

MATH 343 / 643 Homework #2

Natasha Watson

Friday 11th April, 2025

Problem 1

This problem is about OLS estimation in regression. You can assume that

$\mathbf{X} := [\mathbf{1}_n \mid \mathbf{x}_{.1} \mid \dots \mid \mathbf{x}_{.p}]$ with column indices $0, 1, \dots, p$ and row indices $1, 2, \dots, n$

$\mathbf{H} := \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$

$\mathbf{B} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y} = \mathbf{X}\mathbf{B}$

$\mathbf{E} := \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I}_n - \mathbf{H})\mathbf{Y}$

where the entries of \mathbf{X} are assumed fixed and known and the entries of $\boldsymbol{\beta}$ are the unknown parameter).

- (a) [easy] When we “do inference” for the linear model, what is the parameter vector?

$\vec{\boldsymbol{\beta}}$

The purpose of OLS estimation and inference is to estimate and draw conclusions about the components of $\vec{\boldsymbol{\beta}}$. IT is the $p+1$ dimensional vector of true regression coefficients of the model.

- (b) [easy] When we “do inference” for the linear model, what are considered the fixed and known quantities?

The design matrix \mathbf{X} of predictor values, including an intercept column of 1's. The sample size n and the number of features p .

- (c) [easy] When we “do inference” for the linear model, what are considered the random quantities? And what is the notation for their corresponding realizations?

The error term ϵ is random, we consider ϵ_i to be a random variable.

(d) [easy] What is the “core assumption” in which the classic linear model inference follows?

The core assumption is the combination of 3 main assumptions about the ϵ_i 's, and becomes our base case:

- (1) Let $\epsilon_i = u_{1i} + u_{2i} + u_{3i} + \dots$, i.e. the sum of many unknown values is likely independent $\Rightarrow \epsilon_i$ is realized from an independent normal distribution
- (2) The mean of the Random Variable ϵ_i was drawn from is zero (mean centered)
- (3) The variance is equal $\forall i, \sigma^2$ (Homoscedacity Assumption)

The amalgamation of these three assumptions leads us to the core assumption:

$$\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_n \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

Where $\epsilon_1, \epsilon_2, \dots, \epsilon_n$ are its realizations $\leftrightarrow \vec{\mathcal{E}} \sim \mathcal{N}_n(\vec{0}_n, \sigma^2 I_n)$

We are able to say via the core assumption that (1) the relationship between predictors and the response is linear in parameters, the error is normally distributed, there is a constant variance, and the errors are uncorrelated (ind)

(e) [easy] From the core assumption, derive the distribution of \vec{B} .

We recall from Math 340 that if $\vec{X} \sim \mathcal{N}_n(\mu, \Sigma)$ and if $A \in \mathbb{R}^{k \times n}, b \in \mathbb{R}^k$, then $\vec{Y} = \vec{b} + A\vec{X} \sim \mathcal{N}_k(A\vec{\mu} + \vec{b}, A\Sigma A^\top)$

So if $\vec{Y} = \mathbf{X}\vec{\beta} + \vec{\mathcal{E}} \sim \mathcal{N}_n(\mathbf{X}\vec{\beta}, \sigma^2 I_n)$, then we can think of \vec{B} as being a random vector that realizes \vec{b} . If we recall that OLS estimator for $\vec{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{Y}$ we can derive its distribution like so:

$$\vec{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{Y}$$

Recall $\vec{Y} = \mathbf{X}\vec{\beta} + \vec{\mathcal{E}}$

$$\vec{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \vec{\beta} + \vec{\mathcal{E}}$$

Distribute

$$\Rightarrow (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \vec{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{\mathcal{E}}$$

Notice that $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} = \mathbf{I}$ and so this term cancels out. We are then left with:

$$\Rightarrow \vec{\beta} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{\mathcal{E}}$$

And

$$\vec{B} \sim \mathcal{N}_{p+1}(\vec{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$$

- (f) [easy] From this result, derive the distribution of B_j .

Recall that the elements of a multivariate normal random variable are distributed as univariate normals. Therefore,

$$B_j \sim \mathcal{N}(\beta_j, \sigma^2(\mathbf{X}^\top \mathbf{X})_{j,j}^{-1})$$

Where $(\mathbf{X}^\top \mathbf{X})_{j,j}^{-1}$ is the j,j-th index of the covariance matrix, because the diagonal entries give the variances of each component

This distribution is the foundation for confidence intervals for β_j and let's sue perform t-tests once we estimate σ^2 .

- (g) [easy] From this result, derive the distribution of B_j standardized.

B_j standardized:

$$\frac{B_j - \beta_j}{\sqrt{\sigma^2(\mathbf{X}^\top \mathbf{X})_{j,j}^{-1}}} \sim \mathcal{N}(0, 1)$$

Recall, this is when variance is known! If not then this becomes a t-distribution with n-(p+1) degrees of freedom and we use $\hat{\sigma}^2$

- (h) [easy] from the core assumption, derive the distribution of $\hat{\mathbf{Y}}$.

Recall from 342 $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$ where $\mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}\mathbf{X}^\top$ is the hat matrix. If we express $\hat{\mathbf{Y}}$ in terms of ϵ , and we know $\vec{\mathbf{Y}} = \mathbf{X}\vec{\beta} + \vec{\mathcal{E}}$:

$$\hat{\mathbf{Y}} = \mathbf{X}\vec{\beta} = \mathbf{H}\mathbf{Y} = \mathbf{H}(\mathbf{X}\vec{\beta} + \vec{\mathcal{E}})$$

Distribute

$$\mathbf{H}\mathbf{X}\vec{\beta} + \mathbf{H}\vec{\mathcal{E}}$$

Recall \mathbf{H} is in the colsp $[\mathbf{X}]$ and so $\mathbf{H}\mathbf{X} = \mathbf{X}$:

$$\hat{\mathbf{Y}} = \mathbf{X}\vec{\beta} + \mathbf{H}\vec{\mathcal{E}} \sim ?$$

Recall from math 340:

$$\vec{U} \sim \mathcal{N}_n(\vec{0}_n, \mathcal{E}) \quad \mu \in \mathbb{R}^m; \mathbf{A} \in \mathbb{R}^{m \times n} \Rightarrow \vec{\mu} + \mathbf{A}\vec{U} \sim \mathcal{N}_m(\vec{\mu}, \mathbf{A}\mathcal{E}\mathbf{A}^\top)$$

Applying this shift and scale method we see,

$$\hat{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\vec{\beta}, \mathbf{H}(\sigma^2 I_n)\mathbf{H}^\top)$$

If we pull out the scalar σ^2 in the variance term, then we are left with

$$\sigma^2(\mathbf{H}(I_n)\mathbf{H}^\top) \Rightarrow \sigma^2(\mathbf{H}\mathbf{H}^\top)$$

By idempotency and symmetry of the hat matrix, this reduces to:

$$\sigma^2 \mathbf{H}$$

And so the final distribution of $\hat{\mathbf{Y}}$ is:

$$\hat{\mathbf{Y}} \sim \mathcal{N}_n(\mathbf{X}\vec{\beta}, \mathbf{H}\sigma^2)$$

- (i) [easy] From this result, derive the distribution of \hat{Y}_i .

