

Nov 8, 2023.

Reference: Vanderlaet

Maximization

Extremum estimators are called M-estimators

- $\hat{r}(\theta) = \frac{1}{n} \sum_{i=1}^n |x_i - \theta|$: extremum estimator is not differentiable; popⁿ mean would be (Least Absolute Deviation Estimator)

$$X \sim \text{Bernoulli}(p)$$

$$P(X=x) = p^x (1-p)^{1-x}$$

distⁿ.

cannot technically take log.

$$\log \text{ likelihood } f^n = \log P(X=x) = x \log p + (1-x) \log(1-p)$$

$$\hat{l}(p) = \frac{1}{n} \sum_{i=1}^n (x_i \log(p) + (1-x_i) \log(1-p))$$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{L}(p) = \exp(n \hat{l}(p)) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$= \prod_{i=1}^n P(X=x_i) \rightarrow \text{Joint Dist}^n \text{ of the data.}$$

maximises wot to what?

We are maximizing this joint probability \Rightarrow maximizes the probability of our data.

~ MEAN VALUE EXPANSION ~

(*) Empirical Risk

$$\nabla \hat{R}(\hat{\theta}) = \vec{0}$$

Assumed convergence conditions apply.

$$\Rightarrow \hat{\theta} \xrightarrow{P} \theta_0$$

Assumed twice continuity & differentiability

$$\hat{H}(\theta) = \left(\frac{\partial^2 \hat{R}(\theta)}{\partial \theta_i \partial \theta_j} \right) \hat{c}_j$$

$$\nabla \hat{R}(\hat{\theta}) = \nabla \hat{R}(\theta_0) + H(\theta^*)$$

$$\|\theta^* - \theta_0\| < \|\hat{\theta} - \theta_0\|$$

$$\nabla \hat{R}(\hat{\theta}) = \nabla \hat{R}(\theta_0) + H(\theta^*)(\hat{\theta} - \theta_0) = 0 \quad \text{Linearization}$$

$$\hat{\theta} - \theta_0 = -H(\theta^*)^{-1} \nabla \hat{R}(\theta_0)$$

iteration of Newton - Raphson algorithm

$$\theta_{k+1} = \theta_k - H^{-1}(\theta_k) \nabla \hat{R}(\theta_k)$$

$$\left[\begin{array}{l} \text{Mean Value Theorem:} \\ [a, b], c \in [a, b] \\ f(b) - f(a) = f'(c)(b-a) \end{array} \right]$$

$$\hat{\theta} - \theta_0 = -H(\theta^*)^{-1} \nabla \hat{R}(\theta_0)$$

we are trying to find $\hat{\theta}$

$$\begin{aligned} \nabla \hat{R}(\theta_0) &= \nabla \left(\frac{1}{n} \sum_{i=1}^n \ell(Y_i, \theta) \right) \\ &= \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(Y_i, \theta_0) \end{aligned}$$

We know risk is minimized at θ_0 by construction
i.e.

$$\begin{aligned} \nabla R(\theta_0) &= 0 = \nabla E[\ell(Y, \theta_0)] \\ &= E[\nabla_{\theta} \ell(Y, \theta_0)] \end{aligned}$$

(we can't always
do this).

By LLN

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \ell(Y_i, \theta_0) \xrightarrow{P} 0$$

Fubini's Theorem: Allow
you to change the order
of integrals & differentials

$$\text{If } \nabla \hat{R}(\theta_0) \xrightarrow{P} 0$$

$$\Rightarrow \hat{\theta} - \theta_0 = -\hat{H}(\theta^*)^{-1} \nabla \hat{R}(\theta_0) \xrightarrow{P} 0$$

(*) Fubini's Theorem: Allows you to change the order of integrals & differentials

$$\left| \frac{d}{dt} \int_a^b f(x,t) dx = \int_a^b \frac{\partial f(x,t)}{\partial t} dx \right|$$

↓
Leibniz

GENERALIZATION OF CLT:-

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i = \begin{pmatrix} \frac{1}{n} \sum x_i^1 \\ \frac{1}{n} \sum x_i^2 \\ \vdots \\ \frac{1}{n} \sum x_i^d \end{pmatrix}$$

dimension of vector X

Covariance matrix $\Sigma = \left(E[(X^k - E[X^k])(X^p - E[X^p])] \right)_{kp}$

kp element of the matrix

$$\sqrt{n}(\bar{X}_n - \mu) \longrightarrow N(0, \Sigma)$$

$\sqrt{n} \nabla \hat{R}(\theta_0) \xrightarrow{d} N(0, \Sigma) \rightarrow$ Assumption. This is acceptable

'cause

$$\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} l(X_i, \theta) \xrightarrow{P} 0$$

- $\sup_{\theta} |\hat{H}(\theta) - H(\theta)| \xrightarrow{P} 0$ (Assumption)

$$\theta^* \xrightarrow{P} \theta_0 \quad (\text{By construction})$$

$$\Rightarrow \hat{H}(\theta^*) \xrightarrow{P} H(\theta_0)$$

- Random vector $X \sim N(\mu, \Sigma)$ Covariance matrix.

$A \rightarrow$ fixed matrix

$AX \sim (?)$ What would be its distⁿ.

$$\Rightarrow AX \sim N(A\mu, A\Sigma A^T)$$

$$\Sigma = E[(X - E[X])(X - E[X])^T] \rightarrow \text{Outer product}$$

$$\sqrt{n}(\hat{\theta} - \theta_0) = -\hat{H}^{-1}(\theta^*) \sqrt{n} \nabla \hat{L}(\theta_0) \xrightarrow{d} N(0, \hat{H}^{-1}(\theta_0) \Sigma \hat{H}^{-1}(\theta_0)^T)$$