

LECTURE 14: STATISTICAL PROPERTIES OF ORDINARY LEAST SQUARES

MECO 7312.
INSTRUCTOR: DR. KHAI CHIONG
NOVEMBER 20, 2019

1. Multivariate linear regression

Consider the regression model:

$$(1) \quad y_i = \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{iK} + \epsilon_i$$

Using Matrix Algebra:

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1K} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nK} \end{bmatrix}$$

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_K \end{bmatrix}$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

$\mathbf{X}\boldsymbol{\beta}$ is the matrix product of a $n \times K$ matrix with a $K \times 1$ matrix, resulting in a $n \times 1$ matrix.

2. Ordinary Least Squares (OLS) estimator

Our dataset is: $(y_i, x_{i1}, x_{i2}, \dots, x_{ik})$ for $i = 1, \dots, n$. In matrix form, the data matrix is $[\mathbf{y}, \mathbf{X}]$. Data-generating process (DGP): $(y_i, x_{i1}, x_{i2}, \dots, x_{ik}, \epsilon_i)$ are i.i.d random sample from $(Y, X_1, X_2, \dots, X_k, \epsilon)$ such that $Y = \sum_{k=1}^K \beta_k X_k + \epsilon$. We want to estimate the vector of parameters $\boldsymbol{\beta}$.

Make additional assumption on the data-generating process, derive k number of moments conditions, and use Method of Moments: $\mathbb{E}[\epsilon|X_1] = 0$, $\mathbb{E}[\epsilon|X_2] = 0$, \dots , $\mathbb{E}[\epsilon|X_k] = 0$. This is known as the conditional mean independence assumption, or the exogeneity assumption.

Recall that $\mathbb{E}[\epsilon|X] = 0$ implies that $\mathbb{E}[X\epsilon] = 0$, $\text{Cov}(X, \epsilon) = 0$, and also that $\mathbb{E}[\epsilon] = 0$.

The sample moment conditions can be written as: $\sum_{i=1}^n x_{i1}\epsilon_i = 0$, $\sum_{i=1}^n x_{i2}\epsilon_i = 0$, \dots , $\sum_{i=1}^n x_{ik}\epsilon_i = 0$. In matrix notation:

$$(2) \quad \mathbf{X}_1^T \boldsymbol{\epsilon} = 0$$

$$\vdots$$

$$(3) \quad \mathbf{X}_k^T \boldsymbol{\epsilon} = 0$$

Where \mathbf{X}_k is the k -th row of \mathbf{X} , and \mathbf{X}^T is the matrix transpose of \mathbf{X} .

Finally, the moment conditions can be summarized as just:

$$(4) \quad \mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{0}$$

Now, we can derive the OLS (Ordinary Least Square) estimators:

$$(5) \quad \mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{0}$$

$$(6) \quad \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

$$(7) \quad \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$$(8) \quad (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{0}$$

$$(9) \quad (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} - \boldsymbol{\beta} = \mathbf{0}$$

$$(10) \quad \boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

3. The geometry of OLS estimation

Consider the matrix of regressors:

$$(11) \quad \mathbf{X} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_K]$$

Each column vector \mathbf{x}_k is a point in the Euclidean space \mathbb{R}^n , where n is the number of observations (and the length of the vector \mathbf{x}_k). Now consider the set of all possible points (vectors) that can be achieved as a linear combination of $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$. That is, $\mathcal{S}(\mathbf{X}) = \{\mathbf{z} \in \mathbb{R}^n \mid \mathbf{z} = \sum_{k=1}^K \alpha_k \mathbf{x}_k, \alpha_k \in \mathbb{R}\}$. The object $\mathcal{S}(\mathbf{X}) \subseteq \mathbb{R}^n$ is called a subspace of \mathbf{X} (or the column space of \mathbf{X} , or the space span by $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_K$).

