

We first set up some notations. Let \mathbf{y} be an $(n \times 1)$ vector, collecting observations of the response (or dependent) variable. Let \mathbf{X} be an $(n \times k)$ matrix, collecting observations of the explanatory (or independent) variables.

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_n^\top \end{bmatrix}$$

Here, the rows of \mathbf{X} store observations (individuals, countries, etc.) and the columns store variables. The row vector $\mathbf{x}_i^\top = (x_{i1}, x_{i2}, \dots, x_{ik})$ contains the values of each of the k regressors for observation i . Unless otherwise stated, for this section, we assume that $x_{i1} = 1$ is a constant term.

1 Linear Projection

Given observed data $(y_i, \mathbf{x}_i)_{i=1}^n$, we are interested in describing the relationship between y_i and \mathbf{x}_i . An often used device is the conditional expectation function, $\mu(\mathbf{x}_i) \equiv \mathbb{E}[y_i | \mathbf{x}_i]$. Its popularity stems from an oft-touted property of the conditional expectation function is that it has an interpretation as an optimal predictor.

Proposition 1.1. Assume Y_i is square integrable, then $\mu(\mathbf{x}_i) = \arg \min_{m \in \mathcal{M}} \mathbb{E}[(y_i - m(\mathbf{x}_i))]$, where \mathcal{M} is set of all square integrable functions.

Proof. We decompose the square loss function as follows

$$\mathbb{E}[(y_i - m(\mathbf{x}_i))^2] = \mathbb{E}[(y_i - \mu(\mathbf{x}_i)) - (\mu(\mathbf{x}_i) - m(\mathbf{x}_i))]^2 = \mathbb{E}[(y_i - \mu(\mathbf{x}_i))^2] + \mathbb{E}[(\mu(\mathbf{x}_i) - m(\mathbf{x}_i))^2].$$

For the optimization problem, the first term does not matter and the second is 0 when we set $m(\mathbf{x}_i) = \mu(\mathbf{x}_i)$. □

For this reason the conditional expectation is sometimes called the “best” or minimum mean square error predictor. It is worth pointing out that this optimality property is an artifact of the mean square loss function, to which we need not attribute any special importance. Were we to adopt an absolute loss function, the optimal predictor is the conditional median function.

The estimation of $\mu(\mathbf{x}_i)$ can be difficult, since without any restriction on the joint distribution, this object is an infinite dimensional object. Instead we can consider a linear approximation/predictor. Let $\text{Proj}(y_i | \mathbf{x}_i) \equiv \mathbf{x}_i^\top \beta$ with $\beta \equiv \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top]^{-1} \mathbb{E}[\mathbf{x}_i y_i]$ to be the population linear projection, where we assumed the invertibility of $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top]$.

Proposition 1.2. Assume Y_i is square integrable, then $\text{Proj}(y_i | \mathbf{x}_i) = \arg \min_{b \in \mathbb{R}^k} \mathbb{E}[(y_i - \mathbf{x}_i^\top b)^2]$.

Proof. Same as above. □

$\text{Proj}(y_i | \mathbf{x}_i)$ is thus termed as the best linear predictor.

The formulation of β is obtained from the first order condition of the criterion function

$$\mathbb{E}[\mathbf{x}_i (y_i - \mathbf{x}_i^\top \beta)] = 0.$$

To facilitate interpretation, we assume that \mathbf{x}_i include a constant (unless otherwise stated), so then the necessary condition for optimality is a covariance restriction which states that the prediction error $u_i \equiv y_i - \mathbf{x}_i^\top \beta$ such that

$\mathbb{E}[u_i] = 0$ must exhibit no (linear) covariation with the predictor variables \mathbf{x}_i . Assuming that $\mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top]$ is invertible, we obtain

$$\beta = \mathbb{E}[\mathbf{x}_i \mathbf{x}_i^\top]^{-1} \mathbb{E}[\mathbf{x}_i y_i],$$

whose sample analogue is the least square estimator. Before moving onto least square properties, we wish to state two properties of the linear projection operator as it inherits properties from the conditional expectation operator.

