# 2 Basic Approach (theory).

M-estimators are solutions of the vector equation    (iid case)

$\underbrace{\quad\quad}_{\hat{\theta}}$

$$\sum_{i=1}^{n} \psi(\mathbf{Y}_i, \boldsymbol{\theta}) = \mathbf{0}.$$

but what are they estimating? Some true parameter $\theta_0$, where.

$$(*) \quad E_F\left[\psi(Y_i; \theta_0)\right] = \int \psi(y; \theta_0) \, dF(y) = 0 \quad \text{where} \quad Y_i \sim F$$

**Example (Sample Mean, cont'd):** Recall we said $\theta = \overline{Y} = \frac{1}{n}\sum_{i=1}^{n} Y_i$ is an M-estimator for $\psi(Y_i, \theta) = Y_i - \theta$. What is the true parameter?

The true parameter solves    $\int (y - \theta) \, dF(y) = 0$

$$\Rightarrow \quad \underbrace{\int y \, dF(y)}_{\text{population mean}} = \theta$$

Recall the 5-dimensional motivating example.

We said the $\alpha$ which maximizes the pairwise log likelihood seems like it would be a good estimator for $\alpha_0$. We didn't show this.

To do this, we would need to use $(*)$

To arrive at the sandwich estimator, assume $\boldsymbol{Y}_1, \ldots, \boldsymbol{Y}_n \overset{iid}{\sim} F$ and define

$$\boldsymbol{G}_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\psi}(\boldsymbol{Y}_i; \boldsymbol{\theta}).$$

$\uparrow$ depend on $n$.

In the likelihood case:

$$\frac{1}{n} \sum_{i=1}^{n} \underbrace{\psi(Y_i; \theta)}_{\text{score contribution, deriv of log likelihood contribution}}$$

mean derivative of log likelihood contributions.

Taylor expansion of $\boldsymbol{G}_n(\boldsymbol{\theta})$ around $\boldsymbol{\theta}_0$ evaluated at $\hat{\boldsymbol{\theta}}$ yields

$$0 = \underbrace{G_n(\hat{\theta})}_{\substack{\text{definition} \\ \text{of } \hat{\theta}}} = \underbrace{G_n(\theta_0)}_{\substack{\uparrow \\ b \times l}} + \underbrace{G_n'(\theta_0)}_{\substack{\uparrow \\ 5 \times 6 \text{ Jacobian}}} \underbrace{(\hat{\theta} - \theta_0)}_{b \times l} + \underbrace{R_n}_{\substack{\uparrow \\ \text{higher order terms} \\ \text{"residual"}}}$$

Rearranging: $\quad -G_n'(\theta_0)(\hat{\theta} - \theta_0) = G_n(\theta_0) + R_n$

$$\hat{\theta} - \theta_0 = \left\{ -G_n'(\theta_0) \right\}^{-1} G_n(\theta_0) + \underbrace{\left\{ -G_n'(\theta_0) \right\}^{-1} R_n}_{R_n^*}$$

$$\sqrt{n}\left(\hat{\theta} - \theta_0\right) = \underbrace{\left\{ -G_n'(\theta_0) \right\}^{-1}}_{①} \underbrace{\sqrt{n}\, G_n(\theta_0)}_{②} + \underbrace{\sqrt{n}\, R_n^*}_{③}$$

We will look at each piece.

① *  $-G_n'(\theta_0) = \frac{d}{dt} - G_n(\theta_0) = \frac{d}{d\theta}\left[-\frac{1}{n}\sum_{i=1}^{n}\Psi(Y_i,\theta_0)\right] = \frac{1}{n}\sum_{i=1}^{n}-\Psi'(Y_i,\theta_0)$

Define $\boldsymbol{A}(\boldsymbol{\theta}_0) = \mathrm{E}_F[-\boldsymbol{\psi}'(\boldsymbol{Y}_1,\boldsymbol{\theta}_0)]$.

Then $-G_n'(\theta_0) \xrightarrow{p} A(\theta_0)$ by WLLN.

