

# MATH 390.4 / 650.2 Spring 2020 Homework #3

Pizon Shetu

Due noon Friday, March 13, 2020 under the door of KY604

(this document last updated 1:21am on Thursday 19<sup>th</sup> March, 2020)

## Instructions and Philosophy

The path to success in this class is to do many problems. Unlike other courses, exclusively doing reading(s) will not help. Coming to lecture is akin to watching workout videos; thinking about and solving problems on your own is the actual “working out.” Feel free to “work out” with others; **I want you to work on this in groups.**

Reading is still *required*. For this homework set, read Chapters 3-6 of Silver’s book. You should be googling and reading about all the concepts introduced in class online. This is your responsibility to supplement in-class with *your own* readings.

The problems below are color coded: **green** problems are considered *easy* and marked “[easy]”; **yellow** problems are considered *intermediate* and marked “[harder]”, **red** problems are considered *difficult* and marked “[difficult]” and **purple** problems are extra credit. The *easy* problems are intended to be “giveaways” if you went to class. Do as much as you can of the others; I expect you to at least attempt the *difficult* problems.

This homework is worth 100 points but the point distribution will not be determined until after the due date. See syllabus for the policy on late homework.

Up to 7 points are given as a bonus if the homework is typed using L<sup>A</sup>T<sub>E</sub>X. Links to installing L<sup>A</sup>T<sub>E</sub>X and program for compiling L<sup>A</sup>T<sub>E</sub>X is found on the syllabus. You are encouraged to use [overleaf.com](https://www.overleaf.com). If you are handing in homework this way, read the comments in the code; there are two lines to comment out and you should replace my name with yours and write your section. The easiest way to use overleaf is to copy the raw text from hwxx.tex and preamble.tex into two new overleaf tex files with the same name. If you are asked to make drawings, you can take a picture of your handwritten drawing and insert them as figures or leave space using the “\vspace” command and draw them in after printing or attach them stapled.

The document is available with spaces for you to write your answers. If not using L<sup>A</sup>T<sub>E</sub>X, print this document *including this first page* and write in your answers. **I do not accept homeworks which are *not* on this printout.**

NAME: PIZON SHETU

## Problem 1

These are questions about Silver's book, chapters 3-6. For all parts in this question, answer using notation from class (i.e.  $t, f, g, h^*, \delta, \epsilon, e, t, z_1, \dots, z_t, \mathbb{D}, \mathcal{H}, \mathcal{A}, \mathcal{X}, \mathcal{Y}, X, y, n, p, x_1, \dots, x_p, x_1, \dots, x_n$ , etc.

- (a) [easy] What algorithm that we studied in class is PECOTA most similar to?

K Nearest Neighbors

- (b) [easy] Is baseball performance as a function of age a linear model? Discuss.

No, it is difficult to rate players on age alone as the data is noisy, but for most players performances increases from early 20's to 30 and decline as they age.

- (c) [harder] How can baseball scouts do better than a prediction system like PECOTA?

They would have know more features, by that I mean more relevant independent features which might allow them to tune out the noise and make the data fit better.

- (d) [harder] Why hasn't anyone (at the time of the writing of Silver's book) taken advantage of Pitch f/x data to predict future success?

Because that type of data was extremely expensive and at that time, not a lot of people had access to such tools to monitor such data

- (e) [difficult] Chapter 4 is all about predicting weather. Broadly speaking, what is the problem with weather predictions? Make sure you use the framework and notation from class. This is not an easy question and we will discuss in class. Do your best.

Being able to predict the weather with high accuracy and precision is nearly impossible, weather takes into multiple factors in reality, and we would need to be able to predict accurately for each of those factors making weather predictions a multi-level prediction where each predictions is predicated on the supporting factors to be correct for its predictions to be valid. In theoretical sense  $f$  is difficult to understand, we would need to able to properly calculate motion, moisture, and etc but this is difficult so we simplify and make approximations.

- (f) [easy] Why does the weatherman lie about the chance of rain? And where should you go if you want honest forecasts?

Its easier to be wrong on bad news than to be right on good news. If there is a possibility of rain the weatherman will lie about the chances to make sure he has some knowledge of the occurrences but if he says there is no chance of rain even though it will be a sunny day he hedges and comes out good because people don't get upset about wrong news if its good news.