Recall that the elements of a multivariate normal random variable are distributed as univariate normals. Therefore,

$$\hat{Y}_i \sim \mathcal{N}(\vec{x}_i \vec{\beta}, \sigma^2 \mathbf{H}_{i,i})$$

- (j) [easy] From this result, derive the distribution of \hat{Y}_i standardized.

Subtract the mean and divide by the standard deviation:

$$\frac{\hat{Y}_i - \vec{x}_i \vec{\beta}}{\sqrt{\sigma^2 \mathbf{H}_{i,i}}} \sim \mathcal{N}(0, 1)$$

- (k) [easy] from the core assumption, derive the distribution of \mathbf{E} .

Recall we assumed $\mathcal{E} \sim \mathcal{N}_n(\vec{0}_n, \sigma^2 \mathbf{I}_n)$. Rewriting the residuals, we get:

$$\vec{E} = \vec{Y} - \hat{\mathbf{Y}} = \vec{Y} - \mathbf{X}\vec{B}$$

Recall

$$\mathbf{Y} - \mathbf{X}\mathbf{B}$$

Substitute in OLS estimator $\mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$:

$$\mathbf{Y} - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

Factor out \mathbf{Y} :

$$= (\mathbf{I}_n - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{Y}$$

Recognize the hat matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$:

$$= (\mathbf{I}_n - \mathbf{H})\mathbf{Y}$$

Use this:

$$\begin{aligned} (\mathbf{I} - \mathbf{H})(\mathbf{X}\vec{\beta} + \vec{\mathcal{E}}) &= (\mathbf{I} - \mathbf{H})\mathbf{X}\vec{\beta} + (\mathbf{I} - \mathbf{H})\vec{\mathcal{E}} \sim \mathcal{N}_n(\vec{0}_n, (\mathbf{I} - \mathbf{H})\sigma^2\mathbf{I}_n(\mathbf{I} - \mathbf{H})^\top) \\ \Rightarrow \vec{E} &\sim \mathcal{N}_n(\vec{0}_n, \sigma^2(\mathbf{I} - \mathbf{H})) \end{aligned}$$

- (l) [easy] From this result, derive the distribution of E_i .

Recall that the elements of a multivariate normal random variable are distributed as univariate normals. Therefore,

$$E_i \sim \mathcal{N}(0, \sigma^2(1 - \mathbf{H}_{i,i}))$$

- (m) [easy] From this result, derive the distribution of E_i standardized.

Subtract the mean and divide by the SD:

$$\frac{E_i}{\sqrt{\sigma^2(1 - \mathbf{H}_{i,i})}} \sim \mathcal{N}(0, 1)$$

- (n) [easy] From the core assumption, show that $\frac{1}{\sigma^2}\vec{\mathcal{E}}^\top\vec{\mathcal{E}} \sim \chi_n^2$.

$$\text{Recall } \vec{Z} \sim \mathcal{N}_n(\vec{0}_n, \sigma^2\mathbf{I}_n) \Rightarrow \vec{Z}^\top\vec{Z} \sim \chi_n^2$$

This is equivalent to saying $\vec{Z}^\top\vec{Z} = \frac{1}{\sigma^2}\vec{\mathcal{E}}^\top\vec{\mathcal{E}} \sim \chi_n^2$, where $\vec{Z} = \frac{1}{\sigma^2}\vec{\mathcal{E}}$

- (o) [easy] Let $\mathbf{B}_1 = \mathbf{H}$ and let $\mathbf{B}_2 = \mathbf{I}_n - \mathbf{H}$. Justify the use of Cochran's theorem and then find the distributions of $\frac{1}{\sigma^2}\vec{\mathcal{E}}^\top\mathbf{B}_1\vec{\mathcal{E}}$ and $\frac{1}{\sigma^2}\vec{\mathcal{E}}^\top\mathbf{B}_2\vec{\mathcal{E}}$.

We know that $\mathbf{B}_1 + \mathbf{B}_2 = \mathbf{I}_n$ and because \mathbf{H} is the orthogonal projection matrix onto the colsp $[\mathbf{X}]$, then the $\mathbf{rank}[\mathbf{B}_1] = p + 1$ and $\mathbf{rank}[\mathbf{B}_2] = n - (p + 1)$. We can use Cochran's theorem because $\mathbf{rank}[\mathbf{B}_1] + \mathbf{rank}[\mathbf{B}_2] = n$

Therefore,

$$\vec{Z}^\top\mathbf{B}_1\vec{Z} \sim \chi_{p+1}^2$$

and

$$\vec{Z}^\top\mathbf{B}_2\vec{Z} \sim \chi_{n-(p+1)}^2$$

Or,

$$\frac{1}{\sigma^2}\vec{\mathcal{E}}^\top\mathbf{H}\vec{\mathcal{E}} \sim \chi_{p+1}^2$$

and

$$\frac{1}{\sigma^2}\vec{\mathcal{E}}^\top(\mathbf{I} - \mathbf{H})\vec{\mathcal{E}} \sim \chi_{n-(p+1)}^2$$

(p) [easy] Show that $\frac{1}{\sigma^2} \mathbf{E}^\top \mathbf{B}_1 \mathbf{E} = \frac{1}{\sigma^2} \|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2$.

Recall, $\mathbf{B} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$

$$= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{X} \vec{\boldsymbol{\beta}} + \vec{\mathcal{E}}) = \mathbf{I} \vec{\boldsymbol{\beta}} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{\mathcal{E}}$$

$$\Rightarrow \mathbf{B} - \boldsymbol{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{\mathcal{E}}$$

Multiply both sides by \mathbf{X}

$$\Rightarrow \mathbf{X}(\mathbf{B} - \boldsymbol{\beta}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \vec{\mathcal{E}}$$

Notice $(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}$ is the hat matrix, which is in the colsp $[\mathbf{X}]$. Recall $\mathbf{B}_1 = \mathbf{H}$. Take the norm and square both sides:

$$\|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2 = \|\mathbf{B}_1 \vec{\mathcal{E}}\|^2$$

By idempotency and symmetry:

$$\mathbf{E}^\top \mathbf{B}_1 \mathbf{E}$$

Divide both sides by σ^2 :

$$\boxed{\frac{1}{\sigma^2} \|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2 = \frac{1}{\sigma^2} \mathbf{E}^\top \mathbf{B}_1 \mathbf{E}}$$

(q) [harder] Why is the term $\|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2$ used to measure the model's "estimation error"?

Because when we fit a linear model we are trying to estimate the true mean response vector $\mathbf{X}\boldsymbol{\beta}$ by using $\mathbf{X}\mathbf{B}$ and so the difference between the two $\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})$ is the error in our estimate of the true mean response.

