

Setting Up The MLR Model

Population: $E(Y \mid X_1, X_2, \dots, X_{K-1}) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{K-1} X_{K-1}$

Define $\epsilon = Y - \beta_0 - \beta_1 X_1 - \beta_2 X_2 - \dots - \beta_{K-1} X_{K-1}$

X_1, \dots, X_{K-1} refer to $K - 1$ different regressors, not $K - 1$ obs. of a variable X

We have

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{K-1} X_{K-1} + \epsilon, \quad E(\epsilon \mid X_1, \dots, X_{K-1}) = 0$$

Setting Up The MLR Model

Representative iid sample from the population: $\{Y_i, X_{i1}, X_{i2}, \dots, X_{i,K-1}\}_{i=1}^n$

Sample satisfies

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_{K-1} X_{i,K-1} + \epsilon_i$$

$$E(\epsilon_i \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = 0$$

$$E(\epsilon_i^2 \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = \sigma_i^2 \text{ or } \sigma^2$$

$$E(\epsilon_i \epsilon_j \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = 0$$

for all $i, j = 1, \dots, n, i \neq j$

Setting Up The MLR Model

The equations

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_{K-1} X_{i,K-1} + \epsilon_i \quad \text{for all } i = 1, \dots, n$$

can be written

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,K-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,K-1} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{K-1} \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\text{i.e., } y = X\beta + \epsilon$$

X_{ij} is i th observation of the j th regressor

Setting Up The MLR Model

We can partition the regressor matrix by observation:

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,K-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,K-1} \end{bmatrix} = \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix}$$

and write the model as

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix} \beta + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \quad \text{or} \quad Y_i = X_{i*} \beta + \epsilon_i, \quad i = 1, 2, \dots, n$$

Setting Up The MLR Model

It is sometimes helpful to partition the regressor matrix by variable:

$$X = \left[\begin{array}{c|c|c|c|c} 1 & X_{11} & X_{12} & \cdots & X_{1,K-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,K-1} \end{array} \right] = [i_n \quad X_{*1} \quad X_{*2} \quad \cdots \quad X_{*,K-1}]$$

OLS estimation requires no perfect collinearity in the sample:

$$Xc = c_0 + c_1 X_{*1} + c_2 X_{*2} + \cdots + c_{K-1} X_{*,K-1} = 0_{n \times 1} \iff c = 0_K$$

where $c = [c_0 \quad \cdots \quad c_{K-1}]^T$. This assumption can also be written:

$$Xc \neq 0_n \iff c \neq 0_K$$

Setting Up The MLR Model

The assumption

$$E(\epsilon_i \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = 0 \quad \text{for } i = 1, \dots, n$$

can be written as

$$\begin{bmatrix} E(\epsilon_1 \mid X) \\ E(\epsilon_2 \mid X) \\ \vdots \\ E(\epsilon_n \mid X) \end{bmatrix} = 0_n \quad \iff \quad E(\epsilon \mid X) = 0_n$$

Setting Up The MLR Model

The assumptions

$$E(\epsilon_i^2 \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = \sigma_i^2 \quad \text{for } i = 1, \dots, n$$

$$E(\epsilon_i \epsilon_j \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = 0 \quad \text{for } i, j = 1, \dots, n, i \neq j$$

can be written as

$$\begin{bmatrix} E(\epsilon_1^2 \mid X) & E(\epsilon_1 \epsilon_2 \mid X) & \cdots & E(\epsilon_1 \epsilon_n \mid X) \\ E(\epsilon_2 \epsilon_1 \mid X) & E(\epsilon_2^2 \mid X) & \cdots & E(\epsilon_2 \epsilon_n \mid X) \\ \vdots & \vdots & \ddots & \vdots \\ E(\epsilon_n \epsilon_1 \mid X) & E(\epsilon_n \epsilon_2 \mid X) & \cdots & E(\epsilon_n^2 \mid X) \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_2^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n^2 \end{bmatrix}$$

or

$$\text{Var}(\epsilon \mid X) = E(\epsilon \epsilon^T \mid X) = \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2)$$

Setting Up The MLR Model

If the noise terms are homoskedastic, i.e.,

$$E(\epsilon_i^2 \mid X_{11}, \dots, X_{n1}, X_{12}, \dots, X_{n2}, \dots, X_{1,K-1}, \dots, X_{n,K-1}) = \sigma^2 \quad \text{for } i = 1, \dots, n,$$

then we have

$$E(\epsilon\epsilon^T \mid X) = \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 I_n$$

which we will assume for now

OLS Estimation of MLR

In matrix form, the MLR is

$$Y = X\beta + \varepsilon, \quad E(\varepsilon | X) = 0, \quad E(\varepsilon\varepsilon^T | X) = \sigma^2 I_n$$

