

Notation: in this note, y is a scale random variable, and x is a $K \times 1$ random vector.

1 Conditional Expectation Model

A regression model can be written as

$$y = m(x) + \epsilon,$$

where $m(x) = E[y|x]$ is called the *conditional mean function*, and $\epsilon = y - m(x)$ is called the *regression error*. Such an equation holds for (y, x) that follows any joint distribution, as long as $E[y|x]$ exists. The error term ϵ satisfies these properties:

- $E[\epsilon|x] = 0$,
- $E[\epsilon] = 0$,
- $E[h(x)\epsilon] = 0$, where h is a function of x .

The last property implies that ϵ is uncorrelated with any function of x .

If we are interested in predicting y given x , then the conditional mean function $E[y|x]$ is “optimal” in terms of the *mean squared error* (MSE).

As y is not a deterministic function of x , we cannot predict it with certainty. In order to evaluate different methods of prediction, we must therefore propose a criterion for comparison. For an arbitrary prediction method $g(x)$, we employ a *loss function* $L(y, g(x))$ to measure how wrong is the prediction, and the expected value of the loss function is called the

risk $R(y, g(x))$. The *quadratic loss function* is defined as

$$L(y, g(x)) = (y - g(x))^2,$$

and the corresponding risk

$$R(y, g(x)) = E \left[(y - g(x))^2 \right]$$

is called the MSE.

Due to its operational convenience, MSE is one of the most widely used criterion. Under MSE, the conditional expectation function happens to be the best prediction method for y given x . In other words, the conditional mean function $m(x)$ minimizes the MSE.

We can take a guess-and-verify this claim of optimality. For an arbitrary $g(x)$, the risk can be decomposed into three terms

$$\begin{aligned} & E \left[(y - g(x))^2 \right] \\ = & E \left[(y - m(x))^2 \right] + 2E \left[(y - m(x)) (m(x) - g(x)) \right] + E \left[(m(x) - g(x))^2 \right]. \end{aligned}$$

The first term is irrelevant to $g(x)$. The second term $2E[\epsilon(m(x) - g(x))] = 0$ is again irrelevant of $g(x)$. The third term, obviously, is minimized at $g(x) = m(x)$.

2 Linear Projection Model

As discussed in the previous section, we are interested in the conditional mean function $m(x)$. However, remind that

$$m(x) = E[y|x] = \int y f(y|x) dy$$

is a complex function of x , as it depends on the joint distribution of (y, x) .

A particular form of the conditional mean function is

$$m(x) = x'\beta,$$

a linear function of x .

Remark. The linear function is not as restrictive as one might thought. It can be used to generate some nonlinear (in random variables) effect if we re-define x . For example, if

$$y = x_1\beta_1 + x_2\beta_2 + x_1x_2\beta_3 + e,$$

then $\frac{\partial}{\partial x_1}m(x_1, x_2) = \beta_1 + x_2\beta_3$, which is nonlinear in x_1 , while it is still linear in the parameter β if we define a set of new regressors as $(\tilde{x}_1, \tilde{x}_2, \tilde{x}_3) = (x_1, x_2, x_1x_2)$.

Example. If $\begin{pmatrix} y \\ x \end{pmatrix} \sim N\left(\begin{pmatrix} \mu_y \\ \mu_x \end{pmatrix}, \begin{pmatrix} \sigma_y^2 & \rho\sigma_y\sigma_x \\ \rho\sigma_y\sigma_x & \sigma_x^2 \end{pmatrix}\right)$, then

$$E[y|x] = \mu_y + \rho\frac{\sigma_y}{\sigma_x}(x - \mu_x) = \left(\mu_y - \rho\frac{\sigma_y}{\sigma_x}\mu_x\right) + \rho\frac{\sigma_y}{\sigma_x}x.$$

Even though in general $m(x) \neq x'\beta$, the linear form $x'\beta$ is still useful as an approximation, as will be clear soon. Therefore, we may write the linear regression model, or the *linear projection model*, as

$$\begin{aligned} y &= x'\beta + e \\ E[xe] &= 0, \end{aligned}$$

where e is called the *projection error*, to be distinguished from $\varepsilon = y - m(x)$.

Remark. If a constant is included in x as a regressor, we have $E[e] = 0$.

The coefficient β in the linear projection model has a straightforward closed-form. Multiplying x on both sides and taking expectation, we have $E[xy] = E[xx']\beta$. If $E[xx']$ is invertible, we can explicitly solve

$$\beta = (E[xx'])^{-1} E[xy].$$

Now we justify $x'\beta$ as an approximation to $m(x)$. Indeed, $x'\beta$ is the optimal *linear* predictor in terms of MSE; in other words,

$$\beta = \arg \min_{b \in \mathbb{R}^K} E[(y - x'b)^2]. \quad (1)$$

This fact can be verified by taking the first-order condition of the above minimization problem $\frac{\partial}{\partial \beta} E[(y - x'\beta)^2] = 2E[x(y - x'\beta)] = 0$.

In the meantime, $x'\beta$ is also the best *linear* approximation to $m(x)$. If we replace y in (1) by $m(x)$, we solve the minimizer as

$$(E[xx'])^{-1} E[xm(x)] = (E[xx'])^{-1} E[E[xy|x]] = (E[xx'])^{-1} E[xy] = \beta.$$

Therefore β is also the best linear approximation to $m(x)$ in terms of MSE.

