# Advanced Econometrics

## Lecture 4: Conditional Expectation, Linear Projection and Linear Causal Model (Hansen Chapter 2)

Instructor: Ma, Jun

Renmin University of China

October 8, 2021

# Conditional expectation function

- Given the joint density $f_{Y,X}(y, x)$, $X$ has the marginal density

$$f_X(x) = \int_{-\infty}^{\infty} f_{Y,X}(y, x) \, dy.$$

- The conditional density of $Y$ given $X$:

$$f_{Y|X}(y \mid x) = \frac{f_{Y,X}(y, x)}{f_X(x)}.$$

- The CEF of $Y$ given $X = x$ is the mean of the conditional density:

$$m(x) = \int_{-\infty}^{\infty} y f_{Y|X}(y \mid x) \, dy.$$

  $m(x)$ is the mean of for the idealized subpopulation where the conditioning variables are fixed at $x$.

- $\mathbb{E}(Y \mid X = x)$ or $\mathbb{E}(Y \mid x)$ is interpreted as $m(x)$; $\mathbb{E}(Y \mid X)$ is interpreted as $m(X)$.

- We call this the conditional expectation function (CEF). The CEF is a function of $X$.

# Law of iterated expectations

> **Theorem (Simple Law of Iterated Expectations)**
> *If* $\mathbb{E}\,|Y| < \infty$ *then for any random vector* $\boldsymbol{X}$,
>
> $$\mathbb{E}\left(\mathbb{E}\left(Y \mid \boldsymbol{X}\right)\right) = \mathbb{E}\left(Y\right).$$

▶ When $\boldsymbol{X}$ is discrete

$$\mathbb{E}\left(\mathbb{E}\left(Y \mid \boldsymbol{X}\right)\right) = \sum_{j=1}^{\infty} \mathbb{E}\left(Y \mid \boldsymbol{x}_j\right) \Pr\left(\boldsymbol{X} = \boldsymbol{x}_j\right)$$

and when $\boldsymbol{X}$ is continuous

$$\mathbb{E}\left(\mathbb{E}\left(Y \mid \boldsymbol{X}\right)\right) = \int_{\mathbb{R}^k} \mathbb{E}\left(Y \mid \boldsymbol{x}\right) f_{\boldsymbol{X}}\left(\boldsymbol{x}\right) d\boldsymbol{x}.$$

▶ In this course, we ignore conditions such as $\mathbb{E}\,|Y| < \infty$.

Theorem
*For any random vectors $X_1$ and $X_2$,*

$$\mathbb{E}\left(\mathbb{E}\left(Y \mid X_1, X_2\right) \mid X_1\right) = \mathbb{E}\left(Y \mid X_1\right)$$
$$\mathbb{E}\left(\mathbb{E}\left(Y \mid X_1\right) \mid X_1, X_2\right) = \mathbb{E}\left(Y \mid X_1\right)$$

When you condition on a random vector $X$ you can effectively treat it as if it is constant. For example, $\mathbb{E}(X \mid X) = X$ and $\mathbb{E}(g(X) \mid X) = g(X)$ for any function $g(\cdot)$. The general property is known as the Conditioning Theorem.

Theorem (Conditioning Theorem)

$$\mathbb{E}(g(X)Y \mid X) = g(X)\mathbb{E}(Y \mid X)$$
$$\mathbb{E}(g(X)Y) = \mathbb{E}(g(X)\mathbb{E}(Y \mid X)).$$

# CEF error

▶ The CEF error is defined as

$$e = Y - m(X).$$

By construction,

$$Y = m(X) + e.$$

▶ By the linearity of expectations, the definition $m(X) = \mathbb{E}(Y \mid X)$ and the Conditioning Theorem

$$\begin{aligned}
\mathbb{E}(e \mid X) &= \mathbb{E}((Y - m(X)) \mid X) \\
&= \mathbb{E}(Y \mid X) - \mathbb{E}(m(X) \mid X) \\
&= m(X) - m(X) \\
&= 0.
\end{aligned}$$

▶ The unconditional mean is also zero:

$$\mathbb{E}(e) = \mathbb{E}(\mathbb{E}(e \mid X)) = \mathbb{E}(0) = 0.$$

> Theorem
> *Properties of the CEF error*
> *1.* $\mathbb{E}(e \mid X) = 0$.
> *2.* $\mathbb{E}(e) = 0$.
> *3. For any function* $h(\cdot)$, $\mathbb{E}(h(X)e) = 0$.