We take the squared norm as we would in the SSE as it makes mathematical sense (leads to closed form solution of best $\boldsymbol{\beta}$)

(r) [easy] Show that $\frac{1}{\sigma^2} \mathbf{E}^\top \mathbf{B}_2 \mathbf{E} = \frac{1}{\sigma^2} \|\mathbf{E}\|^2$.

We saw $\frac{1}{\sigma^2} \vec{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H}) \vec{\mathcal{E}} \sim \chi_{n-p-1}^2$. By idempotency we can express this as:

$$\frac{1}{\sigma^2} \vec{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H}) (\mathbf{I} - \mathbf{H}) \vec{\mathcal{E}}$$

By symmetry:

$$\frac{1}{\sigma^2} \vec{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H})^\top (\mathbf{I} - \mathbf{H}) \vec{\mathcal{E}}$$

Which can be expressed as the norm squared

$$\Rightarrow \left\| \frac{1}{\sigma^2} (\mathbf{I} - \mathbf{H}) \vec{\mathcal{E}} \right\|^2$$

Recall $\vec{Y} = \mathbf{X}\vec{\beta} + \vec{\mathcal{E}}$ and therefore $\vec{\mathcal{E}} = \vec{Y} - \mathbf{X}\vec{\beta}$. Substitute this information:

$$\Rightarrow \frac{1}{\sigma^2} \left\| (\mathbf{I} - \mathbf{H})(\vec{Y} - \mathbf{X}\vec{\beta}) \right\|^2$$

Distribute:

$$\frac{1}{\sigma^2} \left\| (\mathbf{I} - \mathbf{H})\vec{Y} - (\mathbf{I} - \mathbf{H})\mathbf{X}\vec{\beta} \right\|^2$$

Recall \mathbf{H} keeps anything in the colsp $[\mathbf{X}]$ while $\mathbf{I} - \mathbf{H}$ kills anything in the colsp $[\mathbf{X}]$, therefore, $\mathbf{X}\vec{\beta}$ vanishes when multiplied by $(\mathbf{I} - \mathbf{H})$. Also, because $(\mathbf{I} - \mathbf{H})\vec{Y}$ is the error, it goes to 0. Therefore we are left with:

$$\frac{1}{\sigma^2} \mathcal{E}^\top \mathbf{B}_2 \mathcal{E} = \frac{1}{\sigma^2} \|\mathbf{E}\|^2$$

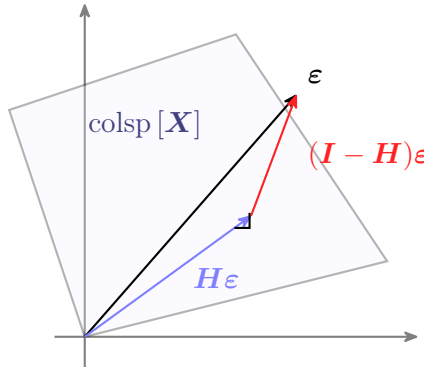
(s) [harder] In what scenarios is $\mathcal{E}^\top \mathbf{B}_1 \mathcal{E} > \mathcal{E}^\top \mathbf{B}_2 \mathcal{E}$?

Recall, we said $\mathbf{B}_1 = \mathbf{H}$, the projection onto the colsp $[\mathbf{X}]$ and $\mathbf{B}_2 = \mathbf{I} - \mathbf{H}$, the projection of the orthogonal complement of the hat matrix (residual space).

If $\vec{\mathcal{E}}$ lies closer to the colsp $[\mathbf{X}]$, then the projection onto that space will be larger and the equality will hold. If $\vec{\mathcal{E}}$ is orthogonal to the colsp $[\mathbf{X}]$ then the opposite holds because $\mathbf{H}\vec{\mathcal{E}}$ will equal 0.

If the former is the case (equality is true), then that means most of the error lies in the model space and a lot of the error is aligned with the predictors i.e. model is mistaking noise for signal and we might have some overfitting (bad model). If the latter is true, then most of the error is orthogonal to the predictors and so the model is not capturing the error (better model)

(t) [harder] Draw an illustration of \mathcal{E} being orthogonally projected onto colsp $[\mathbf{X}]$ via projection matrix \mathbf{H} . Use the previous answers to denote the quantities of the projection and the error of the projection.



- (u) [difficult] A good linear model has a large or a small projection of the error? Discuss.

This is kind of a continuation of part (s). A "good" model has a small projection of the error because if the projection of the error lies in the space which the model is trying to explain, then that means the model is fitting the noise (overfitting). Even if the model has low residuals and appears to fit the data well, this is an indicator that it is actually just capturing random variation as if it is meaningful which will lead to poor generalization to future/new data. i.e. big error projection = bad!

- (v) [easy] Find $\mathbb{E} \left[\frac{1}{\sigma^2} \|\mathbf{E}\|^2 \right]$.

Recall we said in 340 that if $X \sim \chi_k^2$ then $\mathbb{E}[X] = k$.

Because $\frac{1}{\sigma^2} \|\mathbf{E}\|^2 \sim \chi_{n-(p+1)}^2$, $\mathbb{E} \left[\frac{1}{\sigma^2} \|\mathbf{E}\|^2 \right] = n - p - 1$

- (w) [easy] Show that $\frac{\|\mathbf{E}\|^2}{n-(p+1)}$ is an unbiased estimate of σ^2 .

In 342 we saw that the SSE

$$= \vec{e}^\top \vec{e} = \|\vec{e}\|^2$$

. We proved that the MSE is an unbiased estimate for σ^2 , and the

$$\frac{\|\mathbf{E}\|^2}{n - (p + 1)} = \frac{SSE}{n - (p + 1)} = \text{MSE}$$

- (x) [easy] Prove that $\frac{\sqrt{n - (p + 1)}(B_j - \beta_j)}{\|\mathbf{E}\| \sqrt{(\mathbf{X}^T \mathbf{X})_{j,j}^{-1}}} \sim T_{n-(p+1)}$.

We found via Cochran's theorem that:

$$\frac{1}{\sigma^2} \|\mathbf{E}\|^2 \sim \chi_{n-(p+1)}^2$$

And we know from our previous work that:

$$\frac{B_j - \beta_j}{\sqrt{\sigma^2 (\mathbf{X}^{top} \mathbf{X})_{j,j}^{-1}}} \sim \mathcal{N}(0, 1)$$

We know via 340 that dividing a normal by a chi-squared gives us a student-t distribution, with the same degrees of freedom as the chi-squared we divide by. Thus,

$$\frac{\sqrt{n - (p + 1)}(B_j - \beta_j)}{\|\mathbf{E}\| \sqrt{(\mathbf{X}^T \mathbf{X})_{j,j}^{-1}}} \sim T_{n-(p+1)}$$

- (y) [easy] Let $H_0 : \beta_j = 0$. Find the test statistic using the fact from the previous question.
 Let s_e denote $RMSE := \sqrt{MSE} := \sqrt{SSE/(n - (p + 1))} = \sqrt{\|e\|^2 / (n - (p + 1))}$.

$$\hat{t} = \frac{b_j}{s_e \sqrt{(\mathbf{X}^T \mathbf{X})_{j,j}^{-1}}}$$

- (z) [easy] Consider a new parameter of interest $\mu_* := \mathbb{E}[Y_*] = \mathbf{x}_* \boldsymbol{\beta}$, this is the expected response for a unit with measurements given in row vector \mathbf{x}_* whose first entry is 1.