Now observe that $\mathbf{X}\boldsymbol{\beta}$ is just a linear combination of the columns of \mathbf{X} :

$$(12) \quad \mathbf{X}\boldsymbol{\beta} = [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \dots \quad \mathbf{x}_K] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{bmatrix} = \sum_{k=1}^K \beta_k \mathbf{x}_k$$

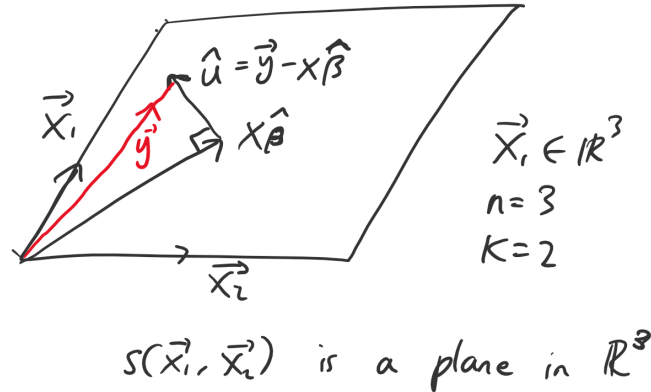
Therefore, the fitted values $\mathbf{X}\hat{\boldsymbol{\beta}}$ lies in the subspace of \mathbf{X} . Recall that the OLS estimator $\hat{\boldsymbol{\beta}}$ is a solution to $\mathbf{X}^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}$, which implies that $\mathbf{x}_k^T(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = 0$ for all $k = 1, \dots, K$, and that \mathbf{x}_k is *orthogonal* to $\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$. We term $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ the vector of fitted residuals. Therefore, we have $\sum_{i=1}^n x_{ik}\hat{u}_i = 0$ for $k = 1, \dots, n$.

Since the vector $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ is orthogonal to all the column vectors of \mathbf{X} , it is orthogonal to the subspace of \mathbf{X} . As such, the vector of fitted value $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ is also orthogonal to $\hat{\mathbf{u}}$, i.e., $\sum_{i=1}^n \hat{y}_i \hat{u}_i = 0$.

Therefore, for a given vector \mathbf{y} , OLS finds $\mathbf{X}\hat{\boldsymbol{\beta}}$, the vector in the subspace of \mathbf{X} , that is the closest to \mathbf{y} . In another words, $\mathbf{X}\hat{\boldsymbol{\beta}}$ is the orthogonal projection of \mathbf{y} onto the subspace of \mathbf{X} .

The implication is that OLS decomposes \mathbf{y} into a systematic part explained by \mathbf{X} , and a residual part that is unrelated/orthogonal to \mathbf{X} . Moreover, \mathbf{y} , $\mathbf{X}\hat{\boldsymbol{\beta}}$ and $\hat{\mathbf{u}}$ forms a right-angled triangle in Euclidean space, and we can apply the Pythagoras' theorem:

$$(13) \quad \|\mathbf{y}\|^2 = \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2 + \|\hat{\mathbf{u}}\|^2$$



In words, the total sum of squares is equal to the explained sum of squares plus the sum of squared residuals. This motivates the use of R^2 as the goodness-of-fit of OLS.

$$(14) \quad R^2 = \frac{\|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{\|\mathbf{y}\|^2} = \frac{(\mathbf{X}\hat{\boldsymbol{\beta}})^T(\mathbf{X}\hat{\boldsymbol{\beta}})}{\mathbf{y}^T\mathbf{y}} = \frac{\hat{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}\hat{\boldsymbol{\beta}}}{\mathbf{y}^T\mathbf{y}}$$

Which always lies between 0 and 1. $\|\mathbf{y}\|^2$ is the Euclidean length of the vector \mathbf{y} .

3.1. Linear transformation of regressors

Suppose we take a regressor \mathbf{x}_k and transform it as $a\mathbf{1} + b\mathbf{x}_k$, where $\mathbf{1}$ is a vector of ones of the same dimension as \mathbf{x}_k .

How does OLS estimation change? *The fitted OLS values and the residuals do not change! (The estimated coefficients would change however, but we can just apply the inverse of the transformation on the estimated coefficients)* This is a consequence of OLS as an orthogonal projection.

