H.W. 0

**1.**

Linear Regression Problem with $n$-training points and $d$-features when $n = d$, the feature matrix becomes $\mathbb{R}^{n \times n}$ where $F$ has max singular value $\alpha$ and a very tiny min. singular value.

for $y = Fw^* + e$

$$\hat{w}_{inv} = F^{-1}y \quad \Rightarrow \quad \|\hat{w}_{inv} - w^*\|_2^2 = 10^{10}$$
$$\downarrow$$

Taking the Gradient Descent Approach,

$$\ell(w) = \frac{1}{2} \|Y - Fw\|^2$$

Starting from $w_0 = 0$

Gradient Descent Update $\Rightarrow$

$$w_t = w_{t-1} - \eta \left( F^T (F w_{t-1} - y) \right)$$

we are interested in minimizing $\|w_k - w^*\|_2^2$

Because

$$F = V \Lambda V^{-1}$$
$$F^{-1} = V \Lambda^{-1} V$$

so,

$$\Lambda^{-1} = \begin{bmatrix} 1/r_1 & 1/r_2 \\ & & \ddots \\ & & & 1/r_n \end{bmatrix}$$
$$\Downarrow$$

Small singular value makes the reciprocal very large!

Show that: for $t > 0$, $\|w_t\|_2 \leq \|w_{t-1}\|_2 + \eta \alpha \|y\|_2$

So, we know that $\cancel{w_{t+1}}$ $w_t = w_{t-1} - \eta F^T (F w_{t-1} - y))$

$$\cdots\cdots\cdots\cdots \quad w_t = w_{t-1} - \eta F^T F w_{t-1} + \eta F^T y$$

$w_t = w_{t-1} (I - \eta F^T F) + \eta F^T y$

take norm (L-2) on both sides,

$$\|w_t\|_2 = \|w_{t-1}(I - \eta F^T F)\|_2 + \eta \|F^T y\|_2$$
$$\|w_t\|_2 \leq \|w_{t-1}\|_2 \|I - \eta F^T F\|_2 + \eta \|F^T y\|_2$$

Recall, for gradient descent to converges we need to look more closely at $\|I - F^T F v\|$

If $F \in \mathbb{R}^{n \times n}$ then, $F = V \Lambda V^T$ where $\max(\Lambda) = \alpha$.

And, $F^T F = (V \Lambda V^T)^T (V \Lambda V^T)$

$$= V \Lambda^T V^T V \Lambda V^T \quad \text{where } V = \text{orthonormal basis}$$

$$F^T F = V \Lambda^2 V^T$$

$$\|F^T F\|_2 = \sqrt{\|\Lambda\|^4} = \|\Lambda\|_2^2.$$

So, $\|I - F^T F \eta\| < 1$

$$\|w_t\|_2 \leq \|w_{t-1}\|_2 + \eta \|F^T y\|_2$$

$$\|w_t\|_2 \leq \|w_{t-1}\|_2 + \eta \|V \Lambda V^T\|_2 \|y\|_2$$

$$\|w_t\|_2 \leq \|w_{t-1}\|_2 + \eta \alpha \|y\|_2 \quad \text{since} \quad \max \|V \Lambda V^T\| \text{ is } \alpha \text{ i.e.,}$$

All diagonal entries have $\alpha$s.

hence shown !

<u>2.</u>    Show that    ①.

$$\boxed{\arg \min_{w} \| \hat{y} - \hat{X} w \|_2^2}$$
has the same solution

as    $$\boxed{\hat{w} = (X^T X + \Sigma^{-1})^{-1} X^T y} \quad - ⓘⓘ.$$

and    $$\hat{X} = \begin{bmatrix} X \\ \Gamma \end{bmatrix} , \quad \hat{y} = \begin{bmatrix} y \\ 0_d \end{bmatrix}$$

