

Why do we run regressions?

Let's say we have two variables y and x . We may regress y on x , because

- ① we want to see the linear association between y and x , OR
 - ② we want to predict y given values of x , OR
 - ③ we want to measure the **causal effect** of x on y .
- What do we mean by the **causal effect**?
 - We will use the definition of SW book:

*A causal effect is defined to be the effect measured in an **ideal randomized controlled experiment (IRCE)**.*

IRCE

- **Ideal:** subjects all follow the treatment protocol, perfect compliance and no errors in reporting.
- **Randomized:** subjects from the population of interest are randomly assigned to a treatment or control group (no confounding factors).
- **Controlled:** having a control group permits measuring the differential effect of the treatment.
- **Experiment:** the treatment is assigned as part of the experiment (no *reverse causality* in which subjects choose the treatment they think will work best).

class size effect on test scores

- Recall the class size example: causal effect of reducing class size (STR) on test scores.
- What would an IRCE be for measuring the effect on Test Score of reducing STR?
 - Students would be randomly assigned to classes, which would have different sizes.
 - Because students are randomly assigned, all student characteristics (ε) would be distributed independently of STR.
 - The **zero conditional mean** assumption holds trivially in an IRCE, i.e.,
 $\mathbb{E}(\varepsilon|STR) = 0$.
- **Observational data** are not collected from IRCEs. They often come from surveys.
- The **zero conditional mean** assumption almost certainly does not hold, resulting in **omitted variables bias**.

Setting

- The linear regression model is given by

$$y_i = \beta_0 \cdot 1 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i \quad (1)$$

for $i = 1, 2, \dots, n$, where $(\beta_0, \dots, \beta_k)'$ are the (unknown) parameters.

- Going back the example, *STR* can be one of the x 's.
- What are the roles of other x 's?
- They are the so-called **control variables**: they control for the omitted variables bias.
- Suppose x_{1i} is the *STR* variable. Then, the role of the control variables can be represented as:

$$\mathbb{E}(\varepsilon_i | STR_i, x_{2i}, \dots, x_{ki}) = \mathbb{E}(\varepsilon_i | x_{2i}, \dots, x_{ki}),$$

for $i = 1, 2, \dots, n$, i.e., **the conditional mean independence** assumption.

Setting

- Hence, controlling for x_2, \dots, x_k , STR does not covary with all student characteristics (as if students were randomly assigned).
- The (almost) causal effect of STR on test scores then can be obtained.
- We will see that if any of the control variables is missing from the specification, the least squares estimator for the slope of STR will generally be biased.
- More specifically, we will see that for the omitted variable bias to occur, the omitted variable(s) will have to covary with STR.

Setting

- The model can be written more compactly

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad (2)$$

for $i = 1, 2, \dots, n$, where $\mathbf{x}_i = (1, x_{1i}, \dots, x_{ki})'$ and $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_k)'$ are the (unknown) parameters.

- Stacking over i 's we can equivalently write

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (3)$$

where $\mathbf{y} = (y_1, \dots, y_n)'$, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$, and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$

The least-squares methodology

- Let $\mathbf{z}_i = (y_i, \mathbf{x}_i')'$ be a random vector, i.e., each element is a random variable.
- A random variable is a mapping, $x(\omega) : \Omega \mapsto \mathbb{R}$.
- As ω varies over Ω , we obtain realizations $x(\omega)$ ranging over \mathbb{R} .
- The terminology random variable is a bit unfortunate, because a random variable is neither random nor a variable.
- x is not a variable, it is a real valued function.
- x is not random, it is fixed. But, $\omega \in \Omega$ is random.
- The realized sample is one of the possible realizations. There is an underlying uncertainty in $x(\omega)$ and $y(\omega)$, and anything derived from them will inherit that uncertainty.

