

Q1

$$E(\text{Score}) = 4 \cdot P(\text{Score}=4) + 3 \cdot P(\text{Score}=3) + 2 \cdot P(\text{Score}=2) + 0.$$

$$4 \cdot P(\text{Score}=4) = 4 \int_0^{\frac{1}{3}} f(x) dx = 4 \cdot \frac{1}{3} = \frac{4}{3}$$

$$3 \cdot P(\text{Score}=3) = 3 \int_{\frac{1}{3}}^{\frac{1}{2}} f(x) dx = 3 \cdot \frac{1}{6} = \frac{1}{2}.$$

$$2 \cdot P(\text{Score}=2) = 2 \int_{\frac{1}{2}}^{\frac{2}{3}} f(x) dx = 2 \cdot \frac{1}{6} = \frac{1}{3}.$$

therefore

Expected score of single shot

$$= \frac{4}{3} + \frac{1}{2} + \frac{1}{3} = \frac{8+3+2}{6} = \frac{13}{6}$$

Q2

n observations  $t_1, t_2, t_3, \dots, t_n$  independent

$$P(X=x_1, X_2=x_2, \dots, X_n=x_n) = P(X=x_1)P(X=x_2) \cdots P(X=x_n)$$

$$= \theta e^{-\theta x_1} \theta e^{-\theta x_2} \cdots \theta e^{-\theta x_n}$$

Want to Maximize this value.  $\theta^n e^{-\theta(\sum x_i)}$

$$\left[ \theta^n e^{-\theta(\sum x_i)} \right]' = n \theta^{n-1} e^{-\theta(\sum x_i)} + \theta^n e^{-\theta(\sum x_i)} \cdot \left( -\sum x_i \right) = 0$$

NR  $n = \theta \sum_{i=1}^n x_i$   $\theta = \frac{n}{\sum_{i=1}^n x_i}$

Q3 a)  $X^T A X = [x_1 \dots x_n] \begin{bmatrix} \sum_{j=1}^n a_{1j} x_j \\ \sum_{j=1}^n a_{2j} x_j \\ \vdots \\ \sum_{j=1}^n a_{nj} x_j \end{bmatrix}$

$$= x_1 \sum_{i=1}^n a_{1i} x_i + x_2 \sum_{i=1}^n a_{2i} x_i + \dots + x_n \sum_{i=1}^n a_{ni} x_i$$

$$= \sum_{j=1}^n \sum_{i=1}^n a_{ji} x_i x_j$$

b) Symmetric then  ~~$A_{ij}$~~   $A_{ij} = A_{ji}$

plug in  $x = e_1, e_2, e_3, \dots$

$$x_i = e_i = \begin{bmatrix} 0 \\ \vdots \\ i \\ 0 \end{bmatrix} \rightarrow i\text{th}$$

$x_i^T A x_i = a_{ii} > 0$  according to the definition of positive matrix

Q4: a) If positive  $X^TAX > 0$

$$X^T(A + \gamma I)X = X^TAX + X^T\gamma IX = X^TAX + \gamma X^TIX$$

from definition.  $X^TAX > 0$  and  $\gamma X^TIX$  also  $> 0$ .

then  $A + \gamma I$  also positive.

b)  $V$  eigenvector and  $\lambda$  the corresponding eigenvalue.  
then  $AV = \lambda V$ .

$$V^TAV = V^T\lambda V = \lambda V^TV > 0 \Leftarrow V^TV > 0$$

and  $V^TV > 0$  then  $\lambda > 0$

$\therefore$  all eigenvalue of  $A$  greater than zero

c)  $A$  positive  $\Rightarrow \forall x \in \mathbb{R}^n \quad X^TAX > 0$ .

then Null space of  $A$  is  $\emptyset$

because if  $\exists x \quad Ax = 0$  then  $X^TAX = 0$   $\text{Ø}$

Null space = 0  $A$  has  $n$  independent column span  $\mathbb{R}^n$   
thus invertible.

d)  $A$  symmetric then  $A$  can be written as  $U\Lambda U^T$ .