Prove that $\frac{\hat{Y}_* - \mu_*}{\sigma \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim \mathcal{N}(0, 1)$.

$$\Rightarrow \hat{\mathbf{Y}} = \vec{X}_* \mathbf{B} \sim ?$$

Recall that we previously showed:

$$\mathbf{B} \sim \mathcal{N}_{p+1}(\vec{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

And so a linear transformation of this will also be distributed normally:

$$\hat{\mathbf{Y}} = \vec{x}_* \mathbf{B} \sim \mathcal{N}(\vec{x}_* \vec{\beta}, \sigma^2 \vec{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \vec{x}_*^\top)$$

Subtracting the mean and dividing by the standard deviation, we get:

$$\frac{\hat{Y}_* - \mu_*}{\sigma \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim \mathcal{N}(0, 1)$$

- (aa) [easy] Prove that $\frac{\sqrt{n - (p + 1)}(\hat{Y}_* - \mu_*)}{\|\mathbf{E}\| \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim T_{n-(p+1)}$.

We know

$$\frac{\hat{Y}_* - \mathbf{x}_* \mathbf{B}}{\sigma \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim \mathcal{N}(0, 1)$$

Where $\mathbf{x}_* \mathbf{B} = \mu_*$. And,

$$\sqrt{\frac{\frac{1}{\sigma^2} \|\vec{E}\|^2}{n - p - 1}} \sim \chi_{n-p-1}^2$$

And that both of these are independent. Dividing a normal by a chi-squared gives us a student-t and so,

$$\frac{\sqrt{n - (p + 1)}(\hat{Y}_* - \mu_*)}{\|\mathbf{E}\| \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim T_{n-(p+1)}$$

- (bb) [easy] Let $H_0 : \mu_* = 17$. Find the test statistic using the fact from the previous question. Let s_e denote the *RMSE*.

$$\hat{t} = \frac{\hat{\mathbf{y}}_* - \mu_{*0}}{s_e \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}}$$

- (cc) [easy] Consider a new parameter of interest $y_* = \mathbf{x}_* \boldsymbol{\beta} + \epsilon_*$, this is the response for a unit with measurements given in row vector \mathbf{x}_* whose first entry is 1. Prove that

$$\frac{\hat{Y}_* - y_*}{\sigma \sqrt{1 + \mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim \mathcal{N}(0, 1).$$

$$\begin{aligned} \mathbf{Y}_* - \hat{\mathbf{Y}}_* &= \mathbf{Y}_* - \vec{x}_* \mathbf{B} = (\vec{x}_* \mathbf{B} + \vec{\mathcal{E}}_*) - (\vec{x}_* \mathbf{B}) \\ &\Rightarrow \vec{x}_* \mathbf{B} + \vec{\mathcal{E}}_* - \vec{x}_* \mathbf{B} \end{aligned}$$

Where $\vec{x}_* \mathbf{B} + \vec{\mathcal{E}}_* \sim \mathcal{N}(\vec{x}_* \mathbf{B}, \sigma^2)$

So,

$$\mathbf{Y}_* \sim \mathcal{N}(\vec{x}_* \mathbf{B} - \vec{x}_* \mathbf{B}, \sigma^2 + \sigma^2 \vec{x}_* ((\mathbf{X}^T \mathbf{X})^{-1}) \vec{x}_*^\top)$$

We can subtract the mean and divide by the standard deviation to get the standard normal:

$$\frac{\hat{Y}_* - y_*}{\sigma \sqrt{1 + \mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim \mathcal{N}(0, 1)$$

- (dd) [easy] Prove that $\frac{\sqrt{n - (p + 1)}(\hat{Y}_* - y_*)}{\|\mathbf{E}\| \sqrt{1 + \mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim T_{n-(p+1)}.$

We again have a standard normal divided by a chi squared with $n-(p+1)$ degrees of freedom.

- (ee) [easy] Let $H_0 : y_* = 37$. Find the test statistic using the fact from the previous question. Let s_e denote the *RMSE*.

$$\hat{t} = \frac{37 - \hat{\mathbf{y}}_*}{s_e \sqrt{1 + \vec{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \vec{x}_*^\top}}$$

- (ff) [difficult] Let $S \subseteq \{1, 2, \dots, p\}$, let $k := |S|$ and let $A = \{0\} \cup S^C$, its complement with zero for the index of the intercept. For convenience, assume you rearrange the columns of the design matrix so that $\mathbf{X} = [\mathbf{X}_A \mid \mathbf{X}_S]$ and the first column is $\mathbf{1}_n$. Let $\mathbf{H}_A := \mathbf{X}_A (\mathbf{X}_A^\top \mathbf{X}_A)^{-1} \mathbf{X}_A^\top$. It is obvious that $\mathbf{H} - \mathbf{H}_A$ is symmetric as both \mathbf{H} and \mathbf{H}_A are symmetric. To prove that $\mathbf{H} - \mathbf{H}_A$ is an orthogonal projection matrix, prove that it is idempotent. Hint: use the Gram-Schmidt decomposition for both matrices

and use block matrix format for \mathbf{H} .

We WTS $\mathbf{H} - \mathbf{H}_A$ is idempotent via the beginning steps of Gram-Schmidt decomposition.

Recall that \mathbf{X}_A is the matrix of the intercept column and the predictors not in the set S , while \mathbf{X}_S contains the predictors in the set S . Recall $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is the full projection matrix. We have $\mathbf{H}_A := \mathbf{X}_A (\mathbf{X}_A^T \mathbf{X}_A)^{-1} \mathbf{X}_A^T$.