2.1 Subvector Regression

Sometimes we are interested in a subvector of β . For example, when we include an intercept and some variables in x , we are often more interested in the slope coefficients—the parameters associated with the random regressors—as they represent the size of effect of these explanatory factors. In such a regression

$$y = \beta_1 + x' \beta_2 + e,$$

we take an expectation to get $E[y] = \beta_1 + E[x]' \beta_2$. Take the difference of the two equations,

$$y - E[y] = (x - E[x])' \beta_2.$$

Therefore, we can explicitly solve β_2 as

$$\beta_2 = (E[(x - E[x])(x - E[x])'])^{-1} E[(x - E[x])(y - E[y])] = (\text{var}(x))^{-1} \text{cov}(x, y),$$

This is a special case of the subvector regression.

To discuss the general case, we need to know *the formula of the inverse of a partitioned matrix*. If $Q = \begin{pmatrix} Q_{11} & Q_{12} \\ Q_{21} & Q_{22} \end{pmatrix}$ is a symmetric and positive definite matrix, then

$$Q^{-1} = \begin{pmatrix} (Q_{11} - Q_{12}Q_{22}^{-1}Q_{21})^{-1} & -(Q_{11} - Q_{12}Q_{22}^{-1}Q_{21})^{-1}Q_{12}Q_{22}^{-1} \\ -(Q_{22} - Q_{21}Q_{11}^{-1}Q_{12})^{-1}Q_{21}Q_{11}^{-1} & (Q_{22} - Q_{21}Q_{11}^{-1}Q_{12})^{-1} \end{pmatrix}.$$

We apply the above formulate to the expression of β , and we obtain $\beta_1 = A_{11 \cdot 2}^{-1} A_{1y \cdot 2}$, where

$$\begin{aligned} A_{11 \cdot 2} &= E[x_1 x_1'] - E[x_1 x_2'] (E[x_2 x_2'])^{-1} E[x_2 x_1'] \\ A_{1y \cdot 2} &= E[x_1 y] - E[x_1 x_2'] (E[x_2 x_2'])^{-1} E[x_2 y]. \end{aligned}$$

This is a brutal force approach for the explicit expression of the subvector β_1 .

Alternatively, we can proceed in two steps. First, we run a multiple regression¹

$$\begin{aligned} x_1 &= x_2' \gamma + u \\ E[x_2 u] &= 0 \end{aligned}$$

so that the regressor error

$$u = x_1 - x_2' \gamma = x_1 - x_2' (E[x_2 x_2'])^{-1} E[x_2 x_1'] = x_1 - E[x_1 x_2'] (E[x_2 x_2'])^{-1} x_2.$$

We then run a simple regression of y on u , and the coefficient is

$$\theta = (E[uu'])^{-1} E[u'y].$$

The nominator is

$$E[u'y] = E[x_1 y] - E[x_1 x_2'] (E[x_2 x_2'])^{-1} E[x_2 y] = A_{1y \cdot 2}$$

¹We do allow x_1 to be a vector. However, one may find it is easier to consider the special case that x_1 is a scalar random variable.

and the denominator is

$$E[uu'] = E\left[\left(x_1 - E[x_1x_2'] (E[x_2x_2'])^{-1} x_2\right) \left(x_1 - E[x_1x_2'] (E[x_2x_2'])^{-1} x_2\right)'\right] = A_{11.2}.$$

It turns out $\beta_2 = \theta$.

While we can derive the expression of β_1 as a subvector of β , why do we come up with the two-step derivation? The latter makes clear that the coefficient represents the *partial effect* of the associate random variable.

2.2 Omitted Variable Bias

We write the *long regression* as

$$y = x_1'\beta_1 + x_2'\beta_2 + e,$$

and the *short regression* as

$$y = x_1'\gamma + u.$$

If β_1 in the long regression is the parameter of interest, omitting x_2 as in the short regression will render *omitted variable bias* (meaning $\gamma \neq \beta_1$) unless x_1 and x_2 are uncorrelated.

We first demean all the variables in the two regressions, which is equivalent as if we project out the effect of the constant. The long regression becomes

$$\tilde{y} = \tilde{x}_1'\beta_1 + \tilde{x}_2'\beta_2 + e,$$

and the short regression becomes

$$\tilde{y} = \tilde{x}_1' \gamma + u,$$

where *tilde* denotes the demeaned variable.

After demeaning, the cross-moment equals to the covariance. The short regression coefficient

$$\begin{aligned} \gamma &= (E[\tilde{x}_1 \tilde{x}_1'])^{-1} E[\tilde{x}_1 \tilde{y}] \\ &= (E[\tilde{x}_1 \tilde{x}_1'])^{-1} E[\tilde{x}_1 (\tilde{x}_1' \beta_1 + \tilde{x}_2' \beta_2 + e)] \\ &= \beta_1 + (E[\tilde{x}_1 \tilde{x}_1'])^{-1} E[\tilde{x}_1 \tilde{x}_2'] \beta_2. \end{aligned}$$

Therefore, $\gamma = \beta_1$ if and only if $E[\tilde{x}_1 \tilde{x}_2'] \beta_2 = 0$, which demands either $E[\tilde{x}_1 \tilde{x}_2'] = 0$ or $\beta_2 = 0$.

Obviously we prefer to run the long regression to attain β_1 if possible. However, sometimes x_2 is simply unobservable so the long regression is infeasible. When only the short regression is available, in some cases we are able to sign the bias, meaning that we know whether γ is bigger or smaller than β_1 .