~~U~~ consists by ~~eigenvectors~~ and  $\Lambda$  is diagonal, eigenvalues.  
normalized eigenvectors

~~then assume~~  $v_1, v_2, v_3, \dots, v_n$  are  $n$  normalized eigenvectors

~~then~~  ~~$A_{ij} = \sum_{k=1}^n v_{ik} \lambda_k v_{kj}$~~

$\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$

$\Lambda$  diagonal and all entries  $> 0$

then  $\begin{vmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_n \end{vmatrix} = \begin{vmatrix} \sqrt{\lambda_1} & 0 & 0 \\ 0 & \sqrt{\lambda_2} & 0 \\ 0 & 0 & \sqrt{\lambda_n} \end{vmatrix}^2 = (\sqrt{\lambda})^2$

therefore  $A = U\Lambda U^T$

$$= (\Lambda U)^T (\Lambda U)$$

are  $n$  eigenvalues  
all of them are positive

Q5.

a).  $x, a \in \mathbb{R}^n$  get  $\frac{\partial(x^T a)}{\partial x}$

$$f(x) = x^T a = x_1 a_1 + x_2 a_2 + x_3 a_3 + x_4 a_4$$

$$\frac{\partial(x^T a)}{\partial x} = \begin{vmatrix} \frac{\partial f(x)}{\partial x_1} & T \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{vmatrix} = a^T$$

$$(a^T)^T = a$$

b).  $A \in \mathbb{R}^{n \times n}$   $x \in \mathbb{R}^n$   $\frac{\partial(x^T A x)}{\partial x}$

$$Ax = \sum_{i=1}^n A_{ii} x_i$$

$$\frac{\partial f(x)}{\partial x} = x^T A x = x_1 \sum_{i=1}^n A_{1i} x_i + x_2 \sum_{i=1}^n A_{2i} x_i + \dots + x_n \sum_{i=1}^n A_{ni} x_i$$

$$\frac{\partial f(x)}{\partial x_k} = x_1 A_{1k} + x_2 A_{2k} + \dots + x_n A_{nk}$$

~~$$\sum_{i=1}^n A_{ki} x_i + 2A_{kk} x_k$$~~

$$= \sum_{i=1}^n A_{ki} x_i + \sum_{i=1}^n A_{ik} x_i$$

~~$$\frac{\partial f(x)}{\partial x} = \left( \sum_{i=1}^n A_{1i} x_i + \sum_{i=1}^n A_{2i} x_i, \sum_{i=1}^n A_{3i} x_i + \sum_{i=1}^n A_{4i} x_i, \dots \right)$$~~

$$= x^T A^T + x^T A$$

$$(x^T A^T + x^T A)^T = A^T x + A x$$

$$c) (XA)_{ij} = \sum_{k=1}^n X_{ik} A_{kj} \quad \text{trace}(XA) = \sum_{i=1}^n (XA)_{ii}$$

$$= \sum_{i=1}^n \sum_{k=1}^n X_{ik} A_{ki}$$

$$\frac{\partial (XA)}{\partial x_{ab}} = Aba \quad \frac{\partial (XA)}{\partial x} = \begin{vmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \dots & A_{nn} \end{vmatrix} = A^T$$

d) want to proof

$$\sqrt{\bar{x}_i^2} \leq \bar{|x_i|} \leq \sqrt{n} \sqrt{\bar{x}_i^2}$$

$$(\sqrt{\bar{x}_i^2})^2 \leq (\bar{|x_i|})^2$$

$$\bar{x}_i^2 \leq \bar{|x_i|^2} + 2 \sum_{j=1}^{n-1} |x_i| \cdot |x_j| \quad \text{if } i$$

greater or equal to zero

therefore

$$\|x\|_2 \leq \|A\|_1$$

$$\text{want to show } \sum |x_i| \leq \sqrt{n} \sqrt{\bar{x}_i^2}$$

$$(\sum |x_i|)^2 \leq \bar{x}_i^2 \sum_{i=1}^n 1^2$$

$$\sum a_i^2 \sum b_i^2 \geq (\sum a_i b_i)^2$$

$$\bar{x}_i^2 \sum 1^2 \geq (\sum |x_i| \cdot 1)^2$$

Done

$$e) \text{ Maximize } g(x) = x^T z \quad \|z\|_1 \leq 1$$

$$x^T z = x_1 z_1 + x_2 z_2 + \dots + x_n z_n \leq \max(|x_i|) (|z_1| + |z_2| + \dots + |z_n|) \leq \max(|x_i|) \cdot 1$$

$$Q6 a). P(X|W_i) = N(\mu_i \beta^2)(x). \quad P(W_i|x) = \frac{P(x|W_i)P(W_i)}{P(x)}$$

$$\begin{aligned} P(x) &= P(x|W_1)P(W_1) + P(x|W_2)P(W_2) \\ &= \frac{1}{2}N(\mu_1 \beta^2)(x) + N(\mu_2 \beta^2)(x) \end{aligned}$$