Let us define $P_S := \mathbf{H} - \mathbf{H}_A$, the projection onto the space in $\text{colsp}[\mathbf{X}]$, but \perp to $\text{colsp}[\mathbf{X}_A]$ (It is the extra contribution of \mathbf{X}_S after adjusting for \mathbf{X}_A)

Then we can say

$$P_S = (\mathbf{I} - \mathbf{H}_A) \mathbf{X}_S [\mathbf{X}_S^T (\mathbf{I} - \mathbf{H}_A) \mathbf{X}_S]^{-1} \mathbf{X}_S^T (\mathbf{I} - \mathbf{H}_A)$$

Here, $(\mathbf{I} - \mathbf{H}_A)$ is the residual maker for the subspace $\text{colsp}[\mathbf{X}_A]$. We are saying take \mathbf{X}_S and remove any part that is in the span \mathbf{X}_A i.e. perform G-S on \mathbf{X}_S relative to \mathbf{X}_A

Let us define $\mathbf{X}_S^\perp := (\mathbf{I} - \mathbf{H}_A) \mathbf{X}_S$

$$\Rightarrow P_S = \mathbf{X}_S^\perp [(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp]^{-1} (\mathbf{X}_S^\perp)^\top$$

Square it:

$$P_S^2 = \mathbf{X}_S^\perp [(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp]^{-1} (\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp [(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp]^{-1} (\mathbf{X}_S^\perp)^\top$$

Notice $(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp [(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp]^{-1} = \mathbf{I}$ We now have:

$$P_S^2 = \mathbf{X}_S^\perp [(\mathbf{X}_S^\perp)^\top \mathbf{X}_S^\perp]^{-1} (\mathbf{X}_S^\perp)^\top$$

Therefore,

$$(\mathbf{H} - \mathbf{H}_A)^2 = \mathbf{H} - \mathbf{H}_A \quad \square$$

(gg) [easy] Let $\hat{\mathbf{Y}}_A := \mathbf{H}_A \mathbf{Y}$, the orthogonal projection onto $\text{colsp}[\mathbf{X}_A]$. Prove that

$$\frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2}{k \left\| \mathbf{E} \right\|^2} \sim F_{k, n-(p+1)}.$$

We know the rank of $\mathbf{H}_A = p + 1 - k$ and rank of $\mathbf{H} - \mathbf{H}_A$ is k , so because the ranks total to $p+1$ we can use Cochran's thm;

$$\frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top \boldsymbol{\varepsilon} = \frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top \mathbf{H}_A \boldsymbol{\varepsilon} + \frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top (\mathbf{H} - \mathbf{H}_A) \boldsymbol{\varepsilon} + \frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{H}) \boldsymbol{\varepsilon}$$

Where

$$\frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top \mathbf{H}_A \boldsymbol{\varepsilon} \sim \chi_{p+1-k}^2$$

independent of

$$\frac{1}{\sigma^2} \boldsymbol{\varepsilon}^\top (\mathbf{H} - \mathbf{H}_A) \boldsymbol{\varepsilon} \sim \chi_k^2$$

and independent of

$$\frac{1}{\sigma^2} \boldsymbol{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H}) \boldsymbol{\mathcal{E}} \sim \chi_{n-(p+1)}^2$$

By idempotency and symmetry

$$\Rightarrow \frac{1}{\sigma^2} \boldsymbol{\mathcal{E}}^\top (\mathbf{H} - \mathbf{H}_A) \boldsymbol{\mathcal{E}} = \frac{1}{\sigma^2} \|\mathbf{H} - \mathbf{H}_A \boldsymbol{\mathcal{E}}\|^2$$

Sub for $\boldsymbol{\mathcal{E}}$

$$\frac{1}{\sigma^2} \|(\mathbf{H} - \mathbf{H}_A)(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2$$

Distribute

$$\frac{1}{\sigma^2} \|\mathbf{H}\mathbf{Y} - \mathbf{H}_A\mathbf{Y} - \mathbf{H}\mathbf{X}\boldsymbol{\beta} + \mathbf{H}_A\mathbf{X}\boldsymbol{\beta}\|^2$$

We know $\mathbf{H}\mathbf{Y} = \hat{\mathbf{Y}}$ and $\mathbf{H}_A\mathbf{Y} = \hat{\mathbf{Y}}_A$

$$\Rightarrow \frac{1}{\sigma^2} \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A - \mathbf{X}\boldsymbol{\beta} + \mathbf{H}_A\mathbf{X}\boldsymbol{\beta} \right\|^2$$

Augment $\mathbf{X} = [\mathbf{X}_A \mid \mathbf{X}_S]$ and $\boldsymbol{\beta} = [\boldsymbol{\beta}_A \mid \boldsymbol{\beta}_S]^\top$. Now perform block multiplication

$$\begin{aligned} \Rightarrow \frac{1}{\sigma^2} \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A - \mathbf{X}_A\boldsymbol{\beta}_A - \mathbf{X}_S\boldsymbol{\beta}_S + \mathbf{X}_A\boldsymbol{\beta}_A + \mathbf{H}_A\mathbf{X}_S\boldsymbol{\beta}_S \right\|^2 \\ \Rightarrow \frac{1}{\sigma^2} \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A - (\mathbf{I} - \mathbf{H}_A)\mathbf{X}_S\boldsymbol{\beta}_S \right\|^2 \end{aligned}$$

Further,

$$\frac{\frac{1}{\sigma^2} \boldsymbol{\mathcal{E}}^\top (\mathbf{H} - \mathbf{H}_A) \boldsymbol{\mathcal{E}}}{\frac{1}{\sigma^2} \boldsymbol{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H}) \boldsymbol{\mathcal{E}}} = \frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2}{k \|\mathbf{E}\|^2} \sim F_{k, n-(p+1)}$$

Recall: this test tells us if at least one of all the predictors has a statistically significant effect.

(hh) [difficult] Let $\hat{\mathbf{E}}_A := (\mathbf{I}_n - \mathbf{H}_A)\mathbf{Y}$, the orthogonal projection onto the $\text{colsp}[\mathbf{X}_{A^\perp}]$. Prove that $\left\| \hat{\mathbf{E}}_A \right\|^2 - \left\| \hat{\mathbf{E}} \right\|^2 = \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2$.

$$\left\| \hat{\mathbf{E}}_A \right\|^2 = \left\| \hat{\mathbf{E}} + (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A) \right\|^2 = \hat{\mathbf{E}}^\top \hat{\mathbf{E}} + 2[\hat{\mathbf{E}}^\top (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A)] + (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A)^\top (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A)$$

$$\left\| \hat{\mathbf{E}}_A \right\|^2 - \left\| \hat{\mathbf{E}} \right\|^2 = \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2 + 2[\hat{\mathbf{E}}^\top (\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A)]$$

Therefore, the equality holds if the dot product of $\hat{\mathbf{E}}$ and $(\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A)$ is zero, which it is! This is because $(\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A) \in \text{colsp}[\mathbf{X}]$ and $\hat{\mathbf{E}} \in \text{colsp}[\mathbf{X}^\perp]$, making the two vectors orthogonal to each other and so their dot product is indeed zero. therefore,

$$\left\| \hat{\mathbf{E}}_A \right\|^2 - \left\| \hat{\mathbf{E}} \right\|^2 = \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2 \quad \square$$

- (ii) [easy] Combining the two previous problems, write the test statistic for $H_0 : \beta_S = \mathbf{0}_k$ where β_S denotes the subvector of β with indices S . Use the notation $\Delta SSE := SSE_A - SSE$ and MSE .

