

MATH 390.4 Spring 2020 Homework #3

Tziporah Horowitz

Monday 16th March, 2020

1 Silver's Book, Chapters 3-6

For all parts in this question, answer using notation from class (i.e. $t, f, g, h^*, \delta, \epsilon, e, t, z_1, \dots, z_t, \mathbb{D}, \mathcal{H}, \mathcal{A}, \mathcal{X}, \mathcal{Y}, X, y, n, p, x_{\cdot 1}, \dots, x_{\cdot p}, x_{1\cdot}, \dots, x_{n\cdot}$, etc.)

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

PECOTA is most similar to k nearest neighbors.

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

Baseball performance as a function of age is not a linear model because it is parabolic, with different players peaking at different ages.

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

Scouts did better than PECOTA because they used a hybrid approach. They used both human judgement and statistics and therefore had more data available to them.

- (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?

At the time, they did have complete 3-dimensional recordings of everything on the baseball field.

- (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.

The problem with predicting weather is that it is too chaotic. The weather behaves like a dynamic system, so f cannot be linear. Slight inaccuracies in prior data (\mathbb{D}) cause much larger inaccuracies in the predictions. In addition, to predict accurately, you have to increase the number of dimensions; however, if you increase the number of dimensions, you're left with exponentially more equations to solve.

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

Weathermen lie about the chance of rain for economic incentives. It gives them room to for fewer complaints when it does rain in a low probability case. To get more "honest" forecasts, you should turn to The National Weather Service.

- (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.

We cannot truly measure X so it becomes too noisy and the models are left with too much residual error, e .

- (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?

Using the color of a lock to match it to known combinations. It is overly specific.

- (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?

Neumann meant that when you model noise, it creates illustrations of reality that are inaccurate. It's important to recognize the difference between necessary features and noise.

- (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.

Economic concepts are hard to predict because they rely on economic policy; once we begin to measure \mathbb{D} , its behavior begins to change. Similar to weather systems, economic data are chaotic. Nevertheless, meteorologists have a stronger understanding of the atmosphere and can therefore make more accurate forecasts, while economists cannot.

- (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.

In order to make good predictions, you need a good X . Since there is no way of knowing the true causal inputs, the definition of X is subjective to your knowledge about y . So, the more you know about how y functions, the better your judgement will be about choosing X .

2 Multivariate Linear Model Fitting Using LS

- (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.

$$\begin{aligned}
 \frac{\partial}{\partial \mathbf{c}} [\mathbf{c}^\top A \mathbf{c}] &= \frac{\partial}{\partial \mathbf{c}} \left[\mathbf{c}^\top \begin{bmatrix} \mathbf{a}_{1\cdot} \\ \vdots \\ \mathbf{a}_{n\cdot} \end{bmatrix} \mathbf{c} \right] \\
 &= \frac{\partial}{\partial \mathbf{c}} \left[\mathbf{c}^\top \begin{bmatrix} \mathbf{a}_{1\cdot} \mathbf{c} \\ \vdots \\ \mathbf{a}_{n\cdot} \mathbf{c} \end{bmatrix} \right] \\
 &= \frac{\partial}{\partial \mathbf{c}} [c_1 \mathbf{a}_{1\cdot} \mathbf{c} + \dots + c_n \mathbf{a}_{n\cdot} \mathbf{c}] \\
 &= \frac{\partial}{\partial \mathbf{c}} [c_1 (a_{11}c_1 + \dots + a_{1n}c_n) + \dots + c_n (a_{n1}c_1 + \dots + a_{nn}c_n)] \\
 &= \frac{\partial}{\partial \mathbf{c}} [(c_1^2 a_{11} + c_1 a_{12}c_2 + \dots + c_1 a_{1n}c_n) + \dots + (c_n a_{n1}c_1 + \dots + c_n^2 a_{nn})] \\
 &\Rightarrow \frac{\partial}{\partial c_1} = (a_{1\cdot} + a_{\cdot 1}) \mathbf{c} \\
 &= \begin{bmatrix} (a_{1\cdot} + a_{\cdot 1}) \mathbf{c} \\ \vdots \\ (a_{n\cdot} + a_{\cdot n}) \mathbf{c} \end{bmatrix} \\
 &= \begin{bmatrix} \mathbf{a}_{1\cdot} \\ \vdots \\ \mathbf{a}_{n\cdot} \end{bmatrix} \mathbf{c} + \begin{bmatrix} \mathbf{a}_{\cdot 1} \\ \vdots \\ \mathbf{a}_{\cdot n} \end{bmatrix} \mathbf{c}
 \end{aligned}$$

- (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.