$$\begin{aligned} P(W_1|x) &= \frac{\frac{1}{2}N(\mu_1 \beta^2)(x)}{\frac{1}{2}[N(\mu_1 \beta^2)(x) + N(\mu_2 \beta^2)(x)]} \\ &\Rightarrow P(W_2|x) = \frac{\frac{1}{2}N(\mu_2 \beta^2)(x)}{\frac{1}{2}[N(\mu_1 \beta^2)(x) + N(\mu_2 \beta^2)(x)]} \end{aligned}$$

$$\mu_1 < \mu_2$$

$$N(\mu_1 \beta^2)(x) = N(\mu_2 \beta^2)(x)$$

$$\frac{1}{N\beta^2\pi} e^{-\frac{(x-\mu_1)^2}{2\beta^2}} = \frac{1}{N\beta^2\pi} e^{-\frac{(x-\mu_2)^2}{2\beta^2}}$$

$$(\mu - \mu_1)^2 = (\mu - \mu_2)^2 \quad \mu = \frac{\mu_1 + \mu_2}{2}$$

Boundary  $\mu = \frac{\mu_1 + \mu_2}{2}$ . When  $\mu \in W_1$  or  $\mu \in W_2$

$$b) P(\text{error}) = P(\text{misclassified } w_2 | w_1) P(w_1) + P(\text{misclassified as } w_1 | w_2) P(w_2)$$

$$X_C = \frac{w_1 + w_2}{2}$$

$$A = \frac{1}{2} \int_{X_C}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-w_1)^2}{2\sigma^2}} dx$$

$$B = \frac{1}{2} \cdot \int_{-\infty}^{X_C} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-w_2)^2}{2\sigma^2}} dx$$

$$Z_1 = \left( \frac{w_1 - w_2}{\sigma} \right) \quad \text{replace } x \text{ in } A \text{ as } Z_1$$

$$A \Rightarrow \frac{1}{2} \int_{\frac{w_1 - w_2}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

$$Z_2 = \left( \frac{w_1 - w_2}{\sigma} \right) \quad \text{replace } x \text{ in } B \text{ as } Z_2$$

$$B \Rightarrow \frac{1}{2} \int_{-\infty}^{\frac{w_1 - w_2}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

$A + B = \text{P(error)}$

~~$A = \frac{1}{2} \int_{-\infty}^{\frac{w_1 - w_2}{\sigma}} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$~~

$$B = \frac{1}{2} \int_{-\frac{w_1 - w_2}{\sigma}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz \quad \text{since } e^{-\frac{z^2}{2}}$$

$$= A$$

$e^{-\frac{z^2}{2}}$  is even function

therefore

$$P(\text{error}) = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{z^2}{2\sigma^2}} dz$$

$$\begin{aligned}
 Q7. \text{ a). } f(w) &= \frac{1}{2} \|Xw - y\|_2^2 + \frac{\lambda}{2} \|w\|_2^2 \\
 &= \frac{1}{2}(Xw - y)^T(Xw - y) + \frac{\lambda}{2} w^T w \\
 \frac{\partial f(w)}{w} &\Rightarrow (Xw - y) = K_{rw} \quad \frac{\partial K_{rw}}{\partial w} = 2K^T \cdot X \\
 \text{therefore } \frac{\partial f(w)}{w} &= \frac{1}{2} \cdot 2(Xw - y)^T \cdot X + \frac{\lambda}{2} \cdot 2w^T \\
 &= (Xw - y)^T X + \lambda w^T
 \end{aligned}$$

$$\begin{aligned}
 \text{b). want derivative } \frac{\partial}{w} &= \lambda w^T + (Xw - y)^T \cancel{f} \\
 &= \lambda w^T + w^T X^T \cancel{y} - y^T X \\
 \text{② } \cancel{y^T X} &= w^T (\lambda I + X^T \cancel{f}) \\
 \cancel{f^T y} &= (\lambda I + X^T \cancel{f})^T w, \\
 \cancel{f^T y} &= (\lambda I + X^T \cancel{f}) \cdot w. \\
 \text{since } \lambda > 0 \text{ therefore } (\lambda I + X^T \cancel{f}) &\text{ is invertible} \\
 \Rightarrow (\lambda I + X^T \cancel{f})^{-1} \cancel{f^T y} &= w.
 \end{aligned}$$

the first term

$$\frac{1}{2} \|Xw - y\|_2^2 = \frac{1}{2}(Xw - y)^T(Xw - y) \quad w = w_n + X^T a.$$

