

### Problem 1

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

$$\begin{aligned}\mathbf{X} &:= [\mathbf{1}_n \mid \mathbf{x}_{\cdot 1} \mid \dots \mid \mathbf{x}_{\cdot p}] \text{ with column indices } 0, 1, \dots, p \text{ and row indices } 1, 2, \dots, n \\ \mathbf{H} &:= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\ \mathbf{Y} &= \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \\ \mathbf{B} &:= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ \hat{\mathbf{Y}} &= \mathbf{HY} = \mathbf{XB} \\ \mathbf{E} &:= \mathbf{Y} - \hat{\mathbf{Y}} = (\mathbf{I}_n - \mathbf{H}) \mathbf{Y}\end{aligned}$$

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?

$$\boldsymbol{\beta}$$

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

$$\mathbf{X}, \mathbf{H}$$

- (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?

$$\mathbf{Y}, \mathbf{y} \quad \text{and} \quad \boldsymbol{\beta}, \hat{\boldsymbol{\beta}}$$

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

mean centered, homoscedasticity, and normal dist

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

$$\begin{aligned}\tilde{\mathbf{B}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tilde{\mathbf{Y}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\varepsilon} \\ &\sim N_{p+1}(\bar{\boldsymbol{\beta}}, 6 \cdot (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-1})^T) \\ &\sim N_{p+1}(\bar{\boldsymbol{\beta}}, 6 \cdot (\mathbf{X}^T \mathbf{X})^{-1})\end{aligned}$$

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

$$\beta_j \sim N(\beta_j, 6 \cdot (\mathbf{X}^T \mathbf{X})^{-1}_{j,j})$$

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

$$\frac{\beta_j - \bar{\beta}}{\sqrt{6 \cdot (\mathbf{X}^T \mathbf{X})^{-1}_{j,j}}} \sim N(0, 1)$$

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

$$\tilde{Y} = X\tilde{\beta} + H\tilde{\epsilon} = H(X\beta + \tilde{\epsilon}) = X\beta + H\tilde{\epsilon} \sim N_n(X\beta, H\sigma^2 I_n H)$$

$$= N_n(X\beta, \sigma^2 H)$$

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

$$\hat{Y}_i \sim N(\hat{x}_i \hat{\beta}, \sigma^2 H_{ii})$$

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

$$\frac{\hat{y}_i - \hat{x}_i \hat{\beta}}{\sqrt{\sigma^2 H_{ii}}}$$

(k) [easy] from the core assumption, derive the distribution of  $E_i$ .

$$\tilde{E}_i = \tilde{Y}_i - X\tilde{\beta} = (I - H)\tilde{Y}_i = (I - H)(X\beta + \tilde{\epsilon}) = (I - H)\tilde{\epsilon} \sim N_n(0_n, (I - H)\sigma^2 I_n (I - H))$$

$$= N_n(0_n, \sigma^2 (I - H))$$

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

$$e_i \sim N(0, \sigma^2 (I - H_{ii}))$$

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

$$\frac{e_i}{\sqrt{\sigma^2 (I - H_{ii})}}$$

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

$$\begin{aligned} z \cdot z &\sim \chi_n^2 \\ \Rightarrow \left( \frac{1}{\sigma} e \right) \left( \frac{1}{\sigma} e \right)^\top &= \frac{1}{\sigma^2} e^\top e \sim \chi_n^2 \end{aligned}$$

- (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  $\frac{1}{\sigma^2} \mathbf{\mathcal{E}}^\top \mathbf{B}_1 \mathbf{\mathcal{E}}$  and  $\frac{1}{\sigma^2} \mathbf{\mathcal{E}}^\top \mathbf{B}_2 \mathbf{\mathcal{E}}$ .

