

# Best Predictors and Linear Regression

Chi-Chao Hung \*

last updated: March 17, 2025

## Contents

<b>1</b>	<b>Notations</b>	<b>2</b>
<b>2</b>	<b>Best Predictor</b>	<b>4</b>
2.1	CEF Error . . . . .	4
2.2	CEF as the Best Predictor . . . . .	5
<b>3</b>	<b>Best Linear Predictor</b>	<b>6</b>
3.1	Deriving the Best Linear Predictor . . . . .	6
3.2	Univariate Case . . . . .	7
3.3	Projection Error . . . . .	8
3.4	Linear CEF . . . . .	9
<b>4</b>	<b>Linear Regression</b>	<b>9</b>
4.1	What is Linear Regression? . . . . .	9
4.2	What is Linear Regression used for? . . . . .	10
<b>5</b>	<b>More Details</b>	<b>10</b>
5.1	Independence, Mean Independence, and Uncorrelatedness . . . . .	10
5.2	Existence and Uniqueness of the BLP . . . . .	11

---

\*This note borrows heavily from [Hansen \(2022\)](#). I was also supported by ChatGPT. All errors are mine. If you find any errors or have any suggestions, please contact me via email at [r13323021@ntu.edu.tw](mailto:r13323021@ntu.edu.tw).

## 1 Notations

Let us begin by defining some notations that will be used throughout this note.

**Definition 1.1** (Random Vector). A **random vector**  $X$  is a collection of random variables arranged as a vector. Formally, an  $k$ -dimensional random vector is written as:

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{bmatrix}$$

where each  $X_i$  is a random variable.

The **expectation** (mean) of  $X$  is a vector:

$$E[X] = \begin{bmatrix} E[X_1] \\ E[X_2] \\ \vdots \\ E[X_k] \end{bmatrix}.$$

**Definition 1.2** (Inner Product). Consider a **random vector**  $X$  and a corresponding **coefficient vector**  $\beta$ . The inner product  $X'\beta$  is given by:

$$\begin{aligned} X'\beta &= \begin{bmatrix} X_1 & X_2 & \dots & X_k \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \\ &= X_1\beta_1 + X_2\beta_2 + \dots + X_k\beta_k. \end{aligned}$$

We use the notation  $X'$  to represent the **transpose** of  $X$ , converting it from a column vector to a row vector. This transformation allows us to express the inner product as a *matrix multiplication*.

Squaring the inner product gives a **quadratic form**:

$$(X'\beta)^2 = (X_1\beta_1 + X_2\beta_2 + \dots + X_k\beta_k)^2.$$

Notice that using matrix notation,  $X'\beta$  and  $\beta'X$  are the same. Thus, squaring  $X'\beta$  gives:

$$(X'\beta)^2 = \underset{(1 \times 1)}{(\beta'X)} \underset{(1 \times 1)}{(X'\beta)} = \underset{(1 \times k)}{\beta'} \underset{(k \times k)}{(XX')} \underset{(k \times 1)}{\beta}.$$

Here,  $XX'$  is an outer product matrix:

$$XX' = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{bmatrix} \begin{bmatrix} X_1 & X_2 & \dots & X_k \end{bmatrix} = \begin{bmatrix} X_1X_1 & X_1X_2 & \dots & X_1X_k \\ X_2X_1 & X_2X_2 & \dots & X_2X_k \\ \vdots & \vdots & \ddots & \vdots \\ X_kX_1 & X_kX_2 & \dots & X_kX_k \end{bmatrix}.$$

Thus,  $XX'$  is an  $k \times k$  *symmetric matrix* containing all possible pairwise products of the components of  $X$ . The **expectation** of  $XX'$  is a matrix:

$$\mathbf{Q}_{XX} := E[XX'] = \begin{bmatrix} E[X_1X_1] & E[X_1X_2] & \dots & E[X_1X_k] \\ E[X_2X_1] & E[X_2X_2] & \dots & E[X_2X_k] \\ \vdots & \vdots & \ddots & \vdots \\ E[X_kX_1] & E[X_kX_2] & \dots & E[X_kX_k] \end{bmatrix}.$$