In the likelihood setting, what is A? curvature! because $\Psi$ is the score function (derivative of log likelihood)
$\Rightarrow \Psi'$ is the 2nd derivative of the log likelihood.

② *  $\sqrt{n}\, G_n(\theta_0) = \sqrt{n}\,\frac{1}{n}\sum_{i=1}^{n}\Psi(Y_i,\theta_0) \xrightarrow{d} N(0, B(\theta_0))$  ← Why? because we have a correctly scaled sum of iid things (CLT).

What is $B(\theta_0)$? Should be the variance of $\Psi$.

$B(\theta_0) = E_F\left\{\Psi(Y_1,\theta_0)\,\Psi(Y_1,\theta_0)^T\right\}$.

③ *  $\sqrt{n}\, R_n^* \xrightarrow{p} 0$

This is the "hard part" to prove. We will skip, see Huber (1967) or Serfling (1980).

So putting ①,②,③ together

$$\sqrt{n}\left(\hat{\theta} - \theta_0\right) \xrightarrow{d} \{A(\theta_0)\}^{-1} N(0, B(\theta_0))  \quad\text{(Slutsky's)}.$$

$$\rightsquigarrow^d N\left(0, A(\theta_0)^{-1} B(\theta_0)\{A(\theta_0)^{-1}\}^T\right)$$

or $\quad \hat{\theta} \sim N\left(\theta_0, \frac{1}{n} A(\theta_0)^{-1} B(\theta_0)\{A(\theta_0)\}^{-1}{}^T\right)$

In practice, we don't know $\theta_0 \Rightarrow$ replace w/ $\hat{\theta}$:

$$\hat{\theta} \sim N\left(\theta_0, \frac{1}{n} A(\hat{\theta})^{-1} B(\hat{\theta})\{A(\hat{\theta})^{-1}\}^T\right)$$

$\qquad\qquad\qquad\uparrow\qquad\qquad\uparrow\qquad\quad\uparrow$
$\qquad\qquad\text{curvature}\quad\text{variance}\quad\text{curvature}.$

$\qquad\qquad\quad\text{bread}\qquad\text{meat}\qquad\text{bread} = \text{sandwich!}$

## 2.1 Estimators for $A, B$

If the data truly come from the assumed parametric family $f(y; \boldsymbol{\theta})$,

Then $A(\underline{\theta}_o) = B(\underline{\theta}_o) = I(\underline{\theta}_o)$

                          Information matrix. where the 2 definitions of $I(\underline{\theta}_o)$ are used.

$\Rightarrow$ the sandwich estimator $A(\underline{\theta}_o)^{-1} B(\underline{\theta}_o) \{A^{-1}(\underline{\theta}_o)\}^T = I(\underline{\theta}_o)^{-1}$

One of the key contributions of M-estimation theory is to point out what happens when the assumed parametric family is not correct.

Then $A(\underline{\theta}_o) \neq B(\underline{\theta}_o)$ and we should use the correct limiting dsn covariance matrix.

$$A(\underline{\theta}_o)^{-1} B(\underline{\theta}_o) \{A^{-1}(\underline{\theta}_o)\}^T$$

We can use empirical estimators of $\boldsymbol{A}$ and $\boldsymbol{B}$:

$$A_n(\underline{y}, \hat{\underline{\theta}}) = \frac{1}{n} \sum_{i=1}^{n} \{-\psi'(y_i, \hat{\theta})\}$$

                         average curvature evaluated at $\hat{\theta}$

$$B_n(\underline{y}, \hat{\underline{\theta}}) = \frac{1}{n} \sum_{i=1}^{n} \psi(y_i, \hat{\theta}) \, \psi(y_i, \hat{\theta})^T$$

                                   variance estimate.

might need to use numeric differentiation to approximate.

Remember, the Hessian in code is $nA_n(\underline{y}, \hat{\underline{\theta}})$.

**Example (Coefficient of Variation):** Let $Y_1, \ldots, Y_n$ be iid from some distribution with finite fourth moment. The coefficient of variation is defined at $\hat\theta_3 = s_n/\overline{Y}$.

