

## 2. Linear Neural Networks

(a) ① If there exists  $W_1, W_2, \dots, W_L$  such that  $\bar{W} = W_0$  then  
 $\text{rank}(W_0) \leq \min_{i=0}^L d_i$ .

Proof:  $\because W_0 = \bar{W} = W_L \times W_{L-1} \times \dots \times W_1$

$$\therefore \text{rank}(W_0) = \min(\text{rank}(W_L), \text{rank}(W_{L-1}), \dots, \text{rank}(W_1))$$

$$\therefore \text{rank}(W_L) \leq \min(d_L, d_{L-1})$$

$$\text{rank}(W_{L-1}) \leq \min(d_{L-1}, d_{L-2})$$

$\vdots$

$$\text{rank}(W_1) \leq \min(d_1, d_0)$$

$$\therefore \text{rank}(W_0) \leq \min_{i=0}^L d_i$$

② If  $\text{rank}(W_0) \leq \min_{i=0}^L d_i$  then there exists  $W_1, W_2, \dots, W_L$  such that  $\bar{W} = W_0$

Proof: Assume  $\text{rank}(W_0) = d^* \leq \min_{i=0}^L d_i$

$$\text{rank}(W_i) = R_i$$

$$W_0 = U_0 \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{d^*} \\ & & 0 & \end{pmatrix} V_0'$$

$p \times p \quad p \times d \quad d \times d$

$$\bar{W} = W_L \times W_{L-1} \times \dots \times W_1$$

$$= U_L \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_L} \\ & & 0 & \end{pmatrix} \underbrace{\left( V_L \times U_{L-1} \right)}_{\substack{d_L \times d_{L-1} \\ d_{L-1} \times d_{L-2}}} \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_{L-1}} \\ & & 0 & \end{pmatrix} \underbrace{\left( V_{L-1} \times U_{L-2} \right)}_{d_{L-1} \times d_{L-2}} \dots$$

$p \times d_{L-1} \quad = I_{d_{L-1}} \quad = I_{d_{L-2}}$

$$\dots \underbrace{\left( V_2' \times U_1 \right)}_{d_2 \times d_1} \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_1} \\ & & 0 & \end{pmatrix} V_1'$$

$= I_{d_1}$

$$= U_L \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_L} \\ & & 0 & \end{pmatrix} \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_{L-1}} \\ & & 0 & \end{pmatrix} \dots \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ 0 & & s_{R_1} \\ & & 0 & \end{pmatrix} V_1'$$

$p \times p \quad p \times d \quad d \times d$

$$= U_L \begin{pmatrix} s_1 & & 0 \\ & \ddots & \\ & & s_d' \\ 0 & & 0 \end{pmatrix} V_L'$$

$$d' = \min_i (R_i) \leq \min_i d_i$$

$$\therefore d^* \leq \min_i d_i$$

$$\therefore d' \text{ can be } d^*$$

$$\therefore \bar{W} \text{ can be equal to } W_0$$

$$\begin{aligned} (b) \quad \hat{P}(f_{\bar{W}, L}) &= \frac{1}{2} \sum_{i=1}^n \|f_{\bar{W}, L}(x_i) - y_i\|_2^2 \\ &= \frac{1}{2} \sum_{i=1}^n \|\bar{W} x_i - y_i\|_2^2 \\ &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^p \|\bar{W}[j] x_i - y_i\|_2^2 \\ &= \frac{1}{2} \|\bar{W} X^T - Y^T\|_F^2 = Q(\bar{W}) \end{aligned}$$

$$Q(\bar{W}) = \frac{1}{2} \|\bar{W} X^T - Y^T\|_F^2$$

$$= \sum_{j=1}^p \frac{1}{2} \|\bar{W}[j] x^T - Y^T[j]\|_2^2$$

$$\nabla_{\bar{W}[j]} \left( \frac{1}{2} \|\bar{W}[j] x^T - Y^T[j]\|_2^2 \right)$$

