

1.1.1

We first consider a simple case with only one data, i.e., let $x_1 = [2, 1]$, $w_0 = [0, 0]$ and $t = 3$. By the formula of mini-batch SGD, we have

$$\begin{aligned} w_1 &= w_0 - \eta \nabla_{w_0} L(x_1, w_0) \\ &= w_0 - 2\eta x_1 (w_0^T x_1 - t_0) \\ &= (2\eta t_0) x_1 \quad (\text{since } w_0 = [0, 0]) \\ &= \alpha x_1 \quad \text{for some } \alpha \in \mathbb{R} \end{aligned}$$

$$\begin{aligned} w_2 &= w_1 - \eta \nabla_{w_1} L(x_1, w_1) \\ &= w_1 - 2\eta x_1 (w_1^T x_1 - t_1) = w_1 - 2\eta x_1 (\alpha x_1^T x_1 - t_1) \\ &= \alpha x_1 - \underbrace{(2\eta \alpha \|x_1\|_2^2 - 2\eta t_1)}_{\text{constant}} x_1 = \alpha' x_1 \quad \text{for some } \alpha' \in \mathbb{R} \end{aligned}$$

Similarly, we have $w_3 = \alpha'' x_1$ for some $\alpha'' \in \mathbb{R}$. As we can see \hat{w} is in the span of X . It can be generalized since for any x_j , $\nabla_{w_t} L(x_j, w_t) = \nabla_{w_t} (w_t^T x_j - t_j)^2 = \underbrace{2x_j (w_t^T x_j - t_j)}_{\text{constant}} = \alpha x_j$ for some $\alpha \in \mathbb{R}$. Since $w_0 = 0$, we have w_t lies in $\text{span}(X)$. Therefore, we have shown that \hat{w} can be written as $\mathcal{B}^T a$ for $a \in \mathbb{R}^b$ in the stationary cond. for mini-batch SGD. Also, as $\mathcal{B} \in \mathbb{R}^{b \times d}$ and $X \in \mathbb{R}^{n \times d}$ with full rank, and $b \leq n < d$, we have \mathcal{B} in the row span of X . Therefore, we can write $\mathcal{B} = bX$ for $b \in \mathbb{R}^{b \times n}$, and so $\hat{w} = \mathcal{B}^T a = X^T b^T a$

$$\begin{aligned} &= X^T (b^T a) \\ &= X^T c \end{aligned}$$

for some $c \in \mathbb{R}^n$. As a result, since w^* & \hat{w} can be written in the same form as $X^T c$ for some $c \in \mathbb{R}^n$, we must have $w^* = \hat{w} = X^T (XX^T)^{-1} t$. To show \hat{w} is indeed the minimum norm, assume there exists some other

solution w , but we have:

$$\begin{aligned}
 (\hat{w} - w)^T \hat{w} &= (\hat{w} - w)^T X^T C \\
 &= [X(\hat{w} - w)]^T C \\
 &= [t - t]^T C = 0
 \end{aligned}$$

which means that $\|w\|_2^2 = \|w - \hat{w} + \hat{w}\|_2^2 = \|w - \hat{w}\|_2^2 + \|\hat{w}\|_2^2 \geq \|\hat{w}\|_2^2$, and so we show \hat{w} is indeed the minimum norm.

1.2.1

Start with $x_1 = [2, 1]$, $w_0 = [0, 0]$ and $t = [2]$, we have:

$$\begin{aligned}
 \nabla_{w_0} L(x_1, w_0) &= 2(w_0^T x_1 - t)x_1 \\
 &= 2(-2)[2, 1] = [-8, -4]
 \end{aligned}$$

Therefore, $\nabla_{w_{0,0}} L(w_{0,0}) = -8$ & $\nabla_{w_{1,0}} L(w_{1,0}) = -4$, and

$$\begin{aligned}
 v_{0,0} &= \beta v_{0,-1} + (1-\beta)(-8)^2 \\
 &= 64(1-\beta)
 \end{aligned}$$

Similarly,

$$v_{1,0} = (1-\beta)(-4)^2 = 16(1-\beta)$$

As a result,

$$\begin{aligned}
 w_{0,1} &= w_{0,0} - \frac{\eta}{\sqrt{v_{0,0}} + \epsilon} (-8) \\
 &= \frac{8\eta}{8\sqrt{1-\beta} + \epsilon}
 \end{aligned}$$

$$\begin{aligned}
 w_{1,1} &= w_{1,0} - \frac{\eta}{\sqrt{v_{1,0}} + \epsilon} (-4) \\
 &= \frac{4\eta}{4\sqrt{1-\beta} + \epsilon}
 \end{aligned}$$