$$Xw = Xw_n + X^T a \quad \text{since } w_n \in \text{Null}(X)$$

$$Xw = X^T a.$$

$$\text{therefore } \frac{1}{2} \|Xw - y\|_2^2 = \frac{1}{2}(X^T a - y)^T(X^T a - y) \quad \text{only depends on } a.$$

$$\text{d7. } \frac{\lambda}{2} w^T w \quad w = w_n + X^T a$$

$$\frac{\lambda}{2} (w_n^T + X^T a)(w_n + X^T a)$$

$$= \frac{\lambda}{2} (w_n^T w_n + \underline{w_n^T X^T a} + \underline{a^T w_n} + a^T \cancel{X^T a}).$$

$$w_n^T \cancel{X^T a} = w_n \cdot \cancel{X^T a}$$

$$a^T w_n = (X a)^T \cdot w_n = \cancel{X^T a} \cdot w_n = 0$$

$$\text{therefore } \frac{1}{2} \|W\|_2^2 = \frac{\lambda}{2} [W_n W_n^T + (\tilde{f}^T a)^T (\tilde{f}^T a)]$$

$\geq 0 \quad \geq 0.$

$$\text{if a fixed } \frac{1}{2} \|W\|_2^2 \geq \frac{\lambda}{2} (\tilde{f}^T a)^T (\tilde{f}^T a).$$

equal when  $W_n = 0$

e7. minimize  $\frac{1}{2} \|X^T a - y\|_2^2 + \frac{\lambda}{2} \|A^T a\|_2^2 \Rightarrow \text{depend only on } a$   
since from c and d. We find that

App. the first term of (1) not depend on  $W_n$  and the second term of (1) is minimized (over  $W_n$ ) when  $W_n = 0$ .

e7. then take derivative of  $\frac{1}{2} \|X^T a - y\|_2^2 + \frac{\lambda}{2} \|A^T a\|_2^2$ . (over  $a$ )

$$\begin{aligned} & \frac{\partial}{\partial a} \left[ \frac{1}{2} (a^T X^T - y^T) (X^T a - y) + \frac{\lambda}{2} (A^T a)^T (A^T a) \right] \\ &= \frac{\partial}{\partial a} \left[ \frac{1}{2} (a^T X^T X^T a - a^T X^T y - y^T X^T a + y^T y) + \frac{\lambda}{2} a^T A^T a \right] \\ &= \frac{1}{2} (a^T X^T X^T a + a^T X^T X^T - y^T X^T a - y^T X^T) + \frac{\lambda}{2} a^T A^T a \end{aligned}$$

$X^T = A$       Asymmetric

$$\Rightarrow \text{then } a^T A^2 - y^T A + \lambda a^T A = 0$$

$$A^T (A^2 + \lambda I) = y^T A \quad A \text{ invertible}$$

$$(A^T A + \lambda I) A = A^T y$$

$$A^T (A^T A + \lambda I) A = A^T y$$

$$\Rightarrow (A^T A + \lambda I) A = y \quad A^T = X^T X \text{ semipositive.}$$

→ Proof

By Q.F. part a.

$$V^T X^T X V = \|XV\|_2^2 \geq 0.$$

$(A^T A + \lambda I) \text{ invertible}$

$$\text{since } A^T = X^T X \text{ semipositive therefore } A = (A^T A + \lambda I)^{-1} y$$

#7

97. find  $\hat{A}x = (X\hat{X}^T \lambda I)^{-1} Y$

$X = n \times d$ .

get  $X\hat{X}^T$  takes  $n \times n \times d$ .

Inverse takes  $n^3$ .

$y = n \times 1$

then takes  $n^2$

$\Theta(n^3 + n^2d) \cancel{\Theta}$

$= \Theta(n^3 + n^2d)$  if  $n, d$  huge

find  $\hat{W}x = (\lambda I - X^T X)^{-1} X^T Y$

$X^T$  takes  $d \times n$

$X^T X$  takes  $d^2n$

Inverse takes  $d^3$ .

