

1 Linear Regression

1.1 Linear regression and Least Square Solution

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where:

- \mathbf{y} is $n \times 1$ vector of outcomes
- \mathbf{X} is $n \times k$ matrix of regressors (full column rank)
- $\boldsymbol{\beta}$ is $k \times 1$ parameter vector
- $\boldsymbol{\epsilon}$ is $n \times 1$ error vector

Assumptions

1. Linear
2. X matrix has full rank. In other words, no multicollinearity.
2. error term has zero mean $\mathbb{E}[\epsilon|X] = 0$
3. Homoscedasticity or equal variance of ϵ . In other words, no autocorrelation between disturbances. $cov(\epsilon_i, \epsilon_j) = 0$.
6. Number of observations n must be greater than the number of parameters.

Least Square Solution

The cost function is given by

$$f(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}$$

Since third term are scalar,

$$\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} = (\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y})^T = \mathbf{y}^T \mathbf{X}\boldsymbol{\beta}$$

$$f(\boldsymbol{\beta}) = \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta} = \mathbf{y}^T \mathbf{y} - 2(\mathbf{X}^T \mathbf{y})^T \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}$$

The first term is a constant and its derivative is zero.

The derivative of 2nd term

Consider the derivative of $\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}$ with respect to $\boldsymbol{\beta}$.

$$\begin{aligned} \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} &= \sum \alpha_i \beta_i \\ \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}}{\partial \beta_i} &= \alpha_i \end{aligned}$$

Write the derivative in matrix form

$$\begin{pmatrix} \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}}{\partial \beta_1} \\ \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}}{\partial \beta_2} \\ \dots \\ \frac{\partial \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y}}{\partial \beta_3} \end{pmatrix} = \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_p \end{pmatrix}$$

So if we let $\alpha = X^T Y$, we have

$$\frac{\partial 2(X^T Y)^T \beta}{\partial \beta} = 2X^T Y$$

The derivative of 3rd term

let $A = X^T X$,

$$\beta^T X^T X \beta = \beta^T \begin{pmatrix} \Sigma_i A_{1k} \beta_k \\ \Sigma_i A_{2k} \beta_k \\ \dots \\ \Sigma_k A_{pk} \beta_k \end{pmatrix} = \Sigma_j \beta_j (\Sigma_k A_{jk} \beta_k)$$

To calculate the derivative of $f(\beta)$, we note there are only 3 cases that the derivative does not vanish

1) $l = j = k$

$$\frac{f(\beta)}{\partial \beta_l} = 2A_{ll} \beta_l$$

2) $l=j, j \neq k$

$$\frac{f(\beta)}{\partial \beta_l} = \Sigma_{k, k \neq l} A_{lk} \beta_k$$

3) $l=k, j \neq k$

$$\frac{f(\beta)}{\partial \beta_l} = \Sigma_{j, j \neq l} A_{jl} \beta_j = \Sigma_{j, j \neq l} A_{lj}^T \beta_j$$

Therefore

$$\begin{aligned} \frac{f(\beta)}{\partial \beta_l} &= A_{ll} \beta_l + \Sigma_{k, k \neq l} A_{lk} \beta_k + A_{ll} \beta_l + \Sigma_{j, j \neq l} A_{lj}^T \beta_j \\ &= \Sigma_k A_{lk} \beta_k + \Sigma_j A_{lj}^T \beta_j \end{aligned}$$

The first term is the l th row of vector $A\beta = X^T X\beta$, and the 2nd term is the l th row of vector $A^T \beta = X^T X\beta$. So we put the whole derivative in matrix form

$$\frac{f(\beta)}{\partial \beta} = -2X^T Y + 2X^T X \beta$$

which is a $p \times 1$ vector with each row corresponding to the derivative with respect to β_i letting the derivative equal to zero yields the **normal equation** and the estimation of β

Normal equation

$$(X^T X) \hat{\beta} = X^T Y$$

Estimator of β

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

Using the rule of matrix multiplication, we can rewrite $X^T X$ as the f

$$\begin{aligned}
X^T X &= \begin{pmatrix} \sum_{i=1}^N x_{1i}^T x_{i1} & \sum_{i=1}^N x_{1i}^T x_{i2} & \dots & \sum_{i=1}^N x_{1i}^T x_{ik} \\ \dots & \dots & \dots & \dots \\ \sum_{i=1}^N x_{ki}^T x_{i1} & \sum_{i=1}^N x_{ki}^T x_{i2} & \dots & \sum_{i=1}^N x_{ki}^T x_{ik} \end{pmatrix} \\
&= \sum_{i=1}^N \begin{pmatrix} x_{i1} x_{i1} & x_{i1} x_{i2} & \dots & x_{i1} x_{ik} \\ \dots & \dots & \dots & \dots \\ x_{ik} x_{i1} & x_{ik} x_{i2} & \dots & x_{ik} x_{ik} \end{pmatrix} \\
&= \sum_{i=1}^N \begin{pmatrix} x_{i1} \\ \dots \\ x_{ik} \end{pmatrix} \begin{pmatrix} x_{i1} & \dots & \dots & x_{ik} \end{pmatrix} \\
&= \sum_{i=1}^N X_i X_i^T
\end{aligned}$$

where X_i is $K \times 1$ vector, and $X_i X_i^T$ is $K \times K$ matrix. Similarly,

$$X^T Y = \sum_{i=1}^N X_i y_i$$

where X_i is $K \times 1$ vector, y_i is a scalar, and $X_i y_i$ is $K \times 1$ vector. So the estimator of β is

$$\hat{\beta} = \left(\sum_{i=1}^N X_i X_i^T \right)^{-1} \left(\sum_{i=1}^N X_i y_i \right)$$

Least Square Estimator for Simple Linear Regression

$$y = \beta_0 + \beta_1 X + \epsilon$$

$$\begin{aligned}
&\begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \\
&= (X^T X)^{-1} X^T Y \\
&= \left(\begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} \\
&= \frac{1}{n \sum x_i^2 - (\sum x_i)^2} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} \begin{pmatrix} \sum y_i \\ -\sum x_i y_i \end{pmatrix}
\end{aligned}$$

So

$$\beta_0 = \frac{\sum x_i^2 \sum y_i - \sum x_i (\sum x_i y_i)}{n \sum x_i^2 - (\sum x_i)^2} \quad (1)$$

$$\beta_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2} \quad (2)$$

β_1 can also be written using the covariance

$$\beta_1 = \frac{\sum_i^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_i^n (x_i - \bar{x})(x_i - \bar{x})} \quad (3)$$

And it is easy to show

$$\begin{aligned} \beta_2 &= \frac{\sum_i^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_i^n (x_i - \bar{x})(x_i - \bar{x})} \\ &= \frac{\sum_i^n (x_i y_i - \bar{x} y_i - x_i \bar{y} + \bar{x} \bar{y})}{\sum_i^n (x_i^2 - 2\bar{x} x_i + (\bar{x})^2)} \\ &= \frac{\sum_i^n x_i y_i - \sum_i^n \bar{x} y_i - \sum_i^n x_i \bar{y} + \sum_i^n \bar{x} \bar{y}}{\sum_i^n x_i^2 - \sum_i^n 2\bar{x} x_i + \sum_i^n (\bar{x})^2} \\ &= \frac{\sum_i^n x_i y_i - (\frac{1}{n} \sum_j^n x_j)(\sum_i^n y_i) - (\sum_i^n x_i)(\frac{1}{n} \sum_j^n y_j) + \sum_i^n (\frac{1}{n} \sum_j^n x_i)(\frac{1}{n} \sum_k^n y_i)}{\sum_i^n x_i^2 - \sum_i^n 2(\frac{1}{n} \sum_j^n x_j)x_i + \sum_i^n (\frac{1}{n} \sum_j^n x_j)^2} \\ &= \frac{\sum_i^n x_i y_i - (\frac{1}{n} \sum_j^n x_j)(\sum_i^n y_i) - (\sum_i^n x_i)(\frac{1}{n} \sum_j^n y_j) + n(\frac{1}{n} \sum_j^n x_i)(\frac{1}{n} \sum_k^n y_k)}{\sum_i^n x_i^2 - \sum_i^n 2(\frac{1}{n} \sum_j^n x_j)x_i + n(\frac{1}{n} \sum_j^n x_j)^2} \\ &= \frac{\sum_i^n x_i y_i - \frac{1}{n}(\sum_i^n x_i)(\sum_j^n y_j) - \frac{1}{n}(\sum_i^n x_i)(\sum_j^n y_j) + \frac{1}{n}(\sum_j^n x_i)(\sum_k^n y_k)}{\sum_i^n x_i^2 - \frac{2}{n}(\sum_j^n x_j)(\sum_i^n x_i) + \frac{1}{n}(\sum_j^n x_j)^2} \\ &= \frac{\sum_i^n x_i y_i - \frac{1}{n}(\sum_i^n x_i)(\sum_j^n y_j)}{\sum_i^n x_i^2 - \frac{1}{n}(\sum_j^n x_j)(\sum_i^n x_i)} \\ &= \frac{n \sum_i^n x_i y_i - (\sum_i^n x_i)(\sum_j^n y_j)}{n \sum_i^n x_i^2 - (\sum_j^n x_j)(\sum_i^n x_i)} \\ &= \frac{n \sum x_i y_i - (\sum x_i)(\sum y_j)}{n \sum x_i^2 - (\sum x_i)^2} \end{aligned}$$

which is the same as Eq.2. We can interpret β as ratio of the covariance of x and y to the variance of x.

1.2 Projection matrix

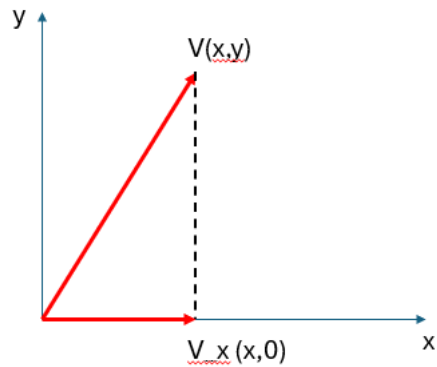
Given $\hat{\beta} = (X^T X)^{-1} X^T Y$, we have the predictor value of $y = X\beta$

$$\hat{y} = X(X^T X)^{-1} X^T y$$

The matrix $P = X(X^T X)^{-1} X^T$ is a projection matrix. It projects the vector of y into the column space of X.

Understand the word projection

Let us understand this first through geometry point of view. Consider a vector on 2 dimensional space, $V_1 = (x_1, y_1)^T$, where x_1 and y_1 are the x and y component, respectively. If we project the vector V into x-line, then apparently we get $V_x = (x_1, 0)^T$, see graph below.



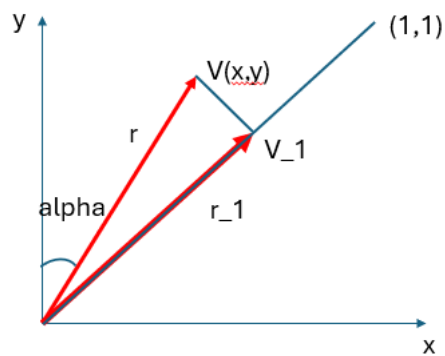
If we have a vector that is along the x axis

$$X = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

The projection matrix of a vector into x line is

$$\begin{aligned} P_x &= x(x^T x)^{-1} x^T \\ &= \begin{pmatrix} 1 \\ 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \end{pmatrix} \\ &= \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \end{aligned}$$

Applying this projection matrix to any 2 dimensional vector V gives $(V_x, 0)^T$. So it projects the vector into x line. Let us take another example. Imagine V_1 is vector if we project V onto the line that has 45 degree angle with x axis. See below.



In order to calculate V_1 , we see

$$r_1 = r \cos(\pi/4 - \alpha) = r \left(\frac{\sqrt{2}}{2} \frac{y}{r} + \frac{\sqrt{2}}{2} \frac{x}{r} \right) = \frac{\sqrt{2}}{2} y + \frac{\sqrt{2}}{2} x$$

$$V_{1x} = r_1 \cos(\pi/4) = \frac{x+y}{2}$$

$$V_{1y} = r_1 \sin(\pi/4) = \frac{x+y}{2}$$

After we understand this using geometry point of view, we can work out from algebra point of view. The vector we want to project onto is

$$i = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

The projection matrix of a vector into x line is

$$\begin{aligned} P_x &= x(x^T x)^{-1} x^T \\ &= \begin{pmatrix} 1 \\ 1 \end{pmatrix} \left(\begin{pmatrix} 1 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right)^{-1} \begin{pmatrix} 1 & 1 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \end{aligned}$$

Therefore we easily see

$$V_1 = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(x+y) \\ \frac{1}{2}(x+y) \end{pmatrix}$$

which is the same as what we get based on geometry. For n dimensional vector y, if our X matrix has rank of k, then the projection matrix P projects the vector y into k dimensional hyperplane. For example, if we define

$$i_N = \begin{pmatrix} 1 \\ 1 \\ \dots \\ 1 \end{pmatrix}$$

The projection matrix P is

$$P = i \frac{1}{N} i^T = \frac{1}{N} \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \dots & \dots & \dots & \dots \\ 1 & 1 & \dots & 1 \end{pmatrix}$$

Projection matrix into null space

If P is a projection matrix, the matrix $I - P$ is also a projection matrix. In linear regression model

$$y = X\beta + \epsilon$$

$$P = X(X^T X)^{-1} X^T$$

Define residual vector $\hat{\epsilon}$

$$\hat{\epsilon} = (I - P)y = (I - X(X^T X)^{-1} X^T)y$$

And it is easy to show $\hat{\epsilon}$ and X are orthogonal.

$$X^T \hat{\epsilon} = X^T (I - P)y = X^T (I - X(X^T X)^{-1} X^T)y = (X^T - X^T X(X^T X)^{-1} X^T)y = 0y = 0$$

For the above example, we define $M = I - \frac{1}{N}ii^T$, and My express the mean deviations of a vector.

