

# M-Estimators

## 1.1 Definition

**Definition 1.1.1.** For any function  $\psi(\mathbf{y}; \boldsymbol{\theta})$ , define the *functional*  $T(F) = \boldsymbol{\theta}_0$ , where  $\boldsymbol{\theta}_0$  is a solution to the equation:

$$\Psi_F(\boldsymbol{\theta}_0) \equiv \int \psi(\mathbf{y}; \boldsymbol{\theta}_0) dF(\mathbf{y}) = 0. \quad (1.1.1)$$

Suppose  $(\mathbf{y}_i)_{i=1}^n$  is a random sample from distribution  $F$ , and let  $\mathbb{F}_n$  denote the empirical distribution function. The **M-estimator** corresponding to  $\psi$  is  $\hat{\boldsymbol{\theta}}_n = T(\mathbb{F}_n)$ , which is a solution to the equation:

$$\Psi_n(\hat{\boldsymbol{\theta}}_n) = \int \psi(\mathbf{y}; \boldsymbol{\theta}_0) d\mathbb{F}_n(\mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \psi(\mathbf{y}_i; \hat{\boldsymbol{\theta}}_n) \stackrel{\text{Set}}{=} \mathbf{0}.$$

■

**Discussion 1.1.1.** Often,  $\psi(\mathbf{y}; \boldsymbol{\theta})$  arises as the gradient of an objective function  $q(\mathbf{y}; \boldsymbol{\theta})$ :

$$\psi(\mathbf{y}; \boldsymbol{\theta}_0) = \frac{\partial}{\partial \boldsymbol{\theta}_0} q(\mathbf{y}; \boldsymbol{\theta}_0),$$

and  $\boldsymbol{\theta}_0$  solving (1.1.1) is a solution to the minimization problem:

$$\boldsymbol{\theta}_0 = \arg \min_{\boldsymbol{\theta} \in \Theta} \int q(\mathbf{y}; \boldsymbol{\theta}) dF(\mathbf{y}).$$

For example, suppose  $F$  is a parametric distribution  $F(\mathbf{y}; \boldsymbol{\theta}_0)$  and the goal is to estimate  $\boldsymbol{\theta}_0$ . Consider specifying the negative log likelihood  $-\ln f(\mathbf{y}; \boldsymbol{\theta})$  as the objective  $q(\mathbf{y}; \boldsymbol{\theta})$ . The gradient  $\psi(\mathbf{y}; \boldsymbol{\theta})$  of  $q(\mathbf{y}; \boldsymbol{\theta})$  is the parametric score equation:

$$\Psi_n(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \boldsymbol{\theta}} \ln f(\mathbf{y}_i; \boldsymbol{\theta}).$$

The maximum likelihood estimator  $\hat{\boldsymbol{\theta}}_n$ , which satisfies  $\Psi_n(\hat{\boldsymbol{\theta}}_n) = \mathbf{0}$ , is therefore an example of an M-estimator. See section (1.3) for additional examples. ♠

## 1.2 Asymptotics

**Theorem 1.2.1 (Consistency).** Suppose  $(\mathbf{y}_i)_{i=1}^n$  is a random sample from  $F$ , and that  $\boldsymbol{\theta}$  belongs to a compact parameter space  $\Theta$ . Assume that:

- i.  $\Psi_F(\boldsymbol{\theta})$  exists for all  $\boldsymbol{\theta} \in \Theta$ , and that  $\boldsymbol{\theta}_0$  is the unique zero of  $\Psi_F(\boldsymbol{\theta})$ .
- ii. Each component of  $\boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta})$  is continuous and  $\boldsymbol{\theta}$  and bounded by an integrable function of  $\mathbf{y}$ , not depending on  $\boldsymbol{\theta}$ .

