

# 24500 HW7

Bin Yu

Feb 27, 2025

## Question 1

(a)

Recall the fact that

$$\frac{\|y - X\hat{\beta}\|^2}{\sigma^2} \sim \chi_{n-p}^2.$$

Thus,

$$\hat{\sigma}^2 = \frac{\|y - X\hat{\beta}\|^2}{n-p} = \sigma^2 \frac{\chi_{n-p}^2}{n-p}.$$

$$E(\hat{\sigma}^2) = E\left(\sigma^2 \cdot \frac{\chi_{n-p}^2}{n-p}\right) = \sigma^2 \frac{E(\chi_{n-p}^2)}{n-p}.$$

Since for  $\chi_{n-p}^2$  we have

$$E(\chi_{n-p}^2) = n-p,$$

it follows that

$$E(\hat{\sigma}^2) = \sigma^2 \frac{n-p}{n-p} = \sigma^2.$$

Hence,  $\hat{\sigma}^2$  is an unbiased estimator of  $\sigma^2$ .

$$\text{Var}(\hat{\sigma}^2) = \text{Var}\left(\sigma^2 \cdot \frac{\chi_{n-p}^2}{n-p}\right) = \sigma^4 \text{Var}\left(\frac{\chi_{n-p}^2}{n-p}\right).$$

We know

$$\text{Var}(\chi_{n-p}^2) = 2(n-p).$$

Hence,

$$\text{Var}\left(\frac{\chi_{n-p}^2}{n-p}\right) = \frac{1}{(n-p)^2} \text{Var}(\chi_{n-p}^2) = \frac{2(n-p)}{(n-p)^2} = \frac{2}{n-p}.$$

Therefore,

$$\text{Var}(\hat{\sigma}^2) = \sigma^4 \frac{2}{n-p} = \frac{2\sigma^4}{n-p}.$$

A chi-square random variable  $\chi_{n-p}^2$  can be viewed as the sum of squares of  $n-p$  i.i.d. standard normal variables  $Z_i \sim N(0, 1)$ :

$$\chi_{n-p}^2 = \sum_{i=1}^{n-p} Z_i^2.$$

For large  $n - p$ , the distribution of this sum can be approximated by a normal distribution via the Central Limit Theorem, we have:

$$Y = \chi_{n-p}^2, \quad \mu = n - p, \quad \text{and} \quad \sigma_Y^2 = 2(n - p),$$

then

$$\frac{Y - \mu}{\sqrt{2(n - p)}} = \frac{\chi_{n-p}^2 - (n - p)}{\sqrt{2(n - p)}} \xrightarrow{d} N(0, 1) \quad \text{as } (n - p) \rightarrow \infty.$$

Since

$$\hat{\sigma}^2 - \sigma^2 = \sigma^2 \left( \frac{\chi_{n-p}^2}{n-p} - 1 \right),$$

we define

$$Z = \frac{\chi_{n-p}^2 - (n - p)}{\sqrt{2(n - p)}} \xrightarrow{d} N(0, 1).$$

Note that

$$\frac{\chi_{n-p}^2}{n-p} - 1 = \frac{\chi_{n-p}^2 - (n - p)}{n - p} = \frac{Z}{\sqrt{n - p}} \sqrt{2}.$$

Thus,

$$\hat{\sigma}^2 - \sigma^2 = \sigma^2 \left( \frac{\chi_{n-p}^2}{n-p} - 1 \right) = \sigma^2 \frac{Z}{\sqrt{n - p}} \sqrt{2}.$$

Multiplying both sides by  $\sqrt{n - p}$  yields

$$\sqrt{n - p} (\hat{\sigma}^2 - \sigma^2) = \sigma^2 \cdot Z \sqrt{2}.$$

As  $Z$  converges in distribution to  $N(0, 1)$ , the right-hand side converges to a normal random variable with mean 0 and variance  $2\sigma^4$ . Hence,

$$\sqrt{n - p} (\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4).$$

the asymptotic distribution is

$$\sqrt{n - p} (\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4)$$

**(b)**

We use the Delta method. From part (a), if

$$\sqrt{n - p} (\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4),$$

then for a differentiable  $g$ ,

$$\sqrt{n - p} (g(\hat{\sigma}^2) - g(\sigma^2)) \xrightarrow{d} N\left(0, [g'(\sigma^2)]^2 \cdot 2\sigma^4\right).$$

We want the limiting variance to be 1, so we need

$$[g'(\sigma^2)]^2 2\sigma^4 = 1$$

$$g'(\sigma^2) = \frac{1}{\sqrt{2}\sigma^2}.$$

Integrating,

$$g(x) = \int \frac{1}{\sqrt{2}x} dx = \frac{1}{\sqrt{2}} \log(x) + C.$$

Ignore the constant  $C$  for inference purposes, so a choice is

$$g(x) = \frac{1}{\sqrt{2}} \log(x).$$

Then

$$\sqrt{n - p} (g(\hat{\sigma}^2) - g(\sigma^2)) \xrightarrow{d} N(0, 1).$$

(c)

