

# 100-CSC413 Homework 1

Name: Kwan Kiu Choy  
Student No: 1005005879

[1.1](#)

[1.2](#)

[2.1.2](#)

[2.2.1](#)

[2.2.2](#)

[2.2.3](#)

[3.2.1](#)

[3.2.2](#)

[3.3.2](#)

[3.3.4](#)

1.1

1.1

Assumption:  $x_1, x_2$  are distinct and positive

$$w^{(1)} = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix} \quad w^{(2)} = \begin{bmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

$$\phi^{(1)}(z) = \begin{matrix} \text{absolute value} \\ \text{function} \end{matrix} \quad \phi^{(2)}(z) = z \quad \left( \begin{matrix} \text{this is allowed} \\ \text{based on @23 on Piazza} \end{matrix} \right)$$

Given input  $z$

$$\phi^{(1)}(z) = |z|$$

Note: we don't have to worry  
about  $|x_1 + x_2|$  changing the sign of  $x_1 + x_2$  since  
both  $x_1, x_2$  are positive.

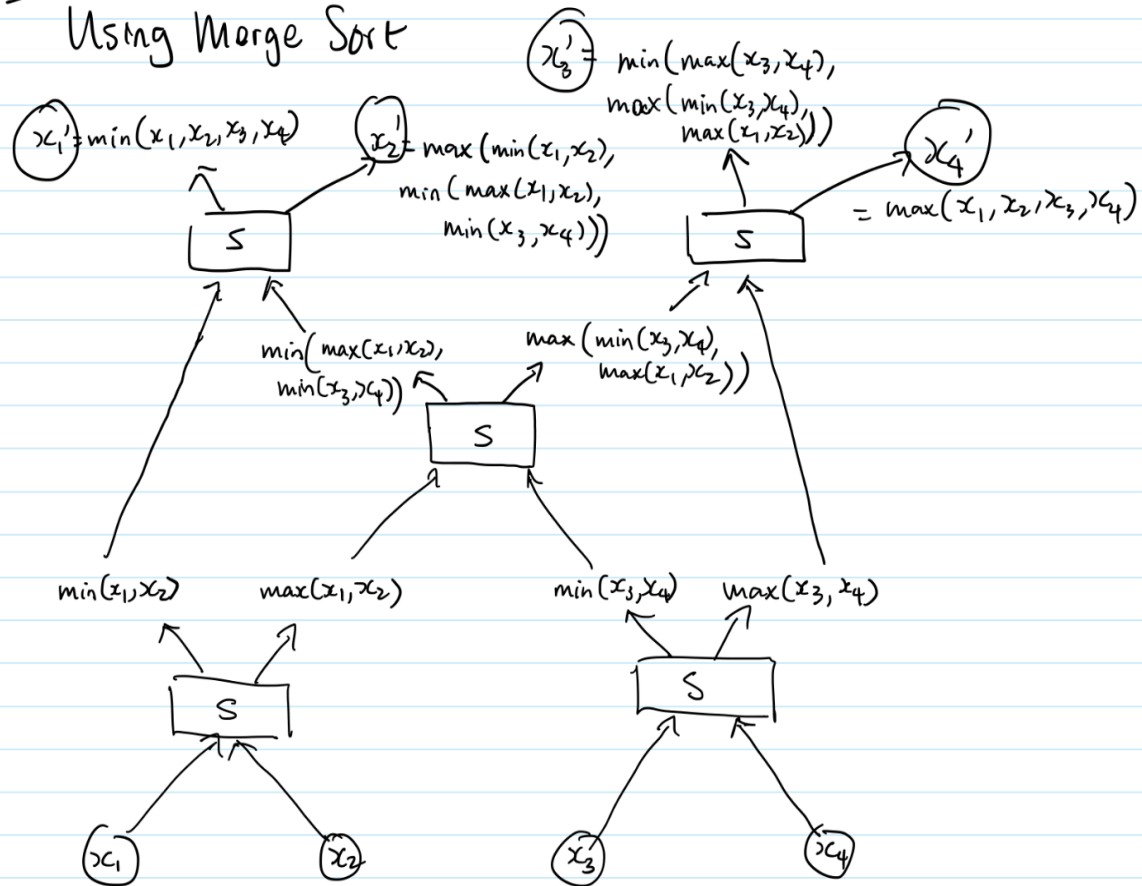
$$b^{(1)} = \vec{0} \quad (\text{zero vector})$$