If  $\Psi_n(\hat{\boldsymbol{\theta}}_n) \xrightarrow{as} \mathbf{0}$ , then  $\hat{\boldsymbol{\theta}}_n \xrightarrow{as} \boldsymbol{\theta}_0$ . □

**Remark 1.2.1.** See Boos and Stefanski (2013), theorem 7.1. ◆

**Theorem 1.2.2 (Asymptotic Normality).** Suppose  $(\mathbf{y}_i)_{i=1}^n$  is a random sample from  $F$ , and  $\boldsymbol{\theta}$  belongs to a compact parameter space  $\Theta$ . Assume that:

- i.  $\Psi_F(\boldsymbol{\theta})$  exists for all  $\boldsymbol{\theta} \in \Theta$ , and that  $\boldsymbol{\theta}_0$  is the unique zero of  $\Psi_F(\boldsymbol{\theta})$ .
- ii.  $\boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta})$  is continuous and twice differentiable with respect to  $\boldsymbol{\theta}$  for all  $\mathbf{y}$  in the support of  $F$  and  $\boldsymbol{\theta}$  in a neighborhood of  $\boldsymbol{\theta}_0$ .
- iii. For  $\boldsymbol{\theta}$  in a neighborhood of  $\boldsymbol{\theta}_0$ , there exists an integrable function  $g(\mathbf{y})$  such that:

$$\left| \frac{\partial^2}{\partial \theta_{j_1} \partial \theta_{j_2}} \psi_{j_3}(\mathbf{y}; \boldsymbol{\theta}) \right| \leq g(\mathbf{y})$$

for  $\forall (j_1, j_2, j_3) \in \{1, \dots, \dim(\boldsymbol{\theta})\}^3$ , and  $\int g(\mathbf{y}) dF(\mathbf{y})$ .

- iv.  $\mathbf{A}(\boldsymbol{\theta}_0)$  exists and is non-singular, where:

$$\mathbf{A}(\boldsymbol{\theta}_0) = -\mathbb{E}_F\{\dot{\boldsymbol{\psi}}(\mathbf{y}; \boldsymbol{\theta}_0)\} = -\int \dot{\boldsymbol{\psi}}(\mathbf{y}; \boldsymbol{\theta}_0) dF(\mathbf{y}).$$

- v.  $\mathbf{B}(\boldsymbol{\theta}_0)$  exists and is finite, where:

$$\mathbf{B}(\boldsymbol{\theta}_0) = \mathbb{E}_F\{\boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta}_0) \otimes \boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta}_0)\} = \int \boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta}_0) \otimes \boldsymbol{\psi}(\mathbf{y}; \boldsymbol{\theta}_0) dF(\mathbf{y}).$$

If  $\Psi_n(\hat{\boldsymbol{\theta}}_n) = o_p(n^{-1/2})$ , then:

$$\sqrt{n}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N\{\mathbf{0}, \boldsymbol{\Omega}(\boldsymbol{\theta}_0)\},$$

$$\boldsymbol{\Omega}(\boldsymbol{\theta}_0) = \mathbf{A}^{-1}(\boldsymbol{\theta}_0) \mathbf{B}(\boldsymbol{\theta}_0) \mathbf{A}^{-T}(\boldsymbol{\theta}_0). \quad \square$$

**Remark 1.2.2.** See Boos and Stefanski (2013), theorem 7.2. ◆

**Proof.** By Taylor expansion:

$$o_p(n^{-1/2}) = \Psi_n(\hat{\theta}_n) = \Psi_n(\theta_0) + \dot{\Psi}_n(\theta_0)(\hat{\theta}_n - \theta_0) + \mathbf{R}_n(\hat{\theta}_n - \theta_0),$$

where the remainder:

$$\mathbf{R}_n = \sum_{j=1}^{\dim(\theta)} \frac{\partial \dot{\Psi}_n}{\partial \theta_j} \Big|_{\theta=\theta_n^*} (\hat{\theta}_{n,j} - \theta_{0,j})$$

and  $\|\theta_n^* - \theta_0\| \leq \|\hat{\theta}_n - \theta_0\|$ . By condition (iii.), the matrix of partial derivatives appearing in the remainder is bounded by a function not depending on  $\theta$ , hence:

$$\|\mathbf{R}_n\| = \mathcal{O}_p(\|\hat{\theta}_n - \theta_0\|).$$

Since the conditions for consistency are contained within the conditions for asymptotic normality,  $\hat{\theta}_n - \theta_0 = o_p(1)$ , which implies  $\mathbf{R}_n = o_p(1)$ . Now:

$$\begin{aligned} o_p(n^{-1/2}) &= \Psi_n(\theta_0) + \{\dot{\Psi}_n(\theta_0) + o_p(1)\}(\hat{\theta}_n - \theta_0), \\ \sqrt{n}(\hat{\theta}_n - \theta_0) &= \{-\dot{\Psi}_n(\theta_0) + o_p(1)\}^{-1} \sqrt{n}\Psi_n(\theta_0). \end{aligned}$$