For any potential estimator $\hat{\beta}$, define

- Fitted values: $\hat{y} = X\hat{\beta}$
- Residuals: $\hat{\varepsilon} = y - \hat{y} = y - X\hat{\beta}$
- Residual Sum of Squares $RSS(\hat{\beta})$: $\sum_{i=1}^n \hat{\varepsilon}_i^2 = \hat{\varepsilon}^T \hat{\varepsilon}$

OLS Estimation of MLR

We have

$$\begin{aligned}RSS(\hat{\beta}) &= \hat{\varepsilon}^T \hat{\varepsilon} = (y - X\hat{\beta})^T (y - X\hat{\beta}) \\&= (y^T - \hat{\beta}^T X^T)(y - X\hat{\beta}) \\&= y^T y - \hat{\beta}^T X^T y - y^T X\hat{\beta} + \hat{\beta}^T X^T X\hat{\beta} \\&= y^T y - 2\hat{\beta}^T X^T y + \hat{\beta}^T X^T X\hat{\beta}\end{aligned}$$

OLS:

$$\begin{aligned}\hat{\beta}^{ols} &= \arg \min_{\hat{\beta}} RSS(\hat{\beta}) \\&= \arg \min_{\hat{\beta}} y^T y - 2\hat{\beta}^T X^T y + \hat{\beta}^T X^T X\hat{\beta}\end{aligned}$$

OLS Estimation of MLR

First and second derivatives are

$$\frac{\partial RSS(\hat{\beta})}{\partial \beta} = -2X^T y + 2X^T X \hat{\beta} \quad \text{and} \quad \frac{\partial^2 RSS(\hat{\beta})}{\partial \beta \partial \beta^T} = 2X^T X$$

$$\text{FOC: } \left. \frac{\partial SSR(\hat{\beta})}{\partial \beta} \right|_{\hat{\beta}^{ols}} = -2X^T y + 2X^T X \hat{\beta}^{ols} = 0_K \text{ which implies}$$

$$\hat{\beta}^{ols} = (X^T X)^{-1} X^T y$$

- $Xc \neq 0 \Leftrightarrow c \neq 0$ means that $X^T X$ is non-singular, and
- $c^T X^T X c = (Xc)^T (Xc) > 0$ (Hessian is pos. def.) so $\hat{\beta}^{ols}$ minimizes $SSR(\hat{\beta})$

OLS Estimation of MLR

Another way of expressing $\hat{\beta}^{ols}$

Writing

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,K-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,K-1} \end{bmatrix} = \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix}$$

which emphasizes observations.

Note X_{i*} is $1 \times K$ vector comprising i th observation of all variables (including '1' for the intercept term)

OLS Estimation of MLR

Then

$$\begin{aligned}
 \hat{\beta}^{ols} &= (X^T X)^{-1} X^T y \\
 &= \left\{ \begin{bmatrix} X_{1*}^T & X_{2*}^T & \dots & X_{n*}^T \end{bmatrix} \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix} \right\}^{-1} \begin{bmatrix} X_{1*}^T & X_{2*}^T & \dots & X_{n*}^T \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \\
 &= \underbrace{\left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}}_{\text{sum of } k \times k \text{ matrices}} \underbrace{\sum_{i=1}^n X_{i*}^T Y_i}_{\text{sum of } k \times 1 \text{ vectors}}
 \end{aligned}$$

OLS Estimation of MLR (Quick Summary)

- MLR: $Y = X\beta + \epsilon$, $E(\epsilon | X) = 0$, $Var(\epsilon | X) = E(\epsilon\epsilon^T | X) = \sigma^2 I_n$
- OLS Estimator for β : $\hat{\beta}^{ols} = (X^T X)^{-1} X^T y$
- OLS Fitted values:

$$\hat{y}^{ols} = X\hat{\beta}^{ols} = X(X^T X)^{-1} X^T y = Py \quad \text{where } P = X(X^T X)^{-1} X^T$$

- OLS Residuals:

$$\begin{aligned}\hat{\epsilon}^{ols} &= y - \hat{y}^{ols} = y - X\hat{\beta}^{ols} = y - X(X^T X)^{-1} X^T y \\ &= (I_n - X(X^T X)^{-1} X^T)y = My \quad \text{where } M = I_n - X(X^T X)^{-1} X^T\end{aligned}$$

Where helpful to do so, I will place “ols” marker in subscript instead of superscript

OLS Estimation of MLR (More on the FOC)

