

# Generalised Least Squares

## References

- Davidson, Russell and James G. MacKinnon (2004). *Econometric Theory and Methods*. Oxford University Press.
- Hayashi, Fumio (2000). *Econometrics*. Princeton University Press.

## 1 Introduction

If the parameters of a regression model are to be estimated efficiently by OLS, the error terms must be uncorrelated and have the same variance. These assumptions are needed to prove the Gauss-Markov Theorem and to show that the least squares estimator is asymptotically efficient. Moreover, the usual estimators of the covariance matrix of the OLS estimators are not valid when these assumptions do not hold, although alternative “sandwich” covariance matrix estimators that are asymptotically valid may be available.

Thus, it is clear that we need new estimation methods to handle regression models with error terms that are heteroskedastic, serially correlated, or both.

Since heteroskedasticity and serial correlation affect both linear and nonlinear regression models in the same way, there is no harm in limiting our attention to the simpler, linear case. We will be concerned with the model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \mathbb{E}[\mathbf{u}\mathbf{u}^\top] = \boldsymbol{\Omega}, \quad (1)$$

where  $\boldsymbol{\Omega}$ , the covariance matrix of the error terms, is a positive definite  $n \times n$  matrix. If  $\boldsymbol{\Omega}$  is equal to  $\sigma^2\mathbf{I}$ , then (1) is just the linear regression model we have always looked at when constructing OLS estimators (error terms are uncorrelated and homoskedastic). If  $\boldsymbol{\Omega}$  is diagonal with non-constant diagonal elements, then the error terms are still uncorrelated, but they are heteroskedastic. If  $\boldsymbol{\Omega}$  is not diagonal, then  $u_i$  and  $u_j$  are correlated whenever  $\Omega_{ij}$ , the  $ij^{th}$  element of  $\boldsymbol{\Omega}$ , is nonzero.

We will obtain an efficient estimator for the vector  $\boldsymbol{\beta}$  in the model (1). This efficient estimator is called the General Least Squares (GLS) estimator.

## 2 The GLS estimator

In order to obtain the GLS estimator, we need to transform the model so that the transformed model satisfies the conditions of the Gauss-Markov theorem. Estimating the transformed model by OLS therefore yields efficient estimates. The

transformation is expressed in terms of an  $n \times n$  matrix  $\Psi$ , which is usually triangular, that satisfies the equation:

$$\Omega^{-1} = \Psi \Psi^\top. \quad (2)$$

Such a matrix can always be found using Crout's algorithm. Premultiplying (1) by  $\Psi^\top$  gives:

$$\Psi^\top \mathbf{y} = \Psi^\top \mathbf{X} \boldsymbol{\beta} + \Psi^\top \mathbf{u}. \quad (3)$$

Because the covariance matrix  $\Omega$  is nonsingular, the matrix  $\Psi$  must be as well, and so the transformed regression model (3) is perfectly equivalent to the original model (1). The OLS estimator of  $\boldsymbol{\beta}$  from regression (3) is:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\text{GLS}} &= (\mathbf{X}^\top \Psi \Psi^\top \mathbf{X})^{-1} \mathbf{X}^\top \Psi \Psi^\top \mathbf{y} \\ &= (\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega^{-1} \mathbf{y}. \end{aligned} \quad (4)$$

This estimator is called the GLS estimator of  $\boldsymbol{\beta}$ . The GLS estimator solves the following problem:

$$\hat{\boldsymbol{\beta}}_{\text{GLS}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X} \boldsymbol{\beta}\|_{\Omega} = (\mathbf{y} - \mathbf{X} \boldsymbol{\beta})^\top \Omega^{-1} (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}),$$

and its first order condition is:

$$-2 \mathbf{X}^\top \Omega^{-1} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}_{\text{GLS}}) = 0.$$

Substituting (1) into (4) yields:

$$\hat{\boldsymbol{\beta}}_{\text{GLS}} - \boldsymbol{\beta} = (\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega^{-1} \mathbf{u}. \quad (5)$$

Since the GLS estimator is simply the OLS estimator from (3), its covariance matrix can be found directly from the standard formula for the OLS covariance matrix:

$$\mathbb{E} \left[ (\hat{\boldsymbol{\beta}}_{\text{GLS}} - \boldsymbol{\beta}) (\hat{\boldsymbol{\beta}}_{\text{GLS}} - \boldsymbol{\beta})^\top | \mathbf{X} \right] = (\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Omega^{-1} \mathbb{E}[\mathbf{u} \mathbf{u}^\top | \mathbf{X}] \Omega^{-1} \mathbf{X} (\mathbf{X}^\top \Omega^{-1} \mathbf{X})^{-1},$$

which simplifies to:

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

## 2.1 Efficiency of the GLS estimator and unbiasedness

It's quite easy to prove that the GLS estimator is unbiased. From (3) we have the following model and assumptions:

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}} \boldsymbol{\beta} + \tilde{\mathbf{u}}, \quad \mathbb{E}[\tilde{\mathbf{u}} | \mathbf{X}] = 0, \quad \mathbb{E}(\tilde{\mathbf{u}} \tilde{\mathbf{u}}^\top | \mathbf{X}) = \sigma^2 \mathbf{I}, \quad (6)$$

where a tilde variable denotes that the variable has been premultiplied by  $\Psi^\top$ , so the GLS estimator is the OLS estimator of the transformed regression (3):

$$\hat{\beta}_{\text{GLS}} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{y}}.$$

To show its unbiasedness:

$$\begin{aligned} \hat{\beta}_{\text{GLS}} &= \beta + (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \tilde{\mathbf{u}} \\ \mathbb{E}(\hat{\beta}_{\text{GLS}} | \mathbf{X}) &= \beta + (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \underbrace{\mathbb{E}[\tilde{\mathbf{u}} | \mathbf{X}]}_{=0}. \end{aligned}$$