Let  $\tilde{X}$  be the de-meanned version of  $X$ :

$$\tilde{X} = X - E[X].$$

**Definition 1.3** (Covariance Matrix). The **covariance matrix** of  $X$  is defined as:

$$\begin{aligned} \text{Var}[X] &= E[\tilde{X}\tilde{X}'] = E[(X - E[X])(X - E[X])'] \\ &= \begin{bmatrix} \text{Var}[X_1] & \text{Cov}[X_1, X_2] & \dots & \text{Cov}[X_1, X_k] \\ \text{Cov}[X_2, X_1] & \text{Var}[X_2] & \dots & \text{Cov}[X_2, X_k] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[X_k, X_1] & \text{Cov}[X_k, X_2] & \dots & \text{Var}[X_k] \end{bmatrix}. \end{aligned}$$

Finally, let  $Y$  be a random variable.  $XY$  is a random vector.

$$\underset{(k \times 1)(1 \times 1)}{X} \underset{(1 \times 1)}{Y} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_k \end{bmatrix} Y = \begin{bmatrix} X_1Y \\ X_2Y \\ \vdots \\ X_kY \end{bmatrix}.$$

The **expectation** of  $XY$  is a vector:

$$\mathbf{Q}_{XY} := E[XY] = \begin{bmatrix} E[X_1Y] \\ E[X_2Y] \\ \vdots \\ E[X_kY] \end{bmatrix}.$$

## 2 Best Predictor

Let  $X$  be a random vector of student characteristics, such as **age**, **past GPA**, and **the number of hours spent studying**. Professor Chen might be interested in *predicting* the midterm exam score  $Y$  of a student based on these characteristics.

Our goal is to find a **predictor** for  $Y$ . A predictor is a function  $g(X)$  that takes the observed characteristics  $X$  and provides a prediction of  $Y$ . For example, we can predict that “*all students studying less than 10 hours would get a score of 50, and all students studying more than 10 hours would get a score of 80.*” This is a very simple predictor.

Some questions naturally arise:

- How do we measure the quality of a predictor?
- What is the best predictor of  $Y$ ?
- What are the properties of the best predictor?

It turns out that the conditional expectation function (CEF) is the best predictor of  $Y$ . So let us discuss some properties of the CEF functions.

**Remark 1.** Here we are studying the **population**. This means we assume that the joint distribution of  $X$  and  $Y$  is known. We are not fitting a line to any data points, nor are we estimating the relationship between  $X$  and  $Y$ . These topics will be covered in Chapters 3 and 4 of [Hansen \(2022\)](#).

### 2.1 CEF Error

Let  $Y$  be a random variable, and let  $X$  be a set of observed variables. The **CEF error** is the deviation of  $Y$  from its conditional expectation:

$$e := Y - E[Y | X].$$

This error term has several important properties:<sup>1</sup>

**Property 2.1.** The mean of CEF error is 0.

$$E[e] = E[Y - E[Y | X]] = E[Y] - E[E[Y | X]] = 0.$$

**Property 2.2 (Mean Independence).** The expectation of the CEF error given  $X$  is zero:

$$E[e | X] = E[Y | X] - E[Y | X] = 0.$$

<sup>1</sup>An implicit assumption is that  $E[|Y|] < \infty$ . See Theorem 2.4 in [Hansen \(2022\)](#) for more details.

This implies that, on average, the error does not systematically deviate from zero for any given  $X$ .

**Property 2.3.** For any function  $h(x)$  such that  $E[h(X)e]$  exists,

$$E[h(X)e] = 0.$$

A special case is when  $h(X) = X$ , which gives:  $E[Xe] = 0$ .

Notice that given property 2.1,  $\text{Cov}(h(X), e) = E[h(X)e]$ . So this property also implies that “the CEF error is **uncorrelated** with any function of  $X$ .”

*Proof.*

$$\begin{aligned} E[h(X)e] &= E[E[h(X)e \mid X]] = E[h(X) E[e \mid X]] \\ &= E[h(X)0] = 0 \end{aligned}$$

□

## 2.2 CEF as the Best Predictor

Now let us return to the question of evaluating the quality of a predictor.

**Definition 2.1** (Mean Squared Error). The mean squared error (MSE) of a predictor  $g(X)$  is given by:

$$\text{MSE}(g) = E[(Y - g(X))^2].$$

A **best predictor** is a function of  $X$  that *minimizes* the MSE. In fact, the CEF,  $E[Y \mid X]$ , is the optimal predictor of  $Y$  in this sense.