The FOC can be written as: $X^T \hat{\epsilon}^{ols} = 0_{K \times 1}$

$$\text{Since } X = \begin{bmatrix} 1 & X_{11} & X_{12} & \cdots & X_{1,K-1} \\ 1 & X_{21} & X_{22} & \cdots & X_{2,K-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n1} & X_{n2} & \cdots & X_{n,K-1} \end{bmatrix} = [i_n \quad X_{*1} \quad X_{*2} \quad \cdots \quad X_{*,K-1}]$$

$$\text{the FOC says: } X^T \hat{\epsilon}_{ols} = \begin{bmatrix} i_n^T \\ X_{*1}^T \\ X_{*2}^T \\ \vdots \\ X_{*,K-1}^T \end{bmatrix} \hat{\epsilon}^{ols} = \begin{bmatrix} i_n^T \hat{\epsilon}^{ols} \\ X_{*1}^T \hat{\epsilon}^{ols} \\ X_{*2}^T \hat{\epsilon}^{ols} \\ \vdots \\ X_{*,K-1}^T \hat{\epsilon}^{ols} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n \hat{\epsilon}_i^{ols} \\ \sum_{i=1}^n X_{i1} \hat{\epsilon}_i^{ols} \\ \sum_{i=1}^n X_{i2} \hat{\epsilon}_i^{ols} \\ \vdots \\ \sum_{i=1}^n X_{i,K-1} \hat{\epsilon}_i^{ols} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

OLS Estimation of MLR (More on the FOC)

The first line $\sum_{i=1}^n \hat{\epsilon}_i^{ols} = 0$ means that $\overline{\hat{\epsilon}^{ols}} = 0$ so

$$\sum_{i=1}^n X_{ij} \hat{\epsilon}_i^{ols} = 0 \iff \text{smp. cov.}(X_{ij}, \hat{\epsilon}_i^{ols}) = 0$$

for each variable X_j , $j = 1, \dots, K - 1$

OLS estimator $\hat{\beta}^{ols}$ makes OLS residuals

- mean zero and
- uncorrelated with each regressor

Goodness of Fit

$TSS = FSS + RSS$ decomposition continues to hold in the general MLR case

Let $M_0 = I_n - i_n(i_n^T i_n)^{-1} i_n^T = I_n - (1/n) i_n i_n^T$ which is symmetric and idempotent

- $M_0^T = M_0$ and $M_0 M_0 = M_0$

We have

- $M_0 y = \begin{bmatrix} Y_1 - \bar{Y} \\ Y_2 - \bar{Y} \\ \vdots \\ Y_n - \bar{Y} \end{bmatrix}$ and $y^T M_0^T M_0 y = y^T M_0 y = \sum_{i=1}^n (Y_i - \bar{Y})^2$

- $M_0 \hat{\varepsilon}^{ols} = \hat{\varepsilon}^{ols}$

Goodness of Fit

Therefore

$$y = \hat{y}_{ols} + \hat{\varepsilon}_{ols}$$

$$M_0 y = M_0 \hat{y}_{ols} + \hat{\varepsilon}_{ols}$$

$$y^T M_0^T M_0 y = (M_0 \hat{y}_{ols} + \hat{\varepsilon}_{ols})^T (M_0 \hat{y}_{ols} + \hat{\varepsilon}_{ols})$$

$$y^T M_0 y = \hat{y}_{ols}^T M_0 \hat{y}_{ols} + \hat{\varepsilon}_{ols}^T \hat{y}_{ols} + \hat{y}_{ols}^T \hat{\varepsilon}_{ols} + \hat{\varepsilon}_{ols}^T \hat{\varepsilon}_{ols}$$

$$y^T M_0 y = \hat{y}_{ols}^T M_0 \hat{y}_{ols} + \hat{\varepsilon}_{ols}^T \hat{\varepsilon}_{ols}$$

$$TSS = FSS + RSS$$

Goodness of Fit

We can use this to define the goodness-of-fit measure:

$$R^2 = 1 - \frac{RSS}{TSS}$$

and “Adjusted R^2 ”

$$\text{Adj.-}R^2 = 1 - \frac{\frac{1}{n-K}RSS}{\frac{1}{n-1}TSS} = 1 - \frac{RSS}{TSS} \frac{n-1}{n-K}$$

OLS Estimator Properties: Unbiasedness

- OLS estimator $\hat{\beta}^{ols}$ is unbiased:

$$\hat{\beta}^{ols} = (X^T X)^{-1} X^T y = (X^T X)^{-1} X^T (X\beta + \varepsilon) = \beta + (X^T X)^{-1} X^T \varepsilon$$

$$E(\hat{\beta}^{ols} | X) = \beta + E((X^T X)^{-1} X^T \varepsilon | X)$$

$$= \beta + (X^T X)^{-1} X^T E(\varepsilon | X)$$

$$= \beta$$

$$E(\hat{\beta}^{ols}) = \beta$$

OLS Estimator Variance (Homoskedasticity Case)

Under homoskedasticity $Var(\varepsilon\varepsilon^T | X) = \sigma^2 I_n$, we have

$$\begin{aligned}Var(\hat{\beta}^{ols} | X) &= Var(\beta + (X^T X)^{-1} X^T \varepsilon | X) \\&= (X^T X)^{-1} X^T Var(\varepsilon | X) X (X^T X)^{-1} \\&= (X^T X)^{-1} X^T (\sigma^2 I_n) X (X^T X)^{-1} \\&= \sigma^2 (X^T X)^{-1} X^T I_n X (X^T X)^{-1} \\&= \sigma^2 (X^T X)^{-1}\end{aligned}$$

but must estimate σ^2

OLS Estimator Variance (Homoskedasticity Case)

Unbiased estimator of σ^2 is $\widehat{\sigma}^2 = \frac{1}{n-K} \sum_{i=1}^n \hat{\epsilon}_{i,ols}^2 = \frac{\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols}}{n-K}$

Proof: $\hat{\epsilon}_{ols} = My$ where $M = I_n - X(X^T X)^{-1} X^T$

note that

- $MX = (I_n - X(X^T X)^{-1} X^T)X = X - X = 0_{n \times K}$
- M is symmetric (Exercise!)
- M is idempotent, i.e., $MM = M$ (Exercise!)

