

Parametric Estimation

Maxim Panov

Skoltech

November, 2021

Skoltech

Outline

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

Outline

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

Problem Statement

Data: $X_1, \dots, X_n \sim f_o$.

Model: we aim to model the data by $f(x; \theta)$ for some value of θ .

Definition

General form of a parametric model:

$$\mathfrak{F} = \left\{ f(x; \theta), \theta \in \Theta \subseteq \mathbb{R}^k \right\},$$

where Θ is the parameter space, $\theta = (\theta_1, \dots, \theta_k)$ is parameter vector and $k \in \mathbb{N}$.

Problem: find an estimate of $T(\theta)$, where T is some function.

Usual assumption: There exists θ_* such that $f_o = f(x; \theta_*)$.

Problem Statement

Model: we aim to model the data by $f(x; \theta)$ for some value of θ .

Problem: find an estimate of $T(\theta)$, where T is some function.

Example

- ▶ Consider a random variable with the distribution $\mathcal{N}(\mu, \sigma^2)$.
- ▶ In this case $\theta = (\mu, \sigma)$.
- ▶ If the task is to estimate just μ , then $\mu = T(\theta)$.
- ▶ In this case σ is called a *nuisance* parameter.

Example

Example

- ▶ Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$.
- ▶ $\theta = (\mu, \sigma)$ is a vector from a parameter space $\Theta = \{(\mu, \sigma) : \mu \in \mathbb{R}, \sigma > 0\}$.
- ▶ Suppose that X_i is a certain integral characteristic of a blood sample.
- ▶ **Problem** is stated as follows: estimate τ is a fraction of the blood samples, for which this value exceeds 1.

$$\begin{aligned}\tau = \mathbb{P}(X > 1) &= 1 - \mathbb{P}(X < 1) = 1 - \mathbb{P}\left(\frac{X - \mu}{\sigma} < \frac{1 - \mu}{\sigma}\right) = \\ &= 1 - \mathbb{P}\left(Z < \frac{1 - \mu}{\sigma}\right) = 1 - \Phi\left(\frac{1 - \mu}{\sigma}\right).\end{aligned}$$

- ▶ $\tau = T(\mu, \sigma) = 1 - \Phi((1 - \mu)/\sigma)$ is our parameter of interest and Z is a standard normal random variable.

One More Example

Example

- ▶ Let $X \sim \text{Gamma}(\alpha, \beta)$, that has the following density:

$$f(x; \alpha, \beta) = \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta}, \quad \text{where } \alpha, \beta, x > 0.$$

- ▶ Here $\theta = (\alpha, \beta)$ is a parameter vector and $\Gamma(\alpha)$ is Gamma function:

$$\Gamma(\alpha) = \int_0^\infty y^{\alpha-1} e^{-y} dy.$$

- ▶ Gamma distribution may be used to model life expectancy.
- ▶ If the task is to estimate the average life expectancy then

$$T(\alpha, \beta) = \mathbb{E}_\theta(X) = \alpha\beta.$$

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

Method of Moments

- ▶ Let $\theta = (\theta_1, \dots, \theta_k)$ be a parameter vector.
- ▶ For $1 \leq j \leq k$ define j -th (non-centered) moment as follows:

$$\alpha_j(\theta) = \mathbb{E}_\theta(X^j) = \int x^j dF_\theta(x) = \int x^j f(x; \theta) dx,$$

- ▶ If we are given the values $\alpha_1, \dots, \alpha_k$ then we can consider the system of equations:

$$\alpha_1(\theta) = \alpha_1,$$

$$\alpha_2(\theta) = \alpha_2,$$

...

$$\alpha_k(\theta) = \alpha_k.$$

- ▶ If the system above has solution, then it can be used as an estimate of θ_* .

Method of Moments

Idea: use δ -method by computing the corresponding j -th sample moment:

$$\hat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n X_i^j$$

and plugging it into the equations.

Definition

$\hat{\theta}_n$ is a method of moments estimate of $\theta = (\theta_1, \dots, \theta_k)$ if

$$\alpha_1(\hat{\theta}_n) = \hat{\alpha}_1,$$

$$\alpha_2(\hat{\theta}_n) = \hat{\alpha}_2,$$

...

$$\alpha_k(\hat{\theta}_n) = \hat{\alpha}_k.$$

Example

Example

- ▶ Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$.
- ▶ Find an estimate of parameter p .

Example

Example

- ▶ Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$.
- ▶ Find an estimate of parameter p .
- ▶ $\alpha_1 = \mathbb{E}_\theta(X) = p$ and $\hat{\alpha}_1 = n^{-1} \sum_{i=1}^n X_i$,
- ▶ from which we get $\hat{p}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

One More Example

Example

Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$. Estimate parameters μ and σ .

