

(Q1.1.1)

assume dimensions:

$$w_t = d \times 1$$

$$x_i = d \times 1$$

Then,

$$L(x_i, w_t) = \|w_t^T x_i - t_i\|_2^2$$

$$\nabla_{w_t} L(x_i, w_t) = \underbrace{2(w_t^T x_i - t_i)x_i}_{\in \mathbb{R}}, \text{ and } x_i \text{ is in rowspace of } X \quad (\star)$$

Assuming that we start from zero weight initialization

$w_0 = 0$, which is a linear combination of X ,
the next weight updates are always in the row span of X ,

since

$$w_{t+1} = \underbrace{w_t}_{\in \text{Span}\{x\}} - \frac{1}{b} \sum_{j \neq i} \underbrace{\nabla_{w_t} L(x_j, w_t)}_{\in \text{Span}\{x\}} \text{ by } (\star)$$

is $\in \text{Span}\{x\}$ by induction.

So we can represent $\hat{w} = X^T a$ for some $a \in \mathbb{R}^n$, i.e. $\hat{w} \in \text{Span}(X)$.

Next, want to show that $\hat{w} = X^T a$ is indeed the minimum norm solution. This would be sufficient to show that $\hat{w} = w^*$ from gradient descent since w^* is the unique mn. norm solution.

take any solution w' , we have

$$(\hat{w} - w')^T \hat{w} = (\hat{w} - w')^T X^T a = (X(\hat{w} - w'))^T a = (X\hat{w} - Xw')^T a = (t - t)^T a = 0$$

since \hat{w} and w' are both solutions.

which means that $\hat{w} - w'$ is orthogonal to \hat{w} .

Then by the Pythagorean theorem,

$$\|w'\|_2^2 = \|\hat{w} - w' - \hat{w}\|_2^2 = \|\hat{w} - w'\|_2^2 + \|\hat{w}\|_2^2 \geq \|\hat{w}\|_2^2$$

since they are orthogonal

so \hat{w} is the min norm solution.
and so $\hat{w} = w^*$ from gradient descent.

Q 1.2.1)

Claim: RMSProp does not always obtain the minimum norm solution.

Counterexample:

It is sufficient to find a counterexample where RMSProp converges after one iteration at a \hat{W} that does not have the minimum norm.

$$\text{take } \mathbf{x}_1 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad \mathbf{w}_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad t = 2,$$

and hyperparameters β, ε, η .

$$\text{Then } \nabla_{\mathbf{w}, t=0} \mathcal{L} = 2 \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} 2 \\ 1 \end{bmatrix} - 2 \right) \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} -8 \\ -4 \end{bmatrix}$$

$$\text{and } \mathbf{w}_{t=1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} \frac{\eta}{\sqrt{\beta(0)+(1-\beta)(-8)^2}} \cdot (-8) \\ \frac{\eta}{\sqrt{\beta(0)+(1-\beta)(-4)^2}} \cdot (-4) \end{bmatrix} = \begin{bmatrix} \frac{\eta}{\sqrt{1-\beta}} \\ \frac{\eta}{\sqrt{1-\beta}} \end{bmatrix} \quad \text{assuming } \varepsilon=0$$

$$\Rightarrow \nabla_{\mathbf{w}, t=1} \mathcal{L} = 2 \left(\begin{bmatrix} \frac{\eta}{\sqrt{1-\beta}} \\ \frac{\eta}{\sqrt{1-\beta}} \end{bmatrix}^T \begin{bmatrix} 2 \\ 1 \end{bmatrix} - 2 \right) \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \left(\frac{4\eta}{\sqrt{1-\beta}} + \frac{2\eta}{\sqrt{1-\beta}} - 4 \right) \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

assuming convergence here \Rightarrow setting $\nabla_{\mathbf{w}, t=1} \mathcal{L} = 0$.