Therefore $\hat{\epsilon}_{ols} = My = M(X\beta + \varepsilon) = M\varepsilon$ and

$$\hat{\epsilon}_{ols}^T \hat{\epsilon}_{ols} = (M\varepsilon)^T M\varepsilon = \varepsilon^T M^T M\varepsilon = \varepsilon^T MM\varepsilon = \varepsilon^T M\varepsilon$$

OLS Estimator Variance (Homoskedasticity Case)

$$\begin{aligned} E(\hat{\varepsilon}_{ols}^T \hat{\varepsilon}_{ols} | X) &= E(\varepsilon^T M \varepsilon | X) \\ &= E(\text{trace}(\varepsilon^T M \varepsilon) | X) \quad \text{because } \varepsilon^T M \varepsilon \text{ is a scalar} \\ &= E(\text{trace}(\varepsilon \varepsilon^T M) | X) \\ &= \text{trace}(E(\varepsilon \varepsilon^T M) | X) \quad \text{since trace is a sum} \\ &= \text{trace}(E(\varepsilon \varepsilon^T | X) M) = \text{trace}(\sigma^2 I_n M) = \sigma^2 \text{trace}(M) \\ &= \sigma^2 \text{trace}(I_n - X(X^T X)^{-1} X^T) = \sigma^2 (\text{trace}(I_n) - \text{trace}(X(X^T X)^{-1} X^T)) \\ &= \sigma^2 (n - \text{trace}((X^T X)^{-1} X^T X)) = \sigma^2 (n - \text{trace}(I_K)) \\ &= \sigma^2 (n - K) \end{aligned}$$

Therefore $E\left(\frac{\hat{\varepsilon}_{ols}^T \hat{\varepsilon}_{ols}}{n - K} \mid X\right) = \sigma^2$, which implies $E\left(\frac{\hat{\varepsilon}_{ols}^T \hat{\varepsilon}_{ols}}{n - K}\right) = \sigma^2$

OLS Estimator Variance (Heteroskedasticity Case)

$$\begin{aligned} \text{Var}(\hat{\beta}^{ols} | X) &= (X^T X)^{-1} X^T \text{Var}(\varepsilon | X) X (X^T X)^{-1} \\ &= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \begin{bmatrix} X_{1*}^T & X_{2*}^T & \dots & X_{n*}^T \end{bmatrix} \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n^2 \end{bmatrix} \begin{bmatrix} X_{1*} \\ X_{2*} \\ \vdots \\ X_{n*} \end{bmatrix} \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \\ &= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \sigma_i^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \end{aligned}$$

Heteroskedasticity-Robust Estimator Variance-Covariance Matrix: replace σ_i^2 with $\hat{\epsilon}_{i,ols}^2$

$$\widehat{\text{Var}}_{HCO}(\hat{\beta}^{ols}) = \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \hat{\epsilon}_{i,ols}^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

OLS (Best Linear Unbiased, under Homoskedasticity)

We have shown that $\hat{\beta}^{ols}$ is unbiased

It is also “linear”

- A linear estimator is one with the form $\tilde{\beta} = Ay$
- Each $\tilde{\beta}_j = \sum_{i=1}^n a_{ji} Y_i$
- OLS estimator is $\hat{\beta}^{ols} = \underbrace{(X^T X)^{-1} X^T}_A y$

Now we show: OLS estimators are BLU (they have the smallest variance among all linear unbiased estimators)

OLS (Best Linear Unbiased, under Homoskedasticity)

BLU in the sense that

$$\text{Var}(c^T \hat{\beta}^{ols} | X) \leq \text{Var}(c^T \tilde{\beta} | X)$$

for all $K \times 1$ vectors c , and for all unbiased estimators of the form $\tilde{\beta} = By$

- each individual $\hat{\beta}_k$ is BLU
- all linear combinations of $\hat{\beta}$ are BLU

Consider predicting Y at the new observation $X_{0*} = [1 \quad X_{01} \quad \dots \quad X_{0,K-1}]$. OLS predictor is

$$\hat{Y}(X_{0*}) = X_{0*} \hat{\beta}^{ols}$$

which is a linear combination of the parameter estimates in $\hat{\beta}^{ols}$, i.e., OLS prediction rule gives us the most precise linear unbiased prediction of Y at X_{0*}