Thus, the GLS estimator is unbiased. To show its efficiency, consider (1) and the OLS estimator:

$$\begin{aligned} \hat{\beta}_{\text{OLS}} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \\ \implies \mathbb{E} \left[ (\hat{\beta}_{\text{OLS}} - \beta) (\hat{\beta}_{\text{OLS}} - \beta)^\top | \mathbf{X} \right] &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \\ \therefore \text{Var}(\hat{\beta}_{\text{OLS}} | \mathbf{X}) &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \end{aligned}$$

When  $\boldsymbol{\Omega} = \sigma^2 \mathbf{I}$ , as in the case for unconditional homoskedasticity, then the variance-covariance matrix of the OLS estimator is  $\sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$ . Also, notice that if  $\boldsymbol{\Omega} = \mathbb{E}[\mathbf{u}\mathbf{u}^\top | \mathbf{X}]$  (conditional homoskedasticity), then we have the sandwich estimator, which we looked at when we analysed heteroskedasticity-consistent variance-covariance matrices.

One can also easily see that

$$\left( \text{Var}(\hat{\beta}_{\text{GLS}}) \right)^{-1} - \left( \text{Var}(\hat{\beta}_{\text{OLS}}) \right)^{-1}$$

results in a positive semidefinite matrix. i.e.,

$$\left( \text{Var}(\hat{\beta}_{\text{GLS}}) \right)^{-1} - \left( \text{Var}(\hat{\beta}_{\text{OLS}}) \right)^{-1} = (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X}) - \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X}$$

is positive semidefinite, implying that the GLS estimator is more efficient than the OLS estimator when we no longer can assume homoskedasticity.

## 2.2 Weighted least squares

It is particularly easy to obtain GLS estimates when the error terms are heteroskedastic but uncorrelated. This implies that the matrix  $\boldsymbol{\Omega}$  is diagonal. Let  $\omega_t^2$  denote the  $t^{\text{th}}$  diagonal element of  $\boldsymbol{\Omega}$ . Then  $\boldsymbol{\Omega}^{-1}$  is a diagonal matrix with  $t^{\text{th}}$  diagonal element  $\omega_t^{-2}$ , and  $\Psi$  can be chosen as the diagonal matrix with  $t^{\text{th}}$  diagonal element  $\omega_t^{-1}$ . Thus, we see that, for a typical observation, regression (3) can be written as:

$$\frac{1}{\omega_t} y_t = \frac{1}{\omega_t} \mathbf{X}_t \beta + \frac{1}{\omega_t} u_t \quad (7)$$

where we are trying to “weigh the observations”. Also, as said:

$$\mathbf{\Omega} = \begin{bmatrix} \omega_1^2 & 0 & \cdots & 0 \\ 0 & \omega_2^2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \omega_n^2 \end{bmatrix} = \text{diag}(\{\omega_t^2\}_{t=1}^n),$$

which implies that

$$\mathbf{\Psi} = \begin{bmatrix} \omega_1^{-1} & 0 & \cdots & 0 \\ 0 & \omega_2^{-1} & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \omega_n^{-1} \end{bmatrix}.$$

In practice, we do not know  $\omega_t^2$ . One choice of skedastic function is

$$\omega_t^2 = \mathbb{E}[u_t^2] = \exp(\delta + Z_t\gamma).$$

Note that for any values of  $\delta$  and  $\gamma$ ,  $\exp(\delta + Z_t\gamma) > 0$ . Hence why many economists like to express the function as an exponential function, this makes later evaluation, using logs, much easier. Also, some or all elements of  $Z_t$  may well belong to  $X_t$ . The function  $\exp(\delta + Z_t\gamma)$  is an example of a skedastic function. In the same way that a regression function determines the conditional mean of a random variable, a skedastic function determines its conditional variance. The skedastic function has the property that it is positive for any parameter  $\gamma$ . This is desirable, since any negative estimated variances would be highly inconvenient. The estimation procedure for  $\omega_t^2$  and  $\beta$  is as follows:

1. Estimate  $\beta$  by OLS, and compute  $\hat{\mathbf{u}} = \mathbf{M}_{\mathbf{X}}\mathbf{y}$ , the OLS residuals.
2. Estimate  $\delta$  and  $\gamma$  from the auxiliary regression

$$\ln \hat{u}_t^2 = \delta + Z_t\gamma + v_t.$$

3. Compute and estimate  $\omega_t$  from

$$\hat{\omega}_t = \left( \exp(\hat{\delta} + Z_t\hat{\gamma}) \right)^{\frac{1}{2}}.$$

4. Estimate  $\beta$  using  $\hat{\omega}_t$  in place of  $\omega_t$ . This is also known as the feasible GLS (FGLS):

$$\hat{\beta}_{\text{FGLS}} = \left( \mathbf{X}^\top \hat{\mathbf{\Omega}}^{-1} \mathbf{X} \right)^{-1} \mathbf{X}^\top \hat{\mathbf{\Omega}}^{-1} \mathbf{y},$$

where

$$\hat{\mathbf{\Omega}} = \mathbf{\Omega}(\hat{\delta}, \hat{\gamma}).$$

### 2.3 Asymptotic distribution of GLS and FGLS

The GLS estimator can be written as

$$\hat{\beta}_{\text{GLS}} = \beta + (\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{u},$$

and

$$\sqrt{n}(\hat{\beta}_{\text{GLS}} - \beta) = \left( n^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X} \right)^{-1} \frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{u}.$$

Now, using the central limit theorem, we have

$$\begin{aligned} \frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{u} &\overset{a}{\sim} N(0, \mathbf{S}_{\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X}}), \\ \text{plim}_{n \rightarrow \infty} \left( n^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X} \right) &= \mathbf{S}_{\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X}}. \end{aligned}$$

Combining the above results we derive

$$\sqrt{n}(\hat{\beta}_{\text{GLS}} - \beta) \overset{a}{\sim} N(0, \mathbf{S}_{\mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X}}^{-1}).$$