$$b^{(2)} = \vec{0} \quad (\text{zero vector})$$

1.2

1.2

Using Merge Sort



## 2.1.2

### 2.1.2 Backward pass

Goal: Compute  $\bar{x} = \frac{\partial J}{\partial x}$

$$\frac{\partial J}{\partial y} = 1 //$$

$$\frac{\partial J}{\partial s} = \frac{\partial J}{\partial y} \frac{\partial y}{\partial s} = \bar{y} \times (-1) = -\bar{y} //$$

$$\begin{aligned} \frac{\partial J}{\partial y'} &= \frac{\partial J}{\partial s} \frac{\partial s}{\partial y'} = \bar{s} \left( \begin{bmatrix} 1 \\ 0 \end{bmatrix} \frac{1}{y_t} \right) \\ &= \begin{bmatrix} 0 \\ 1 \end{bmatrix} \frac{1}{y_t} \bar{s} // \end{aligned}$$

one hot vector where only the  $t$ th position has a 1 and other position is 0.

$$\frac{\partial J}{\partial y} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial y} \frac{\partial y}{\partial y'} = \bar{y} \text{softmax}'(y)$$

$$\frac{\partial J}{\partial g} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial g} \frac{\partial g}{\partial y'} \frac{\partial y'}{\partial y} = \bar{y} W_3$$

$$\frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial g} \frac{\partial g}{\partial y'} \frac{\partial y'}{\partial y} \frac{\partial y}{\partial h_1} = h_2 \circ \bar{y}$$

element-wise multiplication

$$\frac{\partial J}{\partial h_2} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial g} \frac{\partial g}{\partial y'} \frac{\partial y'}{\partial y} \frac{\partial y}{\partial h_2} = h_1 \circ \bar{y}$$

derivative of RELU can be split into cases

$$\text{RELU}'(z) = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \\ \text{undefined} & z = 0 \end{cases}$$

$$\frac{\partial J}{\partial z_1} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial g} \frac{\partial g}{\partial y'} \frac{\partial y'}{\partial y} \frac{\partial y}{\partial z_1} \frac{\partial z_1}{\partial h_1} = \bar{h}_1 \frac{\partial h_1}{\partial z_1} \rightarrow \bar{h}_1 \text{ or } 0 \text{ or undefined} //$$

$$\frac{\partial J}{\partial z_2} = \frac{\partial J}{\partial s} \frac{\partial s}{\partial g} \frac{\partial g}{\partial y'} \frac{\partial y'}{\partial y} \frac{\partial y}{\partial z_2} \frac{\partial z_2}{\partial h_2} = \bar{h}_2 \frac{e^{z_2}}{(1+e^{z_2})^2} = \bar{h}_2 (\sigma(z_2)(1-\sigma(z_2)))$$

$$\begin{aligned} \frac{\partial J}{\partial x} &= \frac{\partial J}{\partial z_1} \frac{\partial z_1}{\partial x} + \frac{\partial J}{\partial z_2} \frac{\partial z_2}{\partial x} + \frac{\partial J}{\partial y} \frac{\partial y}{\partial x} \\ &= \bar{z}_1 W_1 + \bar{z}_2 W_2 + \bar{y} W_4 \end{aligned}$$

$$\text{Therefore, } \bar{x} = \frac{\partial J}{\partial x}^T = W_1^T \bar{z}_1^T + W_2^T \bar{z}_2^T + W_4^T \bar{y}^T //$$

