

Advanced Econometrics

Lecture 2: Conditional Expectation and Projection (Hansen Chapter 2)

Instructor: Ma, Jun

Renmin University of China

Fall 2018

Conditional Expectation Function

- Conditional expectations can be written with the generic notation

$$\mathbb{E}(Y \mid X_1, X_2, \dots, X_k) = m(X_1, X_2, \dots, X_k).$$

We call this the conditional expectation function (CEF). The CEF is a function of (X_1, X_2, \dots, X_k) as it varies with the variables.

- For greater compactness, we will typically write the conditioning variables as a vector in \mathbb{R}^k :

$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{pmatrix}.$$

Given this notation, the CEF can be compactly written as

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

Conditional Expectation Function

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

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

- ▶ For any \mathbf{x} such that $f_X(\mathbf{x}) > 0$ the conditional density of Y given X is defined as

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

- ▶ The CEF of Y given $X = \mathbf{x}$ is the mean of the conditional density

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

Intuitively, $m(\mathbf{x})$ is the mean of for the idealized subpopulation where the conditioning variables are fixed at \mathbf{x} .

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

Law of Iterated Expectations

Theorem (Simple Law of Iterated Expectations)

If $\mathbb{E}|Y| < \infty$ then for any random vector \mathbf{X} ,

$$\mathbb{E}(\mathbb{E}(Y | \mathbf{X})) = \mathbb{E}(Y).$$

When \mathbf{X} is discrete

离散

$$\mathbb{E}(\mathbb{E}(Y | \mathbf{X})) = \sum_{j=1}^{\infty} \mathbb{E}(Y | \mathbf{x}_j) \Pr(\mathbf{X} = \mathbf{x}_j)$$

and when \mathbf{X} is continuous

连续

$$\mathbb{E}(\mathbb{E}(Y | \mathbf{X})) = \int_{\mathbb{R}^k} \mathbb{E}(Y | \mathbf{x}) f_{\mathbf{X}}(\mathbf{x}) d\mathbf{x}.$$

Law of Iterated Expectations

Theorem

If $\mathbb{E}|y| < \infty$ then for any random vectors \mathbf{X}_1 and \mathbf{X}_2 ,

$$\mathbb{E}(\mathbb{E}(Y \mid \mathbf{X}_1, \mathbf{X}_2) \mid \mathbf{X}_1) = \mathbb{E}(Y \mid \mathbf{X}_1)$$

Law of Iterated Expectations

A property of conditional expectations is that when you condition on a random vector \mathbf{x} you can effectively treat it as if it is constant. For example, $\mathbb{E}(\mathbf{X}|\mathbf{X}) = \mathbf{X}$ and $\mathbb{E}(g(\mathbf{X})|\mathbf{X}) = g(\mathbf{X})$ for any function $g(\cdot)$. The general property is known as the Conditioning Theorem.

求关于 \mathbf{X} 的条件期望, 那么
所有有关 \mathbf{X} 的都看成一个常数.

Theorem (Conditioning Theorem)

If $\mathbb{E}|y| < \infty$ then

$$\mathbb{E}(g(\mathbf{X})y|\mathbf{X}) = g(\mathbf{X})\mathbb{E}(y|\mathbf{X})$$

If in addition $\mathbb{E}|g(\mathbf{X})y| < \infty$, then

$$\mathbb{E}(g(\mathbf{X})y) = \mathbb{E}(g(\mathbf{X})\mathbb{E}(y|\mathbf{X})).$$

In this course, we can safely ignore conditions such as $\mathbb{E}|Y| < \infty$ and $\mathbb{E}|g(\mathbf{X})Y| < \infty$.