$$\Rightarrow \left(\frac{4\eta}{\sqrt{1-\beta}} + \frac{2\eta}{\sqrt{1-\beta}} - 4 \right) \begin{bmatrix} 2 \\ 1 \end{bmatrix} = 0$$

$$\Rightarrow \frac{4\eta}{\sqrt{1-\beta}} + \frac{2\eta}{\sqrt{1-\beta}} - 4 = 0$$

$$\Rightarrow \frac{4\eta}{\sqrt{1-\beta}} + \frac{2\eta}{\sqrt{1-\beta}} = 4$$

$$\Rightarrow \frac{6\eta}{\sqrt{1-\beta}} = 4$$

$$\Rightarrow \eta = \frac{4\sqrt{1-\beta}}{6}$$

So if we take $\beta = 0.99$, and $\eta = \frac{1}{15}$

$$\text{then } \mathbf{w}_{t=1} = \begin{bmatrix} \frac{1}{\sqrt{1-0.99}} \\ \frac{1}{\sqrt{1-0.99}} \end{bmatrix} = \begin{bmatrix} \frac{2}{3} \\ \frac{2}{3} \end{bmatrix}, \quad \text{and } \nabla_{\mathbf{w}, t=1} \mathcal{L} = 0$$

$$\Rightarrow \mathbf{w}_{t=2} = \mathbf{w}_{t=1} \quad \text{so we have reached convergence with } \hat{W} = \begin{bmatrix} \frac{2}{3} \\ \frac{2}{3} \end{bmatrix}$$

However, this is not the minimum norm weights since

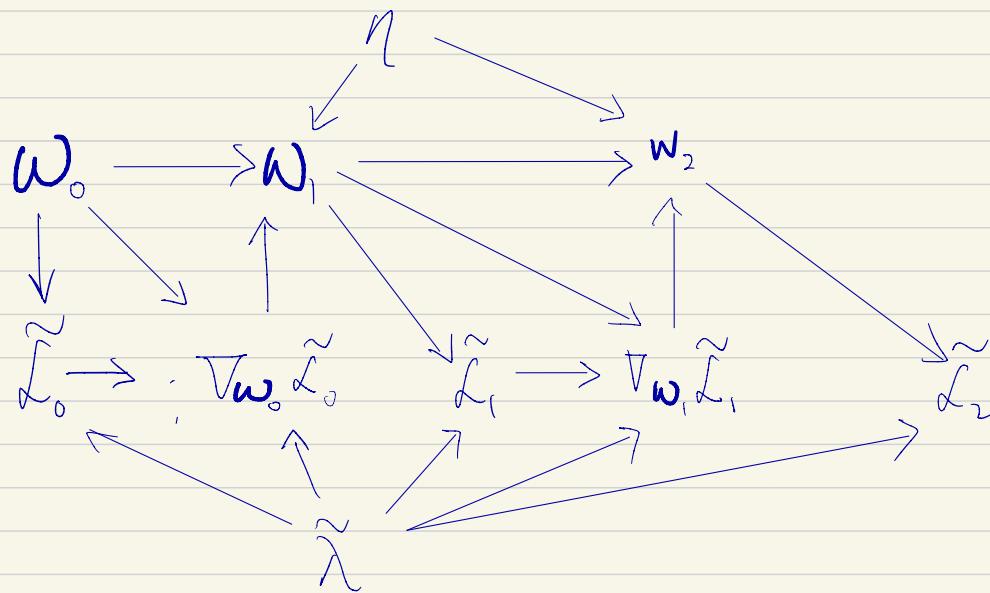
$$\text{the norm of the gradient descent solution } \|\mathbf{w}^*\|_2^2 = \left\| \begin{bmatrix} \frac{4}{3} \\ \frac{2}{3} \end{bmatrix} \right\|_2^2 = \frac{4}{3} < \|\hat{W}\|_2^2 = \left\| \begin{bmatrix} \frac{2}{3} \\ \frac{2}{3} \end{bmatrix} \right\|_2^2 = \frac{8}{9}.$$

This is because the gradient updates for RMSProp do not necessarily lie in the span of X .

2.1.1)

$$\mathcal{L} = \frac{1}{n} \|X\hat{w} - t\|_2^2 + \tilde{\lambda} \|\hat{w}\|_2^2$$

$$\frac{\partial \mathcal{L}}{\partial \hat{w}} = \frac{2}{n} X^T (X\hat{w} - t) + 2\tilde{\lambda} \hat{w} = \frac{2}{n} (X^T X \hat{w} - X^T t) + 2\tilde{\lambda} \hat{w}$$



2.1.2)

forward : O(1)

since each iteration only needs to access results from the previous iteration.

backward : O(t)

since each iteration, $\frac{\partial \mathcal{L}_t}{\partial \eta}$ requires information from all the iterations prior (ex, needs $w_1, w_2 \dots w_t$).