Idempotent property of projection matrix

Consider the previous example that we project a vector V onto x axis, how about we do this projection twice, we would end up the same vector V_x . Using a little matrix algebra, it is easy to prove that for any project matrix P , we have $PP = P$.

1.3 Partitioned Regression and Regression

$$y = X\beta + \epsilon = X_1\beta_1 + X_2\beta_2 + \epsilon$$

The normal equation is

$$\begin{pmatrix} X_1^T X_1 & X_1^T X_2 \\ X_2^T X_1 & X_2^T X_2 \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1^T y \\ X_2^T y \end{pmatrix}$$

If X_1 and X_2 are orthogonal, namely, $X_1^T X_2 = 0$, then

$$\begin{aligned} \hat{\beta}_1 &= (X_1^T X_1)^{-1} X_1^T y \\ \hat{\beta}_2 &= (X_2^T X_2)^{-1} X_2^T y \end{aligned}$$

If X_1 and X_2 are not orthogonal, we can solve for β_2 in the above normal equation set and get β_2

$$\begin{aligned} \hat{\beta}_2 &= [X_2^T (I - X_1(X_1^T X_1)^{-1} X_1^T) X_2]^{-1} [X_2^T (I - X_1(X_1^T X_1)^{-1} X_1^T) y] \\ &= (X_2^T M_1 X_2)^{-1} (X_2^T M_1 y) \end{aligned}$$

Given the fact that M_1 is symmetrical and idempotent, we can rewrite the above expression

$$\begin{aligned} \hat{\beta}_2 &= (X_2^T M_1 M_1 X_2)^{-1} (X_2^T M_1 M_1 y) \\ &= (X_2^T M_1^T M_1 X_2)^{-1} (X_2^T M_1^T M_1 y) \\ &= ((M_1 X_2)^T M_1 X_2)^{-1} ((M_1 X_2)^T M_1 y) \end{aligned} \tag{4}$$

The above uses the property that $M_1^T = M_1$ and $M_1 M_1 = M_1$

The $\hat{\beta}_2$ is also the solution of

$$M_1 Y = M_1 X_2 \beta_2 + \epsilon$$

where $M_1 y$ is the residual of y regressed on X_1 and $M_1 X_2$ is the residual of X_2 regressed on X_1 . For example, in simple linear regression

$$Y = \beta_0 + x\beta_1$$

Where $X_1 = 1_N$, so its projection matrix is $i\frac{1}{N}i^T$, and the corresponding M matrix is $I - \frac{1}{N}ii^T$. We try to calculate β using partition regression.

$$MY = (I - \frac{1}{N}ii^T)Y = Y - \bar{Y}MX = (I - \frac{1}{N}ii^T)Y = X - \bar{X}$$

Then

$$\beta_1 = ((MX)^T(MX))^{-1}((MX)^T(MY)) = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2} \quad (5)$$

which is the same as Eq.3

1.4 Variance componet identity

If we define our mean projection matrix P

$$P = i\frac{1}{N}i^T$$

and silimilarly we define mean deviation project matrix

$$M = I - P = i\frac{1}{N}i^T$$

We have

$$y = \hat{y} + \hat{\epsilon} = X\hat{\beta} + \hat{\epsilon}$$

Multiplying M matrix on the left, we have

$$My = MX\hat{\beta} + M\hat{\epsilon} = MX\hat{\beta} + \hat{\epsilon}$$

$$\begin{aligned} (My)^2 &= (MX\hat{\beta} + \hat{\epsilon})^T(MX\hat{\beta} + \hat{\epsilon}) \\ &= (\beta^T X^T M^T + \hat{\epsilon}^T)(MX\hat{\beta} + \hat{\epsilon}) \\ &= (MX\hat{\beta})^2 + \beta^T X^T M^T \hat{\epsilon} + (\beta^T X^T M^T \hat{\epsilon})^T + (\hat{\epsilon})^2 \end{aligned}$$

The 2nd and 3rd terms are zero because that 1) $\hat{\epsilon}$ has zero mean, so $M^T \hat{\epsilon} = M\hat{\epsilon} = \hat{\epsilon}$ and 2) $X^T \hat{\epsilon} = 0$, so

$$(My)^2 = (MX\hat{\beta})^2 + (\hat{\epsilon})^2$$

Rewriting the above equation using summation, we have

$$\sum_i (y_i - \bar{y})^2 = \sum_i (\bar{y}_i - \bar{\hat{y}})^2 + \sum_i (y_i - \hat{y})^2$$

Define

$$\begin{aligned} SST &= \sum_i (y_i - \bar{y})^2 \\ SSR &= \sum_i (\bar{y}_i - \bar{\hat{y}})^2 \\ SSE &= \sum_i (y_i - \hat{y})^2 \end{aligned}$$

SSE can also be written as

$$\begin{aligned} SSE &= SST - SSR \\ &= SST - (MX\hat{\beta})^2 \\ &= SST - ((MX)((MX)^T(MX))^{-1}(MX)^T(MY))^2 \end{aligned}$$

Let $U = MX$, and $V = MY$ then

$$\begin{aligned} SSE &= SST - (Z(Z^T Z)^{-1} Z^T)^2 \\ &= SST - (U(U^T U)^{-1} U^T V)^T (U(U^T U)^{-1} U^T V) \\ &= SST - V^T U(U^T U)^{-1} U^T U(U^T U)^{-1} U^T V \\ &= SST - V^T U(U^T U)^{-1} U^T V \\ &= SST - (MY)^T (MX) ((MX)^T (MX))^{-1} (MX)^T (MY) \end{aligned}$$

Define

$$\begin{aligned} S_{xx} &= (MX)^T (MX) = \sum_i^N (x_i - \bar{x})^2 \\ S_{xy} &= (MX)^T (MY) = \sum_i^N (x_i - \bar{x})(y_i - \bar{y}) \end{aligned}$$

So

$$SSE = SST - S_{xy}^T S_{xx}^{-1} S_{xy}$$

Then we have

$$SST = SSR + SSE$$

1.5 Variance of $\hat{\beta}$ and σ^2 estimation

$$\begin{aligned} Var(\hat{\beta}) &= Var((X^T X)^{-1} X^T \epsilon) = (X^T X)^{-1} X^T Var(\epsilon) ((X^T X)^{-1} X^T)^T \\ &= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} \end{aligned}$$

The above derivation use the fact that ϵ has a normal distribution with mean 0 and variance σ^2 . For simple linear regression

$$Var(\hat{\beta}) = \frac{\sigma^2}{n\Sigma x_i^2 - (\Sigma x_i)^2} \begin{pmatrix} \Sigma x_i^2 & -\Sigma x_i \\ -\Sigma x_i & n \end{pmatrix}$$

$$Var(\hat{\beta}_0) = \frac{\Sigma x_i^2 \sigma^2}{n\Sigma x_i^2 - (\Sigma x_i)^2}$$

$$Var(\hat{\beta}_1) = \frac{n\sigma^2}{n\Sigma x_i^2 - (\Sigma x_i)^2}$$

Try

$$\begin{aligned} \Sigma(x_i - \bar{x})^2 &= \Sigma(x_i^2 - 2\bar{x}x_i + \bar{x}^2) = \Sigma x_i^2 - 2(\Sigma_j \frac{x_j}{n})x_i + \frac{(\Sigma_j x_j)^2}{n^2} \\ &= \Sigma x_i^2 - \frac{2}{n}(\Sigma_i x_i)^2 + \frac{(\Sigma_i x_i)^2}{n} = \Sigma x_i^2 - \frac{1}{n}(\Sigma x_i)^2 \end{aligned}$$

So

$$Var(\hat{\beta}_0) = \frac{\Sigma x_i^2 \sigma^2}{n\Sigma(x_i - \bar{x})^2}$$

$$Var(\hat{\beta}_1) = \frac{n\sigma^2}{n\Sigma(x_i - \bar{x})^2} = \frac{\sigma^2}{\Sigma(x_i - \bar{x})^2}$$

$$\begin{aligned} SSE &= \Sigma_i (y - \hat{y}_i)^2 \\ &= (Y - X\beta)^T (Y - X\beta) \\ &= (Y - X(X^T X)^{-1} X^T Y)^T (Y - X(X^T X)^{-1} X^T Y) \\ &= (Y - PY)^T (Y - PY) \\ &= Y^T (1 - P)^T (1 - P) Y = Y^T (1 - P) Y \\ &= (X\beta + \epsilon)^T (1 - P) (X\beta + \epsilon) \\ &= \beta^T X^T (1 - P) X \beta + 2\beta^T X^T (1 - P) \epsilon + \epsilon^T (1 - P) \epsilon \end{aligned}$$

$$E[SSE] = E[\epsilon^T (1 - P) \epsilon] = E[\epsilon^T \epsilon] \text{trace}(1 - P) = \sigma^2(n - k)$$

We obtain the unbiased estimator of σ^2

$$\hat{\sigma}^2 = \frac{SSE}{n - k}$$

Therefore the estimator of variance of β

$$\hat{Var}(\hat{\beta}_i) = \hat{\sigma}^2 (X^T X)^{-1}_{ii}$$

and the stanard error of β_i is

$$SE(\hat{\beta}_i) = \sqrt{\hat{\sigma}^2 (X^T X)^{-1}_{ii}}$$

2 Properties of Least Square Estimators

When we have an estimator, we need to evaluate how good our estimator is? A few questions we can ask are 1) how far is the value of our estimator away from the true value, even in the ideal case when the sample size is infinite? 2) when 1) is true, with finite sample size, does the value of our estimator approach the true value as the sample size increases? In other words, does the estimator converge to the true value as sample size goes to infinity? 3) when 1) and 2) are true, as the sample size increases, how fast does our estimator converge to the true value? 4) with 1) 2) and 3), what is the asymptotic distribution of the estimator? If the distribution is normal, it can be used to do interval estimation such as confidence interval. The 1st question defines unbiasedness, the 2nd one defines consistency, and the 3rd one defines efficiency.

2.1 Unbiasness

Unbiased

$$\begin{aligned}\hat{\beta} &= (X^T X)^{-1} X^T Y \\ &= (X^T X)^{-1} X^T (X\beta + \epsilon) \\ &= (X^T X)^{-1} X^T X\beta + (X^T X)^{-1} X^T \epsilon \\ &= \beta + (X^T X)^{-1} X^T \epsilon\end{aligned}$$

Then the expectation of $\hat{\beta}$ condition on X is

$$E[\hat{\beta}|X] = \beta + (X^T X)^{-1} X^T E(\epsilon|X)$$

t

The last term is zero by assumption of linear regression. So

$$E[\hat{\beta}] = \beta$$

The expectation of the estimator is the same as true value, this is called **unbiased**.

Bias due to omission of relevant variables

Suppose we have a model

$$y = X_1\beta_1 + X_2\beta_2 + \epsilon$$

If we regress y on X_1 only, our estimator is

$$\hat{\beta}_1 = (X_1^T X_1)^{-1} X_1^T y = \beta_1 + (X_1^T X_1)^{-1} X_1^T X_2\beta_2 + (X_1^T X_1)^{-1} X_1^T \epsilon$$

On the second term, we see unless 1) X_1 and X_2 are orthogonal, or 2) $\beta_2 = 0$, $\hat{\beta}_1$ is biased.

2.2 Consistency

The unbiasedness gives us a metric of measuring how good our estimator is, from population perspective. In reality, as our sample size is finite, we need ask ourselves does our estimator converges to true value when sample size is sufficiently large. We know

$$\hat{\beta} = \beta + (X^T X)^{-1} X^T \epsilon$$

$$\begin{aligned}\hat{\beta} &= \beta + (X^T X)^{-1} X^T \epsilon \\ &= \beta + (\sum_{i=1}^N X_i X_i^T)^{-1} X^T \epsilon \\ &= \beta + (\sum_{i=1}^N \frac{1}{N} X_i X_i^T)^{-1} (\frac{X^T \epsilon}{N})\end{aligned}$$

To show $\hat{\beta}$ converges to β , we need to show two things:

(1) $\frac{1}{n} \sum_i X_i X_i^T$ converges to Q in probability when N is large. Also the inverse of Q exists.

(2) $\frac{X^T \epsilon}{N}$ converges to zero in probability when N is large.

