**a)** $H$ is symmetric (i.e., $H^T = H$).

$$H^T = \left(X\left(X^TX\right)^{-1}X^T\right)^T$$
$$= X\left(\left(X^TX\right)^{-1}\right)^T X^T \qquad \left[\because (ABC)^T = C^TB^TA^T\right]$$
$$= X\left(\left(X^TX\right)^T\right)^{-1} X^T \qquad \left[\because (X^{-1})^T = (X^T)^{-1}\right]$$
$$= X\left(X^TX\right)^{-1} X^T$$
$$= H.$$

**b)** $H^{k} = H$

$$H^2 = X(X^TX)^{-1}X^T \; X(X^TX)^{-1}X^T$$
$$= X(X^TX)^{-1}\underbrace{(X^TX)(X^TX)^{-1}}_{I}X^T$$
$$= X(X^TX)^{-1}$$

for $k=1$, $H^1 = H$

Suppose, $H^k = H$

Then $H^{k+1} = H^k. H$
$$= H \cdot H = X(X^TX)^{-1}X^T \; X(X^TX)^{-1}X^T$$
$$= X(X^TX)^{-1}\underbrace{(X^TX)(X^TX)^{-1}}_{I}X^T$$
$$= X(X^TX)^{-1}X^T = H.$$

$\left.\begin{array}{l}\end{array}\right\}$ Principle of mathematical induction

**c)** $(I-H)^{k} = (I-H)$.

when $k=1 \Rightarrow (I-H)^1 = I-H$

Suppose, $(I-H)^k = (I-H)$.

Then $(I-H)^{k+1} = (I-H)^k (I-H)$
$$= (I-H)(I-H)$$
$$= I - 2H + H^2$$
$$= I - 2H + H \qquad (\text{from } 1b \because H^2 = H).$$
$$= I - H.$$

$\left.\begin{array}{l}\end{array}\right\}$ principle of mathematical induction.

**d)** trace $(H) = d+1$.

trace $\left(\underbrace{X}_{A}\underbrace{(X^TX)^{-1}X^T}_{B}\right)$ = trace $\left((X^TX)^{-1}X^TX\right)$ $\qquad (\because tr(AB) = tr(BA))$
$$= \text{trace } (I_{(d+1)\times(d+1)}) = d+1.$$

2a) $w_{lin} = (x^T x)^{-1} x^T \vec{y}$

$\quad = (x^T x)^{-1} x^T (Xw^* + \vec{e})$

$\hat{y} = x w_{lin} = X(x^T x)^{-1} x^T (X w^* + \vec{e})$

$\quad = Xw^* + H\vec{e}$

Alternatively, ②

$H = X(x^T x)^{-1} x^T$

$\hat{y} = H\vec{y}$

$\quad = X(x^T x)^{-1} x^T (Xw^* + \vec{e})$

$\quad = Xw^* + H\vec{e}$.

Note: bold y $= \vec{y}$, $X_{N \times (d+1)}$

bold $\epsilon = \vec{e}$. $\quad w_{(d+1) \times 1}$

b) $\hat{y} - \vec{y} = Xw^* + H\vec{e} - (Xw^* + \vec{e})$

$\quad = H\vec{e} - \vec{e} = (H - I)\vec{e}$.

c) $E_{in}(w_{lin}) = \frac{1}{N} \| \hat{y} - \vec{y} \|^2$

$\quad = \frac{1}{N} (\hat{y} - \vec{y})^t (\hat{y} - \vec{y})$

$\quad = \frac{1}{N} [(H - I)\vec{e}]^t [(H - I)\vec{e}]$

$\quad = \frac{1}{N} \vec{e}^t (H - I)^t (H - I)\vec{e}$

$\quad = \frac{1}{N} \vec{e}^t (I - H)^2 \vec{e} \qquad \text{(from 1c)}$

$\quad = \frac{1}{N} \vec{e}^t (I - H) \vec{e}. \qquad \text{(from 1c)}$

d) $E_D [E_{in}(w_{lin})] = E_D \frac{1}{N} [\vec{e}^t (I - H)\vec{e}]$

$\quad = \frac{1}{N} E_D [\vec{e}^t \vec{e} - \vec{e}^t H\vec{e}]$

$\quad = \frac{1}{N} E_D [\sum_{i=1}^{N} e_i^2 - \sum_{i=1}^{N}$

$E_D [\vec{e}^t \epsilon] = E_D [\sum_{i=1}^{N} e_i^2] = \sum E_D [e_i^2]$

$\quad = N\sigma^2$

$\begin{cases} \because E[e_i^2] = var(e_i) \\ \qquad + (E(e_i))^2 \\ E[e_i^2] = \sigma^2 + 0 = \sigma^2 \end{cases}$