We have found that a convenient transform for  $\hat{\sigma}^2$  is

$$g(x) = \frac{1}{\sqrt{2}} \log(x).$$

Under suitable regularity conditions (and for large  $n - p$ ), one obtains

$$\begin{aligned} \sqrt{n-p} (g(\hat{\sigma}^2) - g(\sigma^2)) &= \sqrt{n-p} \left( \frac{1}{\sqrt{2}} \log(\hat{\sigma}^2) - \frac{1}{\sqrt{2}} \log(\sigma^2) \right) \xrightarrow{d} N(0, 1). \\ \frac{\sqrt{n-p}}{\sqrt{2}} (\log(\hat{\sigma}^2) - \log(\sigma^2)) &\xrightarrow{d} N(0, 1). \end{aligned}$$

To construct  $(1 - \alpha)$  confidence interval, begin from the asymptotic result:

$$P\left(-z_{1-\alpha/2} \leq \frac{\sqrt{n-p}}{\sqrt{2}} [\log(\hat{\sigma}^2) - \log(\sigma^2)] \leq z_{1-\alpha/2}\right) = 1 - \alpha.$$

Inside the probability:

$$-z_{1-\alpha/2} \leq \frac{\sqrt{n-p}}{\sqrt{2}} [\log(\hat{\sigma}^2) - \log(\sigma^2)] \leq z_{1-\alpha/2}$$

is equivalent to

$$\log(\hat{\sigma}^2) - z_{1-\alpha/2} \sqrt{\frac{2}{n-p}} \leq \log(\sigma^2) \leq \log(\hat{\sigma}^2) + z_{1-\alpha/2} \sqrt{\frac{2}{n-p}}.$$

Exponentiating each part, we obtain

$$\hat{\sigma}^2 \exp\left(-z_{1-\alpha/2} \sqrt{\frac{2}{n-p}}\right) \leq \sigma^2 \leq \hat{\sigma}^2 \exp\left(z_{1-\alpha/2} \sqrt{\frac{2}{n-p}}\right).$$

Hence, we have the corresponding  $(1 - \alpha)$  confidence interval for  $\sigma^2$ :

$$\left[ \hat{\sigma}^2 \exp\left(-z_{1-\alpha/2} \sqrt{\frac{2}{n-p}}\right), \hat{\sigma}^2 \exp\left(z_{1-\alpha/2} \sqrt{\frac{2}{n-p}}\right) \right]$$

## Question 2

(a)

We want to solve

$$\hat{\beta}_0 = \arg \min_{\beta_0} \sum_{i=1}^n (y_i - \beta_0)^2.$$

Define the objective function:

$$Q(\beta_0) = \sum_{i=1}^n (y_i - \beta_0)^2.$$

To find the minimizer, we take the derivative of  $Q(\beta_0)$  with respect to  $\beta_0$  and set it to zero:

$$\frac{\partial}{\partial \beta_0} \sum_{i=1}^n (y_i - \beta_0)^2 = \sum_{i=1}^n \frac{\partial}{\partial \beta_0} (y_i - \beta_0)^2 = \sum_{i=1}^n (-2)(y_i - \beta_0).$$

Hence,

$$0 = \sum_{i=1}^n (y_i - \beta_0) = \sum_{i=1}^n y_i - n \beta_0.$$

$$\hat{\beta}_0 = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y},$$

the sample mean of the observations.

Because each  $y_i$  is distributed as  $N(\beta_0, \sigma^2)$  and they are i.i.d., we know

$$\bar{y} = \hat{\beta}_0 = \frac{1}{n} \sum_{i=1}^n y_i$$

is also normally distributed with

$$E(\hat{\beta}_0) = \beta_0, \quad \text{Var}(\hat{\beta}_0) = \frac{\sigma^2}{n}.$$

Hence,

$$\hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma^2}{n}\right).$$

$$E[(\hat{\beta}_0 - \beta_0)^2] = \text{Var}(\hat{\beta}_0) + E^2[(\hat{\beta}_0 - \beta_0)] = \frac{\sigma^2}{n} + 0 = \frac{\sigma^2}{n}$$

Thus,

$$E[(\hat{\beta}_0 - \beta_0)^2] = \frac{\sigma^2}{n}.$$

**(b)**

To find the minimizer, we take the derivative of  $Q(\beta_0)$  with respect to  $\beta_0$  and set it to zero:

$$\frac{\partial}{\partial \beta_0} Q(\beta_0, \beta_1) = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-1) = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i).$$

Setting this to zero:

$$0 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = \sum_{i=1}^n y_i - n\beta_0 - \beta_1 \sum_{i=1}^n x_i.$$

Hence,

$$\sum_{i=1}^n y_i = n\beta_0 + \beta_1 \sum_{i=1}^n x_i. \tag{1}$$

For  $\beta_1$ :

$$\frac{\partial}{\partial \beta_1} Q(\beta_0, \beta_1) = \sum_{i=1}^n 2(y_i - \beta_0 - \beta_1 x_i)(-x_i) = -2 \sum_{i=1}^n x_i(y_i - \beta_0 - \beta_1 x_i).$$

Setting this to zero:

$$0 = \sum_{i=1}^n x_i(y_i - \beta_0 - \beta_1 x_i) = \sum_{i=1}^n x_i y_i - \beta_0 \sum_{i=1}^n x_i - \beta_1 \sum_{i=1}^n x_i^2.$$

Thus,

$$\sum_{i=1}^n x_i y_i = \beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2. \tag{2}$$

Solving equations (1) and (2) gives:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ .

In the standard setting of simple linear regression we learned from class:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2), \quad i = 1, \dots, n,$$

with  $x_1, \dots, x_n$  treated as fixed, the well-known OLS result states that:

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

Even if the true slope  $\beta_1 = 0$ , the same formula for  $\text{Var}(\hat{\beta}_0)$  remains valid.

Therefore,

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

Now calculate  $E[\hat{\beta}_1]$ . By linearity of expectation,

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

We examine its expectation. Note that

$$E[\hat{\beta}_1] = \frac{\sum_{i=1}^n (x_i - \bar{x}) E[y_i - \bar{y}]}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Since under the linear model,  $E[y_i] = \beta_0$  and

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_{i=1}^n y_i \\ E[\bar{y}] &= \beta_0, \end{aligned}$$

we have

$$E[y_i - \bar{y}] = \beta_0 - E[\bar{y}] = \beta_0 - \beta_0 = 0.$$

In this case,  $E[\hat{\beta}_1] = 0$ , and so  $\hat{\beta}_1$  is unbiased for the true slope (which is zero).

Next, consider  $\hat{\beta}_0$ . Recall that

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

Taking expectation,

$$E[\hat{\beta}_0] = E[\bar{y}] - E[\hat{\beta}_1] \bar{x} = \beta_0 - 0 = \beta_0$$

Hence,  $\hat{\beta}_0$  is also an unbiased estimator of  $\beta_0$ .

Finally, the mean squared error (MSE) of  $\hat{\beta}_0$  is simply its variance (since it is unbiased):

$$E[(\hat{\beta}_0 - \beta_0)^2] = \text{Var}(\hat{\beta}_0) + E^2[(\hat{\beta}_0 - \beta_0)] = \text{Var}(\hat{\beta}_0) = \text{Var}(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

(c)

In part (a), we estimate  $\beta_0$  by minimizing

$$\sum_{i=1}^n (y_i - \beta_0)^2,$$

which leads to the estimator

$$\hat{\beta}_0^{(a)} = \bar{y}.$$

Because  $y_i \sim N(\beta_0, \sigma^2)$  are i.i.d.,  $\bar{y}$  is unbiased for  $\beta_0$  with

$$\text{Var}(\hat{\beta}_0^{(a)}) = \frac{\sigma^2}{n}.$$

Hence,

$$E[(\hat{\beta}_0^{(a)} - \beta_0)^2] = \frac{\sigma^2}{n}.$$

In part (b), we include a slope parameter  $\beta_1$  in the model

$$y_i = \beta_0 + \beta_1 x_i,$$

and estimate both  $\beta_0$  and  $\beta_1$  by ordinary least squares (OLS), even though the true slope is 0. The resulting estimator for the intercept is

$$\hat{\beta}_0^{(b)} = \bar{y} - \hat{\beta}_1 \bar{x}.$$

From standard regression theory, its variance is given by

$$\text{Var}(\hat{\beta}_0^{(b)}) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

Thus, in this scenario,

$$E[(\hat{\beta}_0^{(b)} - \beta_0)^2] = \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Therefore, we observe that:

- **Both are correct in expectation.** In each method,  $\hat{\beta}_0$  is an unbiased estimator of  $\beta_0$ , so both approaches give valid estimates on average.
- **Different variance (accuracy).** Approach (a) yields  $\text{Var}(\hat{\beta}_0^{(a)}) = \frac{\sigma^2}{n}$ . Approach (b) inflates the variance by an extra term

$$\frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

because it also estimates a slope parameter that, in truth, does not exist ( $\beta_1 = 0$ ).

- **implication.** Therefore, fitting a more complicated model (b) when the simpler model is actually correct will not improve (and can worsen) the precision of estimating  $\beta_0$ . The extra parameter  $\beta_1$  adds uncertainty to the intercept estimate.

Hence, while both approaches are *correct* in the sense of unbiasedness, method (a) has smaller variance for estimating  $\beta_0$  in this particular data-generating scenario.

### Question 3

(a)

We estimate  $\beta_0$  by solving

$$\hat{\beta}_0 = \arg \min_{\beta_0} \sum_{i=1}^n (y_i - \beta_0)^2.$$

From earlier derivations (as in Question 2(a)), we know:

$$\hat{\beta}_0 = \bar{y} = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i + \varepsilon_i) = \beta_0 + \beta_1 \bar{x} + \bar{\varepsilon},$$

where  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$  and  $\bar{\varepsilon} = \frac{1}{n} \sum_{i=1}^n \varepsilon_i$ .

Since we have:

$$E[y_i] = \beta_0 + \beta_1 x_i.$$

Taking  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ , we observe

$$E(\hat{\beta}_0) = \frac{1}{n} \sum_{i=1}^n E[y_i] = \frac{1}{n} \sum_{i=1}^n (\beta_0 + \beta_1 x_i) = \beta_0 + \beta_1 \bar{x},$$

which means the estimator is biased for  $\beta_0$  unless  $\bar{x} = 0$ .