The least-squares methodology

- From hereon, let's denote the unknown true population value of the parameters by $\{\beta'_0, \sigma_0^2\}$ and arbitrary values by $\{\beta', \sigma^2\}$.
- The least-squares methodology tries to find the best value for the unknown β_0 so that the difference between the realized values of the outcome variable and the guesses from the model is the smallest in some sense.
- Let $\hat{\beta}$ the value picked. With this, we can compute the difference between the realized value of the outcome variable and the guess from the model (fitted value) for each i . Call this difference the **residuals**.
- We cannot simply add the residuals over the i 's to find the best value for β , because negatives and positives will negate.
- Instead, the least-squares minimizes the **sum of the squared residuals** to find the best approximation to β_0 .
- The **least-squares estimator** solves

$$\arg \min_{\beta} \sum_{i=1}^n (y_i - x_i' \beta)^2 = \arg \min_{\beta} (\mathbf{y} - \mathbf{X}\beta)' (\mathbf{y} - \mathbf{X}\beta). \quad (4)$$

The least-squares methodology

- Using some matrix algebra, finding the solution is not difficult.
- Let $S(\beta)$ denote $(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta)$.
- The first-order condition for a minimum is

$$\frac{\partial S(\beta)}{\partial \beta} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\beta \stackrel{set}{=} 0. \quad (5)$$

- Then, the solution from (5) follows given \mathbf{X} is full column rank, i.e., $\mathbf{X}'\mathbf{X}$ is invertible,

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \quad (6)$$

- To see that this solution is indeed a minimum, it suffices to note that $S(\beta)$ is convex.

The least-squares methodology

- Equivalently, we can compute the second-order derivatives. The second-order condition for a minimum is

$$\frac{\partial^2 \mathcal{S}(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} = 2\mathbf{X}'\mathbf{X} \quad (7)$$

must be positive definite.

- To see, $\mathbf{X}'\mathbf{X}$ is positive definite, let \mathbf{c} be an arbitrary $(k+1) \times 1$ nonzero vector and let $\mathbf{v} = \mathbf{X}\mathbf{c}$.
- Then,

$$\mathbf{c}'\mathbf{X}'\mathbf{X}\mathbf{c} = \mathbf{v}'\mathbf{v} = \sum_{i=1}^n v_i^2. \quad (8)$$

- (8) equals zero only if $v_i = 0$ for $i = 1, \dots, n$. It can happen only if the columns of \mathbf{X} linearly dependent.
- This is not allowed if \mathbf{X} is full column rank.

Properties

- \mathbf{X} is orthogonal to the residuals by construction, i.e., the inner product of every column of \mathbf{X} and $\hat{\boldsymbol{\varepsilon}}$ is zero (they form a right angle):

$$\begin{aligned}\mathbf{X}'\hat{\boldsymbol{\varepsilon}} &= \mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (9) \\ &= \mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{y} = \mathbf{0}.\end{aligned}$$

- Recall that the first column in \mathbf{X} is a column of 1's. Let \mathbf{l} denote it. The least-squares residuals sum to zero. From (9),

$$\mathbf{X}'\hat{\boldsymbol{\varepsilon}} = (\mathbf{l}, \mathbf{x}_1, \dots, \mathbf{x}_k)' \hat{\boldsymbol{\varepsilon}} = \begin{pmatrix} \mathbf{l}'\hat{\boldsymbol{\varepsilon}} \\ \mathbf{x}_1'\hat{\boldsymbol{\varepsilon}} \\ \vdots \\ \mathbf{x}_k'\hat{\boldsymbol{\varepsilon}} \end{pmatrix} = \mathbf{0}.$$

- Notice that $\mathbf{l}'\hat{\boldsymbol{\varepsilon}}$ is the sum of least squares residuals and equals zero.

Properties

- The regression hyperplane passes through the data means: $\bar{y} = \bar{X}\hat{\beta}$.
- From (9), we have $X' y = X' X \hat{\beta}$. Then,

$$\begin{pmatrix} l' \\ \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_k \end{pmatrix} y = \begin{pmatrix} l' \\ \mathbf{x}'_1 \\ \vdots \\ \mathbf{x}'_k \end{pmatrix} (l, \mathbf{x}_1, \dots, \mathbf{x}_k) \hat{\beta}.$$

- From the first row, $l' y = (l' l, l' x_1, \dots, l' x_k) \hat{\beta}$, hence

$$\bar{y} = \frac{1}{n} l' y = \left(\frac{1}{n} l' l, \frac{1}{n} l' x_1, \dots, \frac{1}{n} l' x_k \right) \hat{\beta} = \bar{X} \hat{\beta}.$$