## 2.2.1

### 2.2.1 Naive Computation



$$\begin{aligned}
 \frac{\partial J}{\partial w_2} &= \left( \frac{\partial J}{\partial y} \frac{\partial y}{\partial h} \frac{\partial h}{\partial w_2} \right)^T & z &= \begin{bmatrix} 1 & 2 & 1 \\ -2 & 1 & 0 \\ 1 & -2 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \\
 &= \left( y \frac{\partial y}{\partial w_2} \right)^T & &= \begin{bmatrix} 8 \\ 1 \\ -6 \end{bmatrix} \\
 &= \left( y h^T \right)^T & h &= \begin{bmatrix} 8 \\ 1 \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} 1 \\ 1 \end{bmatrix} [8 \ 1 \ 0]^T & \frac{\partial J}{\partial h} &= w_2^T y \\
 &= \begin{bmatrix} 8 & 8 & 8 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial J}{\partial z} &= h \circ \text{RELU}'(z) & &= \begin{bmatrix} -2 & 1 & -3 \\ 4 & -2 & 4 \\ 1 & -3 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \\
 &= \begin{bmatrix} -4 \\ 6 \\ 4 \end{bmatrix} \circ \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} & &= \begin{bmatrix} -4 \\ 6 \\ 4 \end{bmatrix} \\
 &= \begin{bmatrix} -4 \\ 6 \\ 0 \end{bmatrix}
 \end{aligned}$$

The Jacobian matrix for  $\frac{\partial J}{\partial w_1}$  and  $\frac{\partial J}{\partial w_2}$  are

$$\begin{aligned}
 \frac{\partial J}{\partial w_1} &= \left( \bar{z} X^T \right)^T & \frac{\partial J}{\partial w_1} &= \begin{bmatrix} -4 & 6 & 0 \\ 12 & 18 & 0 \\ -4 & 6 & 0 \end{bmatrix} & \text{Norm: } \left\| \frac{\partial J}{\partial w_1} \right\|_F^2 &= 572 \\
 &= \begin{bmatrix} -4 \\ 6 \\ 0 \end{bmatrix} [1 \ 3 \ 1]^T & \frac{\partial J}{\partial w_2} &= \begin{bmatrix} 8 & 8 & 8 \\ 1 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} & \text{Norm: } \left\| \frac{\partial J}{\partial w_2} \right\|_F^2 &= 195
 \end{aligned}$$

## 2.2.2

2.2.2 Efficient Computation

$$\begin{aligned}\left\| \frac{\partial J}{\partial w_1} \right\|_F^2 &= (X^T X) (\bar{z}^T \bar{z}) \\ &= [1 \ 3 \ 1] \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \times [4 \ 6 \ 0] \begin{bmatrix} 4 \\ 6 \\ 0 \end{bmatrix} \\ &= 11 \times (16 + 36) \\ &= 11 \times 52 \\ &= 572_{//}\end{aligned}$$

$$\begin{aligned}\left\| \frac{\partial J}{\partial w_2} \right\|_F^2 &= (h^T h) (\bar{y}^T \bar{y}) \\ &= [8 \ 10] \begin{bmatrix} 8 \\ 1 \\ 6 \end{bmatrix} \times [1 \ 11] \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \\ &= 65 \times 3 \\ &= 195_{//}\end{aligned}$$

Both norms computed what we got in 2.1.2, are the same, with what we have.

## 2.2.3

2.2.3

① Forward pass:

i) T-Naive:

Each layer to layer transition

$$W(\text{Input}) = \text{Output}$$

$\begin{matrix} \uparrow & \uparrow & \uparrow \\ D \times D & D \times 1 & D \end{matrix}$

memory needed:  $(D^2 + D)N = O(ND^2 + ND) = O(ND^2)$

For each input, each layer to layer transition we have  $D^2$  scalar multiplications

$$T_{\text{Naive Forward}} = KND^2$$

ii) T-Efficient:

Each layer to layer transition same forward pass with Naive

$$\text{So } T_{\text{Efficient Forward}} = KND^2$$

② Backward pass

i) T-Naive:

Each layer to layer transition

2x no of scalar multiplication in forward pass

$$= 2KND^2$$

memory needed

weight matrix and

$$= O(NKD^2)$$

ii) T-Efficient:

Since we don't want to update the weight/bias we don't need to compute the  $\frac{\partial L}{\partial W_k}$  Jacobian matrix.

$$\text{Hence } T_{\text{Efficient}} = NKD^2$$

③ Compute Norm

T-Naive

$\frac{\partial L}{\partial W_k} \rightarrow$  will be a  $D \times D$  matrix

$$\text{memory needed } (D^2 + D)N = O(ND^2 + ND) = O(ND^2)$$

We have  $N$   $D \times D^2$  matrices

$$\frac{\partial L}{\partial W_k} \frac{\partial L}{\partial W_k} \Rightarrow D^3 \text{ scalar mults} \rightarrow D^2 \text{ effective scalar multiplication}$$

$$\text{memory: } O(ND^2)$$

Total:  $KND^3$  scalar mults

T-Efficient

Each input, each weight matrix

$$\left( \begin{matrix} 1 & \xrightarrow{D} & 0 \end{matrix} \right) \otimes \left( \begin{matrix} 1 & \xrightarrow{D} & 0 \end{matrix} \right)$$