OLS (Best Linear Unbiased, under Homoskedasticity)

Proof of Efficiency: let $\tilde{\beta} = By$ be an unbiased estimator where $B \neq (X^T X)^{-1} X^T$

Let D be such that $B = D + (X^T X)^{-1} X^T$, so that

$$\begin{aligned}\tilde{\beta} &= By = (D + (X^T X)^{-1} X^T)y \\ &= (D + (X^T X)^{-1} X^T)(X\beta + \varepsilon) \\ &= DX\beta + D\varepsilon + \beta + (X^T X)^{-1} X^T \varepsilon\end{aligned}$$

To ensure unbiasedness of $\tilde{\beta}$, we must assume $DX = 0$

OLS (Best Linear Unbiased, under Homoskedasticity)

Then

$$\begin{aligned} \text{Var}(\tilde{\beta} | X) &= E((\tilde{\beta} - \beta)(\tilde{\beta} - \beta)^T | X) \\ &= E((D + (X^T X)^{-1} X^T) \varepsilon \varepsilon^T (D + (X^T X)^{-1} X^T)^T | X) \\ &= (D + (X^T X)^{-1} X^T) E(\varepsilon \varepsilon^T | X) (D + (X^T X)^{-1} X^T)^T \\ &= \sigma^2 (D + (X^T X)^{-1} X^T) (D + (X^T X)^{-1} X^T)^T \\ &= \sigma^2 (D + (X^T X)^{-1} X^T) (D^T + X (X^T X)^{-1}) \\ &= \sigma^2 [DD^T + (X^T X)^{-1} X^T D^T + DX (X^T X)^{-1} + (X^T X)^{-1}] \\ &= \sigma^2 [DD^T + (X^T X)^{-1}] = \sigma^2 DD^T + \sigma^2 (X^T X)^{-1} = \sigma^2 DD^T + \text{Var}(\hat{\beta} | X) \end{aligned}$$

OLS (Best Linear Unbiased, under Homoskedasticity)

Therefore

$$\begin{aligned} \text{Var}(c^T \tilde{\beta} \mid X) &= c^T \text{Var}(\tilde{\beta} \mid X) c \\ &= c^T (\sigma^2 D D^T + \text{Var}(\hat{\beta} \mid X)) c \\ &= \sigma^2 c^T D D^T c + c^T \text{Var}(\hat{\beta} \mid X) c \\ &= \sigma^2 (D^T c)^T D^T c + \text{Var}(c^T \hat{\beta} \mid X) \geq \text{Var}(c^T \hat{\beta} \mid X) \end{aligned}$$

NB: $D^T c$ is a vector, therefore $(D^T c)^T D^T c$ is a sum of squares, which cannot be negative

Example

```
mlr <- function(y, X){
  n <- dim(X)[1]
  K <- dim(X)[2]
  XTXinv <- solve(t(X)%*%X)
  betahat <- XTXinv %*% t(X)%*%y
  yhat <- X %*% betahat
  ehat <- y - yhat
  sigmasqhat <- sum(ehat^2)/(dim(X)[1]-dim(X)[2])
  betahatvar <- sigmasqhat * XTXinv
  betahatse <- sqrt(diag(betahatvar))
  betahat_t <- betahat/betahatse
  results <- cbind(
    betahat, betahatse, betahat_t, 2*pt(-abs(betahat_t), n-K)
  )
  colnames(results) <- c("coef.", "s.e.", "t-stat", "p-value")
  model_return <- list("results"=results, "sigmasqhat"=sigmasqhat, "betahatvar"=betahatvar)
  return(model_return)
}
```

Example

```
y = log(dat1$earn)
X = cbind("intercept"=1, "educ"=dat1$educ, "tenure"=dat1$tenure,
         "age"=dat1$age, "agesq"=dat1$age^2)
model2 <- mlr(y,X)
model2$results %>% round(5)
```

	coef.	s.e.	t-stat	p-value
intercept	-0.14085	0.11333	-1.24284	0.21399
educ	0.12590	0.00379	33.21175	0.00000
tenure	0.01545	0.00108	14.26744	0.00000
age	0.06251	0.00478	13.08669	0.00000
agesq	-0.00067	0.00005	-12.76758	0.00000

Example

Using the `lm()` function

```
model2lm <- lm(log(earn)~educ + tenure + age + I(age^2), data=dat1)
model2lm %>% summary %>% coefficients %>% round(5)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.14085	0.11333	-1.24284	0.21399
educ	0.12590	0.00379	33.21175	0.00000
tenure	0.01545	0.00108	14.26744	0.00000
age	0.06251	0.00478	13.08669	0.00000
I(age^2)	-0.00067	0.00005	-12.76758	0.00000

Example

Estimated var-covariance matrix

