

Advanced Econometrics

Lecture 5: The Algebra of Least Squares (Hansen Chapter 3)

Instructor: Ma, Jun

Renmin University of China

October 10, 2021

Samples

- ▶ Consider the best linear predictor of Y given \mathbf{X} for a pair of random variables $(Y, \mathbf{X}) \in \mathbb{R} \times \mathbb{R}^k$ with joint distribution F and call this the linear projection model. We are interested in estimating the projection coefficients

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{X}\mathbf{X}'))^{-1} \mathbb{E}(\mathbf{X}Y).$$

- ▶ The dataset is $\{(Y_i, \mathbf{X}_i) : i = 1, \dots, n\}$. We call this the **sample** or the **observations**.
- ▶ From the viewpoint of empirical analysis, a dataset is an array of numbers often organized as a table, where the columns of the table correspond to distinct variables and the rows correspond to distinct observations.
- ▶ For empirical analysis, the dataset and observations are fixed in the sense that they are numbers presented to the researcher. For statistical analysis we need to view the dataset as random, or more precisely as a realization of a random process.

Assumption

The observations $\{(Y_1, \mathbf{X}_1), \dots, (Y_i, \mathbf{X}_i), \dots, (Y_n, \mathbf{X}_n)\}$ are identically distributed; they are draws from a common distribution F .

- ▶ In econometric theory, we refer to the underlying common distribution as the population. Some authors prefer the label the **data-generating-process** (DGP).
- ▶ In contrast we refer to the observations available to us $\{(Y_i, \mathbf{X}_i) : i = 1, \dots, n\}$ as the sample or dataset.

We can write the model as

$$Y_i = \mathbf{X}'_i \boldsymbol{\beta} + e_i,$$

where the linear projection coefficient $\boldsymbol{\beta}$ is defined as

$$\boldsymbol{\beta} = \underset{b \in \mathbb{R}^k}{\operatorname{argmin}} S(\mathbf{b}),$$

the minimizer of the expected squared error

$$S(\mathbf{b}) = \mathbb{E} \left((Y_i - \mathbf{X}'_i \mathbf{b})^2 \right),$$

and has the explicit solution

$$\boldsymbol{\beta} = (\mathbb{E}(\mathbf{X}_i \mathbf{X}'_i))^{-1} \mathbb{E}(\mathbf{X}_i Y_i).$$

Moment estimators

- ▶ Suppose that we are interested in the population mean μ of a random variable Y_i : $\mu = \mathbb{E}(Y_i)$.
- ▶ The mean μ is a function of the distribution F . To estimate μ given a sample $\{Y_1, \dots, Y_n\}$ a natural estimator is the sample mean $\hat{\mu} = \bar{Y} = n^{-1} \sum_{i=1}^n Y_i$.
- ▶ Now suppose that we are interested in a set of population means of possibly non-linear functions of a random vector \mathbf{Y} , say $\boldsymbol{\mu} = \mathbb{E}(\mathbf{h}(\mathbf{Y}_i))$. For example, we may be interested in the first two moments of Y_i , $\mathbb{E}(Y_i)$ and $\mathbb{E}(Y_i^2)$. In this case the natural estimator is the vector of sample means, $\hat{\boldsymbol{\mu}} = n^{-1} \sum_{i=1}^n \mathbf{h}(\mathbf{Y}_i)$.
- ▶ For example, $\hat{\mu}_1 = n^{-1} \sum_{i=1}^n Y_i$ and $\hat{\mu}_2 = n^{-1} \sum_{i=1}^n Y_i^2$. We call $\hat{\boldsymbol{\mu}}$ the moment estimator for $\boldsymbol{\mu}$.

- ▶ Now suppose that we are interested in a nonlinear function of a set of moments. For example,

$$\sigma^2 = \text{Var}(Y_i) = \mathbb{E}(Y_i^2) - (\mathbb{E}(Y_i))^2.$$

Many parameters of interest can be written as a function of moments of Y :

$$\beta = g(\mu), \text{ where } \mu = \mathbb{E}(h(Y_i)).$$

- ▶ In this context a natural estimator of β is obtained by replacing μ with $\hat{\mu}$:

$$\hat{\beta} = g(\hat{\mu}), \text{ where } \hat{\mu} = \frac{1}{n} \sum_{i=1}^n h(Y_i).$$

We call $\hat{\beta}$ a moment estimator of β . For example, the moment estimator of σ^2 is

$$\hat{\sigma}^2 = \hat{\mu}_2 - \hat{\mu}_1^2 = \frac{1}{n} \sum_{i=1}^n Y_i^2 - \left(\frac{1}{n} \sum_{i=1}^n Y_i \right)^2.$$

Least squares estimator

- ▶ The moment estimator of $S(\mathbf{b})$ is the sample average:

$$\begin{aligned}\widehat{S}(\mathbf{b}) &= \frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{X}'_i \mathbf{b})^2 \\ &= \frac{1}{n} SSE(\mathbf{b})\end{aligned}$$

where

$$SSE(\mathbf{b}) = \sum_{i=1}^n (Y_i - \mathbf{X}'_i \mathbf{b})^2$$

is called the sum-of-squared-errors function.