2.2.1)

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

$$\begin{aligned}\mathbf{w}_1 &= \mathbf{w}_0 - \eta \nabla_{\mathbf{w}_0} \tilde{\mathcal{L}}_0 \\ &= \mathbf{w}_0 - \eta \left(\frac{2}{n} \mathbf{X}^\top (\mathbf{X} \mathbf{w}_0 - \mathbf{t}) \right) \\ &= \mathbf{w}_0 - \frac{2\eta}{n} \mathbf{X}^\top \mathbf{a} \quad , \text{ where } \mathbf{a} = \mathbf{X} \mathbf{w}_0 - \mathbf{t}\end{aligned}$$

$$\begin{aligned}\mathcal{L}_1 &= \frac{1}{n} \| \mathbf{X} \mathbf{w}_1 - \mathbf{t} \|_2^2 \\ &= \frac{1}{n} \| \mathbf{X} \left(\mathbf{w}_0 - \frac{2\eta}{n} \mathbf{X}^\top \mathbf{a} \right) - \mathbf{t} \|_2^2 \\ &= \frac{1}{n} \| \mathbf{X} \mathbf{w}_0 - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} - \mathbf{t} \|_2^2 \\ &= \frac{1}{n} \| \mathbf{X} \mathbf{w}_0 - \mathbf{t} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \|_2^2 \\ &= \frac{1}{n} \| \mathbf{a} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \|_2^2 \\ &= \frac{1}{n} \left(\mathbf{a} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \right)^\top \left(\mathbf{a} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \right) \\ &= \frac{1}{n} \left(\mathbf{a}^\top - \left(\frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \right)^\top \right) \left(\mathbf{a} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \right) \\ &= \frac{1}{n} \left(\mathbf{a}^\top - \frac{2\eta}{n} \mathbf{a}^\top \mathbf{X}^\top \right) \left(\mathbf{a} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \mathbf{a} \right) \\ &= \frac{1}{n} \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right) \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right) \mathbf{a} \\ &= \frac{1}{n} \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^2 \mathbf{a}\end{aligned}$$

2.2.3).

$$\begin{aligned}\frac{\partial \mathcal{L}_1}{\eta} &= \frac{2}{\eta} \left(\frac{1}{n} \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right)^2 \mathbf{a} \right) \\ &= \frac{2}{n} \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right) \left(-\frac{2}{n} \mathbf{X} \mathbf{X}^\top \right) \mathbf{a} \\ &= \frac{4}{n^2} \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta}{n} \mathbf{X} \mathbf{X}^\top \right) (-\mathbf{X} \mathbf{X}^\top) \mathbf{a}\end{aligned}$$

$\stackrel{\text{set } 0}{=}$

$$\Rightarrow \mathbf{a}^\top \left(\mathbf{I} - \frac{2\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right) (-\mathbf{X} \mathbf{X}^\top) \mathbf{a} = 0$$

$$\Rightarrow \mathbf{a}^\top (-\mathbf{X} \mathbf{X}^\top) \mathbf{a} + \mathbf{a}^\top \left(-\frac{2\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right) (-\mathbf{X} \mathbf{X}^\top) \mathbf{a} = 0$$

$$\Rightarrow -\mathbf{a}^\top (\mathbf{X} \mathbf{X}^\top) \mathbf{a} + \mathbf{a}^\top \left(\frac{2\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right) (\mathbf{X} \mathbf{X}^\top) \mathbf{a} = 0$$

$$\Rightarrow \mathbf{a}^\top \left(\frac{2\eta^*}{n} \mathbf{X} \mathbf{X}^\top \right) (\mathbf{X} \mathbf{X}^\top) \mathbf{a} = \mathbf{a}^\top (\mathbf{X} \mathbf{X}^\top) \mathbf{a}$$

$$\Rightarrow \frac{2\eta^*}{n} \mathbf{a}^\top \mathbf{X} \mathbf{X}^\top \mathbf{X} \mathbf{X}^\top \mathbf{a} = \mathbf{a}^\top \mathbf{X} \mathbf{X}^\top \mathbf{a}$$

$$\Rightarrow \frac{2\eta^*}{n} \|\mathbf{X} \mathbf{X}^\top \mathbf{a}\|_2^2 = \|\mathbf{X}^\top \mathbf{a}\|_2^2$$

$$\Rightarrow \frac{2\eta^*}{n} = \frac{\|\mathbf{X}^\top \mathbf{a}\|_2^2}{\|\mathbf{X} \mathbf{X}^\top \mathbf{a}\|_2^2} \quad (\text{element wise division})$$

$$\Rightarrow \eta^* = \frac{n}{2} \cdot \frac{\|\mathbf{X}^\top \mathbf{a}\|_2^2}{\|\mathbf{X} \mathbf{X}^\top \mathbf{a}\|_2^2} \quad , \text{ where } \mathbf{a} = \mathbf{X} \mathbf{w}_0 - \mathbf{t}$$