**Theorem 2.1** (Conditional Mean as Best Predictor). If the second moment of  $Y$  exist<sup>2</sup>, then for any predictor  $g(X)$  the following holds:

$$E[(Y - m(X))^2] \leq E[(Y - g(X))^2],$$

where  $m(X) := E[Y \mid X]$  is the CEF.

That is, the CEF  $E[Y \mid X]$  minimizes the MSE among all predictors that are functions of  $X$ .

---

<sup>2</sup> $E[Y^2] < \infty$ . See Theorem 2.7 in [Hansen \(2022\)](#).

*Proof.*

$$\begin{aligned}
& \mathbb{E} [(Y - g(X))^2] && \text{(MSE of } g(X)) \\
&= \mathbb{E} [(Y - m(X) + m(X) - g(X))^2] && \text{(Add and Subtract)} \\
&= \mathbb{E} [(e + m(X) - g(X))^2] && \text{(Definition of CEF error)} \\
&= \mathbb{E} [e^2 + 2e(m(X) - g(X)) + (m(X) - g(X))^2] \\
&= \mathbb{E} [e^2] + 2\mathbb{E} [e(m(X) - g(X))] + \mathbb{E} [(m(X) - g(X))^2] \\
&= \mathbb{E} [e^2] + \mathbb{E} [(m(X) - g(X))^2] && \text{(Property 2.3)} \\
&\geq \mathbb{E} [e^2] = \mathbb{E} [(Y - m(X))^2] && \text{(MSE of } m(X))
\end{aligned}$$

□

**Remark 2.** Notice that if  $X$  is a constant, then the CEF is simply  $\mathbb{E}[Y]$ . The best predictor of  $Y$  is then the **unconditional mean**.  $Y = \mathbb{E}[Y] + e$  is called the intercept-only model.

### 3 Best Linear Predictor

#### 3.1 Deriving the Best Linear Predictor

Although we have proven that the CEF is the optimal predictor of  $Y$ , its functional form is generally unknown. Therefore, it is often useful to restrict our focus to a class of predictors with a simpler structure. Here, we consider one of the simplest cases: linear predictors.

The **linear predictor** of  $Y$  given  $X$  takes the form

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

where

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_{k-1} \\ 1 \end{bmatrix}$$

and  $\beta$  is an  $k \times 1$  coefficient vector.

**Definition 3.1** (Best Linear Predictor). A function  $\mathcal{P}(Y | X) = X'\beta$  is the best linear predictor (BLP) of  $Y$  given  $X$  if<sup>3</sup>

$$\beta = \arg \min_{\beta \in \mathbb{R}^k} \mathbb{E} [(Y - X'\beta)^2].$$

<sup>3</sup>See Definition 2.5 in Hansen (2022).

Now let us derive the BLP. First define

$$S(\beta) := E[e^2] = E[(Y - X'\beta)^2].$$

Expanding the square:

$$\begin{aligned} S(\beta) &= E[Y^2 - 2YX'\beta + \beta'XX'\beta] \\ &= E[Y^2] - 2\beta'E[XY] + \beta'E[XX']\beta. \end{aligned}$$

Taking the derivative with respect to  $\beta$ ,<sup>4</sup> the first-order condition is

$$\frac{\partial}{\partial \beta} S(\beta) = \underset{(k \times 1)}{-2E[XY]} + \underset{(k \times k)}{2E[XX']} \underset{(k \times 1)}{\beta} = \underset{(k \times 1)}{0}.$$

Solving for  $\beta$ :

$$\underset{(k \times k)}{E[XX']} \underset{(k \times 1)}{\beta} = \underset{(k \times 1)}{E[XY]}$$

If  $E[XX']$  is invertible, the solution for the best linear predictor is:

$$\underset{(k \times k)}{\beta} = \underset{(k \times k)}{E[XX']}^{-1} \underset{(k \times 1)}{E[XY]}.$$

Using the notation  $\mathbf{Q}_{XX} = E[XX']$  and  $\mathbf{Q}_{XY} = E[XY]$ , we rewrite:

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

The solution  $X'\beta = X'(E[XX'])^{-1}E[XY]$  is often called the **linear projection** of  $Y$  on  $X$ , and  $\beta$  is called the **projection coefficient**.<sup>5</sup> Some authors refer to the BLP as the *population* regression. (Angrist and Pischke, 2009)

### 3.2 Univariate Case

To illustrate the BLP, let us consider the univariate case where  $X_1$  is the only random variable included.

$$X = \begin{bmatrix} 1 \\ X_1 \end{bmatrix}$$

The BLP is

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

<sup>4</sup>See Appendix for technical details.