$$\frac{1}{\sigma^2} \mathbf{\mathcal{E}}^\top \mathbf{\mathcal{E}} = \underbrace{\frac{1}{6} \mathbf{\mathcal{E}}^\top \mathbf{H} \mathbf{\mathcal{E}}}_{\chi^2_{p+1}} + \underbrace{\frac{1}{6} \mathbf{\mathcal{E}}^\top (\mathbf{I} - \mathbf{H}) \mathbf{\mathcal{E}}}_{\chi^2_{n-(p+1)}} \quad \begin{matrix} \text{rank}(\mathbf{H}) = p+1 \\ \text{rank}(\mathbf{I} - \mathbf{H}) = n - (p+1) \end{matrix}$$

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

$$\begin{aligned} \frac{1}{\sigma^2} \mathbf{\mathcal{E}}^\top \mathbf{H} \mathbf{\mathcal{E}} &= \frac{1}{\sigma^2} \mathbf{\mathcal{E}}^\top \mathbf{H} \mathbf{H} \mathbf{\mathcal{E}} = \frac{1}{\sigma^2} \|\mathbf{H} \mathbf{\mathcal{E}}\|^2 = \frac{1}{\sigma^2} \|\mathbf{H}(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2 \\ &= \frac{1}{\sigma^2} \|\mathbf{H}\mathbf{Y} - \mathbf{H}\mathbf{X}\boldsymbol{\beta}\|^2 = \frac{1}{\sigma^2} \|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2 \end{aligned}$$

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

$$\|\mathbf{X}(\mathbf{B} - \boldsymbol{\beta})\|^2 = \underbrace{\|\mathbf{X}\mathbf{B} - \mathbf{X}\boldsymbol{\beta}\|}_{\text{Definition of misspecification}}^2$$

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

$$\begin{aligned} \frac{1}{6} \mathbf{\mathcal{E}}^\top (\mathbf{I}_n - \mathbf{H}) \mathbf{\mathcal{E}} &= \frac{1}{6} \mathbf{\mathcal{E}}^\top (\mathbf{I}_n - \mathbf{H})^\top (\mathbf{I}_n - \mathbf{H}) \mathbf{\mathcal{E}} = \frac{1}{6} \|\mathbf{f}_{n-H} \mathbf{\mathcal{E}}\|^2 \\ &= \frac{1}{6} \|\mathbf{(I}_n - \mathbf{H})(\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})\|^2 = \underbrace{\|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{H}\mathbf{X}\boldsymbol{\beta}\|_2^2} \\ &= \frac{1}{6} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{H}\mathbf{Y} + \mathbf{H}\mathbf{X}\boldsymbol{\beta}\|^2 = \frac{1}{6} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \hat{\mathbf{Y}} + \mathbf{X}\boldsymbol{\beta}\|^2 \\ &= \frac{1}{6} \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 \\ &= \frac{1}{6} \|\mathbf{E}\|^2 \end{aligned}$$

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

$$\mathbf{E}^\top (\mathbf{H}^\top \mathbf{E}) > \mathbf{E}^\top (\mathbf{I} - \mathbf{H}) \mathbf{E}$$

when  $p < n$

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

$$H\mathbf{E} = \|\mathbf{X}(\mathbf{B} - \mathbf{P})\|^2$$
$$\|\mathbf{E} - H(\mathbf{X}(\mathbf{B} - \mathbf{P}))\|^2$$

$$\begin{aligned} H\mathbf{E} &\rightarrow \cancel{\mathbf{E}} \cancel{(\cancel{\mathbf{Y}} - \mathbf{X}\mathbf{P})} \\ &= \mathbf{X}(\mathbf{B} - \mathbf{P}) \end{aligned}$$

[easier] A  $2 \times 1$  linear model has a large or a small projection of the error? Discuss.

Small because ~~over the left the error and proj is small~~  
~~if it were large it would mean in the right part of the error~~  
[easy] Find  $\mathbf{E}^\top \mathbf{E}$  if the colsp of the feature?

$$\mathbb{E}(\frac{1}{n} \|\mathbf{E}\|^2) = \sigma^2(n)$$

[easy] Show that  $\frac{\mathbf{E}^\top \mathbf{E}}{n(n-1)}$  is an unbiased estimate of  $\sigma^2$ .

$$\mathbb{E}\left(\frac{\|\mathbf{E}\|^2}{n(n-1)}\right) = \sigma^2$$

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

$$\frac{\beta_j - \hat{\beta}_j}{\sqrt{(\mathbf{X}^T \mathbf{X})_{j,j}}} \sim N(0, 1)$$

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

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

(y) [easy] Let  $H_0 : \beta_j = 0$ . Find the test statistic using the fact from the previous question.

