

Lecture 38: Estimating Variance

38.1

We know that $\hat{\beta} = (X^T X)^{-1} X^T Y$ gives an unbiased estimator for β_0 , where Y comes from the linear model $Y = X\beta_0 + \epsilon$,

with $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma_0^2 I$.

Our goals for this lecture are: - compute $\text{Var}(\hat{\beta})$
- estimate σ_0^2 .

Warm-up #1

Later it will prove helpful to work with the trace of a matrix, the sum of the diagonal entries.

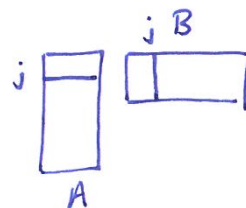
If $A \in \mathbb{R}^{n \times n}$, then

$$\text{trace}(A) = a_{1,1} + a_{2,2} + \dots + a_{n,n}.$$

One of the most useful properties of the trace is that you can commute matrices under the trace:

If $A \in \mathbb{R}^{n \times p}$ and $B \in \mathbb{R}^{p \times n}$, then

$$\begin{aligned} \text{trace}(AB) &= \sum_{j=1}^n (AB)_{j,j} \\ &= \sum_{j=1}^n \sum_{k=1}^p a_{j,k} b_{k,j} \\ &= \sum_{k=1}^p \sum_{j=1}^n b_{k,j} a_{j,k} \\ &= \sum_{k=1}^p (BA)_{k,k} = \text{trace}(BA). \end{aligned}$$



So $\text{trace}(AB) = \text{trace}(BA)$.

Warm-up #2

38

$\hat{q} = (X^T X)^{-1} X^T Y$: Let us focus on the matrix $(X^T X)^{-1} X^T$.

$$\begin{aligned} X^+ &= (X^T X)^{-1} X^T = \left(\begin{array}{|c|} \hline X^T \\ \hline X \\ \hline \end{array} \right)^{-1} \begin{array}{|c|} \hline X^T \\ \hline \end{array} \\ &= \begin{array}{|c|} \hline (X^T X)^{-1} \\ \hline \end{array} \begin{array}{|c|} \hline X^T \\ \hline \end{array} = \hat{p} \begin{array}{|c|} \hline X^+ \\ \hline \end{array} \end{aligned}$$

$X^+ \in \mathbb{R}^{p \times n}$ is the pseudoinverse of $X \in \mathbb{R}^{n \times p}$.

X^+ is a left-inverse of X : $(X^+ X = (X^T X)^{-1} X^T X = I \in \mathbb{R}^{p \times p})$

but it is not generally a right inverse (unless $n=p$):

$$H = X X^+ = X \underset{\substack{\uparrow \\ \text{rank } p}}{(X^T X)^{-1}} \underset{\substack{\uparrow \\ \text{rank } p}}{X^T} \in \mathbb{R}^{n \times n}$$

The product of two rank- p matrices cannot have rank larger than p . Since $H \in \mathbb{R}^{n \times n}$ is an $n \times n$ matrix of rank $\leq p$, it cannot equal $I \in \mathbb{R}^{n \times n}$ (a rank n matrix) unless $n=p$. (If $n=p$, then X is square and $X^+ = X^{-1}$. Usually we have $n \gg p$: many more observations than quantities of interest ("QoIs".))

H still has some nice properties:

$$H^2 = X \underbrace{X^+ X}_{=I} X^+ = X X^+ = H$$

So H is a projector. It is also symmetric:

$$H^T = (X (X^T X)^{-1} X^T)^T = X^{TT} ((X^T X)^{-1})^T X^T = X (X^T X)^{-1} X^T = H.$$

Here we have used the fact that $X^T X$ is symmetric and the inverse of a symmetric (invertible) matrix is also symmetric. (Here's a proof using "the spectral method" from the first part of the Semester. Let $\lambda_1, \dots, \lambda_p$ be the eigenvalues of $X^T X$, with orthonormal eigenvectors v_1, \dots, v_p :

$$X^T X = \begin{bmatrix} | & & | \\ v_1 & \dots & v_p \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_p \end{bmatrix} \begin{bmatrix} -v_1^T \\ \vdots \\ -v_p^T \end{bmatrix} = V \Lambda V^T$$

Then $(V \Lambda V^T)(V \Lambda^{-1} V^T) = V \Lambda \underbrace{V^T V}_{=I \text{ by orthonormality}} \Lambda^{-1} V^T = V \underbrace{\Lambda \Lambda^{-1}}_{=I} V^T = V V^T = I$

Hence $V \Lambda^{-1} V^T = (X^T X)^{-1}$. Now Λ^{-1} is diagonal.

$$[(X^T X)^{-1}]^T = (V \Lambda^{-1} V^T)^T = V (\Lambda^{-1})^T V^T = V \Lambda^{-1} V^T = (X^T X)^{-1}$$

so $(X^T X)^{-1}$ is symmetric.

