Advanced Econometrics I

Jürgen Meinecke Lecture 2 of 12

Research School of Economics, Australian National University

Projections (rinse and repeat) Linear Projections in L_2 – General Case

Ordinary Least Squares Estimation

Given a bunch of random variables X_1, \ldots, X_K, Y , we wanted to express Y as a linear combination in X_1, \ldots, X_K

A fancy way of saying the same thing: We want to project *Y* onto the subspace spanned by *X*₁,...,*X*_{*K*}

That projection is labeled $\mathbb{P}_{\operatorname{sp}(X_1,\ldots,X_K)}Y$ or \hat{Y}

Instead of $\mathbb{P}_{\text{sp}(X_1,...,X_K)}$, we may simply write \mathbb{P}_X , where $X := (X_1,...,X_K)'$

(Aside: the X_i can enter non-linearly, for example $X_2 := X_1^2$)

Viewing X_1, \ldots, X_K, Y as elements of a Hilbert space, we learned the generic characterization using the inner product:

Using the orthonormal basis $\tilde{X}_1, \dots, \tilde{X}_K$ (such that sp $(\tilde{X}_1, \dots, \tilde{X}_K) = sp(X)$)

$$\begin{split} \widehat{Y} &= \mathbb{P}_X Y = \sum_{i=1}^K \langle \widetilde{X}_i, Y \rangle \widetilde{X}_i \\ &= \sum_{i=1}^K \mathbb{E}(\widetilde{X}_i \cdot Y) \widetilde{X}_i \\ &= \sum_{i=1}^K \beta_i^* X_i \end{split}$$

For example, when $X_1 = 1$ (constant term) and K = 2, we saw $\beta_2^* = \frac{\text{Cov}(X_2, Y)}{\text{Var } (X_2)}$ $\beta_1^* = \text{E}Y - \beta_2^* \text{E}X_2$ For general *K*, we use matrices to express $\beta^* := (\beta_1, \dots, \beta_K)'$

Let
$$X := (X_1, X_2, \dots, X_K)'$$
 be a $K \times 1$ vector
 $\widehat{Y} = \mathbb{P}_X Y = \sum_{i=1}^K \beta_i^* X_i \stackrel{(\star)}{=} X' \beta^*$,

where $\beta^* := (E(XX'))^{-1} E(XY)$ is a $K \times 1$ vector

Aside: equality (*) above justified by this bit of linear algebra: $\sum_{i=1}^{K} x_i y_i = x' y = y' x_i$

for generic vectors *x* and *y*:

$$x := \begin{pmatrix} x_1 \\ \vdots \\ x_K \end{pmatrix}, \qquad y := \begin{pmatrix} y_1 \\ \vdots \\ y_K \end{pmatrix},$$

or written compactly as $x := (x_1, \dots, x_K)'$ and $y := (y_1, \dots, y_K)'$

When $X_1 = 1$, β^* can be expressed via covariances

Corollary

When $X = (1, X_2, ..., X_K)'$, then the projection coefficients are $(\beta_2^*, ..., \beta_K^*)' = \Sigma_{XX}^{-1} \Sigma_{XY}$ $\beta_1^* = EY - \beta_2^* EX_2 - \dots - \beta_K^* EX_K$,

where

$$\Sigma_{XX} := \begin{bmatrix} \sigma_2^2 & \sigma_{23} & \dots & \sigma_{2K} \\ \sigma_{32} & \sigma_3^2 & \dots & \sigma_{3K} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{K2} & \sigma_{K3} & \dots & \sigma_K^2 \end{bmatrix} \text{ and } \Sigma_{XY} := \begin{bmatrix} \text{Cov}(X_2, Y) \\ \text{Cov}(X_3, Y) \\ \vdots \\ \text{Cov}(X_K, Y) \end{bmatrix}$$

 Σ_{XX} is matrix that collects variances of X on the diagonal and covariances on the off-diagonal