2.3.1)

$$\tilde{\mathcal{L}} = \frac{1}{n} \|Xw - t\|_2^2 + \tilde{\lambda} \|w\|_2^2, \quad \mathcal{L} = \frac{1}{n} \|Xw - t\|_2^2$$

$$\Rightarrow \frac{\partial \tilde{\mathcal{L}}}{\partial w} = \frac{2}{n} X^T (Xw - t) + 2\tilde{\lambda} w, \quad \frac{\partial \mathcal{L}}{\partial w} = \frac{2}{n} X^T (Xw - t)$$

$$\Rightarrow \textcircled{1} \text{ L2 regularization: } w_1 = w_0 - \eta \nabla_{w_0} \tilde{\mathcal{L}}$$

$$= w_0 - \eta \left(\frac{2}{n} X^T (Xw_0 - t) + 2\tilde{\lambda} w_0 \right)$$

$$= w_0 - \eta \left(\frac{2}{n} X^T a + 2\tilde{\lambda} w_0 \right) \quad , \text{ where } a = Xw_0 - t$$

and \textcircled{2} weight decay: $w_1 = (1-\lambda) w_0 - \eta \nabla_{w_0} \mathcal{L}$

$$= (1-\lambda) w_0 - \eta \left(\frac{2}{n} X^T (Xw_0 - t) \right)$$

$$= (1-\lambda) w_0 - \frac{2\eta}{n} X^T a \quad , \text{ where } a = Xw_0 - t$$

2.3.2)

$$\text{set } w_0 - \eta \left(\frac{2}{n} X^T a + 2\tilde{\lambda} w_0 \right) = (1-\lambda) w_0 - \frac{2\eta}{n} X^T a$$

$$\Rightarrow w_0 - \frac{2\eta}{n} X^T a - 2\eta \tilde{\lambda} w_0 = w_0 - \lambda w_0 - \frac{2\eta}{n} X^T a$$

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

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

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

Q3.1

$$I = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \end{bmatrix} \quad J = \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$$

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

This filter detects horizontal edges. Negative values in output means original image is dark above and light below, while positive values in output means original image is light above and dark below.

(large values) (small values)
 (small values) (large values)

Q3.2)

assuming no bias.

assuming that we don't count input image as # neurons.

CNN

$$\# \text{neurons} : C_1 + P_1 + C_2 + P_2 + C_3 = 1664$$

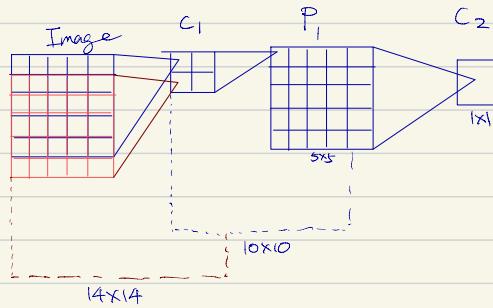
$$\# \text{trainable parameters} : k_1 + k_2 + k_3 = 27$$

FCNN :

$$\# \text{neurons} : FC_1 + P_1 + FC_2 + P_2 + FC_3 = 1664$$

$$\# \text{trainable parameters} : W_1 + W_2 + W_3 = 1,118,208$$

(Q3.3)



Assuming a filter size of 5×5 for the convolution layers, then:

→ every neuron after the second convolution layer has a receptive field of a 5×5 patch of neurons from the result of the first pooling layer.

→ every neuron from the result of the first pooling layer has a receptive field of a 2×2 patch of neurons from the result of the first convolution layer.

↳ since there are 5×5 neurons from the result of the first pooling layer to account for, the receptive field for the result of the first convolution layer must be a 10×10 patch of neurons. (since pooling has stride=2, no overlap).

→ the receptive field of the 10×10 patch of neurons on the result of the first convolution layer, must then have a receptive field of a 14×14 neuron patch since the convolution has filter size 5×5 and stride 1.

$$\text{since } 10 = (W - 5) + 1 \Rightarrow W = 14$$

⇒ So in conclusion, the receptive field of a neuron after the second convolution layer is 14×14 .

Other things that can affect the size of the receptive field:

- ① the stride of the max pooling layer
- ② the stride of the convolution filter in convolution layers