$(n+d) \times d$        $(n+d) \times 1$

from ①, if we differentiate w.r.t. w,

$$2(\hat{y} - \hat{X} w) \cdot \hat{X}^T = 0.$$

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

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

now,

$$\left( \begin{bmatrix} X^T & \Gamma^T \end{bmatrix} \begin{bmatrix} X \\ \Gamma \end{bmatrix} \right)^{-1} \begin{bmatrix} X^T & \Gamma^T \end{bmatrix} \begin{bmatrix} y \\ 0_d \end{bmatrix}$$

$d \times (n \times d)$  $(n \times d) \times d$        $d \times n \times d$  $n \times d \times 1$

Performing blockwise matrix multiplication,

$$\hat{w} = (X^T X + \Gamma^T \Gamma)^{-1} [X^T y + \Gamma^T 0\lambda]$$

$$\hat{w} = (X^T X + \Sigma_j^{-1})^{-1} [X^T y + \bar{0}]$$

$$\hat{w} = (X^T X + \Sigma_j^{-1})^{-1} X^T y \quad \text{— hence proven.}$$

3. VECTOR CALCULUS

(a) show that $\quad \dfrac{\partial}{\partial x} \underset{1 \times n \ \ n \times 1}{(X^T C)} = C^T$

$$= \frac{\partial}{\partial x}\left( \underset{1 \times n}{[X_1 \ X_2 \ \cdots \ X_n]} \begin{bmatrix} C_1 \\ C_2 \\ \vdots \\ C_N \end{bmatrix}_{n \times 1} \right)$$

$$= \frac{\partial}{\partial x}\left[ \underset{1 \times 1}{X_1 C_1 + X_2 C_2 + X_3 C_3 + \cdots X_n C_N} \right]$$

$\longrightarrow$ Now, as $X \in \mathbb{R}^n$,

$$= \frac{\partial}{\partial x}\left[ \frac{\partial X_1 C_1 + X_2 C_2 + \cdots X_n C_N}{\partial x_1} \quad \frac{\partial}{\partial x_2}(\quad) \ \cdots \ \frac{\partial}{\partial x_n}(\quad) \right]$$

$$= \left[ C_1 \ C_2 \ \cdots \ C_N \right] = C^T \quad \begin{array}{l}\text{hence} \\ \text{proven!}\end{array}$$

(b) $\quad \dfrac{\partial}{\partial x} \|x\|_2^2 = \dfrac{\partial}{\partial x}\underset{1 \times n \ \ n \times 1}{\langle x^T x \rangle}$

$$= \frac{\partial}{\partial x}(x_1^2 + x_2^2 + \cdots x_n^2)$$

$$= [2x_1, \ 2x_2, \ \cdots \ 2x_n]$$

$$= 2 x^T$$

hence proven.

(c) $\dfrac{\partial}{\partial x}(Ax) = A$

$$\frac{\partial}{\partial x}\begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1N} \\ A_{21} & A_{22} & \cdots & A_{2N} \\ \vdots & & & \\ A_{N1} & A_{N2} & \cdots & A_{NN} \end{bmatrix}\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \frac{\partial}{\partial x}\begin{bmatrix} A_{11}x_1 + A_{12}x_2 + \cdots A_{1N}x_N \\ A_{21}x_1 + A_{22}x_2 + \cdots A_{2N}x_N \\ \vdots \\ A_{N1}x_1 + A_{N2}x_2 + \cdots A_{NN}x_N \end{bmatrix}$$

$$n \times n \qquad\qquad n \times 1 \qquad\qquad\qquad n \times 1$$

$$= \frac{\partial}{\partial x}Ax = \begin{bmatrix} \frac{\partial}{\partial x_1}(\text{first row}) & \frac{\partial}{\partial x_2}(\text{first row}) & \cdots \\ \frac{\partial}{\partial x_2}(\text{first row}) & & \cdots \\ & & \ddots \end{bmatrix}$$

$$= \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1N} \\ A_{21} & A_{22} & \cdots & A_{2N} \\ \vdots & & & \\ A_{N1} & A_{N2} & \cdots & A_{NN} \end{bmatrix} = A.$$