For (1), we write

$$Q^{(n)} = \frac{1}{n} \sum_i X_i X_i^T$$

which is a $K \times K$ matrix. Its element Q_{kl} is

$$Q_{kl}^{(n)} = \frac{1}{n} \sum_i x_{ik} x_{il}$$

Based on Law of large numbers, $Q_{kl}^{(n)}$ converges to its expectation value $E[x_{ik} x_{il}]$. Let

$$E[x_{ik} x_{il}] = Q_{kl}$$

So each element of $Q^{(n)}$ converges Q_{kl} . Therefore, $Q^{(n)}$ converges to Q. To show the inverse of Q exists, we can show Q is a symmetric positive definite matrix. For any k dimensional vector v, we have

$$v^T Q v = E[v^T x_i x_i^T v]$$

Since $v^T x_i = \sum_k v_k x_{ik} = x_i^T v$ which is a scalar, so

$$v^T Q v = E[v^T x_i x_i^T v] = E[(v^T x_i)^2]$$

if for any i, x_{ik} are linear independent, in other words, no multicollinearity. Then $v^T x_i = \sum_k v_k x_{ik} \neq 0$ and $(v^T x_i)^2 > 0$. Therefore Q is SPD matrix and

its inverse exists.

$$\begin{aligned} & \frac{X^T \epsilon}{N} \\ &= \begin{pmatrix} \frac{1}{N} \sum_1^N x_{i1} \epsilon_i \\ \frac{1}{N} \sum_1^N x_{i2} \epsilon_i \\ \frac{1}{N} \sum_1^N x_{i3} \epsilon_i \\ \dots \\ \frac{1}{N} \sum_1^N x_{ik} \epsilon_i \end{pmatrix} \\ &= \frac{1}{N} \sum_{i=1}^N X_i \epsilon_i = \bar{w} \end{aligned}$$

Where \bar{w} is a $k \times 1$ vector. To see the asymptotical behavior of w , we consider its mean and asymptotical variance. The mean is

$$E[w_i] = E_X[E[w_i|x_i]] = E_X[X_i E[\epsilon|X_i]] = 0$$

$$Var[\bar{w}] = E[Var[\bar{w}|X]] + Var[E[\bar{w}|X]] = E[Var[\bar{w}|X]] + 0 = E[Var[\bar{w}|X]]$$

$$Var[\bar{w}|X] = E[\bar{w}\bar{w}^T|X] = \frac{1}{n} X^T E[\epsilon\epsilon^T] X \frac{1}{n} = \frac{\sigma^2}{n} \frac{X^T X}{n}$$

$$E[Var[\bar{w}|X]] = \frac{\sigma^2}{n} E\left(\frac{X^T X}{n}\right)$$

We have shown that $X^T X = \sum_i x_i x_i^T$, so

$$\frac{X^T X}{n} = \frac{1}{n} \sum_i x_i x_i^T$$

When $\frac{X^T X}{n}$ converges to Q ,

$$E[Var[\bar{w}|X]] = 0$$

So \bar{w} converges to $\mathbf{0}(k \times 1)$ vector. Then when N is sufficiently large, $\hat{\beta}$ converges to β . This is the proof of consistency.

There are certain conditions in which the estimators become inconsistent.

1) X is not full rank, or X has multicollinearity 2) $cov[X, \epsilon] \neq 0$

2.3 Efficiency

The least square estimator has the smallest variance, and this can be proved by Gauss-Markov theorem.

2.4 Multicollinearity

Suppose we have a regression model that contains two parameters

$$y = \beta_0 + X_1\beta_1 + X_2\beta_2$$

From above, we know variance of $\hat{\beta}$ is

$$Var(\hat{\beta}) = \frac{\sigma^2}{(X^T X)^{-1}}$$

When X only contains 2 variables, $X = (X_1, X_2)$

$$Var(\hat{\beta}_1) = \sigma^2 \frac{S_{22}}{S_{11}S_{22} - S_{12}^2} = \frac{1}{S_{11}(1 - \frac{S_{12}^2}{S_{11}S_{22}})} = \frac{1}{S_{11}(1 - r_{12}^2)}$$

$$Var(\hat{\beta}_2) = \sigma^2 \frac{S_{11}}{S_{11}S_{22} - S_{12}^2} = \frac{1}{S_{22}(1 - \frac{S_{12}^2}{S_{11}S_{22}})} = \frac{1}{S_{22}(1 - r_{12}^2)}$$

Where

$$S_{11} = \sum (x_{1i} - \hat{x}_1)^2$$

$$S_{22} = \sum (x_{2i} - \hat{x}_2)^2$$

$$S_{12} = \sum (x_{1i} - \hat{x}_1)(x_{2i} - \hat{x}_2)$$

$$r_{12} = \frac{S_{12}}{\sqrt{S_{11}S_{22}}}$$

r_{12} is the correlation coefficient. In extreme case, when X_1 and X_2 are perfectly correlated, the variance becomes infinite.

3 Model Testing

3.1 Lagrange Multiplier(LM) test

3.1.1 Constrained Maximum Likelihood Estimation

Consider a parametric model with likelihood function $L(\theta)$ for $\theta \in \Theta \subset \mathbb{R}^p$. We wish to test q restrictions:

$$H_0 : g(\theta) = 0,$$

where $g : \mathbb{R}^p \rightarrow \mathbb{R}^q$ is continuously differentiable and $q \leq p$.

The constrained maximum likelihood estimator (CMLE) solves:

$$\max_{\theta} \ell(\theta) \quad \text{s.t.} \quad g(\theta) = 0,$$

where $\ell(\theta) = \log L(\theta)$.

3.1.2 Lagrangian and First-Order Conditions

Form the Lagrangian:

$$\mathcal{L}(\theta, \lambda) = \ell(\theta) - \lambda' g(\theta),$$

where $\lambda \in \mathbb{R}^q$ is the vector of Lagrange multipliers.

The first-order conditions are:

$$\frac{\partial \mathcal{L}}{\partial \theta} = s(\theta) - G(\theta)' \lambda = 0, \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -g(\theta) = 0, \quad (7)$$

where:

- $s(\theta) = \frac{\partial \ell(\theta)}{\partial \theta}$ is the $p \times 1$ score vector.
- $G(\theta) = \frac{\partial g(\theta)}{\partial \theta}$ is the $q \times p$ Jacobian matrix of constraints.

Let $(\tilde{\theta}, \tilde{\lambda})$ denote the solution to (6) and (7).

3.1.3 Statistical Interpretation of Lagrange Multipliers

1. Sample vs. Population Lagrange Multipliers

- $\tilde{\lambda}$ is the **sample Lagrange multiplier**: a numerical value computed from the data.
- Under H_0 , the **population Lagrange multiplier** $\lambda_0 = 0$. Why? If the constraint holds in population, imposing it does not change the optimum, so its "shadow price" is zero.

Thus, we can write:

$$\tilde{\lambda} = \lambda_0 + \text{sampling error} = \text{sampling error} \quad \text{under } H_0.$$

2. Distribution of $\tilde{\lambda}$ Under H_0 From (6), at the constrained estimate:

$$s(\tilde{\theta}) = G(\tilde{\theta})' \tilde{\lambda}. \quad (1)$$

Under regularity conditions and H_0 , as $n \rightarrow \infty$:

- $\tilde{\theta} \xrightarrow{p} \theta_0$, where θ_0 is the true parameter.
- By the Central Limit Theorem for the score:

$$\frac{1}{\sqrt{n}} s(\theta_0) \xrightarrow{d} N(0, \mathcal{I}(\theta_0)),$$

where $\mathcal{I}(\theta_0)$ is the Fisher information matrix.

- Applying the delta method and using (1):

$$\frac{1}{\sqrt{n}} \tilde{\lambda} \xrightarrow{d} N\left(0, [G(\theta_0) \mathcal{I}(\theta_0)^{-1} G(\theta_0)']^{-1}\right).$$

Therefore, under H_0 :

$$\tilde{\lambda} \overset{a}{\approx} N\left(0, \frac{1}{n} [G(\theta_0) \mathcal{I}(\theta_0)^{-1} G(\theta_0)']^{-1}\right).$$

3.1.4 Constructing the LM Test Statistic

1. Quadratic Form in $\tilde{\lambda}$ Since $\tilde{\lambda} \sim N(0, \Sigma_\lambda)$ under H_0 , the quadratic form:

$$\tilde{\lambda}' \Sigma_\lambda^{-1} \tilde{\lambda} \sim \chi_q^2.$$

From the asymptotic variance above:

$$\Sigma_\lambda = \frac{1}{n} [G(\theta_0) \mathcal{I}(\theta_0)^{-1} G(\theta_0)']^{-1}.$$

Hence:

$$n \cdot \tilde{\lambda}' G(\theta_0) \mathcal{I}(\theta_0)^{-1} G(\theta_0)' \tilde{\lambda} \xrightarrow{d} \chi_q^2.$$

2. Connection to the Score Vector Using $s(\tilde{\theta}) = G(\tilde{\theta})' \tilde{\lambda}$ from (1):

$$n \cdot \tilde{\lambda}' G(\theta_0) \mathcal{I}(\theta_0)^{-1} G(\theta_0)' \tilde{\lambda} = n \cdot s(\tilde{\theta})' \mathcal{I}(\theta_0)^{-1} s(\tilde{\theta}) + o_p(1).$$

Replacing $\mathcal{I}(\theta_0)$ with a consistent estimator $\hat{\mathcal{I}}$ (e.g., observed information at $\tilde{\theta}$), we obtain the **Lagrange Multiplier (LM) test statistic**:

$$\boxed{\text{LM} = s(\tilde{\theta})' \hat{\mathcal{I}}^{-1} s(\tilde{\theta}) \xrightarrow{d} \chi_q^2.}$$

3.2 Example: Testing $\beta_2 = 0$ in a Linear Model

Consider the linear regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + u, \quad u \sim N(0, \sigma^2 I)$$

with n observations. We want to test $H_0 : \beta_2 = 0$.

Let $\theta = (\beta_0, \beta_1, \beta_2)' \in \mathbb{R}^3$.

3.2.1 Constraint, Maximum Likelihood and Lagrangian

The constraint is $g(\theta) = \beta_2 = 0$. Here $q = 1$, and:

$$G(\theta) = \frac{\partial g}{\partial \theta} = \begin{pmatrix} 0 & 0 & 1 \end{pmatrix}.$$

The log-likelihood is:

$$\ell(\theta) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i})^2.$$

The Lagrangian is:

$$\mathcal{L}(\theta, \lambda) = \ell(\theta) - \lambda \beta_2,$$

where λ is a scalar Lagrange multiplier.

3.2.2 First-Order Conditions

$$\frac{\partial \mathcal{L}}{\partial \beta_0} = \frac{1}{\sigma^2} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) = 0, \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial \beta_1} = \frac{1}{\sigma^2} \sum_{i=1}^n x_{1i} (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) = 0, \quad (9)$$

$$\frac{\partial \mathcal{L}}{\partial \beta_2} = \frac{1}{\sigma^2} \sum_{i=1}^n x_{2i} (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) - \lambda = 0, \quad (10)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = -\beta_2 = 0. \quad (11)$$

From (11), $\beta_2 = 0$ under the constraint.

3.2.3 Constrained Estimates

Let $\tilde{\theta} = (\tilde{\beta}_0, \tilde{\beta}_1, 0)'$ be the constrained MLE. Substituting $\beta_2 = 0$ into (8) and (9) gives the OLS estimates from the restricted model:

$$y = \beta_0 + \beta_1 x_1 + u.$$

These are simply:

$$\begin{pmatrix} \tilde{\beta}_0 \\ \tilde{\beta}_1 \end{pmatrix} = (W'W)^{-1}W'y,$$

where $W = [\mathbf{1}, x_1]$ is the $n \times 2$ matrix of intercept and x_1 .

3.2.4 Solving for $\tilde{\lambda}$

From (10), with $\beta_2 = 0$:

$$\tilde{\lambda} = \frac{1}{\sigma^2} \sum_{i=1}^n x_{2i} (y_i - \tilde{\beta}_0 - \tilde{\beta}_1 x_{1i}) = \frac{1}{\sigma^2} \sum_{i=1}^n x_{2i} \tilde{u}_i,$$

where \tilde{u}_i are residuals from the restricted regression.

In matrix form:

$$\tilde{\lambda} = \frac{1}{\sigma^2} x_2' \tilde{u},$$

where x_2 is the $n \times 1$ vector of x_{2i} and $\tilde{u} = y - W\tilde{\beta}$.

3.2.5 Distribution of $\tilde{\lambda}$ Under H_0

Under $H_0 : \beta_2 = 0$, the true model is $y = \beta_0 + \beta_1 x_1 + u$. Then:

- $\tilde{\beta}_0, \tilde{\beta}_1$ are consistent for β_0, β_1 .
- $\tilde{u}_i = u_i - (\tilde{\beta}_0 - \beta_0) - (\tilde{\beta}_1 - \beta_1)x_{1i}$.

Since $\mathbb{E}[x_{2i}u_i] = 0$ (by exogeneity) under H_0 :

$$\mathbb{E}[\tilde{\lambda}] = 0.$$

The variance can be computed. Write $M_W = I - W(W'W)^{-1}W'$, the annihilator matrix for W . Then $\tilde{u} = M_W y = M_W u$ under H_0 . Thus:

$$\tilde{\lambda} = \frac{1}{\sigma^2} x_2' M_W u.$$

Since $u \sim N(0, \sigma^2 I)$:

$$\tilde{\lambda} \sim N\left(0, \frac{1}{\sigma^2} x_2' M_W x_2\right).$$

Note $x_2' M_W x_2$ is the residual sum of squares from regressing x_2 on W .

Thus, under H_0 :

$$\frac{\tilde{\lambda}^2}{\frac{1}{\sigma^2} x_2' M_W x_2} = \sigma^2 \frac{\tilde{\lambda}^2}{x_2' M_W x_2} \sim \chi_1^2.$$

3.2.6 Score Vector at Constrained Estimate

The score vector is:

$$s(\theta) = \frac{1}{\sigma^2} \begin{pmatrix} \sum (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) \\ \sum x_{1i} (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) \\ \sum x_{2i} (y_i - \beta_0 - \beta_1 x_{1i} - \beta_2 x_{2i}) \end{pmatrix}.$$

At $\tilde{\theta} = (\tilde{\beta}_0, \tilde{\beta}_1, 0)'$:

$$s(\tilde{\theta}) = \frac{1}{\sigma^2} \begin{pmatrix} 0 \\ 0 \\ \sum x_{2i} \tilde{u}_i \end{pmatrix},$$

because the first two equations are satisfied by $\tilde{\beta}_0, \tilde{\beta}_1$ (from FOCs 8, 9).