How would we get a CI for the coefficient of variation, $\theta_3 = \frac{\sigma}{\mu}$?

Bootstrap? probably.

We'll try M-estimation.

Define a three dimensional $\psi$ so that $\hat\theta_3$ is defined by summing the third component. What is the vector valued function $\psi$ which yields an M-estimator for the coefficient of variation?

$$\underline{\Psi}(Y_i, \underline{\theta}) = \begin{pmatrix} Y_i - \theta_1 \\ (Y_i - \theta_1)^2 - \theta_2 \\ \theta_1 \theta_3 - \sqrt{\theta_2} \end{pmatrix}$$

$$\sum_{i=1}^{n} \underline{\Psi}(Y_i, \underline{\theta}) = \begin{pmatrix} \sum_{i=1}^{n} Y_i - n\theta_1 \\ \sum_{i=1}^{n} (Y_i - \theta_1)^2 - n\theta_2 \\ n\theta_1 \theta_3 - \sqrt{\theta_2} \end{pmatrix} \overset{set}{=} 0$$

$$\implies \theta_1 = \frac{1}{n} \sum_{i=1}^{n} Y_i = \overline{Y}$$

$$\theta_2 = \frac{1}{n} \sum_{i=1}^{n} (Y_i - \theta_1)^2 = \underset{n}{\textcircled{$s^2$}} \leftarrow \text{divide by } n, \text{ not } n\text{-}1$$

$$\theta_3 = \frac{\sqrt{\theta_2}}{\theta_1}$$

What parameter vector is being estimated by the M-estimator?

$$E\left[\Psi_1(Y_i,\theta)\right] = E\left[Y_i - \theta_1\right] = \mu - \theta_1 = 0 \implies \theta_1 = \mu.$$

$$E\left[\Psi_2(Y_i,\theta)\right] = E\left[(Y_i-\theta_1)^2 - \theta_2\right] = 0 \implies \theta_2 = Var(Y_1) = \sigma^2$$

$$E\left[\Psi_3(Y_i,\theta)\right] = E\left[\theta_1\theta_3 - \sqrt{\theta_2}\right] = \theta_1\theta_3 - \sqrt{\theta_2} \overset{Set}{=} 0 \implies \theta_3 = \frac{\sqrt{\theta_2}}{\theta_1}.$$

$$\implies \begin{pmatrix} \mu \\ \sigma^2 \\ \frac{\sigma}{\mu} \end{pmatrix}.$$

What are the matrices **A** and **B**?

$$A = E\left[-\Psi'(Y_i,\theta_0)\right] \qquad \Psi(Y_i,\theta) = \begin{pmatrix} Y_i - \theta_1 \\ (Y_i-\theta_1)^2 - \theta_2 \\ \theta_1\theta_3 - \sqrt{\theta_2} \end{pmatrix}.$$

$$\Psi' = \begin{pmatrix} -1 & 0 & 0 \\ -2(Y_i-\theta_1) & -1 & 0 \\ \theta_3 & -\frac{1}{2}\theta_2^{-1/2} & \theta_1 \end{pmatrix}$$

$$A = E\left[-\Psi'(Y_i,\theta_0)\right] = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -\frac{\sigma}{\mu} & \frac{1}{2\sigma} & -\mu \end{pmatrix}$$

$$B = E\left[\Psi(Y_1,\theta_0)\Psi(Y,\theta_0)^T\right]$$

$$= E\begin{bmatrix} (Y_i-\theta_1)^2 & (Y_i-\theta_1)[(Y_i-\theta_1)^2-\theta_2] & (Y_i-\theta_1)(\theta_1\theta_3-\sqrt{\theta_2}) \\ (Y_i-\theta_1)[(Y_i-\theta_1)^2-\theta_2] & [(Y_i-\theta_1)^2-\theta_2]^2 & [(Y_i-\theta_1)^2-\theta_2](\theta_1\theta_3-\sqrt{\theta_2}) \\ (Y_i-\theta_1)(\theta_1\theta_3-\sqrt{\theta_2}) & [(Y_i-\theta_1)^2-\theta_2](\theta_1\theta_3\sqrt{\theta_2}) & (\theta_1\theta_3-\sqrt{\theta_2}) \end{bmatrix}$$