Properties

- The mean of the fitted values from the regression equals the mean of the actual values: $\bar{\hat{y}} = \bar{y}$. This is easily verified because $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$

$$\begin{aligned}\bar{\hat{y}} &= \frac{1}{n} \mathbf{l}' \mathbf{X} \hat{\boldsymbol{\beta}} = \frac{1}{n} \mathbf{l}' \mathbf{X} \hat{\boldsymbol{\beta}} + \frac{1}{n} \mathbf{l}' \hat{\boldsymbol{\varepsilon}} \\ &= \frac{1}{n} \mathbf{l}' (\mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}) = \frac{1}{n} \mathbf{l}' \mathbf{y} = \bar{y}.\end{aligned}$$

Projection Matrices

- We will define two useful matrices. From the definition of the least squares residual

$$\hat{\varepsilon} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{y} = \mathbb{M}_x\mathbf{y}$$

where $\mathbb{M}_x = (\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$.

- Also, notice that

$$\hat{\mathbf{y}} = \mathbf{y} - \hat{\varepsilon} = \mathbf{y} - \mathbb{M}_x\mathbf{y} = (\mathbf{I}_n - \mathbf{I}_n + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbb{P}_x\mathbf{y}$$

where $\mathbb{P}_x = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$.

- We call \mathbb{M}_x and \mathbb{P}_x respectively as the **orthogonal projection matrix** and the **projection matrix**.
- Both matrices are symmetric and idempotent.

Projection Matrices

■ Using \mathbb{M}_x and \mathbb{P}_x , we have

$$\square \quad \mathbf{y} = \mathbf{X}'\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}} = \mathbb{P}_x\mathbf{y} + \mathbb{M}_x\mathbf{y},$$

$$\square \quad \hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}} = \mathbf{y}'\mathbb{M}_x'\mathbb{M}_x\mathbf{y} = \mathbf{y}'\mathbb{M}_x'\mathbf{y} = \hat{\boldsymbol{\varepsilon}}'\mathbf{y},$$

$$\square \quad \hat{\boldsymbol{\varepsilon}}'\hat{\boldsymbol{\varepsilon}} = (\mathbf{y} - \mathbb{P}_x\mathbf{y})'(\mathbf{y} - \mathbb{P}_x\mathbf{y}) = \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbb{P}_x'\mathbf{y} - \mathbf{y}'\mathbb{P}_x\mathbf{y} + \mathbf{y}'\mathbb{P}_x'\mathbb{P}_x\mathbf{y} = \\ \mathbf{y}'\mathbf{y} - \mathbf{y}'\mathbb{P}_x'\mathbb{P}_x\mathbf{y} = \mathbf{y}'\mathbf{y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}},$$

$$\square \quad \mathbb{P}_x\mathbf{X} = \mathbf{X} \text{ and } \mathbb{P}_x\mathbb{M}_x = \mathbf{0}.$$

Goodness of fit

- A special symmetric idempotent matrix is the following

$$\begin{aligned}\mathbb{M}_0 &= \mathbf{I}_n - \mathbf{l}_n (\mathbf{l}_n' \mathbf{l}_n)^{-1} \mathbf{l}_n' = \mathbf{I}_n - \frac{1}{n} \mathbf{l}_n \mathbf{l}_n' \\ &= \begin{pmatrix} 1 - \frac{1}{n} & -\frac{1}{n} & \cdots & -\frac{1}{n} \\ -\frac{1}{n} & 1 - \frac{1}{n} & \cdots & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & \cdots & 1 - \frac{1}{n} \end{pmatrix}\end{aligned}$$

where \mathbf{l}_n is $n \times 1$ vector of 1's.

- \mathbb{M}_0 is often called the **deviation-from-the-mean** matrix.
- Write the linear regression model in deviation from the mean form:

$$\mathbb{M}_0 \mathbf{y} = \mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbb{M}_0 \hat{\boldsymbol{\varepsilon}} = \mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}.$$

- The second equality follows because $\mathbb{M}_0 \hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\varepsilon}} - \frac{1}{n} \mathbf{l}_n \mathbf{l}_n' \hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\varepsilon}}$ as $\mathbf{l}_n' \hat{\boldsymbol{\varepsilon}} = 0$.