$$= \frac{1}{2} \nabla_{\bar{W}[j]} \left[ (\bar{W}[j] x^T - Y^T[j]) (\bar{W}[j] x^T - Y^T[j])^T \right]$$

$$= \frac{1}{2} \nabla_{\bar{W}[j]} \left( \bar{W}[j] x^T x \bar{W}[j]^T - 2 \bar{W}[j] x^T Y[j] + Y[j] Y[j]^T \right)$$

$$= \frac{1}{2} (2 x^T x \bar{W}[j]^T - 2 x^T Y[j])$$

$$= x^T x \bar{W}[j]^T - x^T Y[j]$$

$$\nabla_{\bar{W}[j]} = x^T x$$

$$\therefore \frac{\partial^2 Q(\bar{W})}{\partial \bar{W}^2} = (x^T x)^T = x^T x$$

This is a PSD matrix

Because  $V^T X^T X V = \|XV\|_2^2 \geq 0$   
for any vector  $V$



$\therefore \frac{1}{2} \| \bar{W} [j] X^T - Y^T [j] \|_2^2$  is a convex function of  $\bar{w}_{[j]}$   
 $\therefore Q(\bar{W})$  is a convex function of  $\bar{W}$

$$(c) \quad Q(\bar{W}) = \frac{1}{2} \times \left\| \begin{pmatrix} \bar{W} [1] \\ \bar{W} [2] \\ \vdots \\ \bar{W} [P] \end{pmatrix} (x_1, x_2, \dots, x_n) - (y_1, y_2, \dots, y_n) \right\|_F^2$$

$$\therefore \frac{\partial Q(\bar{W})}{\partial \bar{W}} = \left( \nabla_{\bar{W}} Q(\bar{W}) \right)^T = (\bar{W} X^T X - Y^T X)^T$$

$$= X^T X \bar{W}^T - X^T Y$$

$$(d) \quad \therefore \min_{i=1,2} d_i = \min \{P, d\}$$

$$W_0 = \bar{W}$$

$$Q(W_0) = \frac{1}{2} \| W_0 X^T - Y^T \|_F^2$$

$$\therefore \nabla_{W_0} Q(W_0) = W_0 X^T X - Y^T X$$

$$= \bar{W} X^T X - Y^T X = 0$$

$\therefore Q(W_0)$  is a convex function of  $W_0$

$\therefore W_0 = \bar{W}$  is a global minimizer of  $Q(\bar{W})$

and  $\hat{R}(f_{w_{1:L}})$

$$(e) \quad \frac{\partial \hat{R}(f_{w_{1:L}})}{\partial W_i} = \bar{W}_{(L+1:H)}^T \frac{\partial \hat{R}(f_{w_{1:L}})}{\partial \bar{W}} \bar{W}_{(1:L)}^T$$

If we cannot guarantee that  $\bar{W} = W_0$  then we cannot guarantee the convexity of  $\hat{R}(f_{w_{1:L}})$

$\therefore$  The critical point of  $\hat{R}$  may not be Global minimum.

(f) The sufficient condition can be All  $W_1, \dots, W_L$  are of full rank.

(a) The relevant code is in the appendix

(b) The relevant code and graph are in the appendix  
 $d=10$  takes more iterations to decrease the loss to 0 than  $d=5$ .

But when  $d=30$ , The loss will oscillate around 1 and will not go to 0 at the end of 5000 iterations.