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

$$\frac{\beta_j - \hat{\beta}_j}{s_e \sqrt{(\mathbf{X}^T \mathbf{X})_{j,j}}}$$

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

$$\hat{Y}_* = \mathbf{x}_* \hat{\boldsymbol{\beta}} \Rightarrow \hat{Y}_* \sim N\left(\underbrace{\mathbf{x}_* \hat{\boldsymbol{\beta}}}_{\mu_*}, \sigma^2 \mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top\right)$$

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

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

$$\begin{aligned} & \frac{\hat{Y}_* - \mu_*}{\sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \sim N(0, 1) \sim T_{n-(p+1)} \\ & \sqrt{\frac{\mathbf{x}_*^\top (\mathbf{B}^\top \mathbf{T}) \mathbf{x}_*}{n-(p+1)}} \sim X_{n-(p+1)}^1 \\ & = \frac{n-(p+1)(\hat{Y}_* - \mu_*)}{\|E\| \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}} \end{aligned}$$

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

$$\frac{\hat{Y}_* - 17}{s_e \sqrt{\mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top}}$$

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

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

$$\begin{aligned} \hat{Y}_* - \hat{Y}_r &= N(0, 6^2 + 6 \cdot \mathbf{x}_* (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*^\top) = N(0, 6^2 (1 + \mathbf{x}_*^\top (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*)) \\ & \frac{\hat{Y}_r - \hat{Y}_*}{\sqrt{6^2 (1 + \mathbf{x}_*^\top (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_*))}} \sim N(0, 1) \end{aligned}$$

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

$$\frac{\sqrt{n-(p+1)}(\hat{Y}_* - y_*)}{\sqrt{\frac{\|\mathbf{E}\|^2}{n-(p+1)}} \sqrt{1 + \mathbf{x}_* (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*^\top}} \stackrel{\mathcal{N}(0, 1)}{\sim} \frac{\sqrt{n-(p+1)}(\hat{Y}_* - \bar{y}_*)}{\|\mathbf{E}\| \sqrt{1 + \mathbf{x}_* (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*^\top}} \stackrel{\sim T_{n-(p+1)}}{\sim}$$

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

$$\frac{37 - \hat{Y}_*}{s_e \sqrt{1 + \mathbf{x}_* (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_*^\top}}$$

(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  $\mathbf{H}$ .

$$\mathbf{X}_A = \mathbf{Q}_A \mathbf{R}_A, \mathbf{X} = \mathbf{Q} \mathbf{R}$$

$$\begin{aligned} (\mathbf{H} - \mathbf{H}_A)(\mathbf{H} - \mathbf{H}_A)^\top &= (\mathbf{Q} \mathbf{Q}^\top - \mathbf{Q}_A \mathbf{Q}_A^\top)(\mathbf{Q} \mathbf{Q}^\top - \mathbf{Q}_A \mathbf{Q}_A^\top)^\top \left( (\mathbf{Q}_A \mathbf{Q}_{-A}^\top) \begin{pmatrix} \mathbf{Q}_A^\top \\ \mathbf{Q}_{-A}^\top \end{pmatrix} - \mathbf{Q}_A \mathbf{Q}_{-A}^\top \right)^\top \\ &\quad \left( (\mathbf{Q}_A \mathbf{Q}_{-A}^\top)^\top \mathbf{Q}_A \mathbf{Q}_{-A}^\top + \mathbf{Q}_A \mathbf{Q}_{-A}^\top - \mathbf{Q}_A \mathbf{Q}_{-A}^\top \right) \\ &= (\mathbf{Q}_{-A} \mathbf{Q}_{-A}^\top)^\top \left( \mathbf{Q}_{-A} \mathbf{Q}_{-A}^\top \right) = (\mathbf{Q}_{-A} \mathbf{Q}_{-A}^\top) = (\mathbf{H} - \mathbf{H}_A) \end{aligned}$$