$$E_D\left[\vec{\epsilon}^{t} H \vec{\epsilon}\right] = E_D\left[\sum \epsilon_i H_{ij} \epsilon_j\right]$$

$$= \sum_{i,j=1}^{N} E_D\left[\epsilon_i H_{ij} \epsilon_j\right] \qquad (\because \text{expectation \& finite sum commute})$$

$$= \sum_{i,j=1}^{N} H_{ij} \, E_D\left[\epsilon_i \epsilon_j\right] \qquad \left(H_{ij} \text{ does not depend on } \epsilon_i, \epsilon_j\right)$$

$$= \sum_{\substack{i,j=1 \\ i=j}}^{N} h_{ii} E_D\left[\epsilon_i^2\right] + \underbrace{\sum_{i \neq j} h_{ij} \, E_D\left[\epsilon_i \epsilon_j\right]}_{= E_D[\epsilon_i] \, E_D[\epsilon_j] \; \left(\text{as } \epsilon_i \& \epsilon_j \text{ are independent}\right)}$$

$$= 0.$$

$$= \sum_{i=1}^{N} h_{ii} \sigma^2 = (d+1)\sigma^2 \qquad (\because tr(H) = d+1)$$

combining, we get

$$\frac{1}{N}\left(N\sigma^2 - (d+1)\sigma^2\right) = \sigma^2\left(1 - \frac{d+1}{N}\right).$$

e) $E_{test}(w_{lin}) = \frac{1}{N}\|\hat{\vec{y}} - \vec{y}\|^2$

$$= \frac{1}{N}\|X w^* + H\vec{\epsilon} - (X w^* - \vec{\epsilon'})\|$$

$$= \frac{1}{N}\|H\vec{\epsilon} - \vec{\epsilon'}\|$$

$$= \frac{1}{N}(H\vec{\epsilon} - \vec{\epsilon'})^{t}(H\vec{\epsilon} - \vec{\epsilon'})$$

$$\cancel{= \frac{1}{N}(\vec{\epsilon}\vec{\epsilon}^{t} - H)}$$

$$= \frac{1}{N}\left(\vec{\epsilon}^{t} H^{t} - \vec{\epsilon'}^{t}\right)(H\vec{\epsilon} - \vec{\epsilon'}).$$

$$= \frac{1}{N}\vec{\epsilon}^{t} H^{t} H \vec{\epsilon} - \vec{\epsilon}^{t} H^{t}\vec{\epsilon'} - \vec{\epsilon'}^{t} H\vec{\epsilon} + \vec{\epsilon'}^{t}\vec{\epsilon'}.$$

$$= \frac{1}{N}\left(\vec{\epsilon}^{t} H \vec{\epsilon} - 2\vec{\epsilon'}^{t} H\vec{\epsilon} + \vec{\epsilon'}^{t}\vec{\epsilon'}\right). \quad \left(\because H^2 = H \atop H^{t} = H\right)$$

Similar to (2d),

$$E\left(\vec{\epsilon'}^{t}\vec{\epsilon'}\right) = N\sigma^2$$

$$E\left(\vec{\epsilon} H \vec{\epsilon}\right) = (d+1)\sigma^2.$$

Now, $E_{D,\vec{\epsilon}'}(\vec{\epsilon}'^{\,t} H \vec{\epsilon})$.