If we apply a linear transformation on \mathbf{X} , the subspace of \mathbf{X} does not change. Therefore the orthogonal projection of \mathbf{y} onto the same subspace would give the same OLS fitted values and residuals.

To see this, let \mathbf{A} be a $k \times k$ non-singular matrix. Then, $\mathbf{X}\mathbf{A}$ is a nonsingular linear transformation. For instance, \mathbf{X} consists of two regressors, $\mathbf{X} = [\mathbf{1}, \mathbf{T}]$, where \mathbf{T} is the recorded temperature in Celsius. The transformation $\mathbf{F} = 32 \cdot \mathbf{1} + \frac{9}{5}\mathbf{T}$ can then be expressed as:

$$(15) \quad [\mathbf{1} \quad \mathbf{F}] = [\mathbf{1} \quad \mathbf{T}] \begin{bmatrix} 1 & 32 \\ 0 & 9/5 \end{bmatrix}$$

Now verify that regressing \mathbf{y} on \mathbf{X} is equivalent to regressing \mathbf{y} on \mathbf{XA} .

Let $\hat{\boldsymbol{\beta}}$ be the OLS estimate of regressing \mathbf{y} on \mathbf{X} , and let $\hat{\boldsymbol{\beta}}_A$ be the OLS estimate of regressing \mathbf{y} on \mathbf{XA} . Show that $\hat{\boldsymbol{\beta}}_A = \mathbf{A}^{-1}\hat{\boldsymbol{\beta}}$, and therefore the fitted values are the same.

$$\begin{aligned} \hat{\boldsymbol{\beta}}_A &= ((\mathbf{XA})^T(\mathbf{XA}))^{-1}(\mathbf{XA})^T\mathbf{y} \\ &= (\mathbf{A}^T\mathbf{X}^T\mathbf{XA})^{-1}\mathbf{A}^T\mathbf{X}^T\mathbf{y} \end{aligned}$$

We have used the fact that $(AB)^T = B^T A^T$. Moreover, suppose we have two (square) invertible matrices A and B . Then, $(AB)^{-1} = B^{-1}A^{-1}$. Since \mathbf{A} and $\mathbf{A}^T\mathbf{X}^T\mathbf{X}$ are both *square* invertible,

$$\begin{aligned} \hat{\boldsymbol{\beta}}_A &= ((\mathbf{XA})^T(\mathbf{XA}))^{-1}(\mathbf{XA})^T\mathbf{y} \\ &= (\mathbf{A}^T\mathbf{X}^T\mathbf{XA})^{-1}\mathbf{A}^T\mathbf{X}^T\mathbf{y} \\ &= \mathbf{A}^{-1}(\mathbf{A}^T\mathbf{X}^T\mathbf{X})^{-1}\mathbf{A}^T\mathbf{X}^T\mathbf{y} \\ &= \mathbf{A}^{-1}(\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{A}^T)^{-1}\mathbf{A}^T\mathbf{X}^T\mathbf{y} \\ &= \mathbf{A}^{-1}\hat{\boldsymbol{\beta}} \end{aligned}$$

Therefore the fitted value is $\mathbf{XA}\mathbf{A}^{-1}\hat{\boldsymbol{\beta}} = \mathbf{X}\hat{\boldsymbol{\beta}}$.

3.2. Orthogonal projection

Define the projection matrix as:

$$(16) \quad P_{\mathbf{X}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

Multiplying $P_{\mathbf{X}}$ with any vector \mathbf{y} results in the orthogonal projection of \mathbf{y} on to the subspace of \mathbf{X} :

$$(17) \quad P_{\mathbf{X}}\mathbf{y} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

Suppose that \mathbf{y} can be perfectly explained by the regressors \mathbf{X} , i.e. $\mathbf{y} = \mathbf{X}\tilde{\boldsymbol{\beta}}$ for some vector $\tilde{\boldsymbol{\beta}}$ without any error term $\boldsymbol{\epsilon}$. then check that $P_{\mathbf{X}}\mathbf{y} = \mathbf{y}$, and so the OLS estimator is $\hat{\boldsymbol{\beta}} = \tilde{\boldsymbol{\beta}}$, and $\hat{\mathbf{u}} = \mathbf{0}$.