- (g) [difficult] Chapter 5 is all about predicting earthquakes. Broadly speaking, what is the problem with earthquake predictions? It is *not* the same as the problem of predicting weather. Read page 162 a few times. Make sure you use the framework and notation from class.

It's difficult to measure actual data since most of it is under ground. Experts really only have access to past results. It is hard to predict having mostly results from the past and not being able to access current  $x$ 's

- (h) [easy] Silver has quite a whimsical explanation of overfitting on page 163 but it is really educational! What is the nonsense predictor in the model he describes?

Having the combination to 3 locks, also knowing the flaws of these locks but using this to lock pick other locks. This is over-fitting as combinations to 3 specific locks have no correlation to other locks nor does knowing the flaws of 3 specific locks.

- (i) [easy] John von Neumann was credited with saying that "with four parameters I can fit an elephant and with five I can make him wiggle his trunk". What did he mean by that and what is the message to you, the budding data scientist?

Essentially a good model can be over-fitted merely by adding 1 unnecessary feature

- (j) [difficult] Chapter 6 is all about predicting unemployment, an index of macroeconomic performance of a country. Broadly speaking, what is the problem with unemployment predictions? It is *not* the same as the problem of predicting weather or earthquakes. Make sure you use the framework and notation from class.

Unemployment predictions are often too late, their predictions become relevant only few months prior to the phenomena. Economic forecasters face 3 fundamental challenges according to Silver, "First, it is very hard to determine cause and effect from economic statistics alone. Second, the economy is always changing, so explanations of economic behavior that hold in one business cycle may not apply to future ones. And third, as bad as their forecasts have been, the data that economists have to work with isn't much good either." In other words  $f$  is extremely hard to find, and not  $x$  fit with the ever changing economy so  $f$  is also always changing thus making the prediction difficult to manufacture.

- (k) [E.C.] Many times in this chapter Silver says something on the order of "you need to have theories about how things function in order to make good predictions." Do you agree? Discuss.

## Problem 2

These are questions about multivariate linear model fitting using the least squares algorithm.

- (a) [difficult] Derive  $\frac{\partial}{\partial \mathbf{c}} [\mathbf{c}^\top A \mathbf{c}]$  where  $\mathbf{c} \in \mathbb{R}^n$  and  $A \in \mathbb{R}^{n \times n}$  but *not* symmetric. Get as far as you can.

$$\begin{aligned}
\mathbf{c}^\top \mathbf{A} \mathbf{c} &= \begin{bmatrix} c_1 & c_2 & \dots & c_n \end{bmatrix} \cdot \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \cdot \begin{bmatrix} c_1 \\ c_2 \\ \dots \\ c_n \end{bmatrix} \\
&= \begin{bmatrix} c_1 & c_2 & \dots & c_n \end{bmatrix} \cdot \begin{bmatrix} c_1 a_{11} + c_2 a_{12} + \dots + c_n a_{1n} \\ c_1 a_{21} + c_2 a_{22} + \dots + c_n a_{2n} \\ \dots \\ c_1 a_{n1} + c_2 a_{n2} + \dots + c_n a_{nn} \end{bmatrix} \\
&= c_1 (c_1 a_{11} + c_2 a_{12} + \dots + c_n a_{1n}) + \dots + c_n (c_1 a_{n1} + c_2 a_{n2} + \dots + c_n a_{nn}) \\
&\text{where} \\
&= \sum_{i=1}^n c_i \left( \sum_{j=1}^n c_j a_{ij} \right) \quad \text{and} \quad \frac{\partial}{\partial \mathbf{c}_i} [\mathbf{c}^\top \mathbf{A} \mathbf{c}] = \frac{\partial}{\partial \mathbf{c}_i} \left[ \sum_{i,j=1}^n c_i c_j a_{ij} \right] \\
&\text{where} \\
&= \sum_{j=1}^n c_j a_{ij} + c_j a_{ji}
\end{aligned}$$