Exercise 1.A

Let $(y_i, \mathbf{x}_i, \mathbf{z}_i)_{i=1}^n$ be observed variables, and assume invertibility of $\mathbb{E}[(\mathbf{x}_i, \mathbf{z}_i)^\top (\mathbf{x}_i, \mathbf{z}_i)]$. Prove that

1. if $\mu(\mathbf{x}_i)$ is linear, then $\text{Proj}(y_i | \mathbf{x}_i) = \mu(\mathbf{x}_i)$;
2. the tower property holds, $\text{Proj}(y_i | \mathbf{x}_i) = \text{Proj}(\text{Proj}(y_i | \mathbf{x}_i, \mathbf{z}_i) | \mathbf{x}_i)$.

From the previous discussion, we can posit that

$$y_i = \mathbf{x}_i^\top \beta + u_i,$$

which is simply a decomposition of the dependent variable into a linear predictor component $\mathbf{x}_i^\top \beta$ and a prediction error term u_i for $u_i \equiv y_i - \text{Proj}(y_i | \mathbf{x}_i)$. And u_i satisfies the following two properties $\mathbb{E}[u_i] = 0$ and $\mathbb{E}[u_i | \mathbf{x}_i] = 0$.

Example 1.3 (Freedman, 2008; Lin, 2013). The estimation of the average treatment effect can also be cast in this regression framework. Let $(y_i(1), y_i(0))_{i=1}^n$ be fixed potential outcomes. Let d_i be a treatment indicator such that $d_i \perp (y_i(1), y_i(0))$ for all i and that

$$y_i = d_i y_i(1) + (1 - d_i) y_i(0).$$

Rearranging the assignment equation, we have

$$y_i = d_i (y_i(1) - y_i(0)) + y_i(0) = d_i \tau_i + \varepsilon_i,$$

where $\tau_i = y_i(1) - y_i(0)$ and $\varepsilon_i = y_i(0)$. This equation still deviates from the regression framework in two fronts. First the coefficient on the explanatory variable d_i is usually fixed whilst $\tau_i \neq \tau_j$ in general. Secondly note that $n^{-1} \sum_{i=1}^n \varepsilon_i = n^{-1} \sum_{i=1}^n y_i(0) \neq 0$. To fit the regression framework, we rewrite the equation as

$$\begin{aligned} y_i &= \beta_0 + \beta_1(d_i - \bar{d}) + u_i, \\ \beta_0 &= \bar{d} \bar{y}(1) + (1 - \bar{d}) \bar{y}(0), \\ \beta_1 &= \bar{y}(1) - \bar{y}(0), \\ \beta_{1i} &= (y_i(1) - \bar{y}(1)) - (y_i(0) - \bar{y}(0)), \\ \beta_{0i} &= \bar{d} (y_i(1) - \bar{y}(1)) + (1 - \bar{d}) (y_i(0) - \bar{y}(0)), \\ u_i &= \beta_{0i} + \beta_{1i} (d_i - \bar{d}), \end{aligned}$$

for $\bar{d} = n^{-1} \sum_{i=1}^n d_i$, $\bar{y}(1) = n^{-1} \sum_{i=1}^n y_i(1)$ and $\bar{y}(0) = n^{-1} \sum_{i=1}^n y_i(0)$. The centering of d_i renders the coefficient β_0 more interpretable. Without the centering, β_0 becomes $(1 - \bar{d}) \bar{y}(0)$. Observe now that $\sum_{i=1}^n \beta_{0i} = \sum_{i=1}^n \beta_{1i} = 0$ by construction, and $\mathbb{E}[u_i | d_1, \dots, d_n] = \beta_{0i}$ and $\mathbb{E}[\sum_{i=1}^n u_i | d_1, \dots, d_n] = 0$, so that this formulation satisfies certain “weak” form of regression conditions. We can now write the least square estimator of y_i on $1, (d_i - \bar{d})$ as

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (d_i - \bar{d}) (y_i - \bar{y})}{\sum_{i=1}^n (d_i - \bar{d})^2} = \beta_1 + \frac{\sum_{i=1}^n (d_i - \bar{d}) u_i}{\sum_{i=1}^n (d_i - \bar{d})^2},$$

where the second term is interpreted as errors from not observing all $y_i(1), y_i(0)$.