$$\begin{aligned}
 \frac{\partial}{\partial \mathbf{w}} [\text{SSE}] &= \frac{\partial}{\partial \mathbf{w}} [\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top X^\top \mathbf{y} + \mathbf{w}^\top X^\top X \mathbf{w}] \\
 &= \mathbf{0}_3 - 2X^\top \mathbf{y} + 2X^\top X \mathbf{w} \stackrel{\text{set}}{=} 0
 \end{aligned}$$

$$\begin{aligned}
 X^\top X \mathbf{w} &= X^\top \mathbf{y} \\
 \mathbf{w} &= (X^\top X)^{-1} X^\top \mathbf{y}
 \end{aligned}$$

- (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 $b_1 = r_{\frac{s_y}{s_x}}$.

$$\begin{aligned}
\mathbf{b} &= (X^\top X)^{-1} X^\top \mathbf{y} \\
&= \left(\begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \\
&= \left(\begin{bmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \right)^{-1} \begin{bmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{bmatrix} \\
&= \begin{bmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} n\bar{y} \\ \sum_{i=1}^n x_i y_i \end{bmatrix} \\
&= \frac{1}{n \sum x_i^2 - n^2 \bar{x}^2} \begin{bmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{bmatrix} \begin{bmatrix} n\bar{y} \\ \sum x_i y_i \end{bmatrix} \\
&= \frac{1}{n(\sum x_i^2 - n\bar{x}^2)} \begin{bmatrix} \sum x_i^2 \sum y_i - n\bar{x} \sum x_i y_i \\ -n\bar{x} \sum y_i + n \sum x_i y_i \end{bmatrix} \\
&= \begin{bmatrix} \bar{y} - \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} \bar{x} \\ \frac{\sum x_i y_i - n\bar{x}\bar{y}}{\sum x_i^2 - n\bar{x}^2} \end{bmatrix} \\
&= \begin{bmatrix} \bar{y} - r_{\frac{s_y}{s_x}} \bar{x} \\ r_{\frac{s_y}{s_x}} \end{bmatrix}
\end{aligned}$$

- (d) [easy] If X is rank deficient, how can you solve for \mathbf{b} ? Explain in English.

If X is rank deficient, $X^\top X$ will not be invertible, so to solve for \mathbf{b} you should remove the linearly dependent columns of X .

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

Definition: $\text{rank}[X] = \dim[X] - N[X]$

- (a) Let $X \in N[A]$,

$$\begin{aligned}
AX &= \mathbf{0} \Rightarrow A^\top AX = \mathbf{0} \\
&\Rightarrow X \in N[A^\top A]
\end{aligned}$$

so, $N[A] \subseteq N[A^\top A]$.

(b) Let $X \in N[A^\top A]$,

$$\begin{aligned} A^\top AX &= \mathbf{0} \\ X^\top A^\top AX &= \mathbf{0} \\ (AX)^\top AX &= \mathbf{0} \Rightarrow AX = \mathbf{0} \\ &\Rightarrow X \in N[A] \end{aligned}$$

so, $N[A^\top A] \subseteq N[A]$.

By (a) and (b), $N[A] = N[A^\top A]$ and $\dim[A] = \dim[A^\top A]$. Therefore, $\text{rank}[X] = \text{rank}[X^\top X]$.

- (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.

$$\begin{aligned} \text{SSE} &= (\mathbf{y} - X\mathbf{b})^\top C(\mathbf{y} - X\mathbf{b}) \\ &= (\mathbf{y}^\top C^{\frac{1}{2}} - \mathbf{b}^\top X^\top C^{\frac{1}{2}})(C^{\frac{1}{2}}\mathbf{y} - C^{\frac{1}{2}}X\mathbf{b}) \\ &= \mathbf{y}^\top C\mathbf{y} - \mathbf{y}^\top CX\mathbf{b} - \mathbf{b}^\top X^\top C\mathbf{y} + \mathbf{b}^\top X^\top CX\mathbf{b} \\ \frac{\partial}{\partial \mathbf{b}} [\text{SSE}] &= \mathbf{0}_{p+1} - 2X^\top C\mathbf{y} + 2X^\top CX\mathbf{b} \stackrel{\text{set}}{=} \mathbf{0}_{p+1} \end{aligned}$$

$$\begin{aligned} X^\top CX\mathbf{b} &= X^\top C\mathbf{y} \\ \mathbf{b} &= (X^\top CX)^{-1}X^\top C\mathbf{y} \end{aligned}$$

- (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.