The FGLS can be written as

$$\sqrt{n}(\hat{\beta}_{\text{FGLS}} - \beta) = \left( n^{-1} \mathbf{X}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{X} \right)^{-1} \frac{1}{\sqrt{n}} \mathbf{X}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{u},$$

which is very similar to the expression derived for the GLS estimator. One can show that

$$\begin{aligned} \frac{1}{\sqrt{n}} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{u} &\overset{a}{=} \frac{1}{\sqrt{n}} \mathbf{X}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{u}, \text{ and,} \\ \text{plim}_{n \rightarrow \infty} \left( n^{-1} \mathbf{X}^\top \boldsymbol{\Omega}^{-1} \mathbf{X} \right) &= \text{plim}_{n \rightarrow \infty} \left( n^{-1} \mathbf{X}^\top \hat{\boldsymbol{\Omega}}^{-1} \mathbf{X} \right), \end{aligned}$$

i.e. the FGLS and GLS are asymptotically equivalent. This result holds for the case where  $\boldsymbol{\Omega}$  is diagonal and  $\mathbf{X}$  is exogenous or predetermined. However, if  $\boldsymbol{\Omega}$  is not diagonal then the OLS estimator is not consistent whenever any element of  $\mathbf{X}_t$  is a lagged dependent variable.

### 2.4 Heteroskedasticity

There are two situations in which the error terms are heteroskedastic but serially uncorrelated. In the first, the form of the heteroskedasticity is completely unknown, while, in the second, the heteroskedastic function is known except for the value of the parameters that can be estimated consistently. The fact that these HCCME's are sandwich covariance matrices makes it clear that, although they are consistent under standard regularity conditions, OLS is inefficient when the error

terms are heteroskedastic.

If the variances of all the error terms are known, at least up to a scale factor, then efficient estimates can be obtained by weighted least squares. For a linear model, we need to multiply all of the variables by  $\omega_t^{-1}$ , the inverse of the standard error of  $u_t$ , and then use OLS. The usual OLS covariance matrix is perfectly valid, although it is desirable to replace  $s^2$  by 1 if the variances are completely known, since in that case  $s^2 \rightarrow 1$  as  $n \rightarrow \infty$ .

If the form of the heteroskedasticity is known, but the skedastic function depends on unknown parameters, then we can use feasible weighted least squares and still achieve asymptotic efficiency. An example of such a procedure was discussed in the previous lecture. As we have seen, it makes no difference asymptotically whether the  $\omega_t$  are known or merely estimated consistently, although it can certainly make a difference in finite samples. Asymptotically, at least, the usual OLS covariance matrix is just as valid with feasible WLS as with WLS.

#### 2.4.1 Testing for Heteroskedasticity

In some cases, it may be clear from the specification of the model that the error terms must exhibit a particular pattern of heteroskedasticity. In many cases, however, we may hope that the error terms are homoskedastic but be prepared to admit the possibility that they are not. In such cases, if we have no information on the form of the skedastic function, it may be prudent to employ an HCCME, especially if the sample size is large.

If we have information on the form of the skedastic function, we might well wish to use WLS. Before doing so, it is advisable to perform a specification test of the null hypothesis that the error terms are homoskedastic against whatever heteroskedastic alternatives may seem reasonable. There are many ways to perform this type of specification test. The simplest approach that is widely applicable, and the only one that we will discuss, involves running an artificial regression in which the regressand is the vector of squared residuals from the model or test.

A reasonably general model of conditional heteroskedasticity is:

$$\mathbb{E}[u_t^2|\Omega_t] = h(\delta + Z_t\gamma) \quad (8)$$

where the skedastic function  $h(\cdot)$  is a non-linear function that can take on only positive values,  $Z_t$  is a  $1 \times r$  vector of observations on exogenous or predetermined variables that belong to the information set  $\Omega_t$ ,  $\delta$  is a scalar parameter, and  $\gamma$  is an

$r$ -vector of parameters. Under the null,  $\gamma = 0$  and the function collapses to  $h(\delta)$ , a constant. A very simple procedure to detect the presence of heteroskedasticity is to estimate:

$$u_t^2 = \alpha + \mathbf{Z}_t\eta + \text{residuals}, \quad (9)$$

and test  $\eta = 0$  by an  $F$ -test, asymptotically distributed as  $F(r, \infty)$ . This expression is obtained from a Taylor expansion of  $h(\delta + Z_t\gamma)$  and imposing restrictions from the null assumption. Alternatively, we could perform an  $nR^2$  test, where we use the centered  $R^2$  from the above regression, which is asymptotically distributed as  $\chi^2(r)$ . In practice, we don't actually observe  $u_t^2$ , so we replace it with  $\hat{u}_t^2$ .

## 2.5 Time Series Models

The error terms for nearby observations may be correlated, or may appear to be correlated, in any sort of regression model, but this phenomenon is most commonly encountered in models estimated with time-series data, where it is known as serial correlation or autocorrelation. In practice, what appears to be serial correlation may instead just be a misspecified model. In some circumstances, though, it is natural to model the serial correlation by assuming the error terms follow some sort of stochastic process.

If there is reason to believe that serial correlation may be present, the first step is usually to test the null hypothesis that the errors are serially uncorrelated against a plausible alternative that involves serial correlation.

### 2.5.1 Autoregressive errors

One of the simplest and most commonly used stochastic processes is the first-order autoregressive process, or AR(1) process. We have already encountered regression models with error terms that follow such a process. Recall from Davidson & MacKinnon Exercise (6.04) that the AR(1) process can be written as:

$$u_t = \rho u_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim \text{IID}(0, \sigma_\varepsilon^2), \quad |\rho| < 1, \quad (10)$$

where

$$y_t = \mathbf{X}_t\boldsymbol{\beta} + u_t.$$

The error at time  $t$  is equal to some fraction  $\rho$  of the error at time  $t - 1$ , with the sign changed if  $\rho < 0$ , plus the innovation  $\varepsilon_t$ .