```
vcov(model2lm)
```

	(Intercept)	educ	tenure	age
(Intercept)	1.284416e-02	-1.880320e-04	1.206658e-05	-4.679158e-04
educ	-1.880320e-04	1.436923e-05	-2.315549e-08	-9.040899e-07
tenure	1.206658e-05	-2.315549e-08	1.172805e-06	-6.690121e-07
age	-4.679158e-04	-9.040899e-07	-6.690121e-07	2.281319e-05
I(age^2)	4.957007e-06	1.083268e-08	2.934586e-09	-2.471022e-07
	I(age^2)			
(Intercept)	4.957007e-06			
educ	1.083268e-08			
tenure	2.934586e-09			
age	-2.471022e-07			
I(age^2)	2.747833e-09			

Hypothesis Testing

A general single linear hypothesis can be written as

$$H_0 : r^T \beta = r_0 \quad \text{vs} \quad r^T \beta \neq r_0.$$

Example: To test $\beta_1 + \beta_2 = 1$ in the regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i,$$

set $r^T = [0 \quad 1 \quad 1]$ and $r_0 = 1$.

Hypothesis Testing

We have

$$r^T \hat{\beta} \mid X \sim \text{Normal}(r^T \beta, r^T (\sigma^2 (X^T X)^{-1}) r).$$

If the null hypothesis $r^T \beta = r_0$ holds, then

$$r^T \hat{\beta} \mid X \sim \text{Normal}(r_0, r^T (\sigma^2 (X^T X)^{-1}) r)$$

and

$$\frac{r^T \hat{\beta} - r_0}{\sqrt{r^T (\sigma^2 (X^T X)^{-1}) r}} \sim \text{Normal}(0, 1).$$

Hypothesis Testing

Furthermore, it can be shown that if we replace σ^2 with $\widehat{\sigma}^2$, then

$$\begin{aligned} t &= \frac{r^T \widehat{\beta} - r_0}{\sqrt{r^T (\widehat{\sigma}^2 (X^T X)^{-1}) r}} \\ &= \frac{r^T \widehat{\beta} - r_0}{\sqrt{r^T \widehat{Var}(\widehat{\beta} | X) r}} \sim t(n - K). \end{aligned}$$

This can be used to test the hypothesis $H_0 : r^T \beta = r_0$ in the usual way.

Hypothesis Testing

To test multiple hypotheses jointly, write the hypotheses as

$$H_0 : \mathcal{R}\beta = r_0 \quad \text{vs} \quad H_A : \mathcal{R}\beta \neq r_0$$

where now \mathcal{R} is a $(J \times K)$ matrix, and r_0 is a $(J \times 1)$ vector.

To test the hypotheses

$$H_0 : \beta_1 + \beta_2 = 1 \text{ and } \beta_3 = 0 \quad \text{vs} \quad H_A : \beta_1 + \beta_2 \neq 1 \text{ or } \beta_3 \neq 0 \text{ (or both),}$$

in the regression

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \varepsilon_i$$

set the matrices \mathcal{R} and r_0 to

$$\mathcal{R} = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad r_0 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}.$$

Hypothesis Testing

It turns out that in practice, one does not actually have to compute the restricted regression. It can be shown that

$$\hat{\varepsilon}_{rls}^T \hat{\varepsilon}_{rls} - \hat{\varepsilon}^T \hat{\varepsilon} = (\mathcal{R}\hat{\beta} - r_0)^T (\mathcal{R}(X^T X)^{-1} \mathcal{R}^T)^{-1} (\mathcal{R}\hat{\beta} - r_0)$$

where $\hat{\beta}$ is the unrestricted OLS estimators (for proof, see Econometrics Notes)

Furthermore, denominator of F -statistic is $\widehat{\sigma^2}$, therefore

$$\begin{aligned} F &= (\mathcal{R}\hat{\beta} - r_0)^T (\mathcal{R}(\widehat{\sigma^2}(X^T X)^{-1})\mathcal{R}^T)^{-1} (\mathcal{R}\hat{\beta} - r_0) / J \\ &= (\mathcal{R}\hat{\beta} - r_0)^T (\mathcal{R} \widehat{Var}(\hat{\beta} | X) \mathcal{R}^T)^{-1} (\mathcal{R}\hat{\beta} - r_0) / J \\ &\sim F(J, n - K) \end{aligned}$$

Hypothesis Testing: Example

We use the following example

$$\ln \text{earn} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{black} + \beta_3 \text{female} + \beta_4 \text{black} \cdot \text{female} + \epsilon$$

Use t-test to test: $H_0 : \beta_2 = \beta_3$ vs $H_A : \beta_2 \neq \beta_3$

Use F-test to test: $H_0 : \beta_2 = \beta_3$ and $\beta_4 = 0$ versus $H_0 : \beta_2 \neq \beta_3$ or $\beta_4 \neq 0$