Projection matrix $P_{\mathbf{X}}$ is idempotent: $P_{\mathbf{X}}P_{\mathbf{X}} = P_{\mathbf{X}}$.

Define another projection matrix $M_{\mathbf{X}} = \mathbf{I} - P_{\mathbf{X}}$, where \mathbf{I} is the $n \times n$ identity matrix. While $P_{\mathbf{X}}$ creates the fitted value, $M_{\mathbf{X}}$ is the residual maker. We know from the last section that:

$$(18) \quad \mathbf{y} = P_{\mathbf{X}}\mathbf{y} + M_{\mathbf{X}}\mathbf{y}$$

Hence $P_{\mathbf{X}}\mathbf{y}$ is the vector in the subspace of \mathbf{X} closest to \mathbf{y} , while $M_{\mathbf{X}}\mathbf{y}$ is the vector that is unrelated to \mathbf{X} .

Consider the regression model $\mathbf{y} = \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{X}_2\boldsymbol{\beta}_2 + \mathbf{u}$, where we have partitioned \mathbf{X} into two groups of regressors, \mathbf{X}_1 and \mathbf{X}_2 . For example, \mathbf{X}_1 is a $n \times k_1$ matrix while \mathbf{X}_2 is a $n \times k_2$ matrix, with $K = k_1 + k_2$, and \mathbf{X} is a $n \times K$ matrix.

Consider the projection matrix $P_{\mathbf{X}_1} = \mathbf{X}_1(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T$ and $P_{\mathbf{X}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$. As a shorthand, $P_1 = P_{\mathbf{X}_1}$. A very useful property of projection matrices is that:

$$(19) \quad P_1P_{\mathbf{X}} = P_{\mathbf{X}}P_1 = P_1$$

Therefore if we fit \mathbf{y} to \mathbf{X} , then take the fitted value $P_{\mathbf{X}}\mathbf{y}$ and fit it to \mathbf{X}_1 , we get $P_1\mathbf{y}$, which is the fitted value of \mathbf{y} on \mathbf{X}_1 . Similarly, if we fit \mathbf{y} to \mathbf{X}_1 , then take the fitted value $P_1\mathbf{y}$ and fit it to \mathbf{X} , we get $P_1\mathbf{y}$.

Proof:

$$\begin{aligned} P_{\mathbf{X}}P_1 &= P_{\mathbf{X}}\mathbf{X}_1(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T \\ &= \mathbf{X}_1(\mathbf{X}_1^T\mathbf{X}_1)^{-1}\mathbf{X}_1^T \\ &= P_1 \end{aligned}$$

$P_{\mathbf{X}}\mathbf{X}_1 = \mathbf{X}_1$ because all columns of \mathbf{X}_1 belong to the subspace of \mathbf{X} , so the orthogonal projection of \mathbf{X}_1 onto the subspace of \mathbf{X} is just itself. In another words, if we regress (each column of) \mathbf{X}_1 on \mathbf{X} , then the fitted value is just \mathbf{X}_1 since \mathbf{X}_1 can be explained perfectly by $\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2]$.

Finally, we have:

$$\begin{aligned} P_1^T &= (P_{\mathbf{X}}P_1)^T \\ P_1 &= P_1^T P_{\mathbf{X}}^T \\ P_1 &= P_1 P_{\mathbf{X}} \end{aligned}$$

It is straightforward to verify that $P_1 P_X = P_X P_1 = P_1$ implies that $M_1 M_X = M_X M_1 = M_X$.

4. Frisch-Waugh-Lovell (FWL) theorem

Consider the regression model $\mathbf{y} = \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \text{residuals}$, where we have partitioned \mathbf{X} into two groups of regressors, \mathbf{X}_1 and \mathbf{X}_2 .

$$(20) \quad \mathbf{y} = \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \text{residuals}$$

$$(21) \quad M_{\mathbf{X}_1} \mathbf{y} = M_{\mathbf{X}_1} \mathbf{X}_2 \boldsymbol{\beta}_2 + \text{residuals}$$

Where $M_{\mathbf{X}_1} = \mathbf{I} - P_{\mathbf{X}_1}$, $P_{\mathbf{X}_1} = \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$.