(c) (i)  $\nabla R(p) = \nabla E_x \{ \max\{0, 1 - f_p(x) \cdot (-1)^{y^T x}\} \}$

$$= E_x \{ \nabla_p \{ \max\{0, 1 - f_p(x) \cdot (-1)^{y^T x}\} \} \}$$

$$= E_x \{ \nabla_p [(-1)^{y^T x} - f_p(x)] \cdot 1_{(f_p(x) \cdot (-1)^{y^T x})} \}$$

$$= E_x \{ \nabla_p [y - f_p(x)] \cdot 1_{(f_p(x) \cdot y)} \}$$

$$= E_{x,y} \{ \nabla_p [y \cdot 1_{(f_p(x) \cdot y)}] - \nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot y)}] \}$$

$$= \underbrace{-E_{x,y} [\nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot y)}]}_{\downarrow} + E_{x,y} \left[ y \underbrace{\nabla_p [1_{(f_p(x) \cdot y)}]}_{\downarrow} \right]$$

$$= -E_{x,y} [\nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot y)}]] + E_x \{ (-1)^{y^T x} \cdot \underbrace{\nabla_p [1_{(f_p(x) \cdot (-1)^{y^T x})}]}_{\downarrow g(x)} \}$$

$$= -E_{x,y} [\nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot y)}]]$$

$$= -E_x [\nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot (-1)^{y^T x})}]]$$

$$= - \text{When } (-1)^{y^T x} = 1:$$

$$1_{(f_p(x) \cdot (-1)^{y^T x})} = 1_{(f_p(x))}$$

$$\text{When } (-1)^{y^T x} = -1:$$

$$1_{(f_p(x) \cdot (-1)^{y^T x})} = 1_{(-f_p(x))}$$

$$\therefore -E_x [\nabla_p [f_p(x) \cdot 1_{(f_p(x) \cdot (-1)^{y^T x})}]]$$

$$= -E_x \left[ \frac{1_{(f_p(x))} - 1_{(-f_p(x))}}{1 - (-1)} \cdot \nabla_p [f_p(x)] \right]$$

$$= -E_x \left[ \frac{1_{(f_p(x))} - 1_{(-f_p(x))}}{2} \cdot \nabla_p [f_p(x)] \right]$$

$\downarrow$   
 $a$

$$\therefore \nabla R(p) = a + E_x [(-1)^{y^T x} g(x)] = \nabla_p [1_{(f_p(x) \cdot (-1)^{y^T x})}]$$

where  $g(x)$



C (i'i')

$$g(x) = \nabla_p \mathbb{I}_1 [f_p(x) \cdot (-1)^{v^T x}]$$

$$E_v [1 \otimes R(p) - a]^2 = \frac{1}{2^d} \sum_v (E_x ((-1)^{v^T x} g(x)))^2 \leq \frac{E_x [g(x)^2]}{2^d}$$

$v$  is  $(0,1)^d$   
 $\therefore 2^d$  possible  $v$

$$g(x) = \nabla_p \mathbb{I}_1 [f_p(x) \cdot (-1)^{v^T x}]$$

$$= -\frac{1}{2} \nabla_p f_p(x) \cdot [\mathbb{I}_1(f_p(x)) + \mathbb{I}_1(-f_p(x))]$$

$$g(x)^2 \leq (\nabla_p f_p(x) \cdot \frac{\mathbb{I}_1(f_p(x)) + \mathbb{I}_1(-f_p(x))}{2})^2$$

$$\because \mathbb{I}_1(f_p(x)) + \mathbb{I}_1(-f_p(x)) \leq 2$$

$$\therefore g(x)^2 \leq (\nabla_p f_p(x))^2 = \sum_j (\nabla_{p_j} f_p(x)_j)^2$$

$$\leq d C_p^2$$

$$\therefore E_v [1 \otimes R(p) - a]^2 \leq \frac{E_x [g(x)^2]}{2^d} \leq \frac{d C_p^2}{2^d} \therefore C_p = d C_p^2$$

(d) As we go deeper, the gradient will be the product of more gradients, which will be more unstable.

And for sigmoid activation function, it has 0 gradient at top and bottom.

As a result, we will highly likely stuck at 0 gradient.

Merely because there are more gradients time together, easier to get 0 gradient even when  $d$  is small.