Memory Space:

$$2ND = O(ND)$$

$$= 2D + 1$$

$$\text{Total: } NK(2D + 1) = 2NKD + NK$$

Summary

	T-Naive	T-Efficient	M-Naive	M-Efficient
Forward Pass	$KND^2$	$KND^2$	$O(ND^2)$	$O(ND^2)$
Backward Pass	$2KND^2$	$KND^2$	$O(ND^2)$	$O(ND^2)$
Compute Norm	$KND^2$	$2NKD + NK$	$O(ND^2)$	$O(ND)$

### 3.2.1

3.2.1

$$J = \frac{1}{n} \|X\hat{w} - t\|_2^2$$

$$\frac{\partial J}{\partial \hat{w}} = \frac{2}{n} X^T (X\hat{w} - t)$$

For the gradient descent process to terminate,  $\frac{\partial J}{\partial \hat{w}} = 0$

$$\frac{\partial J}{\partial \hat{w}} = 0$$

$$\frac{2}{n} X^T (X\hat{w} - t) = 0$$

$$X^T X \hat{w} - X^T t = 0$$

$$X^T X \hat{w} = X^T t$$

$$\hat{w} = (X^T X)^{-1} X^T t$$

Since  $d < n$ , we know that  $X^T X$  is invertible and hence we get  $\hat{w} = (X^T X)^{-1} X^T t$ .

### 3.2.2

3.2.2

$$\begin{aligned}
 \text{Error} &= \frac{1}{n} \|X\hat{w} - t\|_2^2 \\
 &= \frac{1}{n} \|X(X^T X)^{-1} X^T t - t\|_2^2 \\
 &= \frac{1}{n} \|X(X^T X)^{-1} X^T t - I\|_2^2 \\
 &= \frac{1}{n} \|X(X^T X)^{-1} X^T (Xw^* + \epsilon)\|_2^2 \quad \xrightarrow{\text{By Associative property of matrix multiplication}} \\
 &= \frac{1}{n} \|X(X^T X)^{-1} X^T Xw^* + (X(X^T X)^{-1} X^T - I)\epsilon\|_2^2 \\
 &= \frac{1}{n} \|\cancel{Xw^*} - \cancel{Xw^*} + (X(X^T X)^{-1} X^T - I)\epsilon\|_2^2 \\
 &= \frac{1}{n} \|(X(X^T X)^{-1} X^T - I)\epsilon\|_2^2
 \end{aligned}$$

Expected Error:

$$\text{Error} = \text{tr}(\text{Error})$$

$$\mathbb{E}(\text{Error}) = \text{tr}(\mathbb{E}(\text{Error}))$$

$$\rightarrow \text{Let } V = (X(X^T X)^{-1} X^T - I)\epsilon$$

$$\begin{aligned}
 \mathbb{E}(\text{Error}) &= \text{tr}(\mathbb{E}(\frac{1}{n} \|V\|_2^2)) \\
 &= \frac{1}{n} \text{tr}(\mathbb{E}(V^T V)) \\
 &= \frac{1}{n} \text{tr}(\mathbb{E}(\epsilon^T (X(X^T X)^{-1} X^T - I)^T (X(X^T X)^{-1} X^T - I) \epsilon))
 \end{aligned}$$

Making use of the cyclic property of trace

$$= \frac{1}{n} \text{tr}((X(X^T X)^{-1} X^T - I)^T (X(X^T X)^{-1} X^T - I) \mathbb{E}(\epsilon \epsilon^T))$$

