

COL 341 Homework1

Chinmay Mittal

TOTAL POINTS

19.75 / 20

QUESTION 1

1 Question 1,2,3 **19.75 / 20**

✓ + **10 pts** Correct

✓ + **5 pts** Click here to replace this description.

✓ + **2 pts** Click here to replace this description.

✓ + **2 pts** Click here to replace this description.

+ **1 pts** Click here to replace this description.

✓ + **0.5 pts** Click here to replace this description.

✓ + **0.25 pts** Click here to replace this description.

Chinmay Mittal
(2020CS10336)

Homework 1 | COL341
Machine Learning

February 3, 2023

Question 1

1 mark

a)

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

To show H is symmetric we compute H^T

$$H^T = (X(X^T X)^{-1} X^T)^T = ((X^T)^T ((X^T X)^{-1})^T X^T)^T \text{ using } (AB)^T = B^T A^T,$$

$$\Rightarrow H^T = (X((X^T X)^T)^{-1} X^T), \text{ using } (A^{-1})^T = (A^T)^{-1}$$

$$\Rightarrow H^T = (X(X^T (X^T)^T)^{-1} X^T) \text{ using } (AB)^T = B^T A^T$$

$$\Rightarrow H^T = (X(X^T X)^{-1} X^T), \text{ using } (A^T)^T = A$$

$$\Rightarrow H^T = H$$

$\Rightarrow H$ is symmetric

b)

To show that $H^K = H$ for any positive integer K , we use induction

Base case is trivially true, for $K=1$, $H^K = H^1 = H$

The induction hypothesis is that $H^{K-1} = H$ and we need to show that $H^K = H$,

1 mark

We have $H^K = H^{K-1} H = H H = H^2$ from the induction hypothesis

$$\Rightarrow H^K = (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T)$$

$$\Rightarrow H^K = (X(X^T X)^{-1} ((X^T X)(X^T X)^{-1}) X^T) \text{ from associativity of matrix multiplication}$$

$$\Rightarrow H^K = (X(X^T X)^{-1} I) X^T = H \text{ which proves our result by induction}$$

c)

To show that $(I - H)^K = (I - H)$ we use induction and as before the base case is trivially true.

The induction hypothesis is that $(I - H)^{K-1} = (I - H)$ and we need to show $(I - H)^K = (I - H)$

$$\Rightarrow (I - H)^K = (I - H)^{K-1} (I - H) = (I - H)^2$$

$$\Rightarrow (I - H)^K = (I - H)^2 = (I - H)(I - H) = I^2 - 2H + H^2$$

$$\Rightarrow (I - H)^K = I - 2H + H^2 = I - 2H + H = I - H \text{ following the above 1b)}$$

1 mark

This proves our induction hypothesis

d)

We need to show that $\text{trace}(H) = d + 1$

Consider $A = X(X^T X)^{-1}$ and $B = X^T$, thus we have $H = AB$

1 mark

Now $\text{trace}(H) = \text{trace}(AB) = \text{trace}(BA) = \text{trace}(X^T X (X^T X)^{-1})$

$$\Rightarrow \text{trace}(H) = \text{trace}(I_{d+1}) = d + 1 \text{ since } X^T X \text{ is } d+1 \text{ dimensional which proves the required result}$$

Question 2

1 mark

a)

The in-sample estimate of y is given by $\hat{y} = X w_{lin}$ where w_{lin} are the learnt linear weights.

