} & e_{vv} \end{bmatrix} \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} + [e_u, e_v] \begin{bmatrix} \Delta u \\ \Delta v \end{bmatrix} + e$$

$$= \frac{1}{2} x^T H x + \nabla E^T x + e$$

$$\text{When optimum } \frac{\partial \hat{E}_2}{\partial x} = 0 = Hx + \nabla E$$

$$\therefore x = -H^{-1} \nabla E$$

$$= -(\nabla^2 E(u,v))^{-1} \nabla E(u,v) \quad \neq$$

#### Problem 4

(16)

In multiclass logistic regression,  
the goal is to maximize probability:

$$\prod_{n=1}^N h_y(x_n) = \prod_{n=1}^N \frac{\exp(W_{y_n}^T x_n)}{\sum_{k=1}^K \exp(W_k^T x_n)}$$

which can be considered as minimizing loss:

$$\begin{aligned} -\frac{1}{N} \sum_{n=1}^N \ln h_y(x_n) &= -\frac{1}{N} \sum_{n=1}^N \ln \left( \frac{\exp(W_{y_n}^T x_n)}{\sum_{k=1}^K \exp(W_k^T x_n)} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \left[ \ln \left( \sum_{k=1}^K \exp(W_k^T x_n) \right) - W_{y_n}^T x_n \right] \end{aligned}$$

#### Problem 5

$$\textcircled{I} \quad \min_W \frac{1}{N+K} \left( \sum_{n=1}^N (y_n - w^T x_n)^2 + \sum_{k=1}^K (y_k - w^T x_k)^2 \right)$$

$$\equiv \min_W \left[ (XW - y)^2 + (\tilde{X}W - \tilde{y})^2 \right]$$

$$\text{Let } E = (XW - y)^2 + (\tilde{X}W - \tilde{y})^2$$

$$\text{when optimum, } \frac{\partial E}{\partial W} = 2(\tilde{X}^T XW - \tilde{X}^T y) + 2(\tilde{X}^T \tilde{X}W - \tilde{X}^T \tilde{y}) = 0$$

$$\Rightarrow (X^T X + \tilde{X}^T \tilde{X}) W = X^T y + \tilde{X}^T \tilde{y}$$

$$\therefore W = (X^T X + \tilde{X}^T \tilde{X})^{-1} (X^T y + \tilde{X}^T \tilde{y}) \quad \#$$

### Problem 6

(6)

$$\text{Let } E_{\text{aug}} = \lambda \|w\|^2 + \|Xw - y\|^2$$

$$\frac{\partial E_{\text{aug}}}{\partial w} = 2\lambda w + 2(X^T X w - X^T y)$$

$$\text{When minimum } E_{\text{aug}}, (\lambda I + X^T X)w - X^T y = 0$$

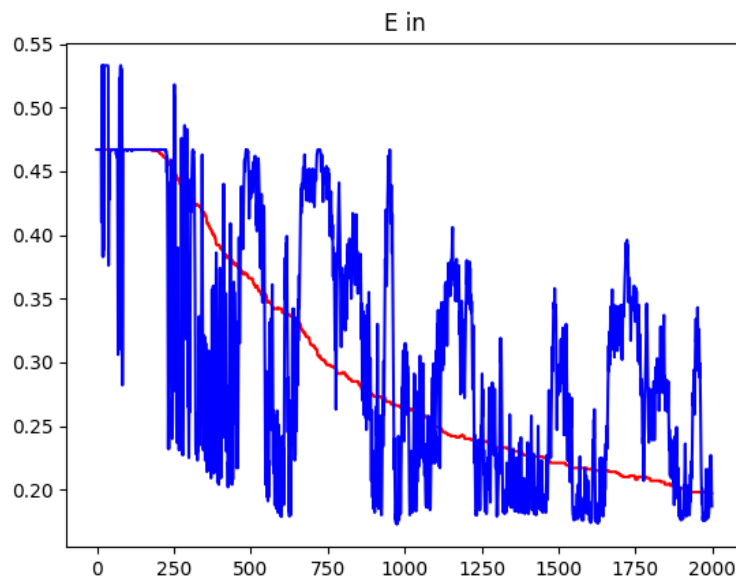
$$w = (X^T X + \lambda I)^{-1} X^T y$$

$$(\text{compared to } \textcircled{5}) = (X^T X + \tilde{X}^T \tilde{X})^{-1} (X^T y + \tilde{X}^T \tilde{y})$$

$$\therefore \tilde{X}^T \tilde{X} = \lambda I, \quad \tilde{X} = \sqrt{\lambda} I$$