- ▶ Since the projection coefficient minimizes $S(\mathbf{b})$, the OLS estimator minimizes $\widehat{S}(\mathbf{b})$:

$$\widehat{\boldsymbol{\beta}} = \underset{\mathbf{b} \in \mathbb{R}^k}{\operatorname{argmin}} \widehat{S}(\mathbf{b}) = \underset{\mathbf{b} \in \mathbb{R}^k}{\operatorname{argmin}} SSE(\mathbf{b}).$$

Solving for least squares with one regressor

- ▶ Consider the case $k = 1$ so that the coefficient β is a scalar. Then

$$\begin{aligned}SSE(\beta) &= \sum_{i=1}^n (Y_i - X_i\beta)^2 \\ &= \left(\sum_{i=1}^n Y_i^2 \right) - 2\beta \left(\sum_{i=1}^n X_i Y_i \right) + \beta^2 \left(\sum_{i=1}^n X_i^2 \right).\end{aligned}$$

- ▶ The minimizer of $SSE(\beta)$ is

$$\hat{\beta} = \frac{\sum_{i=1}^n X_i Y_i}{\sum_{i=1}^n X_i^2}.$$

- ▶ The intercept-only model: $X_i = 1$ and

$$\hat{\beta} = \frac{\sum_{i=1}^n 1Y_i}{\sum_{i=1}^n 1^2} = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}.$$

Solving for least squares with multiple regressors

- ▶ Expand SSE to find

$$SSE(\mathbf{b}) = \sum_{i=1}^n Y_i^2 - 2\mathbf{b}' \sum_{i=1}^n \mathbf{X}_i Y_i + \mathbf{b}' \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \mathbf{b}.$$

- ▶ The first-order condition is

$$0 = \frac{\partial}{\partial \mathbf{b}} SSE(\hat{\boldsymbol{\beta}}) = -2 \sum_{i=1}^n \mathbf{X}_i Y_i + 2 \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \hat{\boldsymbol{\beta}},$$

which is actually a system of k equations with k unknowns.

- ▶ We find an explicit formula for the OLS:

$$\hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{X}_i Y_i \right).$$

- ▶ Alternatively, we can write the projection coefficient β as an explicit function of the moments $\mathbf{Q}_{XY} = \mathbb{E}(XY)$ and $\mathbf{Q}_{XX} = \mathbb{E}(XX')$. Their moment estimators are

$$\hat{\mathbf{Q}}_{XY} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i Y_i \text{ and } \hat{\mathbf{Q}}_{XX} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i'.$$

- ▶ The moment estimator of β replaces the population moments with the sample moments:

$$\begin{aligned} \hat{\beta} &= \hat{\mathbf{Q}}_{XX}^{-1} \hat{\mathbf{Q}}_{XY} \\ &= \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i Y_i \right) \\ &= \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i' \right)^{-1} \left(\sum_{i=1}^n \mathbf{X}_i Y_i \right) \end{aligned}$$

which is identical with OLS.

Least squares residuals

- ▶ Define the fitted value $\hat{Y}_i = \mathbf{X}'_i \hat{\boldsymbol{\beta}}$ and the residual $\hat{e}_i = Y_i - \hat{Y}_i = Y_i - \mathbf{X}'_i \hat{\boldsymbol{\beta}}$.
- ▶ Note that $Y_i = \hat{Y}_i + \hat{e}_i$ and $Y_i = \mathbf{X}'_i \hat{\boldsymbol{\beta}} + \hat{e}_i$.
- ▶ e_i is called error and \hat{e}_i is called residual. The OLS first-order condition implies $\sum_{i=1}^n \mathbf{X}_i \hat{e}_i = \mathbf{0}$.
- ▶ Alternatively,

$$\begin{aligned} \sum_{i=1}^n \mathbf{X}_i \hat{e}_i &= \sum_{i=1}^n \mathbf{X}_i (Y_i - \mathbf{X}'_i \hat{\boldsymbol{\beta}}) \\ &= \sum_{i=1}^n \mathbf{X}_i Y_i - \sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i \hat{\boldsymbol{\beta}} \\ &= \sum_{i=1}^n \mathbf{X}_i Y_i - \sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i \left(\sum_{i=1}^n \mathbf{X}_i \mathbf{X}'_i \right)^{-1} \left(\sum_{i=1}^n \mathbf{X}_i Y_i \right) \\ &= \mathbf{0}. \end{aligned}$$