(d)

$$\frac{\partial}{\partial x}\left(x^T A x\right) = x^T(A + A^T)$$

$$1 \times n \;\; n \times n \;\; n \times 1$$

Using the definition of fundamental theorem:

$$f(x + \Delta) = (x + \Delta)^T A (x + \Delta)$$

$$= x^T A x + \Delta^T A x + x^T A \Delta + \Delta^T A \Delta$$
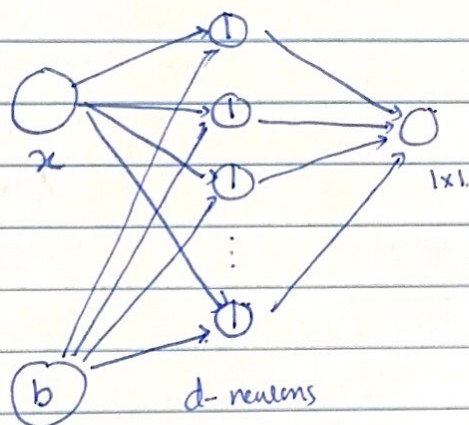
$$= f(x) + (x^T A^T + x^T A)\Delta$$

which yields the derivative as $x^T A^T + x^T A$
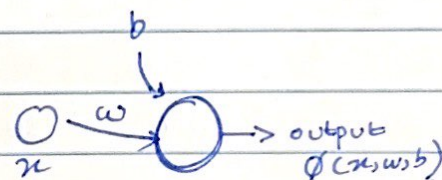
$$\Rightarrow (A^T + A)x^T$$

(e) If $A^T A$ is symmetric.

4. RELU Elbow Update Under SGD.

$$\hat{f}(x) = \underset{1\times1}{W^{(2)}} \underset{1\times d}{\Phi}( \underset{d\times1}{W^{(1)}} \underset{1\times1}{x} + \underset{d\times1}{b})$$



$$\text{loss } \frac{1}{2}\| y - \hat{f}(x)\|_2^2$$

$1\times1$

d- neurons

(a) Starting with ONE RELU: →



output
$\Phi(x, w, b)$

$$\Phi(x) = \begin{cases} wx+b & wx+b > 0 \\ 0 & \text{o.w.} \end{cases}$$

$$\text{loss} = \frac{1}{2}\|\Phi(x) - y\|_2^2$$

(i) Location of the **elbow** $e$ of the function where it transitions from $0$ to something else.

$$\Phi(x) = 0 = wx+b$$
$$x = -\frac{b}{w}$$

(ii)  $\dfrac{\partial l}{\partial \phi}$  $\begin{cases} (\phi(x) - y) & \phi(x) > 0 \\ 0 & o.w. \end{cases}$

(iii)  $l(x, w, b) = \dfrac{1}{2} \| \phi(x) - y \|^2.$

$\dfrac{\partial l}{\partial w} = \dfrac{\partial l}{\partial \phi} \cdot \dfrac{\partial \phi}{\partial w} = \begin{cases} (\phi(x) - y) \circ x & wx + b > 0 \\ 0 & o.w. \end{cases}$

(iv)  $\dfrac{\partial l}{\partial b} = \begin{cases} (\phi(x) - y) & wx + b > 0 \\ 0 & o.w. \end{cases}$

(b)  $\phi(x) - y = 1$

(i)  $\phi(x) = 0$ 　　　　Gradient Descent:

$$w_t = w_{t-1} - \eta (\nabla_w l)$$

If $\phi(x) = 0 \Rightarrow \dfrac{\partial l}{\partial w} = 0 \Rightarrow$ Gradient remains unchanged.

So, slope remains unchanged and since $\dfrac{\partial l}{\partial \phi} = 0$ thus no change in elbow either.

