

Gauss-Markov Theorem for OLS is the best linear unbiased estimator

Andersen Ang

First created: 2014. Last update : 2017-Feb

1 The regression problem

$$y = X\beta + \varepsilon$$

Gauss-Markov assumptions on ε

- $\mathbb{E}\varepsilon = 0$
- $\text{Cov}(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij}$

It means, the noise is zero-mean

The variance $\text{Var}(\varepsilon_i) = \text{Cov}(\varepsilon_i, \varepsilon_i) = \sigma^2 < \infty$, all noise have same variance.

The covariance $\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \forall i \neq j$, all noise are uncorrelated

Then under these assumptions, OLS estimator is the best linear unbiased estimator. That means, for all unbiased estimators for this regression, OLS estimator has the smallest variance.

Note. The assumptions only requires the noise to be zero-mean and $\text{Cov}(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij}$. It does not require the noise to be Gaussian noise.

2 The proof

Recall that $\hat{\beta}^{OLS} = (X^T X)^{-1} X^T y$, $\mathbb{E}[\hat{\beta}^{OLS}] = \beta_0$, $\text{Var}[\hat{\beta}^{OLS}] = \sigma^2 (X^T X)^{-1}$

Now let $\hat{\beta}^{Other} = ((X^T X)^{-1} X^T + D)y$ where D is a non-zero matrix.

Then

$$\mathbb{E}[\hat{\beta}^{Other}] = \mathbb{E}\left[\left((X^T X)^{-1} X^T + D\right)y\right]$$

$$\mathbb{E}[\hat{\beta}^{Other}] = \mathbb{E}\left[(X^T X)^{-1} X^T y + Dy\right]$$

$$\mathbb{E}[\hat{\beta}^{Other}] = \mathbb{E}\left[(X^T X)^{-1} X^T y\right] + \mathbb{E}[Dy]$$

$$\mathbb{E}[\hat{\beta}^{Other}] = \mathbb{E}[\hat{\beta}^{OLS}] + \mathbb{E}[D(X\beta_0 + \varepsilon)]$$

$$\mathbb{E}[\hat{\beta}^{Other}] = \beta_0 + \mathbb{E}[DX]\beta_0 + \underbrace{\mathbb{E}[D\varepsilon]}_{D\mathbb{E}[\varepsilon]=0}$$

$$\mathbb{E} [\hat{\beta}^{Other}] = [I + DX] \beta_0$$

So $\hat{\beta}^{Other}$ is unbiased iff $DX = 0$.

$$Var [\hat{\beta}^{Other}] = Var \left[\left((X^T X)^{-1} X^T + D \right) y \right]$$

$$Var [\hat{\beta}^{Other}] = \mathbb{E} \left[\left\{ \left((X^T X)^{-1} X^T + D \right) y \right\} \left\{ \left((X^T X)^{-1} X^T + D \right) y \right\}^T \right]$$

$$Var [\hat{\beta}^{Other}] = \mathbb{E} \left[\left((X^T X)^{-1} X^T + D \right) y y^T \left((X^T X)^{-1} X^T + D \right)^T \right]$$

$$Var [\hat{\beta}^{Other}] = \left((X^T X)^{-1} X^T + D \right) \mathbb{E} [y y^T] \left((X^T X)^{-1} X^T + D \right)^T$$

Since $\mathbb{E} [y y^T] = Var (\varepsilon) = \sigma^2 I$

$$Var [\hat{\beta}^{Other}] = \left((X^T X)^{-1} X^T + D \right) \sigma^2 I \left((X^T X)^{-1} X^T + D \right)^T$$

$$Var [\hat{\beta}^{Other}] = \sigma^2 \left((X^T X)^{-1} X^T + D \right) \left(\left((X^T X)^{-1} X^T \right)^T + D^T \right)$$

$$Var [\hat{\beta}^{Other}] = \sigma^2 \left[(X^T X)^{-1} X^T \left((X^T X)^{-1} X^T \right)^T + (X^T X)^{-1} X^T D^T + D \left((X^T X)^{-1} X^T \right)^T + D D^T \right]$$

Since $DX = 0$ so $X^T D^T = (DX)^T = 0^T = 0$

$$Var [\hat{\beta}^{Other}] = \sigma^2 \left[\left((X^T X)^{-1} \right)^T + (X^T X)^{-1} \underbrace{X^T D^T}_0 + \underbrace{DX}_0 \left((X^T X)^{-1} \right)^T + D D^T \right]$$

$$\left((X^T X)^{-1} \right)^T = \left((X^T X)^T \right)^{-1} = \left(X^T (X^T)^T \right)^{-1} = (X^T X)^{-1}$$

$$Var [\hat{\beta}^{Other}] = \sigma^2 \left[(X^T X)^{-1} + D D^T \right]$$

$$Var [\hat{\beta}^{Other}] = \underbrace{\sigma^2 (X^T X)^{-1}}_{Var [\hat{\beta}^{OLS}]} + \sigma^2 D D^T$$

$$Var [\hat{\beta}^{Other}] = Var [\hat{\beta}^{OLS}] + \sigma^2 D D^T$$

Since $DD^T > 0$ (DD^T is positive semi-definite matrix)

$$Var [\hat{\beta}^{Other}] \geq Var [\hat{\beta}^{OLS}]$$

Therefore the OLS estimator has the lowest variance, and hence it is the best estimator.