▶ The equations

$$Y = m(X) + e$$
$$\mathbb{E}(e \mid X) = 0$$

together imply that $m(X)$ is the CEF of $Y$ given $X$.

▶ It is important to understand that this is not a restriction. These equations hold true by definition.

► The equation $\mathbb{E}(e \mid X) = 0$ is called a conditional mean restriction, since the conditional mean of the error is restricted to equal zero.

► The property is also called **mean independence**, for the conditional mean of is 0 and thus independent of $X$. However, it does not imply that the distribution of $e$ is independent of $X$.

► As a simple example of a case where $X$ and $e$ are mean independent yet dependent, let $e = X\epsilon$ where $X$ and $e$ are independent $N(0, 1)$. Then conditional on $X$ the error $e$ has the distribution $N(0, x^2)$. Thus $\mathbb{E}(e \mid X) = 0$ and $e$ is mean independent of $X$, yet $e$ is not fully independent of $X$. Mean independence does not imply full independence.

# Regression variance

- An important measure of the dispersion about the CEF function is the unconditional variance of the CEF error $e$. We write this as

$$\sigma^2 = \text{Var}\,(e) = \mathbb{E}\left((e - \mathbb{E}e)^2\right) = \mathbb{E}\left(e^2\right).$$

- We can call $\sigma^2$ the regression variance or the variance of the regression error. The magnitude of $\sigma^2$ measures the amount of variation in $Y$ which is not "explained" or accounted for in the conditional mean $\mathbb{E}\,(Y \mid X)$.

► The regression variance depends on the regressors. Consider two regressions

$$Y = \mathbb{E}(Y \mid \boldsymbol{X}_1) + e_1$$
$$Y = \mathbb{E}(Y \mid \boldsymbol{X}_1, \boldsymbol{X}_2) + e_2.$$

► The simple relationship we now derive shows that the variance of this unexplained portion decreases when we condition on more variables. This relationship is monotonic in the sense that increasing the amount of information always decreases the variance of the unexplained portion.

Theorem

$$\operatorname{Var}(Y) \geq \operatorname{Var}(Y - \mathbb{E}(Y \mid \boldsymbol{X}_1)) \geq \operatorname{Var}(Y - \mathbb{E}(Y \mid \boldsymbol{X}_1, \boldsymbol{X}_2)).$$

# Best predictor

- ▶ Suppose that given a realized value of $X$, we want to create a prediction or forecast of $Y$.

- ▶ We can write any predictor as a function $g(X)$ of $X$. A non-stochastic measure of the magnitude of the prediction error $Y - g(X)$ is the expectation of its square $\mathbb{E}\left((Y - g(X))^2\right)$.

- ▶ What function is the best predictor? It turns out that the answer is the CEF:

$$
\begin{aligned}
\mathbb{E}\left((Y - g(X))^2\right) &= \mathbb{E}\left((e + m(X) - g(X))^2\right) \\
&= \mathbb{E}\left(e^2\right) + 2\mathbb{E}\left(e\left(m(X) - g(X)\right)\right) + \mathbb{E}\left((m(X) - g(X))^2\right) \\
&= \mathbb{E}\left(e^2\right) + \mathbb{E}\left((m(X) - g(X))^2\right) \\
&\geq \mathbb{E}\left(e^2\right) \\
&= \mathbb{E}\left((Y - m(X))^2\right).
\end{aligned}
$$