CEF Error

- The CEF error is defined as the difference between and the CEF evaluated at the random vector X :

$$\text{CEF error} \leftarrow e = Y - \underline{m(X)} \rightarrow \text{条件期望函数(CEF)}$$

By construction, this yields the formula

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

- A key property of the CEF error is that it has a conditional mean of zero. To see this, by the linearity of expectations, the definition $m(X) = \mathbb{E}(Y|X)$ and the Conditioning Theorem

$$E(X+Y|Z) = E(X|Z) + E(Y|Z)$$

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

- The unconditional mean is also zero:

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

$$X \perp\!\!\!\perp Y \Rightarrow E(X|Y) = E(X)$$

$$\Rightarrow \text{cov}(X, Y) = 0$$



$$E(XY) = E(X)E(Y)$$

$$= E E(XY|Y)$$

$$= E[Y \cdot E(X|Y)]$$

$$= E(Y)E(X)$$

$$E(X|Y) = E(X)$$

$$E(X^2|Y) = \text{function of } Y$$

CEF Error

Theorem

Properties of the CEF error

If $\mathbb{E}|y| < \infty$ then

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

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

3. For any function $h(X)$ such that $\mathbb{E}|h(x)e| < \infty$ then

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

$$\begin{aligned}\mathbb{E}(h(X)e) &= \mathbb{E}(\mathbb{E}(h(X)e | X)) \\ &= \mathbb{E}(h(X)\mathbb{E}(e | X)) \\ &= \mathbb{E}(0) = 0\end{aligned}$$

The equations

$$y = m(X) + e$$

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

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.

CEF Error

- ▶ The equation $\mathbb{E}(e | X) = 0$ is sometimes called a conditional mean restriction, since the conditional mean of the error is restricted to equal zero. The property is also sometimes 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 $\mathbb{E}(e|X) = 0$ and e is mean independent of X , yet e is not fully independent of X . Mean independence does not imply full independence.

$$\epsilon \perp X \quad \epsilon \sim N(0, 1)$$

$$X \sim N(0, 1)$$

$$e = X \cdot \epsilon$$

$$\mathbb{E}(e|X) = \mathbb{E}(X\epsilon|X) = X \mathbb{E}(\epsilon|X) = 0$$

$$e|X \sim N(0, X^2)$$

(e 关于 X 的条件分布, X 可以看作一个常数. ϵ 是正态的, e 也是正态的)

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

$$\begin{aligned}\sigma^2 &= \text{Var}(e) = \mathbb{E}\left((e - \mathbb{E}e)^2\right) = \mathbb{E}\left(e^2\right). \\ &= \mathbb{E}(e^2 - 2e\mathbb{E}e + (\mathbb{E}e)^2) = \mathbb{E}e^2 - \underbrace{(\mathbb{E}e)^2}_{=0} = \mathbb{E}e^2\end{aligned}$$

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

$(Y, X) \in \mathbb{R} \times \mathbb{R}^k$

$$\left. \begin{aligned} m(X) &= \mathbb{E}(Y|X) \\ \underline{e = Y - m(X)} \end{aligned} \right\} \Rightarrow \mathbb{E}(e|X) = 0$$

\uparrow CEF error \downarrow $\text{Cov}(e, X) = 0$

只能被 X 解释的是 $m(X)$.

Regression Variance

- ▶ The regression variance depends on the regressors. Consider two regressions

$$Y = \mathbb{E}(Y \mid \mathbf{X}_1) + e_1$$

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

} \rightarrow 一定有 $e_2 \leq e_1$

- ▶ 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.

Jensen 不等式

$$g(\mathbb{E}(X)) \leq \mathbb{E}[g(X)]$$

if $g(\cdot)$ 是凸的.

eg. $g(x) = x^2$

Theorem

If $\mathbb{E}(y^2) < \infty$ then

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

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

这里把 $g(X)$ 作为 Y 的 predictor 的误差的大小.

$$\mathbb{E} \left((Y - g(X))^2 \right). \quad (1)$$

- We can define the best predictor as the function $g(X)$ which minimizes (1). What function is the best predictor? It turns out that the answer is the CEF. This holds regardless of the joint distribution of (Y, X) :

$m(X)$ 是所有预测结果中误差最小的一项.