- (b) [easy] Given matrix  $X \in \mathbb{R}^{n \times (p+1)}$ , full rank and first column consisting of the  $\mathbf{1}_n$  vector, rederive the least squares solution  $\mathbf{b}$  (the vector of coefficients in the linear model shipped in the prediction function  $g$ ). No need to rederive the facts about vector derivatives.

$$\begin{aligned}
SSE &= \sum (\vec{y}_i - \hat{\vec{y}})^2 = (\vec{y} - \hat{\vec{y}})^T (\vec{y} - \hat{\vec{y}}) = (\vec{y}^T - \hat{\vec{y}}^T) (\vec{y} - \hat{\vec{y}}) \\
&= \vec{y}^T \vec{y} - \vec{y}^T \hat{\vec{y}} - \hat{\vec{y}}^T \vec{y} + \hat{\vec{y}}^T \hat{\vec{y}} \\
&= \vec{y}^T \vec{y} - 2\hat{\vec{y}}^T \vec{y} + \hat{\vec{y}}^T \hat{\vec{y}} \\
&= \vec{y}^T \vec{y} - 2(X\vec{b})^T \vec{y} + (X\vec{b})^T (X\vec{b}) \\
&= \vec{y}^T \vec{y} - 2\vec{b}^T X^T \vec{y} + \vec{b}^T X^T X \vec{b}
\end{aligned}$$

$$\frac{\partial}{\partial \vec{b}} [\vec{y}^T \vec{y} - 2\vec{b}^T X^T \vec{y} + \vec{b}^T X^T X \vec{b}] = -2X^T \vec{y} + 2X^T X \vec{b} = 0$$

$$(X^T X)^{-1} (X^T X) \vec{b} = (X^T X)^{-1} X^T \vec{y}$$

$$\vec{b} = (X^T X)^{-1} X^T \vec{y}$$

- (c) [harder] Consider the case where  $p = 1$ . Show that the solution for  $\mathbf{b}$  you just derived is the same solution that we proved for simple regression in Lecture 8. That is, the first element of  $\mathbf{b}$  is the same as  $b_0 = \bar{y} - r \frac{s_y}{s_x} \bar{x}$  and the second element of  $\mathbf{b}$  is  $b_1 = r \frac{s_y}{s_x}$ .

$$\vec{b} = (X^T X)^{-1} X^T \vec{y}$$

$$X \in \mathbb{R}^{n \times 2} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \cdot \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

$$\sum x_i = n\bar{x}$$

$$(X^T X)^{-1} = \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{bmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{bmatrix}$$

$$(X^T X)^{-1} (X^T \vec{y}) = \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{bmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{bmatrix} \cdot \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix}$$

$$= \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{bmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{bmatrix} \cdot \begin{bmatrix} n\bar{y} \\ \sum y_i x_i \end{bmatrix}$$

$$= \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{bmatrix} n\bar{y} \sum x_i^2 - n \sum y_i x_i \\ -n^2 \bar{x} \bar{y} + n \sum y_i x_i \end{bmatrix}$$

$$b_0 = \frac{\bar{y} (\sum x_i^2 - n\bar{x}^2) - \bar{x} (\sum y_i x_i - n\bar{x}\bar{y})}{\sum x_i^2 - n\bar{x}^2}$$

$$b_1 = \frac{\sum y_i x_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2}$$

(d) [easy] If  $X$  is rank deficient, how can you solve for  $\mathbf{b}$ ? Explain in English. Make  $X$  full-rank by making all columns linearly independent

(e) [difficult] Prove  $\text{rank}[X] = \text{rank}[X^T X]$ .

$$\text{rank}[x] = \text{Dim of whole space} - N(X)$$

Prove that  $N(A) = N(A^T A)$ :

$$N(A) \subset N(A^T A)$$

$$x \in N(A)$$

$$Ax = 0$$

$$A^T Ax = A^T 0 = 0 \quad \Rightarrow \quad x \in N(A^T A)$$

$$\text{so } N(A) \subset N(A^T A)$$

$$N(A^T A) \subset N(A)$$

$$\begin{aligned}
x &\in N(A^T A) \\
A^T A x &= 0 \\
x^T A^T A x &= x^T 0 = 0 \\
(Ax)^T (Ax) &= 0 \\
\|Ax\|^2 &= 0 \\
Ax &= 0 \\
\text{so } N(A^T A) &\subset N(A)
\end{aligned}$$