Goodness of fit

- We can write the sum of the squared deviations of the elements of \mathbf{y} from their mean

$$\begin{aligned}(\mathbb{M}_0 \mathbf{y})' \mathbb{M}_0 \mathbf{y} &= \mathbf{y}' \mathbb{M}_0 \mathbf{y} = (\mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}})' (\mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}) \\&= \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\beta}}' \mathbf{X}' \hat{\boldsymbol{\varepsilon}} + \hat{\boldsymbol{\varepsilon}}' \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} \\&= \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbb{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}} \\&= \hat{\mathbf{y}}' \mathbb{M}_0 \hat{\mathbf{y}} + \hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}.\end{aligned}$$

- Recall that $\hat{\bar{y}} = \bar{y}$ and $\hat{\mathbf{y}} - \mathbf{l}_n \hat{\bar{y}} = \hat{\mathbf{y}} - \mathbf{l}_n \bar{y}$ when an intercept is included in \mathbf{X} .
- Let $\mathbf{y}' \mathbb{M}_0 \mathbf{y}$, $\hat{\mathbf{y}}' \mathbb{M}_0 \hat{\mathbf{y}}$ and $\hat{\boldsymbol{\varepsilon}}' \hat{\boldsymbol{\varepsilon}}$ be denoted respectively by total sum of squares (TSS), explained sum of squares (ESS) and residual sum of squares (RSS).
- Clearly, $TSS = ESS + RSS$, where ESS is the portion of the variation in \mathbf{y} explained by the model, and the RSS is the portion of the variation in \mathbf{y} unexplained by the model.

Goodness of fit

- A natural measure of goodness-of-fit is

$$R^2 = \frac{ESS}{TSS} = \frac{\hat{\beta}' \mathbf{X}' \mathbb{M}_0 \mathbf{X} \hat{\beta}}{\mathbf{y}' \mathbb{M}_0 \mathbf{y}} = 1 - \frac{\hat{\epsilon}' \hat{\epsilon}}{\mathbf{y}' \mathbb{M}_0 \mathbf{y}} = 1 - \frac{RSS}{TSS}$$

which is the square of the sample correlation between \mathbf{y} and $\mathbb{P}_x \mathbf{y}$.

- We will refer to this measure as the **coefficient of determination**. It lies between 0 and 1.
- One serious drawback of R^2 is that it never decreases as we include more explanatory variables in \mathbf{X} , even when they are irrelevant (superfluous) in explaining \mathbf{y} .
- Hence, we modify R^2 so that it penalizes for superfluous regressors.
- The modified measure is called the **adjusted** R^2 and is given by

$$\bar{R}^2 = 1 - \frac{\hat{\epsilon}' \hat{\epsilon} / (n - k - 1)}{\mathbf{y}' \mathbb{M}_0 \mathbf{y} / (n - 1)}.$$

Goodness of fit

- As k increases relative to the number of observations the last term increases and \bar{R}^2 falls.
- One can show that

$$\bar{R}^2 = 1 - \left(\frac{n-1}{n-k-1} \right) (1 - R^2).$$

- In fact, \bar{R}^2 increases if and only if the t -statistic of a newly added regressor is greater than one in absolute value and \bar{R}^2 may even get negative values.
- For linear regression models without an intercept, the interpretation of R^2 is not easily available and R^2 should not be used for model comparison.
- Instead, one can use the following measures for model comparison (including nonlinear models):

- Akaike Information Criterion:

$$AIC = \log \left(\frac{\hat{\epsilon}' \hat{\epsilon}}{n} \right) + \frac{2(k+1)}{n},$$

- Bayesian (or Schwartz) Information Criterion:

$$BIC = \log \left(\frac{\hat{\epsilon}' \hat{\epsilon}}{n} \right) + \frac{k \log(n)}{n}.$$

Assumptions

- Characterizing the **finite sample** distribution of the least-squares estimator is feasible, which is often not the case for many other estimators.
- We make the following assumptions:
 - ❶ $\mathbb{E}(\boldsymbol{\varepsilon}|\mathbf{X}) = \mathbf{0}$ a.s.,
 - ❷ $\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X}) = \sigma_0^2 \mathbf{I}_n$ a.s.,
 - ❸ $\text{rank}(\mathbf{X}) = k + 1$ a.s.
- These assumptions correspond to the statistical model of random regressors with independent sampling.
- The first two assumptions mostly rule out models of time-series data.
- By their nature, time series data almost always exhibit dependence of data points.