Theorem
*Conditional Mean as Best Predictor*
*For any predictor $g(X)$,*

$$\mathbb{E}\left((Y - g(X))^2\right) \geq \mathbb{E}\left((Y - m(X))^2\right)$$

*where $m(X) = \mathbb{E}(Y \mid X)$*

# Conditional variance

> **Definition**
> The conditional variance of $W$ given $X$ is
> $$\text{Var}(W \mid X) = \mathbb{E}\left((W - \mathbb{E}(W \mid X))^2 \mid X\right)$$

> **Definition**
> The conditional variance of the regression error $e$ is
> $$\sigma^2(X) = \text{Var}(e \mid X) = \mathbb{E}\left(e^2 \mid X\right)$$

- Generally, $\sigma^2(X) > 0$ is a function of $X$.
- Notice as well that $\sigma^2(X) = \text{Var}(Y \mid X)$ so it is equivalently the conditional variance of the dependent variable.
- We define the **conditional standard deviation** as its square root $\sigma(X) = \sqrt{\sigma^2(X)}$.
- The unconditional error variance and the conditional variance are related by the law of iterated expectations

$$\sigma^2 = \mathbb{E}\left(e^2\right) = \mathbb{E}\left(\mathbb{E}\left(e^2 \mid X\right)\right) = \mathbb{E}\left(\sigma^2(X)\right).$$

▶ Given the conditional variance, we can define a rescaled error

$$\varepsilon = \frac{e}{\sigma(X)}.$$

▶ We can calculate that since $\sigma(X)$ is a function of $X$

$$\mathbb{E}(\varepsilon \mid X) = \mathbb{E}\left(\frac{e}{\sigma(X)} \mid X\right) = \frac{1}{\sigma(X)}\mathbb{E}(e \mid X) = 0$$

and

$$\text{Var}(\varepsilon \mid X) = \mathbb{E}\left(\varepsilon^2 \mid X\right) = \mathbb{E}\left(\frac{e^2}{\sigma^2(X)} \mid X\right) = \frac{1}{\sigma^2(X)}\mathbb{E}\left(e^2 \mid X\right) = 1$$

Thus $\varepsilon$ has a conditional mean of zero, and a conditional variance of 1.

# Homoskedasticity and heteroskedasticity

Definition
The error is homoskedastic if $\mathbb{E}\left(e^2 \mid X\right) = \sigma^2$ does not depend on $X$.

Definition
The error is heteroskedastic if $\mathbb{E}\left(e^2 \mid X\right) = \sigma^2\left(X\right)$ depends on $X$.

- ► Older textbooks also tend to describe homoskedasticity as a component of a correct regression specification, and describe heteroskedasticity as an exception or deviance.

- ► The correct view is that heteroskedasticity is generic and "standard", while homoskedasticity is unusual and exceptional. The default in empirical work should be to assume that the errors are heteroskedastic.

- ► We will still frequently impose homoskedasticity when making theoretical investigations. In many cases homoskedasticity greatly simplifies the theoretical calculations.

# Regression derivative

- When a regressor $X_1$ is continuously distributed, we define the **marginal effect** of a change in $X_1$, holding the variables $X_2, ..., X_k$ fixed, as the partial derivative of the CEF $\frac{\partial}{\partial X_1} m(X_1, \ldots, X_k)$.

- When $X_1$ is discrete we define the marginal effect as a discrete difference. For example, if $X_1$ is binary, then the marginal effect of $X_1$ on the CEF is

$$m(1, X_2, ..., X_k) - m(0, X_2, ..., X_k).$$

- We can unify the continuous and discrete cases with the notation

$$\nabla_1 m(X) = \begin{cases} \frac{\partial}{\partial X_1} m(X_1, \ldots, X_k), & \text{if } x_1 \text{ is continuous} \\ m(1, X_2, ..., X_k) - m(0, X_2, ..., X_k), & \text{if } x_1 \text{ is binary.} \end{cases}$$