The OLS estimates $\boldsymbol{\beta}_2$ from (20) and (21) are identical.

When we apply $M_{\mathbf{X}_1}$ to \mathbf{y} , it removes and purges \mathbf{X}_1 from \mathbf{y} , as such $M_{\mathbf{X}_1} \mathbf{y}$ is orthogonal and unrelated to \mathbf{X}_1 .

FWL theorem is very important, as it means that we can first regress \mathbf{y} on \mathbf{X}_1 , collect the residual as $\tilde{\mathbf{y}}$. Then regress each of the column vector of \mathbf{X}_2 on \mathbf{X}_1 , collect the residuals as $\tilde{\mathbf{X}}_2$. Finally, regress $\tilde{\mathbf{y}}$ on $\tilde{\mathbf{X}}_2$ as if \mathbf{X}_1 has been controlled for.

For example, a regression involving time-series data that have been de-trend and de-seasonalize is equivalent to a regression involving the original time-series data plus seasonality dummy and time trend variables. FWL also underlies the 2-stage least square (2SLS) procedure in instrumental variable estimation.

As a shortcut, let $M_1 = M_{\mathbf{X}_1} \mathbf{X}_2$. Projection matrices have the property that $M_1 M_1 = M_1$.

The OLS from the second regression is:

$$\begin{aligned} \tilde{\boldsymbol{\beta}}_2 &= ((M_1 \mathbf{X}_2)^T M_1 \mathbf{X}_2)^{-1} (M_1 \mathbf{X}_2)^T M_1 \mathbf{y} \\ &= (\mathbf{X}_2^T M_1^T M_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^T M_1^T M_1 \mathbf{y} \\ &= (\mathbf{X}_2^T M_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^T M_1 \mathbf{y} \end{aligned}$$

Let $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ be the OLS from the first regression. As such,

$$\begin{aligned}
\mathbf{y} &= \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + M_X \mathbf{y} \\
\mathbf{X}_2^T M_1 \mathbf{y} &= \mathbf{X}_2^T M_1 \mathbf{X}_1 \hat{\boldsymbol{\beta}}_1 + \mathbf{X}_2^T M_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{X}_2^T M_1 M_X \mathbf{y} \\
&= \mathbf{X}_2^T M_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + \mathbf{X}_2^T M_1 M_X \mathbf{y} \\
&= \mathbf{X}_2^T M_1 \mathbf{X}_2 \hat{\boldsymbol{\beta}}_2 + (M_X M_1 \mathbf{X}_2)^T \mathbf{y}
\end{aligned}$$

Now $M_X M_1 \mathbf{X}_2 = \mathbf{0}$. Recall from the last section that $M_X M_1 = M_X$. As such, $M_X \mathbf{X}_2 = \mathbf{0}$ because \mathbf{X}_2 can be perfectly explained by \mathbf{X} and so there is no residual left.

Therefore $\hat{\boldsymbol{\beta}}_2 = (\mathbf{X}_2^T M_1 \mathbf{X}_2)^{-1} \mathbf{X}_2^T M_1 \mathbf{y} = \tilde{\boldsymbol{\beta}}_2$

5. Unbiasedness of OLS estimators

Is the OLS estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ unbiased? What does unbiasedness mean here? We need a ground truth, and say that it is unbiased with respect to a data-generating process.

DGP: $(y_i, x_{i1}, x_{i2}, \dots, x_{ik}, \epsilon_i)$, for $i = 1, \dots, n$, are generated from some joint distribution that obeys the equation $\mathbf{y} = \mathbf{X} \boldsymbol{\beta}_0 + \boldsymbol{\epsilon}$. We can be agnostic about this joint distribution, in particular, ϵ_i may not even be i.i.d across i .