For the variance, since  $y_i$  are i.i.d. with variance  $\sigma^2$ ,

$$\text{Var}(\hat{\beta}_0) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(y_i) = \frac{n \sigma^2}{n^2} = \frac{\sigma^2}{n}.$$

Because

$$\text{Bias}(\hat{\beta}_0) = E[\hat{\beta}_0] - \beta_0 = \beta_1 \bar{x},$$

the MSE is

$$E[(\hat{\beta}_0 - \beta_0)^2] = \text{Var}(\hat{\beta}_0) + [\text{Bias}(\hat{\beta}_0)]^2 = \frac{\sigma^2}{n} + \beta_1^2 \bar{x}^2.$$

**(b)**

Now we include the slope term and solve:

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2.$$

From OLS theory (see Question 2(b)), we get

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

From the standard OLS estimates we have known that:

$$E[\hat{\beta}_0] = \beta_0,$$

and the variance is

$$\text{Var}(\hat{\beta}_0) = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

Therefore,

$$E[(\hat{\beta}_0 - \beta_0)^2] = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

which is the same as the result in Question 2(b).

(c)

Suppose the true data-generating process is

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad \beta_1 \neq 0.$$

However, if we fit the simpler model from part (a), we are assuming:

$$y_i \approx \beta_0 \quad (\text{no slope term}),$$

which ignores the covariates  $\{x_i\}$ . Our results show:

If  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \neq 0$ , then the simpler model's estimator

$$\hat{\beta}_0^{(\text{mean})} = \bar{y}$$

will be a biased estimator of  $\beta_0$ . Indeed,

$$\text{Bias}(\hat{\beta}_0^{(\text{mean})}) = \beta_1 \bar{x}.$$

This arises because the slope term  $\beta_1 x_i$  is ignored, shifting it away from the true  $\beta_0$ .

By contrast, part (b) incorporates the slope parameter  $\beta_1$ , matching the true data distribution. The estimator  $\hat{\beta}_0$  in the full regression model is unbiased for  $\beta_0$ , and its MSE depends only on its variance.

Also, the mean squared error (MSE) of the simpler estimator is

$$\text{MSE}_{\text{simple}} = \frac{\sigma^2}{n} + \beta_1^2 \bar{x}^2.$$

and the full model's MSE is

$$\text{MSE}_{\text{full}} = \sigma^2 \left[ \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right].$$

To compare the two, consider the difference:

$$\begin{aligned} \text{MSE}_{\text{simple}} - \text{MSE}_{\text{full}} &= \left[ \frac{\sigma^2}{n} + \beta_1^2 \bar{x}^2 \right] - \left[ \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ &= \beta_1^2 \bar{x}^2 - \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \bar{x}^2 \left( \beta_1^2 - \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right). \end{aligned}$$

Thus, the simpler model yields a higher MSE than the full model (i.e.,  $\text{MSE}_{\text{simple}} > \text{MSE}_{\text{full}}$ ) if and only if

$$\beta_1^2 > \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

Equivalently, if

$$|\beta_1| > \frac{\sigma}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}.$$

Therefore, when the magnitude of the slope  $\beta_1$  exceeds this threshold, the bias incurred by omitting the slope leads to a larger MSE for the simpler model compared to the full model.

## Question 4

(a)

Let us denote by  $S^c$  the complement of  $S$  in  $\{0, 1, 2, \dots, p-1\}$ . Hence,  $\beta_{S^c}$  is the vector of parameters  $\beta_j$  with  $j \in S^c$ , and  $X_{S^c}$  is the corresponding submatrix of  $X$  consisting of those columns indexed by  $S^c$ .



Under the null hypothesis  $H_0 : \beta_S = 0$ , the regression model only involves the covariates in  $S^c$ , i.e.,

$$y \sim N(X_{S^c} \beta_{S^c}, \sigma^2 I_n)$$

Imposing the restriction  $\beta_S = 0$  yields the least squares estimator (LSE) for the remaining parameters:

$$\hat{\beta}_{H_0} = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top y.$$

where,

$$X_{S^c} = \begin{pmatrix} | & & | \\ x_{:,j_1} & \cdots & x_{:,j_{p-s}} \\ | & & | \end{pmatrix},$$

with  $j_1, \dots, j_{p-s} \in S^c$ .

We can compute the mean and variance of  $\hat{\beta}_{H_0}$  by noting that  $\hat{\beta}_{H_0}$  is a linear function of  $y$ . Specifically,

$$\hat{\beta}_{H_0} = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top y,$$

$$E[\hat{\beta}_{H_0}] = E[(X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top y] = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top E[y] = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top X_{S^c} \beta_{S^c} = \beta_{S^c}.$$

Hence,  $\hat{\beta}_{H_0}$  is unbiased for  $\beta_{S^c}$  under  $H_0$ .