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

$$\begin{aligned} &\mathbb{E} \left(e(m(X) - g(X)) \right) \\ &= \mathbb{E} \left(\mathbb{E} \left(e(m(X) - g(X)) \mid X \right) \right) \\ &= \mathbb{E} \left((m(X) - g(X)) \underbrace{\mathbb{E}(e \mid X)}_{=0} \right) \\ &= 0 \end{aligned}$$

Best Predictor

$$m(\cdot) = \arg \min_{g(\cdot)} \mathbb{E} (Y - g(X))^2$$

Theorem

Conditional Mean as Best Predictor

If $\mathbb{E}(Y^2) < \infty$, then for any predictor $g(X)$,

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

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

Conditional Variance

$$\text{Var}(W) = E(W - E(W))^2$$

把所有期望换成关于 X 的条件期望

$$\text{Var}(W|X) = E((W - E(W|X))^2 | X)$$

Definition

If $E(W^2) < \infty$, the conditional variance of W given X is

$$\text{Var}(W | X) = E\left((W - E(W | X))^2 | X\right)$$

Definition 回归误差 e 的条件方差

If $E(e^2) < \infty$, the conditional variance of the regression error e is

$$\sigma^2(X) = \text{Var}(e | X) = E(e^2 | X)$$

Conditional Variance

- ▶ Generally, $\sigma^2(\mathbf{X})$ is a non-trivial function of \mathbf{X} and can take any form subject to the restriction that it is non-negative.
- ▶ Notice as well that $\sigma^2(\mathbf{X}) = \text{Var}(Y|\mathbf{X})$ so it is equivalently the conditional variance of the dependent variable.
- ▶ We define the **conditional standard deviation** as its square root $\sigma(\mathbf{X}) = \sqrt{\sigma^2(\mathbf{X})}$.
- ▶ The unconditional error variance and the conditional variance are related by the law of iterated expectations

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

Conditional Variance

$$\mathbb{E}(g(X)e|X) = g(X) \mathbb{E}(e|X)$$

Conditioning Theorem

- 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

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

and

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

Thus ε has a conditional mean of zero, and a conditional variance of 1.

Homoskedasticity and Heteroskedasticity

Definition

同方差

是一个常数, 和 X 无关.

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

Definition

异方差

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

Some older or introductory textbooks describe heteroskedasticity as the case where “the variance of e varies across observations”. This is a poor and confusing definition. It is more constructive to understand that heteroskedasticity means that the conditional variance $\sigma^2(X)$ depends on observables.

Homoskedasticity and Heteroskedasticity

- ▶ 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, not the converse.
- ▶ We will still frequently impose the homoskedasticity assumption when making theoretical investigations into the properties of estimation and inference methods. The reason is that in many cases homoskedasticity greatly simplifies the theoretical calculations, and it is therefore quite advantageous for teaching and learning.

都认为同方差是正常的, 异方差是特例.

但现代观点认为异方差才是一般的情况, 同方差才是特例.

实证研究中都假设误差是异方差.

现在常用 White heteroskedasticity robust se. 怀特异方差稳健标准差

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, \dots, X_k fixed, as the partial derivative of the CEF $\frac{\partial}{\partial X_1} m(X_1, \dots, X_k)$.
- ▶ 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

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

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

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

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

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

Linear CEF

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

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

$$= \beta_0 + X_1\beta_1 + X_2\beta_2 + \cdots + X_k\beta_k$$

- An easy way to do so is to augment the regressor vector X by listing the number “1” as an element. We call this the “constant” and the corresponding coefficient is called the “intercept”:

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \\ 1 \end{pmatrix}.$$

← 截距项

With this redefinition, the CEF is

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

线性条件期望模型

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

Linear CEF

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

Linear CEF Model

$$y = \mathbf{X}'\boldsymbol{\beta} + e$$

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

Homoskedastic Linear CEF Model

$$y = \mathbf{X}'\boldsymbol{\beta} + e$$

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

$$\mathbb{E}(e^2 \mid \mathbf{X}) = \sigma^2$$

Linear CEF with Nonlinear Effects

- ▶ We can include as regressors nonlinear transformations of the original variables. In this sense, the linear CEF framework is flexible and can capture many nonlinear effects.
- ▶ The CEF could take the quadratic form

$$m(X_1, X_2) = X_1\beta_1 + X_2\beta_2 + X_1^2\beta_3 + X_2^2\beta_4 + X_1X_2\beta_5 + \beta_6.$$

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 β_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 .