<sup>5</sup>Loosely speaking, “linear projection” means to find the random variable linear in  $X$  that has the smallest distance to  $Y$ . Here MSE can be seen as a measure of distances between two random variables.

where

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} 1 & E[X_1] \\ E[X_1] & E[X_1^2] \end{bmatrix}^{-1} \begin{bmatrix} E[Y] \\ E[X_1 Y] \end{bmatrix}$$

Since

$$\begin{bmatrix} 1 & E[X_1] \\ E[X_1] & E[X_1^2] \end{bmatrix}^{-1} = \frac{1}{E[X_1^2] - E[X_1]^2} \begin{bmatrix} E[X_1^2] & -E[X_1] \\ -E[X_1] & 1 \end{bmatrix}$$

we can derive the projection coefficients  $\beta_1$  and  $\beta_0$ .

$$\beta_1 = \frac{(E[X_1 Y] - E[X_1] E[Y])}{E[X_1^2] - E[X_1]^2} = \frac{\text{Cov}(X_1, Y)}{\text{Var}(X_1)}$$

and

$$\begin{aligned} \beta_0 &= \frac{E[Y] E[X_1^2] - E[X_1 Y] E[X_1]}{\text{Var}(X_1)} \\ &= E[Y] - \frac{\text{Cov}(X_1, Y)}{\text{Var}(X_1)} E[X_1] \\ &= E[Y] - \beta_1 E[X_1] \end{aligned}$$

The BLP of  $Y$  is then

$$\mathcal{P}(Y | X_1) = \beta_0 + \beta_1 X_1$$

### 3.3 Projection Error

The **projection error** is the deviation of  $Y$  from its linear projection:

$$e := Y - \mathcal{P}(Y | X).$$

This error term has several important properties.

**Property 3.1.**  $E[Xe] = 0$ .

*Proof.*

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

□

This is sometimes called the **orthogonality condition**.

**Property 3.2.**  $E[e] = 0$  if  $X$  includes a constant. Notice that property 3.1 is a set of  $k$  equations:

$$E[X_j e] = 0 \text{ for all } j = 1, \dots, k.$$

Whenever the constant 1 is included in  $X$ , it follows immediately from property 3.1 that  $E[e] = 0$ . This also implies that  $\text{Cov}(X_j, e) = 0$  since  $\text{Cov}(X_j, e) = E[X_j e] = E[X_j] E[e]$ .



### 3.4 Linear CEF

In general, BLP is not CEF. However, if we consider the special case that we know (or assume) that CEF is linear in  $X$ , then the BLP and the CEF coincide.

**Theorem 3.1** (Linear CEF). If the CEF is linear in  $X$ , then the CEF is the BLP.

The proof is trivial since we already know that the CEF is the best predictor among *all* predictors. Here we present another way to understand this result.

*Proof.* Let  $e := Y - E[Y | X]$ . Under the linear CEF assumption, we have

$$Y = E[Y | X] + e = X'\beta + e$$

Using property 2.1 of the CEF error, we have

$$E[XY] = E[X(X'\beta + e)] = E[XX']\beta + E[Xe] = E[XX']\beta.$$

As a result,  $\beta = E[XX']^{-1} E[XY]$  as in the BLP. □

**Definition 3.2** (Linear CEF Model).

$$E[Y | X] = X'\beta \tag{1}$$

$$Y = E[Y | X] + e \tag{2}$$

$$E[e | X] = 0 \tag{3}$$

**Remark 3.** Notice that a linear CEF is a *model assumption*. If the CEF is not linear, we say the linear CEF model is **misspecified**. However, the BLP is still the best linear approximation of the CEF.<sup>a</sup> This is a motivation for using the BLP.

<sup>a</sup>See Section 2.25 in Hansen (2022).