$X^T Y$  takes  $d \times n$

Multiply  $(\lambda I - X^T X)^{-1}$  and  $X^T Y$   
 $d \times d$                                     $d \times 1$

takes  $-d^2$

Therefore  $\Theta(d^2 + nd + d^3 + d^2n)$  if  $d, n$  large  
 $= \Theta(d^3 + d^2n)$

Compare

$\hat{A}x: \Theta(n^3 + n^2d)$  and  $\hat{W}x: \Theta(d^3 + d^2n)$

if  $n \geq d$  Should choose  $\hat{W}x$

if  $n \ll d$   $\hat{A}x$

(Q8) a7.  $W^*$

$$\sum_{i=0}^n (\bar{y}_i^T W - \bar{y}_i^T)(W^T \bar{x}_i - \bar{y}_i) + \lambda \sum_{j=1}^{k-1} \sum_{i=1}^n w_{ij} = \text{tr}(W^T W)$$

$$= \sum_{i=0}^n [\bar{x}_i^T W W^T \bar{x}_i + \bar{x}_i^T W \bar{y}_i - \bar{y}_i^T W^T \bar{x}_i + \bar{y}_i^T \bar{y}_i] + \text{tr}(W^T W)$$

Take derivative (over  $W$ ).  $\Rightarrow$

$$= \left\{ \sum_{i=0}^n \left[ \cancel{2W^T \bar{x}_i \bar{x}_i^T} - 2W^T \bar{x}_i \bar{x}_i^T + 2\lambda W^T \right] \right\}^T$$

$$= \left\{ 2W^T \sum_{i=1}^n \bar{x}_i \bar{x}_i^T - 2 \sum_{i=1}^n \bar{x}_i \bar{x}_i^T + 2\lambda W^T \right\}^T$$

$$\begin{aligned} \frac{\partial}{\partial W} \text{tr}(W^T W) &= \frac{\partial}{\partial W} \text{tr}(W^T \bar{x}_i \bar{x}_i^T) \\ \sum_{i=1}^n \bar{x}_i \bar{x}_i^T &= \sum_{i=1}^n \bar{x}_i \bar{x}_i^T + \sum_{i=1}^n \bar{x}_i \bar{x}_i^T \end{aligned}$$

$$= X X^T \quad X = [x_1, x_2, x_3, \dots]$$

$$\sum Y_i X_i^T = Y X \quad \text{defn.} \\ Y = [y_1, y_2, y_3, \dots, y_n]$$

$$K \times n \quad \rightarrow T$$

therefore gradient =  $\left\{ 2W^T X X^T - 2W^T X^T + 2\lambda W^T \right\}^T$

Want gradient  $\geq 0$  then  $2W^T X X^T + 2\lambda W^T = 2Y X^T$

$$W^T (2X X^T + 2\lambda I) = 2Y X^T$$

if  $\lambda > 0$   $(2X X^T + 2\lambda I)$  positive

thus invertible.

$$(2X X^T + 2\lambda I) \cancel{W} = 2Y X^T$$

$$W = (2X X^T + 2\lambda I)^{-1} \cdot 2Y X^T \\ = (X X^T + \lambda I)^{-1} \cdot Y X^T$$

b)  $\rightarrow$  python code

c) to compute  $(XX^T + \lambda I)^{-1} \cdot X^T$

$X: d \times n \rightarrow n \times c$

$XX^T$  takes  $d^2n$

inverse takes  $d^3$

$X^T$  takes  $d \times k$

therefore  ~~$O(d^3n + d^3 + dk^2)$~~

and  $(XX^T + \lambda I)^{-1} = X^T$

takes  $d^2k$

therefore  $O(d^2k + dk + d^3 + d^2n)$

e)  $W = [w_1, w_2, \dots, w_k]$

$X: d \times n$

original classifier

$$\sum_{i=0}^n \|Wx_i - y_i\|^2 + \lambda \|W\|_F^2$$

$$W^T x_i = \begin{vmatrix} W_1 \cdot x_i \\ W_2 \cdot x_i \\ \vdots \\ W_k \cdot x_i \end{vmatrix} \quad \|W\|_F^2 = \|W_1\|^2 + \|W_2\|^2 + \|W_3\|^2 + \dots + \|W_k\|^2$$

$$X = [x_1, x_2, \dots, x_n]$$

so we can choose  $w_i$  orthogonal to all  $x_i$

$$\text{then } X^T w_i = 0.$$

and we can take  $W$  apart  $W_L$  and  $W_R$

$W_L$  is consisted by  $w_i$   $X^T w_i = 0$

$W_R$  is consisted by  $w_i$  column space of  $X$ .

$$\text{then } W = W_L + X A \quad A = n \times c \text{ matrix}$$

We then plug this in the equation.