$$= \begin{bmatrix} \sigma^2 & \mu_3 & 0 \\ \mu_3 & \mu_4-\sigma^4 & 0 \\ 0 & 0 & 0 \end{bmatrix} \qquad \text{where } \mu_j = E(Y_1-\theta_1)^j$$
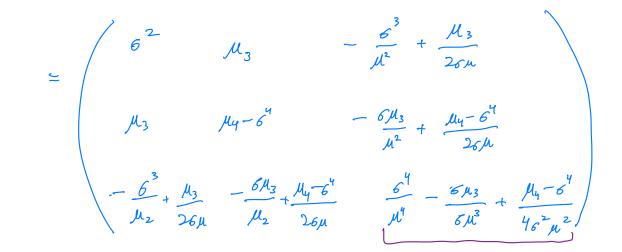
Write out the asymptotic variance, **V**.

$$V = A^{-1} B (A^{-1})^T$$

Using row operations (not shown):

$$A^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -\frac{\sigma}{\mu^2} & \frac{1}{2\sigma\mu} & -\frac{1}{\mu} \end{pmatrix}$$

$$A^{-1}B = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -\frac{\sigma}{\mu^2} & \frac{1}{2\sigma\mu} & -\frac{1}{\mu} \end{pmatrix}\begin{pmatrix} \sigma^2 & \mu_3 & 0 \\ \mu_3 & \mu_4 - \sigma^4 & 0 \\ 0 & 0 & 0 \end{pmatrix} = \begin{pmatrix} \sigma^2 & \mu_3 & 0 \\ \mu_3 & \mu_4 - \sigma^4 & 0 \\ -\frac{\sigma^3}{\mu_2} + \frac{\mu_3}{2\sigma\mu} & -\frac{\sigma\mu_3}{\mu_2} + \frac{\mu_4 - \sigma^4}{2\mu} & 0 \end{pmatrix}$$

$$\Rightarrow A^{-1}B(A^{-1})^T = \quad \ldots\ldots$$

$$\simeq \begin{pmatrix} \sigma^2 & \mu_3 & -\frac{\sigma^3}{\mu^2} + \frac{\mu_3}{2\sigma\mu} \\[2em] \mu_3 & \mu_4 - \sigma^4 & -\frac{\sigma\mu_3}{\mu^2} + \frac{\mu_4 - \sigma^4}{2\sigma\mu} \\[2em] -\frac{\sigma^3}{\mu_2} + \frac{\mu_3}{2\sigma\mu} & -\frac{\sigma\mu_3}{\mu_2} + \frac{\mu_4 - \sigma^4}{2\sigma\mu} & \frac{\sigma^4}{\mu^4} - \frac{\sigma\mu_3}{\sigma\mu^3} + \frac{\mu_4 - \sigma^4}{4\sigma^2\mu^2} \end{pmatrix}$$

Assume $Y_i$ are iid from a normal distribution with mean 10 and standard deviation 1. Calculate $V_{3,3}$. Assume you have a sample of size 25 and you get an estimated coefficient of variation of 0.11. Give the asymptotic 95% confidence interval.

$$V_{3,3} = \frac{\sigma^4}{\mu^4} - \frac{\sigma\mu_3}{2\sigma\mu^3} + \frac{\mu_4 - \sigma^4}{4\sigma^2\mu^2}$$

If $Y \sim N(10,1)$, $\mu_3 = 0$, $\mu_4 = 3$ (looked up).

$$\Rightarrow V_{33} = \frac{1}{10^4} - 0 + \frac{3-1}{4\cdot1\cdot10^2} = \frac{1}{10000} + \frac{1}{200} = .0051.$$

$$n = 25 \Rightarrow Var(\hat{\theta_3}) = \frac{.0051}{25} = .0002^{04}.$$

$$CI: 0.11 \pm 1.96\sqrt{2.04e^{-4}}$$

$$(0.082, 0.138)$$