线性指的是对参数 β 是线性的.

交互项 当 $\beta_5 > 0$, X_1 的边际效应对 X_2 是单调的.

教育 \times 性别 — 教育回报在性别间存在不平等.

Best Linear Predictor 最优线性预测

A linear predictor for is a function of the form $X'\beta$ for some $\beta \in \mathbb{R}^k$.

The mean squared prediction error is

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

Definition

The Best Linear Predictor of Y given X is

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

where β minimizes the mean squared prediction error

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

The minimizer

$$\beta = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} S(b) \quad \text{线性投影}$$

is called the Linear Projection Coefficient.

Best Linear Predictor

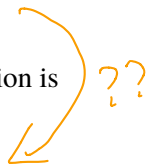
- By calculations,

$$S(\beta) = \mathbb{E}(Y^2) - 2\beta' \mathbb{E}(XY) + \beta' \mathbb{E}(XX') \beta.$$

$$\mathbb{E}(X' \beta)^2 = \mathbb{E}(\beta' X X' \beta) = \beta' \mathbb{E}(XX') \beta$$

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

$$\mathbf{0} = \frac{\partial}{\partial \beta} S(\beta) = -2\mathbb{E}(XY) + 2\mathbb{E}(XX') \beta.$$



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

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

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

最优线性近似

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

Best Linear Predictor

- 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

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

Best Linear Predictor

Theorem (Properties of Linear Projection Model)

1. *The Linear Projection Coefficient equals*

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

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

$$\mathcal{P}(Y | \mathbf{X}) = \mathbf{X}' (\mathbb{E}(\mathbf{X}\mathbf{X}'))^{-1} \mathbb{E}(\mathbf{X}Y).$$

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

投影误差 $\mathbb{E}(\mathbf{X}e) = \mathbf{0}.$

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

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

Best Linear Predictor

Linear Projection Model

$$y = \mathbf{x}'\boldsymbol{\beta} + e$$

$$\mathbb{E}(\mathbf{x}e) = \mathbf{0}$$

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{x}\mathbf{x}'))^{-1} \mathbb{E}(\mathbf{x}y)$$

不考

Regression Sub-Vectors

- Let the regressors be partitioned as

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}.$$

We can write the projection of Y on X as

$$\begin{aligned} y &= X'\beta + e \\ &= X_1'\beta_1 + X_2'\beta_2 + e \\ \mathbb{E}(Xe) &= \mathbf{0}. \end{aligned}$$

- Partition:

$$Q_{XX} = \begin{bmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{bmatrix} = \begin{bmatrix} \mathbb{E}(X_1 X_1') & \mathbb{E}(X_1 X_2') \\ \mathbb{E}(X_2 X_1') & \mathbb{E}(X_2 X_2') \end{bmatrix}$$

and

$$Q_{XY} = \begin{bmatrix} Q_{1Y} \\ Q_{2Y} \end{bmatrix} = \begin{bmatrix} \mathbb{E}(X_1 Y) \\ \mathbb{E}(X_2 Y) \end{bmatrix}.$$

$$X = \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \quad \begin{matrix} k \times 1 \\ k \times 1 \end{matrix} \quad \begin{matrix} k_1 + k_2 = k \end{matrix}$$

$$X' = \begin{pmatrix} X_1' & X_2' \end{pmatrix} \quad \begin{matrix} 1 \times k \\ 1 \times k_1 \\ 1 \times k_2 \end{matrix}$$

$$Q_{XX} = \mathbb{E}(XX')_{k \times k}$$

$$Q_{XY} = \mathbb{E}(XY)_{k \times 1}$$

Regression Sub-Vectors

- By the partitioned matrix inversion formula,

$$\mathbf{Q}_{XX}^{-1} = \begin{bmatrix} \mathbf{Q}_{11} & \mathbf{Q}_{12} \\ \mathbf{Q}_{21} & \mathbf{Q}_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \mathbf{Q}_{11.2}^{-1} & -\mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22.1}^{-1} \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22.1}^{-1} \end{bmatrix},$$

where $\mathbf{Q}_{11.2} = \mathbf{Q}_{11} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{21}$ and $\mathbf{Q}_{22.1} = \mathbf{Q}_{22} - \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} \mathbf{Q}_{12}$.