Note that the third element equals $\tilde{\lambda}$, consistent with the general relation $s(\tilde{\theta}) = G(\tilde{\theta})' \tilde{\lambda}$.

3.2.7 Information Matrix

The information matrix for this normal linear model is:

$$\mathcal{I}(\theta) = \frac{1}{\sigma^2} \begin{pmatrix} n & \sum x_{1i} & \sum x_{2i} \\ \sum x_{1i} & \sum x_{1i}^2 & \sum x_{1i} x_{2i} \\ \sum x_{2i} & \sum x_{1i} x_{2i} & \sum x_{2i}^2 \end{pmatrix} = \frac{1}{\sigma^2} X' X,$$

where $X = [\mathbf{1}, x_1, x_2]$.

3.2.8 LM Statistic

The LM statistic is:

$$\text{LM} = s(\tilde{\theta})' \mathcal{I}(\tilde{\theta})^{-1} s(\tilde{\theta}).$$

Since only the third element of $s(\tilde{\theta})$ is non-zero, and $\mathcal{I}^{-1} = \sigma^2 (X' X)^{-1}$, we have:

$$\text{LM} = \frac{1}{\sigma^4} (0, 0, \sum x_{2i} \tilde{u}_i) \cdot \sigma^2 (X' X)^{-1} \cdot \begin{pmatrix} 0 \\ 0 \\ \sum x_{2i} \tilde{u}_i \end{pmatrix}.$$

Let $v = (0, 0, 1)'$ and note $\sum x_{2i}\tilde{u}_i = x_2'\tilde{u}$. Then:

$$\text{LM} = \frac{1}{\sigma^2} (x_2'\tilde{u}) \cdot v'(X'X)^{-1}v \cdot (x_2'\tilde{u}).$$

But $v'(X'X)^{-1}v = [(X'X)^{-1}]_{33}$, the (3,3) element of $(X'X)^{-1}$, which is $(x_2'M_W x_2)^{-1}$ (the inverse of the residual variance from regressing x_2 on W).

Thus:

$$\text{LM} = \frac{(x_2'\tilde{u})^2}{\sigma^2 \cdot (x_2'M_W x_2)} = \frac{\tilde{\lambda}^2}{\frac{1}{\sigma^2} x_2'M_W x_2},$$

exactly the quadratic form in $\tilde{\lambda}$ derived earlier.

3.2.9 Connection to Usual t-Test

The unconstrained OLS estimate of β_2 is:

$$\hat{\beta}_2 = \frac{x_2' M_W y}{x_2' M_W x_2}.$$

Under H_0 , $\hat{\beta}_2 \sim N\left(0, \frac{\sigma^2}{x_2' M_W x_2}\right)$.

The t-statistic is:

$$t = \frac{\hat{\beta}_2}{\text{SE}(\hat{\beta}_2)} = \frac{\hat{\beta}_2}{\sqrt{\hat{\sigma}^2 / (x_2' M_W x_2)}},$$

where $\hat{\sigma}^2$ is the error variance estimator from the unconstrained regression.

Note that:

$$x_2'\tilde{u} = x_2' M_W y = (x_2' M_W x_2)\hat{\beta}_2.$$

Thus:

$$\tilde{\lambda} = \frac{1}{\sigma^2} (x_2' M_W x_2)\hat{\beta}_2.$$

The LM statistic becomes:

$$\text{LM} = \frac{(x_2' M_W x_2)^2 \hat{\beta}_2^2 / \sigma^4}{\frac{1}{\sigma^2} x_2' M_W x_2} = \frac{(x_2' M_W x_2) \hat{\beta}_2^2}{\sigma^2}.$$

If we use the true σ^2 , then $\text{LM} = \hat{\beta}_2^2 / \text{Var}(\hat{\beta}_2) \sim \chi_1^2$, and:

$$\text{LM} = t^2.$$

In practice, we estimate σ^2 from the constrained regression as $\tilde{\sigma}^2 = \frac{1}{n} \sum \tilde{u}_i^2$, and the LM statistic becomes nR^2 from the auxiliary regression of \tilde{u} on x_1 and x_2 .

3.2.10 Summary

For testing $H_0 : \beta_2 = 0$:

- The Lagrange multiplier $\tilde{\lambda} = \frac{1}{\sigma^2} \sum x_{2i}\tilde{u}_i$ measures the marginal increase in log-likelihood from relaxing $\beta_2 = 0$.
- Under H_0 , $\tilde{\lambda} \sim N\left(0, \frac{1}{\sigma^2} x_2' M_W x_2\right)$.

- The LM statistic $LM = \frac{\tilde{\lambda}^2}{\text{Var}(\tilde{\lambda})} \sim \chi_1^2$.
- This equals $s(\tilde{\theta})' \mathcal{I}(\tilde{\theta})^{-1} s(\tilde{\theta})$, connecting constrained optimization to hypothesis testing.
- When σ^2 is known, $LM = t^2$, showing equivalence to the usual t-test.

This example illustrates the general principle: LM tests check whether the score (gradient of the likelihood) at the constrained estimate is statistically zero, which is equivalent to testing whether the Lagrange multipliers (shadow prices of constraints) are zero.

3.2.11 Example: Testing constraint $R\beta = r$ in Linear Model

Consider $y = X\beta + u$, $u \sim N(0, \sigma^2 I)$, with constraint $R\beta = r$.

1. Constrained MLE The Lagrangian is:

$$\mathcal{L}(\beta, \lambda) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y - X\beta)'(y - X\beta) - \lambda'(R\beta - r).$$

FOCs:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta} &= \frac{1}{\sigma^2} X'(y - X\beta) - R'\lambda = 0, \\ \frac{\partial \mathcal{L}}{\partial \lambda} &= -(R\beta - r) = 0. \end{aligned}$$

The constrained estimator is $\tilde{\beta}$ with $R\tilde{\beta} = r$, and:

$$\tilde{\lambda} = \frac{1}{\sigma^2} [R(X'X)^{-1}R']^{-1} (R\hat{\beta}_{OLS} - r),$$

where $\hat{\beta}_{OLS} = (X'X)^{-1}X'y$.

2. Distribution Under H_0 Under $H_0 : R\beta = r$, and assuming σ^2 known:

$$\tilde{\lambda} \sim N\left(0, \frac{1}{\sigma^2} [R(X'X)^{-1}R']^{-1}\right).$$

Thus:

$$\sigma^2 \tilde{\lambda}' R(X'X)^{-1} R' \tilde{\lambda} \sim \chi_q^2.$$

3. Equivalence to Score Form The score at $\tilde{\beta}$ is:

$$s(\tilde{\beta}) = \frac{1}{\sigma^2} X'(y - X\tilde{\beta}) = R'\tilde{\lambda}.$$

The information matrix is $\mathcal{I}(\beta) = \frac{1}{\sigma^2} X'X$.

Then:

$$LM = s(\tilde{\beta})' \mathcal{I}(\tilde{\beta})^{-1} s(\tilde{\beta}) = \sigma^2 \tilde{\lambda}' R(X'X)^{-1} R' \tilde{\lambda},$$

identical to the quadratic form in $\tilde{\lambda}$.

3.2.12 Key Insights

1. Why Test $\lambda = 0$?
 - λ measures the marginal increase in log-likelihood from relaxing the constraint.
 - If H_0 is true, relaxing the constraint should not improve the likelihood significantly $\Rightarrow \lambda$ should be near zero.
 - A large $|\tilde{\lambda}|$ suggests the constraint is costly, evidence against H_0 .
2. Advantages of the Score Form While we *could* directly test $\tilde{\lambda} = 0$, the score form $s(\tilde{\theta})$ is preferred because:
 - (a) **Invariance:** $LM = s(\tilde{\theta})' \hat{\mathcal{I}}^{-1} s(\tilde{\theta})$ is invariant to reparameterization of constraints, while $\tilde{\lambda}$ is not.
 - (b) **Computational simplicity:** Often $s(\tilde{\theta})$ is easier to compute (e.g., via auxiliary regressions like White's test).
 - (c) **Unified theory:** The score vector appears in other contexts (Cramér-Rao bound, MLE asymptotics).
3. The Big Picture The LM test elegantly connects:
 - **Constrained optimization** (Lagrange multipliers)
 - **Asymptotic distribution theory** (CLT for scores)
 - **Quadratic forms** (chi-square distributions)

It tests whether the gradient of the likelihood at the constrained estimate is statistically zero—which is exactly what we expect if the constraints are true in population.

Suppose we have two models, one is restricted, the other is unrestricted:
 Restricted (R): $y = X_1\beta_1 + \epsilon$
 Unrestricted (U): $y = X_1\beta_1 + X_2\beta_2 + \epsilon$
 Given the unrestricted model, the likelihood function is

$$L(\beta_1, \beta_2, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (y - X_1\beta_1 - X_2\beta_2)' (y - X_1\beta_1 - X_2\beta_2)\right)$$

$$S_2 = \frac{\partial L}{\partial \beta_2} = \frac{1}{\sigma^2} X_2^T (y - X_1\beta_1 - X_2\beta_2)$$

When $\beta_2 = 0$, define

$$M_1 = I - X_1(X_1^T X_1)^{-1} X_1^T$$

and $M_1 X_1 = 0$.

$$S_2 = \frac{1}{\sigma^2} X_2^T (y - X_1\hat{\beta}_1) = \frac{1}{\sigma^2} X_2^T M_1 y = \frac{1}{\sigma^2} X_2^T M_1 (X_1\beta_1 + \epsilon) = \frac{1}{\sigma^2} X_2^T M_1 \epsilon$$

The last equal sign uses the fact $M_1 X_1 = 0$.

$$\begin{aligned} \text{Var}(X_2^T M_1 \epsilon) &= \text{Var}(X_2^T M_1 \epsilon) \\ &= X_2^T M_1 \text{Var}(\epsilon) (X_2^T M_1)^T \\ &= X_2^T M_1 M_1^T X_2 \text{Var}(\epsilon) \\ &= \sigma^2 X_2^T M_1 X_2 \end{aligned}$$

Define

$$V = X_2^T M_1 X_2$$

Then

$$\text{Var}(X_2^T M_1 \epsilon) = \sigma^2 V$$

So $X_2^T M_1 \epsilon$ follows normal distribution with mean 0 and variance $\sigma^2 X_2^T M_1 X_2$.
Define

$$Z = \frac{X_2^T M_1 \epsilon}{\sqrt{\sigma^2 X_2^T M_1 X_2}} = \frac{S_2}{\sqrt{\sigma^2 V}}$$

then Z follows standard normal distribution. The **Lagrange Multiplier (LM) test** is defined

$$LM = Z^2 = \frac{(X_2^T M_1 \epsilon)^2}{\sigma^2 X_2^T M_1 X_2} = \frac{S_2^2}{\sigma^2 V}$$

which follows χ^2 distribution with degree of freedom 1.

F test

We define

$$\begin{aligned} SSE_U &= \|y - X_1 \hat{\beta}_1 - X_2 \hat{\beta}_2\|^2 \\ SSE_R &= \|y - X_1 \hat{\beta}_1\|^2 \end{aligned}$$

F test is defined as

$$F = \frac{\frac{\text{Extra explained variation}}{\text{Degree of Freedom}}}{\frac{\text{Remaining unexplained variation}}{\text{Degree of Freedom}}} = \frac{SSE_R - SSE_U}{\frac{SSE_U}{n-1}}$$

Let $X = (X_1, X_2)$, and we define two projection matrices

$$\begin{aligned} P_U &= X(X^T X)^{-1} X^T \\ P_R &= X_1(X_1^T X_1)^{-1} X_1^T \end{aligned}$$

$$SSR_R - SSR_U = y^T P_R y - y^T P_U y = y^T (P_R - P_U) y$$

recall

$$\hat{\beta}_2 = (X_2^T M_1 X_2)^{-1} X_2^T M_1 y$$

The corresponding projection matrix is

$$M_1 X_2 (X_2^T M_1 X_2)^{-1} X_2^T M_1$$

The $SSR_R - SSR_U$ is the additional variance explained by X_2 after removing the linear space of X_1 on X_2 . This means the projection matrix corresponding to β_2 is $P_U - P_R$. So we can get $P_U - P_R$ using the interpretation of projection matrix instead solving for the projection matrix itself.

$$P_U - P_R = M_1 X_2 (X_2^T M_1 X_2)^{-1} X_2^T M_1$$

The extra explained sum of squares by the unrestricted model is

$$SSR_R - SSR_U = y^T (P_R - P_U) y = (X_2^T M_1 y)^T (X_2^T M_1 X_2)^{-1} X_2^T M_1 y$$

with 1 degree of freedom as X_2 only contains 1 parameter.

$$F = \frac{SSR_R - SSR_U}{\frac{SSR_U}{n-k}} = \frac{(X_2 M_1 y)^T (X_2^T M_1 y)}{\hat{\sigma}^2 (X_2^T M_1 X_2)} = LM$$

We see that F test and LM test are equivalent.