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

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

$$\begin{aligned} & \frac{1}{6} \cdot \vec{\epsilon}^\top (\mathbf{H} - \mathbf{H}_A \mathbf{H}_A^\top) \vec{\epsilon} \sim \chi^2_k \\ & = \frac{1}{6} \|\vec{\mathbf{Y}} - \vec{\mathbf{Y}}_A + (\mathbf{I} - \mathbf{H}_A) \vec{\mathbf{Y}}, \vec{\beta}_S\|^2 \\ & \Rightarrow \frac{1}{6} \left\| \frac{\vec{\mathbf{Y}} - \vec{\mathbf{Y}}_A}{k} \right\|^2 = \frac{(n - (p+1)) \|\vec{\mathbf{Y}} - \vec{\mathbf{Y}}_A\|^2}{k \|\vec{\mathbf{B}}\|^2} \sim F_{k, n-(p+1)} \end{aligned}$$

under  
H

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

$$\text{Prove that } \|\hat{\mathbf{E}}_A\|^2 - \|\hat{\mathbf{E}}\|^2 = \|\hat{\mathbf{Y}} - \hat{\mathbf{Y}}_A\|^2.$$

$$\begin{aligned} \|(I_n - H_A)Y\|^2 - \|Y - \hat{Y}\|^2 &= \|Y - \hat{Y}_A\|^2 - \|Y - \hat{Y}\|^2 \\ &= (Y - \hat{Y}_A)^\top (Y - \hat{Y}_A) - (Y - \hat{Y})^\top (Y - \hat{Y}) \\ &= Y^\top Y - Y^\top \hat{Y}_A - \hat{Y}_A^\top Y + \hat{Y}_A^\top \hat{Y}_A - Y^\top Y - \hat{Y}^\top Y + \hat{Y}^\top Y - \hat{Y}^\top \hat{Y} \\ &= (\hat{Y} - \hat{Y}_A)^\top (\hat{Y} - \hat{Y}_A) = \|\hat{Y} - \hat{Y}_A\|^2 \end{aligned}$$

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

$$\hat{F} = \frac{\underline{SSE_A} - SSE}{\frac{k}{MSE}}$$

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

$$\sqrt{\frac{SSE_A - SSR}{MSB}} = \sqrt{\frac{SSE_A - SSR}{MSE}} = \sqrt{\frac{SSE_A - SSR}{Se}} = \frac{\sqrt{SSE_A - SSR}}{Se} = \frac{\sqrt{\|Y - \hat{Y}\|^2}}{Se}$$

$$= \frac{\sqrt{\|Y - \bar{Y}_A\|^2}}{Se} = \frac{\sqrt{\|X\beta - X^b\|^2}}{Se} = \frac{\sqrt{\|X^b - X_A b\|^2}}{Se}$$

$$= \frac{\sqrt{\Delta SSE}}{Se} = \frac{b_j}{\sqrt{(X^T X)^{-1}}} = \frac{b_j}{Se \sqrt{(X^T X)^{-1}}}$$

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

$$\frac{(n - (p+1)) \|\hat{Y} - \bar{Y}_A\|^2}{p \|\mathbf{E}\|^2} = \frac{(n - (p+1)) \|\hat{Y} - \bar{Y}\|^2}{p \|\mathbf{E}\|^2}$$

$\hat{Y}_A = \bar{Y} \mathbf{1}_n$  under null

$$\hat{Y}_A = X\beta$$

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

$$\frac{\|\hat{Y} - \bar{Y}\|^2}{\frac{\|\mathbf{E}\|^2}{n-(p+1)}} = \frac{\frac{SSR}{p}}{\frac{SSB}{n-(p+1)}} = \frac{MSR}{MSE}$$

- (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  $\beta$  is  $\hat{\beta}$ , the OLS estimator.