$w_{lin} = (X^T X)^{-1} X^T y \Rightarrow \hat{y} = (X(X^T X)^{-1} X^T y) = H y$, follows from the proof in class of learned linear weights

$$\Rightarrow \hat{y} = H y = H(X w^* + \epsilon) = (H X) w^* + H \epsilon$$

Also we have, $H X = X(X^T X)^{-1} X^T X = X$

$$\Rightarrow \hat{y} = X w^* + H \epsilon \text{ as required}$$

b)

$$\hat{y} - y = (X w^* + H \epsilon) - (X w^* + H \epsilon) = (H - I_N) \epsilon$$

1 mark

Thus the matrix is $H - I_N$

c)

The insample error is given by $E_{in} = \frac{1}{N}(\hat{\mathbf{y}} - \mathbf{y})^T(\hat{\mathbf{y}} - \mathbf{y})$

$$\begin{aligned} &\Rightarrow E_{in} = \frac{1}{N}((H - I)\boldsymbol{\epsilon})^T((H - I)\boldsymbol{\epsilon}) \quad \text{1mark} \\ &\Rightarrow E_{in} = \frac{1}{N}\boldsymbol{\epsilon}^T(H - I)^T(H - I)\boldsymbol{\epsilon} \\ &\Rightarrow E_{in} = \frac{1}{N}\boldsymbol{\epsilon}^T(H^T - I^T)(H - I)\boldsymbol{\epsilon} \\ &\Rightarrow E_{in} = \frac{1}{N}\boldsymbol{\epsilon}^T(H - I)^2\boldsymbol{\epsilon} \text{ since H is symmetric} \\ &\Rightarrow E_{in} = \frac{1}{N}\boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon} \text{ since } (H - I)^2 = (I - H)^2 = I - H \text{ from 1c)} \end{aligned}$$

d)

$$\begin{aligned} &\text{From 1c) We have that } E_{in} = \frac{1}{N}\boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon} \\ &\Rightarrow \mathbb{E}_{\mathcal{D}}[E_{in}(\mathbf{w}_{lin})] = \mathbb{E}_{\mathcal{D}}[\frac{1}{N}\boldsymbol{\epsilon}^T(I - H)\boldsymbol{\epsilon}] \quad \text{2mark} \\ &\Rightarrow \mathbb{E}_{\mathcal{D}}[E_{in}(\mathbf{w}_{lin})] = \frac{1}{N}(\mathbb{E}_{\mathcal{D}}[\boldsymbol{\epsilon}^T\boldsymbol{\epsilon}] - \mathbb{E}_{\mathcal{D}}[\boldsymbol{\epsilon}^TH\boldsymbol{\epsilon}]) \\ &= \frac{1}{N}(\mathbb{E}_{\mathcal{D}}[\sum_{i=1}^N(\epsilon_i^2)] - \mathbb{E}_{\mathcal{D}}[\boldsymbol{\epsilon}^TH\boldsymbol{\epsilon}]) = \frac{1}{N}(\sum_{i=1}^N(\mathbb{E}_{\mathcal{D}}[\epsilon_i^2]) - \mathbb{E}_{\mathcal{D}}[\boldsymbol{\epsilon}^TH\boldsymbol{\epsilon}]) \\ &= \frac{1}{N}(N\sigma^2 - \mathbb{E}_{\mathcal{D}}[\boldsymbol{\epsilon}^TH\boldsymbol{\epsilon}]) \text{ since } \epsilon_i \text{ is normally distributed with variance } \sigma^2 \\ &= \sigma^2 - \frac{1}{N}\mathbb{E}_{\mathcal{D}}(\sum_{i=1}^N \sum_{j=1}^N (c_{ij}\epsilon_i\epsilon_j)) \text{ from the definition of matrix multiplication} \\ &= \sigma^2 - \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^N (\mathbb{E}_{\mathcal{D}}(c_{ij}\epsilon_i\epsilon_j)) \\ &= \sigma^2 - \frac{1}{N} \sum_{i=1}^N (\mathbb{E}_{\mathcal{D}}(c_{ii}\epsilon_i\epsilon_i)), \text{ because } \mathbb{E}_{\mathcal{D}}[\epsilon_i\epsilon_j] = \mathbb{E}_{\mathcal{D}}[\epsilon_i]\mathbb{E}_{\mathcal{D}}[\epsilon_j] = 0 \text{ since } \epsilon_i \text{ and } \epsilon_j \text{ are independent.} \\ &\text{Also Note here that } c_{ii} = H_{ii} \text{ from the definition of matrix multiplication} \\ &= \sigma^2 - \frac{1}{N}\mathbb{E}_{\mathcal{D}}[\epsilon_i^2] \sum_{i=1}^N H_{ii} \\ &\Rightarrow \sigma^2 - \frac{\sigma^2}{N}\text{trace}(H) \\ &= \sigma^2(1 - \frac{d+1}{N}) \text{ from 1d)} \end{aligned}$$