Since  $y$  has covariance matrix  $\sigma^2 I_n$  under the model, let  $A = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top$ , then  $\hat{\beta}_{H_0} = Ay$

$$\text{Var}(\hat{\beta}_{H_0}) = A \text{Var}(y) A^\top = A \sigma^2 I_n A^\top = \sigma^2 A A^\top.$$

$$A A^\top = (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top X_{S^c} [(X_{S^c}^\top X_{S^c})^{-1}]^\top = (X_{S^c}^\top X_{S^c})^{-1},$$

using the fact that  $(X_{S^c}^\top X_{S^c})$  is symmetric and invertible. Hence

$$\text{Var}(\hat{\beta}_{H_0}) = \sigma^2 (X_{S^c}^\top X_{S^c})^{-1}.$$

Because  $\hat{\beta}_{H_0}$  is a linear transformation of the Gaussian vector  $y$ ,  $\hat{\beta}_{H_0}$  itself is normally distributed. Putting it all together,

$$\hat{\beta}_{H_0} \sim N(\beta_{S^c}, \sigma^2 (X_{S^c}^\top X_{S^c})^{-1}).$$

(b)

We can write

$$\begin{aligned} y - y_{H_0} &= y - \hat{y} + \hat{y} - y_{H_0}. \\ y - y_{H_0} &= (y - \hat{y}) + (\hat{y} - y_{H_0}). \end{aligned}$$

$$\|y - y_{H_0}\|^2 = \|y - \hat{y} + \hat{y} - y_{H_0}\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - y_{H_0}\|^2 + 2(y - \hat{y})^\top (\hat{y} - y_{H_0})$$

For the cross product term  $(y - \hat{y})^\top (\hat{y} - y_{H_0})$ : section\*Algebraic Proof of  $(y - \hat{y})^\top (\hat{y} - y_{H_0}) = 0$  Using Projection Matrices

Let

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

We claim  $H$  is the projection matrix onto the column space of  $X$ . To verify:

- $H$  is symmetric:

$$H^\top = [X (X^\top X)^{-1} X^\top]^\top = X [(X^\top X)^{-1}]^\top X^\top = X (X^\top X)^{-1} X^\top = H,$$

since  $(X^\top X)^{-1}$  is symmetric.

- $H$  is idempotent:

$$H^2 = [X (X^\top X)^{-1} X^\top] [X (X^\top X)^{-1} X^\top] = X (X^\top X)^{-1} X^\top X (X^\top X)^{-1} X^\top = X (X^\top X)^{-1} X^\top = H.$$

Hence  $H$  is projection matrix that projects onto  $\mathcal{C}(X)$ , the column space of  $X$ .

Let  $S^c$  be the complement of  $S$  in  $\{0, 1, \dots, p-1\}$ , and let  $X_{S^c} \in R^{n \times (p-s)}$  be the submatrix of  $X$  containing the columns indexed by  $S^c$ . Define

$$H_c = X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top.$$

Similarly, one checks that:

$$H_c^\top = H_c, \quad H_c^2 = H_c,$$

so  $H_c$  is the projection matrix onto  $\mathcal{C}(X_{S^c})$ . Because  $X_{S^c}$  consists of a subset of the columns of  $X$ , we have

$$\mathcal{C}(X_{S^c}) \subseteq \mathcal{C}(X).$$

In other words, projecting first onto  $\mathcal{C}(X_{S^c})$  is a stricter (or smaller) projection than projecting onto  $\mathcal{C}(X)$ .

Because  $\mathcal{C}(X_{S^c}) \subseteq \mathcal{C}(X)$ , we have the nesting of subspaces, and we have:

$$H H_c = H_c, \quad H_c H = H_c.$$

Here is the verification:

$$\begin{aligned} H H_c &= [X (X^\top X)^{-1} X^\top] [X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top] \\ &= X (X^\top X)^{-1} X^\top X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top. \end{aligned}$$

Because  $X_{S^c}$  consists of a subset of the columns of  $X$ , each column of  $X_{S^c}$  lies in the column space of  $X$ . Thus there exists some matrix  $B$  such that

$$X_{S^c} = X B.$$

Consequently,

$$X^\top X_{S^c} = X^\top (X B) = (X^\top X) B.$$

$$\begin{aligned} H H_c &= [X (X^\top X)^{-1} X^\top] [X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top] \\ &= X (X^\top X)^{-1} \underbrace{X^\top X_{S^c}}_{\text{cross-term}} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top \\ &= X (X^\top X)^{-1} [(X^\top X) B] (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top \quad (\text{since } X^\top X_{S^c} = (X^\top X) B). \end{aligned}$$

Thus

$$H H_c = X [(X^\top X)^{-1} (X^\top X) B] (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top = X B (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top.$$

and  $X_{S^c} = X B$ , therefore,

$$H H_c = X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top = H_c$$

Similarly,

$$H_c H = [X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top] [X (X^\top X)^{-1} X^\top].$$

Again substitute  $X_{S^c} = X B$ :

$$X_{S^c}^\top X = (X B)^\top X = B^\top \underbrace{X^\top X}_{\text{used in the same manner}},$$

and proceed as above. The same pattern emerges, and one obtains  $H_c H = H_c$ .

In the full model, the fitted values are  $\hat{y} = H y$ , and hence the residuals are  $y - \hat{y} = (I - H) y$ .