The test statistic for the partial F-Test is:

$$\frac{\frac{\Delta SSE}{k}}{\frac{SSE}{n-(p+1)}}$$

- (jj) [difficult] Prove that the square root of the test statistic in (ii) is the same as t-test statistic from (y) when $k = 1$.

$$\frac{\frac{\Delta SSE}{k}}{\frac{SSE}{n-(p+1)}} = \frac{SSE_A - SSE}{MSE}$$

Recall the T-test

$$\hat{t} = \frac{b_j}{s_e \sqrt{(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} \sim T_{n-(p+1)} \Rightarrow \frac{b_j^2}{s_e^2 (\mathbf{X}^T \mathbf{X})_{jj}^{-1}} \sim F_{1, n-(p+1)}$$

And s_e denotes $RMSE := \sqrt{MSE} := \sqrt{SSE/(n-(p+1))} = \sqrt{\|e\|^2/(n-(p+1))}$.
And $b_j^2 = \Delta SSE$

$$\begin{aligned} F &= \frac{\Delta SSE}{MSE} \\ &= \frac{SSE_A - SSE}{MSE} \\ &= \frac{B_j^2}{(\mathbf{X}^\top \mathbf{X})_{jj}^{-1} \cdot MSE} \end{aligned}$$

$$\begin{aligned} \sqrt{F} &= \sqrt{\frac{B_j^2}{(\mathbf{X}^\top \mathbf{X})_{jj}^{-1} \cdot MSE}} \\ &= \frac{|B_j|}{\sqrt{(\mathbf{X}^\top \mathbf{X})_{jj}^{-1} \cdot \sqrt{MSE}}} \\ &= \left| \frac{B_j}{SE(B_j)} \right| \\ &= |t| \end{aligned}$$

- (kk) [harder] The point of this exercise is to demonstrate that the estimator used for the omnibus / global / overall F-test is nothing but a special case of the main result from (gg). Let $S = \{1, 2, \dots, p\}$ and thus $k = p$ and $A = \{0\}$. Using the result from (gg),

show that
$$\frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \bar{y} \mathbf{1}_n \right\|^2}{p \left\| \mathbf{E} \right\|^2} \sim F_{p, n-(p+1)}.$$

Let $\mathbf{X}_A = \mathbf{1}$. Find the fitted values under the reduced model:

$$\mathbf{H}_A = \mathbf{X}_A (\mathbf{X}_A^\top \mathbf{X}_A)^{-1} \mathbf{X}_A^\top$$

Substitute $\mathbf{X}_A = \mathbf{1}_n$:

$$\mathbf{H}_A = \mathbf{1}_n (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} \mathbf{1}_n^\top$$

Since:

$$\mathbf{1}_n^\top \mathbf{1}_n = \sum_{i=1}^n 1 = n, \quad \text{so} \quad (\mathbf{1}_n^\top \mathbf{1}_n)^{-1} = \frac{1}{n}$$

Thus:

$$\mathbf{H}_A = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top$$

Apply \mathbf{H}_A to \mathbf{Y} :

$$\hat{\mathbf{Y}}_A = \mathbf{H}_A \mathbf{Y} = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top \mathbf{Y}$$

$$\mathbf{1}_n^\top \mathbf{Y} = \sum_{i=1}^n y_i = n\bar{y}$$

$$\hat{\mathbf{Y}}_A = \frac{1}{n} \mathbf{1}_n (n\bar{y}) = \bar{y} \cdot \mathbf{1}_n$$

Therefore, the fitted values under the reduced model are $\hat{\mathbf{Y}}_A = \mathbf{H}_A \mathbf{Y} = \bar{y} \mathbf{1}$

Plug in:

$$\frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A \right\|^2}{k \left\| \mathbf{E} \right\|^2} \sim F_{k, n-(p+1)} \Rightarrow \frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \bar{y} \mathbf{1}_n \right\|^2}{p \left\| \mathbf{Y} - \hat{\mathbf{Y}} \right\|^2} \sim F_{p, n-(p+1)}$$

- (ll) [easy] Prove that the omnibus / global / overall F-test statistic is $\hat{\hat{F}} = MSR/MSE$ by using the result from (kk).

We know:

$$\frac{(n - (p + 1)) \left\| \hat{\mathbf{Y}} - \bar{y} \mathbf{1}_n \right\|^2}{p \left\| \mathbf{Y} - \hat{\mathbf{Y}} \right\|^2} \sim F_{p, n-(p+1)}$$

We can see that $\|\hat{\mathbf{Y}} - \bar{y}\mathbf{1}_n\|^2$ is the sum of square residuals due to regression:

$$SSR := \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \|\hat{\mathbf{Y}} - \bar{y}\mathbf{1}_n\|^2$$

Therefore,

$$\frac{\|\hat{\mathbf{Y}} - \bar{y}\mathbf{1}_n\|^2}{p} = \frac{SSR}{p} := MSR$$

And we saw earlier that the denominator is the MSE:

$$\frac{\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2}{n - (p + 1)} = \frac{SSE}{n - (p + 1)} := MSE$$

So we have,

$$\hat{F} = \frac{MSR}{MSE} \sim F_{p, n-(p+1)}$$

(mm) [difficult] [MA] Prove that the distribution that realizes the R^2 metric (the proportion of response variance explained by the model) is distributed as Beta $\left(\frac{p}{2}, \frac{n-(p+1)}{2}\right)$. This amounts to proving a fact found at the bottom of the F distribution's Wikipedia page .

(nn) [easy] Prove that the maximum likelihood estimate for β is \mathbf{b} , the OLS estimator.

From 342, we know that $\vec{b} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ which came from $\arg \min \{ \|\mathbf{y} - \mathbf{X}\vec{w}\|^2 \}$. In 343, $\vec{\hat{\theta}} = \vec{b}$ drawn from $\vec{\theta} = \mathbf{B}$. Recall 341's use of the method of moments and the maximum likelihood estimate to find estimators. Let us attempt this here; derive $\vec{\hat{\theta}}^{\text{MLE}}$ under core assumption:

$$\mathcal{E} \sim \mathcal{N}_n(\vec{0}_n, \sigma^2 \mathbf{I}_n) \Rightarrow \mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$$

Use pdf of this and find the likelihood:

$$\Rightarrow \frac{1}{(2\pi)^{n/2} (\sigma^2)^{n/2}} e^{-\frac{1}{2} [(\mathbf{Y} - \mathbf{X}\beta)^\top (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{Y} - \mathbf{X}\beta)]}$$

Notice, $-\frac{1}{2} [(\mathbf{Y} - \mathbf{X}\beta)^\top (\sigma^2 \mathbf{I}_n)^{-1} (\mathbf{Y} - \mathbf{X}\beta)] = -\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2$

$$\Rightarrow \ell(\beta, \sigma^2; \mathbf{y}; \mathbf{X}) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) \left[-\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{X}\beta\|^2 \right]$$

Differentiate with respect to β

$$\Rightarrow -\frac{1}{2\sigma^2} \frac{\partial}{\partial \beta} \left[\sum_{i=1}^n (y_i - \mathbf{X}_i \beta)^2 \right] \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \vec{\hat{\beta}} = \vec{b} \text{ from OLS} \quad \square$$

(oo) [harder] Prove that the maximum likelihood estimate for σ^2 is SSE/n .