Therefore,  $N(X) = N(X^T X)$  and  $\text{rank}[X] = \text{rank}[X^T X]$

- (f) [difficult] Given matrix  $X \in \mathbb{R}^{n \times (p+1)}$ , full rank and first column consisting of the  $\mathbf{1}_n$  vector, now consider cost multiples (“weights”)  $c_1, c_2, \dots, c_n$  for each mistake  $e_i$ . As an example, previously the mistake for the 17th observation was  $e_{17} := y_{17} - \hat{y}_{17}$  but now it would be  $e_{17} := c_{17}(y_{17} - \hat{y}_{17})$ . Derive the weighted least squares solution  $\mathbf{b}$ . No need to rederive the facts about vector derivatives. Hints: (1) show that SSE is a quadratic form with the matrix  $C$  in the middle (2) Split this matrix up into two pieces i.e.  $C = C^{\frac{1}{2}} C^{\frac{1}{2}}$ , distribute and then foil (3) note that a scalar value equals its own transpose and (4) use the vector derivative formulas.

$$\begin{aligned}
SSE &= (\mathbf{y} - X\vec{b})^T C (\mathbf{y} - X\vec{b}) \\
&= (\mathbf{y}^T C \mathbf{y} - \mathbf{y}^T C X \vec{b} - \vec{b}^T X^T C \mathbf{y} + \vec{b}^T X^T C X \vec{b})
\end{aligned}$$

$$\frac{\partial SSE}{\partial \vec{b}} = -2X^T C \mathbf{y} + 2X^T C X \vec{b} = 0$$

$$\vec{b} = (X^T C X)^{-1} X^T C \mathbf{y}$$

- (g) [difficult] If  $p = 1$ , prove  $r^2 = R^2$  i.e. the linear correlation is the same as proportion of sample variance explained in a least squares linear model.

$$\begin{aligned}
b_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \\
b_0 &= \bar{y} - b_1 \bar{x}
\end{aligned}$$

$$\hat{\mathbf{y}}_i = b_0 + b_1 x_i = \bar{y} - \left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right] \bar{x} + \left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right] x_i$$

$$\begin{aligned}
R^2 &= \frac{SSR}{SST} \\
&= \frac{\sum (\hat{\mathbf{y}}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\
&= \frac{\sum \left( \bar{y} - \left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right] \bar{x} + \left[ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right] x_i - \bar{y} \right)^2}{\sum (y_i - \bar{y})^2}
\end{aligned}$$

$$\begin{aligned}
&= \frac{\sum \left( \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \right)^2 (x_i - \bar{x})^2}{\sum (y_i - \bar{y})^2} \\
&= \frac{[\sum (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2} \\
&= \frac{[Cov(x, y)]^2}{s_x^2 s_y^2} \\
&= r^2
\end{aligned}$$

(h) [harder] Prove that  $g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p]) = \bar{y}$  in OLS.

Recall that

$$\bar{y} = b_0 + b_1 \bar{x}_1 + b_2 \bar{x}_2 + \dots + b_{p+1} \bar{x}_{p+1}$$

and

$$\hat{\mathbf{y}}^* = g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p])$$

since

$$\begin{aligned}
\hat{\mathbf{y}}^* &= b_0 + b_1 \bar{x}_1 + b_2 \bar{x}_2 + \dots + b_{p+1} \bar{x}_{p+1} \\
\hat{\mathbf{y}}^* &= \frac{1}{n} \sum b_0 + \frac{1}{n} \sum b_1 x_{i1} + \frac{1}{n} \sum b_2 x_{i2} + \dots + \frac{1}{n} \sum b_{p+1} x_{ip+1}
\end{aligned}$$

and we know that

$$\hat{\mathbf{y}}^* = \frac{1}{n} \sum \hat{\mathbf{y}}_i \text{ since } \hat{\mathbf{y}}_i = b_0 + b_1 x_{i1} + b_2 x_{i2} + \dots + b_{p+1} x_{ip+1}$$

Also recall  $\frac{1}{n} \sum (\mathbf{y}_i - \hat{\mathbf{y}}_i) = \frac{1}{n} \sum e_i = 0$   
 So  $\frac{1}{n} \sum \hat{\mathbf{y}}_i = \frac{1}{n} \sum \mathbf{y}_i = \bar{\mathbf{y}}$

Hence,

$$\hat{\mathbf{y}}^* = g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p]) = \bar{\mathbf{y}}$$

(i) [harder] Prove that  $\bar{e} = 0$  in OLS.