- ▶ When \mathbf{X}_i contains a constant, $\sum_{i=1}^n \hat{e}_i = 0$.

Demeaned regressors

- ▶ Sometimes it is useful to separate the constant from the other regressors: $Y_i = \mathbf{X}'_i \boldsymbol{\beta} + \alpha + e_i$, where α is the intercept and \mathbf{X}_i does not contain a constant.
- ▶ The least-squares estimates and residuals can be written as $Y_i = \mathbf{X}'_i \hat{\boldsymbol{\beta}} + \hat{\alpha} + \hat{e}_i$.
- ▶ Then $\sum_{i=1}^n \hat{e}_i = 0$ and $\sum_{i=1}^n \mathbf{X}_i \hat{e}_i = \mathbf{0}$ can be written as

$$0 = \sum_{i=1}^n (Y_i - \mathbf{X}'_i \hat{\boldsymbol{\beta}} - \hat{\alpha})$$

$$\mathbf{0} = \sum_{i=1}^n \mathbf{X}_i (Y_i - \mathbf{X}'_i \hat{\boldsymbol{\beta}} - \hat{\alpha}).$$

- ▶ Inserting $\hat{\alpha} = \bar{Y} - \bar{X}'\hat{\beta}$ into the second equation:

$$\sum_{i=1}^n X_i \left((Y_i - \bar{Y}) - (X_i - \bar{X})' \hat{\beta} \right) = \mathbf{0}.$$

- ▶ Solving for $\hat{\beta}$ we find

$$\hat{\beta} = \left(\sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})' \right)^{-1} \left(\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \right).$$

Model in matrix notation

- ▶ We can stack these n equations together

$$Y_1 = \mathbf{X}'_1 \boldsymbol{\beta} + e_1$$

$$Y_2 = \mathbf{X}'_2 \boldsymbol{\beta} + e_2$$

$$\vdots$$

$$Y_n = \mathbf{X}'_n \boldsymbol{\beta} + e_n.$$

- ▶ Define

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} \mathbf{X}'_1 \\ \mathbf{X}'_2 \\ \vdots \\ \mathbf{X}'_n \end{pmatrix}, \mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}.$$

\mathbf{Y} and \mathbf{e} are $n \times 1$ vectors and \mathbf{X} is an $n \times k$ matrix.

- ▶ The system of n equations can be written as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}.$$

- ▶ Sample sums can be written in matrix notation:

$$\sum_{i=1}^n X_i X_i' = X'X \text{ and } \sum_{i=1}^n X_i Y_i = X'Y.$$

- ▶ Therefore the least-squares estimator can be written as

$$\hat{\beta} = (X'X)^{-1} (X'Y).$$

- ▶ The residual vector is $\hat{e} = Y - X\hat{\beta}$. We can write $\sum_{i=1}^n X_i \hat{e}_i = \mathbf{0}$ as $X'\hat{e} = \mathbf{0}$.

Projection matrix

- ▶ Define

$$P = X (X'X)^{-1} X'.$$

Observe

$$PX = X (X'X)^{-1} X'X = X.$$

This is a property of a **projection** matrix.

- ▶ For any matrix Z which can be written as $Z = X\Gamma$ for some matrix Γ ,

$$PZ = PX\Gamma = X (X'X)^{-1} X'X\Gamma = X\Gamma = Z.$$

- ▶ If we partition the matrix X into two matrices X_1 and X_2 so that

$$X = \begin{bmatrix} X_1 & X_2 \end{bmatrix}$$

then $PX_1 = X_1$.

