

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, \dots, X_K, Y , we wanted to express Y as a linear combination in X_1, \dots, X_K

A fancy way of saying the same thing:

We want to project Y onto the subspace spanned by X_1, \dots, X_K

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

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

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

Viewing X_1, \dots, 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 $\text{sp}(\tilde{X}_1, \dots, \tilde{X}_K) = \text{sp}(X)$)

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

For example, when $X_1 = 1$ (constant term) and $K = 2$, we saw

$$\beta_2^* = \frac{\text{Cov}(X_2, Y)}{\text{Var}(X_2)}$$

$$\beta_1^* = EY - \beta_2^* EX_2$$

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

$$\hat{Y} = \mathbb{P}_X Y = \sum_{i=1}^K \beta_i^* X_i = X' \beta^*,$$

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

Aside: the last equality above is justified by the following result from linear algebra:

$$\sum_{i=1}^K x_i y_i = x' y = y' x,$$

for generic vectors x and y :

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

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 $X = (1, X_2, \dots, X_K)'$, then the projection coefficients are

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

$$\beta_1^* = EY - \beta_2^* EX_2 - \dots - \beta_K^* EX_K,$$

where

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

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 :

$$\begin{aligned} Y &= \mathbb{P}_X Y + (Y - \mathbb{P}_X Y) \\ &=: \mathbb{P}_X Y + u \\ &= X' \beta^* + u, \end{aligned}$$

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

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

Proof that u is orthogonal to $\text{sp}(X)$, that is, $E(Xu) = 0$:

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

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) $X_1, \dots, X_K, Y \in L_2$
- (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 Estimator

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

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 $\gamma(P) = \int Z dP = 0$
(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, \dots, Z_N are called a **random sample of size N** from the population P if Z_1, \dots, Z_N are mutually independent and all have probability distribution P .

The joint distribution of Z_1, \dots, Z_N is P^N by independence

We sometimes say that Z_1, \dots, Z_N are iid copies of Z

We sometimes say that Z_1, \dots, Z_N are iid random variables

By the way: Z_i could be vectors or matrices too

Definition (Statistic)

A **statistic** is any function $g : \mathbb{R}^N \rightarrow \mathbb{R}$ 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: $g(Z_1, \dots, Z_N) = \sum_{i=1}^N Z_i / N =: \bar{Z}_N$

sample variance: $g(Z_1, \dots, Z_N) = \sum_{i=1}^N (Z_i - \bar{Z}_N)^2 / N$

sample min: $g(Z_1, \dots, Z_N) = \min \{Z_1, \dots, Z_N\}$

answer to everything: $g(Z_1, \dots, Z_N) = 42$

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 $\gamma(P) = EZ = \int Z dP$

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)

Analogy Principle

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

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

$$\int Z dP_N = \sum_{i=1}^N Z_i / N =: \bar{Z}_N$$

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 $X := (X_1, \dots, X_K)'$, and $X_1, \dots, X_K, Y \in L_2$

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...

The random sample offers us N iid copies of the random vector X and the random variable Y

Letting i run from 1 to N , I define $X_i = (X_{i1}, \dots, X_{iK})'$

This gives us a collection of K -vectors, X_1, \dots, X_N

Notice that this effectively overwrites earlier notation:

- previously: X_1 denoted the first random variable in a collection of K random variables that span the subspace for the projection of Y ;
- from now on: X_1 is a K -dimensional column vector collecting iid copies of the K random variables that hitherto were labelled X_1, \dots, X_K

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

With this new notation, the random sample is $(X_1, Y_1), \dots, (X_N, Y_N)$

Or simpler: $(X_i, Y_i), i = 1, \dots, N$ is a random sample

These are iid copies of the ordered pair (X, Y)

Given the random sample $(X_i, Y_i), i = 1, \dots, N$ we can write the linear projection representation as

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

Because $E(u|X) = 0$ we have $E(u_i X_i) = 0$

Combining findings from last lecture and assignment 1:

$$\begin{aligned} \beta^* &= \operatorname{argmin}_{b \in \mathbb{R}^K} E \left((Y - X'b)^2 \right) \\ &= \operatorname{argmin}_{b \in \mathbb{R}^K} E \left((Y_i - X_i'b)^2 \right) \end{aligned} \quad (1)$$

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

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:

$$\beta^* := \operatorname{argmin}_{b \in \mathbb{R}^K} \mathbb{E} \left((Y_i - X_i' b)^2 \right),$$

then the analogy principle suggests the estimator

$$\operatorname{argmin}_{b \in \mathbb{R}^K} \sum_{i=1}^N (Y_i - X_i' b)^2$$

This seems very sensible and deserves a famous definition

Definition (Ordinary Least Squares (OLS) Estimator)

The ordinary least squares estimator is

$$\hat{\beta}^{\text{OLS}} := \operatorname{argmin}_{b \in \mathbb{R}^K} \sum_{i=1}^N (Y_i - X_i' b)^2$$

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 $\hat{\beta}^{\text{OLS}} = \hat{\beta}^{\text{MM}}$

I'll simply refer to it as the OLS estimator

The OLS estimator does have a compact matrix representation

Recall that $X_i := (X_{i1}, X_{i2}, \dots, X_{iK})'$ is the K -dimensional column vector that collects the K 'regressors' for observation i