e)

$$\begin{aligned} &\mathbf{y}_{test} - \hat{\mathbf{y}}_{test} = (X\mathbf{w}^* + \boldsymbol{\epsilon}') - (X\mathbf{w}^* + H\boldsymbol{\epsilon}) = \boldsymbol{\epsilon}' - H\boldsymbol{\epsilon} \quad \text{1.75mark} \\ &E_{test}(\mathbf{w}_{lin}) = \frac{1}{N}(\mathbf{y}_{test} - \hat{\mathbf{y}}_{test})^T(\mathbf{y}_{test} - \hat{\mathbf{y}}_{test}) = \frac{1}{N}(H\boldsymbol{\epsilon} - \boldsymbol{\epsilon}')^T(H\boldsymbol{\epsilon} - \boldsymbol{\epsilon}') \\ &= \frac{1}{N}((H\boldsymbol{\epsilon})^T - \boldsymbol{\epsilon}'^T)(H\boldsymbol{\epsilon} - \boldsymbol{\epsilon}') = \frac{1}{N}(\boldsymbol{\epsilon}^TH^T - \boldsymbol{\epsilon}'^T)(H\boldsymbol{\epsilon} - \boldsymbol{\epsilon}') \\ &= \frac{1}{N}(\boldsymbol{\epsilon}^TH^TH\boldsymbol{\epsilon} + \boldsymbol{\epsilon}'^2 - \boldsymbol{\epsilon}'^TH\boldsymbol{\epsilon} - \boldsymbol{\epsilon}^TH^T\boldsymbol{\epsilon}') = \frac{1}{N}(\boldsymbol{\epsilon}^THH\boldsymbol{\epsilon} + \boldsymbol{\epsilon}'^2 - \boldsymbol{\epsilon}'^TH\boldsymbol{\epsilon} - \boldsymbol{\epsilon}^TH^T\boldsymbol{\epsilon}') = \frac{1}{N}(\boldsymbol{\epsilon}^TH\boldsymbol{\epsilon} + \boldsymbol{\epsilon}'^2 - \boldsymbol{\epsilon}'^TH\boldsymbol{\epsilon} - \boldsymbol{\epsilon}^TH^T\boldsymbol{\epsilon}') \\ &\Rightarrow \mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[E_{test}(\mathbf{w}_{lin})] = \frac{1}{N}(\sigma^2(d+1) + N\sigma^2) - \mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[\boldsymbol{\epsilon}'^TH\boldsymbol{\epsilon} + \boldsymbol{\epsilon}^TH^T\boldsymbol{\epsilon}'] \text{ which follows from 2d)} \\ &= \frac{1}{N}(\sigma^2(d+1) + N\sigma^2) + \mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[\sum_{i=1}^N \sum_{j=1}^N (c_{ij})\epsilon_i\epsilon'_j] \text{ for some constant } c_{ij} \text{ derived from H following definition of} \\ &\text{matrix multiplication} \\ &= \frac{1}{N}(\sigma^2(d+1) + N\sigma^2) + \sum_{i=1}^N \sum_{j=1}^N c_{ij}\mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[\epsilon_i]\mathbb{E}_{\mathcal{D}, \boldsymbol{\epsilon}'}[\epsilon'_j] = \frac{1}{N}(\sigma^2(d+1) + N\sigma^2) + 0 \\ &= \sigma^2(1 + \frac{d+1}{N}) \end{aligned}$$