 Σ_{XY} is vector that collects covariances between X and Y

Linear projection representation of *Y*:

$$Y = \mathbb{P}_X Y + (Y - \mathbb{P}_X Y)$$

=: $\mathbb{P}_X Y + u$
= $X'\beta^* + u$,

The element Y of a Hilbert space can be reached by adding two elements:

- an element $X'\beta^*$ from the subspace sp(X) ;
- an element u that is orthogonal to sp(X)

Proof that u is orthogonal to sp(X), that is, E(Xu) = 0:

$$\begin{split} \mathsf{E}(X(Y - \mathbb{P}_X Y)) &= \mathsf{E}(X(Y - X' \cdot \mathsf{E}(XX')^{-1}\mathsf{E}(XY))) \\ &= \mathsf{E}(XY - XX' \cdot \mathsf{E}(XX')^{-1}\mathsf{E}(XY)) \\ &= \mathsf{E}(XY) - \mathsf{E}(XX')\mathsf{E}(XX')^{-1}\mathsf{E}(XY) \\ &= 0 \end{split}$$

Using the linear projection representation $Y = X'\beta^* + u$

Once you learn that E(Xu) = 0 you know that β^* must be the projection coefficient

You have learned that it exists and is unique

It is important to understand that the definition of the linear projection model is not restrictive

In particular, E(uX) = 0 is not an assumption, it is *definitional*

To drive home this point, suppose I claim

 $Y = X'\theta + w$

Next I tell you that E(wX) = 0

You therefore conclude that $\theta = \beta^* = (E(XX'))^{-1} E(XY)$

In summary

Definition (Linear Projection Model)

Given

(i) $X_1, \ldots, X_K, Y \in L_2$

(ii) E(XX') > 0 (positive definite)(aka, no perfect multicollinearity)

Then the **linear projection model** is given by $Y = X'\beta^* + u$,

where E(uX) = 0 and $\beta^* = (E(XX'))^{-1} E(XY)$.

We accept and understand now that the unique projection coefficient exists

Let's say we're interested in knowing the value of β^*

We just learned that $(\beta_2^*, \dots, \beta_K^*)' = \Sigma_{XX}^{-1} \Sigma_{XY}$

Do we know the objects on the rhs?

These are **population** variances and covariances We don't know these, therefore we don't know β^* How else could we quantify β^* ?

Projections (rinse and repeat)

Ordinary Least Squares Estimation The Problem of Estimation Definition of the OLS Estimator Basic Asymptotic Theory (part 1 of 2) Large Sample Properties of the OLS Estimato Let's indulge ourselves and take a short detour to think about estimation in an abstract way

This subsection is based on Stachurski A Primer in Econometric Theory chapters 8.1 and 8.2

We're dealing with a random variable Z with distribution P

We're interested in a *feature* of P

Definition (Feature)

Let $Z \in L_2$ and $P \in \mathcal{P}$ where \mathcal{P} is a class of distributions on Z. A **feature** of P is an object of the form $\gamma(P)$ for some $\gamma : \mathcal{P} \to S$.

Here S is an arbitrarily flexible space (usually \mathbb{R})

Examples of features: means, moments, variances, covariances

For some reason we are interested in $\gamma(P)$

If we knew P then we may be able to derive $\gamma(P)$

Example: P is standard normal and $\gamma(P) = \int ZdP = 0$ (mean of the standard normal distribution)

But we typically don't know P

If all we're interested in is $\gamma(P)$ then we may not need to know P (unless the feature we're interested in is P itself)

Instead, we use a random sample to make an inference about a feature of *P*

Definition (Random Sample)

The random variables Z_1, \ldots, Z_N are called a **random sample of** size N from the population P if Z_1, \ldots, Z_N are mutually independent and all have probability distribution P.

The joint distribution of $Z_1, ..., Z_N$ is P^N by independence We sometimes say that $Z_1, ..., Z_N$ are iid copies of ZWe sometimes say that $Z_1, ..., Z_N$ are iid random variables By the way: Z_i could be vectors or matrices too

Definition (Statistic)

A statistic is any function $g: \mathbb{R}^N \to \mathbb{R}$ that maps the sample data somewhere.