$$\begin{aligned} R^2 &= \frac{\text{SSR}}{\text{SST}} \\ &= \frac{\sum(\hat{y}_i - \bar{y})^2}{\sum(y_i - \bar{y})^2} \\ &= \frac{\sum \hat{y}_i^2 - n\bar{y}^2}{\sum(y_i - \bar{y})^2} \end{aligned}$$

$$\begin{aligned}
&= \frac{\sum (b_0 + b_1 x_i)^2 - n \bar{y}^2}{(n-1)s_y^2} \\
&= \frac{\sum (b_0^2 + 2b_0 b_1 x_i + b_1^2 x_i^2) - n \bar{y}^2}{(n-1)s_y^2} \\
&= \frac{nb_0^2 + 2b_0 b_1 \sum x_i + b_1^2 \sum x_i^2 - n \bar{y}^2}{(n-1)s_y^2} \\
&= \frac{(\bar{y} - b_1 \bar{x})^2 n + 2(\bar{y} - b_1 \bar{x}) b_1 n \bar{x} + b_1^2 \sum x_i^2 - n \bar{y}^2}{(n-1)s_y^2} \\
&= \frac{b_1^2 \sum x_i^2 - b_1^2 \bar{x}^2 n}{(n-1)s_y^2} \\
&= \frac{b_1^2 (\sum x_i^2 - \bar{x}^2 n)}{(n-1)s_y^2} \\
&= \frac{r^2 \frac{s_y^2}{s_x^2} \sum (x_i - \bar{x})^2}{(n-1)s_y^2} \\
&= r^2
\end{aligned}$$

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

$$\begin{aligned}
\bar{y} &= \frac{1}{n} \sum y_i \\
&= \frac{1}{n} \sum (b_0 + b_1 x_{1i} + \dots + b_p x_{pi} + e_i) \\
&= \frac{1}{n} \sum b_0 + \frac{1}{n} \sum b_1 x_{1i} + \dots + \frac{1}{n} \sum b_p x_{pi} + \frac{1}{n} \sum e_i \\
&= b_0 + b_1 \bar{x}_1 + \dots + b_p \bar{x}_p + \bar{e} \\
&= b_0 + b_1 \bar{x}_1 + \dots + b_p \bar{x}_p \\
&= g([1 \ \bar{x}_1 \ \bar{x}_2 \ \dots \ \bar{x}_p])
\end{aligned}$$

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

$$\begin{aligned}
\bar{e} &= \frac{1}{n} \sum (y_i - \hat{y}_i) \\
&= \frac{1}{n} \sum y_i - \frac{1}{n} \sum (b_0 + b_1 x_{1i} + \dots + b_p x_{pi}) \\
&= \bar{y} - (b_0 + b_1 \bar{x}_1 + \dots + b_p \bar{x}_p) \\
&= \bar{y} - \bar{y} \\
&= 0
\end{aligned}$$

- (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.

Show $\mathbf{b} = (X^\top X)^{-1} X^\top \mathbf{y}$ is \bar{y}_A if $x = A$, \bar{y}_B if $x = B$ and \bar{y}_C if $x = C$.

$$\begin{aligned} X^\top \mathbf{y} &= \begin{matrix} A \\ B \\ C \end{matrix} \begin{bmatrix} 1 & \dots & 1 & 0 & \dots & \dots & \dots & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \\ &= \begin{bmatrix} \sum y \mathbf{1}_{A=1} \\ \sum y \mathbf{1}_{B=1} \\ \sum y \mathbf{1}_{C=1} \end{bmatrix} \end{aligned}$$

$$\begin{aligned} X^\top X &= \begin{matrix} A \\ B \\ C \end{matrix} \begin{bmatrix} 1 & \dots & 1 & 0 & \dots & \dots & \dots & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 & 1 & \dots & 1 \end{bmatrix} \begin{matrix} A & B & C \\ \begin{bmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix} \end{matrix} \end{aligned}$$

$$\begin{aligned} &= \begin{bmatrix} n_A & 0 & 0 \\ 0 & n_B & 0 \\ 0 & 0 & n_C \end{bmatrix} \\ (X^\top X)^{-1} &= \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix} \\ \mathbf{b} = (X^\top X)^{-1} X^\top \mathbf{y} &= \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix} \begin{bmatrix} \sum y \mathbf{1}_{A=1} \\ \sum y \mathbf{1}_{B=1} \\ \sum y \mathbf{1}_{C=1} \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{n_A} \sum y \mathbf{1}_{A=1} \\ \frac{1}{n_B} \sum y \mathbf{1}_{B=1} \\ \frac{1}{n_C} \sum y \mathbf{1}_{C=1} \end{bmatrix} \\ &= \begin{bmatrix} \bar{y}_A \\ \bar{y}_B \\ \bar{y}_C \end{bmatrix} \end{aligned}$$

3 Orthogonal Projection and QR Decomposition

(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 $\forall n$.