The condition in equation (10) that  $|\rho| < 1$  is called a stationary condition, because it is necessary for the AR(1) process to be stationary. There are several

definitions of stationarity in time series analysis.

**Definition:** A series of typical element  $u_t$  is stationary if the unconditional expectation  $\mathbb{E}[u_t]$  and the unconditional variance  $\text{Var}(u_t)$  exist and are independent of  $t$ , and if the covariance  $\text{Cov}(u_t, u_{t-j})$  is also, for any given  $j$ , independent of  $t$ . This definition of stationarity is sometimes referred to as covariance stationarity, or wide sense stationarity.

Consider:

$$\mathbb{E} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix} \begin{bmatrix} u_1 & u_2 & \cdots & u_n \end{bmatrix}$$

and

$$\mathbb{E} \begin{bmatrix} \mathbb{E}[u_1^2] & \mathbb{E}[u_1 u_2] & \cdots & \mathbb{E}[u_1 u_n] \\ \mathbb{E}[u_2 u_1] & \mathbb{E}[u_2^2] & & \vdots \\ \vdots & & \ddots & \mathbb{E}[u_{n-1} u_n] \\ \mathbb{E}[u_n u_1] & \cdots & \mathbb{E}[u_n u_{n-1}] & \mathbb{E}[u_n^2] \end{bmatrix} = \mathbf{\Omega}.$$

We can see that:

$$u_t = \varepsilon_t + \rho \varepsilon_{t-1} + \rho^2 \varepsilon_{t-2} + \rho^3 \varepsilon_{t-3} + \dots,$$

and as the innovations are independent and uncorrelated:

$$\begin{aligned} \sigma_u^2 &= \text{Var}(u_t) = \text{Var}(\varepsilon_t + \rho \varepsilon_{t-1} + \rho^2 \varepsilon_{t-2} + \rho^3 \varepsilon_{t-3} + \dots) \\ &= \text{Var}(\varepsilon_t) + \text{Var}(\rho \varepsilon_t) + \text{Var}(\rho^2 \varepsilon_{t-2}) + \dots \\ &= \text{Var}(\varepsilon_t) + \rho^2 \text{Var}(\varepsilon_{t-1}) + \rho^4 \text{Var}(\varepsilon_{t-2}) + \dots \\ &= (1 + \rho^2 + \rho^4 + \dots) \sigma_\varepsilon^2 \\ &= \frac{1}{1 - \rho^2} \sigma_\varepsilon^2 = \frac{\sigma_\varepsilon^2}{1 - \rho^2} = \sigma_u^2 \end{aligned}$$

We say that the residuals follow an AR(1) process. By imposing  $|\rho| < 1$  we guarantee that the variance of  $u_t$  does not explode when  $t \rightarrow \infty$ .

In order to find our  $\mathbf{\Psi}$  matrix, we need to compute the following results (for the off diagonals in our  $\mathbf{\Omega}$  matrix):



$$\begin{aligned}
\mathbb{E}[u_t u_{t-1}] &= \mathbb{E}[(\rho u_{t-1} + \varepsilon_t) u_{t-1}] \\
&= \rho \mathbb{E}[u_{t-1}^2] + \mathbb{E}[\varepsilon_t u_{t-1}] \\
&= \rho \frac{\sigma_\varepsilon^2}{(1 - \rho^2)}, \\
\mathbb{E}[u_t u_{t-2}] &= \mathbb{E}[(\rho^2 u_{t-2} + \rho \varepsilon_{t-1} + \varepsilon_t) u_{t-2}] \\
&= \rho^2 E(u_{t-2}^2) + \rho E(\varepsilon_{t-1} u_{t-2}) + E(\varepsilon_t u_{t-1}) \\
&= \rho^2 \left( \frac{\sigma_\varepsilon^2}{1 - \rho^2} \right), \\
\Rightarrow E[u_t u_{t-s}] &= \rho^s \left( \frac{\sigma_\varepsilon^2}{1 - \rho^2} \right).
\end{aligned}$$

Then we find:

$$\mathbf{\Omega} = \begin{bmatrix} \frac{\sigma_\varepsilon^2}{1-\rho^2} & \rho \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & \rho^2 \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & \cdots & \rho^{n-1} \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) \\ \rho \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & \frac{\sigma_\varepsilon^2}{1-\rho^2} & & & \vdots \\ \rho^2 \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & & \ddots & & \\ \vdots & & & \ddots & \rho \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) \\ \rho^{n-1} \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & \cdots & & \rho \left( \frac{\sigma_\varepsilon^2}{1-\rho^2} \right) & \frac{\sigma_\varepsilon^2}{1-\rho^2} \end{bmatrix},$$

and

$$\mathbf{\Omega} = \frac{\sigma_\varepsilon^2}{1 - \rho^2} \begin{bmatrix} 1 & \rho & \cdots & \rho^{n-1} \\ \rho & 1 & & \vdots \\ \vdots & & \ddots & \rho \\ \rho^{n-1} & \cdots & \rho & 1 \end{bmatrix},$$

and so:

$$\mathbf{\Psi}(\rho) = \begin{bmatrix} (1 - \rho)^{-\frac{1}{2}} & -\rho & 0 & \cdots & 0 \\ 0 & 1 & -\rho & \cdots & 0 \\ \vdots & 0 & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & -\rho \\ 0 & \cdots & \cdots & 0 & 1 \end{bmatrix},$$

thus we can yield:

$$\begin{aligned} \Psi^\top \mathbf{u} &= \boldsymbol{\varepsilon} \\ \begin{bmatrix} (1-\rho)^{-\frac{1}{2}} & 0 & \cdots & \cdots & 0 \\ -\rho & 1 & 0 & \ddots & \vdots \\ 0 & -\rho & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & 0 \\ 0 & 0 & \cdots & -\rho & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ \vdots \\ u_n \end{bmatrix} &= \begin{bmatrix} (1-\rho^2)u_1 \\ u_2 - \rho u_1 \\ u_3 - \rho u_2 \\ \vdots \\ u_n - \rho u_{n-1} \end{bmatrix} \Rightarrow \begin{matrix} \varepsilon_1 \\ \varepsilon_2 \\ \varepsilon_3 \\ \vdots \\ \varepsilon_n \end{matrix} \end{aligned}$$

The idea is similar to the WLS estimates:

1. Estimate  $\boldsymbol{\beta}$  by OLS, and compute  $\hat{\mathbf{u}} = \mathbf{M}_{\mathbf{X}}\mathbf{y}$ , the OLS residuals (Note:  $X_t$  cannot have lagged dependent variables!).