The definition of a statistic is deliberately broad

It is a function that maps the sample data somewhere Where to? Depends on the feature $\gamma(P)$ you're interested in There are countless examples

Illustration: let K = 1 (i.e., univariate) sample mean: $g(Z_1, ..., Z_N) = \sum_{i=1}^N Z_i/N =: \overline{Z}_N$ sample variance: $g(Z_1, ..., Z_N) = \sum_{i=1}^N (Z_i - \overline{Z}_N)^2/N$ sample min: $g(Z_1, ..., Z_N) = \min \{Z_1, ..., Z_N\}$ answer to everything: $g(Z_1, ..., Z_N) = 42$ A statistic becomes an estimator when linked to a feature $\gamma(P)$

Definition (Estimator)

An **estimator** $\hat{\gamma}$ is a statistic used to infer some feature $\gamma(P)$ of an unknown distribution *P*.

In other words: an estimator is a statistic with a purpose

Earlier example: P is the standard normal distribution (but let's pretend we don't know this, as is usually the case) So $Z \sim N(0, 1)$

And we're interested in EZ so we set $\gamma(P) = EZ = \int ZdP$ We have available a random sample $\{Z_1, \dots, Z_N\}$ Each $Z_i \sim N(0,1)$, but we don't know this But we do know: all Z_i are iid So they must all have the same mean EZ_i What would be an estimator for EZ? Aside: there are infinitely many What would be a *good* estimator for EZ_i ? (perhaps not so many anymore)

A good way to create estimators is the **analogy principle**

Goldberger explains the main idea of it:

the analogy principle of estimation...proposes that population parameters be estimated by sample statistics which have the same property in the sample as the parameters do in the population (Goldberger, 1968, as cited in Manski, 1988)

That is very unspecific, of course

Manski (1988) wrote an entire book on analog estimation and explains the analogy principle precisely and comprehensively

But we can illustrate it using our earlier framework

Definition (Empirical Distribution)

The **empirical distribution** P_N of the sample $\{Z_1, ..., Z_N\}$ is the discrete distribution that puts equal probability 1/N on each sample point Z_i , i = 1, ..., N.

Definition (Analogy Principle)

To estimate $\gamma(P)$ use $\hat{\gamma} := \gamma(P_N)$.

How do we use this in our example?

We wanted to estimate $\gamma(P) := \int Z dP$

According to the analogy principle, we should use $\int Z dP_N$

By definition, the empirical distribution is discrete therefore $\int ZdP_N = \sum_{i=1}^N Z_i/N =: \bar{Z}_N$

This is, of course, the sample average and we use the conventional notation \bar{Z}_{N}

The analogy principle results in the estimator $\hat{\gamma} = \sum_{i=1}^{N} Z_i / N$

How can we use the analogy principle to estimate β^* ?

Projections (rinse and repeat)

Ordinary Least Squares Estimation The Problem of Estimation Definition of the OLS Estimator

Basic Asymptotic Theory (part 1 of 2)

Large Sample Properties of the OLS Estimator

Recall linear projection representation $Y = X'\beta^* + u,$

where $X := (X_1, ..., X_K)'$, and $X_1, ..., X_K, Y \in L_2$ We saw that E(uX) = 0 implied $\beta^* = (E(XX'))^{-1} E(XY)$ In other words: β^* is the projection coefficient We want to estimate β^* using a random sample The use of a random sample necessitates some changes of notation... Until now, the symbol X represented the K-vector of random variables X_1, \ldots, X_K

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_K \end{pmatrix}$$

where, for example, X_1 was the constant (degenerate) random variable, X_2 was schooling, and so forth

Now we're given N copies of each of these, necessitating a double-subscript X_{ik} , where i = 1, ..., N and k = 1, ..., K

For example, the second copy (when i = 2) is

$$X_2 := \begin{pmatrix} X_{21} \\ X_{22} \\ \vdots \\ X_{2K} \end{pmatrix}$$

Do you see how this overloads the X_2 symbol?