## 4 Linear Regression

### 4.1 What is Linear Regression?

Regression is a method that allows researchers to summarize how predictions or average values of an outcome vary across individuals defined by a set of predictors. (Gelman et al., 2021) Generally, a linear regression model takes the form

$$Y = \underbrace{X'\beta}_{\text{predictor}} + \underbrace{e}_{\text{error}}$$

In [Hansen \(2022\)](#), the term “linear regression” refers to the linear CEF model. In some other cases,  $X'\beta$  is viewed as the best linear predictor (may not be CEF) and  $e$  is the projection error. An important thing to remember is that in both cases the error term  $e$  has no life of its own. It is always defined as the difference between the outcome and the predictor.

Table 1: Comparison of the CEF error and Projection Error

	CEF Error	Projection Error
Definition	$e := Y - E[Y   X]$	$e := Y - \mathcal{P}(Y   X)$
Properties	$E[e   X] = 0$ $E[h(X)e] = 0$ $E[e] = 0$	$E[Xe] = 0$ $E[e] = 0$ when 1 is included in X

**Remark 4.** In older textbooks, the assumption  $E[e | X] = 0$  is called the “exogeneity” assumption. This is a very *misleading* terminology since “exogeneity” is a causal concept while the CEF itself has no causal interpretation.

## 4.2 What is Linear Regression used for?

1. **Prediction:** Either linear CEF or linear projection serves as the best linear predictor.
2. **Comparison:** Linear regression can be used to compare the average values of  $Y$  across different groups defined by  $X$ . This is especially useful when  $X$  is binary or categorical.
3. **Exploring Assotiations:** Linear regression can be used to explore the association between  $Y$  and  $X$ .
4. **Causal Inference:** So far we said nothing about causality. In fact, linear regression is not a causal model. The linear function  $Y = X'\beta$  does not tell us wether it is  $X$  causing  $Y$  or  $Y$  causing  $X$ . We will discuss the **Potential Outcome Framework** at the end this semester, which provides a framework for causal inference. Once a causal model is established, linear regression may be used as a tool.

## 5 More Details

### 5.1 Independence, Mean Independence, and Uncorrelatedness

1.  $X$  and  $Y$  are **statistically independent** if and only if their joint distribution equals the product of their marginals.

2.  $Y$  is **mean independent** of  $X$  if and only if  $E[Y | X] = E[Y]$ .
3.  $X$  and  $Y$  are **uncorrelated** if and only if  $\text{Cov}(X, Y) = 0$ . If  $E[X]$  or  $E[Y]$  is zero, then  $E[XY] = 0$ .

**Theorem 5.1.** If  $X$  and  $Y$  are statistically independent, then  $X$  is mean independent of  $Y$  and  $Y$  is mean independent of  $X$ . Moreover,  $X$  and  $Y$  are uncorrelated. The converse is not true.

**Theorem 5.2.** If  $X$  is mean independent of  $Y$ , then  $X$  and  $Y$  are uncorrelated. The converse is not true.

The proofs is left for the readers. See [DiTraglia \(2023\)](#) for more details.

## 5.2 Existence and Uniqueness of the BLP

Though we made no assumptions on the CEF when deriving the BLP, the existence of the linear projection requires some conditions.

1.  $E[Y^2] < \infty$
2.  $E\|X\|^2 < \infty$
3.  $\mathbf{Q}_{XX} = E[XX']$  is invertible.

The first two conditions says that both  $X$  (vector) and  $Y$  (scalar) have finite second moments. The third condition is the most important. It says that the predictors in  $X$  are not perfectly collinear. If  $\mathbf{Q}_{XX}$  is not invertible, then the unique linear projection does not exist.

## References

**Angrist, Joshua David and Jörn-Steffen Pischke**, *Mostly Harmless Econometrics: An Empiricist's Companion*, Princeton: Princeton University Press, 2009.

**DiTraglia, Frank**, “Why Econometrics Is Confusing Part II: The Independence Zoo,” <https://www.econometrics.blog/post/why-econometrics-is-confusing-part-ii-the-independence-zoo/> January 2023.

**Gelman, Andrew, Jennifer Hill, and Aki Vehtari**, *Regression and Other Stories* Analytical Methods for Social Research, Cambridge: Cambridge University Press, 2021.

**Hansen, Bruce E.**, *Econometrics*, Princeton: Princeton University Press, 2022.

## Acronyms

BLP    best linear predictor. [6](#), [7](#)  
CEF    conditional expectation function. [4–6](#)  
MSE    mean squared error. [5](#)