(j) [difficult] If you model  $\mathbf{y}$  with one categorical nominal variable that has levels  $A, B, C$ , prove that the OLS estimates look like  $\bar{y}_A$  if  $x = A$ ,  $\bar{y}_B$  if  $x = B$  and  $\bar{y}_C$  if  $x = C$ . You can choose to use an intercept or not. Likely without is easier.

$$\vec{x} = \begin{bmatrix} A \\ B \\ A \\ C \\ \dots \end{bmatrix}, \mathbb{X} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ \cdot & , & , \end{bmatrix}, \vec{b} = (X^T X)^{-1} X^T \vec{y}$$

$$(X^T X) = \begin{bmatrix} 1 & 0 & 1 & 0 & \dots \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} n_A & 0 & 0 \\ 0 & n_B & 0 \\ 0 & 0 & n_C \end{bmatrix}$$

$$(X^T X)^{-1} = \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix}, (X^T \vec{y}) = \begin{bmatrix} 1 & 0 & 1 & 0 & \dots \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 0 & 0 & \dots \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_n \end{bmatrix} = \begin{bmatrix} \sum_{i=A} Y_i \\ \sum_{i=B} Y_i \\ \sum_{i=C} Y_i \end{bmatrix}$$

$$\vec{b} = (X^T X)^{-1} (X^T \vec{y}) = \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix} \begin{bmatrix} \sum_{i=A} Y_i \\ \sum_{i=B} Y_i \\ \sum_{i=C} Y_i \end{bmatrix} = \begin{bmatrix} \bar{y}_A \\ \bar{y}_B \\ \bar{y}_C \end{bmatrix}$$

### Problem 3

These are questions related to the concept of orthogonal projection, QR decomposition and its relationship with least squares linear modeling.

(a) [difficult] [MA] Prove that if a square matrix is both symmetric and idempotent then it must be an orthogonal projection matrix.

(b) [easy] Prove that  $I_n$  is an orthogonal projection matrix  $\forall n$ .

Assume  $I_n$  is an orthogonal projection matrix then  $I_n$  is symmetrical and idempotent

$$I_n = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \text{ then } I_n^T = I_n \text{ meaning } \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

And finally the idempotency  $I_n I_n = I_n$

$$\begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots \\ 0 & 0 & 1 & 0 & \dots \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

(c) [easy] What subspace does  $I_n$  project onto?

It projects onto the colspace of  $I_n$

(d) [easy] Consider least squares linear regression using a design matrix  $X$  with rank  $p+1$ . What are the degrees of freedom in the resulting model? What does this mean?



Since there are  $p+1$  degrees of freedom thus  $\hat{\mathbf{y}} = w_0 + w_1x_1 + \dots + w_px_p$  would have  $p+1$  weight parameters which can be adjusted

- (e) [harder] If you are orthogonally projecting the vector  $\mathbf{y}$  onto the column space of  $X$  which is of rank  $p+1$ , derive the formula for  $\text{Proj}_{\text{colsp}[X]}[\mathbf{y}]$ . Is this the same as in OLS?

$$\begin{aligned}\text{Proj}_{\text{colsp}[X]}[\mathbf{y}] &= X\vec{w} \\ X^T(\mathbf{y} - X\vec{w}) &= 0 \text{ Because of orthogonality} \\ X^T\mathbf{y} - X^TX\vec{w} &= 0 \\ X^T\mathbf{y} &= X^TX\vec{w} \\ \vec{w} &= (X^TX)^{-1}X^T\mathbf{y} \\ \text{Proj}_{\text{colsp}[X]}[\mathbf{y}] &= X\vec{w} = X(X^TX)^{-1}X^T\mathbf{y} = H\mathbf{y} \\ \text{Since this is the same as OLS, yes it is.}\end{aligned}$$

