

This material comes from Hansen Chapter 2. I think that this material should be mostly review from ECON 8070, so I will cover it rather quickly.

## Linear Regression Preliminary Notes: Conditional expectation and projection

H: 2.5

We will spend much time this semester thinking about **conditional expectations**, that is,  $E[Y|X = x]$ . This is the mean of  $Y$  conditional on  $X$  taking the particular value  $x$ . The book sometimes uses the short-hand notation  $m(x) := E[Y|X = x]$ . You can think of  $E[Y|X = x]$  as a function — that is, when you plug in a new value of  $x$ , the value of the conditional expectation function can change. For example, if  $Y$  is a person's earnings and  $X$  is their years of education, the  $E[Y|X = 12]$  could differ (perhaps substantially) from, say,  $E[Y|X = 16]$ . Sometimes it will be useful to view  $E[Y|X]$  as a function of the random variable  $X$ ; in this case,  $E[Y|X]$  is itself random (because  $X$  is random). This is different from when you plug in a particular value of  $x$ ,  $E[Y|X = x]$  is no longer random; it is equal to some number (though, in most cases, it is unlikely that we know the value of this sort of population quantity).

### Notation:

I'll follow the convention in the book by writing

$$E[Y|X = x] = x_1\beta_1 + x_2\beta_2 + \cdots + x_{k-1}\beta_{k-1} + \beta_k$$

so that the “intercept” is in the last position. More specifically,

$$X = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_{k-1} \\ 1 \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}$$

so that  $X$  and  $\beta$  are both  $k \times 1$  vectors. And  $Y$ , the outcome, is a scalar.

When referring to particular observations, I'll use the notation

$$X_i = \begin{pmatrix} X_{i1} \\ X_{i2} \\ \vdots \\ X_{ik} \end{pmatrix} \quad \text{and} \quad Y_i$$

where  $X_i$  is  $k \times 1$  and  $Y_i$  is a scalar.

## Law of iterated expectations

H: 2.7

One of the most useful tools that you likely will recall from ECON 8070 is the law of iterated expectations. We'll provide a simple version and a general version (both require the regularity condition that  $E|Y| < \infty$ ). The simple version is

$$E[Y] = E[E[Y|X]]$$

In words, this says that the expectation of the conditional expectation is equal to the unconditional expectation. Here is an example. Continue to suppose that  $Y$  is a person's earnings and  $X$  is their years of education.  $E[Y|X = x]$  can vary arbitrarily for different values of  $x$ . But if you know  $E[Y|X = x]$  for all possible values of  $x$ , this will pin down the value of  $E[Y]$  (i.e., if you know mean earnings for all years of education, then you also should be able to recover the overall mean value of earnings). And, in particular, the law of iterated expectations says that the overall mean is equal to the mean of the conditional expectations (i.e., it puts more "weight" on conditional expectations of relatively common values of  $X$ ).

A more general version of the law of iterated expectations is the following: for any two random vectors  $X_1$  and  $X_2$ ,

$$E[Y|X_1] = E[E[Y|X_1, X_2]|X_1]$$

The inside expectation is an expectation of  $Y$  conditional on  $X_1$  and  $X_2$ , the outside expectation is over the distribution of  $X_2$  conditional on  $X_1$ .

## CEF Error

H: 2.8

We define the CEF error as the difference between  $Y$  and  $E[Y|X]$ . That is,

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

Rearranging terms also implies the following expression that we will use often

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

That is, the actual outcome  $Y$  is equal to the CEF plus the CEF error. Notice that (besides regularity conditions that expectations exist), we are not making any assumptions here — we can

essentially always write that  $Y$  is equal to its conditional expectation plus CEF error.

A fundamental property of the CEF error is that  $E[e|X] = 0$ . Let us show why this holds

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

From the law of iterated expectations, this also implies that

$$E[e] = E[\underbrace{E[e|X]}_{=0}] = 0$$

The condition that  $E[e|X] = 0$  is called **mean independence**. It is weaker than “full” independence. To give an example, even if  $E[e|X] = 0$ , it could be the case that  $E[e^2|X]$  varies with  $X$  which would imply that  $e$  and  $X$  are not independent.