$$\begin{aligned} \vec{y} &\sim N_n(\vec{x}\vec{\beta}, \sigma^2 I_n) \\ \ell(\vec{\beta}, \sigma^2; \vec{y}, \vec{x}) &= \frac{1}{(2\pi)^n \sqrt{\det(\sigma^2 I_n)}} e^{-\frac{1}{2}(\vec{y} - \vec{x}\vec{\beta})^T (\sigma^2 I_n)^{-1} (\vec{y} - \vec{x}\vec{\beta})} \\ \ell(\vec{\beta}, \sigma^2; \vec{y}, \vec{x}) &= -\frac{n}{2} \ln(\sigma^2) - \frac{n}{2} \ln(r^2) - \frac{1}{2\sigma^2} \|\vec{y} - \vec{x}\vec{\beta}\|^2 \\ \frac{\partial}{\partial \vec{\beta}} [\ell] &= -\frac{1}{2\sigma^2} \frac{\partial}{\partial \vec{\beta}} [\|\vec{y} - \vec{x}\vec{\beta}\|^2] \stackrel{\text{set}}{=} 0 \\ \frac{\partial}{\partial \vec{\beta}} &\Rightarrow \hat{\beta}^{\text{MLE}} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y} = \vec{b} \end{aligned}$$

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

$$\begin{aligned} \frac{\partial}{\partial \sigma^2} [\ell] &= -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)} \|\vec{y} - \vec{x}\vec{\beta}\|^2 \stackrel{\text{set}}{=} 0 \Rightarrow -n + \frac{1}{\sigma^2} \|\vec{y} - \vec{x}\vec{\beta}\|^2 = 0 \\ \hat{\sigma}^2 &= \frac{1}{n} \|\vec{y} - \vec{x}\vec{\beta}\|^2 = \frac{1}{n} \|\vec{y} - \vec{x}(\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}\|^2 \\ &= \frac{1}{n} \|(\vec{I} - \vec{H})\vec{y}\|^2 \\ &= \frac{1}{n} \|e\|^2 = \frac{1}{n} \sum e_i^2 = \frac{1}{n} SSE \end{aligned}$$

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

$$\begin{aligned}\text{Bias}(\hat{\sigma}^2) &= E\left[\frac{1}{n} \text{sse}\right] - E\left[\frac{\text{sse}}{n-(p+1)}\right] \\ &= E\left[\frac{1}{n} E\left(\frac{1}{n} - \frac{1}{n-(p+1)}\right)\right] \\ &\stackrel{n}{\rightarrow} -\frac{1}{n-(p+1)} E\left(\frac{1}{n}\right)^2 \\ &= \left(\frac{n-(p+1)}{n} - 1\right) \hat{\sigma}^2 = \left(\frac{p+1}{n}\right) \hat{\sigma}^2\end{aligned}$$

### 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  $\beta$  which yields the ridge estimates.

$$f(\vec{\beta}) = M_{\mu_0}(\vec{\alpha}_0, \tau^2 I_p) \quad f(\sigma^2) \propto \frac{1}{\sigma^2}$$

