Linear Regression

Nuno Carvalho

June 3, 2025

1 Generalized Least Squares (GLS)

Largely adapted from the respective Wikipedia article.

1.1 Model definition

In Generalized Least Squares (GLS), we define an outcome's mean to be a linear function of a set of predictors:

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

Where:

- y: outcome vector. $N \times 1$
 - -N: number of samples
- **X**: design matrix. $N \times K$ matrix
 - K: number of predictors
- $-\beta$: coefficients vector. $K \times 1$ vector
- $-\epsilon$: error term. $N \times 1$ vector

Because of our conditional mean definition earlier, the error term, ϵ , has a mean of zero for a given set of predictor values (X):

$$E[\boldsymbol{\epsilon}|\mathbf{X}] = \mathbf{0}$$

We also assume that the variance of the error term given **X** is described by an invertible $N \times N$ covariance matrix, Ω :

$$\operatorname{Cov}[\boldsymbol{\epsilon}|\mathbf{X}] = \mathbf{\Omega}$$

We further assume that the error term follows a multivariate normal distribution with mean 0 and covariance Ω :

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{\Omega})$$

And so it follows that:

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Omega})$$

1.2 Least squares estimate of β

We denote a candidate estimate for the β vector as **b** and define its residual vector as $\mathbf{y} - \mathbf{X}\mathbf{b}$. The goal of GLS is to find the estimate of β that maximizes the likelihood of the data given the above model, which we can calculate using the probability density function of a multivariate normal distribution:

$$\frac{1}{\sqrt{(2\pi)^K |\mathbf{\Sigma}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

We can plug in $\Sigma = \Omega$, $\mathbf{x} = \mathbf{y}$, and $\boldsymbol{\mu} = \mathbf{X}\mathbf{b}$. The **b** that maximizes this likelihood function will be the same that minimizes the inside of the exponent (without the negative), which is also the squared Mahalanobis length of the residual vector:

$$\hat{\boldsymbol{\beta}} = \operatorname*{arg\,min}_{b} \, \left(\mathbf{y} - \mathbf{X} \mathbf{b} \right)^{\intercal} \boldsymbol{\Omega}^{-1} (\mathbf{y} - \mathbf{X} \mathbf{b})$$

Through matrix algebra, this is equivalent to:

$$\hat{\boldsymbol{\beta}} = \underset{b}{\operatorname{arg \, min}} \ (\mathbf{y} - \mathbf{X}\mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} (\mathbf{y} - \mathbf{X}\mathbf{b})$$

$$= \underset{b}{\operatorname{arg \, min}} \ \mathbf{y}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y} - \mathbf{y}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X}\mathbf{b} - (\mathbf{X}\mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y} + (\mathbf{X}\mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X}\mathbf{b}$$

$$= \underset{b}{\operatorname{arg \, min}} \ \mathbf{y}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y} + (\mathbf{X}\mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} (\mathbf{X}\mathbf{b}) - 2(\mathbf{X}\mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y}$$

We can use calculus to solve for the \mathbf{b} that minimizes this expression by taking the partial derivative with respect to \mathbf{b} and solving for 0:

$$\mathbf{0} = \frac{\partial}{\partial \mathbf{b}} [\mathbf{y}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y} + (\mathbf{X} \mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} (\mathbf{X} \mathbf{b}) - 2(\mathbf{X} \mathbf{b})^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y}]$$
$$= 2\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X} \hat{\boldsymbol{\beta}} - 2\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y}$$

Which yields:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y}$$

1.2.1 Specific case of Ordinary Least Squares

Note that in Ordinary Least Squares (OLS), the covariance matrix is an identity matrix, $\Omega = I$. That is, the residuals are uncorrelated with each other. This simplifies the above equation to:

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{y}$$
$$= (\mathbf{X}^{\mathsf{T}} I \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} I^{-1} \mathbf{y}$$
$$= (\mathbf{X}^{\mathsf{T}} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{y}$$

1.3 Variance of $\hat{\beta}$

To get the variance of the $\hat{\boldsymbol{\beta}}$ estimate, we only need to focus on the variance of \mathbf{y} , as all other terms are not random variables. $\operatorname{Var}[\mathbf{y}] = \mathbf{\Omega}$ since $\boldsymbol{\epsilon}$ is independent of $\mathbf{X}\boldsymbol{\beta}$, as shown below:

$$\begin{aligned} \operatorname{Var}[\mathbf{y}] &= \operatorname{Var}[\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}] \\ &= \operatorname{Var}[\mathbf{X}\boldsymbol{\beta}] + \operatorname{Var}[\boldsymbol{\epsilon}] \\ &= \mathbf{0} + \operatorname{E}[\operatorname{Var}[\boldsymbol{\epsilon}|\mathbf{X}]] + \operatorname{Var}[\operatorname{E}[\boldsymbol{\epsilon}|\mathbf{X}]] \\ &= \operatorname{E}[\boldsymbol{\Omega}] + \operatorname{Var}[\mathbf{0}] \\ &= \boldsymbol{\Omega} \end{aligned}$$

Let's define $\mathbf{A} = (\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X})^{-1} \mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1}$, such that $\hat{\boldsymbol{\beta}} = \mathbf{A} \mathbf{y}$. Furthermore, note that since $\mathbf{\Omega}$ is a covariance matrix, it (and its inverse) is symmetric: $\mathbf{\Omega} = \mathbf{\Omega}^{\mathsf{T}}$. The same symmetry property applies to $\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X}$. We can then solve:

$$\begin{aligned} \operatorname{Var}[\hat{\boldsymbol{\beta}}] &= \operatorname{Var}[(\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{y}] \\ &= \operatorname{Var}[\mathbf{A}\mathbf{y}] \\ &= \mathbf{A}\operatorname{Var}[\mathbf{y}]\mathbf{A}^{\mathsf{T}} \\ &= \mathbf{A}\boldsymbol{\Omega}\mathbf{A}^{\mathsf{T}} \\ &= (\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{\Omega}((\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1})^{\mathsf{T}} \\ &= (\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\boldsymbol{\Omega}\boldsymbol{\Omega}((\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1})^{\mathsf{T}} \\ &= (\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1}(\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})(\mathbf{X}^{\mathsf{T}}\boldsymbol{\Omega}^{-1}\mathbf{X})^{-1} \end{aligned}$$

We can define $\mathbf{B} = \mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X}$, allowing us to simplify the above equation, $\mathbf{B}^{-1} \mathbf{B} \mathbf{B}^{-1} = \mathbf{B}^{-1}$, yielding:

$$\mathrm{Var}[\hat{\boldsymbol{\beta}}] = (\mathbf{X}^{\mathsf{T}} \mathbf{\Omega}^{-1} \mathbf{X})^{-1}$$