Wald Test

Recall the estimator for β_2 in Eq.4,

$$\hat{\beta}_2 = (X_2^T M_1 X_2)^{-1} (X_2^T M_1 y)$$

Substitue

$$y = X_1 \beta_1 + X_2 \beta_2 + \epsilon$$

we get

$$\hat{\beta}_2 = (X_2^T M_1 X_2)^{-1} (X_2^T M_1 \epsilon)$$

Since $\epsilon \sim N(0, \sigma^2 I)$, we obtain

$$\beta_2 \sim N(0, \sigma^2 (X_2^T M_1 X_2)^{-1})$$

Thus, scaling by $1/\sigma^2$, we arrive at

$$\frac{1}{\sigma^2} X_2^T M_1 \epsilon \sim N(0, X_2^T M_1 X_2)$$

Construct Wald test W

$$W = \frac{\hat{\beta}_2}{\sqrt{\hat{Var}(\hat{\beta}_2)}}$$

W follows t distribution. We now show W test is equivalent to LM test. Consider W^2

$$\begin{aligned} W^2 &= \hat{\beta}^T (Var(\hat{\beta}))^{-1} \hat{\beta} = \frac{1}{\sigma^2} \hat{\beta}^T V \hat{\beta} = \frac{1}{\sigma^2} (V^{-1} S_2)^T V V^{-1} S_2 = \frac{1}{\sigma^2} S_2^T V^{-1} V V^{-1} S_2 \\ &= \frac{1}{\sigma^2} S_2^T V^{-1} S_2 \\ &= LM \end{aligned}$$

4 Generalized Regression Model

4.1 Generalized Least Squares (GLS)

4.1.1 Problem Setting and Motivation

Recall the linear regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where:

- \mathbf{y} is $n \times 1$ vector of outcomes
- \mathbf{X} is $n \times k$ matrix of regressors (full column rank)
- $\boldsymbol{\beta}$ is $k \times 1$ parameter vector
- $\boldsymbol{\epsilon}$ is $n \times 1$ error vector

The classical OLS assumptions include spherical errors:

$$\mathbb{E}[\boldsymbol{\epsilon} \mid \mathbf{X}] = \mathbf{0}, \quad \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}' \mid \mathbf{X}] = \sigma^2 \mathbf{I}_n.$$

When the second assumption fails (heteroskedasticity and/or autocorrelation), OLS remains unbiased but is no longer efficient. The Generalized Least Squares (GLS) framework provides the BLUE under a known non-spherical covariance structure.

4.1.2 Assumptions

Assume:

1. $\mathbb{E}[\boldsymbol{\epsilon} \mid \mathbf{X}] = \mathbf{0}$
2. $\mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}' \mid \mathbf{X}] = \sigma^2 \boldsymbol{\Omega}$, where $\boldsymbol{\Omega}$ is a known $n \times n$ symmetric positive definite matrix.
3. \mathbf{X} has full column rank k , and $\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X}$ is invertible.

4.1.3 Derivation via Transformation

Since $\boldsymbol{\Omega}$ is positive definite, there exists a nonsingular matrix \mathbf{P} such that:

$$\mathbf{P}\boldsymbol{\Omega}\mathbf{P}' = \mathbf{I}_n.$$

One common choice is $\mathbf{P} = \boldsymbol{\Omega}^{-1/2}$ (the inverse of the Cholesky factor).

Transform the model by premultiplying by \mathbf{P} :

$$\mathbf{P}\mathbf{y} = \mathbf{P}\mathbf{X}\boldsymbol{\beta} + \mathbf{P}\mathbf{u}.$$

Define $\mathbf{y}^* = \mathbf{P}\mathbf{y}$, $\mathbf{X}^* = \mathbf{P}\mathbf{X}$, $\mathbf{u}^* = \mathbf{P}\mathbf{u}$. Then:

$$\mathbb{E}[\mathbf{u}^* \mid \mathbf{X}] = \mathbf{P}\mathbb{E}[\mathbf{u} \mid \mathbf{X}] = \mathbf{0},$$

$$\mathbb{E}[\mathbf{u}^*\mathbf{u}^{*'} \mid \mathbf{X}] = \mathbf{P}(\sigma^2 \boldsymbol{\Omega})\mathbf{P}' = \sigma^2 \mathbf{I}_n.$$

Applying OLS to the transformed model gives the GLS estimator:

$$\hat{\boldsymbol{\beta}}_{\text{GLS}} = (\mathbf{X}^{*'}\mathbf{X}^*)^{-1}\mathbf{X}^{*'}\mathbf{y}^* = (\mathbf{X}'\mathbf{P}'\mathbf{P}\mathbf{X})^{-1}\mathbf{X}'\mathbf{P}'\mathbf{P}\mathbf{y}.$$

Since $\mathbf{P}'\mathbf{P} = \boldsymbol{\Omega}^{-1}$, we obtain the canonical form:

$$\boxed{\hat{\boldsymbol{\beta}}_{\text{GLS}} = (\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Omega}^{-1}\mathbf{y}}.$$

4.1.4 Properties

Under assumptions 1-3, $\hat{\beta}_{\text{GLS}}$ is the Best Linear Unbiased Estimator (BLUE) of β .

Let $\tilde{\beta} = \mathbf{A}\mathbf{y}$ be any linear unbiased estimator. Write $\mathbf{A} = (\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{\Omega}^{-1} + \mathbf{D}$. Unbiasedness requires $\mathbf{D}\mathbf{X} = \mathbf{0}$. Then:

$$\text{Var}(\tilde{\beta} \mid \mathbf{X}) = \sigma^2 [(\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X})^{-1} + \mathbf{D}\mathbf{\Omega}\mathbf{D}'] .$$

Since $\mathbf{D}\mathbf{\Omega}\mathbf{D}'$ is positive semidefinite, $\hat{\beta}_{\text{GLS}}$ has minimal variance.

The asymptotic distribution (under appropriate regularity conditions) is:

$$\sqrt{n}(\hat{\beta}_{\text{GLS}} - \beta) \xrightarrow{d} N(\mathbf{0}, \sigma^2 \mathbf{Q}^{-1}) ,$$

where $\mathbf{Q} = \text{plim } \frac{1}{n}\mathbf{X}'\mathbf{\Omega}^{-1}\mathbf{X}$.

4.2 Feasible Generalized Least Squares (FGLS)

4.2.1 The Problem of Unknown $\mathbf{\Omega}$

In practice, $\mathbf{\Omega}$ is unknown. Feasible GLS (FGLS) replaces $\mathbf{\Omega}$ with a consistent estimator $\hat{\mathbf{\Omega}}$.

4.2.2 Two-Step Procedure

1. **First step:** Obtain initial OLS estimator $\hat{\beta}_{\text{OLS}}$ and residuals $\hat{\mathbf{u}} = \mathbf{y} - \mathbf{X}\hat{\beta}_{\text{OLS}}$.
2. **Estimate $\mathbf{\Omega}$:** Model $\mathbf{\Omega}$ as a function of parameters $\boldsymbol{\theta}$, e.g.:
 - For heteroskedasticity: $\Omega_{ii} = h(\mathbf{z}_i'\boldsymbol{\theta})$, where \mathbf{z}_i may include \mathbf{x}_i .
 - For AR(1) errors: $\Omega_{ij} = \rho^{|i-j|}$.

Estimate $\boldsymbol{\theta}$ from residuals (e.g., regress $\log(\hat{u}_i^2)$ on \mathbf{z}_i for multiplicative heteroskedasticity).

3. **Second step:** Compute FGLS estimator:

$$\boxed{\hat{\beta}_{\text{FGLS}} = (\mathbf{X}'\hat{\mathbf{\Omega}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{\Omega}}^{-1}\mathbf{y}} .$$

4.2.3 Large-Sample Properties

Under regularity conditions (consistency of $\hat{\mathbf{\Omega}}$ and finite moments):

$$\sqrt{n}(\hat{\beta}_{\text{FGLS}} - \beta) \xrightarrow{d} N(\mathbf{0}, \sigma^2 \mathbf{Q}^{-1}) ,$$

the same asymptotic distribution as GLS. Thus, FGLS is **asymptotically efficient**.

Important: In finite samples, FGLS may not outperform OLS if $\mathbf{\Omega}$ is poorly estimated. Iterated FGLS (updating residuals and $\hat{\mathbf{\Omega}}$ until convergence) is sometimes used.

4.3 Weighted Least Squares (WLS)

4.3.1 Special Case of GLS

WLS arises when $\mathbf{\Omega}$ is diagonal (heteroskedasticity with no autocorrelation):

$$\mathbf{\Omega} = \text{diag}(\omega_1, \omega_2, \dots, \omega_n), \quad \omega_i > 0.$$

Then $\mathbf{\Omega}^{-1} = \text{diag}(\omega_1^{-1}, \dots, \omega_n^{-1})$.

4.3.2 Derivation

The GLS estimator simplifies to:

$$\hat{\boldsymbol{\beta}}_{\text{WLS}} = \left(\sum_{i=1}^n \omega_i^{-1} \mathbf{x}_i \mathbf{x}_i' \right)^{-1} \left(\sum_{i=1}^n \omega_i^{-1} \mathbf{x}_i y_i \right).$$

Equivalently, $\hat{\boldsymbol{\beta}}_{\text{WLS}}$ minimizes the weighted sum of squares:

$$\sum_{i=1}^n \omega_i^{-1} (y_i - \mathbf{x}_i' \boldsymbol{\beta})^2.$$

If we define weights $w_i = \omega_i^{-1}$, the estimator is called Weighted Least Squares.

4.3.3 Common Weight Specifications

- **Known weights:** From survey data where ω_i reflects sampling probabilities.
- **Variance proportional to a known function:** Suppose $\text{Var}(u_i | \mathbf{x}_i) = \sigma^2 v(\mathbf{x}_i)$. Then $\omega_i = v(\mathbf{x}_i)$.
- **Estimated weights:** A two-step FGLS procedure where:
 1. Regress y on \mathbf{X} by OLS, get residuals \hat{u}_i .
 2. Regress \hat{u}_i^2 (or $\log \hat{u}_i^2$) on functions of \mathbf{x}_i to estimate $v(\mathbf{x}_i)$.
 3. Use $\hat{\omega}_i = \hat{v}(\mathbf{x}_i)$ in WLS.

4.4 Relationships and Comparison

Method	$\mathbf{\Omega}$ known?	Structure	Efficiency
OLS	$\mathbf{\Omega} = \mathbf{I}$	Spherical errors	Inefficient if $\mathbf{\Omega} \neq \mathbf{I}$
WLS	Diagonal known/estimated	Heteroskedastic only	Efficient within diagonal class
GLS	Known	General	BLUE
FGLS	Estimated	General	Asymptotically efficient

Table 1: Comparison of Least Squares Methods

4.4.1 Key Theorems

Equivalence of GLS and OLS on transformed data

GLS is numerically equivalent to OLS on data transformed by $\mathbf{P} = \mathbf{\Omega}^{-1/2}$.

Invariance to choice of \mathbf{P}

If \mathbf{P}_1 and \mathbf{P}_2 satisfy $\mathbf{P}_i \mathbf{\Omega} \mathbf{P}_i' = \mathbf{I}$, the GLS estimator is unchanged.

Asymptotic equivalence of FGLS and GLS

If $\hat{\mathbf{\Omega}} \xrightarrow{p} \mathbf{\Omega}$ and certain moment conditions hold, $\hat{\beta}_{\text{FGLS}}$ has the same asymptotic distribution as $\hat{\beta}_{\text{GLS}}$.

4.4.2 Empirical Considerations

4.5 Testing for Heteroskedasticity/Autocorrelation

Before applying FGLS, test OLS residuals:

- Breusch-Pagan / White test for heteroskedasticity
- Durbin-Watson / Breusch-Godfrey test for autocorrelation

4.6 Robust Inference

Even with FGLS, it is common to report **heteroskedasticity-robust standard errors** (White, 1980) because:

1. The $\mathbf{\Omega}$ model may be misspecified.
2. Robust standard errors provide valid inference under weaker conditions.

The robust asymptotic variance estimator for FGLS is:

$$\widehat{\text{Var}}_{\text{robust}}(\hat{\beta}_{\text{FGLS}}) = (\mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1} \left(\sum_{i=1}^n \hat{w}_i^2 \hat{u}_i^2 \mathbf{x}_i \mathbf{x}_i' \right) (\mathbf{X}' \hat{\mathbf{\Omega}}^{-1} \mathbf{X})^{-1},$$

where \hat{w}_i are the FGLS weights and \hat{u}_i the FGLS residuals.

