

MATH 390.4 / 650.2 Spring 2020 Homework #3

Tyron Samaroo

Saturday 9th May, 2020

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?

PECOTA seems to be most similar to KNN(k nearest neighbors) based on his premises of comparing players by points and finding the difference.

- (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 since some players can perform really well early on and better than more experience players.

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

Baseball scouts were able to do better than a prediction system like PECOTA because scouts use a hybrid approach. They have access to more information than statistics alone.

- (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?

They had other tools like radar guns but didn't have completely functional three-dimensional cameras. Also Pitch f/x seems expensive and Moneyball could do the same.

- (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 weather prediction is that it can add a lot of dimensions and make it very difficult and messy since its a dynamic system. There is also many inaccuracies in the data \mathbb{D}

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

Weatherman lie about the chance of rain because the public won't believe then either way. Weatherman would rather get chance rain wrong because people wont curse them hence economic incentive. If you want honest forecast you should go to The Weather Channel since they can use all of the government's raw data.

- (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 do not know the true function t with its causal inputs z 's. We also don't have enough data in X and knowledge of features x_i to predict earthquakes especially ones with higher magnitude. We end up building a model that over-fits and will then have huge errors and bad prediction.

- (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?

The nonsense predictor in model Silver describe is someone who assumes that since the locks with certain color unlocks with a certain code then all locks with same color will unlock with the same code. This is not true since other locks will not unlock with the same code and have the same color. So this is over fitting. You assume you learned something but in fact you did not.

- (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?

Adding more features to the model can create noise and not do what is intended to do. We end up fitting noise than the signal and perform badly in real world.

- (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.

It's much harder to find something that identifies the signal; variables that are leading indicators in one economic cycle will not be in the next. There are other decisions such political decisions you have to worry about other than economic ones. Its too difficult to understand the true function with its causal inputs and to even estimate it.

- (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.

$$A\vec{c} = \begin{bmatrix} a_{11}c_1 + \dots + a_{1n}c_n \\ a_{21}c_1 \dots a_{2n}c_n \\ \vdots \vdots \vdots \\ a_{n1}c_1 \dots a_{nn}c_n \end{bmatrix}$$

$$\frac{\partial}{\partial \mathbf{c}} [\mathbf{c}^\top A \mathbf{c}] = \begin{bmatrix} \frac{\partial}{\partial \mathbf{c}} [A\vec{c}_1] \\ \vdots \\ \frac{\partial}{\partial \mathbf{c}} [A\vec{c}_n] \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial \mathbf{c}} [A\vec{c}_1] \\ \vdots \\ \frac{\partial}{\partial \mathbf{c}} [A\vec{c}_n] \end{bmatrix} = 2 \begin{bmatrix} \vec{a}_1 \\ \vdots \\ \vec{a}_n \end{bmatrix} \vec{c} = 2A\vec{c}$$

- (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.

Proof. We know that $\vec{\hat{y}} = X\vec{w}$

$$\begin{aligned} SSE &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= (\vec{y} - \vec{\hat{y}})^T (\vec{y} - \vec{\hat{y}}) \\ &= (\vec{y}^T - \vec{\hat{y}}^T) (\vec{y} - \vec{\hat{y}}) \\ &= \vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \vec{w}^T X^T \vec{w} \end{aligned}$$

Then we can do

$$\begin{aligned} \frac{\partial}{\partial \vec{w}} [SSE] &= \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \vec{w}^T X^T \vec{w}] \\ &= \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y}] - 2 \frac{\partial}{\partial \vec{w}} [\vec{w}^T X^T \vec{y}] + \frac{\partial}{\partial \vec{w}} [\vec{w}^T X^T X \vec{w}] \\ &= \vec{0} - 2X^T \vec{y} + 2X^T X \vec{w} \\ X^T \vec{y} &= X^T X \vec{w} \\ \vec{w} = \vec{b} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y} \end{aligned}$$

□

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

When $p=1$ we use the lost function to fit the b's. We can use the SSE or sum of squared error formula and do some manipulation.

$$\hat{y} = b_0 + b_1(x_i)$$

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Substitute \hat{y}

$$\begin{aligned} &= \sum_{i=1}^n y_i - (b_0 + b_1 x_i)^2 \\ &= \sum_{i=1}^n y_i^2 + b_0^2 + b_1^2 x_i^2 - 2y_i b_0 - 2y_i b_1 x_i + 2b_0 b_1 x_i \\ &= \sum_{i=1}^n y_i^2 + n b_0^2 + b_1^2 \sum x_i^2 - 2n \bar{y} b_0 - 2b_1 \sum x_i y_i + 2b_0 b_1 n \bar{x} \end{aligned}$$

Then we take partial of b_0

$$\frac{\partial}{\partial b_0} [SSE] = 2n b_0 - 2n \bar{y} = 0$$

$$b_0 = \bar{y} - b_1 \bar{x}$$

Then we take partial of b_1

$$\frac{\partial}{\partial b_1} [SSE] = 2b_1 \sum x_i^2 - 2 \sum x_i y_i + 2b_0 n \bar{x}$$