Remarks

- σ_0^2 is a constant not depending on \mathbf{X} . $\sigma_0^2 \mathbf{I}_n$ is the conditional covariance of ε , but since it does not depend on \mathbf{X} , it is also the unconditional covariance of ε .
- From the first two assumptions, we have $\mathbb{E}(\varepsilon_i) = 0$, $\text{Var}(\varepsilon_i) = \sigma_0^2$ for $i = 1, \dots, n$ and $\text{cov}(\varepsilon_i, \varepsilon_j) = \mathbb{E}(\varepsilon_i \varepsilon_j) = 0$ for $i \neq j$ (unconditionally).
- The independent random sampling provides the uncorrelatedness.
- However, uniformity of the variances is simply of convenience, but otherwise, it is necessary to know what the variances are.
- The last assumption is needed for the existence of the least squares estimator. Failure of this assumption is known as **perfect collinearity**.
- Near failure of this assumption results in $\mathbf{X}'\mathbf{X}$ being poorly conditioned for inversion. It is known as **multicollinearity**.

Unbiasedness

Claim

Assumptions 1–3 imply that $\hat{\beta}$ is unbiased.

Proof.

We can write

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\beta_0 + \varepsilon) = \beta_0 + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon.$$

Hence, the expectation of $\hat{\beta}$ conditional on \mathbf{X} is,

$$\begin{aligned}\mathbb{E}(\hat{\beta}|\mathbf{X}) &= \beta_0 + \mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\varepsilon|\mathbf{X}] \\ &= \beta_0 + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbb{E}(\varepsilon|\mathbf{X}) \\ &= \beta_0,\end{aligned}$$

where we used Assumption 1. Furthermore, by the LIE, the unconditional expectation of $\hat{\beta}$ is

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}[\mathbb{E}(\hat{\beta}|\mathbf{X})] = \mathbb{E}(\beta_0) = \beta_0.$$



Unbiasedness

- In repeated samples (i.e., sample infinitely many times), the distribution of the estimator is centered on the true value β_0 .
- In repeated samples, the estimator captures the truth on average.
- This is a desirable property for an estimator, but by no means a sufficient condition for the estimator to be useful.
- Next, we look at the bias property of the least-squares estimator of σ_0^2 .
- Recall that the least squares estimator simply uses the sample variance of the residuals, $\hat{\sigma}^2 = \hat{\epsilon}'\hat{\epsilon}/n$.
- Recall that

$$\hat{\epsilon} = \mathbb{M}_x \mathbf{y} = \mathbb{M}_x (\mathbf{X}\beta + \epsilon) = \mathbb{M}_x \epsilon$$

since $\mathbb{M}_x \mathbf{X} = \mathbf{0}$.

Unbiasedness

- Then, $\hat{\epsilon}'\hat{\epsilon} = (\mathbb{M}_x\epsilon)'(\mathbb{M}_x\epsilon) = \epsilon'\mathbb{M}_x\epsilon$.
- Using the fact that trace of scalar equals itself and by Assumption 2, we have

$$\begin{aligned}\mathbb{E}(\epsilon'\mathbb{M}_x\epsilon|\mathbf{X}) &= \mathbb{E}[\text{tr}(\epsilon'\mathbb{M}_x\epsilon)|\mathbf{X}] = \mathbb{E}[\text{tr}(\mathbb{M}_x\epsilon\epsilon')|\mathbf{X}] \\ &= \text{tr}[\mathbb{M}_x\mathbb{E}(\epsilon\epsilon'|\mathbf{X})] = \sigma_0^2\text{tr}(\mathbb{M}_x).\end{aligned}$$

- Furthermore, $\text{tr}(\mathbb{M}_x) = \text{tr}(\mathbf{I}_n) - \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] =$
 $\text{tr}(\mathbf{I}_n) - \text{tr}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}] = \text{tr}(\mathbf{I}_n) - \text{tr}(\mathbf{I}_{k+1}) = n - k - 1$.
- Hence,

$$\mathbb{E}(\hat{\sigma}^2|\mathbf{X}) = \frac{1}{n}\mathbb{E}(\hat{\epsilon}'\hat{\epsilon}|\mathbf{X}) = \frac{1}{n}\mathbb{E}(\epsilon'\mathbb{M}_x\epsilon|\mathbf{X}) = \left(\frac{n-k-1}{n}\right)\sigma_0^2.$$

- Also, by the LIE, this is true unconditionally.