Question 3

a)

$$\begin{aligned} &\text{We have a test point with } \mathbf{y} = \mathbf{x}^T\mathbf{w}^* + \epsilon \\ &\text{The learn weights are } \mathbf{w}_{lin} = (X^TX)^{-1}X^T\mathbf{y}_{train} \text{ where } \mathbf{y}_{train} = X\mathbf{w}^* + \boldsymbol{\epsilon} \\ &\text{Thus we get } g(\mathbf{x}) = \mathbf{x}^T\mathbf{w}_{lin} = \mathbf{x}^T(X^TX)^{-1}X^T\mathbf{y}_{train} = \mathbf{x}^T(X^TX)^{-1}X^T(X\mathbf{w}^* + \boldsymbol{\epsilon}) \\ &\Rightarrow g(\mathbf{x}) = \mathbf{x}^T(X^TX)^{-1}X^TX\mathbf{w}^* + \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon} = \mathbf{x}^T(I)\mathbf{w}^* + \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon} = \mathbf{x}^T\mathbf{w}^* + \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon} \\ &\Rightarrow \mathbf{y} - g(\mathbf{x}) = (\mathbf{x}^T\mathbf{w}^* + \epsilon) - (\mathbf{x}^T\mathbf{w}^* + \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}) \\ &\Rightarrow \mathbf{y} - g(\mathbf{x}) = \epsilon - \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon} \quad \text{1mark} \end{aligned}$$

b)

$$\begin{aligned} &\text{From the previous question we can conclude } E_{out} = \mathbb{E}_{\mathbf{x}, \epsilon}[(\mathbf{y} - g(\mathbf{x}))^2] \\ &E_{out} = \mathbb{E}_{\mathbf{x}, \epsilon}[(\epsilon - \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon})^2] = \mathbb{E}_{\mathbf{x}, \epsilon}[\epsilon^2 - 2\epsilon(\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}) + \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}(\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon})^T] \\ &= \mathbb{E}_{\mathbf{x}, \epsilon}[\epsilon^2] - 2\mathbb{E}_{\mathbf{x}, \epsilon}[\epsilon]\mathbb{E}_{\mathbf{x}, \epsilon}[\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}] + \mathbb{E}_{\mathbf{x}, \epsilon}[\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}\boldsymbol{\epsilon}^TX(X^TX)^{-1}\mathbf{x}] \\ &= \sigma^2 - 0(\mathbb{E}_{\mathbf{x}, \epsilon}[\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}]) + \mathbb{E}_{\mathbf{x}, \epsilon}[\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}\boldsymbol{\epsilon}^TX(X^TX)^{-1}\mathbf{x}] = \sigma^2 + \mathbb{E}_{\mathbf{x}, \epsilon}[\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}\boldsymbol{\epsilon}^TX(X^TX)^{-1}\mathbf{x}] \\ &= \sigma^2 + \mathbb{E}_{\mathbf{x}, \epsilon}[\text{trace}(\mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}\boldsymbol{\epsilon}^TX(X^TX)^{-1}\mathbf{x})], \text{ since the second term is a scalar taking the trace doesn't} \\ &\text{make a difference} \\ &\text{Now let } A = \mathbf{x}^T(X^TX)^{-1}X^T\boldsymbol{\epsilon}\boldsymbol{\epsilon}^TX(X^TX)^{-1} \text{ and } B = \mathbf{x} \text{ such that } = \sigma^2 + \mathbb{E}_{\mathbf{x}, \epsilon}[\text{trace}(AB)] = \sigma^2 + \mathbb{E}_{\mathbf{x}, \epsilon}[\text{trace}(BA)] \end{aligned}$$

2marks

$$\begin{aligned}
&= \sigma^2 + \mathbb{E}_{x,\epsilon}[\text{trace}(xx^T(X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})] = \text{trace}(\mathbb{E}_{x,\epsilon}[xx^T(X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1}]) \text{ since trace} \\
&\text{and expectation commute} \\
&= \sigma^2 + \text{trace}(\mathbb{E}_{x,\epsilon}[xx^T]((X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})) \\
&= \sigma^2 + \text{trace}(\sum (X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})
\end{aligned}$$

c)