Summary

- **GLS** is the theoretical benchmark with known $\mathbf{\Omega}$.
- **FGLS** is the feasible version using estimated $\hat{\mathbf{\Omega}}$; asymptotically equivalent to GLS.
- **WLS** is a special case for diagonal $\mathbf{\Omega}$ (pure heteroskedasticity).
- In practice, FGLS is commonly implemented with:
 1. Model specification for the error structure
 2. Two-step estimation
 3. Robust inference to guard against misspecification

5 Panel Data Model

We can view panel data as a "two dimensional" data set in which the sample does not only come from different individuals, but also same individual across different time point. We can write the regression model as

$$y_{it} = \alpha_{it} + \sum_k x_{itk} \beta_{itk} + u_{it}$$

where $1 < i < N$, $1 < t < T$, and $1 < k < K$. The equation has total sample size of NT with total number of parameter $NT(K+1)$, therefore it is not estimable. So we will make the following few assumptions

	$\alpha_{it} = \alpha_{is}$	$\alpha_{it} = \alpha_{jt}$	$\beta_{it} = \beta_{is}$	$\beta_{itk} = \beta_{jtk}$
Pooled	yes	yes	yes	yes
Fixed Effect	yes	no	yes	yes
Unrestricted	yes	no	yes	no

5.1 The unrestriced model

$$y_{it} = \alpha_i + \sum_k x_{itk} \beta_{ik} + u_{it}$$

The above equation can be written in matrix form:
for $i = 1$,

$$\begin{pmatrix} y_{11} \\ y_{12} \\ \dots \\ y_{1T} \end{pmatrix} = \begin{pmatrix} 1 & x_{111} & x_{112} & \dots & x_{11K} \\ 1 & x_{121} & x_{122} & \dots & x_{12K} \\ 1 & \dots & \dots & \dots & \dots \\ 1 & x_{1T1} & x_{1T2} & \dots & x_{1TK} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \beta_{11} \\ \beta_{12} \\ \dots \\ \beta_{1K} \end{pmatrix} + \begin{pmatrix} u_{11} \\ u_{12} \\ \dots \\ u_{1T} \end{pmatrix}$$

which we can also write as for $i = 2$,

$$\begin{pmatrix} y_{21} \\ y_{22} \\ \dots \\ y_{2T} \end{pmatrix} = \begin{pmatrix} 1 & x_{211} & x_{212} & \dots & x_{21K} \\ 1 & x_{221} & x_{222} & \dots & x_{22K} \\ 1 & \dots & \dots & \dots & \dots \\ 1 & x_{2T1} & x_{2T2} & \dots & x_{2TK} \end{pmatrix} \begin{pmatrix} \alpha_2 \\ \beta_{21} \\ \beta_{22} \\ \dots \\ \beta_{2K} \end{pmatrix} + \begin{pmatrix} u_{21} \\ u_{22} \\ \dots \\ u_{2T} \end{pmatrix}$$

So for each i , we can write

$$Y_i = 1_T \alpha_i + X_i \beta_i + U_i$$

where $Y_i = (y_{i1}, y_{i2}, \dots, y_{iT})^T$, 1_T is a one vector of length T , X_i is $K \times T$ matrix, $\beta_i = (\beta_{i1}, \beta_{i2}, \dots, \beta_{iK})^T$, and $U_i = (u_{i1}, u_{i2}, \dots, u_{iT})^T$.

If we consolidate equation set for all the value of i

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \dots \\ Y_N \end{pmatrix} = \begin{pmatrix} 1_T & 0 & 0 & \dots & 0 \\ 0 & 1_T & 0 & \dots & 0 \\ 0 & 0 & 1_T & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_{NT} \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_N \end{pmatrix} + \begin{pmatrix} X_1 & 0 & 0 & \dots & 0 \\ 0 & X_2 & 0 & \dots & 0 \\ 0 & 0 & X_3 & \dots & 0 \\ 0 & 0 & 0 & \dots & X_N \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_N \end{pmatrix} + \begin{pmatrix} U_1 \\ U_2 \\ \dots \\ U_N \end{pmatrix}$$

To solve for β_i , we can use the strategy of partition regression. β_i is the solution of the the regression

$$MY_i = MX_i\beta + U$$

where

$$\begin{aligned} M &= I - \frac{1}{T}1_T1_T' \\ MY_i &= y_{it} - \bar{y}_i. \\ MX_i &= x_{it} - \bar{x}_i. \end{aligned}$$

Here to avoid duplicate notation, we denote the transposed matrix using $'$. The estimate of β is

$$\hat{\beta}_i = ((MX_i)^T(MX_i))^{-1}((MX_i)^T(MY_i)) = W_{xx,i}^{-1}W_{xy,i}$$

where

$$\begin{aligned} W_{xy,i} &= \sum_t^T (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i) \\ W_{xx,i} &= \sum_t^T (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)^T \end{aligned}$$

5.2 The pooled model

$$y_{it} = \alpha + \sum_k x_{itk}\beta_k + u_{it}$$

The above equation can be written in matrix form:
for $i = 1$,

$$\begin{pmatrix} y_{11} \\ y_{12} \\ \dots \\ y_{1T} \\ y_{21} \\ y_{22} \\ \dots \\ y_{2T} \\ \dots \\ y_{NT} \end{pmatrix} = \begin{pmatrix} 1 & x_{111} & x_{112} & \dots & x_{11K} \\ 1 & x_{121} & x_{122} & \dots & x_{12K} \\ 1 & \dots & \dots & \dots & \dots \\ 1 & x_{1T1} & x_{1T2} & \dots & x_{1TK} \\ 1 & x_{211} & x_{212} & \dots & x_{21K} \\ 1 & x_{221} & x_{222} & \dots & x_{22K} \\ 1 & \dots & \dots & \dots & \dots \\ 1 & x_{2T1} & x_{2T2} & \dots & x_{2TK} \\ 1 & \dots & \dots & \dots & \dots \\ 1 & x_{NT1} & x_{NT2} & \dots & x_{NTK} \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \\ \dots \\ \beta_K \end{pmatrix} + \begin{pmatrix} u_{11} \\ u_{12} \\ \dots \\ u_{1T} \\ u_{21} \\ u_{22} \\ \dots \\ u_{2T} \\ \dots \\ u_{NT} \end{pmatrix}$$

Similarly, using the solution of β from Eq.5, the estimated β can be written as

$$M = I - \frac{1}{NT}1_{NT}1_{NT}'$$

$$\begin{aligned}
\hat{\beta} &= ((MX)^T(MX))^{-1}((MX)^T(MY_i)) \\
&= \left(\sum_i^N \sum_t^T (x_{it} - \bar{x}_{..})(x_{it} - \bar{x}_{..})^T \right)^{-1} \sum_i^N \sum_t^T (x_{it} - \bar{x}_{..})(y_{it} - \bar{y}_{..}) \\
&= T_{xx}^{-1} T_{xy}
\end{aligned}$$

where

$$\begin{aligned}
T_{xy} &= \sum_i^N \sum_t^T (x_{it} - \bar{x}_{..})(y_{it} - \bar{y}_{..}) \\
T_{xx} &= \sum_i^N \sum_t^T (x_{it} - \bar{x}_{..})(x_{it} - \bar{x}_{..})^T \\
T_{yy} &= \sum_i^N \sum_t^T (y_{it} - \bar{y}_{..})^2
\end{aligned}$$

We call T_{xx} , T_{yy} and T_{xy} the total sum square of x, total sum square of y, and total sum of cross product. The sum of square error for the pooled model is

$$SSE_{pooled} = T_{yy} - T_{xy}' T_{xx}^{-1} T_{xy}$$

with N-1-K degrees of freedom.

5.3 The fixed effect model

5.3.1 Model formulation and estimator

$$y_{it} = \alpha_i + \sum_k x_{itk} \beta_k + u_{it}$$

The above equation can be written in matrix form:

$$\begin{pmatrix} y_{11} \\ y_{12} \\ \dots \\ y_{1T} \\ y_{21} \\ y_{22} \\ \dots \\ y_{2T} \\ \dots \\ y_{NT} \end{pmatrix} = \begin{pmatrix} 1_T & 0 & 0 & \dots & 0 \\ 0 & 1_T & 0 & \dots & 0 \\ 0 & 0 & 1_T & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_T \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_N \end{pmatrix} \\
+ \begin{pmatrix} x_{111} & x_{112} & \dots & x_{11K} \\ x_{121} & x_{122} & \dots & x_{12K} \\ \dots & \dots & \dots & \dots \\ x_{1T1} & x_{1T2} & \dots & x_{1TK} \\ x_{211} & x_{212} & \dots & x_{21K} \\ x_{221} & x_{222} & \dots & x_{22K} \\ \dots & \dots & \dots & \dots \\ x_{2T1} & x_{2T2} & \dots & x_{2TK} \\ \dots & \dots & \dots & \dots \\ x_{NT1} & x_{NT2} & \dots & x_{NTK} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_K \end{pmatrix} + \begin{pmatrix} u_{11} \\ u_{12} \\ \dots \\ u_{1T} \\ u_{21} \\ u_{22} \\ \dots \\ u_{2T} \\ \dots \\ u_{NT} \end{pmatrix}$$

Where 1_T is a $1 \times T$ vector. Similarly, using the solution of β from Eq.5, the M matrix is

$$M = I - \frac{1}{T} \begin{pmatrix} 1_{T \times T} & 0 & 0 & \dots & 0 \\ 0 & 1_{T \times T} & 0 & \dots & 0 \\ 0 & 0 & 1_{T \times T} & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_{T \times T} \end{pmatrix}$$

Define

$$\tilde{x}_{itk} = (MX)_{itk} = x_{itk} - \bar{x}_{i.k}$$

which is an NT x K matrix,

$$\tilde{y}_{it} = (MY)_{it} = y_{it} - \bar{y}_i.$$

which is an NT x 1 vector. The estimated β can be written as

$$\hat{\beta} = ((MX)'(MX))^{-1}((MX)^T(MY_i))$$

and

$$\begin{aligned} ((MX)'(MX))_{kl} &= \sum_i \sum_t \tilde{x}_{itk} \tilde{x}_{itl} \\ ((MX)'(MY))_k &= \sum_i \sum_t x_{itk} \tilde{y}_{it} \end{aligned}$$

Let $\tilde{X}_{it} = (x_{it1} - \bar{x}_{i.1}, x_{it2} - \bar{x}_{i.2}, \dots, x_{itk} - \bar{x}_{i.k})$, $\tilde{y}_{it} = y_{it} - \bar{y}_i$. Then

$$\begin{aligned} ((MX)'(MX)) &= \sum_i \sum_t \tilde{X}_{it} \tilde{X}_{it}' \\ ((MX)'(MY)) &= \sum_i \sum_t \tilde{X}_{it} \tilde{y}_{it} \end{aligned}$$

$$\begin{aligned} \hat{\beta} &= [\sum_i \sum_t \tilde{X}_{it} \tilde{X}_{it}']^{-1} [\sum_i \sum_t \tilde{X}_{it} \tilde{y}_{it}] \\ &= W_{xx}^{-1} W_{xy} \end{aligned}$$

where

$$\begin{aligned} W_{xy} &= \sum_i \sum_t \tilde{X}_{it} \tilde{y}_{it} \\ W_{xx} &= \sum_i \sum_t \tilde{X}_{it} \tilde{X}_{it}' \\ W_{yy} &= \sum_i \sum_t \tilde{y}_{it}^2 \end{aligned}$$

We call W_{xx} , W_{yy} and W_{xy} the within-groups sum square of x, the within-groups sum square of y, and within-groups of cross product. The name within-groups means they utilized the variation within group i. The sum of square error is

$$SSE_{fix} = W_{yy} - W_{xy}' W_{xx}^{-1} W_{xy}$$

with NT - N - K degrees of freedom.

If we remind ourselves of the concept of partial regression, we let

$$Y^* = MYX^* = MXU^* = MU$$

then the fix effect model can also be written as

$$Y^* = X^* \beta + U^*$$

or

$$y_{it} - \bar{y}_i. = (X_{it} - \bar{X}_{i.})\beta + (u_{it} - \bar{u}_{i.})$$

5.3.2 Properties of Fixed effect estimator

1) Unbiasedness

$$\begin{aligned}
\hat{\beta} &= [\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T (\tilde{X}_{it})(\tilde{y}_{it})] \\
&= [\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T (\tilde{X}_{it})(\tilde{X}'_{it}\beta + \tilde{u}_{it})] \\
&= [\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T (\tilde{X}_{it}\tilde{X}'_{it}\beta + \tilde{X}_{it}\tilde{u}_{it})] \\
&= [\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T \tilde{X}_{it}\tilde{X}'_{it}]\beta + [\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T \tilde{X}_{it}\tilde{u}_{it}]
\end{aligned}$$

So

$$E[\hat{\beta}] = \beta + E[\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T \tilde{X}_{it}\tilde{u}_{it}]$$

As u_{it} is uncorrelated with x_{it} , it is easy to prove the 2nd term above is zero. Therefore, the within group estimator is unbiased.

2) Consistency

To show consistency, we need to prove

$$[\sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\sum_i^N \sum_t^T \tilde{X}_{it}\tilde{u}_{it}] = [\frac{1}{NT} \sum_i^N \sum_t^T \tilde{X}_{it} X'_{it}]^{-1} [\frac{1}{NT} \sum_i^N \sum_t^T \tilde{X}_{it}\tilde{u}_{it}] = 0$$

Similar to our previous consistency proof, the above converges to zero when 1) either T or N goes to infinity, 2) \tilde{x}_{it} and \tilde{u}_{it} are not correlated, 3) x_{it} and u_{it} has finite second moment.

5.3.3 Connections between fixed effect model estimator and pooled model estimator when T is large

To see the connection between fixed effect estimator and pooled model estimator, we first introduce between-groups estimator, which allows us to easily understand the relationship among different estimators.