- Thus,

$$\begin{aligned} \beta &= \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} \\ &= \begin{bmatrix} \mathbf{Q}_{11.2}^{-1} & -\mathbf{Q}_{11.2}^{-1} \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \\ -\mathbf{Q}_{22.1}^{-1} \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} & \mathbf{Q}_{22.1}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{Q}_{1Y} \\ \mathbf{Q}_{2Y} \end{bmatrix} \\ &= \begin{pmatrix} \mathbf{Q}_{11.2}^{-1} \left(\mathbf{Q}_{1Y} - \mathbf{Q}_{12} \mathbf{Q}_{22}^{-1} \mathbf{Q}_{2Y} \right) \\ \mathbf{Q}_{22.1}^{-1} \left(\mathbf{Q}_{2Y} - \mathbf{Q}_{21} \mathbf{Q}_{11}^{-1} \mathbf{Q}_{1Y} \right) \end{pmatrix} \end{aligned}$$

$\beta = \mathbf{Q}_{XX}^{-1} \mathbf{Q}_{XY}$

Coefficient Decomposition

分块回归

- ▶ $\beta_1 \in \mathbb{R}$ and

$$Y = X_1\beta_1 + X_2'\beta_2 + e.$$

X_1 是第一个变量

X_2 是其他变量的向量

- ▶ Now consider the projection of X_1 on X_2 :

$$X_1 = X_2'\gamma_2 + U_1$$

$$\mathbb{E}(X_2 U_1) = \mathbf{0}.$$

γ_2 是 X_1 在 X_2 上的线性投影.

$$\gamma_2 = \underbrace{\mathbb{E}(X_2 X_2')^{-1}}_{Q_{22}} \underbrace{\mathbb{E}(X_2 X_1)}_{Q_{21}}$$

- ▶ $\gamma_2 = Q_{22}^{-1} Q_{21}$ and

$$\begin{aligned}\mathbb{E}(U_1^2) &= \mathbb{E}\left((X_1 - X_2'\gamma_2)^2\right) \\ &= \mathbb{E}(X_1^2) - 2\mathbb{E}(X_1 X_2')\gamma_2 + \gamma_2'\mathbb{E}(X_2 X_2')\gamma_2 \\ &= Q_{11} - 2Q_{12}Q_{22}^{-1}Q_{21} + Q_{12}Q_{22}^{-1}Q_{22}Q_{22}^{-1}Q_{21} \\ &= Q_{11} - Q_{12}Q_{22}^{-1}Q_{21}.\end{aligned}$$

Coefficient Decomposition

- ▶ Also calculate

$$\mathbb{E}(U_1 Y) = \mathbb{E}((X_1 - \gamma_2' X_2) Y) = \boldsymbol{Q}_{1Y} - \boldsymbol{Q}_{12} \boldsymbol{Q}_{22}^{-1} \boldsymbol{Q}_{2Y}.$$

- ▶ We found

$$\beta_1 = \boldsymbol{Q}_{11 \cdot 2}^{-1} \left(\boldsymbol{Q}_{1Y} - \boldsymbol{Q}_{12} \boldsymbol{Q}_{22}^{-1} \boldsymbol{Q}_{2Y} \right) = \frac{\mathbb{E}(u_1 y)}{\mathbb{E}(u_1^2)}$$

the coefficient from the simple regression of Y on U_1 .

Omitted Variable Bias

- Consider the projection of Y on X_1 only:

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

- Typically, $\beta_1 \neq \gamma_1$:

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

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

$$\mathbb{E}(eX) = 0 \Leftrightarrow \mathbb{E}\begin{pmatrix} eX_1 \\ eX_2 \end{pmatrix} = 0$$

KX1

$$\Rightarrow \mathbb{E}(eX_1) = 0$$

重要

Best Linear Approximation

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

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

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

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

- It turns out that the best linear predictor and the best linear approximation are identical: by law of iterated expectation,

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

min $E((Y - X'b)^2)$
 $b \in \mathbb{R}^k$
 $m(X)$

$$\begin{aligned} E(XY) &= E(E(XY|X)) \\ &= EX E(Y|X) \\ &= E(Xm(X)) \end{aligned}$$

$$\mathcal{P}(Y|X) = X'E(XX')^{-1}E(XY)$$

① Y 对 X 的投影

② 对 $E(Y|X)$ 的最佳线性近似

Random Coefficient Model

- ▶ A linear random coefficient model takes the form

$$Y = X'\eta,$$

where the individual-specific coefficient η is random and independent of X .