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_n) \Rightarrow \mathbf{Y} \sim \mathcal{N}_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$$

$$\ell(\boldsymbol{\beta}, \sigma^2 \mid \mathbf{Y}) = \frac{1}{(2\pi\sigma^2)^{n/2}} e^{(-\frac{1}{2\sigma^2} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2)}$$

Let $\mathbf{b} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ be the OLS estimator of $\boldsymbol{\beta}$, and define:

$$SSE := \|\mathbf{Y} - \mathbf{X}\mathbf{b}\|^2$$

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} \ln(\ell) &= -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} SSE \\ -\frac{n}{2\sigma^2} + \frac{SSE}{2(\sigma^2)^2} &= 0 \Rightarrow -n\sigma^2 + SSE \stackrel{\text{set}}{=} 0 \\ \Rightarrow \hat{\sigma}_{\text{MLE}}^2 &= \frac{SSE}{n} \end{aligned}$$

(pp) [harder] Find the bias of the maximum likelihood estimator for σ^2 using your answers from (w) and (oo).

Recall the bias of an estimator $\hat{\theta}$ for a parameter θ is:

$$\text{Bias}[\theta] = \mathbb{E}[\hat{\theta}] - \theta$$

In this case it is:

$$\text{Bias}[\sigma_{\text{MLE}}^2] = \mathbb{E}[\sigma_{\text{MLE}}^2] - \sigma^2$$

Recall from the previous problem,

$$\hat{\sigma}_{\text{MLE}}^2 = \frac{SSE}{n}$$

Therefor,

$$\mathbb{E}[\sigma_{\text{MLE}}^2] = \mathbb{E}\left[\frac{SSE}{n}\right] = \frac{1}{n} \mathbb{E}[SSE]$$

But we know

$$\frac{1}{n} \mathbb{E}[SSE] = \frac{n - (p + 1)}{n} \sigma^2$$

So,

$$\begin{aligned} \text{Bias}[\sigma_{\text{MLE}}^2] &= \frac{n - (p + 1)}{n} \sigma^2 - \sigma^2 \\ &\Rightarrow \frac{n\sigma^2 - p\sigma^2 - \sigma^2}{n} - \frac{n\sigma^2}{n} \\ &\Rightarrow \text{Bias}[\sigma_{\text{MLE}}^2] = -\frac{(p + 1)\sigma^2}{n} \end{aligned}$$

The MLE does not account for the fact that we estimated $p+1$ parameters from the data. Therefore, it systematically underestimates the variance and the unbiased estimator corrects this by dividing $n - (p+1)$ instead of just n .

Problem 2

This problem is about two types of Bayesian estimation of the slope parameters in linear regression which lead to the ridge and lasso estimates.

- (a) [easy] Write the prior assumption about β which yields the ridge estimates.

$$f(\vec{\beta}) = \mathcal{N}_{p+1} \left(\vec{0}_{p+1}, \tau^2 \mathbf{I}_{p+1} \right)$$

$$f(\sigma^2) = \frac{1}{\sigma^2}$$

- (b) [easy] Using the prior and core assumption (which implies a likelihood function for \mathbf{B}), derive the ridge estimates.

$$f(\vec{\beta}; \sigma^2 \mid \vec{y}, X) \propto f(\vec{y} \mid \vec{\beta}, \sigma^2, X) f(\vec{\beta}) f(\sigma^2)$$

$$\propto (\sigma^2)^{-\frac{n}{2}-1} e^{-1/2\sigma^2 \|\vec{y} - \mathbf{X}\vec{\beta}\|^2} e^{-1/2\tau^2 \|\vec{\beta}\|^2}$$

This is not a known distribution. But if we want just point estimation, we don't need the distribution, we can simply just use the MAP:

$$\hat{\vec{\beta}}^{\text{MAP}} = \arg \max \left\{ (-n/2 - 1) \ln(\sigma^2) - \frac{1}{2\sigma^2} \|\vec{y} - \mathbf{X}\vec{\beta}\|^2 - \frac{1}{2\tau^2} \|\vec{\beta}\|^2 \right\}$$

End up with:

$$\Rightarrow \arg \min \left\{ SSE + \lambda \|\vec{\beta}\|^2 \right\}$$

Where $\lambda = \frac{\sigma^2}{\tau^2}$

- (c) [easy] Write the prior assumption about β which yields the lasso estimates.

We use a Laplace prior:

$$f(\vec{\beta}) = \frac{1}{2\tau^2} e^{-\frac{|\beta|}{\tau^2}}$$

And Jeffries for variance:

$$f(\sigma^2) = \frac{1}{\sigma^2}$$

- (d) [easy] Using the prior and core assumption (which implies a likelihood function for \mathbf{B}), derive the lasso estimates to the point where you need to use a computer to run the optimization.

$$\begin{aligned} f(\vec{\beta}; \sigma^2 \mid \vec{y}, X) &\propto f(\vec{y} \mid \vec{\beta}, \sigma^2, X) f(\vec{\beta}) f(\sigma^2) \\ &\propto (\sigma^2)^{-\frac{n}{2}-1} e^{-1/2\sigma^2 \|\vec{y} - \mathbf{X}\vec{\beta}\|^2} e^{-1/2\tau^2 \sum_{j=0}^p |\beta_j|} \end{aligned}$$

Again, we do not know this distribution, so we will find the MAP estimate for point estimation:

$$\hat{\vec{\beta}}^{\text{MAP}} = \arg \max \left\{ (-n/2 - 1) \ln(\sigma^2) - \frac{1}{2\sigma^2} \|\vec{y} - \mathbf{X}\vec{\beta}\|^2 - \frac{1}{2\tau^2} \sum_{j=0}^p |\beta_j| \right\}$$

Which becomes:

$$\arg \min \left\{ SSE + \lambda \sum_{j=0}^p |\beta_j| \right\}$$

Where $\lambda = \frac{2\sigma^2}{\tau^2}$. There is no closed form solution and we need an optimizer.

- (e) [easy] Both ridge and lasso shrink the estimate of β towards what vector value?

$$\vec{0}_{p+1}$$

In both cases, ridge and lasso, the prior mean is zero so the MAP estimate is pulled toward zero. When $\lambda = 0$ it becomes the OLS estimate and when λ goes to infinity the $\arg \min \{\|\cdot\|^2\} \vec{\beta} \rightarrow \vec{0}_{p+1}$ in Ridge regression.

Recall, the shrinkage estimator is used to pull the estimated parameters toward a target value to reduce variance at the risk of introducing bias, which can be useful if (a) p is very large or (b) data is noisy. We did this initially because we want to regress $p = 20,000$ genes but we only have $n = 300$ people.

Ridge regression shrinks all coefficients towards zero, but keeps them nonzero - does not perform variable selection i.e. keeps all predictors.

Lasso regression shrinks some coefficients exactly to zero, it does both shrinkage and variable selection which can be helpful when you believe only a small number of predictors matter.

- (f) [easy] Describe what the prestep called “variable selection” is within the modeling enterprise.