$$= E_D\left[E_{\vec{\epsilon}'|D}(\vec{\epsilon}'^{\,t} H \vec{\epsilon})\right]$$

$$= E_D\left[E_{\vec{\epsilon}'|D}(\vec{\epsilon}'^{\,t})\cdot H\vec{\epsilon}\right] \qquad \left[\begin{array}{l}\because H\vec{\epsilon} \text{ is independent}\\ \qquad \text{of } \vec{\epsilon}'\end{array}\right]$$

$$= E_D\left[0\cdot H\vec{\epsilon}\right]. \qquad \left(as \; E[\epsilon_i'] = 0\right)$$

$$= 0 \qquad \left(\Rightarrow E_{\vec{\epsilon}'|D}\circ(\vec{\epsilon}'^{\,t}) = \vec{0}\right)$$

combining,

$$N\sigma^2 + (d+1)\sigma^2 = \frac{1}{N}\left(N\sigma^2 + (d+1)\sigma^2\right)$$

$$= \sigma^2\left(1 + \frac{d+1}{N}\right).$$

## Question 3

a) $\epsilon \to$ noise realization of test point

$\vec{\epsilon} \to$ bold epsilon i.e., noise realization of all train data points

$$y = w^{*T}x + \epsilon$$
$$= x^T w^* + \epsilon$$

$$w_{lin} = (X^TX)^{-1}X^T y_{data}$$
$$= (X^TX)^{-1}X^T(Xw^* + \vec{\epsilon})$$
$$= w^* + (X^TX)^{-1}X^T\vec{\epsilon}.$$

$$\Rightarrow g(x) = x^t\left(w^* + (X^TX)^{-1}X^T\vec{\epsilon}\right).$$
$$= x^t w^* + x^t(X^TX)^{-1}X^T\vec{\epsilon}$$

$$\Rightarrow y - g(x) = \epsilon - x^t(X^TX)^{-1}X^T\vec{\epsilon}.$$

b) Take joint expectation w.r.t $x \& \epsilon$.

$$E_{out} = E_{x, \epsilon} \left[ (y - g(x) \right]^2$$

$$= E_{x, \epsilon} \left[ (y - g(x))^t (y - g(x)) \right]$$

$$= E_{x, \epsilon} \left[ (\epsilon - x^t (x^T x)^{-1} x^T \vec{\epsilon})^t (\epsilon - x^t (x^T x)^{-1} x^T \vec{\epsilon}) \right]$$

$$= E_{x, \epsilon} \left[ (\epsilon - \vec{\epsilon}^t x (x^T x)^{-1} x) (\epsilon - x^t (x^T x)^{-1} x^T \vec{\epsilon}) \right]$$

$$\left( \text{as } \epsilon \text{ is scalar} \atop \epsilon^t = \epsilon \right)$$

$$= E_{x, \epsilon} \left[ \epsilon^2 - \epsilon \, x^t (x^T x)^{-1} x^T \vec{\epsilon} - \vec{\epsilon}^t x (x^T x)^{-1} x \epsilon + \right.$$
$$\left. \vec{\epsilon}^t x (x^T x)^{-1} x x^t (x^T x)^{-1} x^T \vec{\epsilon} \right]$$

$$= E_{x, \epsilon} [\epsilon^2] - E_{x, \epsilon} [2 \epsilon \, x^t (x^T x)^{-1} x^T \vec{\epsilon}] + E_{x, \epsilon} [\vec{\epsilon}^t x (x^T x)^{-1} x x^t (x^T x)^{-1} x^T \vec{\epsilon}]$$

$$\left( \text{as 2nd \& 3rd terms are scalars \&} \atop \text{for scalar } a, \, a = a^t \right)$$

---

**1st term**  $= E_{x, \epsilon}$

$$E_{x, \epsilon} [\epsilon^2] = \cancel{\phi \phi \phi \phi} E_x \, E_{\epsilon | x} [\epsilon^2] = E_x [\sigma^2] = \sigma^2$$

**2nd term**

$$E_{x, \epsilon} [2 \epsilon \, x^t (x^T x)^{-1} x^T \vec{\epsilon}] = 2 \, E_x \left[ E_{\epsilon | x} \, \epsilon \, \underbrace{x^t (x^T x)^{-1} x^T \vec{\epsilon}}_{\text{independent of} \atop \epsilon \, \& \, x} \right]$$

$$= 2 \, E_x \left[ \underbrace{E_{\epsilon | x} [\epsilon]}_{0}^{\;0} \; \cancel{\otimes} \; \epsilon \, x^t (x^T x)^{-1} x^T \vec{\epsilon} \right]$$

$$= 0$$

**3rd term :-**

$$E_{x, \epsilon} [\vec{\epsilon}^t x (x^T x)^{-1} x x^t (x^T x)^{-1} x^T \vec{\epsilon}]$$