(a) Symmetry: $I_n^\top = I_n$

(b) Idempotency: $I_n I_n = I_n$

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

\mathbb{R}^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?

There are $p+1$ degrees of freedom, i.e. $p+1$ independent parameters.

(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?

Yes,

$$\begin{aligned} \text{Proj}_{\text{colsp}[X]}[\mathbf{y}] &= X\mathbf{w} \\ \Rightarrow \begin{cases} \mathbf{x}_1^\top(\mathbf{y} - X\mathbf{w}) = 0 \\ \vdots \\ \mathbf{x}_n^\top(\mathbf{y} - X\mathbf{w}) = 0 \end{cases} \end{aligned}$$

$$\begin{aligned} \begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix} (\mathbf{y} - X\mathbf{w}) &= \mathbf{0}_n \\ X^\top(\mathbf{y} - X\mathbf{w}) &= \mathbf{0}_n \\ X^\top \mathbf{y} &= X^\top X \mathbf{w} \\ \mathbf{w} &= (X^\top X)^{-1} X^\top \mathbf{y} \end{aligned}$$

- (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}_{new} . Would \mathbf{e}_{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.

No, the projection mapping yields the least squares error. Any further projection onto the $\text{colsp}[\mathbf{X}]$ would only yield the same result.

- (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.

$$\begin{aligned}\mathbf{Q}^\top &= \mathbf{Q}^{-1} \Rightarrow \mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}^{-1} \mathbf{Q} = \mathbf{I}_n \\ \mathbf{Q}^\top \mathbf{Q} &= \begin{bmatrix} \mathbf{q}_1^\top \\ \vdots \\ \mathbf{q}_n^\top \end{bmatrix} [\mathbf{q}_1 \quad \dots \quad \mathbf{q}_n] \\ &= \begin{bmatrix} \mathbf{q}_1^\top \mathbf{q}_1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & \mathbf{q}_n^\top \mathbf{q}_n \end{bmatrix} = \mathbf{I}_n\end{aligned}$$

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

$$\begin{aligned}\text{Proj}_V[\mathbf{a}] &= \mathbf{H} \mathbf{a}, \text{ Let } \|\mathbf{v}_1\| = \dots = \|\mathbf{v}_d\| = 1 \\ \mathbf{H} &= [\mathbf{v}_1 \mathbf{v}_1^\top \quad \dots \quad \mathbf{v}_d \mathbf{v}_d^\top] \\ &= [v_{11} \mathbf{v}_1 \quad \dots \quad v_{1n} \mathbf{v}_1] + [v_{21} \mathbf{v}_2 \quad \dots \quad v_{2n} \mathbf{v}_2] + \dots + [v_{d1} \mathbf{v}_d \quad \dots \quad v_{dn} \mathbf{v}_d] \\ &= [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_d] \begin{bmatrix} \mathbf{v}_1^\top \\ \vdots \\ \mathbf{v}_d^\top \end{bmatrix} \\ &= \mathbf{Q} \mathbf{Q}^\top\end{aligned}$$

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

$$\begin{aligned}
 \text{Proj}_Q[a] &= Q(Q^\top Q)^{-1}Q^\top a \\
 &= QI_n Q^\top a \\
 &= QQ^\top a \\
 &= \sum q_i q_i^\top a \\
 &= \sum \text{Proj}_{q_i}[a]
 \end{aligned}$$

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

$$\begin{aligned}
 \text{SST} &= \text{SSR} + \text{SSE} \\
 \text{SST}_* &= \text{SSR}_* + \text{SSE}_* \\
 &= (\text{SSR} + k) + (\text{SSE} - k) \\
 &= \text{SSR} + \text{SSE} + k - k \\
 &= \text{SSR} + \text{SSE} \\
 &= \text{SST}
 \end{aligned}$$

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

4 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 $\boldsymbol{\beta} = [\beta_0 \ \beta_1 \ \dots \ \beta_p]^\top$, 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 σ^2 is known.

- (a) [E.C.] Show that $\mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$.
- (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. $\boldsymbol{\mathcal{E}}$. 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 σ^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 $\mathbf{U} \sim \mathcal{N}_n(n, \mathbf{I}_n)$ independent of $\mathbf{V} \sim \mathcal{N}_n(n, \mathbf{I}_n)$. Let θ be the r.v. model of the angle between \mathbf{U} and \mathbf{V} . How is θ distributed?