In the restricted model  $H_0$ , the fitted values are  $y_{H_0} = H_c y$ . Hence the difference between the two fitted values is  $\hat{y} - y_{H_0} = (H - H_c) y$ .

$$\begin{aligned} (y - \hat{y})^\top (\hat{y} - y_{H_0}) &= \left[ (I - H) y \right]^\top \left[ (H - H_c) y \right] \\ &= y^\top (I - H)^\top (H - H_c) y. \end{aligned}$$

Because  $H$  is symmetric,  $(I - H)^\top = I - H$ . So:

$$\begin{aligned} &= y^\top (I - H) (H - H_c) y. \\ (I - H) (H - H_c) &= (I - H) H - (I - H) H_c. \end{aligned}$$

for each term:

$$(I - H) H = H - H^2.$$

But  $H^2 = H$  (idempotent), so

$$H - H^2 = 0.$$

Hence  $(I - H) H = 0$ . and,

$$(I - H) H_c = H_c - H H_c.$$

But from above,  $H H_c = H_c$ . Therefore

$$H_c - H H_c = 0.$$

Hence  $(I - H) H_c = 0$ .

Putting these together,

$$(I - H) (H - H_c) = (I - H) H - (I - H) H_c = 0 - 0 = 0.$$

Thus

$$y^\top (I - H) (H - H_c) y = y^\top 0 y = 0.$$

Hence

$$(y - \hat{y})^\top (\hat{y} - y_{H_0}) = 0,$$

We obtain:

$$\|y - y_{H_0}\|^2 = \|y - \hat{y} + \hat{y} - y_{H_0}\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - y_{H_0}\|^2 + 0 = \|y - \hat{y}\|^2 + \|\hat{y} - y_{H_0}\|^2$$

Thus we have shown the desired identity:

$$\|y - y_{H_0}\|^2 = \|y - \hat{y}\|^2 + \|\hat{y} - y_{H_0}\|^2,$$

**(c)**

$\|y - \hat{y}\|^2$  and  $\|\hat{y} - y_{H_0}\|^2$  can each be written as a quadratic form in  $y$ . Specifically,

$$\|y - \hat{y}\|^2 = \|(I - H) y\|^2, \quad \|\hat{y} - y_{H_0}\|^2 = \|(H - H_c) y\|^2,$$

where

$$H = X (X^\top X)^{-1} X^\top, \quad H_c = X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top.$$

Since both  $(I - H)y$  and  $(H - H_c)y$  are linear transformations of the Gaussian vector  $y$ . Hence each is a Gaussian vector in  $R^n$ . To show the corresponding sums of squares are independent, we only need to show these two Gaussian

vectors are independent as random vectors. For jointly Gaussian vectors, uncorrelatedness implies independence. Hence it suffices to show their covariance is zero.

From part(b), we already know  $(I_n - H)(H - H_c) = 0$ , and  $H - H_c$  is also a projection matrix, indicating that  $(H - H_c)^T = (H - H_c)$

$$\text{Cov}((I_n - H)y, (H - H_c)y) = (I_n - H) \text{Var}[y] (H - H_c)^T = \sigma^2 (I_n - H)(H - H_c) = 0_n.$$

Therefore, zero covariance implies  $\|y - \hat{y}\|^2$  and  $\|\hat{y} - y_{H_0}\|^2$  are independent.

(d)

Under  $H_0$ , the true model becomes

$$y = X_{S^c} \beta_{S^c} + \sigma Z,$$

where  $Z \sim N(0, I_n)$ . The restricted fit is  $\hat{y}_{H_0} = H_c y$ . Therefore,

$$\begin{aligned} y - \hat{y}_{H_0} &= (I_n - H_c)y = (I_n - H_c)(X_{S^c} \beta_{S^c} + \sigma Z). \\ &= [X_{S^c} - X_{S^c}(X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top X_{S^c}] \beta_{S^c} + \sigma (I_n - H_c) Z \\ &= [X_{S^c} - X_{S^c}] \beta_{S^c} + \sigma (I_n - H_c) Z \\ &= \sigma (I_n - H_c) Z. \end{aligned}$$

Thus,

$$\|y - \hat{y}_{H_0}\|^2 = \sigma^2 \|(I_n - H_c) Z\|^2.$$

Since  $Z \sim N(0, I_n)$ , the vector  $(I_n - H_c) Z$  is a Gaussian in  $R^n$  whose dimension equals the rank of  $(I_n - H_c)$ . By lemma from class,

$$\|(I_n - H_c) Z\|^2 \sim \chi_{\text{rank}(I_n - H_c)}^2.$$

Note:

$$\text{rank}(I_n - H_c) = \text{Tr}(I_n - H_c) = \text{Tr}(I_n) - \text{Tr}(H_c) = n - \text{Tr}(H_c).$$

The  $\text{Tr}(H_c)$  equals the dimension of the column space of  $X_{S^c}$ , which is  $p - s$  if  $X$  is full rank ( $p \leq n$ ) and  $|S| = s$ . Hence

$$\text{Tr}(H_c) = p - s \implies \text{rank}(I_n - H_c) = n - (p - s) = n - p + s.$$