The matrix P is symmetric and idempotent:

$$\begin{aligned} P' &= \left(X (X'X)^{-1} X' \right)' \\ &= (X')' \left((X'X)^{-1} \right)' (X)' \\ &= X ((X'X)')^{-1} X' \\ &= X ((X)' (X')')^{-1} X' \\ &= P \end{aligned}$$

$$\begin{aligned} PP &= P X (X'X)^{-1} X' \\ &= X (X'X)^{-1} X' \\ &= P. \end{aligned}$$

- ▶ The matrix P has the property that it creates the fitted values in a least-squares regression:

$$PY = X (X'X)^{-1} X'Y = X\hat{\beta} = \hat{Y}.$$

- ▶ A special example of a projection matrix occurs when $X = \mathbf{1}$ is an n -vector of ones. Then

$$\begin{aligned} P_1 &= \mathbf{1} (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}' \\ &= \frac{1}{n} \mathbf{1}\mathbf{1}'. \end{aligned}$$

- ▶ Note

$$\begin{aligned} P_1 Y &= \mathbf{1} (\mathbf{1}'\mathbf{1})^{-1} \mathbf{1}' Y \\ &= \mathbf{1} \bar{Y} \end{aligned}$$

creates an n -vector whose elements are the sample mean \bar{Y} of Y_i

Theorem

The i^{th} diagonal element of $\mathbf{P} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$ is

$$h_{ii} = \mathbf{X}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}_i.$$

$$\sum_{i=1}^n h_{ii} = \text{tr}\mathbf{P} = k$$

and $0 \leq h_{ii} \leq 1$.

$$\begin{aligned}\text{tr}\mathbf{P} &= \text{tr}\left(\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\right) \\ &= \text{tr}\left((\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}\right) \\ &= \text{tr}(\mathbf{I}_k) \\ &= k\end{aligned}$$

One implication is that the rank of \mathbf{P} is k .

Orthogonal projection

- ▶ Define

$$\begin{aligned} \mathbf{M} &= \mathbf{I}_n - \mathbf{P} \\ &= \mathbf{I}_n - \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'. \end{aligned}$$

- ▶ Note

$$\mathbf{MX} = (\mathbf{I}_n - \mathbf{P})\mathbf{X} = \mathbf{X} - \mathbf{PX} = \mathbf{X} - \mathbf{X} = \mathbf{0}.$$

Thus \mathbf{M} and \mathbf{X} are orthogonal. We call \mathbf{M} the **orthogonal projection matrix**.

- ▶ If $\mathbf{Z} = \mathbf{X}\mathbf{\Gamma}$, then

$$\mathbf{MZ} = \mathbf{Z} - \mathbf{PZ} = \mathbf{0}.$$

For example, $\mathbf{MX}_1 = \mathbf{0}$ for any subcomponent \mathbf{X}_1 of \mathbf{X} and $\mathbf{MP} = \mathbf{0}$.

- ▶ \mathbf{M} is symmetric ($\mathbf{M}' = \mathbf{M}$) and idempotent ($\mathbf{MM} = \mathbf{M}$).
 $\text{tr}\mathbf{M} = n - k$. The rank of \mathbf{M} is $n - k$.

- ▶ M creates least-square residuals:

$$MY = Y - PY = Y - X\hat{\beta} = \hat{e}.$$

- ▶ When $X = \mathbf{1}$,

$$\begin{aligned} M_1 &= I_n - P_1 \\ &= I_n - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}' \end{aligned}$$

and M_1 creates demeaned values $M_1Y = Y - \mathbf{1}\bar{Y}$.

- ▶ We find

$$\hat{e} = MY = M(X\beta + e) = Me.$$

Estimation of error variance

- ▶ If e_i were observed, we would estimate $\sigma^2 = \mathbb{E}(e_i^2)$ by $\tilde{\sigma}^2 = n^{-1} \sum_{i=1}^n e_i^2$. This is infeasible as e_i is not observed.
- ▶ The feasible estimator: $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n \hat{e}_i^2$. In matrix notation, $\tilde{\sigma}^2 = n^{-1} \mathbf{e}'\mathbf{e}$ and $\hat{\sigma}^2 = n^{-1} \hat{\mathbf{e}}'\hat{\mathbf{e}}$.
- ▶ Since $\hat{\mathbf{e}} = \mathbf{M}\mathbf{Y} = \mathbf{M}\mathbf{e}$,

$$\hat{\sigma}^2 = n^{-1} \hat{\mathbf{e}}'\hat{\mathbf{e}} = n^{-1} \mathbf{Y}'\mathbf{M}\mathbf{M}\mathbf{Y} = n^{-1} \mathbf{Y}'\mathbf{M}\mathbf{Y} = n^{-1} \mathbf{e}'\mathbf{M}\mathbf{e}.$$