By the LLN:

$$-\dot{\Psi}_n(\theta_0) + o_p(1) \xrightarrow{p} -E\{\dot{\psi}(\mathbf{y}; \theta_0)\} = \mathbf{A}(\theta_0)$$

By continuous mapping:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \sqrt{n}\mathbf{A}^{-1}(\theta_0)\Psi_n(\theta_0) + o_p(1) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \{\mathbf{A}^{-1}(\theta_0)\psi(\mathbf{y}_i; \theta_0)\} + o_p(1).$$

Identify  $\varphi(\mathbf{y}_i; \theta_0) = \mathbf{A}^{-1}(\theta_0)\psi(\mathbf{y}_i; \theta_0)$  as the *influence function* of  $\hat{\theta}_n$ . By the CLT,

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} N\{\mathbf{0}, \Omega(\theta_0)\},$$

where:

$$\Omega(\theta_0) = \text{Var}\{\varphi(\mathbf{y}_i; \theta_0)\} = \mathbf{A}^{-1}(\theta_0)\mathbb{E}\{\psi(\mathbf{y}_i; \theta_0) \otimes \psi(\mathbf{y}_i; \theta_0)\}\mathbf{A}^{-T}(\theta_0).$$

■

## 1.3 Examples

**Example 1.3.1 (Ordinary Least Squares).** Consider the model:

$$Y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i,$$

where  $(\epsilon_i)_{i=1}^n$  are independent, random residuals with  $\mathbb{E}(\epsilon_i|\mathbf{x}_i) = 0$  and  $\text{Var}(\epsilon_i|\mathbf{x}_i) = \sigma^2$ . Define the squared error objective function:

$$Q(\boldsymbol{\beta}) = \frac{1}{2n} \sum_{i=1}^n (Y_i - \mathbf{x}_i' \boldsymbol{\beta})^2.$$

Differentiating to obtain the estimating equation:

$$\boldsymbol{\Psi}_n(\boldsymbol{\beta}) = -\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = -\frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i' \boldsymbol{\beta}) \mathbf{x}_i.$$

Identify  $\boldsymbol{\psi}(Y_i, \mathbf{x}_i; \boldsymbol{\beta}) = (Y_i - \mathbf{x}_i' \boldsymbol{\beta}) \mathbf{x}_i$ . Let  $\hat{\boldsymbol{\beta}}$  denote a solution to  $\boldsymbol{\Psi}_n(\hat{\boldsymbol{\beta}}) = \mathbf{0}$ , then under the assumptions of theorem (1.2.2),  $\hat{\boldsymbol{\beta}}$  is an  $M$ -estimator.

The partial of  $\boldsymbol{\psi}$  with respect to  $\boldsymbol{\beta}'$  is:

$$\frac{\partial \boldsymbol{\psi}}{\partial \boldsymbol{\beta}'} = -\mathbf{x}_i \otimes \mathbf{x}_i.$$

$\mathbf{A}(\boldsymbol{\theta}_0)$  takes the form:

$$\mathbf{A}(\boldsymbol{\theta}_0) \equiv -\mathbb{E}_0 \left( \frac{\partial \boldsymbol{\psi}}{\partial \boldsymbol{\beta}'} \right) = \mathbb{E}_0(\mathbf{x}_i \otimes \mathbf{x}_i).$$

The empirical estimate of  $\mathbf{A}(\boldsymbol{\theta}_0)$  is:

$$\hat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \otimes \mathbf{x}_i.$$

The outer product of  $\boldsymbol{\psi}$  is:

$$\boldsymbol{\psi} \otimes \boldsymbol{\psi} = (Y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \mathbf{x}_i \otimes \mathbf{x}_i.$$