- ▶ It is interesting to discover that the linear random coefficient model implies a linear CEF. Let

$$\beta = \mathbb{E}(\eta)$$

$$\Sigma = \text{Var}(\eta)$$

and decompose $\eta = \beta + U$. Now U is distributed independently of X with mean zero and covariance matrix Σ .

- ▶ Then

$$\mathbb{E}(Y | X) = X'\mathbb{E}(\eta | X) = X'\mathbb{E}(\eta) = X'\beta$$

so the CEF is linear in X , and the coefficients β equal the mean of the random coefficient η .

Random Coefficient Model

Theorem

In the linear random coefficient model $Y = X'\boldsymbol{\eta}$ with $\boldsymbol{\eta}$ independent of \mathbf{x} , then

$$\mathbb{E}(Y \mid \mathbf{X}) = \mathbf{X}'\boldsymbol{\beta}$$

$$\text{Var}(e \mid \mathbf{X}) = \mathbf{X}'\boldsymbol{\Sigma}\mathbf{X}$$

where $\boldsymbol{\beta} = \mathbb{E}(\boldsymbol{\eta})$, $\boldsymbol{\Sigma} = \text{Var}(\boldsymbol{\eta})$ and $e = Y - \mathbf{X}'\boldsymbol{\beta}$. So the error is conditionally heteroskedastic with its variance a quadratic function of \mathbf{X} .

Causal Effects 因果效应

因果关系是单向的.

- ▶ 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:

$$Y = h(X_1, X_2, U),$$

Y 的数据生成过程. X_1, X_2 被解释. U 观测不到的因素.

where X_1 and X_2 are observed variables, U is some unobserved random factor and h is a functional relationship.

- ▶ This framework is called the **potential outcomes** framework.

Causal Effects

Definition

In the model (2.52) the causal effect of x_1 on y is

$$\text{因果效应} \Rightarrow C(x_1, \mathbf{x}_2, \mathbf{u}) = \nabla_1 h(x_1, \mathbf{x}_2, \mathbf{u})$$

the change in y due to a change in x_1 , holding \mathbf{x}_2 and \mathbf{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 \mathbf{X}_2 and \mathbf{U} constant.

Causal Effects

- ▶ A popular example arises in the analysis of treatment effects with a binary regressor X_1 . Let $X_1 = 1$ indicate treatment (e.g. a medical procedure) and $X_1 = 0$ indicate non-treatment.
- ▶ In this case: $Y(0) = h(0, X_2, U)$ and $Y(1) = h(1, X_2, U)$.
- ▶ $Y(0)$ and $Y(1)$ are the latent outcomes associated with non-treatment and treatment. That is, for a given individual, $Y(0)$ is the health outcome if there is no treatment, and $Y(1)$ is the health outcome if there is treatment.
- ▶ The causal effect of treatment for the individual is the change in their health outcome due to treatment — the change in as we hold both X_2 and U constant:

$$C(X_2, U) = Y(1) - Y(0).$$

- ▶ In a sample, we cannot observe both outcomes from the same individual, we only observe the realized value

$$Y = \begin{cases} Y(0) & \text{if } X_1 = 0 \\ Y(1) & \text{if } X_1 = 1. \end{cases}$$

treatment effect 就是 Cause effect 中的一个特例。只不过 X_1 是离散的。

不能同时观测到 $Y(0)$ 、 $Y(1)$ ，是同一个体的两种不同状态。

Average Causal Effect

$$C(X_1, X_2, U) = \nabla_1 h(X_1, X_2, U)$$

$$\mathbb{E}[C(X_1, X_2, U) | X_1, X_2]$$

$$= ACE(X_1, X_2) \text{ 平均因果效应}$$

Definition

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

$$\begin{aligned} ACE(X_1, X_2) &= \mathbb{E}(C(X_1, X_2, U) | X_1, X_2) \\ &= \int \nabla_1 h(X_1, X_2, u) f_{U|X_1, X_2}(u | X_1, X_2) du \end{aligned}$$

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

What is the relationship between the average causal effect $ACE(X_1, X_2)$ and the regression derivative $\nabla_1 m(X_1, X_2)$?