We define the between-groups sum of square is

$$\begin{aligned}
B_{xx} &= \sum_i^N T(\bar{x}_{i.} - \bar{x}_{...})(\bar{x}_{i.} - \bar{x}_{...})' \\
B_{yy} &= \sum_i^N T(\bar{y}_{i.} - \bar{y}_{...})(\bar{y}_{i.} - \bar{y}_{...})' \\
B_{xy} &= \sum_i^N T(\bar{x}_{i.} - \bar{x}_{...})(\bar{y}_{i.} - \bar{y}_{...})
\end{aligned}$$

It is easy to show that

$$\begin{aligned}T_{xx} &= W_{xx} + B_{xx} \\T_{yy} &= W_{yy} + B_{yy} \\T_{xy} &= W_{xy} + B_{xy}\end{aligned}$$

And similarly we can define between-groups estimator

$$\beta_{between} = \frac{B_{xy}}{B_{xx}}$$

We rewrite the estimator of β for the pooled model

$$\begin{aligned}\beta_{pooled} &= \frac{T_{xy}}{T_{xx}} \\&= \frac{B_{xy}}{T_{xx}} + \frac{W_{xy}}{T_{xx}} \\&= \frac{B_{xx}}{T_{xx}} \frac{B_{xy}}{B_{xx}} + \frac{W_{xx}}{T_{xx}} \frac{W_{xy}}{W_{xx}}\end{aligned}$$

because

$$T_{xx} = W_{xx} + B_{xx}$$

So if we define

$$\omega = \frac{W_{xx}}{T_{xx}}$$

then

$$\beta_{pooled} = \omega \frac{W_{xy}}{W_{xx}} + (1 - \omega) \frac{B_{xy}}{B_{xx}}$$

When $T \rightarrow \infty$, how do W_{xx} and B_{xx} behave? remember

$$\begin{aligned}B_{xx} &= \sum_i^N T(\bar{x}_{i.} - \bar{x}_{...})(\bar{x}_{i.} - \bar{x}_{...})' \\W_{xx} &= \sum_i^N \sum_t^T (x_{it} - \bar{x}_{i.})(x_{it} - \bar{x}_{i.})^T\end{aligned}$$

When T increase, $\bar{x}_{i.}$ would move close to $E[x_{i.}]$. so its variance would decrease, so B_{xx} grows sub-linealy. In practice, most of the panel data noise comes from transitory noise. so W_{xx} grows faster than B_{xx} when T is large. Then $\omega \rightarrow 1$. So in large T , the pooled estimator converges to fix effect estimator.

5.3.4 Limitations of fixed effect model

1) Time-invariant variables are dropped

For a given individual i and regressor k , and any different time t and $s (t \neq s)$, if $x_{itk} = x_{isk}$, this x variable will be dropped during estimation. Because this condition leads to $x_{it} = \bar{x}_i$.

2) Loss degree of freedom

In fixed effect, each individual adds up one parameter, so total N individuals need N parameters. This can cause problem in short panel.

5.3.5 F test for fixed effect model

$$\begin{aligned} F &= \frac{(SSE_{pooled} - SSE_{fix}) / ((NT - 1 - K) - (NT - N - K))}{SSE_{fix} / (NT - N - K)} \\ &= \frac{(SSE_{pooled} - SSE_{fix}) / (N - 1)}{SSE_{fix} / (NT - N - K)} \end{aligned}$$

5.4 The random effect model

5.4.1 Model formulation and estimator

In fixed effect model, we treat individual mean α_i is a constant. In random effect model, we treat α_i as a random variable. We write our random effect model as

$$y_{it} = \sum_k x_{itk} \beta_k + \alpha_i + u_{it}$$

Where $\alpha_i \in Normal(0, \sigma_\alpha^2)$, $u_{it} \in Normal(0, \sigma_u^2)$, $E(\alpha_i u_{it}) = 0$. Define

$$v_{it} = \alpha_i + u_{it}$$

The variance of v_{it} for fixed i is

$$E(v_{it} v_{is}) = E(\alpha_i + u_{it})(\alpha_i + u_{is}) = \sigma_\alpha^2 + \delta_{ts} \sigma_u^2$$

The last term is non-zero only when $t = s$. Thus, the covariance matrix for individual i is

$$V_i = \begin{pmatrix} \sigma_\alpha^2 + \sigma_u^2 & \sigma_\alpha^2 & \dots & \sigma_\alpha^2 \\ \sigma_\alpha^2 & \sigma_\alpha^2 + \sigma_u^2 & \dots & \sigma_\alpha^2 \\ \dots & \dots & \dots & \dots \\ \sigma_\alpha^2 & \sigma_\alpha^2 & \dots & \sigma_\alpha^2 + \sigma_u^2 \end{pmatrix}$$

where V_i is $T \times T$ matrix. Stack together for all individuals, the whole covariance matrix is

$$V = \begin{pmatrix} V_1 & 0 & \dots & 0 \\ 0 & V_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & V_N \end{pmatrix}$$

Based on the solution of GLS, we define

$$V_i^{-1/2} = \frac{1}{\sigma_u} [I - \frac{\theta}{T} 1_T 1_T']$$

where

$$\theta = 1 - \frac{\sigma_u}{\sqrt{\sigma_u^2 + T\sigma_\alpha^2}}$$

and the transformation of X_i and y_i

$$\begin{aligned}\tilde{y}_i &= V_i^{-1/2} y_i = y_i - \theta \bar{y}_i \\ \tilde{X}_i &= V_i^{-1/2} X_i = X_i - \theta \bar{X}_i\end{aligned}$$

where y_i is $T \times 1$ vector, and X_i is $T \times K$ matrix. The estimator of β becomes

$$\hat{\beta}_{RE} = (\tilde{X}' V^{-1} \tilde{X})^{-1} (\tilde{X}' V^{-1} \tilde{Y})$$

With a little derivation, we can prove $\hat{\beta}_{RE}$ is a combination of estimator of pooled model and fixed effect model:

$$\hat{\beta}_{RE} = (1 - \omega) \hat{\beta}_{pooled} + \omega \hat{\beta}_{FE}$$

where

$$\omega = \frac{T\sigma_\alpha^2}{T\sigma_\alpha^2 + \sigma_u^2}$$

5.4.2 Connections to estimator of pooled model and fixed effect model

1) when $\sigma_\alpha \gg \sigma_u$, then $\omega \rightarrow 1$, $\theta \rightarrow 1$, this leads to fixed-effect model. In fixed effect model

a. Large σ_α means large variation of α for different i , in other words, α_i is very different.

b. Since Y_{it} for fixed i is center around α_i , this means centers of Y_i given different i are far apart.

c. The smaller σ_u compared to σ_α means the variation of Y_{it} is small. Therefore, the probability density functions of Y_i for different i barely overlap, and their tails almost do not touch each other at all.

2) when $\sigma_\alpha \ll \sigma_u$, then $\omega \rightarrow 0$, $\theta \rightarrow 0$, this leads to pooled model. In pooled model,

a. Small σ_α means α_i s are almost identical.

b. Since Y_{it} for fixed i is center around α_i , this means centers of Y_i given different i are very close to each other.

c. The larger σ_u compared to σ_α means the variation of Y_{it} is large. Therefore, the probability density functions of Y_i for different i are completely overlapping. The difference between different individual i is hardly visible. So it is unnecessary to model the individual mean.

5.4.3 Estimation of σ_u^2 and σ_α^2

Estimation of σ_u^2 can be based on the regression formula

$$y_{it} - \bar{y}_i = (X_{it} - \bar{X}_{i.})\beta + (u_{it} - \bar{u}_i.)$$

and σ_u^2 can be estimated using sum of square error and its degrees of freedom.

$$\hat{\sigma}_u^2 = \frac{\sum_{i=1}^N \sum_{t=1}^T [(y_{it} - \bar{y}_i) - \hat{\beta}_{fix}(x_{it} - \bar{x}_i)]^2}{N(T-1) - K}$$

σ_α^2 can be estimated by first taking the average of random effect model.

$$\bar{y}_i = \mu + \bar{X}_i\beta + \alpha_i + \bar{u}_i$$

Then because α and u_i are independent,

$$Var(\bar{y}_i - \mu - \bar{X}_i\beta) = Var(\alpha) + Var(\bar{u}_i) = Var(\alpha) + \frac{1}{T}Var(u)$$

Then

$$\sigma_\alpha^2 = \frac{\sum_{i=1}^N (\bar{y}_i - \hat{\mu} - \bar{X}_i\hat{\beta})^2}{N - (K+1)} - \frac{1}{T}\hat{\sigma}_u^2$$

5.4.4 Coping with the Limitation of Random Model

The random model, like linear model in general, assumes α does not correlate with variable X . Mundlak introduced auxiliary regression where we write

$$\alpha_i = \sum_k \bar{x}_{i.k} a_k + \omega_i$$

where $\omega_i \in Normal(0, \sigma_\omega^2)$. So

$$y_{it} = \sum_k x_{itk}\beta_k + \sum_k \bar{x}_{i.k}a_k + \omega_i + u_{it}$$

If we define the error term

$$v_{it} = \omega_i + u_{it}$$

The above equation can be written in matrix form:

$$\begin{pmatrix} y_{11} \\ y_{12} \\ \dots \\ y_{1T} \\ y_{21} \\ y_{22} \\ \dots \\ y_{2T} \\ \dots \\ y_{NT} \end{pmatrix} = \begin{pmatrix} 1_T & 0 & 0 & \dots & 0 \\ 0 & 1_T & 0 & \dots & 0 \\ 0 & 0 & 1_T & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_T \end{pmatrix} \begin{pmatrix} \bar{x}_{1.1} & \bar{x}_{1.2} & \dots & \bar{x}_{1.K} \\ \bar{x}_{2.1} & \bar{x}_{2.2} & \dots & \bar{x}_{2.K} \\ \bar{x}_{3.1} & \bar{x}_{3.2} & \dots & \bar{x}_{3.K} \\ \dots & \dots & \dots & \dots \\ \bar{x}_{N.1} & \bar{x}_{N.2} & \dots & \bar{x}_{N.K} \end{pmatrix} \begin{pmatrix} \bar{a}_1 \\ \bar{a}_2 \\ \bar{a}_3 \\ \dots \\ \bar{a}_K \end{pmatrix}$$

$$+ \begin{pmatrix} x_{111} & x_{112} & \dots & x_{11K} \\ x_{121} & x_{122} & \dots & x_{12K} \\ \dots & \dots & \dots & \dots \\ x_{1T1} & x_{1T2} & \dots & x_{1TK} \\ x_{211} & x_{212} & \dots & x_{21K} \\ x_{221} & x_{222} & \dots & x_{22K} \\ \dots & \dots & \dots & \dots \\ x_{2T1} & x_{2T2} & \dots & x_{2TK} \\ \dots & \dots & \dots & \dots \\ x_{NT1} & x_{NT2} & \dots & x_{NTK} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_K \end{pmatrix} + \begin{pmatrix} v_{11} \\ v_{12} \\ \dots \\ v_{1T} \\ v_{21} \\ v_{22} \\ \dots \\ v_{2T} \\ \dots \\ v_{NT} \end{pmatrix}$$

Similarly to the original random model, the variance matrix of v_{it} is

$$V_i = \begin{pmatrix} \sigma_\omega^2 + \sigma_u^2 & \sigma_\omega^2 & \dots & \sigma_\omega^2 \\ \sigma_\omega^2 & \sigma_\omega^2 + \sigma_u^2 & \dots & \sigma_\omega^2 \\ \dots & \dots & \dots & \dots \\ \sigma_\omega^2 & \sigma_\omega^2 & \dots & \sigma_\omega^2 + \sigma_u^2 \end{pmatrix}$$

The way to estimator β should be GLM, but what happens if we just run OLS on this model by assuming v_{it} follows standard normal distribution? Remember the method of doing partial regression. Running OLS on this model is equivalently running the **residual** of Y after regressing on \bar{X}_i , on the residual of X after regressing on \bar{X}_i .

Obviously, regressing X_{it} on X_i , leads the residual

$$X_{it}^* = X_{it} - \bar{X}_i.$$

Where each term is a K by 1 vector Let us regress Y_{it} on \bar{X}_i , and compute the residual. The regression is

$$y_{it} = \mu + \lambda \bar{x}_i + \eta_{it}$$

We realize that μ and \bar{x}_i are constant for a given i. So the fitted value on a constant is the mean of the y_{it} with respect t. We can write

$$\bar{y}_i = \sum_k \beta_k \bar{x}_{i.k} + \sum_k \bar{x}_{i.k} a_k + \bar{v}_i = \sum_k (\beta_k + a_k) \bar{x}_{i.k} + \bar{v}_i$$

We know $E[v_i] = 0$. So when we regress y on \bar{x}_i , the fitted value becomes

$$\hat{y}_{it} = \sum_k (\beta_k + a_k) \bar{x}_{i.k}$$

Then the residual is

$$\begin{aligned}
y_{it} - \hat{y}_{it} &= \sum_k x_{ik} \beta_k + \sum_k \bar{x}_{i.k} a_k + \epsilon_{it} - \left(\sum_k (\beta_k + a_k) \bar{x}_{i.k} \right) \\
&= \sum_k \beta_k (x_{ik} - \bar{x}_{i.k}) + \epsilon_{it} - \bar{\epsilon}_i. \\
&= y_{it} - \bar{y}_i.
\end{aligned}$$

So the OLS we try to run is regress $y_{it} - \bar{y}_i$ on $x_{it} - \bar{x}_i$, which is equivalent to fixed effect estimation.