- Collecting the $k$ effects into one $k \times 1$ vector, we define we define the regression derivative to be

$$\boldsymbol{\nabla} m(X) = \begin{bmatrix} \nabla_1 m(X) \\ \nabla_2 m(X) \\ \vdots \\ \nabla_k m(X) \end{bmatrix}.$$

# Linear CEF

▶ An important special case is when the CEF $m(X) = \mathbb{E}(Y \mid X)$ is linear in $X$:

$$m(X) = X_1\beta_1 + X_2\beta_2 + \cdots + X_k\beta_k + \beta_{k+1}.$$

▶ An easy way to do so is to augment the regressor vector $X$ by listing the number "1" as an element. The corresponding coefficient is called the "intercept":

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_{k-1} \\ 1 \end{pmatrix}.$$

▶ With this redefinition, the CEF is

$$\begin{aligned} m(X) &= X_1\beta_1 + X_2\beta_2 + \cdots + X_k\beta_k + \beta_{k+1} \\ &= X'\boldsymbol{\beta} \end{aligned}$$

where $\boldsymbol{\beta} = (\beta_1, ..., \beta_{k+1})'$. This is the **linear CEF model**. It is also often called the **linear regression** model.

- In the linear CEF model, the regression derivative is simply the coefficient vector: $\nabla m(X) = \beta$. The coefficients have simple and natural interpretations as the marginal effects of changing one variable, holding the others constant.

- Linear CEF model:

$$Y = X'\beta + e$$
$$\mathbb{E}(e \mid X) = 0.$$

- Homoskedastic linear CEF model:

$$Y = X'\beta + e$$
$$\mathbb{E}(e \mid X) = 0$$
$$\mathbb{E}\left(e^2 \mid X\right) = \sigma^2.$$

# Linear CEF with nonlinear effects

- We can include as regressors nonlinear transformations of the original variables.
- The CEF could take the quadratic form

$$m(X_1, X_2) = X_1\beta_1 + X_2\beta_2 + X_1^2\beta_2 + X_2^2\beta_4 + X_1X_2\beta_5 + \beta_6.$$

This is also a linear CEF in the sense of being linear in the coefficients.

- The regression derivatives:

$$\frac{\partial}{\partial X_1}m(X_1, X_2) = \beta_1 + 2X_1\beta_3 + X_2\beta_5$$
$$\frac{\partial}{\partial X_2}m(X_1, X_2) = \beta_2 + 2X_2\beta_4 + X_1\beta_5.$$

We typically call $\beta_5$ the **interaction effect**. If $\beta_5 > 0$ then the regression derivative with respect to $X_1$ is increasing in the level of $X_2$.

# Best linear predictor

A linear predictor for is a function of the form $X'b$ for some $b \in \mathbb{R}^k$.
The mean squared prediction error is

$$S(b) = \mathbb{E}\left((Y - X'b)^2\right).$$

Definition
The Best Linear Predictor of $Y$ given $X$ is

$$\mathcal{P}(Y \mid X) = X'\beta$$

where $\beta$ minimizes the mean squared prediction error

$$S(b) = \mathbb{E}\left((Y - X'b)^2\right)$$

The minimizer

$$\beta = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} S(b)$$

is called the Linear Projection Coefficient.

▶ By calculations,

$$S(\boldsymbol{b}) = \mathbb{E}\left(Y^2\right) - 2\boldsymbol{b}'\mathbb{E}(XY) + \boldsymbol{b}'\mathbb{E}\left(XX'\right)\boldsymbol{b}.$$

▶ By matrix calculus, the first-order condition for minimization is

$$\boldsymbol{0} = \frac{\partial}{\partial \boldsymbol{b}} S(\boldsymbol{b}) = -2\mathbb{E}(XY) + 2\mathbb{E}\left(XX'\right)\boldsymbol{b}.$$

Solving for the first-order condition, $\boldsymbol{\beta} = \boldsymbol{Q}_{XX}^{-1}\boldsymbol{Q}_{XY}$ where $\boldsymbol{Q}_{XY} = \mathbb{E}(XY)$ is $k \times 1$ and $\boldsymbol{Q}_{XX} = \mathbb{E}(XX')$ is $k \times k$.