(ii) $\quad w > 0, \; x > 0, \; \phi(x) > 0.$

now, $\quad \phi(x) \geqslant 0$ which means $\quad wx + b > 0$

$$\frac{\partial l}{\partial \phi} = \phi(x) - y = 1$$

Gradient Descent $\Rightarrow$ $W_{new} = W_{old} - \eta \frac{\partial l}{\partial w}.$

$$W_{new} = W_{old} - \eta \, (\phi(x) - y) \cdot x$$

since $x > 0, \; w > 0,$

$$W_{new} = W_{old} - \eta \, (1) \cdot x$$

when $\quad W_{new} < W_{old},$

$\quad wx + b$ would decrease

so slope will decrease!

$$\text{eff elbow} = \frac{-(b + \Delta b)}{(w - \Delta w)} \, .$$

$$= \frac{-\left(b - \frac{\partial l}{\partial b}\right)}{w - \frac{\partial l}{\partial w}}$$

(iii) $\quad w > 0, \; x < 0, \; \phi(x) > 0$

$$\frac{\partial l}{\partial \phi} = \phi(x) - y = 1$$

$$\frac{\partial l}{\partial w} = (\phi(x) - y)(x) = 1(x) < 0$$

so, $\quad W_{new} = W_{old} - \eta \, (1)(x)$

$$W_{old} < W_{new}.$$

$\qquad\qquad (wx + b)$

$W_{new} x \cancel{+b} + b > W_{old} + b$

so, slope increases.

$$e' = \frac{-b + \frac{\partial \theta}{\partial b}}{W - \gamma \frac{\partial \theta}{\partial w}} \quad \Big\} \quad \begin{array}{l} \text{elbow should move} \\ \text{to left} \\ \text{as error decreases.} \end{array}$$

(iv)    $w < 0, \ x > 0, \ \phi(x) > 0.$

$W_{new} = W_{old} - \gamma (\phi(x) - y) \cdot x$

$W_{new} = W_{old} - \gamma (1)(x)$

$W_{new} < W_{old}$

$\hookrightarrow$ becomes more -ve.

So, slope gets smaller

elbow moves to left.

(c)    $\hat{f}(x) = W^{(2)} \Phi (w^{(1)} x + b)$

$$Loss = \frac{1}{2} \| \hat{f}(x) - y \|_2^2.$$

$$Loss = \frac{1}{2} \left[ W_i^2 \Phi(w_i^1 x + b) \right]^2$$

$$\frac{\partial Loss_i}{\partial w_i^2} = \frac{\cancel{2}}{\cancel{2}} \frac{\Phi(w_i^1 x + b)^2}{}$$

$$\frac{\partial Loss}{\partial \phi_i} = (w_i^1 x + b), \quad \frac{\partial Loss}{\partial w_i^1} = (w_i^1 x + b) x.$$

elbow for $i^{th}$ neuron    would be

$$e_i = \frac{-b}{w_i'} \qquad w_i' x + b = 0$$
$$x = \frac{-b}{w_i'}$$

(中)                                                       $x = \frac{-b}{w_i'}$

(d)    $e_i^{i*} \Rightarrow$ new elbow

$$= \frac{-(b - \Delta b)}{(w_i^* - \Delta w_i')}$$

we need to know what is $\Delta b$ and $\Delta w_i'$ after one SGD iteration.

$\Delta b = \frac{\partial l}{\partial b}$      $Loss = \frac{1}{2} \| w^2 \Phi(w' x + b) - y \|^2$

$$\frac{\partial L}{\partial w^2} = (w^2 \Phi(w' x + b) - y) \cdot \Phi(w x + b)$$

$$\frac{\partial L}{\partial w'} = (w^2 \Phi(w' x + b) - y) \cdot x$$

$$\frac{\partial L}{\partial w_i'} = w_i^2 [\Phi(w_i' x + b) - y] \cdot x.$$

$$\frac{\partial L}{\partial b} = w_i^2 [\Phi(w_i' x + b) - y]$$

$$e_i = \frac{-(b_i - \nu w_i^2 (\Phi(w_i' x + b) - y))}{w_i' - \nu w_i^2 x [\Phi(w_i' x + b) - y]}.$$

6    Homework Process & study Group.

(a)    sources =>    eecs127_reader.pdf , EECS16B SVD Notes

(b)    worked Individually

(c)    Number of hours => 5