So we have $w_1 = \left[\frac{8\eta}{8\sqrt{1-\beta} + \epsilon}, \frac{4\eta}{4\sqrt{1-\beta} + \epsilon} \right]$. Now, we compute $w_{0,2}$:

$$\therefore \nabla_{w_1} L(x_1, w_1) = 2(w_1^T x_1 - t)x_1$$

$$\begin{aligned}
&= 2 \left[\frac{16\eta}{8\sqrt{1-\beta}+\varepsilon} + \frac{4\eta}{4\sqrt{1-\beta}+\varepsilon} - 2 \right] [2, 1] \\
&= \left[\frac{16\eta}{8\sqrt{1-\beta}+\varepsilon} + \frac{4\eta}{4\sqrt{1-\beta}+\varepsilon} - 2 \right] [4, 2]
\end{aligned}$$

We denote $\nabla_{w_{0,1}} L(w_{0,1}) = a_1$ & $\nabla_{w_{1,1}} L(w_{1,1}) = a_2$, we have

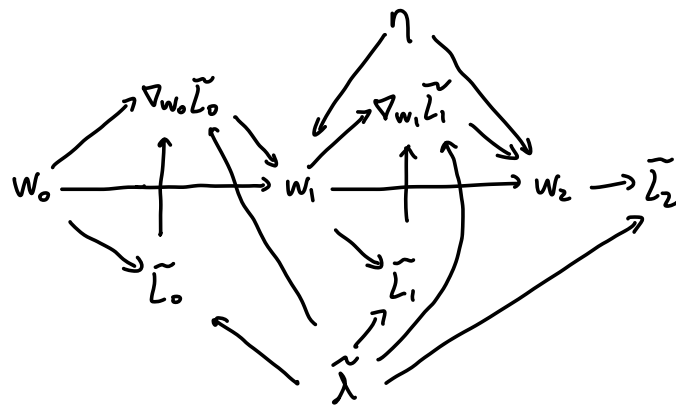
$$v_{0,1} = \beta v_{0,0} + (1-\beta) a_1^2 = (1-\beta)(64\beta + a_1^2)$$

$$\begin{aligned}
w_{0,2} &= w_{0,1} - \frac{\eta}{\sqrt{v_{0,1}} + \varepsilon} a_1 \\
&= \frac{8\eta}{8\sqrt{1-\beta} + \varepsilon} - \frac{\eta a_1}{\sqrt{(64\beta + a_1^2)(1-\beta)} + \varepsilon}
\end{aligned}$$

However, this component is not a linear function of the first component of x_1 , meaning that w stops to be in the span of X since span is the linear combination of rows in X , here it's x_1 . Therefore, it violates the condition required in 1.1.1, leading to the consequence that w is not always the minimum norm. This still conv. to some other pt though.

2.1.1

The computation graph looks like follows:



2.1.2

Memory complexity for the forward-propagation is $O(1)$ since we only need to save the w_i at a time for $i=1, \dots, t$. For the backward propagation to compute $\nabla_{\eta} \tilde{L}_t$ is $O(t)$ since from the computation graph we know η is involved in every w_i for $i=1, \dots, t$, and so storing all w_i requires $O(t)$

2.2.1

$$\begin{aligned} \because w_1 &= w_0 - \eta \nabla_{w_0} L_0 = w_0 - \frac{2\eta}{n} X^T (Xw_0 - t) \\ &= w_0 - \frac{2\eta}{n} X^T a \quad \text{for } a = Xw_0 - t \end{aligned}$$

$$\therefore L_1 = \frac{1}{n} \|Xw_1 - t\|_2^2$$

$$\begin{aligned} &= \frac{1}{n} \|X(w_0 - \frac{2\eta}{n} X^T a) - t\|_2^2 = \frac{1}{n} \left\| -\frac{2\eta}{n} X X^T a + Xw_0 - t \right\|_2^2 \\ &= \frac{1}{n} \left\| -\frac{2\eta}{n} X X^T a + a \right\|_2^2 \\ &= \frac{1}{n} \left\| \left(I - \frac{2\eta}{n} X X^T \right) a \right\|_2^2 \quad (I \text{ is an } n \times n \text{ identity matrix}) \\ &= \frac{1}{n} \left[a^T \left(I - \frac{2\eta}{n} X X^T \right)^T \left(I - \frac{2\eta}{n} X X^T \right) a \right] \\ &= \frac{1}{n} a^T \left(I - \frac{2\eta}{n} X X^T \right)^2 a \end{aligned}$$