In words: the new X₂ is the second copy of the original *K* random variables collecting all regressors

We'll also overload the X symbol like so:

$$X := \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1K} \\ X_{21} & X_{22} & \cdots & X_{2K} \\ X_{31} & X_{32} & \cdots & X_{3K} \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} & X_{N2} & \cdots & X_{NK} \end{pmatrix}$$

This is a big $N \times K$ matrix

Each row collects one copy of the original *K* random variables Example: the second row mirrors the transpose of our new *X*₂ Notice that

$$X := \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1K} \\ X_{21} & X_{22} & \cdots & X_{2K} \\ X_{31} & X_{32} & \cdots & X_{3K} \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} & X_{N2} & \cdots & X_{NK} \end{pmatrix} = \begin{pmatrix} X'_1 \\ X'_2 \\ X'_3 \\ \vdots \\ X'_N \end{pmatrix}$$

Alternatively, $X := (X_1, X_2, \dots, X_N)'$

Similarly, Y_1, \ldots, Y_N are iid copies of Y

Define $Y := (Y_1, Y_2, ..., Y_N)'$ to be the $N \times 1$ vector collecting all Y_i With this new notation, the random sample is $(X_1, Y_1), ..., (X_N, Y_N)$ Or simpler: $(X_i, Y_i), i = 1, ..., N$ is a random sample These are iid copies of the ordered pair (X, Y)Notice: I'm not pedantic enough to write $(X'_i, Y_i)'$ Given the random sample (X_i, Y_i) , i = 1, ..., N we can write the linear projection representation as

$$Y_i = X_i'\beta^* + u_i,$$

where we have $E(u_iX_i) = 0$

Combining findings from last lecture and assignment 1:

$$\beta^* = \operatorname*{argmin}_{b \in \mathbb{R}^K} \mathbb{E}\left(\left(Y_i - X'_i b\right)^2\right) \tag{1}$$

$$= \mathsf{E}(X_i X_i')^{-1} \mathsf{E}(X_i Y_i) \tag{2}$$

Equations (1) and (2) motivate two succinct analog estimators for β^* :

- (1) the ordinary least squares estimator;
- (2) the method of moments estimator

Let's look at both

If we define β^* like so: $\beta^* := \operatorname*{argmin}_{b \in \mathbb{R}^K} \mathbb{E}\left(\left(Y_i - X'_i b\right)^2\right),$

then the analogy principle suggests the estimator $\operatorname*{argmin}_{b\in\mathbb{R}^{K}}\sum_{i=1}^{N}\left(Y_{i}-X_{i}^{\prime}b\right)^{2}$

This seems very sensible and deserves a famous definition

Definition (Ordinary Least Squares (OLS) Estimator)

The ordinary least squares estimator is

$$\hat{\beta}^{\text{OLS}} := \underset{b \in \mathbb{R}^{K}}{\operatorname{argmin}} \sum_{i=1}^{N} \left(Y_{i} - X_{i}' b \right)^{2}$$

It is obvious how this estimator obtained its name

When you solve this you get

$$\hat{\beta}^{\text{OLS}} = \left(\frac{1}{N}\sum_{i=1}^{N} X_{i}X_{i}'\right)^{-1} \left(\frac{1}{N}\sum_{i=1}^{N} X_{i}Y_{i}\right)$$

Most people, when writing vectors, use the default column notation, meaning that if I tell you that X_i is a K-dimensional vector, you automatically know it is a $K \times 1$ vector

The second way of defining an estimator for β^* , via: $\beta^* = E(X_i X_i')^{-1} E(X_i Y_i)$

The analogy principle suggests the estimator

$$\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}'\right)^{-1}\frac{1}{N}\sum_{i=1}^{N}X_{i}Y_{i}$$

This also seems very sensible and deserves a familiar name:

Definition (Method of Moments (MM) Estimator)

Applying the analogy principle results in

$$\hat{\beta}^{\text{MM}} = \left(\frac{1}{N}\sum_{i=1}^{N} X_i X_i'\right)^{-1} \frac{1}{N}\sum_{i=1}^{N} X_i Y_i$$