- (f) [difficult] We saw that the perceptron is an *iterative algorithm*. This means that it goes through multiple iterations in order to converge to a closer and closer  $\mathbf{w}$ . Why not do the same with linear least squares regression? Consider the following. Regress  $\mathbf{y}$  using  $\mathbf{X}$  to get  $\hat{\mathbf{y}}$ . This generates residuals  $\mathbf{e}$  (the leftover piece of  $\mathbf{y}$  that wasn't explained by the regression's fit,  $\hat{\mathbf{y}}$ ). Now try again! Regress  $\mathbf{e}$  using  $\mathbf{X}$  and then get new residuals  $\mathbf{e}_{\text{new}}$ . Would  $\mathbf{e}_{\text{new}}$  be closer to  $\mathbf{0}_n$  than the first  $\mathbf{e}$ ? That is, wouldn't this yield a better model on iteration #2? Yes/no and explain.

$H = X(X^TX)^{-1}X^T$  and  $H \cdot H = H$  Since the projection onto the  $X$  gives the least square error on the 1st iteration. Due to the idem-potency of the projection matrix, you would get the same thing over and over again.

- (g) [harder] Prove that  $\mathbf{Q}^\top = \mathbf{Q}^{-1}$  where  $\mathbf{Q}$  is an orthonormal matrix such that  $\text{colsp}[\mathbf{Q}] = \text{colsp}[\mathbf{X}]$  and  $\mathbf{Q}$  and  $\mathbf{X}$  are both matrices  $\in \mathbb{R}^{n \times (p+1)}$ . Hint: this is purely a linear algebra exercise.

$$\begin{bmatrix} \leftarrow & q_{\cdot 1} & \rightarrow \\ \leftarrow & q_{\cdot 2} & \rightarrow \\ & \dots & \\ \leftarrow & q_{\cdot n} & \rightarrow \end{bmatrix} \cdot \begin{bmatrix} \uparrow & \uparrow & & \uparrow \\ q_{\cdot 1} & q_{\cdot 2} & \dots & q_{\cdot n} \\ \downarrow & \downarrow & & \downarrow \end{bmatrix} = \begin{bmatrix} q_{\cdot 1}^\top q_{\cdot 1} & q_{\cdot 1}^\top q_{\cdot 2} & \dots & q_{\cdot 1}^\top q_{\cdot n} \\ q_{\cdot 2}^\top q_{\cdot 1} & q_{\cdot 2}^\top q_{\cdot 2} & \dots & q_{\cdot 2}^\top q_{\cdot n} \\ \dots & \dots & \dots & \dots \\ q_{\cdot n}^\top q_{\cdot 1} & q_{\cdot n}^\top q_{\cdot 2} & \dots & q_{\cdot n}^\top q_{\cdot n} \end{bmatrix}$$

$q_{\cdot i}^\top q_{\cdot i} = ||q_{\cdot i}||^2$ , and  $q_{\cdot i}^\top q_{\cdot j} = 0$  when  $i \neq j$  because of orthonormality of  $\mathbf{Q}$ .

$$\mathbf{Q}^\top \mathbf{Q} = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{bmatrix} = I$$

- (h) [harder] Prove that the least squares projection  $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{Q}^T$ .

Since  $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = (\mathbf{Q} \mathbf{R}) \left( (\mathbf{Q} \mathbf{R})^T \mathbf{Q} \mathbf{R} \right)^{-1} (\mathbf{Q} \mathbf{R})^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{R} (\mathbf{R}^T \mathbf{Q}^T \mathbf{Q} \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{R} (\mathbf{R}^T \mathbf{I} \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{R} (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \mathbf{Q}^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{R} \mathbf{R}^{-1} (\mathbf{R}^T)^{-1} \mathbf{R}^T \mathbf{Q}^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{I} \mathbf{I} \mathbf{Q}^T$$

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{Q} \mathbf{Q}^T.$$

- (i) [harder] Prove that an orthogonal projection onto the colsp  $[\mathbf{Q}]$  is the same as the sum of the projections onto each column of  $\mathbf{Q}$ .