We can also define the variance of the CEF error as

$$\sigma^2 = \text{var}(e) = E[(e - E[e])^2] = E[e^2]$$

$\sigma^2$  measures the amount of variation in  $Y$  that is not accounted for by  $E[Y|X]$ .

## Best Predictor

H: 2.11

Next, we will show that the conditional expectation function is the “best” predictor of  $Y$  given  $X$ . We can write any predictor as a function  $g(X)$ ; i.e., a function that takes values of  $X$  and makes predictions about what the outcome will be. In order to evaluate how well a predictor makes predictions, we need some criteria. The most common criteria is **mean squared prediction error**. This is given by

$$E[(Y - g(X))^2]$$

The squared difference between  $Y$  and its prediction  $g(X)$  is a measure of the distance between  $Y$  and  $g(X)$  — in particular, it is always non-negative and gets larger as when  $Y$  and  $g(X)$  are further away from each other. The outside expectation averages this distance over the distribution of  $Y$  and  $X$ .

Next, we show that  $m(X)$  minimizes mean squared prediction error. To this end, notice that

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

where the first equality holds by adding and subtracting  $m(X)$ , the second equality holds by the definition of  $e$ , and the last equality holds by squaring the main term and pushing the expectation through the sum. Now, let's consider each term. First,  $\mathbb{E}[e^2] = \sigma^2$  which does not depend on our prediction function  $g(X)$ . Next, consider the middle term

$$\begin{aligned}
\mathbb{E}[e(m(X) - g(X))] &= \mathbb{E}[(m(X) - g(X))\mathbb{E}[e|X]] \\
&= 0
\end{aligned}$$

where the first equality holds by the law of iterated expectations and the second equality holds because  $\mathbb{E}[e|X] = 0$ . Finally, the third term is minimized at 0 by setting  $g(X) = m(X)$ . This implies that setting  $g(X) = m(X)$  minimizes mean squared prediction error.

## Linear CEF

H: 2.15

An important special case is when  $\mathbb{E}[Y|X] = X'\beta$ ; that is, when the CEF is linear. Importantly, unlike the previous generic discussion of CEFs, in many cases, it may be a strong assumption to impose a linear CEF.

## Best Linear Predictor

H: 2.18

In many cases, we may be hesitant (or not have a good reason to believe) that the CEF is actually linear. Even in this case, we can still “run a regression” of  $Y$  on  $X$ .

In this section, we'll consider the best *linear* predictor of  $Y$  given  $X$ . A linear predictor for  $Y$  is a function  $X'b$  for some  $b \in \mathbb{R}^k$ . We will choose the best possible value of  $b$ . For this section, we will make the following regularity assumptions (Assumption 2.1 in the textbook)

1.  $\mathbb{E}[Y^2] < \infty$
2.  $\mathbb{E}[\|X\|^2] < \infty$
3.  $\mathbb{E}[XX']$  is positive definite

where  $\|x\| = (x'x)^{1/2}$  is the Euclidean length of the vector  $x$ . The first two assumptions imply that  $Y$  and  $X$  have finite second moments (and, therefore, finite means, variances, and covariances). As earlier, we will try to minimize mean squared prediction error; that is,

$$\beta = \underset{b}{\operatorname{argmin}} S(b)$$

where we define

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

We can use tools from calculus to solve this. As a first step, it is helpful to notice that

$$\mathbb{E}[(Y - X'b)^2] = \mathbb{E}[Y^2] - 2\beta' \mathbb{E}[XY] + \beta' \mathbb{E}[XX']\beta$$

As a side-comment, notice that the function that we are minimizing is a scalar (we will likely see other parameters/estimators this semester that minimize or maximize some objective function — these objective functions always return a scalar). Second, the expansion in the previous equation may not be obvious; in general, I think you can do quite well at linear algebra by keeping track of the dimensions of the terms that you are working with. In particular, you would know that you were making a mistake if the dimension of any term above were not scalar. Finally, we combined the two terms  $\mathbb{E}[YX']\beta$  and  $\beta' \mathbb{E}[XY]$  for the middle term above — these are equal to the transpose of each other and since they are both scalars, they are exactly equal to each other.