The formulation in this section implies that the least square procedure reduces the potentially very complicated structure of the data into a lower dimensional parametric model capable of being directly interpreted. Even if this low

dimensional model comes at a cost of misspecification, it is still guaranteed to provide certain types of approximations which we can interpret mechanically without knowing the exact misspecification.

2 Regression Algebra

Given a sample of $(y_i, \mathbf{x}_i)_{i=1}^n$, we could estimate the predictor coefficient β by

$$\hat{\beta} = \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i y_i \right) = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}),$$

which is solution to the following optimization problem

$$\hat{\beta} = \arg \min_{b \in \mathbb{R}^k} (\mathbf{y} - \mathbf{X}b)^T (\mathbf{y} - \mathbf{X}b) = \arg \min_{b \in \mathbb{R}^k} \sum_{i=1}^n (y_i - \mathbf{x}_i^T b)^2.$$

We now define two matrices for understanding matrix algebra, the projection matrix $\mathbf{P}_X \equiv \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$, and the annihilator matrix $\mathbf{M}_X = \mathbf{I}_n - \mathbf{P}_X$.¹ These two matrices have special properties.

Exercise 2.A

Prove that \mathbf{P}_X and \mathbf{M}_X are idempotent and symmetric, and $\mathbf{P}_X \mathbf{M}_X = 0$.

It is immediate to see that $\mathbf{P}_X \mathbf{X} = \mathbf{X}$ and $\mathbf{M}_X \mathbf{X} = 0$. More importantly, we have

$$\mathbf{P}_X \mathbf{y} = \mathbf{X} \hat{\beta} = \hat{\mathbf{y}} \text{ and } \mathbf{M}_X \mathbf{y} = \mathbf{y} - \hat{\mathbf{y}} = \hat{\mathbf{u}},$$

which gives respectively the fitted values and the residuals from the regression.

2.1 Frisch-Waugh-Lovell Theorem

Let us partition the regressor matrix $\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 \end{bmatrix}$ into two parts \mathbf{X}_1 , $n \times k_1$ matrix, and \mathbf{X}_2 , $n \times k_2$ matrix, such that $k = k_1 + k_2$. The regression model can now be written as

$$\mathbf{y} = \mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 + \mathbf{u}. \quad (2.1)$$

Define $\mathbf{P}_1 \equiv \mathbf{X}_1 (\mathbf{X}_1^T \mathbf{X}_1)^{-1} \mathbf{X}_1^T$ and $\mathbf{P}_2 \equiv \mathbf{X}_2 (\mathbf{X}_2^T \mathbf{X}_2)^{-1} \mathbf{X}_2^T$, and the associated annihilator matrix \mathbf{M}_1 and \mathbf{M}_2 . If we pre-multiply the regression equation by \mathbf{M}_1 , we have

$$\mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \beta_2 + \mathbf{M}_1 \mathbf{u}. \quad (2.2)$$

We can now state the Frisch-Waugh-Lovell theorem.

Theorem 2.1. *The regression coefficients $\hat{\beta}_2$ from Equation (2.1) are identical to the regression coefficients $\tilde{\beta}_2$ from Equation (2.2). The residuals obtained from the respective regressions are identical.*

Proof. We first write out the equation for $\tilde{\beta}_2$ as

$$\tilde{\beta}_2 = \left(\mathbf{X}_2^T \mathbf{M}_1 \mathbf{X}_2 \right)^{-1} \mathbf{X}_2^T \mathbf{M}_1 \mathbf{y}.$$

By regression algebra we can express \mathbf{y} as

$$\mathbf{y} = \mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \hat{\mathbf{u}} = \mathbf{X}_1 \hat{\beta}_1 + \mathbf{X}_2 \hat{\beta}_2 + \mathbf{M}_X \mathbf{y}.$$

¹ The motivation of the name of \mathbf{P}_X can be seen as the result of orthogonal projection. Alternatively we can use singular value decomposition to under \mathbf{P}_X . By SVD, we write $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T$, where \mathbf{U} is $n \times k$ matrix, Σ is $k \times k$ diagonal matrix of singular values of \mathbf{X} and \mathbf{V} is $k \times k$ matrix, and that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_k$. Using this formulation, we see that $\mathbf{P}_X = \mathbf{U} \mathbf{U}^T$. The projection matrix takes in any vector $\mathbf{y} \in \mathbb{R}^n$ and apply $\mathbf{U}^T \mathbf{y}$ to find coefficients that could express \mathbf{y} using the basis of the column space of \mathbf{X} and then return it as an element in the column space of \mathbf{X} .