Therefore,

$$\frac{\|y - \hat{y}_{H_0}\|^2}{\sigma^2} = \|(I_n - H_c) Z\|^2 \sim \chi_{n-p+s}^2.$$

For distribution of  $\|\hat{y} - \hat{y}_{H_0}\|^2/\sigma^2$ :

The difference between the full fit  $\hat{y} = H y$  and the restricted fit  $\hat{y}_{H_0} = H_c y$  is

$$\hat{y} - \hat{y}_{H_0} = (H - H_c)y.$$

Under  $H_0$ , write again  $y = X_{S^c} \beta_{S^c} + \sigma Z$ . Then

$$(H - H_c)y = (H - H_c)(X_{S^c} \beta_{S^c} + \sigma Z).$$

Therefore,

$$\begin{aligned} \hat{y} - \hat{y}_{H_0} &= (H - H_c)y = (H - H_c)(X_{S^c} \beta_{S^c} + \sigma Z) \\ &= \left[ X(X^\top X)^{-1} X^\top - X_{S^c}(X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top \right] (X_{S^c} \beta_{S^c} + \sigma Z) \end{aligned}$$

Since in (c) we already have:

$$X_{S^c} = X B.$$

Consequently,

$$X^\top X_{S^c} = X^\top (XB) = (X^\top X) B.$$

Substitute:

$$\begin{aligned} & \left[ X (X^\top X)^{-1} X^\top - X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top \right] (X_{S^c} \beta_{S^c} + \sigma Z) \\ = & \left[ X (X^\top X)^{-1} X^\top XB - X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top X_{S^c} \right] \beta_{S^c} + \sigma \left[ X (X^\top X)^{-1} X^\top - (X_{S^c} - X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top) \right] Z \\ = & \left[ XB - X_{S^c} \right] \beta_{S^c} + \sigma \left[ X (X^\top X)^{-1} X^\top - (X_{S^c} - X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top) \right] Z \\ = & (X_{S^c} - X_{S^c}) \beta_{S^c} + \sigma \left[ X (X^\top X)^{-1} X^\top - (X_{S^c} - X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top) \right] Z \\ = & 0 \cdot \beta_{S^c} + \sigma (H - H_c) Z = \sigma (H - H_c) Z. \end{aligned}$$

We get

$$\hat{y} - \hat{y}_{H_0} = \sigma (H - H_c) Z.$$

$$\|\hat{y} - \hat{y}_{H_0}\|^2 = \sigma^2 \|(H - H_c) Z\|^2.$$

By the same lemma,

$$\frac{\|\hat{y} - \hat{y}_{H_0}\|^2}{\sigma^2} = \|(H - H_c) Z\|^2 \sim \chi_{\text{rank}(H - H_c)}^2.$$

Finally, we compute

$$\text{rank}(H - H_c) = \text{Tr}(H - H_c) = \text{Tr}(H) - \text{Tr}(H_c).$$

Since  $H$  is the full projection of rank  $p$  (assuming  $X$  is  $n \times p$  with full column rank),

$$\text{Tr}(H) = p, \quad \text{Tr}(H_c) = p - s.$$

Hence

$$\text{rank}(H - H_c) = p - (p - s) = s.$$

Therefore,

$$\frac{\|\hat{y} - \hat{y}_{H_0}\|^2}{\sigma^2} \sim \chi_s^2.$$

Therefore, under the null hypothesis  $H_0 : \beta_S = 0$ , we have:

$$\frac{\|y - \hat{y}_{H_0}\|^2}{\sigma^2} \sim \chi_{n-p+s}^2, \quad \frac{\|\hat{y} - \hat{y}_{H_0}\|^2}{\sigma^2} \sim \chi_s^2.$$

(e)

Given the model

$$y \sim N(X\beta, \sigma^2 I_n),$$

we want to test

$$H_0 : \beta_S = 0 \quad \text{vs.} \quad H_1 : \beta_S \neq 0,$$

where  $S \subseteq \{0, 1, \dots, p-1\}$  is a subset of indices, of size  $s = |S|$ , and  $X_S$  (respectively  $X_{S^c}$ ) is the submatrix of  $X$  corresponding to those columns (respectively, the complementary columns).

- **Restricted fit:**  $\hat{y}_{H_0} = H_c y$ , where

$$H_c = X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top.$$

The residual sum of squares (RSS) under  $H_0$  is

$$\|y - \hat{y}_{H_0}\|^2.$$

- **Full fit:**  $\hat{y} = H y$ , where

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

The residual sum of squares (RSS) under the full model is

$$\|y - \hat{y}\|^2.$$

from part(d) we have

$$\|y - \hat{y}_{H_0}\|^2 \sim \sigma^2 \chi_{n-p+s}^2 \quad \text{and} \quad \|\hat{y} - \hat{y}_{H_0}\|^2 \sim \sigma^2 \chi_s^2 \quad (\text{independently}),$$

under  $H_0$ . Also

$$\|y - \hat{y}\|^2 = \|y - \hat{y}_{H_0}\|^2 - \|\hat{y} - \hat{y}_{H_0}\|^2.$$

Hence  $\|y - \hat{y}\|^2 / \sigma^2 \sim \chi_{n-p}^2$

The “extra sum of squares” contributed by the columns in  $S$  is

$$\|\hat{y} - \hat{y}_{H_0}\|^2 = \|y - \hat{y}_{H_0}\|^2 - \|y - \hat{y}\|^2.$$

From the previous results:

$$F = \frac{[\|\hat{y} - \hat{y}_{H_0}\|^2 / \sigma^2] / s}{[\|y - \hat{y}\|^2 / \sigma^2] / (n-p)}.$$

$$F = \frac{\|\hat{y} - \hat{y}_{H_0}\|^2 / s}{\|y - \hat{y}\|^2 / (n-p)} \sim F_{s, n-p} \quad (\text{under } H_0).$$