$\epsilon \epsilon^T$ is a matrix with the entries $[\epsilon_i \epsilon_j]$
 $\mathbb{E}[\epsilon_i \epsilon_i] = \sigma^2$ and $\mathbb{E}[\epsilon_i \epsilon_j] = \mathbb{E}[\epsilon_i] \mathbb{E}[\epsilon_j] = 0$ when $i \neq j$ 2marks
Hence we get that $\mathbb{E}_\epsilon[\epsilon \epsilon^T] = \sigma^2 \mathbb{I}_N$

d)

From 3b) We have that $E_{out} = \sigma^2 + \text{trace}(\sum (X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})$
 $\Rightarrow \mathbb{E}_\epsilon[E_{test}] = \mathbb{E}_\epsilon[\sigma^2 + \text{trace}(\sum (X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})] = \sigma^2 + \mathbb{E}_\epsilon[\text{trace}(\sum (X^T X)^{-1}X^T \epsilon \epsilon^T X(X^T X)^{-1})]$
 $= \sigma^2 + \text{trace}(\sum (X^T X)^{-1}X^T \mathbb{E}_\epsilon[\epsilon \epsilon^T] X(X^T X)^{-1})$ since other parts are independent of the expectation
 $= \sigma^2 + \text{trace}(\sum (X^T X)^{-1}X^T (\sigma^2 \mathbb{I}) X(X^T X)^{-1})$ from 3c) 2marks
 $= \sigma^2 + \sigma^2 \text{trace}(\sum (X^T X)^{-1}X^T X(X^T X)^{-1})$
 $= \sigma^2 + \sigma^2 \text{trace}(\sum ((X^T X)^{-1}(X^T X))(X^T X)^{-1})$
 $= \sigma^2 + \sigma^2 \text{trace}(\sum (X^T X)^{-1}) = \sigma^2 + \sigma^2 (\sum \frac{1}{N} (\frac{1}{N} (X^T X))^{-1})$
 $\Rightarrow E_{out} = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\sum (\frac{1}{N} (X^T X)^{-1}))$
Also we have that $\frac{1}{N} (X^T X)$ is the N sample estimate of \sum thus if we approximate $\sum \approx \frac{1}{N} (X^T X)^{-1}$ we get
 $E_{out} \approx \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\sum \sum^{-1}) = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\mathbb{I}_{d+1}) \approx \sigma^2 + \frac{\sigma^2}{N} (d+1)$ on average

e)

Consider the Random Variable (Matrix) defined as $Y = \sum \frac{1}{N} (X^T X)^{-1}$, here the randomness is over the values of X that is the features.

The random variable Y has a true mean (which is a matrix) this true mean is the same as taking $N \rightarrow \infty$ from which we can see $\mathbb{E}[Y] = \mathbb{I}_{d+1}$. 2marks

Hoeffding's inequality can be applied to each entry in the matrix which is itself a random variable from this, we can say that with a high probability $\sum \frac{1}{N} (X^T X)^{-1} < \mathbb{I}_{d+1} + c$ where c is a constant matrix and the inequality is applied element-wise (this is just an alternate form of the Law of Large Numbers).

Thus using the Law of Large Numbers in this form of Hoeffding's inequality we get that with high probability:

$$\begin{aligned}
&E_{out} = \sigma^2 + \frac{\sigma^2}{N} \text{trace}(\mathbb{I}_{d+1} + c) \text{ where } c \text{ is a constant matrix} \\
&\Rightarrow E_{out} = \sigma^2 + \frac{\sigma^2}{N} (d+1 + k) \text{ with high probability where } k \text{ is constant such that } k = \text{trace}(c) \\
&\Rightarrow E_{out} = \sigma^2 (1 + \frac{d+1}{N} + \frac{k}{N}) = \sigma^2 (1 + \frac{d+1}{N} + o(\frac{1}{N})) \text{ from the definition of the little-}(o) \text{ notation which proves our result.}
\end{aligned}$$