$$= E_{x, \epsilon} \left[ tr \left( \underbrace{\vec{\epsilon}^t x (x^T x)^{-1}}_{A} \, \underbrace{x x^t (x^T x)^{-1} x^T \vec{\epsilon}}_{B} \right) \right] \quad \left( \text{as the term} \atop \text{inside trace} \atop \text{is scalar} \right)$$

$$= E_{\eta, \epsilon} \left\{ tr \left[ \eta \eta^t (X^TX)^{-1} X^T \vec{\epsilon} \vec{\epsilon}^t X (X^TX)^{-1} \right] \right\} \qquad \boxed{tr(AB) = tr(BA)}$$

$$= tr \; E_{\eta, \epsilon} \left[ \qquad " \qquad \right] \qquad \left[ \begin{array}{l} \text{as expectation } \& \\ \text{trace commute} \end{array} \right]$$

$$= tr \left\{ E_{\eta} \left[ E_{\epsilon | \eta} (\eta \eta^t) \right] \underbrace{(X^TX)^{-1} X^T \vec{\epsilon} \vec{\epsilon}^t X (X^TX)^{-1}}_{\text{independent of both } \epsilon \& \eta.} \right\}$$

$$= tr \left[ E_{\eta} (\eta \eta^t) \left[ (X^TX)^{-1} X^T \vec{\epsilon} \vec{\epsilon}^t X (X^TX)^{-1} \right] \right]$$

$$= tr \left[ \Sigma (X^TX)^{-1} X^T \vec{\epsilon} \vec{\epsilon}^t X (X^TX)^{-1} \right]$$

c) $E_{\vec{\epsilon}} \left[ \vec{\epsilon} \vec{\epsilon}^t \right] = E_{\vec{\epsilon}} \begin{bmatrix} \epsilon_1^2 & \epsilon_1 \epsilon_2 & \cdots & \epsilon_1 \epsilon_n \\ \epsilon_2 \epsilon_1 & \epsilon_2^2 & \cdots & \epsilon_2 \epsilon_n \\ \vdots & & & \\ \epsilon_n^2 & \epsilon_n \epsilon_1 & \cdots & \epsilon_n \epsilon_n \end{bmatrix}$

$$E_{\vec{\epsilon}} \left[ \epsilon_i^2 \right] = \sigma^2$$

$$E_{\vec{\epsilon}} \left[ \epsilon_i \epsilon_j \right] = E_{\epsilon_i} \left[ \epsilon_i \right] E_{\epsilon_j} \left[ \epsilon_j \right] = 0,$$

$$\text{so,} \quad \sigma^2 I_{N \times N}.$$

as $\epsilon$ is of dimension $N \times 1$.

$$E_{\vec{\epsilon}} \left[ \vec{\epsilon} \vec{\epsilon}^t \right] = \sigma^2 I_{N \times N}.$$

as $\vec{e}$ is of dimension $N \times 1$, ~~the resulting matrix~~

② we have $E_{\vec{e}} \left[ \vec{e} \; \vec{e}^{\,t} \right] = \sigma^2 I_{N \times N}$

d) From ©,

$$E_{out} = \sigma^2 + trace \left[ \Sigma \, (X^T X)^{-1} X^T \vec{e} \vec{e}^{\,T} X \, (X^T X)^{-1} \right]$$

$$\Rightarrow E_{\vec{e}} \left[ E_{out} \right] = \sigma^2 + E_{\vec{e}} \; tr \left[ \qquad \right]$$

$$= \sigma^2 + tr \; E_{\vec{e}} \left[ \qquad \right]$$

$$= \sigma^2 + tr \left( \underbrace{\Sigma \, (X^T X)^{-1} X^T}_{① P} \; E_{\vec{e}} \left[ \vec{e} \vec{e}^{\,T} \right] \; \underbrace{X \, (X^T X)^{-1}}_{②} \right).$$

P & Q are independent of $\vec{e}$.