```
library(car)
dat2<-read_csv("data\\earnings2019.csv",show_col_types=FALSE) %>%
  mutate(female=1-male,
         white=if_else(race=="White", 1,0),
         black=if_else(race=="Black", 1,0),
         other=if_else(race=="Other", 1,0)) %>% select(-race)
mdl_unres <- lm(log(earn) ~ educ + black*female, data=dat2)
```

Hypothesis Testing: Example

```
mdl_unres %>% summary %>% coefficients %>% round(5)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.53686	0.05709	26.92194	0.00000
educ	0.12786	0.00388	32.95769	0.00000
black	-0.26005	0.02696	-9.64677	0.00000
female	-0.28066	0.01959	-14.32381	0.00000
black:female	0.08730	0.03549	2.46010	0.01392

Hypothesis Testing: Example

```
linearHypothesis mdl_unres, c("black=female", "black:female=0"))
```

Linear hypothesis test:

black - female = 0

black:female = 0

Model 1: restricted model

Model 2: log(earn) ~ educ + black * female

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	4943	1606.5				
2	4941	1603.5	2	2.9507	4.5461	0.01065 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Hypothesis Testing: Example

For illustration, we compute the two tests using the formulas derived in this section

```
y = log(dat2$earn)
X = cbind("intercept"=1, "educ"=dat2$educ, "black"=dat2$black, "female"=dat2$female,
         "black.female" = dat2$black * dat2$female)
n = dim(X)[1]
K = dim(X)[2]
model1 <- mlr(y,X)
model1$results %>% round(5)
```

	coef.	s.e.	t-stat	p-value
intercept	1.53686	0.05709	26.92194	0.00000
educ	0.12786	0.00388	32.95769	0.00000
black	-0.26005	0.02696	-9.64677	0.00000
female	-0.28066	0.01959	-14.32381	0.00000
black.female	0.08730	0.03549	2.46010	0.01392

Hypothesis Testing: Example

```
# t-test of "black=female"
r = matrix(c(0,0,1,-1,0), ncol=1)
r0 = 0
betahat <- as.matrix(model1$results[,1])
tstat <- (t(r) %*% betahat - r0)/sqrt(t(r) %*% model1$betahatvar %*% r)
tstat_pval <- 2*pt(-abs(tstat), n-K)
cat("\n Test: b2=b3")
cat("\n tstat:", round(tstat,5), "    p-value:", round(tstat_pval,5), "\n")
```

Test: b2=b3

tstat: 0.76067 p-value: 0.44689

Hypothesis Testing: Example

```
# F-test of "black=female" and "black.female=0"
R <- matrix(c(0,0,1,-1,0,
              0,0,0, 0,1), nrow=2, byrow=TRUE)
J <- dim(R)[1]
r0 <- matrix(c(0,0), nrow = 2)
Fstat <- t(R**%betahat - r0) **%
        solve(R **% model1$betahatvar **% t(R)) **%
        (R **% betahat - r0)/J
Fstat_pval <- 1-pf(Fstat, J, n-K)
cat("\n Test: b2=b3, b4=0")
cat("\n Fstat:", round(Fstat,5), "   p-value:", round(Fstat_pval,5),"\n")
```

Test: b2=b3, b4=0

Fstat: 4.54611 p-value: 0.01065

Asymptotic Properties of OLS Estimators

Multivariate versions of LLN and CLT:

If $\{Z_i\}_{i=1}^n$ iid vectors of random variables with $E(Z_i) = 0$ and $Var(Z_i) = \Omega$, for all i , then

$$\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i = \begin{bmatrix} (1/n) \sum_{i=1}^n Z_{1i} \\ (1/n) \sum_{i=1}^n Z_{2i} \\ \vdots \\ (1/n) \sum_{i=1}^n Z_{ki} \end{bmatrix} \xrightarrow{p} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_k \end{bmatrix} = \mu$$

$$\sqrt{n} \bar{Z} = \frac{1}{\sqrt{n}} \sum_{i=1}^n Z_i = \begin{bmatrix} (1/\sqrt{n}) \sum_{i=1}^n Z_{1i} \\ (1/\sqrt{n}) \sum_{i=1}^n Z_{2i} \\ \vdots \\ (1/\sqrt{n}) \sum_{i=1}^n Z_{ki} \end{bmatrix} \xrightarrow{d} \text{Normal}_k(0, \Omega)$$

Asymptotic Properties

To talk about limiting distributions, we have to scale $\hat{\beta}$. Use

$$\sqrt{n}(\hat{\beta}^{ols} - \beta) = \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*}^T \epsilon_i \right)$$

Our assumptions and the CLT imply

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*} \epsilon_i \xrightarrow{d} \text{Normal}_{k+1}(0, S)$$

therefore

$$\sqrt{n}(\hat{\beta}^{ols} - \beta) = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}}_{\xrightarrow{p} \Sigma_{xx}^{-1}} \underbrace{\left(\frac{1}{\sqrt{n}} \sum_{i=1}^n X_{i*}^T \epsilon_i \right)}_{\xrightarrow{p} \text{Normal}_{k+1}(0, S)} \xrightarrow{d} \text{Normal}_{k+1}(0, \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1})$$