$$W = WL + XA$$

the first term

$$\sum_{i=0}^n \|((WL^T + AT^T)X_i - y_i)\|^2 = \sum_{i=0}^n \|AT^T X_i - y_i\|^2$$

$$= \sum_{i=1}^n |x_i^T X A A^T X_i - 2 x_i^T X A y_i + y_i^T y_i|$$

$$= \text{trace}(X^T X A A^T X) - 2 \text{trace}[X^T X A T] + \text{trace}(Y^T Y)$$

not about  $F$

$$\text{tail: } \|W\|_F^2 = \|WL + XA\|_F^2 = \text{trace}(WL + XA)^T (WL + XA)$$

$$\text{trace}(AT^T X A)$$

$$(WL^T + AT^T)(WL + XA) = WL^T WL + AT^T WL \\ + WL^T XA + AT^T X A$$

the first term does not depend on  $WL$   
and for the second term, when  $WL = 0$  has minimum value

We take derivative

$$\frac{\partial}{\partial A} \left( \text{trace}(X^T X A A^T X) - 2 \text{trace}[X^T X A T] + \text{trace}(AT^T X A) \right) = 0$$

$$\text{then } \cancel{\lambda} \cdot 2 A^T \underline{X^T X} + 2 A^T X^T \cancel{X^T X} = 2 \cancel{\lambda} X^T X \quad \text{Assume } X^T X \text{ invertible}$$

$$\text{then } \cancel{\lambda} \cdot 2 A^T + 2 A^T X^T X = 2 Y$$

$$A^T (\cancel{\lambda I} + X^T X) = Y \quad \text{Since } (\cancel{\lambda I} + X^T X) \text{ invertible}$$

$$A = (X^T X + \lambda I)^{-1} Y^T$$

Complexity:  $X^T X$  takes  $n^2 d$ . More  $n^3$  times  $Y^T$   $n^2 k$ .  
 $O(n^2 d + n^3 + n^2 k)$  if  $n \ll d$  compute  $A$  instead of  $W$ .

# 2.

a)  $X = AZ + b$      $Z = \begin{pmatrix} N(0, 1) \\ N(0, 1) \\ N(0, 1) \\ \vdots \end{pmatrix}$

$$\begin{aligned} S_{ij} &= E[(X_i - b_i)(X_j - b_j)] \text{ since } b = \bar{X} \\ &= E\left[\left(\sum_{k=1}^n A_{ik} Z_k\right)\left(\sum_{p=1}^n A_{jp} Z_p\right)\right] \\ &= \sum_{k=1}^n A_{ik} \bar{Z}_k \sum_{p=1}^n A_{jp} \bar{Z}_p \underbrace{E(Z_k Z_p)}_{\hookrightarrow = V_{kj}} \\ &\quad \hookrightarrow = V_{kj} = \begin{cases} 0 & k \neq j \\ 1 & k=j \end{cases} \end{aligned}$$

$$\bar{Z}_x = AA^T \text{ invertible} \quad = (AA^T)_{ij}$$

iff  $AA^T$  invertible

$\det(AA^T) = \det(A)\det(A^T) \Leftrightarrow A$  invertible

$\bar{Z}_x^{-1}$  not exist  $\Leftrightarrow A$  not invertible

$X = AZ + b \Rightarrow x_1, x_2, x_3, \dots, x_n$  Linearly dependant.

We can get  $A$  from  $X \sqrt{\bar{Z}_x} = A$ .  $\bar{Z}_x$  semipositive therefore  $\sqrt{\bar{Z}_x}$  exist  
then we do row reduce of  $A$  get  $A' = \begin{array}{c|ccc} & \cdots & \cdots & \cdots \\ \text{n rows all zero} & | & 0 & 0 \\ & | & 0 & 0 \\ & | & 0 & 0 \end{array}$

then we act the last  $n$  term in  $X$

convert  $x$  to  $x'$ . Since the last  $n$  are linearly dependent term  
we do not lose information.

b)  $S$  symmetric and positive defined.

$\bar{Z}^{-1}$  also symmetric and positive defined.

$$\bar{Z}^{-1} = U \Lambda U^T \quad \Lambda \text{ positive} \quad \Lambda = \sqrt{\Lambda} \cdot \sqrt{\Lambda}$$

$$X^T \sqrt{\Lambda} \cdot \sqrt{\Lambda} X = \| \sqrt{\Lambda} X \|_2^2$$

$$A = \sqrt{\Lambda} U^T \quad \Lambda = \text{diagonal matrix}$$

after diagonalizing  $\bar{Z}^{-1}$

c). ~~How~~ ~~if~~ ~~it's~~ ~~even~~ sample from

Convert  $X$  back to ~~independently~~  $Z = \begin{bmatrix} N(0,1) \\ N(0,1) \\ N(0,1) \end{bmatrix}$