If we do a projection on each column of  $\mathbf{Q}$ :  $\text{Proj}_{q_i} [\vec{a}] = \frac{q_i q_i^T}{\|q_i\|^2} \vec{a}$

Recall  $\mathbf{Q}$  is orthonormal,  $\|q_i\|^2 = 1$ , and  $\text{Proj}_{q_i} [\vec{a}] = q_i q_i^T \vec{a}$

So

$$\sum_{i=1}^{p+1} \text{Proj}_{q_i} [\vec{a}] = \sum_{i=1}^{p+1} q_i q_i^T \vec{a} = \mathbf{Q} \mathbf{Q}^T \vec{a}$$

Due to orthogonality of all the columns in  $\mathbf{Q}$

- (j) [easy] Prove that adding a new column to  $\mathbf{X}$  results in SST remaining the same.

- (k) [difficult] [MA] Prove that  $\text{rank}[\mathbf{H}] = \text{tr}[\mathbf{H}]$ . Hint: you will need to use facts about eigenvalues and the eigendecomposition of projection matrices that we learned in class.

## Problem 4

All of these are extra credit. This is for students who want to get a taste of a first year linear model theory class at the graduate level. The prereq to do these problems is Math 368/621. Only attempt these if you have time!

In linear modeling,  $\mathcal{H} = \{\mathbf{x}\mathbf{w} : \mathbf{w} \in \mathbb{R}^{p+1}\}$  where  $\mathbf{x} = [1 \ x_1 \ \dots \ x_p]$ , a row vector. Thus, there is a best function  $h^*(\mathbf{x}) = \mathbf{x}\boldsymbol{\beta}$  where  $\boldsymbol{\beta} = [\beta_0 \ \beta_1 \ \dots \ \beta_p]^T$ , a column vector and  $y = h^*(\mathbf{x}) = \mathbf{x}\boldsymbol{\beta} + \mathcal{E}$ . Imagine that for all  $n$  observations in  $\mathbb{D}$ , the  $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mathcal{E}}$  where  $\boldsymbol{\mathcal{E}} \sim \mathcal{N}_n(\mathbf{0}_n, \sigma^2 \mathbf{I}_n)$  and  $\mathbf{Y}$  is a random vector with dimension  $n$  modeling the responses of which  $\mathbf{y}$  is a random realization. Assume  $\sigma^2$  is known.

- (a) [E.C.] Show that  $\mathbf{Y} \sim \mathcal{N}_n(X\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$ .
- (b) [E.C.] Let  $\mathbf{B} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$ , i.e. the r.v. that represents the OLS estimator of which  $\mathbf{b}$  is one realization which changes based on the realizations of the error-vector r.v.  $\boldsymbol{\varepsilon}$ . Find the distribution of  $\mathbf{B}$  and once this is done, its expectation and variance-covariance matrix. Do the entries in  $\mathbf{B}$  have dependence?
- (c) [E.C.] Find the distribution of  $\hat{\mathbf{Y}}$ , the vector r.v. of predictions.
- (d) [E.C.] Find the distribution of  $\mathbf{E}$ , the vector r.v. of residuals.
- (e) [E.C.] Find the distribution of  $SST$ .
- (f) [E.C.] Find the distribution of  $SSE$ .
- (g) [E.C.] Find the distribution of  $SSR$ .
- (h) [E.C.] Find the distribution of  $R^2$ .
- (i) [E.C.] Now let  $\sigma^2$  be unknown. Use the MSE as its estimate. What is the distribution of  $\mathbf{B}$  now?
- (j) [E.C.] What is the distribution of MSE?
- (k) [E.C.] What is the distribution of  $R^2$ ?
- (l) [E.C.] Let  $\mathbf{U} \sim \mathcal{N}_n(\mathbf{0}_n, \mathbf{I}_n)$  independent of  $\mathbf{V} \sim \mathcal{N}_n(\mathbf{0}_n, \mathbf{I}_n)$ . Let  $\theta$  be the r.v. model of the angle between  $\mathbf{U}$  and  $\mathbf{V}$ . How is  $\theta$  distributed?