You immediately see that $\hat{\beta}^{\rm OLS}=\hat{\beta}^{\rm MM}$

I'll simply refer to it as the OLS estimator

The OLS estimator does have a compact matrix representation

Using the $N \times K$ -matrix X and the $N \times 1$ -vector Y, we can replace the summation operator like so:

$$\sum X_i X'_i = X' X$$
$$\sum X_i Y_i = X' Y$$

It follows that $\hat{\beta}^{OLS}$ has a nice and short matrix representation: $\hat{\beta}^{OLS} = (X'X)^{-1}X'Y$

Now let's turn to the question: How good is $\hat{\beta}^{OLS}$?

What is goodness?

In the next few weeks we'll consider things such as

- bias
- variance (small sample and large sample)
- consistency
- distribution (large sample)

Projections (rinse and repeat)

Ordinary Least Squares Estimation The Problem of Estimation Definition of the OLS Estimator Basic Asymptotic Theory (part 1 of 2) Large Sample Properties of the OLS Est

Definition (Convergence in Probability)

A sequence of random variables $Z_1, Z_2, ...$ converges in probability to a real constant c if for every $\delta > 0$,

$$\lim_{N \to \infty} P(|Z_N - c| > \delta) = 0.$$

We say that *c* is the **probability limit** of Z_N and write $Z_N \xrightarrow{p} c$.

Often c=0 and the definition simplifies to $\lim_{N\to\infty}P(|Z_N|>\delta)=0$ for every $\delta>0$

Contrast this with the following definition

Definition (Bounded in Probability)

A sequence of random variables $Z_1, Z_2, ...$ is **bounded in** probability if there exists a finite real number $\delta > 0$ such that $\lim_{N \to \infty} P(|Z_N| > \delta) = 0.$ Let's unpack the previous definition:

- we're studying a sequence of probabilities $P(|Z_1| > \delta), \, P(|Z_2| > \delta), \ldots$
- we want to know its limit, does it get close to zero?
- this means, does there exist $\delta > 0$ and integer N_{ε} such that $P(|Z_N| > \delta) < \epsilon$ for every $\epsilon > 0$, and $N > N_{\varepsilon}$?

Let Z_1, Z_2, \dots be a sequence of iid standard normal random variables

- Recall that Φ is the cdf, and that $\Phi(-1.96)=0.025$
- · Also, by symmetry, $P(|Z_1| > \delta) = 2\Phi(-\delta)$
- Let's pick a small value for ϵ , say $\epsilon=0.05$
- Pick $\delta = 1.96 + 0.01 = 1.97$ which results in $P(|Z_1| > 1.97) < 0.05$
- I can make ϵ smaller and smaller, and I can always find δ that can bound the sequence of probabilities
- It follows that a sequence of iid standard normal random variables is indeed bounded in probability

That same sequence does not converge in probability to zero. Why?

Here's some new notation:

- a sequence Z_N is at most of order N^{λ} if $\frac{Z_N}{N^{\lambda}}$ is bounded in probability; we write $Z_N = O_p(N^{\lambda})$
- a sequence Z_N is order smaller than N^{λ} if $\frac{Z_N}{N^{\lambda}} \xrightarrow{p} 0$ (converges in probability to zero); we write $Z_N = o_p(N^{\lambda})$

The case in which $\lambda = 0$ occurs quite often

- if Z_N is bounded in probability we write $Z_N = O_p(1)$ and say that Z_N "is big Oh-p-one"
- if Z_N converges in probability to zero we write $Z_N = o_p(1)$ and say that Z_N "is little Oh-p-one"

The 'order' (Bachmann-Landau) notation is quite handy

Here some useful rules how to work with the new notation:

Lemma

Let c be a real constant. Then

$$c + o_p(1) = O_p(1)$$
$$c \cdot o_p(1) = o_p(1).$$

Lemma

Let
$$W_N = o_p(1)$$
, $X_N = o_p(1)$, $Y_N = O_p(1)$, and $Z_N = O_p(1)$.

$$\begin{split} W_N + X_N &= o_p(1) \qquad W_N + Y_N = O_p(1) \qquad Y_N + Z_N = O_p(1) \\ W_N \cdot X_N &= o_p(1) \qquad W_N \cdot Y_N = o_p(1) \qquad Y_N \cdot Z_N = O_p(1) \end{split}$$

All results here are quite intuitive

We've got a few more tricks up our sleeves

Theorem (Slutsky Theorem) If $Z_N = c + o_p(1)$ and $g(\cdot)$ is continuous at c then $g(Z_N) = g(c) + o_p(1)$.

