

ABSTRACT. We consider the *population* and put the linear model in perspective of its probabilistic framework. The goal is to clarify the target of OLS estimation.

Let (Y, X) be $L^2(\Omega, \mathcal{F}, P)$ random variables. All we could possibly want to know about these quantities is contained in the joint distributions $F_{Y,X}$. This can be thought of the following collection of probabilities

$$\{P\{(Y, X) \in A\} ; A \text{ rectangle in } \mathbb{R}^{p+1}\}.$$

However our interest is usually asymmetric in that we are primarily interested in the random variable Y , and use random vector X only as a tool to reduce uncertainty about Y . In this case, we don't need to know probabilities of events that allow for uncertainty in both Y and X , but only those events that prescribe a fixed value x to X , whatever that value might happen to be. This information is contained in the collection of conditional distributions F_x ¹ that describe uncertainty about Y , when X is known to have realized as x . Equivalently,

$$\{P(Y \in A|X = x) ; A \text{ rectangle in } \mathbb{R}^1, x \in \mathbb{R}^p\}.$$

A distribution of Y for each value x is a lot of information. The regression function $f(x)$ that assigns to each x the mean of the conditional distribution F_x is a nice way to summarize and visualize the dependence of randomness in Y on X :²

$$f(x) := E[Y|X = x].$$

There is a very nice geometric interpretation of regression that relies on the notion of orthogonality in linear spaces. If we decompose

$$L^2(\Omega, \mathcal{F}, P) = L^2(\Omega, \sigma(X), P) \oplus L^2(\Omega, \sigma(X), P)^\perp,$$

then regression $E[Y|X]$ is the orthogonal projection of Y onto $L^2(\Omega, \sigma(X), P)$:

$$Y \equiv E[Y|X] + \varepsilon,$$

where the residual ε is orthogonal to *all* functions of X : for every $Z \in L^2(\Omega, \sigma(X), P)$ we have $E\varepsilon Z = 0$. We make one more use of orthogonality to express the regression function as an infinite vector of coordinates. Let $\{\psi_j\}_{j \geq 1}$ be an orthonormal basis for $L^2(\Omega, \sigma(X), P)$ chosen in the way that ψ_1, \dots, ψ_p span all linear combinations of random variables X_1, \dots, X_p , and decompose $f(X)$ into its Fourier series

$$f(X) = \sum_{j=1}^p \alpha_j \psi_j + \sum_{j>p} \alpha_j \psi_j. \tag{1}$$

If we think of ψ as appropriately normalized polynomials in X , then the first p terms is the linear part of f and the rest of the series is the nonlinear part of the regression function. Finally, when we estimate a regression, we project Y only onto the first p terms of above decomposition which

¹As a side remark, it is at all not straightforward to give meaning to the above set. The difficulty stems from the fact that the collection of events we must assign conditional probabilities to is uncountable, whereas measure theory allows for only countable operations. The construction is called regular conditional probabilities and is possible only under topological conditions that allow to reconcile this gap.

²Continuing with the side remark, conditional probabilities are defined in terms of the conditional expectation, but this construction does not guaranty that $A \mapsto P(Y \in A|X = x)$ is a probability measure for any given x .

can be called the linear regression function.

In the above discussion we made extensive use of orthogonality, yet regression coefficients are traditionally defined through a minimization problem. In fact, the two definitions are equivalent! This is a very important and useful characterization of orthogonal decomposition. First observe that eq. (1) is another orthogonal decomposition:

$$L^2(\Omega, \sigma(X), P) = \text{span}(X_1, \dots, X_p) \oplus \text{span}(\psi_j)_{j>p}$$

where we have replaced the first p terms of the orthonormal basis with the original linear terms that are not necessarily orthogonal. In this new basis, the linear part of the regression function has coordinates β . We will see how this effects the coefficients of linear regression shortly. Define the linear regression coefficient

$$b := \arg \min_{b \in \mathbb{R}^p} E[f(X) - b'X]^2 \quad (2)$$

and residual

$$\epsilon(b) := f(X) - b'X.$$

Claim: b solves (2) iff $\epsilon(b)$ is orthogonal to $\text{span}(X_1, \dots, X_p)$. This follows immediately from the following decomposition of the norm of $f(X)$:³

$$\|f(X) - b'X\|^2 = \|f(X) - \beta'X\|^2 + 2\langle f(X) - \beta'X, \beta'X - b'X \rangle + \|\beta'X - b'X\|^2 \quad (3)$$

First recall that β is defined as the orthogonal projection of $f(X)$ onto $\text{span}(X_1, \dots, X_n)$, so that the residual $f(X) - \beta'X$ is orthogonal to $\text{span}(X_1, \dots, X_p)$ by construction. Thus, the middle term in eq. (3) is zero for any choice of coefficient b because $\beta'X - b'X \in \text{span}(X_1, \dots, X_p)$, and we see that $\|\epsilon(b)\| \geq \|\epsilon(\beta)\|$. Conversely, if b solves (2), decomposition eq. (3) implies that $\beta'X = b'X$ (the middle term is still zero by definition of β !) and therefore $\epsilon(b) = \epsilon(\beta)$ is orthogonal to $\text{span}(X_1, \dots, X_p)$.

Putting all of the above together, we have produced the following orthogonal decomposition of random variable Y

$$Y = \beta'X + \sum_{j>p} \alpha_j \psi_j + \varepsilon. \quad (4)$$

The target of linear regression is the first term on the right. It is characterized by minimizing the size of the residual $\epsilon + \varepsilon$ and, equivalently, by producing a residual that is orthogonal to all of $\text{span}(X_1, \dots, X_p)$.

Above we used X_1, \dots, X_p as basis to define linear regression coefficient β . This choice of basis is natural because it is how we observe data, but does not immediately allow us to interpret coefficients β . If we partition $X = (D, W)$ and $\beta = (\beta_1, \beta_2)$, we can project orthogonally any variable of interest

³We use $\|Z\|^2$ to denote EZ^2 , and $\langle Y, Z \rangle$ to denote $E[YZ]$.

onto $\text{span}(W)$:

$$\begin{aligned} Y &= W'\gamma_{YW} + \tilde{Y}, & \tilde{Y} &\text{ orthogonal to } \text{span}(W) \\ D &= W'\gamma_{DW} + \tilde{D}, & \tilde{D} &\text{ orthogonal to } \text{span}(W) \\ \epsilon + \varepsilon &= 0 + \epsilon + \varepsilon \\ W &= W + 0, \end{aligned}$$

third line above follows by recalling that ϵ and ε are orthogonal to $\text{span}(X)$. Now substitute above decomposition with respect to subspace spanned by W into the regression equation (4):

$$W'\gamma_{YW} + \tilde{Y} = \beta_1'W'\gamma_{DW} + \beta_1'\tilde{D} + \beta_2'W + \epsilon + \varepsilon$$

rearrange

$$\tilde{Y} = \beta_1'\tilde{D} + \left[\beta_1'W'\gamma_{DW} + \beta_2'W + \epsilon + \varepsilon - W'\gamma_{YW} \right]$$

and note that the term in brackets is orthogonal to $\text{span}(\tilde{D})$! By our equivalence result it follows that β_1 is the projection coefficient of \tilde{Y} onto $\text{span}(\tilde{D})$. So we can think of β_1 as the linear term of regression of $[Y - \text{linear part in } W]$ on $[D - \text{linear part in } W]$.

MIT OpenCourseWare
<https://ocw.mit.edu>

14.382 Econometrics
Spring 2017

For information about citing these materials or our Terms of Use, visit: <https://ocw.mit.edu/terms>.