Unbiasedness

- This result shows that the least-squares estimator $\hat{\sigma}^2$ underestimates σ_0^2 (biased towards zero), but the bias is small for large samples.
- The unbiased least-squares estimator of σ_0^2 is

$$\tilde{\sigma}^2 = \frac{\hat{\varepsilon}'\hat{\varepsilon}}{n - k - 1}$$

where $\tilde{\sigma}$ is called the **standard error of the regression**.

- The sampling variance of $\hat{\beta}$ can be driven as follows:

$$\begin{aligned}\text{Var}(\hat{\beta}|\mathbf{X}) &= \mathbb{E} \left[(\hat{\beta} - \beta_0)(\hat{\beta} - \beta_0)' | \mathbf{X} \right] \\ &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbb{E}(\varepsilon\varepsilon' | \mathbf{X}) \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\ &= \sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1}.\end{aligned}$$

- Unlike the conditional mean of $\hat{\beta}$, conditional variance of $\hat{\beta}$ depends on \mathbf{X} .
- Hence, $\sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1}$ is not the unconditional variance of $\hat{\beta}$, unless \mathbf{X} is deterministic.

Sampling variance of the LSE

- Applying the LIE, we get $\text{Var}(\hat{\beta}) = \mathbb{E}(\text{Var}(\hat{\beta})|\mathbf{X}) = \sigma_0^2 \mathbb{E} \left[(\mathbf{X}'\mathbf{X})^{-1} \right]$.
- Also, note that $\mathbb{E} \left[(\mathbf{X}'\mathbf{X})^{-1} \right] \neq \left[\mathbb{E}(\mathbf{X}'\mathbf{X}) \right]^{-1}$.
- We must assume that $\mathbb{E} \left[(\mathbf{X}'\mathbf{X})^{-1} \right]$ exists, otherwise the variance, conditional or unconditional, is not well-defined.¹
- For randomly drawn \mathbf{X} in repeated samples from an experiment, the dispersion of the distribution of $\hat{\beta}$ is given by $\sigma_0^2 \mathbb{E} \left[(\mathbf{X}'\mathbf{X})^{-1} \right]$, i.e., **the average case**.
- But for practical purposes, the conditional variance can be the relevant measure of dispersion.
- Since \mathbf{X} is observed, in repeated samples one can consider confining attention to those samples with this \mathbf{X} and the relevant measure of dispersion becomes the conditional variance.
- But, note that variance (of the distribution of an estimator) is used to compare the efficiency of estimators, hence the average case is the natural choice in that case.

¹For any random variables y and x , $\mathbb{E}(y|x)$ is well-defined only if $\mathbb{E}|y| < \infty$.

Efficiency of the LSE

- Let $\mathbf{L} \in \mathbb{R}^{(k+1) \times n}$ some matrix such that $\boldsymbol{\beta}^\dagger = \mathbf{L}\mathbf{y}$.
- Here, $\boldsymbol{\beta}^\dagger$ denotes a class of linear estimators of $\boldsymbol{\beta}_0$.
- The Gauss-Markov Theorem claims that the least-squares estimator $\hat{\boldsymbol{\beta}}$ is the best in the sense that within the class of linear estimators of $\boldsymbol{\beta}_0$ the least squares estimator has the smallest variance.
- In other words, it makes the most efficient use of the information in a given sample.
- Before, we proceed to state the theorem and its proof, we make the following assumptions:
 - ④ $\mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{X}, \mathbf{L}) = \mathbf{0}$ a.s.
 - ⑤ $\mathbb{E}(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}' | \mathbf{X}, \mathbf{L}) = \sigma_0^2 \mathbf{I}_n$ a.s.
 - ⑥ $\mathbf{L}\mathbf{X} = \mathbf{I}_{k+1}$ a.s.