In short: $g(c + o_p(1)) = g(c) + o_p(1)$

That's a reason to like the plim, it passes through nonlinear functions (which is not true for expectation operators)

Corollary

 $1/(c + o_p(1)) = 1/c + o_p(1)$ whenever $c \neq 0$.

All the definitions on the previous slides also apply element by element to sequences of random vectors or matrices

Theorem (Weak Law of Large Numbers (WLLN))

Let $Z_1, Z_2, ...$ be independent and identically distributed random variables with $EZ_i = \mu_Z$ and $Var Z_i = \sigma_Z^2 < \infty$. Then $\frac{1}{N} \sum_{i=1}^N Z_i = \mu_Z + o_p(1).$

Of course, $\bar{Z}_N := \frac{1}{N} \sum_{i=1}^N Z_i$ is the sample mean or sample average WLLN in words:

sample mean converges in probability to population mean

Proving the WLLN is easy, using Chebyshev's inequality

Notice that we want to show: $\lim_{N\to\infty} P\left(|\bar{Z}_N - \mu_Z| > \delta\right) = 0$, in other words: the sequence of probabilities $P\left(|\bar{Z}_1 - \mu_Z| > \delta\right)$, $P\left(|\bar{Z}_2 - \mu_Z| > \delta\right)$, $P\left(|\bar{Z}_3 - \mu_Z| > \delta\right)$, ... approaches zero

Lemma (Chebyshev's Inequality)

Let Z be a random variable with $EZ^2 < \infty$. Then for every c > 0 $P(|Z - \mu_Z| \ge c) \le \frac{Var(Z)}{c^2}$.

Idea: tail probabilities can be bounded via the variance

Recall from undergrad metrics: $E(\bar{Z}_N) = \mu_Z$ and $Var \bar{Z}_N = \sigma_Z^2/N$ We're interested in the sequence of probabilities $P(|\bar{Z}_N - \mu_Z| > \delta)$ Applying Chebyshev's inequality,

$$P\left(|\bar{Z}_N - \mu_Z| > \delta\right) \le \frac{\operatorname{Var} Z_N}{\delta^2} = \frac{\sigma_Z^2}{N \cdot \delta^2}$$

which converges to zero as $N \rightarrow \infty$

This takes us back to the analogy principle

Remember earlier:

We wanted to estimate the feature $\gamma(P) := EZ = \int ZdP$

According to the analogy principle, we should use $\int Z dP_N$

This led to the estimator $\hat{\gamma} = \sum_{i=1}^{N} Z_i / N$

Immediately by the WLLN: $\hat{\gamma} \xrightarrow{p} \gamma(P)$

Definition (Consistency of an Estimator)

An estimator $\hat{\gamma}$ for $\gamma := \gamma(P)$ is called **consistent** if $\hat{\gamma} \xrightarrow{p} \gamma$.

Intuition: if the sample size is large, sample mean is almost equal to population mean

So there is some hope that the analogy principle leads to consistent estimators

Projections (rinse and repeat)

Ordinary Least Squares Estimation The Problem of Estimation Definition of the OLS Estimator Basic Asymptotic Theory (part 1 of 2) Large Sample Properties of the OLS Estimator Let's first show that the QLS estimator is consistent

$$\hat{\beta}^{\text{OLS}} := \left(\sum_{i=1}^{N} X_i X'_i\right)^{-1} \sum_{i=1}^{N} X_i Y_i = \beta^* + \left(\frac{1}{N} \sum_{i=1}^{N} X_i X'_i\right)^{-1} \left(\frac{1}{N} \sum_{i=1}^{N} X_i u_i\right),$$

where we have used $Y_i = X'_i \beta^* + u_i$

Big picture to establish consistency: want to show that second term on rhs is close to zero (in a probabilistic sense)