- ▶ An implication:

$$\tilde{\sigma}^2 - \hat{\sigma}^2 = n^{-1} \mathbf{e}'\mathbf{e} - n^{-1} \mathbf{e}'\mathbf{M}\mathbf{e} = n^{-1} \mathbf{e}'\mathbf{P}\mathbf{e} \geq 0.$$

Analysis of variance

- ▶ Write

$$Y = PY + MY = \hat{Y} + \hat{e},$$

where

$$\hat{Y}'\hat{e} = (PY)'(MY) = Y'PMY = 0.$$

- ▶ Then

$$Y'Y = \hat{Y}'\hat{Y} + 2\hat{Y}'\hat{e} + \hat{e}'\hat{e} = \hat{Y}'\hat{Y} + \hat{e}'\hat{e}$$

or

$$\sum_{i=1}^n Y_i^2 = \sum_{i=1}^n \hat{Y}_i^2 + \sum_{i=1}^n \hat{e}_i^2.$$

- Since $Y = \hat{Y} + \hat{e}$,

$$Y - \mathbf{1}\bar{Y} = \hat{Y} - \mathbf{1}\bar{Y} + \hat{e},$$

where

$$\left(\hat{Y} - \mathbf{1}\bar{Y}\right)' \hat{e} = \hat{Y}' \hat{e} - \bar{Y} \mathbf{1}' \hat{e} = 0.$$

- Then

$$(Y - \mathbf{1}\bar{Y})' (Y - \mathbf{1}\bar{Y}) = \left(\hat{Y} - \mathbf{1}\bar{Y}\right)' \left(\hat{Y} - \mathbf{1}\bar{Y}\right) + \hat{e}' \hat{e}$$

or

$$\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 + \sum_{i=1}^n \hat{e}_i^2.$$

This is commonly called the analysis-of-variance formula for least squares regression.

- ▶ A commonly reported statistic is the R^2 :

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n \hat{e}_i^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}.$$

This is a measure of goodness of regression fit.

- ▶ One deficiency with R^2 is that it increases when regressors are added to a regression so the “fit” can be always increased by increasing the number of regressors.

Regression components

- ▶ Partition

$$X = [X_1 \quad X_2]$$

and

$$\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}.$$

- ▶ Then the regression model can be rewritten as

$$Y = X_1\beta_1 + X_2\beta_2 + e$$

and

$$Y = X\hat{\beta} + \hat{e} = X_1\hat{\beta}_1 + X_2\hat{\beta}_2 + \hat{e}.$$

- ▶ We show that

$$\begin{aligned} \hat{\beta}_1 &= (X_1' M_2 X_1)^{-1} (X_1' M_2 Y) \\ \hat{\beta}_2 &= (X_2' M_1 X_2)^{-1} (X_2' M_1 Y). \end{aligned}$$

Residual regression

► Note

$$\begin{aligned}\hat{\beta}_2 &= (X_2' M_1 X_2)^{-1} (X_2' M_1 Y) \\ &= (X_2' M_1 M_1 X_2)^{-1} (X_2' M_1 M_1 Y) \\ &= (\tilde{X}_2' \tilde{X}_2)^{-1} (\tilde{X}_2' \tilde{e}_1)\end{aligned}$$

where

$$\tilde{X}_2 = M_1 X_2 \text{ and } \tilde{e}_1 = M_1 Y.$$

- The estimate $\hat{\beta}_2$ is algebraically equal to the least-squares regression of \tilde{e}_1 on \tilde{X}_2 . \tilde{e}_1 is the least-squares residuals from a regression of Y on X_1 . The columns of \tilde{X}_2 are the least-squares residuals from the regressions of the columns of X_2 on X_1 .

Theorem (Frisch-Waugh-Lovell (FWL))

The OLS estimator of β_2 and the OLS residuals \hat{e} may be equivalently computed by either the OLS regression or via the following algorithm:

1. Regress Y on X_1 , obtain residuals \tilde{e}_1 ;
2. Regress X_2 on X_1 , obtain residuals \tilde{X}_2 ;
3. Regress \tilde{e}_1 on \tilde{X}_2 , obtain OLS estimates $\hat{\beta}_2$ and residuals \hat{e} .

► To check (3), note

$$\begin{aligned} \left(I_n - \tilde{X}_2 \left(\tilde{X}_2' \tilde{X}_2 \right)^{-1} \tilde{X}_2' \right) \tilde{e}_1 &= M_1 Y - M_1 X_2 \hat{\beta}_2 = \\ M_1 \left(X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 + \hat{e} \right) - M_1 X_2 \hat{\beta}_2 &= M_1 \hat{e} = \hat{e}. \end{aligned}$$