5.5 Fixed Effect Model Generalization

We can generalize fixed effect model by including more individual specific variables which vary across different individuals and but do not vary over time. We can write the model as

$$\begin{aligned}
y_{it} &= \mu + \alpha_i + z_{i1}\gamma_1 + z_{i2}\gamma_2 + \dots + z_{ip}\gamma_p + \sum_k x_{itk}\beta_k + u_{it} \\
&= \mu + \alpha_i + \sum_p z_{ip}\gamma_p + \sum_k x_{itk}\beta_k + u_{it}
\end{aligned}$$

$$\begin{aligned}
&\begin{pmatrix} y_{11} \\ y_{12} \\ \dots \\ y_{1T} \\ y_{21} \\ y_{22} \\ \dots \\ y_{2T} \\ \dots \\ y_{NT} \end{pmatrix} = \begin{pmatrix} 1_T & 0 & 0 & \dots & 0 \\ 0 & 1_T & 0 & \dots & 0 \\ 0 & 0 & 1_T & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_T \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \dots \\ \alpha_N \end{pmatrix} \\
&+ \begin{pmatrix} 1_T & 0 & 0 & \dots & 0 \\ 0 & 1_T & 0 & \dots & 0 \\ 0 & 0 & 1_T & \dots & 0 \\ 0 & 0 & 0 & \dots & 1_T \end{pmatrix} \begin{pmatrix} z_{11} & z_{12} & \dots & z_{1p} \\ z_{21} & z_{22} & \dots & z_{2p} \\ \dots & \dots & \dots & \dots \\ z_{N1} & z_{N2} & \dots & z_{Np} \end{pmatrix} \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ \dots \\ \gamma_N \end{pmatrix} \\
&+ \begin{pmatrix} x_{111} & x_{112} & \dots & x_{11K} \\ x_{121} & x_{122} & \dots & x_{12K} \\ \dots & \dots & \dots & \dots \\ x_{1T1} & x_{1T2} & \dots & x_{1TK} \\ x_{211} & x_{212} & \dots & x_{21K} \\ x_{221} & x_{222} & \dots & x_{22K} \\ \dots & \dots & \dots & \dots \\ x_{2T1} & x_{2T2} & \dots & x_{2TK} \\ \dots & \dots & \dots & \dots \\ x_{NT1} & x_{NT2} & \dots & x_{NTK} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_K \end{pmatrix} + \begin{pmatrix} u_{11} \\ u_{12} \\ \dots \\ u_{1T} \\ u_{21} \\ u_{22} \\ \dots \\ u_{2T} \\ \dots \\ u_{NT} \end{pmatrix}
\end{aligned}$$

This model suffers multicollinearity so α and γ are not both estimable. We can set

$$\alpha_i^* = \alpha_i + \sum_p z_{ip} \gamma_p$$

Then the model is effectively the same as the original fixed effect model, so we can use the same trick of solving fixed effect model to estimate β . To estimate γ , we note

$$\bar{y}_i - \sum_k \bar{x}_{i.k} \beta_k = \sum_p z_{ip} \gamma_p + \alpha_i + \bar{u}_i$$

Let $\epsilon_i = \alpha_i + \bar{u}_i$ and by minimizing $\sum_i \epsilon_i^2$, we obtain

$$\hat{\gamma} = \left[\sum_{i=1}^N (z_i - \bar{z})(z_i - \bar{z})' \right]^{-1} \left\{ \sum_{i=1}^N (z_i - \bar{z}) [(\bar{y}_i - \bar{y}) - (\bar{x}_i - \bar{x})' \beta] \right\}$$

Testing of σ_α We define the null and alternative hypothesis

$$H_0 : \sigma_\alpha^2 = 0$$

$$H_1 : \sigma_\alpha^2 > 0$$

We construct LM test, which uses the score function (gradient of log-likelihood) with respect to σ_α^2 evaluated at $\sigma_\alpha^2 = 0$:

$$\begin{aligned} LM &= \frac{(\partial l / \partial \sigma_\alpha^2)^2}{\text{Var}(\partial l / \partial \sigma_\alpha^2)} \Big|_{\sigma_\alpha^2=0} \\ &= \frac{nT}{2(T-1)} \left[\frac{[\sum_{i=1}^n (\bar{e}_i)^2]}{(\frac{1}{nT} \sum_{i=1}^n \sum_{t=1}^T e_{it}^2)} - 1 \right]^2 \end{aligned}$$

6 Panel Data Time-Dependent Models

6.1 Introduction and Model Specification

Panel data (or longitudinal data) consist of observations on N cross-sectional units (individuals, firms, countries, etc.) over T time periods. Time dependence in panel models arises when the error terms or dependent variables exhibit correlation across time for the same unit. A general linear dynamic panel model can be written as:

$$y_{it} = \alpha y_{i,t-1} + \mathbf{x}_{it}' \beta + \eta_i + \varepsilon_{it}, \quad i = 1, \dots, N, \quad t = 1, \dots, T,$$

where:

- y_{it} is the dependent variable for unit i at time t ,
- $y_{i,t-1}$ is the lagged dependent variable (introducing **dynamics**),
- \mathbf{x}_{it} is a $k \times 1$ vector of exogenous regressors,
- η_i is the **unobserved individual effect** (fixed or random),

- ε_{it} is the idiosyncratic error term.

Key features of time dependence:

1. **State dependence:** Past outcomes affect current outcomes ($y_{i,t-1}$ term).
2. **Serial correlation:** ε_{it} may follow an AR(p) process.
3. **Heterogeneity:** Individual-specific effects η_i that may be correlated with regressors.

6.2 Challenges in Estimation

- **Dynamic panel bias:** The lagged dependent variable $y_{i,t-1}$ is correlated with the individual effect η_i , making standard estimators (OLS, fixed effects) inconsistent for fixed T as $N \rightarrow \infty$.
- **Nickell bias:** The within (fixed effects) estimator is biased of order $O(1/T)$, which diminishes only when T is large.
- **Initial conditions problem:** The process may not be stationary, and the initial observation y_{i0} may be correlated with η_i .

6.3 Estimation Methods

6.3.1 Generalized Method of Moments (GMM)

The most common approach for dynamic panels with small T and large N .

First-difference transformation to eliminate η_i :

$$\Delta y_{it} = \alpha \Delta y_{i,t-1} + \Delta \mathbf{x}'_{it} \boldsymbol{\beta} + \Delta \varepsilon_{it}, \quad t = 2, \dots, T.$$

The problem: $\Delta y_{i,t-1}$ is correlated with $\Delta \varepsilon_{it}$.

Arellano-Bond (1991) GMM:

- Use **lagged levels** of y_{it} as instruments for $\Delta y_{i,t-1}$:
 - For Δy_{i2} , instrument: y_{i0}
 - For Δy_{i3} , instruments: y_{i0}, y_{i1}
 - For Δy_{iT} , instruments: $y_{i0}, y_{i1}, \dots, y_{i,T-2}$
- Moment conditions: $\mathbb{E}[y_{i,t-s} \Delta \varepsilon_{it}] = 0$ for $s \geq 2, t = 2, \dots, T$.
- Efficient GMM estimator:

$$\hat{\boldsymbol{\theta}}_{\text{GMM}} = (\mathbf{Z}' \mathbf{X} \mathbf{W}_N \mathbf{X}' \mathbf{Z})^{-1} \mathbf{Z}' \mathbf{X} \mathbf{W}_N \mathbf{Z}' \mathbf{y},$$

where \mathbf{Z} is the instrument matrix, \mathbf{X} the regressor matrix, and \mathbf{W}_N a weight matrix.

Blundell-Bond (1998) System GMM:

- Combines equations in differences with equations in levels.
- Additional moment conditions: Use lagged differences as instruments for levels equations.

- More efficient when α is close to 1 or when series are persistent.
- Moment conditions:

$$\begin{aligned}\mathbb{E}[y_{i,t-s}\Delta\varepsilon_{it}] &= 0 \quad (\text{difference equations}) \\ \mathbb{E}[\Delta y_{i,t-s}(\eta_i + \varepsilon_{it})] &= 0 \quad (\text{level equations})\end{aligned}$$

for appropriate s .

6.3.2 Maximum Likelihood Estimation (MLE)

For known distribution of errors, MLE can be used:

$$L(\boldsymbol{\theta}) = \prod_{i=1}^N f(y_{i1}, \dots, y_{iT} | \mathbf{x}_{i1}, \dots, \mathbf{x}_{iT}; \boldsymbol{\theta}).$$

- **Conditional MLE:** Treat initial conditions y_{i0} as given.
- **Unconditional MLE:** Model the distribution of y_{i0} (e.g., assume stationarity).
- Implementation requires specifying the joint distribution of $(\eta_i, \varepsilon_{i1}, \dots, \varepsilon_{iT})$.

6.3.3 Long-T Panels: Bias-Corrected Estimators

When T is moderately large (e.g., $T > 10$):

- **Bias-corrected fixed effects:** Correct the Nickell bias using analytical or bootstrap methods.
- **Least Squares Dummy Variable (LSDV)** with correction:

$$\hat{\alpha}_{\text{corrected}} = \hat{\alpha}_{\text{LSDV}} + \frac{1 + \hat{\alpha}_{\text{LSDV}}}{T - 1}.$$

6.4 Hypothesis Testing

6.4.1 Specification Tests

1. Serial Correlation Tests:

- **Arellano-Bond test for autocorrelation:**
 - Test H_0 : No autocorrelation in first-differenced errors at order m .
 - Based on sample autocovariances of differenced residuals.
 - Critical: Test for AR(1) in $\Delta\varepsilon_{it}$ (expected) and AR(2) (should be zero under null).
- **Breusch-Godfrey test for panel data:** LM test for serial correlation.

2. Overidentification/Specification Tests:

- **Sargan-Hansen J-test:**

$$J = \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{z}'_i \hat{\mathbf{u}}_i \right)' \hat{\mathbf{W}} \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N \mathbf{z}'_i \hat{\mathbf{u}}_i \right) \xrightarrow{d} \chi^2_{q-k},$$

where q is number of instruments, k number of parameters.

- Tests validity of moment conditions (instruments).
- Difference-in-Sargan test: Compare subsets of instruments.

3. Hausman Tests:

- Compare fixed effects vs. random effects estimators:

$$H = (\hat{\beta}_{\text{FE}} - \hat{\beta}_{\text{RE}})' [\widehat{\text{Var}}(\hat{\beta}_{\text{FE}}) - \widehat{\text{Var}}(\hat{\beta}_{\text{RE}})]^{-1} (\hat{\beta}_{\text{FE}} - \hat{\beta}_{\text{RE}}) \sim \chi^2_k.$$

- For dynamic models, compare GMM estimators with different instrument sets.

4. Unit Root Tests for Panel Data:

- **Levin-Lin-Chu (LLC) test:** Assumes common unit root process.
- **Im-Pesaran-Shin (IPS) test:** Allows for individual unit root processes.
- **Maddala-Wu Fisher test:** Combines p-values from individual ADF tests.

6.4.2 Inference on Parameters

Wald tests for linear restrictions $H_0 : \mathbf{R}\theta = \mathbf{r}$:

$$W = (\mathbf{R}\hat{\theta} - \mathbf{r})' [\mathbf{R}\widehat{\text{Var}}(\hat{\theta})\mathbf{R}']^{-1} (\mathbf{R}\hat{\theta} - \mathbf{r}) \sim \chi^2_q,$$

where q is number of restrictions.

For GMM, use robust variance estimator:

$$\widehat{\text{Var}}(\hat{\theta}_{\text{GMM}}) = (\mathbf{X}'\mathbf{Z}\mathbf{W}_N\mathbf{Z}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Z}\mathbf{W}_N\hat{\mathbf{S}}\mathbf{W}_N\mathbf{Z}'\mathbf{X}(\mathbf{X}'\mathbf{Z}\mathbf{W}_N\mathbf{Z}'\mathbf{X})^{-1},$$

where $\hat{\mathbf{S}}$ estimates $\text{Var}(N^{-1/2} \sum \mathbf{z}'_i \mathbf{u}_i)$.

6.5 Practical Considerations

1. **Instrument proliferation:** Too many instruments in GMM can overfit endogenous variables and bias results. Use:
 - Collapsing instruments (combining lags)
 - Limiting lag depth
 - Principal components of instruments
2. **Weak instruments:** When lagged levels are poor predictors of differences (high persistence). Check with:

- First-stage F-statistics
- Stock-Yogo critical values
- Anderson-Rubin tests

3. Choice between difference and system GMM:

- Use difference GMM when series are stationary
- Use system GMM for persistent data or when α is near 1
- Compare with Hansen test for additional moments

4. Time effects: Include time dummies to account for common shocks:

$$y_{it} = \alpha y_{i,t-1} + \mathbf{x}'_{it}\boldsymbol{\beta} + \eta_i + \delta_t + \varepsilon_{it}.$$

5. Nonlinear dynamic panels: For binary, count, or other nonlinear outcomes:

- Dynamic probit/logit with correlated random effects
- Conditional maximum likelihood for logit
- GMM for exponential regression models

6.6 Software Implementation

Stata:

```
* Arellano-Bond difference GMM
xtabond y x1 x2, lags(1)
estat abond // Test autocorrelation
estat sargan // Test overidentifying restrictions
```

```
* Blundell-Bond system GMM
xtdpdsys y x1 x2, lags(1)
```

R:

```
library(plm)
# Difference GMM
gmm_diff <- pgmm(y ~ lag(y, 1) + x1 + x2 | lag(y, 2:99),
                 data = pdata, effect = "individual",
                 model = "twosteps")
summary(gmm_diff)

# System GMM
gmm_sys <- pgmm(y ~ lag(y, 1) + x1 + x2 | lag(y, 2:99),
                data = pdata, effect = "individual",
                model = "onestep", transformation = "ld")
summary(gmm_sys)
```

6.7 Summary

Time-dependent panel models capture important dynamic relationships but require careful estimation due to:

- The inherent endogeneity from lagged dependent variables
- Unobserved heterogeneity correlated with regressors
- Potential serial correlation in errors

GMM methods, particularly Arellano-Bond and Blundell-Bond estimators, are standard for short panels. For longer panels, bias-corrected fixed effects or MLE may be appropriate. Comprehensive specification testing is crucial for valid inference