2.2.3

$$\nabla_{\eta} L_1 = \frac{2}{n} a^T \left(I - \frac{2\eta}{n} X X^T \right) \left(-\frac{2}{n} X X^T \right) a \quad (\text{by chain rule, and since } L_1 \text{ \& } \eta \text{ are both numbers, the derivative should be number. too})$$

$$\text{Set } \nabla_{\eta} L_1 = 0, \text{ we have } a^T \left(-\frac{2}{n} X X^T + \frac{2\eta}{n} X X^T \right) a = 0$$

$$\Rightarrow \frac{2}{n} \eta = \frac{a^T X X^T a}{a^T (X X^T)^2 a}$$

$$\Rightarrow \eta = \frac{n}{2} \frac{a^T X X^T a}{a^T (X X^T)^2 a}$$

But $a^T x x^T a = \|x^T a\|_2^2$, and $a^T (x x^T)^2 a = a^T x x^T x x^T a$
 $= \|x x^T a\|_2^2$

So after simplification, $\eta = \frac{\eta}{2} \frac{\|x^T a\|_2^2}{\|x x^T a\|_2^2}$

2.3.1

For \tilde{L} : $w_1 = w_0 - \eta \left[\frac{2}{n} x^T (x w_0 - t) + 2 \tilde{\lambda} w_0 \right]$

For L : $w_1 = (1 - \lambda) w_0 - \frac{2\eta}{n} x^T (x w_0 - t)$

2.3.2

Since for \tilde{L} , we have: $w_1 = w_0 - \frac{2\eta}{n} x^T (x w_0 - t) - 2 \tilde{\lambda} \eta w_0$, and for L , we have

$w_1 = w_0 - \frac{2\eta}{n} x^T (x w_0 - t) - \lambda w_0$, we compare two equations and we can have

$$2 \tilde{\lambda} \eta w_0 = \lambda w_0$$

$$\Rightarrow \tilde{\lambda} = \frac{\lambda}{2\eta}$$

3.1

After flipping the filter, we have $\text{flip}(J) = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$, and after computation, we have

$$I * J = \begin{bmatrix} 0 & -1 & -2 & -3 & -2 \\ -2 & -3 & -3 & -2 & -1 \\ -1 & -1 & -1 & 1 & 1 \\ 2 & 2 & 2 & 1 & 1 \\ 1 & 2 & 3 & 2 & 1 \end{bmatrix}$$

This filter detects the horizontal edge.

3.2

total # of neurons:

① CNN: Image: 32×32

Conv 1: 32×32

pool 1: 16×16

Conv 2: 16×16

pool 2: 8×8

Conv 3: 8×8

total: 2688

② FCNN: Image: 1024

FC 1: 1024

pool 1: 256

FC 2: 256

pool 2: 64

FC 3: 64

total: 2688

Total # of weights:

① CNN: Image to Conv1: 3×3

pool1 to conv2: 3×3

pool 2 to conv 2: 3×3

total: 27

② FCNN: Im to FC1: 1024×1024

pool1 to FC2: 256×256

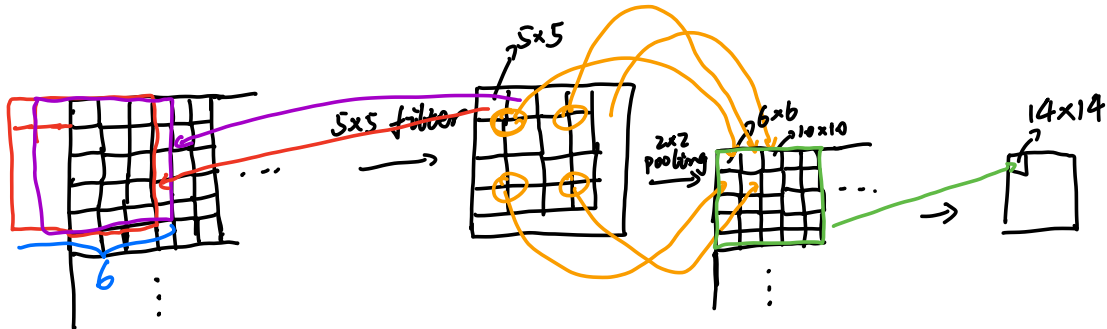
pool 2 to FC: 64×64

total: 1118208

Disadvantage of having more trainable parameters: too high memory cost

3.3.

A rough analysis's guide is following graph:



The receptive field of a neuron after the second conv. layer is 14×14 .

Two other things that can affect the size of the receptive field:

① max-pooling layer size ② the stride when applying the convolutional layer