Efficiency of the LSE

- Given these assumptions we have

$$\mathbb{E}(\beta^\dagger | \mathbf{X}, \mathbf{L}) = \mathbf{LX}\beta_0 + \mathbb{E}(\mathbf{L}\epsilon | \mathbf{X}, \mathbf{L}) = \beta_0,$$

and using the LIE we obtain

$$\mathbb{E}(\beta^\dagger) = \mathbb{E}(\mathbb{E}(\beta^\dagger | \mathbf{X}, \mathbf{L})) = \mathbb{E}(\mathbf{LX})\beta_0 + \mathbb{E}(\mathbf{L}\mathbb{E}(\epsilon | \mathbf{X}, \mathbf{L})) = \beta_0.$$

- So, β^\dagger is unbiased.
- Furthermore, the sampling variance of β^\dagger is given by

$$\text{Var}(\beta^\dagger | \mathbf{X}, \mathbf{L}) = \mathbb{E} \left[(\beta^\dagger - \beta) (\beta^\dagger - \beta)' | \mathbf{X}, \mathbf{L} \right] = \mathbb{E} \left(\mathbf{L}\epsilon\epsilon' \mathbf{L}' | \mathbf{X}, \mathbf{L} \right) = \sigma_0^2 \mathbf{L}\mathbf{L}'$$

and by the LIE we have

$$\text{Var}(\beta^\dagger) = \sigma_0^2 \mathbb{E}(\mathbf{L}\mathbf{L}').$$

Efficiency of the LSE

Theorem (Gauss-Markov)

The difference between $\text{Var}(\beta^\dagger)$ and $\text{Var}(\hat{\beta})$ is a positive semi-definite matrix for every \mathbf{L} satisfying Assumptions 4–6.

- Note that we have not used the terminology **BLUE** here.
- If indeed we had deterministic regressors, we would simply state the theorem as **the least-squares estimator is BLUE**.
- However, this is not necessarily true for a regression model with stochastic regressors.

Efficiency of the LSE

Proof.

Let $D = L - (X'X)^{-1}X'$, so that $DX = 0$ by assumption. Then, note that

$$\begin{aligned} LL' &= (X'X)^{-1}X'X(X'X)^{-1} + (X'X)^{-1}X'D' + DX(X'X)^{-1} + DD' \\ &= (X'X)^{-1} + DD'. \end{aligned}$$

It follows that

$$\text{Var}(\beta^\dagger | X, L) = \text{Var}(\hat{\beta} | X, L) + \sigma_0^2 DD',$$

where DD' positive semi-definite (well-known result). Then, taking unconditional expectations

$$\text{Var}(\beta^\dagger) = \text{Var}(\hat{\beta}) + \sigma_0^2 \mathbb{E}(DD'),$$

and noting that for arbitrary fixed $(k+1)$ -vector c ,

$$c' \mathbb{E}(DD') c = \mathbb{E}(c' DD' c) \geq 0$$

since $c' DD' c$ is non-negative (a sum of squares) for any D . □

Efficiency of the LSE

- Note again that in the non-stochastic \mathbf{X} and \mathbf{L} case, Assumption 6 holds for every member of the class β^\dagger .
- Hence $\hat{\beta}$ is said to be BLUE, best (minimum variance) in the class of linear unbiased estimators.
- But, in the stochastic case, note that

$$\mathbb{E}(\beta^\dagger) = \mathbb{E}(\mathbf{LX})\beta_0 + \mathbb{E}(\mathbf{L}\varepsilon)$$

and given Assumption 4, it is sufficient for unbiasedness if $\mathbb{E}(\mathbf{LX}) = \mathbf{I}_{k+1}$.

- This is weaker than Assumption 6.
- If $\mathbf{LX} \neq \mathbf{I}_{k+1}$ with positive probability, the estimator is conditionally biased.
- However, under $\mathbb{E}(\mathbf{LX}) = \mathbf{I}_{k+1}$, the bias terms average out over the distribution of \mathbf{LX} .