Let $\hat{\boldsymbol{\beta}}$ be the OLS estimator.

$$\begin{aligned}
(22) \quad \mathbb{E}[\hat{\boldsymbol{\beta}}] &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\
(23) \quad &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \boldsymbol{\beta}_0 + \boldsymbol{\epsilon})] \\
(24) \quad &= \mathbb{E}[\boldsymbol{\beta}_0 + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon}] \\
(25) \quad &= \boldsymbol{\beta}_0 + \mathbb{E}[\mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\epsilon} | \mathbf{X}]] \\
(26) \quad &= \boldsymbol{\beta}_0 + \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\boldsymbol{\epsilon} | \mathbf{X}]]
\end{aligned}$$

The Law of Iterated Expectation is applied in the last two equations. It is clear that a sufficient condition for the unbiasedness of OLS estimator is that $\mathbb{E}[\boldsymbol{\epsilon} | \mathbf{X}] = \mathbf{0}$. This expression means that $\mathbb{E}[\epsilon_i | \mathbf{X}] = 0$ for all $i = 1, \dots, n$. Further unpacking, it means that ϵ_i for each $i = 1, \dots, n$ is (conditionally mean) independent from the entire matrix \mathbf{X} , i.e. $\mathbb{E}[\epsilon_i | x_{11}, \dots, x_{ik}, \dots, x_{nk}] = 0$.

Two possible ways to satisfy $\mathbb{E}[\epsilon_i | \mathbf{X}] = 0$.

(1) $(\epsilon_i, x_{i1}, x_{i2}, \dots, x_{iK})$ are independently and identically distributed across i from some probability distributions, *and*, $\mathbb{E}[\epsilon_i | x_{i1}, x_{i2}, \dots, x_{iK}] = 0$. In the i.i.d case, we

can drop the i subscript, and write $\mathbb{E}[\epsilon|X_1, X_2, \dots, X_K] = 0$. Now $\mathbb{E}[\epsilon|X_1, X_2, \dots, X_K] = 0$ implies that $\mathbb{E}[\epsilon|X_k] = 0$ for $k = 1, \dots, K$. This is what we assumed when we use the Method of Moments to derive the OLS estimator. This is a sufficient but not a necessary condition for unbiasedness.

(2)* $(\epsilon_i, x_{i1}, x_{i2}, \dots, x_{iK})$ are *not necessarily* i.i.d across i , but $\mathbb{E}[\epsilon_i|\mathbf{X}] = 0$ for each i . Therefore, in the context of time-series where i.i.d does not hold true, OLS can still be unbiased if x_{tk} is independent of ϵ_t for all $t = 1, \dots, T$. There can be no correlation between the error term at time t and your covariates at time $t + 1, t + 2, \dots$.

Exogeneity alone guarantees unbiasedness. Exogeneity can be violated under many circumstances – whenever the regressor is correlated with the error term. Let us consider cross-sectional data, so that we are in the i.i.d scenario (the first scenario above). The cross-sectional data consist of gas stations' quantities sold (number of gallons) and prices, across multiple gas stations at a given hour. If we regress quantities sold on price, we expect the coefficient to be negative, i.e. high price causes lower demand. But often this regression gives us positive price coefficient! This is because price is an endogenous variable.

Prices are set strategically by firms. Suppose some gas station locations are popular because of their more friendly staffs. Customers are willing to pay more for service friendliness, and so firms respond by charging higher prices. Now it is highly unlikely that we can ever observe or even measure service friendliness. Therefore, service friendliness becomes part of the error term ϵ that can explain demand. It follows that ϵ correlates with price here (positive correlation according to our story). In general, even if we can measure and observe service friendliness, there could be some unobserved preference shocks that somehow drive higher demand at one location, and which correlates with prices.

If price is set randomly (experimentation or A/B testing), then it would be exogenous. This is the foremost concern in any empirical research. The branch of statistics/econometrics dealing with this concern is called *causal inference*. Tools that fall under causal inference include (1) Instrumental Variable approach, (2) Difference-in-difference, (3) Regression discontinuity, (4) Propensity score matching, (5) experimentation and A/B testing.