2. Estimate  $\rho$  from the auxiliary regression:

$$\hat{u}_t = \rho \hat{u}_{t-1} + \varepsilon_t^* \quad (11)$$

we have  $n-1$  elements (because  $\hat{u}_1$  is not included on the LHS).

3. An estimator for  $\Psi(\rho)$  is  $\Psi(\hat{\rho})$ .

The feasible GLS estimator is:

$$\hat{\boldsymbol{\beta}}_{\text{FGLS}} = (\mathbf{X}^\top \Delta(\hat{\rho})^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \Delta(\hat{\rho})^{-1} \mathbf{y},$$

where we define  $\Psi\Psi^\top = \Delta^{-1}$ , and it has the following variance estimator:

$$\text{Var}(\hat{\boldsymbol{\beta}}_{\text{FGLS}}) = s^2 (\mathbf{X}^\top \Delta(\hat{\rho})^{-1} \mathbf{X})^{-1}.$$

$s^2$  is an estimator of the variance of the residual  $\Psi^\top \mathbf{u}$  obtained from the transformed regression:

$$\Psi^\top \mathbf{y} = \Psi^\top \mathbf{X}\boldsymbol{\beta} + \Psi^\top \mathbf{u}.$$

A simple test to verify if the residuals follow an AR(1) process is to test  $H_0 : \rho = 0$  in the regression (11). However the best known test, albeit quite dated, is the  $d$ -test proposed by Durbin and Watson, also known as the DW statistic:

$$d = \frac{\sum_{t=2}^n (\hat{u}_t - \hat{u}_{t-1})^2}{\sum_{t=1}^n \hat{u}_t^2} = \frac{n^{-1} \hat{u}_1^2 + 2n^{-1} \hat{\mathbf{u}}^\top \hat{\mathbf{u}}_1}{n^{-1} \hat{\mathbf{u}}^\top \hat{\mathbf{u}}} \quad (12)$$

Asymptotically, one can show that the  $d$ -test tends to:  $2 - 2\rho$ . Hence a value of  $d \cong 2$  corresponds to the absence of serial correlation in the residuals, while values

of  $d$  less than 2 corresponds to  $\hat{\rho} > 0$ , and values of greater than 2 correspond to  $\hat{\rho} < 0$ .

Note: The DW statistic is not valid when there are lagged dependent variables among the regressors.

### 2.5.2 Moving average errors

Autoregressive processes are not the only way to model stationary time series. Another type of stochastic process the moving average (MA) process. The simplest of these is the first order moving average (MA(1)) process

$$u_t = \epsilon_t + \alpha_1 \epsilon_{t-1}, \quad \epsilon_t \sim \text{IID}(0, \sigma_\epsilon^2) \quad (13)$$

in which the error term  $u_t$  is a weighted average of two successive innovations,  $\epsilon_t$  and  $\epsilon_{t-1}$ . This is also a covariance-stationary process:

$$\begin{aligned} \text{Var}(u_t) &= \text{Var}(\epsilon_t + \alpha_1 \epsilon_{t-1}) = \sigma_\epsilon^2 + \alpha_1^2 \sigma_\epsilon^2 = (1 + \alpha_1^2) \sigma_\epsilon^2 \\ \mathbb{E}[u_t u_{t-1}] &= \mathbb{E}[(\epsilon_t + \alpha_1 \epsilon_{t-1})(\epsilon_{t-1} + \alpha_1 \epsilon_{t-2})] = \alpha_1 \sigma_\epsilon^2 \\ \mathbb{E}[u_t u_{t-2}] &= \mathbb{E}[(\epsilon_t + \alpha_1 \epsilon_{t-1})(\epsilon_{t-2} + \alpha_1 \epsilon_{t-3})] = 0. \end{aligned}$$

The variance-covariance matrix  $\mathbf{\Omega}$  is

$$\mathbf{\Omega} = \sigma_\epsilon^2 \begin{bmatrix} (1 + \alpha_1^2) & \alpha_1 & 0 & \cdots & 0 \\ \alpha_1 & \ddots & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \ddots & \alpha_1 \\ 0 & \cdots & 0 & \alpha_1 & (1 + \alpha_1^2) \end{bmatrix},$$

where it is evident that there is no correlation between error terms which are more than one period apart.

## 2.6 Panel data

Many data sets are measured across two dimensions. One dimension is time, and the other is usually called the cross-section dimension. Data of this type are often referred to as panel data. It is likely that the error terms for a model using panel data will display certain types of dependence, which should be taken into account when we estimate such a model. Consider the following linear regression model:

$$y_{it} = \mathbf{X}_{it} \boldsymbol{\beta} + u_{it}, \quad i = 1, \dots, m, \quad t = 1, \dots, T, \quad (14)$$

where  $\mathbf{X}_{it}$  is a  $1 \times k$  vector of observations on explanatory variables. There is assumed to be  $m$  cross-sectional units and  $T$  time periods, for a total of  $n = m \times T$

observations. If each  $u_{it}$  has expectation zero conditional on its corresponding  $\mathbf{X}_{it}$ , we can estimate (14) via OLS. But the OLS estimator is not efficient if the  $u_{it}$  are not IID, and the IID assumption is rarely realistic with panel data. For example, for (14) consider the following

$$u_{it} = v_i + e_t + \epsilon_{it},$$

$$\begin{aligned}\mathbb{E}[v_i^2] &= \sigma_v^2, \quad \mathbb{E}[v_i v_j] = 0, \\ \mathbb{E}[e_t^2] &= \sigma_e^2, \quad \mathbb{E}[e_t e_{t-s}] = 0 \\ \mathbb{E}[\epsilon_{it}] &= \sigma_\epsilon^2, \quad \mathbb{E}[\epsilon_{it} \epsilon_{jt-s}] = 0.\end{aligned}$$

The above specification for the error term is an example of an ‘error-components model’, where the error term  $u_{it}$  consists of two or three independent shocks. In order to estimate an error-components model, the  $e_t$  and  $v_i$  can be regarded as being either fixed or random. If the  $e_t$  and  $v_i$  are thought of as fixed effects, then they are treated as parameters to be estimated, and it turns out that they can then be estimated by OLS using dummy variables.

