

Estimation of the noise variance in OLS

Andersen Ang

First created: 2014. Last update : 2017-Feb

1 The problem

For the regression problem

$$y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \varepsilon_{n \times 1}$$

$\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$, the noise is zero-mean, uncorrelated and having same variance σ^2 : $\mathbb{E}[\varepsilon_i \varepsilon_j] = \sigma^2 \delta_{ij}$. This is the Gaussian-Markov assumption.

We want to estimate β for the case $n = p$ by the Ordinary Least Square Estimator $\hat{\beta}^{OLS} = \arg \min \|y - X\beta\|_2^2$

it can be shown that $\hat{\beta}^{OLS} = (X^T X)^{-1} X^T y$ and it has the distribution $\hat{\beta}^{OLS} \sim \mathcal{N}(\beta_0, \sigma^2 (X^T X)^{-1})$ for $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$ and β_0 is the true value.

We know that the noise is zero mean , but we don't know the value of σ^2 , so we want to estimate it.

2 Estimation of σ^2

The error e (residue) can be expressed as

$$e = y - X\hat{\beta}$$

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

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

Let $G = \left(I - X (X^T X)^{-1} X^T \right)$, note that

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

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

By $(ABC)^T = A[BC]^T = [BC]^T A^T = C^T B^T A^T$

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

And

$$G^2 = \left(I - X (X^T X)^{-1} X^T \right) \left(I - X (X^T X)^{-1} X^T \right)$$

$$G^2 = I - X (X^T X)^{-1} X^T - X (X^T X)^{-1} X^T + X (X^T X)^{-1} \underbrace{X^T X (X^T X)^{-1}}_I X^T$$

$$G^2 = I - X (X^T X)^{-1} X^T \underbrace{-X (X^T X)^{-1} X^T + X (X^T X)^{-1} X^T}_0 = I - X (X^T X)^{-1} X^T = G$$

And

$$GX = \left(I - X (X^T X)^{-1} X^T \right) X = X - X = 0$$

So

G is symmetric, idempotent and orthogonal to X

Then

$$\begin{aligned} e &= Gy \\ e &= G(X\beta + \varepsilon) \\ e &= \underbrace{GX}_0 \beta + G\varepsilon \\ e &= \left(I - X (X^T X)^{-1} X^T \right) \varepsilon \end{aligned}$$

Therefore

$$\mathbb{E}[e] = \mathbb{E} \left[\left(I - X (X^T X)^{-1} X^T \right) \varepsilon \right] = \left(I - X (X^T X)^{-1} X^T \right) \underbrace{\mathbb{E}[\varepsilon]}_0 = 0$$

$$\text{Var}(e) = \mathbb{E}[ee^T] = \mathbb{E}[(G\varepsilon)(G\varepsilon)^T]$$

Since G is symmetric $G = G^T$

$$\text{Var}(e) = \mathbb{E}[G\varepsilon\varepsilon^T G] = G\mathbb{E}[\varepsilon\varepsilon^T]G$$

Since $\mathbb{E}[\varepsilon\varepsilon^T] = \sigma^2 I$ and $G^2 = G$

$$\text{Var}(e) = \sigma^2 G$$

Thus the theoretical equation to estimate σ^2 is thus

$$\mathbb{E}[ee^T] = \sigma^2 \left(I - X (X^T X)^{-1} X^T \right)$$

3 Practical way to estimate σ^2 in OLS

The above equation cannot be used directly in application since it is not “very useful”.

Consider $\mathbb{E}[ee^T]$, it is a square matrix, one way to compute this value is to use the trace

$$\mathbb{E}[ee^T] = \mathbb{E}[\text{Tr}(ee^T)]$$

Since trace and expectation operator can be interchanged

$$\mathbb{E}[ee^T] = \text{Tr}(\mathbb{E}[ee^T])$$

Since $\mathbb{E}[ee^T] = \sigma^2 G = \sigma^2 \left(I - X (X^T X)^{-1} X^T \right)$, thus

$$\mathbb{E}[ee^T] = \text{Tr}(\mathbb{E}[\sigma^2 G])$$

$$\mathbb{E}[ee^T] = \sigma^2 \text{Tr}(G)$$

$$\mathbb{E}[ee^T] = \sigma^2 \text{Tr} \left(I - X (X^T X)^{-1} X^T \right)$$

Since $\text{Tr}(A + B) = \text{Tr}(A) + \text{Tr}(B)$

$$\mathbb{E}[ee^T] = \sigma^2 \left(\text{Tr}I - \text{Tr} \left(X (X^T X)^{-1} X^T \right) \right)$$

Since $\text{Tr}(AB^T) = \text{Tr}(B^T A)$

$$\mathbb{E}[ee^T] = \sigma^2 \left(\text{Tr}I - \text{Tr} \left(\underbrace{X^T X (X^T X)^{-1}}_I \right) \right)$$

Notice that the first identity matrix is $n \times n$ and the second one is $p \times p$

$$\mathbb{E}[ee^T] = \sigma^2 (\text{Tr}I_{n \times n} - \text{Tr}I_{p \times p})$$

$$\mathbb{E}[ee^T] = \sigma^2 (n - p)$$

Therefore

$$\sigma^2 = \frac{\mathbb{E}[ee^T]}{n - p}$$