▶ We now have an explicit expression for the best linear predictor:

$$\mathcal{P}(Y \mid X) = X'\left(\mathbb{E}\left(XX'\right)\right)^{-1}\mathbb{E}(XY).$$

This expression is also referred to as the **linear projection** of $Y$ on $X$.

- ▶ The **projection error** is

$$e = Y - X'\beta.$$

- ▶ Rewriting, we obtain a decomposition of $Y$ into linear predictor and error

$$Y = X'\beta + e.$$

- ▶ An important property of the projection error is

$$\begin{aligned}
\mathbb{E}(Xe) &= \mathbb{E}\left(X\left(Y - X'\beta\right)\right) \\
&= \mathbb{E}(XY) - \mathbb{E}\left(XX'\right)\left(\mathbb{E}\left(XX'\right)\right)^{-1}\mathbb{E}(XY) \\
&= \mathbf{0}.
\end{aligned}$$

Theorem (Properties of Linear Projection Model)

1. *The Linear Projection Coefficient equals*

$$\boldsymbol{\beta} = \left(\mathbb{E}\left(\boldsymbol{X}\boldsymbol{X}'\right)\right)^{-1}\mathbb{E}\left(\boldsymbol{X}Y\right).$$

2. *The best linear predictor of $Y$ given $\boldsymbol{X}$ is*

$$\mathcal{P}\left(Y \mid \boldsymbol{X}\right) = \boldsymbol{X}'\left(\mathbb{E}\left(\boldsymbol{X}\boldsymbol{X}'\right)\right)^{-1}\mathbb{E}\left(\boldsymbol{X}Y\right).$$

3. *The projection error $e = Y - \boldsymbol{X}'\boldsymbol{\beta}$ satisfies*

$$\mathbb{E}\left(\boldsymbol{X}e\right) = \boldsymbol{0}.$$

4. *If $\boldsymbol{X}$ contains an constant, then*

$$\mathbb{E}\left(e\right) = 0.$$

# Best linear approximation

▶ We start by defining the mean-square approximation error of $X'b$ to $m(X)$ as the expected squared difference between $X'b$ and the conditional mean $m(X)$:

$$d(b) = \mathbb{E}\left((m(X) - X'b)^2\right) = \int_{\mathbb{R}^k} \left(m(x) - x'b\right)^2 f_X(x)\, dx.$$

▶ We can then define the best linear approximation to the conditional mean $m(X)$ as the function $X'\beta$ obtained by selecting $\beta$ to minimize $d(b)$:

$$\beta = \underset{b \in \mathbb{R}^k}{\mathrm{argmin}}\, d(b).$$

▶ It turns out that the best linear predictor and the best linear approximation are identical:

$$
\begin{aligned}
\beta &= (\mathbb{E}(XX'))^{-1}\, \mathbb{E}(Xm(X)) \\
&= (\mathbb{E}(XX'))^{-1}\, \mathbb{E}(XY).
\end{aligned}
$$

# Causal effects

- ► Consider the effect of schooling on wages. The causal effect is the actual difference a person would receive in wages if we could change their level of education holding all else constant.

- ► The causal effect is unobserved because the most we can observe is their actual level of education and their actual wage, but not the counterfactual wage if their education had been different.

- ► A variable $X_1$ can be said to have a causal effect on the response variable if the latter changes when all other inputs are held constant.

- ► A full model:

$$Y = h\left(X_1, \boldsymbol{X}_2, \boldsymbol{U}\right),$$

where $X_1$ and $\boldsymbol{X}_2$ are observed variables, $\boldsymbol{U}$ is some unobserved random factor and $h$ is a functional relationship.