$$\cancel{= \sigma^2 + tr \left( \Sigma \, (X^T X)^{-1} X^T \, X \, (X^T X)^{-1} \right)}$$

$$= \sigma^2 + tr \left[ \Sigma \, (X^T X)^{-1} X^T \sigma^2 I \; X (X^T X)^{-1} \right]$$

$$= \sigma^2 + \sigma^2 \; tr \left[ \Sigma \, (X^T X)^{-1} (X^T X) \, (X^T X)^{-1} \right]$$

$$= \sigma^2 + \sigma^2 \; tr \left( \Sigma \, (X^T X)^{-1} \right)$$

$$= \sigma^2 + \frac{\sigma^2}{N} \; tr \left( \Sigma \left( \frac{X^T X}{N} \right)^{-1} \right).$$

If $\dfrac{X^T X}{N} \approx \Sigma$ then

$$= \sigma^2 + \frac{\sigma^2}{N} \; tr \left( \Sigma \, \Sigma^{-1} \right)$$

$$= \sigma^2 + \frac{\sigma^2}{N} (d+1) = \sigma^2 \left( 1 + \frac{d+1}{N} \right).$$

e) $X^TX$ is an N-Sample estimate of $\Sigma$.

By the law of large number $\dfrac{X^TX}{N} \xrightarrow{p} \Sigma$.

using continuity at $\Sigma$, $\left(\dfrac{X^TX}{N}\right)^{-1} \xrightarrow{p} \Sigma^{-1}$

$$\Rightarrow \quad \Sigma \left(\dfrac{X^TX}{N}\right)^{-1} \xrightarrow{p} \Sigma\Sigma^{-1} = I$$

$$\Rightarrow trace\left(\Sigma \left(\dfrac{X^TX}{N}\right)^{-1}\right) \xrightarrow{p} d+1 \quad \left(\therefore \text{ as trace is continuous at } I\right)$$

$$\Rightarrow trace\left(\Sigma \dfrac{(X^TX)^{-1}}{N}\right) = (d+1) + \epsilon. \quad \text{for som small scalar } \epsilon.$$

$$\Rightarrow \sigma^2 + \dfrac{\sigma^2}{N} trace\left(\Sigma \left(\dfrac{1}{N}X^TX\right)^{-1}\right) = \sigma^2\left(1 + \dfrac{d+1}{N} + \dfrac{\epsilon}{N}\right)$$

$$\Rightarrow E_{out} = \sigma^2\left(1 + \dfrac{d+1}{N} + o\left(\dfrac{1}{N}\right)\right).$$

# COL341 Spring 2023
## Homework 1
## Grading guidelines

## 1 Question1

- Everyone has got good marks in this.

- Question 1a, mention the result $((AB)^{-1})^T = ((AB)^T)^{-1}$ i.e., transpose and inverse commute.

- One can formalize the proof in 1b and 1c by using mathematical induction properly.

## 2 Question2

- You should make the distinction between $\epsilon$ (noise associated with the test input) and $\boldsymbol{\epsilon}$ (noise realizations of the training examples)

- In Q2d, mention the independence of $\epsilon_i$ and $\epsilon_j$, $\mathbb{E}[\epsilon] = 0$ and $\mathbb{E}[\epsilon^2] = var(\epsilon) - [\mathbb{E}[\epsilon]]^2 = \sigma^2$, where you are using these results.

- Mention the independence of $\epsilon$ and $\epsilon'$.

- Mention which random variable are you taking the expectation over, and break the join distribution appropriately, In 2e, $\mathbb{E}_{\mathcal{D},\epsilon'}$ should be broken into $\mathbb{E}_{\mathcal{D}}[\mathbb{E}_{\epsilon'|\mathcal{D}}]$.

## 3 Question3

- 

- You should make the distinction between $\epsilon$ (noise associated with the test input) and $\boldsymbol{\epsilon}$ (noise realizations of the training examples)

- In 3c, mention teh dimension of the matrix.

- Mention in which step did you use, $a = trace(a)$ for any scalar and $trace(AB) = trace(BA)$ and what is your $A$ and $B$.

- Mention trace and expectation commute.

- Some have mentioned that test noise $\epsilon$ is independent of test point $x$, which is not true.

- In 3b, $\mathbb{E}_{\boldsymbol{x},\epsilon'}$ should be broken into $\mathbb{E}_{\boldsymbol{x}}[\mathbb{E}_{\epsilon'|\boldsymbol{x}}]$.

- Some have not done the last bit of question 3d and 0.25 have been deducted for that.

- Additionally, mention how can you take certain terms out of the expectation.

- In 3e, mention the continuity of trace at the identity matrix $\mathbb{I}$.