$$\rightarrow \text{Let } A = X(X^T X)^{-1} X^T - I$$

$$= \frac{1}{n} \text{tr}(A^T A \sigma^2)$$

$$= \frac{\sigma^2}{n} \text{tr}(A^T A) \quad \text{--- } (*)$$

$$\rightarrow \text{Let } B = X(X^T X)^{-1} X^T$$

$$A = B - I$$

$$A^T A = (B - I)^T (B - I)$$

$$= (B^T - I^T) (B - I)$$

$$= B^T B - B - B^T + I$$

$$= \cancel{X(X^T X)^{-1} X^T X(X^T X)^{-1} X^T} - X(X^T X)^{-1} X^T - X(X^T X)^{-1} X^T + I$$

$$= I - B^T$$

$$\text{tr}(A^T A) = \text{tr}(I - B^T)$$

$$= \text{tr}(I) - \text{tr}(B^T)$$

$$= \text{rank}(I) - \text{rank}(B^T)$$

$$= n - d \quad \text{--- } (1)$$

Substituting (1) into (\*), we get

$$\mathbb{E}(\text{Error}) = \frac{\sigma^2}{n} (\text{tr}(A^T A))$$

$$= \frac{\sigma^2}{n} (n - d) //$$



## 3.3.2

3.3.2Let  $\hat{w} = X^T A$  where  $A \in \mathbb{R}^w$ 

$$J = \frac{1}{n} \|X \hat{w} - t\|_2^2$$

$$\frac{\partial J}{\partial w} = \frac{2}{n} X^T \|X \hat{w} - t\|_2^2$$

$$\frac{2}{n} X^T \|X \hat{w} - t\|_2^2 = 0$$

$$\frac{2}{n} X^T (X X^T A - t) = 0$$

$$\frac{2}{n} X X^T (X X^T A - t) = 0$$

$$X X^T A - t = 0$$

$$X X^T A = t$$

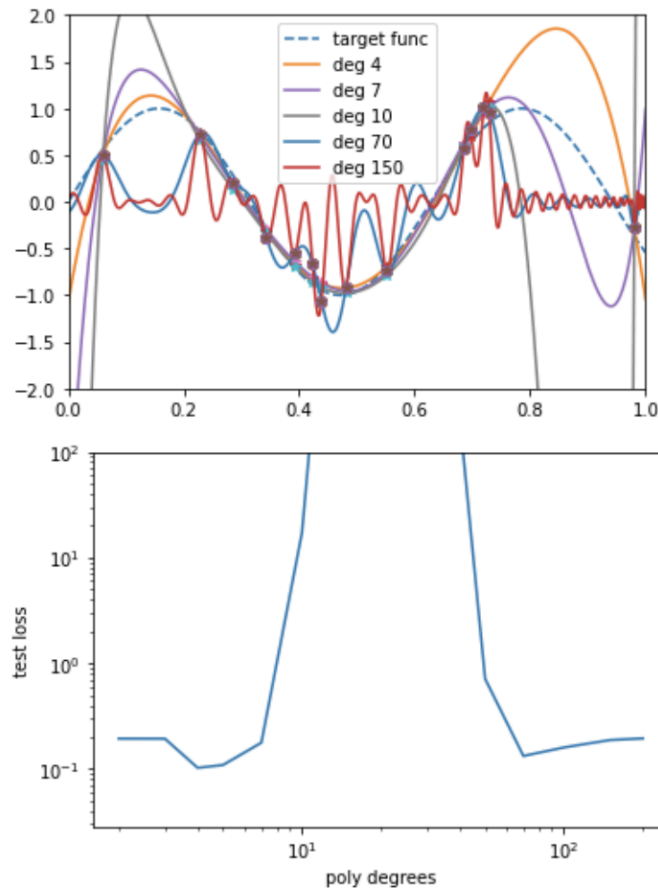
$$A = (X X^T)^{-1} t$$

$$\hat{w} = X^T (X X^T)^{-1} t$$

Since  $d > n$ , we know that  $X X^T$  is invertible and hence  $\hat{w} = X^T (X X^T)^{-1} t$

### 3.3.4

```
def fit_poly(X, d,t):
    X_expand = poly_expand(X, d=d, poly_type = poly_type)
    n = X.shape[0]
    if d > n:
        inv = np.linalg.inv(np.dot(X_expand, X_expand.T))
        head = np.dot(X_expand.T, inv)
        W = np.dot(head, t)
    else:
        inv = np.linalg.inv(np.dot(X_expand.T, X_expand))
        head = np.dot(inv, X_expand.T)
        W = np.dot(head, t)
    return W
```



#### Analysis:

From this diagram above, we know that high degrees don't always lead to test error. Based on the test loss vs poly degrees diagram, we can see a decrease in test loss after degrees exceed a certain value.