Next, let's take the derivative of the previous equation, set it equal to 0 and solve for  $\beta$ . Before we do this, let's be clear about exactly what we are doing. We are taking the derivative of a function  $S(b) : \mathbb{R}^k \rightarrow \mathbb{R}$  (that is a function that takes in a  $k$  dimensional vector and returns a scalar). And, in particular, this vector derivative is given by

$$\frac{\partial S(b)}{\partial b} = \begin{pmatrix} \frac{\partial S(b)}{\partial b_1} \\ \vdots \\ \frac{\partial S(b)}{\partial b_k} \end{pmatrix}$$

which is a  $k \times 1$  vector. As a side-comment, I follow the convention (which is also used in the textbook) of taking first derivatives of scalar-valued functions with respect to a vector “down” (so that the first derivative is  $k \times 1$  vector rather than a  $1 \times k$  vector) and (if needed) second derivatives “across” (so that the second derivative would result in a  $k \times k$  matrix). For more details about taking derivatives with respect to a vector, see the discussion on p.38 of the textbook and in Appendix A.20 (I consider this to be review material, but it is definitely material you should know).

Returning to our present problem, notice that

$$0 = \left. \frac{\partial S(b)}{\partial b} \right|_{b=\beta} = -2\mathbb{E}[XY] + 2\mathbb{E}[XX']\beta$$

In terms of mechanics, this derivative is just like the scalar case (the first term is a linear term and the second term is a quadratic term) except you just need to make sure that you get a  $k \times 1$  vector (rather than, especially, a  $1 \times k$  vector).

We can immediately solve the previous equation for  $\beta$ . The regularity conditions above imply that all of the moments here exist and that the matrix  $E[XX']$  is invertible.

$$\beta = E[XX']^{-1}E[XY]$$

In this context, we'll refer to  $\beta$  as the **linear projection coefficient** and  $X'\beta$  as the **best linear predictor**

We can also define the **projection error**

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

which is the difference between the actual value of  $Y$  and the best linear predictor of  $Y$  given  $X$ .

An important property of the projection error is that  $E[Xe] = 0$ . To see this, notice that

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

Notice that  $E[Xe]$  is a  $k \times 1$  vector. When  $X$  includes an intercept (so that  $X_k = 1$ ), this implies that  $E[e] = 0$ .

## Best linear approximation

H 2.25

Next, we show another interesting property/interpretation for the linear projection coefficient  $\beta$ . Suppose that we are interested in learning about the conditional expectation function  $m(x) = E[Y|X = x]$ , but we have no reason to suppose that it is linear. Therefore, we might be interested in trying to construct the best linear approximation to  $m(x)$ . That is, let's consider choosing  $\beta$  in the following way

$$\beta = \underset{b}{\operatorname{argmin}} E[(m(X) - X'b)^2]$$

Following roughly the same strategy as earlier, notice that

$$E[(m(X) - X'b)^2] = E[m(X)^2] - 2\beta'E[XY] + \beta'E[XX']\beta$$

Taking the derivative and setting equal to 0 implies that

$$0 = -2E[Xm(X)] + 2E[XX']\beta$$

which further implies that

$$\begin{aligned}\beta &= E[XX']^{-1}E[Xm(X)] \\ &= E[XX']^{-1}E[XE[Y|X]] \\ &= E[XX']^{-1}E[XY]\end{aligned}$$

where the second equality holds by the definition of  $m(X)$  and the last equality holds by the law of iterated expectations.

This is exactly the same expression for  $\beta$  as we derived earlier under the motivation of best linear predictor. This implies that  $X'\beta$  can additionally be interpreted as the best linear approximation to the underlying CEF — even if the CEF is nonlinear. This is a nice property for the linear projection model to have. That being said, even the best linear approximation to a nonlinear CEF can sometimes be quite poor. See Section 2.28 for an example and some discussion.