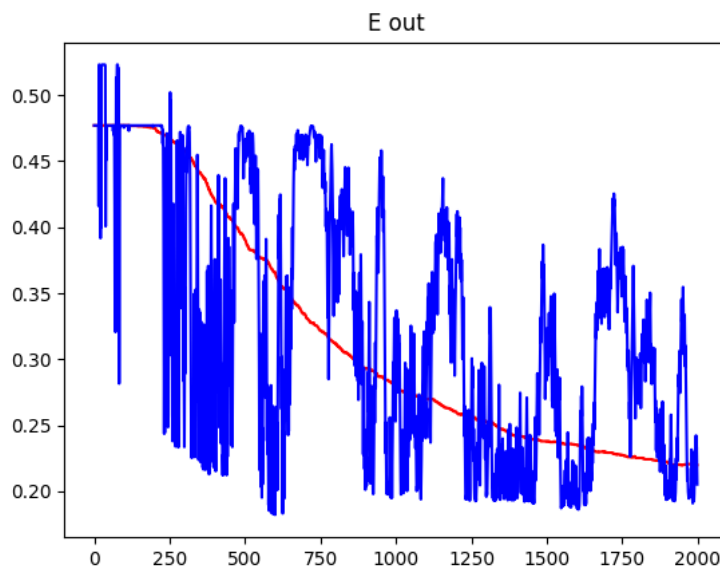
$$\tilde{X}^T \tilde{y} = 0, \quad \tilde{y} = 0 \quad \#$$

### Problem 7



The red line represents the curve for gradient decent, and the blue one is for stochastic gradient decent. We can see that the red one is much more stable than blue one. On the other hand, blue line might reach lowest  $E_{in}$  in some iteration, but it's very unstable.

### Problem 8



The red line represents the curve for gradient decent, and the blue one is for stochastic gradient decent. We can find similar observation with previous question. Fortunately, the model generalize well on testing data with both GD and SGD.



## Problem 9

⑨

$$\begin{aligned}
 (a) \quad X^T X W_{\min} &= X^T X (V \Gamma^{-1} U^T y) \\
 &= X^T (U \Gamma V^T) (V \Gamma^{-1} U^T) y \\
 &= X^T U U^T y \quad (\because V^T V = I, \Gamma \Gamma^{-1} = I) \\
 &= (V \Gamma U^T) U U^T y = V \Gamma U^T y \quad (\because U^T U = I) \\
 &= (U \Gamma V^T)^T y = X^T y
 \end{aligned}$$

(b)  $X^T X W = X^T y$ , by SVD decomposition,  $X = U \Gamma V^T$ .

$$(V \Gamma U^T) (U \Gamma V^T) W = (V \Gamma U^T) y$$

$$V \Gamma^2 V^T W = V \Gamma U^T y \quad \text{multiply } \Gamma^{-2} V^T \text{ on both sides}$$

$$V^T W = \Gamma^{-1} U^T y. \quad (\because V^T V = I) \quad \text{let } \Gamma^{-1} U^T y = a = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix}$$

then it's equivalent to solve  $V^T W = a$ .  $V \in \mathbb{R}^{(d+1) \times p}$

$V$  is composed of  $p$  orthogonal basis in  $(d+1)$ -dimension space, that is,

$$V = [v_1, v_2, v_3, \dots, v_p], \quad \|v_i\| = 1 \quad \forall i, \quad v_i^T v_j = 0 \quad \forall i \neq j$$

let  $W = c_1 v_1 + c_2 v_2 + \dots + c_p v_p + c_k v_k$ , where  $v_k^T v_i = 0, i = 1, 2, \dots, p$ ,  $c_i \in \mathbb{R} \quad \forall i$

$$V^T W = \begin{bmatrix} v_1^T W \\ v_2^T W \\ \vdots \\ v_p^T W \end{bmatrix} = \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ \vdots \\ c_p \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_p \end{bmatrix} \quad \therefore W = a_1 v_1 + a_2 v_2 + \dots + a_p v_p + c_k v_k$$

$$\|W\|^2 = \|a_1 v_1\|^2 + \|a_2 v_2\|^2 + \dots + \|c_k v_k\|^2 \quad \text{reach min when } c_k = 0$$

$$\therefore \arg \min_W \|W\| = a_1 v_1 + a_2 v_2 + \dots + a_p v_p = V a = V \Gamma^{-1} U^T y = W_{\min}$$

which is shortest  $W$  that satisfy  $X^T X W = X^T y$

$$\|W_{\min}\| \leq \|W\| \quad \forall W$$