If they are thought of as random effects, then we must figure out the covariance matrix of the  $u_{it}$  as functions of the variances of the  $e_t$ ,  $v_i$ , and  $\epsilon_{it}$ , and use FGLS. Each of these approaches can be appropriate in some circumstances but may be inappropriate in others.

For simplicity, let us eliminate  $e_t$  and consider  $u_{it} = v_i + \epsilon_{it}$ . In matrix notation, the linear panel data model is

$$\begin{bmatrix} y_{11} \\ \vdots \\ y_{1T} \\ y_{21} \\ \vdots \\ y_{2T} \\ \vdots \\ y_{m1} \\ \vdots \\ y_{mT} \end{bmatrix}_{n \times 1} = \begin{bmatrix} \mathbf{X}_{11} \\ \vdots \\ \mathbf{X}_{1T} \\ \mathbf{X}_{21} \\ \vdots \\ \mathbf{X}_{2T} \\ \vdots \\ \mathbf{X}_{m1} \\ \vdots \\ \mathbf{X}_{mT} \end{bmatrix}_{n \times k} \boldsymbol{\beta} + \begin{bmatrix} 1 & 0 & \cdots & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & \cdots & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 1 & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & \cdots & \cdots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & \cdots & 0 & 1 \end{bmatrix}_{n \times m} \boldsymbol{\eta} + \begin{bmatrix} \epsilon_{11} \\ \vdots \\ \epsilon_{1T} \\ \epsilon_{21} \\ \vdots \\ \epsilon_{2T} \\ \vdots \\ \epsilon_{m1} \\ \vdots \\ \epsilon_{mT} \end{bmatrix}_{n \times 1}, \quad \boldsymbol{\eta} = \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix},$$

$$\Leftrightarrow \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \underbrace{\mathbf{D}\boldsymbol{\eta}}_{\mathbf{u}} + \boldsymbol{\epsilon}, \quad \mathbb{E}[\boldsymbol{\epsilon}\boldsymbol{\epsilon}^\top] = \sigma_\epsilon^2 \mathbf{I}_n, \quad (15)$$

where  $\mathbf{y}$  and  $\boldsymbol{\epsilon}$  are  $n$ -vectors with typical elements  $y_{it}$  and  $\epsilon_{it}$ , respectively, and  $\mathbf{D}$  is an  $n \times m$  matrix of dummy variables, constructed in such a way that the

element in the row corresponding to observation  $it$ , for  $i = 1, \dots, m$  and  $t = 1, \dots, T$ , and column  $j$ , for  $j = 1, \dots, m$ , is equal to 1 if  $i = j$  and equal to 0 otherwise. The  $m$ -vector  $\boldsymbol{\eta}$  has typical element  $v_i$ , and so it follows that the  $n$ -vector  $\mathbf{D}\boldsymbol{\eta}$  has element  $v_i$  in the row corresponding to observation  $it$ . Note that there is exactly one element on  $\mathbf{D}$  equal to 1 in each row, which implies that the  $n$ -vector  $\boldsymbol{\iota}$  with each element equal to 1 is a linear combination of the columns of  $\mathbf{D}$ . Consequently, in order to avoid collinear regressors, the matrix  $\mathbf{X}$  should not contain a constant.

The vector  $\boldsymbol{\eta}$  essentially plays the role of a parameter vector, and it is in this sense the  $v_i$  are called fixed effects (though they can be random, as we will later show). The essential thing is that they must be independent of the error terms,  $\epsilon_{it}$ . They may, however, be correlated with the explanatory variables,  $\mathbf{X}_{it}$ .

We are going to use this specification to investigate two cases:

1. Fixed-effect estimation.
2. Random-effect estimation.

### 2.6.1 Fixed effect estimation

Also known as the “within-group” estimator or dummy variable estimator, because we assume that  $e_t = 0$  (no time component shock). Let us assume that the moment conditions

$$\begin{aligned}\mathbb{E}[y_{it} - \mathbf{X}_{it}\boldsymbol{\beta} - v_i] &= 0, \\ \mathbb{E}[\mathbf{X}_{it}^\top (y_{it} - \mathbf{X}_{it}\boldsymbol{\beta} - v_i)] &= 0,\end{aligned}$$

are satisfied, which is similar to the OLS regression model assumptions. Note that we are allowing that  $\mathbb{E}[\mathbf{X}_{it}^\top v_i] \neq 0$ . By applying the FWL theorem, we have

$$\hat{\boldsymbol{\beta}}_{\text{WG}} = (\mathbf{X}^\top \mathbf{M}_{\mathbf{D}} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}_{\mathbf{D}} \mathbf{y}, \quad (16)$$

where

$$\mathbf{M}_{\mathbf{D}} = \mathbf{I}_n - \mathbf{P}_{\mathbf{D}} = \mathbf{I}_n - \mathbf{D}(\mathbf{D}^\top \mathbf{D})^{-1} \mathbf{D}^\top.$$