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

$$\begin{aligned}f(\vec{\beta}, \sigma^2 | \vec{x}, \vec{y}) &\propto (\sigma^2)^{-\frac{n}{2}-1} e^{-\frac{1}{2\sigma^2} \|(\vec{y} - \vec{x}\vec{\beta})\|^2} \cdot \frac{1}{2\tau^2} \|\vec{\beta}\|^2 \\ \vec{\beta}^{\text{ridge}} &= \underset{\vec{\beta}}{\operatorname{argmax}} \left\{ (-\frac{n}{2} - 1) D_{\text{ridge}} + \frac{1}{2\sigma^2} \|(\vec{y} - \vec{x}\vec{\beta})\|^2 + \frac{1}{2\tau^2} \|\vec{\beta}\|^2 \right\} \\ &= \underset{\vec{\beta}}{\operatorname{argmax}} \left\{ -\frac{1}{2\sigma^2} \|(\vec{y} - \vec{x}\vec{\beta})\|^2 + \frac{6^2}{\tau^2} \|\vec{\beta}\|^2 \right\} \\ &= \underset{\vec{\beta}}{\operatorname{argmin}} \left\{ \|(\vec{y} - \vec{x}\vec{\beta})\|^2 + \frac{\sigma^2}{\tau^2} \|\vec{\beta}\|^2 \right\} \\ &= \frac{\partial}{\partial \vec{\beta}} \left[ (\vec{y} - \vec{x}\vec{\beta})^T (\vec{y} - \vec{x}\vec{\beta}) + \vec{\beta}^T \vec{\beta} \right] \\ &= \frac{\partial}{\partial \vec{\beta}} \left[ \vec{y}^T \vec{y} - 2\vec{\beta}^T \vec{x}^T \vec{y} + \vec{\beta}^T \vec{x}^T \vec{x} \vec{\beta} + \vec{\beta}^T (\vec{x}^T \vec{x} + \lambda I_p) \vec{\beta} \right] \\ &= \frac{\partial}{\partial \vec{\beta}} \left[ -2\vec{\beta}^T \vec{x}^T \vec{y} + \vec{\beta}^T (\vec{x}^T \vec{x} + \lambda I_p) \vec{\beta} \right] \\ &- 2\vec{x}^T \vec{y} + 2(\vec{x}^T \vec{x} + \lambda I_p) \vec{\beta} \stackrel{\text{set}}{=} 0 \\ \vec{x}^T \vec{y} - (\vec{x}^T \vec{x} + \lambda I_p) \vec{\beta} &\Rightarrow \vec{\beta}_{\text{ridge}} = (\vec{x}^T \vec{x} + \lambda I_p)^{-1} \vec{x}^T \vec{y}\end{aligned}$$

- (h) [easy] Describe why Lasso estimation has the added bonus of being able to perform variable selection and ridge does not.

Because Lasso has a minimum at  $\beta = 0$   
due to the Laplace prior

### Problem 3

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

- (a) [easy] If we only know that the errors  $\varepsilon_1, \dots, \varepsilon_n$  are independent, what tried and true method can we employ to get asymptotically valid inference for  $\beta$ ?

Bootstrap

- (b) [easy] If we know that the errors  $\varepsilon_1, \dots, \varepsilon_n$  are iid with expectation zero and variance  $\sigma^2$  for all values of  $x$  (i.e. the errors are “homoskedastic”) but the errors are not necessarily normal, what is the asymptotic distribution of  $B$ ?

$$\tilde{\beta} \sim N(\hat{\beta}, \sigma^2 (X^\top X)^{-1})$$

- (c) [easy] If we know that the errors  $\varepsilon_1, \dots, \varepsilon_n \stackrel{ind}{\sim} \mathcal{N}(0, \sigma_i^2)$  which means the errors are “heteroskedastic”, what is the asymptotic distribution of  $B$  using the Huber-White estimator?

$$\tilde{\beta} \sim N(\hat{\beta}, (X^\top X)^{-1} X^\top \hat{\Omega} X (X^\top X)^{-1})$$

- (d) [easy] If we know that the errors  $\varepsilon_1, \dots, \varepsilon_n$  are independent with expectation zero and variance  $\sigma_i^2$  which means the errors are “heteroskedastic”, what is the asymptotic distribution of  $B$  using the Huber-White estimator?

$$\tilde{\beta} \sim N(\hat{\beta}, (X^\top X)^{-1} X^\top \hat{\Omega} X (X^\top X)^{-1})$$

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

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

$$\beta_0 \sim \beta_p \stackrel{\text{ind}}{\sim} \text{Laplace}(0, \tau^2) = \frac{1}{2\tau} \cdot e^{-\frac{|\beta_j|}{\tau}}$$