Substitute in b_0

$$\begin{aligned} &= b_1 \sum x_i^2 - \sum x_i y_i + (\bar{y} - b_1 \bar{x}) n \bar{x} \\ &= b_1 \sum x_i^2 - \sum x_i y_i + n \bar{x} \bar{y} - b_1 n \bar{x}^2 \\ b_1 (\sum x_i^2 - n \bar{x}^2) &= \sum x_i y_i - n \bar{x} \bar{y} \\ b_1 &= \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \\ b_1 &= \frac{(n-1) S_x y}{(n-1) S_x^2} \\ &= r \frac{S_y}{S_x} \end{aligned}$$

Now since $b_1 = r \frac{S_y}{S_x}$ and $b_0 = \bar{y} - b_1 \bar{x}$

we substitute in b_0 to get $b_0 = \bar{y} - r \frac{S_y}{S_x} \bar{x}$

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

If X is rank deficient we can't solve for \mathbf{b} because it will not be invertible. In order to solve this problem you had to remove anything that is linearly dependent in X

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

$X^T X$ has the dimension of $\mathbb{R}^{p+1 \times p+1}$ When $X^T X$ is full rank it is invertible and rank of X is the same, $p+1$. All the columns are all linear independent so the rank of $X^T X$ and X are both $p+1$

- (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.
- (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.

$$\begin{aligned}
 R^2 &= \frac{\text{SSR}}{\text{SST}} \\
 &= \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} \\
 &= \frac{\sum \hat{y}_i^2 - n\bar{y}^2}{\sum (y_i - \bar{y})^2} \\
 &\text{Substitute in for } \hat{y} = b_0 + b_1 x_i \\
 &= \frac{\sum (b_0 + b_1 x_i)^2 - n\bar{y}^2}{(n-1)s_y^2} \\
 &= \frac{\sum (b_0^2 + 2b_0 b_1 x_i + b_1^2 x_i^2) - n\bar{y}^2}{(n-1)s_y^2} \\
 &= \frac{nb_0^2 + 2b_0 b_1 \sum x_i + b_1^2 \sum x_i^2 - n\bar{y}^2}{(n-1)s_y^2} \\
 &= \frac{(\bar{y} - b_1 \bar{x})^2 n + 2(\bar{y} - b_1 \bar{x}) b_1 n \bar{x} + b_1^2 \sum x_i^2 - n\bar{y}^2}{(n-1)s_y^2} \\
 &= \frac{b_1^2 \sum x_i^2 - b_1^2 \bar{x}^2 n}{(n-1)s_y^2} \\
 &= \frac{b_1^2 (\sum x_i^2 - \bar{x}^2 n)}{(n-1)s_y^2}
 \end{aligned}$$

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

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

Since \bar{y} is the average of all the y_i

$$\bar{y} = \frac{1}{n} \sum y_i$$

y_i can be $b_0 + b_1x_1 + b_2x_2 \dots b_px_p$ where each x is offset by some b

$$\begin{aligned}
&= \frac{1}{n} \sum b_0 + b_1x_1 + \dots + b_px_p \\
&= \frac{1}{n} \sum b_0 + \frac{1}{n} \sum b_1x_1 + \dots + \frac{1}{n} \sum b_px_p \\
&= b_0 + b_1\bar{x}_1 + \dots + b_p\bar{x}_p
\end{aligned}$$

Now b_0 will be some constant vector $\vec{1}$.

So the b 's are apart of some function g that gives

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

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

$$\begin{aligned}
\bar{e} &= \frac{1}{n} \sum (y_i - \hat{y}_i) \\
&= \frac{1}{n} \sum y_i - \frac{1}{n} \sum b_0 + b_1x_1 + \dots + b_px_p
\end{aligned}$$

Multiplying each term by $\frac{1}{n}$ would make the x 's become average \bar{x} 's

$$= \frac{1}{n} \sum y_i - b_0 + b_1\bar{x}_1 + \dots + b_p\bar{x}_p$$

From part h we can see that $\bar{y} = b_0 + b_1\bar{x}_1 + \dots + b_p\bar{x}_p$

$$= \bar{y} - \bar{y}$$

$$= 0$$

(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.

$$\mathbf{X}^T \mathbf{X} = \begin{matrix} & \begin{matrix} A & B & C \end{matrix} \\ \begin{matrix} A \\ B \\ C \end{matrix} & \begin{bmatrix} 1 & \dots & 1 & 0 & \dots & \dots & \dots & \dots & 0 \\ 0 & \dots & 0 & 1 & \dots & 1 & 0 & \dots & 0 \\ 0 & \dots & \dots & \dots & \dots & 0 & 1 & \dots & 1 \end{bmatrix} \end{matrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix}$$

This only leaves the diagonal

$$= \begin{bmatrix} n_A & 0 & 0 \\ 0 & n_B & 0 \\ 0 & 0 & n_C \end{bmatrix}$$

Now the inverse $(\mathbf{X}^T \mathbf{X})^{-1}$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix}$$