$\mathbf{B}(\boldsymbol{\beta}_0)$  takes the form:

$$\mathbf{B}(\boldsymbol{\beta}_0) = \mathbb{E}_0(\boldsymbol{\psi} \otimes \boldsymbol{\psi}) = \mathbb{E}_0\{(Y_i - \mathbf{x}_i' \boldsymbol{\beta})^2 \mathbf{x}_i \otimes \mathbf{x}_i\} = \sigma^2 \mathbb{E}_0(\mathbf{x}_i \otimes \mathbf{x}_i).$$

The empirical estimate of  $\mathbf{B}(\boldsymbol{\beta}_0)$  is:

$$\hat{\mathbf{B}} = \hat{\sigma}^2 \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \otimes \mathbf{x}_i = \hat{\sigma}^2 \mathbf{X}' \mathbf{X},$$

where  $\hat{\sigma}^2$  is any consistent estimate of the residual variance.

The asymptotic approximation to the distribution of  $\hat{\boldsymbol{\beta}}$  is:

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}_0, n^{-1} \hat{\boldsymbol{\Omega}}),$$

where  $\hat{\boldsymbol{\Omega}} = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} \hat{\mathbf{A}}^{-T}$ .



**Example 1.3.2 (Non-linear Least Squares).** Consider the model:

$$Y_i = g(\mathbf{x}_i, \boldsymbol{\beta}) + \epsilon_i,$$

where  $g$  is a known differentiable function, and  $(\epsilon_i)_{i=1}^n$  are independent, random residuals with  $\mathbb{E}(\epsilon_i | \mathbf{x}_i) = 0$  and  $\text{Var}(\epsilon_i | \mathbf{x}_i) = \sigma_i^2$ . Define the squared error objective function:

$$Q(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \{Y_i - g(\mathbf{x}_i, \boldsymbol{\beta})\}^2.$$

Differentiating to obtain the estimating equation:

$$\boldsymbol{\Psi}_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \{Y_i - g(\mathbf{x}_i, \boldsymbol{\beta})\} \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}).$$

Identify  $\boldsymbol{\psi}(Y_i, \mathbf{x}_i; \boldsymbol{\beta}) = \{Y_i - g(\mathbf{x}_i, \boldsymbol{\beta})\} \dot{g}(\mathbf{x}_i, \boldsymbol{\beta})$ . Let  $\hat{\boldsymbol{\beta}}$  denote a solution to  $\boldsymbol{\Psi}_n(\hat{\boldsymbol{\beta}}) = \mathbf{0}$ , then under the assumptions of theorem (1.2.2),  $\hat{\boldsymbol{\beta}}$  is an  $M$ -estimator.

$\mathbf{A}(\boldsymbol{\beta}_0)$  takes the form:

$$\mathbf{A}(\boldsymbol{\beta}_0) = -\mathbb{E}_0 \{ \dot{\boldsymbol{\psi}}(Y_i, \mathbf{x}_i; \boldsymbol{\beta}_0) \} = \mathbb{E}_0 \{ \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \otimes \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \}.$$

The empirical estimate of  $\mathbf{A}(\boldsymbol{\beta}_0)$  is:

$$\hat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \dot{g}(\mathbf{x}_i, \hat{\boldsymbol{\beta}}) \otimes \dot{g}(\mathbf{x}_i, \hat{\boldsymbol{\beta}}).$$

$\mathbf{B}(\boldsymbol{\beta}_0)$  takes the form:

$$\mathbf{B}(\boldsymbol{\beta}_0) = \mathbb{E}_0 \left[ \{Y_i - g(\mathbf{x}_i, \boldsymbol{\beta}_0)\}^2 \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \otimes \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \right] = \sigma_i^2 \mathbb{E}_0 \{ \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \otimes \dot{g}(\mathbf{x}_i, \boldsymbol{\beta}_0) \}.$$

The empirical estimate of  $\mathbf{B}(\boldsymbol{\beta}_0)$  is:

$$\hat{\mathbf{B}} = \frac{1}{n} \sum_{i=1}^n \{Y_i - g(\mathbf{x}_i, \hat{\boldsymbol{\beta}})\}^2 \dot{g}(\mathbf{x}_i, \hat{\boldsymbol{\beta}}) \otimes \dot{g}(\mathbf{x}_i, \hat{\boldsymbol{\beta}}).$$