Aside....

Variance of \hat{q}

Now we can determine how the noise ϵ that pollutes the observations Y filters through to affect the unbiased estimator \hat{q} .

$$\begin{aligned} \text{Var}(\hat{q}) &= \mathbb{E}((\hat{q} - \mathbb{E}(\hat{q}))(\hat{q} - \mathbb{E}(\hat{q}))^T) \\ &= \mathbb{E}((\hat{q} - q_0)(\hat{q} - q_0)^T) \end{aligned}$$

using the definition of variance of a vector variable

since $\mathbb{E}(\hat{q}) = q_0$.

Substitute in our formula for \hat{q} :

38.1

$$\begin{aligned}\hat{q} - q_0 &= X^+ Y - q_0 & Y &= X q_0 + \epsilon \\ &= X^+ (X q_0 + \epsilon) - q_0 \\ &= \underbrace{X^+ X}_{=I} q_0 + X^+ \epsilon - q_0 \\ &= q_0 + X^+ \epsilon - q_0 = X^+ \epsilon.\end{aligned}$$

Thus

$$\begin{aligned}\text{Var}(\hat{q}) &= \mathbb{E}((\hat{q} - q_0)(\hat{q} - q_0)^T) \\ &= \mathbb{E}((X^+ \epsilon)(X^+ \epsilon)^T) \\ &= \mathbb{E}(X^+ \epsilon \epsilon^T (X^+)^T) \\ &= X^+ \mathbb{E}(\epsilon \epsilon^T) (X^+)^T\end{aligned}$$

Where this last step used the linearity of the expected value and the fact that X is a constant matrix.

$$\text{Now } \mathbb{E}(\epsilon \epsilon^T) = \mathbb{E}((\epsilon - \underbrace{\mathbb{E}(\epsilon)}_{=0})(\epsilon - \underbrace{\mathbb{E}(\epsilon)}_{=0})^T) = \text{Var}(\epsilon)$$

$$\begin{aligned}\text{Thus } \text{Var}(\hat{q}) &= X^+ (\sigma_0^2 I) (X^+)^T & &= \sigma_0^2 I. \\ &= \sigma_0^2 X^+ (X^+)^T\end{aligned}$$

$$\begin{aligned}\text{Look at } X^+ (X^+)^T &= (X^T X)^{-1} X^T ((X^T X)^{-1} X^T)^T \\ &= (X^T X)^{-1} X^T \underbrace{X^T}_{X} \underbrace{((X^T X)^{-1})^T}_{\text{Symmetric}} \\ &= (X^T X)^{-1} X^T X (X^T X)^{-1} \\ &= (X^T X)^{-1}.\end{aligned}$$

Thus we can conclude

38.5

$$\text{Var}(\hat{q}) = \sigma_0^2 (X^T X)^{-1}$$

↑ ↑
Var(ϵ) = $\sigma_0^2 \mathbb{I}$ scaling factor.

So $(X^T X)^{-1}$ is a "magnification factor" that describes how much the Variance of the noise ϵ affects \hat{q} .

Unbiased estimate for σ_0^2 .

Now we would like to estimate this variance factor σ_0^2 . Our route to this estimate travels through the residual vector $\hat{R} = Y - X\hat{q}$, which is the mismatch between our observations Y and our unbiased estimate $X\hat{q}$ for these observations.

$$\begin{aligned}\hat{R} &= Y - X\hat{q} = Y - X X^+ Y \\ &= (I - X X^+) Y \\ &= (I - H) Y \\ &= (I - H)(X q_0 + \epsilon) \\ &= (I - H) X q_0 + (I - H) \epsilon\end{aligned}$$

Note that

$$\begin{aligned}(I - H) X q_0 &= X q_0 - H X q_0 = X q_0 - X \underbrace{(X^T X)^{-1} X^T X}_{=I} q_0 \\ &= X q_0 - X q_0 = 0.\end{aligned}$$

Hence

38.1

$$\hat{R} = \underbrace{(I-H)Xq_0}_{=0} + (I-H)\varepsilon = (I-H)\varepsilon.$$

Thus, we can compute an expected value for the square of the norm of the mismatch:

$$\begin{aligned} \mathbb{E}(\|\hat{R}\|^2) &= \mathbb{E}(\hat{R}^T \hat{R}) \\ &= \mathbb{E}(((I-H)\varepsilon)^T ((I-H)\varepsilon)) \\ &= \mathbb{E}(\varepsilon^T \underbrace{(I-H)^T (I-H)}_{=I-H \quad (H=H^T)} \varepsilon) \\ &= \mathbb{E}(\varepsilon^T \underbrace{(I-H)(I-H)}_{(I-H)(I-H) = I-H-H+H^2 = I-H-H+H = I-H} \varepsilon) \\ &= \mathbb{E}(\varepsilon^T (I-H)\varepsilon) \end{aligned}$$