Now the $\mathbf{X}^T y$

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

$$= \begin{bmatrix} \sum y \mathbf{1}_{A=1} \\ \sum y \mathbf{1}_{B=1} \\ \sum y \mathbf{1}_{C=1} \end{bmatrix}$$

Now put $(\mathbf{X}^T \mathbf{X})^{-1}$ and $\mathbf{X}^T y$ together

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T y = \begin{bmatrix} \frac{1}{n_A} & 0 & 0 \\ 0 & \frac{1}{n_B} & 0 \\ 0 & 0 & \frac{1}{n_C} \end{bmatrix} \begin{bmatrix} \sum y \mathbf{1}_{A=1} \\ \sum y \mathbf{1}_{B=1} \\ \sum y \mathbf{1}_{C=1} \end{bmatrix}$$

$$= \begin{bmatrix} \frac{1}{n_A} \sum y \mathbf{1}_{A=1} \\ \frac{1}{n_B} \sum y \mathbf{1}_{B=1} \\ \frac{1}{n_C} \sum y \mathbf{1}_{C=1} \end{bmatrix}$$

$$= \begin{bmatrix} \bar{y}_A \\ \bar{y}_B \\ \bar{y}_C \end{bmatrix}$$

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 $\forall n$.

I_n is an orthogonal project is it is both symmetric and idempotent. If $I_n^T = I_n$ then it is symmetric. If $I_n I_n = I_n$ then it is idempotent.

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

I_n projects onto the subspace \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?

The resulting model will have $p+1$ degrees of freedom.

- (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?

- (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.

- (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.

$$\begin{aligned}\mathbf{Q}^\top &= \mathbf{Q}^{-1} \Rightarrow \mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}^{-1} \mathbf{Q} = I_n \\ \mathbf{Q}^\top \mathbf{Q} &= \begin{bmatrix} q_1^\top \\ \vdots \\ q_n^\top \end{bmatrix} \begin{bmatrix} q_1 & \dots & q_n \end{bmatrix} \\ &= \begin{bmatrix} q_1^\top q_1 & 0 & \dots & 0 \\ 0 & \ddots & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & 0 & q_n^\top q_n \end{bmatrix} = I_n\end{aligned}$$

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

$$\text{Proj}_v[a] = \frac{\vec{v}\vec{v}^T}{\|\vec{v}\|^2} a = H\vec{a}$$

If $\text{Proj}_v[a]$ is in the column space of \mathbf{V} where $V = \{\vec{v}_1 \dots \vec{v}_k\} \in \mathbb{R}^{n \times k}$

$\text{Proj}_v[a] \in$ the column space of V then $\text{Proj}_v[a] = w_1\vec{v}_1 + \dots + w_k\vec{v}_k = V\vec{w}$

Now if V is orthogonal

$$\begin{aligned} \text{Proj}_v[a] &= \text{Proj}_{v_1}[a] + \text{Proj}_{v_2}[a] \\ &= \left(\frac{\vec{v}_1\vec{v}_1^T}{\|\vec{v}_1\|^2} + \frac{\vec{v}_2\vec{v}_2^T}{\|\vec{v}_2\|^2} \right) \vec{a} \end{aligned}$$

If V is orthonormal length = 1 then

$$\begin{aligned} \text{Proj}_v[a] &= (\vec{v}_1\vec{v}_1^T + \vec{v}_2\vec{v}_2^T) \vec{a} \\ &= \left(\begin{bmatrix} v_{11}^2 v_1 & v_{11}v_{12} & \dots & v_{11}v_{1n} \\ & \ddots & & \vdots \\ \vdots & & \ddots & \\ v_{1n}v_1 & \dots & & v_1^2 v_n \end{bmatrix} + \dots + \begin{bmatrix} v_{n1}^2 v_1 & v_{n1}v_{n2} & \dots & v_{n1}v_{nn} \\ & \ddots & & \vdots \\ \vdots & & \ddots & \\ v_{nn}v_{n1} & \dots & & v_n^2 v_n \end{bmatrix} \right) \vec{a} \\ &= \begin{bmatrix} \vec{v}_1 & \vec{v}_2 & \dots & \vec{v}_n \end{bmatrix} \begin{bmatrix} \vec{v}_1^T \\ \vdots \\ \vec{v}_n^T \end{bmatrix} \vec{a} \\ &= \mathbf{Q}\mathbf{Q}^T \vec{a} \end{aligned}$$

- (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} .

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

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

$$SST = SSR + SSE$$

Adding another column does not change the formula since it is another y value "i"

SSR is sum of the squared regression or how far the prediction is from the \bar{y} and SSE is sum of

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

- (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 $\boldsymbol{\beta} = [\beta_0 \ \beta_1 \ \dots \ \beta_p]^\top$, 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 $\mathbf{U} \sim \mathcal{N}_n(\mathbf{0}_n, \mathbf{I}_n)$ independent of $\mathbf{V} \sim \mathcal{N}_n(\mathbf{0}_n, \mathbf{I}_n)$. Let θ be the r.v. model of the angle between \mathbf{U} and \mathbf{V} . How is θ distributed?