Definition
The causal effect of $X_1$ on $Y$ is

$$C(X_1, \boldsymbol{X}_2, \boldsymbol{U}) = \nabla_1 h(X_1, \boldsymbol{X}_2, \boldsymbol{U})$$

the change in $Y$ due to a change in $X_1$, holding $\boldsymbol{X}_2$ and $\boldsymbol{U}$ constant.

▶ We define the causal effect of $X_1$ within this model as the change in due to a change in $X_1$ holding the other variables $\boldsymbol{X}_2$ and $\boldsymbol{U}$ constant.

# Average causal effect

Definition
The average causal effect of $X_1$ on $Y$ conditional on $X_2$ is

$$ACE(X_1, X_2) = \mathbb{E}\left(C(X_1, X_2, U) \mid X_1, X_2\right)$$
$$= \int \nabla_1 h(X_1, X_2, u) \, f_{U|X_1,X_2}(u \mid X_1, X_2) \, du$$

where $f_{U|X_1,X_2}(u \mid x_1, x_2)$ is the conditional density of $U$ given $X_1, X_2$.

# Linear causal model

- Suppose that $h(X_1, X_2, U) = g(X_1, X_2) + U$ ($U$ is one-dimensional). Then, $C(X_1, X_2, U) = \nabla_1 g(X_1, X_2)$ and $ACE(X_1, X_2) = \nabla_1 g(X_1, X_2)$.

- The linear causal model specifies that $g(X_1, X_2)$ is linear in parameters. E.g., $g(X_1, X_2) = \gamma X_1 + X_2'\beta$ ($\nabla_1 g(X_1, X_2) = \gamma$) or $g(X_1, X_2) = \gamma X_1 + X_2'\beta + X_1 X_2'\delta$ ($\nabla_1 g(X_1, X_2) = \gamma + X_2'\delta$), where $\gamma, \beta, \delta$ are unknown parameters.

- Without additional assumptions, it is impossible to recover the parameters $\gamma, \beta, \delta$.

- If we assume the stronger condition $\mathbb{E}(U \mid X_1, X_2) = 0$, then $g(X_1, X_2) = \mathbb{E}(Y \mid X_1, X_2)$.

- If $g(X_1, X_2) = \gamma X_1 + X_2'\beta$ and a weaker condition $\mathbb{E}(X_1 U) = 0$ and $\mathbb{E}(X_2 U) = 0$ holds, then $g(X_1, X_2) = \mathcal{P}(Y \mid X_1, X_2)$.

- Under either stronger or weaker condition, the unobservable factor $U$ is uncorrelated with the observable factors.

# Omitted variable bias

▶ Consider the projection of $Y$ on $X_1$ only:

$$Y = X_1'\gamma_1 + e$$
$$\mathbb{E}(X_1 e) = \mathbf{0}.$$

▶ Suppose that the linear causal model holds:
$Y = X_1'\beta_1 + X_2'\beta_2 + U$, with $\mathbb{E}(X_1 U) = 0$ and $\mathbb{E}(X_2 U) = 0$.

▶ Typically, $\beta_1 \neq \gamma_1$:

$$\begin{aligned}
\gamma_1 &= \left(\mathbb{E}(X_1 X_1')\right)^{-1} \mathbb{E}(X_1 Y) \\
&= \left(\mathbb{E}(X_1 X_1')\right)^{-1} \mathbb{E}\left(X_1 \left(X_1'\beta_1 + X_2'\beta_2 + U\right)\right) \\
&= \beta_1 + \left(\mathbb{E}(X_1 X_1')\right)^{-1} \mathbb{E}(X_1 X_2') \beta_2 \\
&\neq \beta_1
\end{aligned}$$

unless $\left(\mathbb{E}(X_1 X_1')\right)^{-1} \mathbb{E}(X_1 X_2') = \mathbf{0}$ or $\beta_2 = \mathbf{0}$.