Variable selection is the process of deciding which predictors to include in the model before fitting the model. This is important because if there is a large number of predictors, they may not all be important and the inclusion of unnecessary variables can create overfitting in your model.

- (g) [easy] Describe why Lasso estimation has the added bonus of being able to perform variable selection and ridge does

Lasso has an ℓ_1 penalty that is not differentiable at zero, so this encourages exactly zero coefficients and will force some of the weights for certain features to zero. Ridge has an ℓ_2 penalty which is smooth and differentiable everywhere so it will shrink all the coefficients continuously towards zero, but they will never equal zero unless $\lambda = \infty$ (which is impossible).

Problem 3

This problem is about the specific robust regression methods we studied.

- (a) [easy] If we only know that the errors $\mathcal{E}_1, \dots, \mathcal{E}_n$ are independent, what tried and true method can we employ to get asymptotically valid inference for β ?

Perform bootstrap! Sample n rows from \mathbb{D} with replacement to get the bootstrap sample (a bunch of tuples), and the find their OLD estimate. Then, for the j -th estimate, plot the bootstrap distribution.

Ask, does $0 \in [Q[b_{b_j}, \frac{\alpha}{2}], Q[b_{b_j}, 1 - \frac{\alpha}{2}]]$

- (b) [easy] If we know that the errors $\mathcal{E}_1, \dots, \mathcal{E}_n$ are iid with expectation zero and variance σ^2 for all values of \mathbf{x} (i.e. the errors are “homoskedastic”) but the errors are not necessarily normal, what is the asymptotic distribution of \mathbf{B} ?

$$\mathbf{B} \dot{\sim} \mathcal{N}_{p+1} \left(\vec{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- (c) [easy] If we know that the errors $\mathcal{E}_1, \dots, \mathcal{E}_n \stackrel{ind}{\sim} \mathcal{N}(0, \sigma_i^2)$ which means the errors are “heteroskedastic”, what is the asymptotic distribution of \mathbf{B} using the Huber-White estimator?

$$\mathbf{B} \dot{\sim} \mathcal{N}_{p+1} \left(\vec{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

Where residuals are in place of variance and

$$\hat{D} = \begin{bmatrix} e_1^2 & & & \\ & e_2^2 & & \\ & & \ddots & \\ & & & e_n^2 \end{bmatrix}$$

- (d) [easy] If we know that the errors $\mathcal{E}_1, \dots, \mathcal{E}_n$ are independent with expectation zero and variance σ_i^2 which means the errors are “heteroskedastic”, what is the asymptotic distribution of \mathbf{B} using the Huber-White estimator?

same as previous, except even more approximate

$$\mathbf{B} \sim \mathcal{N}_{p+1} \left(\vec{\beta}, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{D} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

Where residuals are in place of variance and

$$\hat{D} = \begin{bmatrix} e_1^2 & & & \\ & e_2^2 & & \\ & & \ddots & \\ & & & e_n^2 \end{bmatrix}$$

- (e) [easy] Is the F-tests we derived under the core assumption valid in any of the four above scenarios? Yes/no

No, because the F-tests derived from the core assumption rely on homoscedacity, normality of error, and IID

Problem 4

This problem is about inference for the generalized linear model (glm).

- (a) [harder] Let $Y_i \stackrel{ind}{\sim} \text{Bernoulli}(\theta_i)$ for $i = 1, \dots, n$ where $\theta_i = \phi(\mathbf{x}_i \boldsymbol{\beta})$ and $\mathbf{x}_i \in \mathbb{R}^{p+1}$ whose first entry is always 1. For the link function, use the complementary log-log (i.e. the standard Gumbel CDF). Write out the full likelihood below. No need to simplify or take the log.

$$\prod_{i=1}^n (e^{-e^{-x_i \beta}})^{Y_i} (1 - e^{-e^{-x_i \beta}})^{1-Y_i}$$

- (b) [harder] Given the assumptions in (a), write the likelihood ratio estimate for the omnibus test of $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$.

$$\hat{\text{LR}} = \frac{\prod_{i=1}^n (e^{-e^{-x_i \beta}})^{Y_i} (1 - e^{-e^{-x_i \beta}})^{1-Y_i}}{\prod_{i=1}^n (e^{-e^{-x_i \tilde{b}_0}})^{Y_i} (1 - e^{-e^{-x_i \tilde{b}_0}})^{1-Y_i}}$$

- (c) [harder] Let $Y_i \stackrel{ind}{\sim} \text{Poisson}(\theta_i)$ for $i = 1, \dots, n$ where $\theta_i = e^{x_i \boldsymbol{\beta}}$ and $\mathbf{x}_i \in \mathbb{R}^{p+1}$ whose first entry is always 1. Write out the likelihood ratio when testing $H_0 : \beta_2 = \beta_3 = 0$.

$$\hat{\text{LR}} = \frac{\prod_{i=1}^n e^{-e^{-x_i \beta}} e^{x_i \beta Y_i}}{\prod_{i=1}^n e^{-e^{-x_i \tilde{b}}} e^{x_i \tilde{b} Y_i}}$$

Where $\tilde{b} = \arg \max \{ \sum_{i=1}^n -e^{\tilde{x} \tilde{w}} + \tilde{x} \tilde{w} y_i \mid w_2 = w_3 = 0 \}$

- (d) [harder] Let $Y_i \stackrel{ind}{\sim} \text{Weibull}(k, \theta_i)$ for $i = 1, \dots, n$ where $\theta_i = e^{\mathbf{x}_i \boldsymbol{\beta}}$ and $\mathbf{x}_i \in \mathbb{R}^{p+1}$ whose first entry is always 1. This uses the alternate parameterization so that $\mathbb{E}[Y_i] = \theta_i \Gamma(1 + 1/k)$. There is a censoring vector \mathbf{c} which is 1 when censored on the right (meaning the real y_i is \geq to the observed y_i) and 0 when not censored. Write out the likelihood ratio when testing $H_0 : \beta_2 = \beta_3 = 0$.

$$\prod_{i:c_i=1} \frac{s}{e^{\vec{x}_i \vec{w}}} \left(\frac{y_i}{e^{\vec{x}_i \vec{w}}} \right)^{s-1} e^{-\left(\frac{y_i}{e^{\vec{x}_i \vec{w}}} \right)^s} \prod_{i:c_i=0} e^{-\left(\frac{y_i}{e^{\vec{x}_i \vec{w}}} \right)^s}$$

- (e) [difficult] [MA] Let $Y_i \stackrel{ind}{\sim} \mathcal{N}(\theta_i, \sigma^2)$ for $i = 1, \dots, n$ where $\theta_i = \mathbf{x}_i \boldsymbol{\beta}$ and $\mathbf{x}_i \in \mathbb{R}^{p+1}$ whose first entry is always 1. So far, this is the vanilla linear model. However, consider now a wrinkle: there is a censoring vector \mathbf{c} which is 1 when censored on the right (meaning the real y_i is \geq to the observed y_i) and 0 when not censored. This is called the Tobit model. Write the likelihood ratio estimate for the omnibus test of $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$.