Average Causal Effect

- Since $Y = h(X_1, X_2, U)$, the CEF is

$$\begin{aligned} m(X_1, X_2) &= \mathbb{E}(h(X_1, X_2, U) \mid X_1, X_2) \\ &= \int h(X_1, X_2, u) f_{U|X_1, X_2}(u \mid X_1, X_2) du \end{aligned}$$

the average causal equation, averaged over the conditional distribution of the unobserved component U .

- The regression derivative is

$$\begin{aligned} \nabla_1 m(X_1, X_2) &= \int \nabla_1 h(X_1, X_2, u) f_{U|X_1, X_2}(u \mid X_1, X_2) du \\ &\quad + \int h(X_1, X_2, u) \nabla_1 f_{U|X_1, X_2}(u \mid X_1, X_2) du \\ &= ACE(X_1, X_2) \\ &\quad + \int h(X_1, X_2, u) \nabla_1 f_{U|X_1, X_2}(u \mid X_1, X_2) du. \end{aligned}$$

- The regression derivative and ACE equal in the special case when $\nabla_1 f_{U|X_1, X_2}(u \mid X_1, X_2) = 0$, that is, when the conditional density of U given X_1, X_2 : $f_{U|X_1, X_2}(u \mid x_1, x_2)$ does not depend on x_1 .

$$\begin{aligned} \frac{\partial m(X_1, X_2)}{\partial X_1} &= \frac{\partial}{\partial X_1} \int h(X_1, X_2, u) f(u \mid X_1, X_2) du \\ &= \int \frac{\partial}{\partial X_1} (h(X_1, X_2, u) f(u \mid X_1, X_2)) du \\ &= \int \left\{ \frac{\partial h(X_1, X_2, u)}{\partial X_1} f(u \mid X_1, X_2) + h(X_1, X_2, u) \frac{\partial f}{\partial X_1} \right\} du \end{aligned}$$

Average Causal Effect 条件独立假设

Definition (Conditional Independence Assumption(CIA))

Conditional on X_2 , the random variables X_1 and U are statistically independent.

- ▶ Like (unconditional) independence ($f_{U|X_1} = f_U$), the conditional independence means $f_{U|X_1, X_2} = f_{U|X_2}$ and thus $\nabla_1 f_{U|X_1, X_2}(u | x_1, x_2) = 0$.
- ▶ Thus CIA implies $\nabla_1 m(X_1, X_2) = ACE(X_1, X_2)$.
- ▶ CIA is weaker than full independence of U from the regressors X_1, X_2 : $f_{U|X_1, X_2} = f_U \implies f_{U|X_1, X_2} = f_{U|X_2}$.

$$X \perp\!\!\!\perp W \Leftrightarrow f_{X|W} = f_X$$

$$X \perp\!\!\!\perp W | Z \Leftrightarrow \underbrace{f_{X|W, Z}}_{\text{Swartz}} = f_{X|Z}$$

$$f_{U|X_1, X_2} = f_U \implies f_{U|X_2} = f_U \implies \underbrace{f_{U|X_2} = f_{U|X_1, X_2}}_{\text{CIA}}$$

$$\begin{aligned}
 & \int \\
 f_{u, x_2}(u, x_2) &= \int f_{u, x_1, x_2}(u, x_1, x_2) dx_1 \\
 &= \int f_u(u) f_{x_1, x_2}(x_1, x_2) dx_1 \\
 &= f_u(u) \underbrace{\int f_{x_1, x_2}(x_1, x_2) dx_1}_{f_{x_2}(x_2)}
 \end{aligned}$$