## Omitted variable bias

H: 2.24

To conclude this section, let's briefly talk about omitted variable bias. Let's partition  $X$  as follows  $X = (X'_1, X'_2)'$  and likewise partition  $\beta$  into  $\beta = (\beta'_1, \beta'_2)'$ . Suppose that we are interested in  $\beta_1$  from the linear projection of  $Y$  onto  $X_1$  and  $X_2$ :

$$Y = X'_1\beta_1 + X'_2\beta_2 + e \tag{1}$$

Since this is a linear projection, it implies that  $E[Xe] = 0$ .

However, let's suppose that  $X_2$  is not observed, so that it is infeasible to run a regression of  $Y$  on  $X_1$  and  $X_2$ . In this section, we consider properties of the following **short regression**

$$Y = X'_1\gamma_1 + u$$

which is the linear projection of  $Y$  on  $X_1$  only. Since this is a linear projection, we also have that  $E[X_1u] = 0$  and that

$$\begin{aligned}\gamma_1 &= E[X_1X'_1]^{-1}E[X_1Y] \\ &= E[X_1X'_1]^{-1}E[X_1(X_1\beta_1 + X_2\beta_2 + e)] \\ &= \beta_1 + E[X_1X'_1]^{-1}E[X_1X_2]\beta_2 \\ &= \beta_1 + \Gamma_{12}\beta_2\end{aligned}$$

where the first equality holds by the definition of linear projection of  $Y$  on  $X_1$ , the second equality holds by substituting for  $Y$ , the third equality combines and cancels terms and also holds because  $E[X_1e] = 0$  (since  $E[Xe] = 0$ ), and the last equality holds because we define  $\Gamma_{12} = E[X_1X_1']^{-1}E[X_1X_2]$ . Notice that  $\Gamma_{12}$  is the coefficient from the linear projection of  $X_2$  on  $X_1$ .

Importantly, the previous expression implies that  $\gamma_1$  is not generally equal to  $\beta_1$ ; that is, in general, we are not able to recover the parameter of interest  $\beta_1$  from the feasible regression of  $Y$  on  $X_1$ . This is probably not surprising — otherwise, our lives would be much easier! The difference between  $\gamma_1$  and  $\beta_1$  is called **omitted variable bias** and is a very important concern in many applications.

The only case where  $\gamma_1 = \beta_1$  is when  $\Gamma_{12}\beta_2 = 0$  which can happen when either  $\Gamma_{12} = 0$  or  $\beta_2 = 0$ .  $\Gamma_{12} = 0$  if  $E[X_1X_2] = 0$  which would be the case if  $X_1$  and  $X_2$  are uncorrelated.  $\beta_2 = 0$  occurs when the coefficient on  $X_2$  in Equation 1 is equal to 0. In words, the cases where you can recover  $\beta_1$  while only using the short regression are (i) if the omitted variables are uncorrelated with the included variables or (ii) if the omitted variables have no effect on the outcome.

**Side-Comment:** There are a number of cases where you might be able to figure out the sign of the omitted variable bias. The textbook gives the following simple example. Consider the case where  $Y$  is a person's earnings,  $X_1$  is a person's years of education, and  $X_2$  is a person's "ability", and where you are interested in  $\beta_1$  (the coefficient on years of education). However, suppose that ability is not observed. In this case, it might be reasonable to suppose that  $\beta_2 > 0$  (i.e., that, conditional on years of education, individuals with higher ability tend to have higher earnings) and that  $\Gamma_{12} > 0$  (i.e., that higher ability is positively correlated with more education). Under these conditions, it would be the case that  $\gamma_1 > \beta_1$ . This discussion suggests that a regression that only includes years of education would overestimate the effect of years of education relative to a model that included both education and ability. This sort of argument is quite common in applied work — something like: "even though we are not able to control for some important variable, it's correlation with the observed variable of interest and likely sign in the long regression indicate that the estimate of our coefficient of interest is likely a lower (or upper) bound."