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

$$\begin{aligned} f(\vec{\beta}, \sigma^2 | \mathbf{x}, \mathbf{y}) &\propto (\sigma^{-1})^{n-1} e^{-\frac{1}{2\sigma^2} \| \mathbf{y} - \mathbf{x}\vec{\beta} \|^2} e^{-\frac{1}{\tau^2} \sum_{j=0}^p |\beta_j|} \\ \hat{\vec{\beta}} &= \arg \min \left\{ \left( -\frac{n}{2} - 1 \right) \ln \sigma^2 + \frac{1}{2\sigma^2} \| \mathbf{y} - \mathbf{x}\vec{\beta} \|^2 + \frac{1}{\tau^2} \sum_{j=0}^p |\beta_j| \right\} \\ &= \arg \min \left\{ -\frac{1}{2\sigma^2} \left( \| \mathbf{y} - \mathbf{x}\vec{\beta} \|^2 + \frac{2\sigma^2}{\tau^2} \sum_{j=0}^p |\beta_j| \right) \right\} \\ &= \arg \min \left\{ \| \mathbf{y} - \mathbf{x}\vec{\beta} \|^2 + \frac{2\sigma^2}{\tau^2} \sum_{j=0}^p |\beta_j| \right\} \\ \hat{\beta}_{\text{lasso}} &= \arg \min \left\{ \| \mathbf{y} - \mathbf{x}\vec{\beta} \|^2 + \lambda \sum_{j=0}^p |\beta_j| \right\} \end{aligned}$$

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

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

*filtering out variables that have little impact  
on the model*

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

### 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.

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{x}) = \prod_{i=1}^n \left( \frac{1}{1 + e^{-\tilde{y}_i \tilde{b}}} \right)^{y_i} \left( \frac{1}{1 + e^{\tilde{y}_i \tilde{b}}} \right)^{1-y_i}$$

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

$$\hat{LR} = \frac{\prod_{i=1}^n \left( \frac{1}{1 + e^{-\tilde{y}_i \tilde{b}}} \right)^{y_i} \left( \frac{1}{1 + e^{\tilde{y}_i \tilde{b}}} \right)^{1-y_i}}{\prod_{i=1}^n \left( \frac{1}{1 + e^{-\tilde{b}_0}} \right)^{y_i} \left( \frac{1}{1 + e^{\tilde{b}_0}} \right)^{1-y_i}}$$

$$\tilde{b}_0 = \arg\max \left( \prod_{i=1}^n \left( \frac{1}{1 + e^{-v}} \right)^{y_i} \left( \frac{1}{1 + e^v} \right)^{1-y_i} \right)$$

- (c) [harder] Let  $Y_i \stackrel{ind}{\sim} \text{Poisson}(\theta_i)$  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. Write out the likelihood ratio when testing  $H_0 : \beta_2 = \beta_3 = 0$ .

$$\hat{LR} = \frac{\prod_{i=1}^n e^{(\tilde{y}_i \tilde{b})} Y_i! e^{-e^{\tilde{y}_i \tilde{b}}}}{y_i!}$$

$$\frac{e^{\tilde{b}(y_1)} e^{\tilde{b}}}{y_1!} + \frac{\prod_{i=3}^n e^{(\tilde{y}_i \tilde{b})} Y_i! e^{-e^{\tilde{y}_i \tilde{b}}}}{y_i!}$$

- (d) [harder] Let  $Y_i \stackrel{\text{ind}}{\sim} \text{Weibull}(k, \theta_i)$  for  $i = 1, \dots, n$  where  $\theta_i = e^{\mathbf{x}_i \beta}$  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  $H_0 : \beta_2 = \beta_3 = 0$ .

$$\hat{\mathcal{L}} = \frac{\prod_{\{i: c_i=0\}} \left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^{k-1} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k} \prod_{\{i: c_i=1\}} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k}}{\prod_{c_i=0} \left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^{k-1} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k} \prod_{i: c_i=1} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k}}$$

where  $\hat{\beta} = \arg\max_{\beta} \left\{ \prod_{c_i=0} \left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^{k-1} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k} \prod_{i: c_i=1} e^{-\left( \frac{y_i}{e^{\mathbf{x}_i \cdot \beta}} \right)^k} \right\}$  under  $H_0$

- (e) [difficult] [MA] Let  $Y_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\theta_i, \sigma^2)$  for  $i = 1, \dots, n$  where  $\theta_i = \mathbf{x}_i \beta$  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  $H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$ .