One can show that (with some butchering of notation):

$$\mathbf{M}_D \mathbf{X} = \begin{bmatrix} \mathbf{X}_{11} - \bar{\mathbf{X}}_1 \\ \vdots \\ \mathbf{X}_{1T} - \bar{\mathbf{X}}_1 \\ \mathbf{X}_{21} - \bar{\mathbf{X}}_2 \\ \vdots \\ \mathbf{X}_{2T} - \bar{\mathbf{X}}_2 \\ \vdots \\ \mathbf{X}_{m1} - \bar{\mathbf{X}}_m \\ \vdots \\ \mathbf{X}_{mT} - \bar{\mathbf{X}}_m \end{bmatrix}, \quad \bar{\mathbf{X}}_i = \frac{1}{T} \sum_{t=1}^T \mathbf{X}_{it}.$$

Since all the variables in (16) are premultiplied by  $\mathbf{M}_D$ , it follows that this estimator makes use only of the information in the variation around the mean for each of the  $m$  groups. Hence,  $\hat{\beta}_{WG}$  is called the within-groups estimator.

Since the OLS and fixed-effect estimation assumptions are the same, the  $\hat{\beta}_{WG}$  is BLUE.

$\hat{\beta}_{WG}$  has advantages and disadvantages. It is easy to compute, even when  $m$  is large, because it is never necessary to make direct use of the  $n \times n$  matrix  $\mathbf{M}_D$ . We just need to compute the  $m$  group means for each variable. Also, the estimates of  $\hat{\eta}$  may themselves be of interest.

However, the estimator cannot be used with explanatory variables that take on the same value for all the observations of each group, because such a column would be collinear with the columns of  $\mathbf{D}$ . More generally, if the explanatory variables in  $\mathbf{X}$  are well explained by the dummy variables in  $\mathbf{D}$ , then  $\beta$  is not estimated at all precisely.

### 2.6.2 Random effects estimation

The random-effects estimator uses the following assumption.  $\mathbf{X}$  and the cross-sectional errors,  $v_i$ , should both be independent of the  $\epsilon_{it}$ , but this does not rule out the possibility of a correlation between them. Thus,

$$\mathbb{E}[u_{it}|\mathbf{X}] = \mathbb{E}[v_i + \epsilon_{it}|\mathbf{X}] = 0, \quad (17)$$

i.e., both  $v_i$  and  $\epsilon_{it}$  are independent of  $\mathbf{X}$ . Condition (17) is precisely the condition which ensures that OLS estimation of the model (14),

$$y_{it} = \mathbf{X}_{it}\beta + u_{it},$$

yields unbiased estimates, rather than the model (15),

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

But the estimation of (14) is, in general, not efficient, because the  $u_{it}$  are not IID. We can calculate the variance-covariance matrix of the  $u_{it}$  by assuming that the  $v_i$  are IID  $\sim (0, \sigma_v^2)$ . This is why we call this “random effects” estimation.

Remember that we have  $e_t = 0$  for all  $t$ , so:

$$\begin{aligned}\text{Var}(u_{it}) &= \text{Var}(v_i) + \text{Var}(\epsilon_{it}) = \sigma_v^2 + \sigma_\epsilon^2, \\ \text{Cov}(u_{it}, u_{it-s}) &= \sigma_v^2, \forall s, \\ \text{Cov}(u_{it}, u_{js}) &= 0, \forall i \neq j.\end{aligned}$$

Therefore the  $n \times n$  variance-covariance matrix,  $\mathbb{E}[\mathbf{uu}^\top | \mathbf{X}] = \boldsymbol{\Omega}$ , where the data is ordered by the cross-sectional units, is

$$\boldsymbol{\Omega} = \begin{bmatrix} \boldsymbol{\Sigma} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \boldsymbol{\Sigma} \end{bmatrix},$$

where

$$\boldsymbol{\Sigma} = \sigma_v^2 \boldsymbol{\nu}_T \boldsymbol{\nu}_T^\top + \sigma_\epsilon^2 \mathbf{I}_T,$$

is a  $T \times T$  matrix with  $\sigma_v^2 + \sigma_\epsilon^2$  along its diagonal and  $\sigma_v^2$  everywhere else. Here,  $\boldsymbol{\nu}_T$  is the  $T$ -unit vector. One can show that

$$\boldsymbol{\Sigma}^{-\frac{1}{2}} = \frac{1}{\sigma_\epsilon} (\mathbf{I}_T - \lambda \mathbf{P}_\boldsymbol{\nu}),$$

where

$$\lambda = 1 - \left( T \frac{\sigma_v^2}{\sigma_\epsilon^2} + 1 \right)^{-\frac{1}{2}}.$$

Also, note that  $\mathbf{P}_\boldsymbol{\nu}$  is the diagonal of  $\mathbf{P}_\mathbf{D}$ . Since  $\frac{1}{\sigma_\epsilon}$  is a scalar, we need to just find an estimator for the transformation matrix

$$\mathbf{I}_n - \lambda \mathbf{P}_\mathbf{D},$$

which depends on  $\sigma_\epsilon^2$  and  $\sigma_v^2$ .

To get the GLS estimates of  $\boldsymbol{\beta}$ , we would need to know the values of  $\sigma_\epsilon^2$  and  $\sigma_v^2$ , or, at least, the value of their ratio. This is because GLS estimation

requires only that  $\Omega$  be specified up to a factor. To get FGLS estimates, we need a consistent estimate of that ratio. Notice that we haven't discussed too much about asymptotic concepts thus far. This is because in order to get definite results, we need to specify what happens to both  $m$  and  $T$  when  $n = mT \rightarrow \infty$ .