$\bar{Z} = VV^T$        $X = VZ + b \rightarrow b = 0$

$\Sigma^{-1} = V^T V^{-1}$        $= VZ$

$X^T \Sigma^{-1} X = X^T V^T V^{-1} X$        $f = VZ + b$ .  $x, z$  are samples from  $X, Z$

$= Z^T V^T V^{-1} V^T VZ = Z^T Z = \|Ax\|^2$

# change Sample space.

d).  $\|Ax\|^2 = X^T \Sigma^{-1} X = X^T U \Lambda U^T X$ .  $U^T$  change base

therefore

$$\|U^T X\|_2 = 1$$

Max = Max<sup>num</sup> entry of  $\Lambda$

Min = Min<sup>entry</sup> Minimum entry of  $\Lambda$

entries of  $\Lambda$  are eigenvalues of  $\Sigma^{-1}$

if  $X_i \perp\!\!\! \perp X_j$  then  $\Sigma$  diagonal matrix entries correspond to variance of  $X_i$

$$\text{Maximum of } \|Ax\|^2 = \frac{1}{a} \quad a = \text{Minimum variance among } X_i$$

$$\text{Minimum of } \|Ax\|^2 = \frac{1}{b} \quad b = \text{Maximum variance among } X_i$$

choose  $X^*$  makes  $\|Ax\|^2 = \frac{1}{b}$   $b = \text{Max Variance}$

#4. a7. code

and use form solution from HW1.

$$w = (\lambda I + X^T X)^{-1} X^T y$$

b7.  $W_{t+1} = W_t - \alpha (2X X^T W_t - 2X T + 2\lambda W_t)$ .

c7.  $W_{t+1} = W_t - \alpha (2(X_i X_i^T W_t - 2x_i y_i^T + 2\lambda W_t))$ .

In order to help the iterative algorithm converge.  $\alpha$  should decrease generally can use exp or  $(1 - \frac{i}{\# \text{iteration}})$  --

d7. gradient descent:  $\alpha = \text{default value}$  reg = 0.1 (default).

Stochastic gradient descent:  $\alpha = \text{default}$  reg = 0.1  
and  $\alpha_i = \alpha * (1 - \frac{i}{\# \text{iteration}})$

From the Graph we can see that the error rate of gradient descent is always decreasing and the curve is very smooth. Because we compute the real derivative and with a proper  $\alpha$  the update direction is always correct.

The graph of stochastic gradient descent is a little different. Although it tends to decrease, we can see a lot of noise during the process because we randomly choose vector to compute the "expected" value of derivative, some times the update direction is not correct.

e7 My score is 0.94460

#5.

a).  $A(A^T A + M I) = (A A^T + M I) A \Leftarrow A A^T A + M I A = A A^T A + A M I$ .

and  $\begin{matrix} \swarrow \text{invertible} \\ A \end{matrix} \Leftarrow \begin{matrix} \text{Both positive} \\ A A^T + M I \end{matrix}$

$$(A A^T + M I)^{-1} A (A^T A + M I) (A^T A + M I)^{-1} = (A A^T + M I)^{-1} (A A^T + M I) A (A^T A + M I)^{-1}$$

$\Downarrow$

$$(A A^T + M I)^{-1} A = A (A^T A + M I)^{-1} \quad \square$$

b). From Hw1 we have  $W_* = W = (\lambda I + X^T X)^{-1} X^T y$

from part a we write  $as = X^T (X X^T + M I)^{-1} y$   
 $= X^T K \quad K = (X X^T + \lambda I)^{-1} y$   
 $= \sum_{i=1}^n a_i x_i$

$$a_i = [(X X^T + \lambda I)^{-1} y]_i$$

~~and which can also indicate that  $a_i$  unique  
since  $(X X^T + \lambda I)^{-1}$  is unique~~

~~$W_*$  is unique since  $W_* = (\lambda I + X^T X)^{-1} X^T y$~~

and  $(\lambda I + X^T X)$  invertible then  $W_*$  uniquely defined.

c).  $W = W_* + W_L \quad W_* = \sum a_i x_i \quad \text{and } W_L \in \text{Null}(X)$

$$\frac{1}{n} \sum_{i=1}^n \underbrace{\|W^T x_i - y_i\|^2}_{B} + M \|W\|^2 \leftarrow B \quad X^T W_L = 0 \quad x_i^T W_L = 0.$$