Collecting all N of these vectors in an $N \times K$ matrix:

$$X := (X_1, X_2, \dots, X_N)' = \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1K} \\ X_{21} & X_{22} & \cdots & X_{2K} \\ X_{31} & X_{32} & \cdots & X_{3K} \\ \vdots & \vdots & \ddots & \vdots \\ X_{N1} & X_{N2} & \cdots & X_{NK} \end{pmatrix}$$

Notice again that this effectively overwrites earlier notation:

- previously: X denoted the K -dimensional column vector that collected the K random variables representing the regressors;
- from now on: X is an $N \times K$ -dimensional matrix collecting in each row the K regressors for each observation $i = 1, \dots, N$

Similarly, from now on $Y := (Y_1, Y_2, \dots, Y_N)'$ is the $N \times 1$ vector collecting all Y_i

The new matrix X and the new vector Y let us replace sums:

$$\sum X_i X_i' = X'X \text{ and } \sum X_i Y_i = X'Y$$

It follows that $\hat{\beta}^{\text{OLS}}$ has a nice and short matrix representation:

$$\hat{\beta}^{\text{OLS}} = (X'X)^{-1}X'Y$$

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

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 Estimator

Definition (Convergence in Probability)

A sequence of random variables Z_1, Z_2, \dots **converges in probability** to a random variable Z if for all $\epsilon > 0$,

$$\lim_{N \rightarrow \infty} P(|Z_N - Z| > \epsilon) = 0.$$

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

Often times the probability limit Z is a degenerate random variable that takes on a constant value everywhere

Definition (Bounded in Probability)

A sequence of random variables Z_1, Z_2, \dots is **bounded in probability** if for all $\epsilon > 0$, there exists $b_\epsilon \in \mathbb{R}$ and an integer N_ϵ such that

$$P(|Z_N| \geq b_\epsilon) < \epsilon \text{ for all } N \geq N_\epsilon.$$

Here's some new notation:

- for sequences that are bounded in probability we write $Z_N = O_p(1)$;
- for sequences that converge to zero in probability we write $Z_N = o_p(1)$.

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

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

Lemma

If $Z_N = c + o_p(1)$ then $Z_N = O_p(1)$ for a real constant c .

Proposition

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

$$W_N + X_N = o_p(1) \quad W_N + Y_N = O_p(1) \quad Y_N + Z_N = O_p(1)$$

$$W_N \cdot X_N = o_p(1) \quad W_N \cdot Y_N = o_p(1) \quad Y_N \cdot Z_N = O_p(1)$$

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 four 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, \dots be independent and identically distributed random variables with $EZ_i = \mu_Z$ and $\text{Var } Z_i = \sigma_Z^2 < \infty$. Then

$$\frac{1}{N} \sum_{i=1}^N Z_i = \mu_Z + o_p(1).$$

Of course, $\frac{1}{N} \sum_{i=1}^N Z_i$ 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

Lemma (Chebyshev's Inequality)

Let Z be a random variable with $EZ^2 < \infty$ and let $g(\cdot)$ be a nonnegative function. Then for any $c > 0$

$$P(g(Z) \geq c) \leq \frac{E(g(Z))}{c}.$$

Let $\bar{Z}_N := \frac{1}{N} \sum_{i=1}^N Z_i$

Here we're interested in bounding $\lim_{N \rightarrow \infty} P(|\bar{Z}_N - \mu_Z| > \epsilon)$

$$\begin{aligned} P(|\bar{Z}_N - \mu_Z| > \epsilon) &= P((\bar{Z}_N - \mu_Z)^2 > \epsilon^2) \\ &\leq \frac{E(\bar{Z}_N - \mu_Z)^2}{\epsilon^2} = \frac{\text{Var } \bar{Z}_N}{\epsilon^2} = \frac{\sigma_Z^2}{N \cdot \epsilon^2} \end{aligned}$$

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

We have used the fact $E(\bar{Z}_N) = \mu_Z$ and $\text{Var } \bar{Z}_N = \sigma_Z^2/N$
(we remember this from undergrad metrics)

This takes us back to the analogy principle

Remember earlier:

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

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 OLS estimator is consistent

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

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:

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

By WLLN

$$\frac{1}{N} \sum_{i=1}^N X_i X_i' = E(X_i X_i') + o_p(1)$$

and for its inverse:

$$\begin{aligned} \left(\frac{1}{N} \sum_{i=1}^N X_i X_i' \right)^{-1} &= \left(E(X_i X_i') + o_p(1) \right)^{-1} \\ &= E(X_i X_i')^{-1} + o_p(1) = O_p(1) \end{aligned}$$

using Slutsky's theorem, and a matrix version of the earlier Lemma that $c + o_p(1) = O_p(1)$, and assuming that $E(X_i X_i')$ is positive definite (inverse exists)

Copy and paste again:

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

For the other factor on the rhs:

$$\frac{1}{N} \sum_{i=1}^N X_i u_i = E(X_i u_i) + o_p(1) = 0 + o_p(1) = o_p(1)$$

It follows

$$\begin{aligned} \hat{\beta}^{\text{OLS}} &= \beta^* + O_p(1) \cdot o_p(1) \\ &= \beta^* + o_p(1) \end{aligned}$$

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

This means $\hat{\beta}^{\text{OLS}}$ 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 $\hat{\beta}^{\text{OLS}}$?

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