

Linear Regression

- Let $f(\tilde{\mathbf{x}}; \tilde{\boldsymbol{\theta}}) = \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}}$, where $\tilde{\mathbf{X}} \in \mathbb{R}^{N_y \times P}$ + $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^P$.

- Let $\tilde{\boldsymbol{\theta}}^0$ be the true, unknown value of parameters.

- Lets assume the following:

$$\left. \begin{array}{l} \text{i)} E[\epsilon_i] = 0, \quad i=1, \dots, N_y \\ \text{ii)} \text{Var}[\epsilon_i] = \sigma_0^2, \quad i=1, \dots, N_y \\ \text{iii)} \text{Cov}[\epsilon_i, \epsilon_j] = 0, \quad i \neq j \end{array} \right\} \sim N(0, \sigma_0^2)$$

Goal get estimators, $\hat{\tilde{\boldsymbol{\theta}}} + \hat{\sigma}^2$, for $\tilde{\boldsymbol{\theta}}^0 + \sigma_0^2$, + then estimates, $\tilde{\boldsymbol{\theta}}_{\text{ols}} + S^2$, with their sampling distributed.

OLS

$$(f(\tilde{\mathbf{x}}; \tilde{\boldsymbol{\theta}}) = \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}})$$

- Assuming $\epsilon_i \sim N(0, \sigma^2)$, we seek $\tilde{\boldsymbol{\theta}}$ s.t.

$$J(\tilde{\boldsymbol{\theta}}) = (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}})^T (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}})$$

- We want to minimize $J(\tilde{\boldsymbol{\theta}})$.

- For vector valued $\tilde{\boldsymbol{\theta}} \in \mathbb{R}^P$, we can use the gradient to minimize $J(\tilde{\boldsymbol{\theta}})$.

$$\nabla_{\tilde{\boldsymbol{\theta}}} J = \nabla_{\tilde{\boldsymbol{\theta}}} [(\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}})^T (\tilde{\mathbf{y}} - \tilde{\mathbf{X}} \tilde{\boldsymbol{\theta}})]$$

$$\nabla_{\vec{\theta}} J = \nabla_{\vec{\theta}} \left[(\vec{y} - \tilde{X}\vec{\theta})^T (\vec{y} - \tilde{X}\vec{\theta}) \right]$$

$$= -2\tilde{X}^T [\vec{y} - \tilde{X}\vec{\theta}] = -2[\tilde{X}^T \vec{y} - \tilde{X}^T \tilde{X} \vec{\theta}]$$

- Setting $\nabla_{\vec{\theta}} J = 0$

$$\Rightarrow \boxed{\vec{\theta}_{OLS} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \vec{y}}$$

• Note: if $\tilde{X}^T \tilde{X}$ is ill-conditioned, we can use

$$\tilde{X}^T \tilde{X} \approx (\tilde{X}^T \tilde{X} + \alpha I)$$

where α regularizes the problem.

Estimator Properties

- Once we get estimates for $\vec{\theta}$, we can look at its sampling distribution. (note that $\hat{\cdot}$ notation \Rightarrow estimator)

$$i) E[\hat{\theta}] = E[(\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \vec{y}] = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T E[\vec{y}]$$

$$= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T E[\tilde{X}\vec{\theta} + \vec{\epsilon}]$$

$$= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T \tilde{X} \vec{\theta} + 0$$

$$= \vec{\theta}^0$$

$$ii) \text{Var}[\hat{\theta}] = E[(\hat{\theta} - \vec{\theta}^0)(\hat{\theta} - \vec{\theta}^0)^T], \text{ let } A = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T$$

$$ii) \text{Var}[\hat{\theta}] = E[(\hat{\theta} - \vec{\theta}^0)(\hat{\theta} - \vec{\theta}^0)'] , \text{ let } A = (X^T X)^{-1} X^T$$

$$(\text{Note: } \hat{\theta} = AY = A(X\theta + \vec{e}))$$

$$\begin{aligned} \Rightarrow E[(\vec{\theta}^0 + A\vec{e} - \vec{\theta}^0)(\vec{\theta}^0 + A\vec{e} - \vec{\theta}^0)'] \\ = E[A\vec{e}\vec{e}'A^T] = A E[\vec{e}\vec{e}'] A^T \\ = A \sigma_e^2 A^T = \boxed{\sigma_e^2 (X^T X)^{-1}} \end{aligned}$$

iii) We need an estimator for σ_e^2 .