Asymptotic Properties

That is, $\hat{\beta}^{ols}$ is consistent, with asymptotic variance $Avar(\hat{\beta}^{ols}) = \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1}$. This result justifies the approximation

$$Var(\hat{\beta}^{ols}) \approx (1/n) \Sigma_{xx}^{-1} S \Sigma_{xx}^{-1}.$$

An obvious estimator for Σ_{xx} is

$$\hat{\Sigma}_{xx} = \frac{1}{n} \sum_{i=1}^N X_{i*}^T X_{i*} = \frac{1}{n} X^T X$$

Some additional assumptions (see advanced econometrics textbooks) guarantee

$$\hat{S} = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \xrightarrow{p} S.$$

Asymptotic Properties

This allows us to consistently estimate the asymptotic variance of $\hat{\beta}$ by

$$\widehat{Avar}(\hat{\beta}^{ols}) = \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1}$$

and justifies the use of

$$\begin{aligned} \widehat{Var}_{HCO}(\hat{\beta}^{ols}) &= \frac{1}{n} \widehat{Avar}(\hat{\beta}^{ols}) \\ &= \frac{1}{n} \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\frac{1}{n} \sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \\ &= \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \left(\sum_{i=1}^n \hat{\epsilon}_i^2 X_{i*}^T X_{i*} \right) \left(\sum_{i=1}^n X_{i*}^T X_{i*} \right)^{-1} \end{aligned}$$

Example

We adapt previously written `mlr` function to give HC0 standard errors:

```
mlr_hc0 <- function(y, X){
  n <- dim(X)[1]; K <- dim(X)[2]
  XTXinv <- solve(t(X)%*%X)
  betahat <- XTXinv %*% t(X)%*%y; yhat <- X %*% betahat; ehat <- y - yhat
  sigmasqhat <- sum(ehat^2)/(dim(X)[1]-dim(X)[2])
  hatS = 0
  for (i in 1:n){
    xi = as.matrix(X[i,])
    hatS = hatS + ehat[i]^2*xi%*%t(xi)
  }
  betahatvar_hc0 <- XTXinv %*% hatS %*% XTXinv ## XTXinv computed earlier
  betahatse_hc0 <- sqrt(diag(betahatvar_hc0))
  betahat_t_hc0 <- betahat/betahatse_hc0
  results <- cbind(betahat, betahatse_hc0, betahat_t_hc0, 2*pnorm(-abs(betahat_t_hc0), 0, 1))
  colnames(results) <- c("coef.", "s.e.(hc0)", "t-stat", "p-value")
  model_return <- list("results"=results, "betahatvar"=betahatvar_hc0)
  return(model_return)
}
```

Example

$$\ln \text{earn} = \beta_0 + \beta_1 \text{educ} + \beta_2 \text{black} + \beta_3 \text{female} + \beta_4 \text{black.female} + \epsilon$$

```
y = log(dat1$earn)
X = cbind("intercept"=1, "educ"=dat2$educ, "black"=dat2$black, "female"=dat2$female,
         "black.female"=dat2$black * dat2$female)
model2_hc0 <- mlr_hc0(y,X)
cat("Estimation results:\n"); model2_hc0$results %>% round(6)
```

Estimation results:

	coef.	s.e.(hc0)	t-stat	p-value
intercept	1.536856	0.056516	27.193109	0.00000
educ	0.127858	0.003927	32.555323	0.00000
black	-0.260050	0.026926	-9.657793	0.00000
female	-0.280661	0.020170	-13.915056	0.00000
black.female	0.087299	0.034640	2.520150	0.01173

Example

- `vcovHC()` function from the `sandwich` package.

```
model3 <- lm(log(earn) ~ educ + black*female, data=dat2)
betahatvar_HCO_sando = sandwich::vcovHC(model3,type="HCO")
round(betahatvar_HCO_sando,6)
```

	(Intercept)	educ	black	female	black:female
(Intercept)	0.003194	-0.000214	-0.000290	-0.000070	0.000188
educ	-0.000214	0.000015	0.000005	-0.000011	0.000003
black	-0.000290	0.000005	0.000725	0.000224	-0.000724
female	-0.000070	-0.000011	0.000224	0.000407	-0.000400
black:female	0.000188	0.000003	-0.000724	-0.000400	0.001200

Example

```
linearHypothesis(model3, c("black=female", "black:female=0"),  
                  vcov=betahatvar_HCO_sando, test="Chisq")
```

Linear hypothesis test:

black - female = 0

black:female = 0

Model 1: restricted model

Model 2: $\log(\text{earn}) \sim \text{educ} + \text{black} * \text{female}$

Note: Coefficient covariance matrix supplied.

	Res.Df	Df	Chisq	Pr(>Chisq)
1	4943			
2	4941	2	9.6196	0.00815 **