Consider the fixed-effects model, (15). If  $m$  remains fixed as  $T \rightarrow \infty$ , then the number of regressors also remains fixed as  $n \rightarrow \infty$ , and standard asymptotic theory applies. But if  $T$  remains fixed as  $m \rightarrow \infty$ , then the number of parameters to be estimated tends to infinity, and the  $m$ -vector,  $\hat{\eta}$ , of estimates of the fixed effects is not consistent.<sup>1</sup>

But it is always possible to find a consistent estimate of  $\sigma_\epsilon^2$  by estimating the fixed-effects model, (15). This is because no matter how  $m$  and  $T$  may behave as  $n \rightarrow \infty$ , there are  $n$  residuals. Thus, if we divide the SSR from (15) by  $n - m - k$ , we obtain an unbiased and consistent estimate of  $\sigma_\epsilon^2$ , since the error terms of this model are just the  $\epsilon_{it}$ .

So, let us first estimate  $\sigma_\epsilon^2$ . Our linear regression model is

$$\mathbf{y} = \mathbf{X}\beta + \underbrace{\mathbf{D}\eta}_{\mathbf{u}} + \epsilon,$$

and via the FWL theorem the OLS residuals for  $\epsilon$  from the above regression and from

$$\mathbf{M}_D \mathbf{y} = \mathbf{M}_D \mathbf{X}\beta + \mathbf{M}_D \epsilon,$$

are numerically the same. So we can compute  $\hat{\epsilon}$  as

$$\hat{\epsilon} = \mathbf{M}_D \mathbf{y} - \mathbf{M}_D \mathbf{X} \hat{\beta}_{\text{FE}},$$

and the estimator for  $\sigma_\epsilon^2$  is:

$$\hat{\sigma}_\epsilon^2 = \frac{\hat{\epsilon}^\top \hat{\epsilon}}{n - (k + m)}. \quad (18)$$

But getting an estimate for  $\sigma_v^2$  is not as simple. The natural estimator of  $\sigma_v^2$ , namely the sample variance of the  $m$  elements of  $\hat{\eta}$ , is not consistent unless  $m \rightarrow \infty$ . Therefore, we should not use the random-effects estimator if  $m$  is small. For estimating  $\sigma_v^2$  we will need the auxiliary regression

$$\mathbf{P}_D \mathbf{y} = \mathbf{P}_D \mathbf{X}\beta + \underbrace{\mathbf{P}_D \mathbf{u}}_{\mathbf{r}},$$

---

<sup>1</sup>Interestingly, one can show that  $\hat{\beta}$  is still consistent under these circumstances, however. Sufficient conditions for this are that  $\text{plim}_{m \rightarrow \infty} m^{-1} \sum_{i=1}^m \mathbf{X}_i^\top \mathbf{M}_i \mathbf{X}_i$  is positive definite (for invertibility) and that  $\text{plim}_{m \rightarrow \infty} m^{-1} \sum_{i=1}^m \mathbf{X}_i^\top \mathbf{M}_i \epsilon_i = \mathbf{0}$ .



where  $\mathbf{P}_D \equiv \mathbf{I} - \mathbf{M}_D$ , so as to obtain the “between-groups” estimator

$$\hat{\beta}_{BG} = (\mathbf{X}^\top \mathbf{P}_D \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_D \mathbf{y}. \quad (19)$$

Although  $\mathbf{P}_D \mathbf{u}$  is  $n \times 1$  in dimension (remember,  $n = mT$ ), we really only have  $m$  observations:

$$\mathbf{P}_D \mathbf{y} = \begin{bmatrix} \bar{\mathbf{u}}_1 \\ \vdots \\ \bar{\mathbf{u}}_1 \\ \bar{\mathbf{u}}_2 \\ \vdots \\ \bar{\mathbf{u}}_2 \\ \vdots \\ \bar{\mathbf{u}}_m \\ \vdots \\ \bar{\mathbf{u}}_m \end{bmatrix}, \quad \bar{\mathbf{u}}_i = \frac{1}{T} \sum_{t=1}^T u_{it},$$

i.e., this is because the regressand and all the regressors are the same for every observation in each group. The estimator  $\hat{\beta}_{BG}$  uses only the variation among the group means, hence why it is termed as the between-groups estimator. Note that if  $m < k$ , then (19) does not exist, since  $\mathbf{X}^\top \mathbf{P}_D \mathbf{X}$  can have rank at most  $m$ .

Next,

$$\text{Var}(\bar{\mathbf{u}}_i) = \text{Var}\left(\frac{1}{T} \sum_{t=1}^T v_i + \frac{1}{T} \sum_{t=1}^T \epsilon_{it}\right) = \sigma_v^2 + \frac{1}{T} \sigma_\epsilon^2 = \sigma_r^2,$$

hence we can estimate  $\sigma_r^2$  by

$$\hat{\sigma}_r^2 = \frac{\hat{\mathbf{r}}^\top \hat{\mathbf{r}}}{(m - k)},$$

where

$$\hat{\mathbf{r}} = \mathbf{P}_D \mathbf{y} - \mathbf{P}_D \mathbf{X} \hat{\beta}_{BG}.$$

We note that in this case that the OLS estimator is a matrix-weighted average estimator:

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top (\mathbf{P}_D + \mathbf{M}_D) \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{P}_D \mathbf{X} \hat{\beta}_{BG} + (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{M}_D \mathbf{X} \hat{\beta}_{WG}. \end{aligned} \quad (20)$$

Finally, we estimate  $\sigma_v^2$  as

$$\hat{\sigma}_v^2 = \hat{\sigma}_r^2 - \frac{1}{T} \hat{\sigma}_\epsilon^2,$$

and the FGLS “random effects” estimator is the OLS estimator of the following transformed regression

$$(\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D)\mathbf{y} = (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D)\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D)\mathbf{u}.$$

In other words

$$\hat{\boldsymbol{\beta}}_{\text{RE}} = \left( \mathbf{X}^\top (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D)^\top (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D) \mathbf{X} \right)^{-1} \mathbf{X}^\top (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D)^\top (\mathbf{I}_n - \hat{\lambda}\mathbf{P}_D) \mathbf{y}, \quad (21)$$

where

$$\hat{\lambda} = 1 - \left( T \frac{\hat{\sigma}_v^2}{\hat{\sigma}_\epsilon^2} + 1 \right)^{-\frac{1}{2}}.$$