$$\text{Let } \hat{R} = \vec{y} - X\hat{\theta}$$

$$\text{Note that } \vec{y} - X\hat{\theta} = \vec{y} - X[(X^T X)^{-1} X^T \vec{y}]$$

$$= (\underline{I} - H)\vec{y}, \quad H = X(X^T X)^{-1} X^T$$

Note: Since $H = X(X^T X)^{-1} X^T$, $X \in \mathbb{R}^{N_y \times P} \Rightarrow H \in \mathbb{R}^{N_y \times N_y}$

$$\Rightarrow H^T = H \quad (\text{sym})$$

$$H^2 = H \quad (\text{Idempotent})$$

$$(\underline{I} - H)^2 = (\underline{I} - H)$$

$$(\underline{I} - H)X = 0$$

$$\hat{R} = (\underline{I} - H)\vec{y} = (\underline{I} - H)(X\vec{\theta} + \vec{e}) = \cancel{(\underline{I} - H)X\vec{\theta}} + (\underline{I} - H)\vec{e} = 0 + (\underline{I} - H)\vec{e}$$

$$\Rightarrow \hat{\vec{R}}_{\sim} = (\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{y} = (\vec{I}_{\sim} - \vec{H}_{\sim}) (\vec{X}_{\sim} \vec{\theta} + \vec{\varepsilon}) = \cancel{(\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{X}_{\sim} \vec{\theta}} + (\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{\varepsilon}$$

$$\Rightarrow \hat{\vec{R}}_{\sim} = (\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{\varepsilon}$$

- Then $\hat{\vec{R}}_{\sim}^+ \hat{\vec{R}}_{\sim} = \vec{\varepsilon}^T (\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{\varepsilon}$

$$\begin{aligned} \Rightarrow E[\hat{\vec{R}}_{\sim}^+ \hat{\vec{R}}_{\sim}] &= E[\vec{\varepsilon}^T (\vec{I}_{\sim} - \vec{H}_{\sim}) \vec{\varepsilon}] \\ &= \sum_i \sum_j h_{ij} \text{cov}(\varepsilon_i, \varepsilon_j) \\ &= \sum_i h_{ii} \text{var}(\varepsilon_i) \\ &= \sigma_0^2 \text{trace}(\vec{I}_{\sim} - \vec{H}_{\sim}) \end{aligned}$$

- Note: $\text{trace}(\vec{I}) = N_y$

$$\text{trace}(\vec{H}) = \text{trace}(\vec{X} (\vec{X}^T \vec{X})^{-1} \vec{X}^T)$$

- Note: $\text{trace}(\vec{A} \vec{B}) = \text{trace}(\vec{B} \vec{A}) = \sum_{i=1}^N \sum_{j=1}^M a_{ij} b_{ji}$

$$\Rightarrow \text{trace}(\underbrace{\vec{X}}_{\vec{A}} \underbrace{(\vec{X}^T \vec{X})^{-1} \vec{X}^T}_{\vec{B}}) = \text{trace}(\underbrace{(\vec{X}^T \vec{X})^{-1} \vec{X}^T}_{\vec{B}} \underbrace{\vec{X}}_{\vec{A}}) = \text{trace}(\vec{I}_M) = P.$$

$$= \sigma_0^2 \cdot (N_y - P)$$

$$\Rightarrow \hat{\sigma}^2 = \frac{1}{N_y - P} \hat{\vec{R}}_{\sim}^+ \hat{\vec{R}}_{\sim}$$

- So in total, given OLS problem with $\vec{y} = \underset{\sim}{X} \vec{\theta} + \vec{e}$

$$\hat{\theta} \sim N(\vec{\theta}^0, \sigma_0^2 (\underset{\sim}{X}^T \underset{\sim}{X})^{-1}), \quad \sigma_0^2 \approx \frac{1}{N_y - p} \underset{\sim}{R}^T \underset{\sim}{R}$$

Frequentist Confidence Intervals

- For $N_y \rightarrow \infty$, the Law of Large numbers provides an asymptotic framework for $\hat{\theta} + \hat{\sigma}^2$

- Since $\hat{\theta} \sim N(\vec{\theta}^0, \sigma^2 (\underset{\sim}{X}^T \underset{\sim}{X})^{-1})$, then the R.V. T_k is

$$T_k = \frac{\hat{\theta}_k - \theta_k^0}{\sqrt{\sigma^2 (\underset{\sim}{X}^T \underset{\sim}{X})^{-1}_{kk}}}$$

- By law of large numbers, $T_k \sim t$ -distribution with $N_y - p$ degrees of freedom

- Then, $\hat{\theta}_k$ has a $1 - \alpha$ confidence interval given by

$$P(\mathcal{P}_-(\hat{\theta}_k) < \hat{\theta}_k < \mathcal{P}_+(\hat{\theta}_k)) = 1 - \alpha$$

$$\mathcal{P}_{\pm} = \hat{\theta}_k \pm t_{N_y - p}^{1 - \alpha/2} \cdot \sqrt{\sigma^2 (\underset{\sim}{X}^T \underset{\sim}{X})^{-1}_{kk}}$$