The asymptotic approximation to the distribution of  $\hat{\boldsymbol{\beta}}$  is:

$$\hat{\boldsymbol{\beta}} \sim N(\boldsymbol{\beta}_0, n^{-1} \hat{\boldsymbol{\Omega}}),$$

where  $\hat{\boldsymbol{\Omega}} = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} \hat{\mathbf{A}}^{-T}$ . ♠

**Example 1.3.3 (Robust Regression).** Consider again the model:

$$Y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i,$$

where  $(\epsilon_i)_{i=1}^n$  are independent, random residuals with  $\mathbb{E}(\epsilon_i|\mathbf{x}_i) = 0$  and  $\text{Var}(\epsilon_i|\mathbf{x}_i) = \sigma_i^2$ .

The ordinary least squares estimating equations are:

$$\frac{1}{n} \sum_{i=1}^n (Y_i - \mathbf{x}_i' \boldsymbol{\beta}) \mathbf{x}_i.$$

A robust, alternative set of estimating equations is:

$$\boldsymbol{\Psi}_n(\boldsymbol{\beta}) = \sum_{i=1}^n \psi_\tau(Y_i - \mathbf{x}_i' \boldsymbol{\beta}) \mathbf{x}_i.$$

Here  $\psi_\tau$  is a bounded loss function, such as Huber's function:

$$\psi_\tau(t) = \begin{cases} -\tau, & t < -\tau, \\ x, & -\tau < t < \tau, \\ \tau, & t > \tau. \end{cases}$$

or Tukey's biweight function:

$$\psi_\tau(t) = \begin{cases} t \left(1 - \frac{t^2}{\tau^2}\right)^2, & |t| < \tau, \\ 0, & |t| > \tau. \end{cases}$$

Let  $\hat{\boldsymbol{\beta}}$  denote a solution to  $\boldsymbol{\Psi}_n(\hat{\boldsymbol{\beta}}) = \mathbf{0}$ , then under the assumptions of theorem (1.2.2),  $\hat{\boldsymbol{\beta}}$  is an  $M$ -estimator.

$\mathbf{A}(\boldsymbol{\beta}_0)$  takes the form:

$$\mathbf{A}(\boldsymbol{\beta}_0) = -\mathbb{E}_0 \{ \dot{\psi}_\tau(Y_i - \mathbf{x}_i' \boldsymbol{\beta}_0) \mathbf{x}_i \otimes \mathbf{x}_i \}.$$

The empirical estimate of  $\mathbf{A}(\boldsymbol{\beta}_0)$  is:

$$\hat{\mathbf{A}} = \frac{1}{n} \sum_{i=1}^n \dot{\psi}_\tau(Y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \mathbf{x}_i \otimes \mathbf{x}_i.$$

$\mathbf{B}(\boldsymbol{\beta}_0)$  takes the form:

$$\mathbf{B}(\boldsymbol{\beta}_0) = \mathbb{E}_0 \{ \psi_\tau^2(Y_i - \mathbf{x}_i' \boldsymbol{\beta}_0) \mathbf{x}_i \otimes \mathbf{x}_i \}.$$

The empirical estimate of  $\mathbf{B}(\boldsymbol{\beta}_0)$  is:

$$\hat{\mathbf{B}} = \frac{1}{n} \sum_{i=1}^n \psi_\tau^2(Y_i - \mathbf{x}_i' \hat{\boldsymbol{\beta}}) \mathbf{x}_i \otimes \mathbf{x}_i.$$

The asymptotic approximation to the distribution of  $\hat{\boldsymbol{\beta}}$  is:

$$\hat{\boldsymbol{\beta}} \dot{\sim} N(\boldsymbol{\beta}_0, n^{-1} \hat{\boldsymbol{\Omega}}),$$

where  $\hat{\boldsymbol{\Omega}} = \hat{\mathbf{A}}^{-1} \hat{\mathbf{B}} \hat{\mathbf{A}}^{-T}$ .