Thus the F-test statistic is

$$F = \frac{[\|\hat{y} - \hat{y}_{H_0}\|^2 / s]}{[\|y - \hat{y}\|^2 / (n-p)]},$$

and we reject  $H_0$  (i.e. conclude that at least one of  $\beta_j, j \in S$  is nonzero) if:

$$\frac{\|\hat{y} - \hat{y}_{H_0}\|^2 / s}{\|y - \hat{y}\|^2 / (n-p)} > F_{s, n-p, 1-\alpha}$$

(f)

In this scenario, we are testing:

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_{p-1} = 0$$

Hence the submatrix  $X_S$  is all columns of  $X$  except the intercept column  $X_0$ . Equivalently,  $X_{S^c}$  consists only of the intercept column. so

$$X_{S^c} = 1_n,$$

where  $1_n$  denotes the  $n$ -vector of all ones.

Under  $H_0$ , the restricted least squares fit is given by

$$\hat{y}_{H_0} = H_c y, \quad \text{with} \quad H_c = X_{S^c} (X_{S^c}^\top X_{S^c})^{-1} X_{S^c}^\top.$$

Since

$$X_{S^c} = 1_n \quad \text{and} \quad X_{S^c}^\top X_{S^c} = 1_n^\top 1_n = n,$$

it follows that

$$H_c = 1_n \frac{1}{n} 1_n^\top = \frac{1}{n} 1_n 1_n^\top.$$

Thus, the restricted fitted values are

$$\hat{y}_{H_0} = H_c y = \frac{1}{n} 1_n 1_n^\top y,$$

which is the  $n$ -vector with every entry equal to the sample mean  $\bar{y}$ .

From our earlier results (see part (e)), we have shown that

$$\hat{y} - \hat{y}_{H_0} = \sigma (H - H_c) Z,$$

where

$$H = X (X^\top X)^{-1} X^\top \quad \text{and} \quad Z \sim N(0, I_n).$$

Thus,

$$\|\hat{y} - \hat{y}_{H_0}\|^2 = \sigma^2 \|(H - H_c)Z\|^2.$$

By a standard lemma,

$$\frac{\|\hat{y} - \hat{y}_{H_0}\|^2}{\sigma^2} \sim \chi_{\text{rank}(H - H_c)}^2.$$

Now, the rank of  $(H - H_c)$  is given by

$$\text{rank}(H - H_c) = \text{Tr}(H - H_c) = \text{Tr}(H) - \text{Tr}(H_c).$$

Since  $X$  is  $n \times p$  with full column rank, the projection matrix  $H$  has

$$\text{Tr}(H) = p.$$

Moreover, as we have just shown,

$$H_c = \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top,$$

so

$$\text{Tr}(H_c) = \frac{1}{n} \text{Tr}(\mathbf{1}_n \mathbf{1}_n^\top) = \frac{1}{n} (\mathbf{1}_n^\top \mathbf{1}_n) = \frac{n}{n} = 1.$$

Thus,

$$\text{rank}(H - H_c) = p - 1.$$

It follows that

$$\frac{\|\hat{y} - \hat{y}_{H_0}\|^2}{\sigma^2} \sim \chi_{p-1}^2.$$

The results in (e) shows that in the full model the residual sum of squares (RSS) is

$$\|y - \hat{y}\|^2,$$

,under the full model, has  $n - p$  degrees of freedom:

$$\frac{\|y - \hat{y}\|^2}{\sigma^2} \sim \chi_{n-p}^2.$$

Thus, the standard  $F$ -statistic for testing

$$H_0 : \beta_1 = \beta_2 = \cdots = \beta_{p-1} = 0$$

is given by

$$F = \frac{[\|\hat{y} - \hat{y}_{H_0}\|^2 / (p - 1)]}{[\|y - \hat{y}\|^2 / (n - p)]} \sim F_{p-1, n-p}.$$

Noting that  $\hat{y}_{H_0} = \bar{y} \mathbf{1}_n$  the above  $F$ -statistic becomes

$$F = \frac{[\|\hat{y} - \bar{y} \mathbf{1}_n\|^2 / (p - 1)]}{[\|y - \hat{y}\|^2 / (n - p)]} \sim F_{p-1, n-p},$$

which is exactly the standard result taught in class. Therefore, we reject  $H_0$  at significance level  $\alpha$  if

$$F > F_{p-1, n-p, 1-\alpha}$$