Let $B = I - H \in \mathbb{R}^{n \times n}$. Then

$$\varepsilon^T (I-H)\varepsilon = \varepsilon^T B \varepsilon = \sum_{j=1}^n \sum_{k=1}^n \varepsilon_j b_{j,k} \varepsilon_k$$

So by linearity of the expected value,

$$\begin{aligned} \mathbb{E}(\|\hat{R}\|^2) &= \mathbb{E}(\varepsilon^T B \varepsilon) \\ &= \sum_{j=1}^n \sum_{k=1}^n b_{j,k} \mathbb{E}(\varepsilon_j \varepsilon_k) \\ &= \sum_{j=1}^n \sum_{k=1}^n b_{j,k} \mathbb{E}(\varepsilon_j \varepsilon_k - \underbrace{\mathbb{E}(\varepsilon_j)}_{=0} \underbrace{\mathbb{E}(\varepsilon_k)}_{=0}) \\ &= \sum_{j=1}^n \sum_{k=1}^n b_{j,k} \text{Cov}(\varepsilon_j, \varepsilon_k) \end{aligned}$$

since $\mathbb{E}(\varepsilon) = 0$

Recall that $\text{Var}(\varepsilon) = \sigma_0^2 I$

38.7

which means $\text{Cov}(\varepsilon_j, \varepsilon_k) = \begin{cases} \sigma_0^2, & j=k \\ 0, & j \neq k \end{cases}$

$$\begin{aligned} \text{Hence } E(\|\hat{R}\|^2) &= \sum_{j=1}^n \sum_{k=1}^n b_{jik} \text{Cov}(\varepsilon_j, \varepsilon_k) \\ &= \sum_{j=1}^n b_{jij} \sigma_0^2 \\ &= \sigma_0^2 \text{trace}(B). \end{aligned}$$

$$\begin{aligned} \text{Now } \text{trace}(B) &= \text{trace}(I - H) \\ &= \underbrace{\text{trace}(I)}_{I \in \mathbb{R}^{n \times n}} - \underbrace{\text{trace}(H)}_{H = XX^T} \end{aligned} \quad \left. \vphantom{\begin{aligned} \text{Now } \text{trace}(B) &= \text{trace}(I - H) \\ &= \text{trace}(I) - \text{trace}(H) \end{aligned}} \right\} \begin{array}{l} \text{trace} = \text{sum of} \\ \text{diagonals is} \\ \text{linear} \end{array}$$
$$\begin{aligned} &= n - \text{trace}(XX^T) \\ &= n - \text{trace}(X^T X) \end{aligned} \quad \begin{array}{l} \text{KEY: use Worm-up \#1} \\ \\ \end{array}$$
$$\begin{aligned} &= n - p \\ &= n - p \end{aligned}$$

$$\begin{array}{c} \boxed{X^T} \quad \boxed{X} \\ \begin{array}{c} 1 \quad p \\ \hline n \end{array} \end{array} = X^T X \in \mathbb{R}^{p \times p}$$

$$\begin{aligned} \text{Thus } E(\|\hat{R}\|^2) &= \sigma_0^2 \text{trace}(B) \\ &= \sigma_0^2 (n-p). \end{aligned}$$

$$\text{Define } \hat{\sigma}_0^2 = \frac{\|\hat{R}\|^2}{n-p}.$$

$$\text{Then } E(\hat{\sigma}_0^2) = \frac{E(\|\hat{R}\|^2)}{n-p} = \frac{\sigma_0^2 (n-p)}{n-p} = \sigma_0^2.$$

38.8

Thus $\hat{\sigma}_0^2$ is an unbiased estimator for σ_0^2 .

In summary:

① $\hat{\beta} = X^+ Y = (X^T X)^{-1} X^T Y =$ unbiased estimator for β_0 .

② $\text{Var}(\hat{\beta}) = \sigma_0^2 (X^T X)^{-1}$

③ $\hat{\sigma}_0^2 = \frac{\|\hat{R}\|^2}{n-p} = \frac{\|Y - X\hat{\beta}\|^2}{n-p} =$ unbiased estimator for σ_0^2 .

See `ls_demo1.m` and `ls_demo2.m` on the website for examples with

$$X = \begin{bmatrix} 1 & t_1 & t_1^2 \\ \vdots & \vdots & \vdots \\ 1 & t_n & t_n^2 \end{bmatrix} \in \mathbb{R}^{n \times 3}$$

One last question: What if $n=p$? Does ③ break down?

In this case, $X^+ = X^{-1}$, so $\hat{\beta} = X^{-1} Y$, and

$$\hat{R} = Y - X X^{-1} Y = 0, \text{ so } \textcircled{3} \text{ reduces to a } 0/0$$

Indeterminate form; we learn nothing about σ_0^2

Since $\mathbb{E}(\|\hat{R}\|^2) = \mathbb{E}(0) = 0$.