$A \rightarrow$  takes parameter  $W^T x_i$  and  $y_i$

$$W^T x_i = (W_*^T + W_L^T) x_i = W_*^T x_i$$

therefore  $W^T x_i$ 's value not depends on  $W_L$   
the value of part A not depends on  $W_L$

$$\min_{W \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \text{Loss}(W^T x_i, y_i) + \frac{\mu}{2} \|W\|_2^2$$

Loss function  
 convex. for  $k_1 \leq k_2$   
 $\lambda \text{Loss}(k_1, y_i) + (1-\lambda) \text{Loss}(k_2, y_i) \geq \text{Loss}(\lambda k_1 + (1-\lambda) k_2, y_i).$

$$\text{For Part B } = \frac{1}{2} \|W\|_2^2 = \frac{1}{2} (W^T W)^T (W^T + W) \\ = \frac{1}{2} (W^T W^T + W^T W)$$

if fix  $W^*$ , has minimum value when  $W^T = 0$

therefore to ensure  $\frac{1}{n} \sum_{i=1}^n \text{Loss}(W^T x_i, y_i) + \frac{\mu}{2} \|W\|_2^2$  has

minimum value,  $W^T$  should be zero

$$W^* = \sum_{i=1}^n x_i \alpha_i$$

We never use the condition "Loss function is convex"

therefore the proof above is a general proof.

Loss not convex, optimal solution still has the form

$$W^* = \sum_{i=1}^n \alpha_i x_i$$

#6. P (see those data)

$$TCP(y_i|x_i)$$

$$\begin{aligned} &= \mathcal{N}(w_0 + w_1 x_i, \sigma^2) (y_i) \\ &= \pi \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - (w_0 + w_1 x_i))^2}{2\sigma^2}} \\ &= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{-\frac{[(y_1 - (w_0 + w_1 x_1))^2 + (y_2 - (w_0 + w_1 x_2))^2 + \dots]}{2\sigma^2}} \end{aligned}$$

Want to Maximize  $TCP(y_i|x_i)$

then Maximize  $\sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2 = f(w_0, w_1)$

$$f_{w_0} = \sum_{i=1}^n 2(y_i - (w_0 + w_1 x_i)) = 0$$

$$f_{w_1} = \sum_{i=1}^n 2(y_i - (w_0 + w_1 x_i))(-x_i) = 0$$

$$\begin{aligned} \frac{\partial f_{w_0}}{\partial w_0} &= \sum_{i=1}^n y_i - \cancel{w_0} + \cancel{w_1} \sum_{i=1}^n x_i = 0 \\ &\quad + 2n w_0 \\ w_0 &= \bar{y} - \bar{w}_1 \bar{x} \end{aligned}$$

plug in  $w_0$ .

$$\begin{aligned} \cancel{\sum_{i=1}^n y_i - \cancel{w_0} + \cancel{w_1} \sum_{i=1}^n x_i = 0} \\ \cancel{+} \\ \cancel{\bar{y} - \bar{w}_1 \bar{x}} \end{aligned}$$

$$\cancel{f_{w_1} = 0}$$

$$= \sum_{i=1}^n y_i x_i + \cancel{w_0 \bar{x} t_i} + \cancel{w_1 \bar{x}^2} = 0.$$

$$\cancel{- \bar{x} \sum_{i=1}^n y_i + \bar{y} \bar{x}^2} + w_1 \bar{x} \bar{x} + w_1 \bar{x}^2 = 0$$

$$w_1 = \frac{\sum_{i=1}^n x_i y_i - \bar{y} \bar{x} n}{\bar{x}^2 - \bar{x} \bar{x}} = \frac{\sum_{i=1}^n x_i y_i - \bar{y} \bar{x} n}{\bar{x} \bar{t}_i^2 - n \bar{x}^2}$$

# 7. a).  $(0, 1) \quad (0, -1) \quad (-1, 0) \quad (1, 0)$  all  $\frac{1}{4}$

$$\text{Cov}(X, Y) = E(XY) = 0 \quad \text{uncorrelated}$$

But  $P(X=0, Y=0) = 0 \neq P(X=0) \cdot P(Y=0)$  not independent.

b). Consider  $X, Y$ .

$$P(X=1) = \frac{1}{2} \quad P(Y=1) = \frac{1}{2} \quad P(X=1, Y=1) = \frac{1}{4} = P(X=1) \cdot P(Y=1)$$

for all other cases  $P(X, Y) = P(X)P(Y)$

for  $X \neq Y$  pair due to symmetry. Symmetric property they are all ~~pairwise~~ pairwise independent

But they are not mutually independent

$$P(X=1, Y=1, Z=1) = 0 \neq \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2}$$