Premultiply both sides by $\mathbf{X}_2^T \mathbf{M}_1$ to get

$$\mathbf{X}_2^T \mathbf{M}_1 \mathbf{y} = \mathbf{X}_2^T \mathbf{M}_1 \mathbf{X}_2 \hat{\beta}_2 + \mathbf{X}_2^T \mathbf{M}_1 \mathbf{X}_2 \mathbf{y} = \mathbf{X}_2^T \mathbf{M}_1 \mathbf{X}_2 \hat{\beta}_2,$$

where the second line follows from that fact that \mathbf{M}_2 annihilates \mathbf{X}_2 . The first claim that $\hat{\beta}_2 = \tilde{\beta}_2$ now follows.

For the second claim, observe that

$$\mathbf{M}_1 \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \tilde{\beta}_2 + \mathbf{M}_1 \mathbf{M}_X \mathbf{y} = \mathbf{M}_1 \mathbf{X}_2 \tilde{\beta}_2 + \mathbf{M}_X \mathbf{y},$$

so the residuals are $\mathbf{M}_X \mathbf{y}$ which is identical to those residuals from the regression of Equation (2.1). \square

The Frisch-Waugh-Lovell theorem is an elegant description of the anatomy of multiple regression. Loosely speaking, the regression coefficient $\hat{\beta}_2$ is interpreted as “the effect” of \mathbf{X}_2 on \mathbf{y} , holding \mathbf{X}_1 constant (or controlling for \mathbf{X}_1). This theorem makes clear of this statement through a two step procedure:

- **Step 1:** Regress the columns of \mathbf{X}_2 onto the columns of \mathbf{X}_1 :

$$\mathbf{X}_2 = \mathbf{X}_1 \gamma + \eta$$

From the usual OLS formulas, we know that $\hat{\gamma} = (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{X}_2$ and that the residuals from this regression are given by:

$$\begin{aligned} \hat{\mathbf{u}}_2 &= \mathbf{X}_2 - \mathbf{X}_1 \hat{\gamma} \\ &= (\mathbf{I} - \mathbf{P}_1) \mathbf{X}_2. \end{aligned}$$

- **Step 2:** Regress \mathbf{y} onto $\hat{\mathbf{u}}_2$:

$$\mathbf{y} = \hat{\mathbf{u}}_2 \delta + \nu$$

From this second regression, we get the following formula for $\hat{\delta}$:

$$\begin{aligned} \hat{\delta} &= (\hat{\mathbf{u}}_2' \hat{\mathbf{u}}_2)^{-1} \hat{\mathbf{u}}_2' \mathbf{y} \\ \hat{\delta} &= (\mathbf{X}_2' (\mathbf{I} - \mathbf{P}_1) \mathbf{X}_2)^{-1} \mathbf{X}_2' (\mathbf{I} - \mathbf{P}_1) \mathbf{y} \end{aligned}$$

The Frisch-Waugh-Lovell theorem states that $\hat{\beta}_2$ is equal to $\hat{\delta}$.

So, what this says is that if we want to estimate the effect of \mathbf{X}_2 on \mathbf{y} , controlling for \mathbf{X}_1 , we first partial out the effect of \mathbf{X}_1 on \mathbf{X}_2 , and we regress \mathbf{y} on only the component of \mathbf{X}_2 that is uncorrelated with \mathbf{X}_1 .

References

- Freedman, David A.** 2008. “On regression adjustments in experiments with several treatments.” *Ann. Appl. Stat.*, 2(1): 176–196.
- Lin, Winston.** 2013. “Agnostic notes on regression adjustments to experimental data: Reexamining Freedman’s critique.” *Ann. Appl. Stat.*, 7(1): 295–318.