Efficiency of the LSE

- Therefore, the estimator is unbiased in repeated sampling under the statistical model with stochastic regressors.
- But, we cannot conclude from $\mathbb{E}(\mathbf{LX}) = \mathbf{I}_{k+1}$ that $\mathbb{E}(\mathbf{DX}(\mathbf{X}'\mathbf{X})^{-1}) = \mathbf{0}$, therefore the proof above breaks down.
- Hence, in such a case, although the estimator is unbiased, we cannot show that the least squares estimator is at least as efficient as this linear estimator.
- Therefore, for a statistical model with random regressors, it is not correct to claim that $\hat{\beta}$ is BLUE (except in the context of the conditional distribution).
- In practice, in a repeated sampling experiment where \mathbf{X} is held fixed and a new ϵ is sampled in each replication of the experiment, the least squares has the variance of $\sigma_0^2(\mathbf{X}'\mathbf{X})^{-1}$ and therefore the most efficient in the class of linear unbiased estimators. However, if both ϵ and \mathbf{X} are sampled in each replication of the experiment, the least squares has the variance of $\sigma_0^2\mathbb{E}[(\mathbf{X}'\mathbf{X})^{-1}]$, and we cannot claim that the least squares estimator is the best in the linear unbiased class. But, there is no obvious candidate for a more efficient estimator, and this result is more formal than practical significance.

Finite sample properties of the LSE

- To sum, Assumptions 1–3 are sufficient to characterize the mean and variance of the finite sample distribution of the least-squares estimator.
- Note that we have not specified any distribution for $\mathbf{y}|\mathbf{X}$, such as the normal distribution or some other distribution.
- Indeed, if we further make the normality assumption, we can characterize the **exact distribution** of the least-squares estimator in finite samples.
- In other words, if we additionally assume that $\mathbf{y}|\mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}_0, \sigma_0^2 \mathbf{I}_n)$, then

$$\hat{\boldsymbol{\beta}}|\mathbf{X} \sim N(\boldsymbol{\beta}_0, \sigma_0^2 (\mathbf{X}'\mathbf{X})^{-1}).$$

- Note again that this is the exact distribution of $\hat{\boldsymbol{\beta}}|\mathbf{X}$. We are not utilizing any approximation methods (large sample methods).

Finite sample properties of the LSE

- You can already see that this result can be used for inference.
- In significance testing exercises, we are interested in whether a given regressor belongs to the model.
- Say, we'd like to test if the j th regressor belongs to the model, i.e.,
 $H_0 : \beta_{j0} = 0$.

- We can compute

$$\frac{\hat{\beta}_j - 0}{\sqrt{(\sigma_0^2(\mathbf{X}'\mathbf{X})^{-1})_{jj}}} \sim N(0, 1).$$

where jj subscript denotes the j th row j th column of the matrix.

- This is the t -test statistic you are familiar with. The only problem is that it is not feasible because we do not know σ_0^2 .
- We replace the unknown σ_0^2 with its unbiased estimator $\tilde{\sigma}^2$ to make it feasible:

$$\frac{\hat{\beta}_j - 0}{\sqrt{(\tilde{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1})_{jj}}} \sim t(n - k - 1).$$

Remarks

- For Gauss-Markov theorem to hold the normality assumption is not necessary.
- Note that we have not characterized the **large sample** distribution of the least-squares estimator.
- The behavior of the least-squares estimator can be studied in large samples (as n grows without a bound) using the so-called large sample tools such as LLNs and CLTs.
- Indeed, it is straightforward to show that the LSE in large samples is approximately distributed as normal, and the larger n is the better the approximation is.
- Also, there is a large sample counterpart of the Gauss-Markov theorem.
- The t -test statistic is also approximately normal in large samples, and the larger n is the better the approximation is.

Hypothesis testing

- The term **hypotheses testing** refers to the process of trying to decide the truth or falsity of hypotheses on the basis of experimental/observational evidence.
- Experimental/observational measurements are subject to random error.
- Therefore, any decision about the truth or falsity of the hypothesis, based on experimental/observational evidence, also is subject to error.
- It will not be possible to avoid an occasional decision error, but it will be possible to construct tests so that such errors occur infrequently at some prescribed rate.
 - ① Type I error: Reject a true H_0 .
 - ② Type II error: Fail to reject a false H_0 .
- For a simple null hypothesis, the probability of rejecting a true H_0 , is referred to as the **significance level** of the test.

Hypothesis testing

Hypothesis testing involve the following steps:

- Decide on your tolerance level (or level of significance) for making a Type I error, i.e., out of 100 times how many times you're willing to reject a true null hypothesis?
- Calculate a test statistic \mathcal{T}_n .
- Obtain the critical value c and
 - ☐ reject H_0 if $\mathcal{T}_n > c$,
 - ☐ fail to reject H_0 if $\mathcal{T}_n \leq c$.