Let's take a look

Copy and paste from previous slide:

$$\hat{\beta}^{\text{OLS}} = \beta^* + \left(\frac{1}{N}\sum_{i=1}^N X_i X_i'\right)^{-1} \left(\frac{1}{N}\sum_{i=1}^N X_i u_i\right)$$

Let's separately deal with $\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}'\right)^{-1}$ and $\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}u_{i}\right)$ By WLLN: $\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}' = E(X_{i}X_{i}') + o_{p}(1)$

But need to consider the asymptotic behavior of $\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}'\right)^{-1}$ Use Slutsky's theorem and pay attention:

$$\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}'\right)^{-1} = \left(\mathsf{E}(X_{i}X_{i}') + \mathsf{o}_{p}(1)\right)^{-1} = \mathsf{E}(X_{i}X_{i}')^{-1} + \mathsf{o}_{p}(1)$$

Pay attention: $E(X_i X'_i)^{-1}$ might not exist

This is analogous to a division-by-zero problem in the scalar world

A technical sufficient condition would be that $E(X_iX'_i)$ is positive definite; this guarantees invertibility, and results in:

Lemma

If $E(X_iX'_i)$ is positive definite then $E(X_iX'_i)^{-1} = O_p(1)$.

A more relatable sufficient condition, in words: none of the K - 1 regressors can be written as a linear function of the others (provided that a constant is included); aka, no perfect multicollinearity

Whenever you derive a consistency result, make sure to justify invertibility!

Copy and paste from two slides ago:

$$\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}X_{i}'\right)^{-1} = \left(\mathsf{E}(X_{i}X_{i}') + \mathsf{o}_{p}(1)\right)^{-1} = \mathsf{E}(X_{i}X_{i}')^{-1} + \mathsf{o}_{p}(1)$$

$$\stackrel{(\star)}{=} \mathsf{O}_{p}(1) + \mathsf{o}_{p}(1) = \mathsf{O}_{p}(1)$$

where (\star) is justified if $E(X_iX_i')$ is positive definite

Continuing our consistency proof, recall

$$\hat{\beta}^{\text{OLS}} = \beta^* + \left(\frac{1}{N}\sum_{i=1}^N X_i X_i'\right)^{-1} \left(\frac{1}{N}\sum_{i=1}^N X_i u_i\right)$$

Let's deal with $\left(\frac{1}{N}\sum_{i=1}^{N}X_{i}u_{i}\right)$

But this is easy, using the WLLN:

$$\frac{1}{N}\sum_{i=1}^{N} X_{i}u_{i} = \mathbb{E}(X_{i}u_{i}) + o_{p}(1) = 0 + o_{p}(1) = o_{p}(1)$$

Combining, we get:

$$\hat{\beta}^{\text{OLS}} = \beta^* + \left(\frac{1}{N}\sum_{i=1}^N X_i X_i'\right)^{-1} \left(\frac{1}{N}\sum_{i=1}^N X_i u_i\right)$$
$$= \beta^* + O_p(1) \cdot O_p(1)$$
$$= \beta^* + O_p(1)$$

In words: $\hat{\beta}^{OLS}$ converges in probability to β^*

This means $\hat{\beta}^{\rm OLS}$ is a consistent estimator for the projection coefficient β^*

It illustrates the benefit of the analogy principle when it works

But what is the distribution of $\hat{\beta}^{OLS}$?

- \cdot that's a tricky one
- $\hat{\beta}^{OLS} = \beta^* + (X'X)^{-1}X'u$, what's the distribution of the second term on the rhs?
- short answer: we have no idea
- there's some suspicion that $\hat{\beta}^{\rm OLS}$ may have an exact normal distribution if u is normally distributed
- but we don't know what the distribution of u is