One More Example

Example

Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$. Estimate parameters μ and σ .

$$\alpha_1 = \mathbb{E}_\theta(X_1) = \mu,$$

$$\alpha_2 = \mathbb{E}_\theta(X_1^2) = \mathbb{V}_\theta(X_1) + (\mathbb{E}_\theta(X_1))^2 = \sigma^2 + \mu^2,$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i,$$

$$\hat{\sigma}^2 + \hat{\mu}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2.$$

Solving the system of equations we get that

$$\hat{\mu}_n = \overline{X}_n \quad \text{and} \quad \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \overline{X}_n)^2.$$

Properties of Method of Moments

Theorem

If $\hat{\theta}_n$ is a method of moments estimate of θ then (given certain assumptions about the distribution of the sample) the following properties hold:

1. $\hat{\theta}_n$ exists with probability 1;
2. $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta_*$ when $n \rightarrow \infty$;
3. estimate is asymptotically normal:

$$\sqrt{n}(\hat{\theta}_n - \theta_*) \rightsquigarrow \mathcal{N}(0, \Sigma),$$

where $\Sigma = g\mathbb{E}(YY^T)g^T$, $Y = (X, X^2, \dots, X^k)^T$,
 $g = (g_1, \dots, g_k)$ and $g_j = \partial \alpha_j^{-1}(\theta_*) / \partial \theta$.

Remark: the last property can be used to derive standard errors and confidence intervals.

Method of Moments: Comments

- ▶ not optimal;
- ▶ easy to use;
- ▶ estimates can be used as initial values for more “accurate” methods.

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

Deriving MLE

- If we know the true density, then we can estimate θ in the following way:

$$\theta^*(f_o) = \arg \min_{\theta} \int_{\mathbb{R}^n} f_o(z) \log \frac{f_o(z)}{f(z; \theta)} dz.$$

- It easy to see that

$$\theta^*(f_o) = \arg \max_{\theta} \int_{\mathbb{R}^n} f_o(z) \log[f(z; \theta)] dz.$$

- Let us note that the knowledge of $f_o(z)$ is needed.
- However, importantly, the functional

$$L[f_o, \theta] = \int_{\mathbb{R}^n} f_o(z) \log[f(z; \theta)] dz$$

is linear in $f_o(\cdot)$!

Deriving MLE

Idea: let's estimate the functional $L[f_\circ, \theta]$ based on X^n via δ -method.

- ▶ It leads to estimation of the *logarithm of likelihood*:

$$\bar{L}(\theta; X^n) = \int_{\mathbb{R}^n} \delta(z - X^n) \log[f(z; \theta)] dz = \log[f(X^n; \theta)].$$

- ▶ Maximum likelihood method:

$$\hat{\theta}(X^n) = \arg \max_{\theta} \{\log[f(X^n; \theta)]\}.$$

- ▶ Maximum likelihood method was suggested by Ronald Fisher when he was 22 years old.

Maximum Likelihood Estimation

For the linear model

$$Y_i = \mu + \epsilon_i, \quad i = 1, \dots, n$$

we obviously have

$$\bar{L}(\mu; Y^n) = \sum_{i=1}^n \log[f(Y_i - \mu)].$$

Thus, MLE reads as

$$\hat{\mu}(Y^n) = \arg \max_{\mu} \left\{ \sum_{i=1}^n \log[f(Y_i - \mu)] \right\}.$$

Maximum Likelihood Estimation

As random variables Y_i are i.i.d., then

$$\mathbb{E}|\log[f(Y_1 - \tilde{\mu})]| < \infty,$$

for any $\tilde{\mu}$, then by Law of Large Numbers for fixed $\tilde{\mu}$ and $n \rightarrow \infty$

$$\frac{1}{n} \sum_{i=1}^n \log[f(Y_i - \tilde{\mu})] \rightarrow \mathbb{E} \log[f(Y_1 - \tilde{\mu})] = \int_{\mathbb{R}^1} f_{\circ}(x) \log[f(x - \tilde{\mu})] dx,$$

where $f_{\circ}(x)$ is the true density of Y_1 .

That's why for large n :

$$\hat{\mu}(Y^n) \rightarrow \arg \max_{\tilde{\mu}} \mathbb{E} \log[f(Y_1 - \tilde{\mu})] = \arg \min_{\tilde{\mu}} \mathbb{E} \log \frac{f_{\circ}(Y_1)}{f(Y_1 - \tilde{\mu})}.$$

Maximum Likelihood Estimation

Definition

Consider an *i.i.d.* estimate $X_1, \dots, X_n \sim F$ and the distribution has a density $f(x; \theta)$. The *likelihood function* is defined as follows:

$$\mathcal{L}_n(\theta) = \prod_{i=1}^n f(X_i; \theta).$$

Log-likelihood is just taking the logarithm of the expression above:

$$\ell_n(\theta) = \log \mathcal{L}_n(\theta).$$

- ▶ We will view likelihood as a function of the model's parameters $\mathcal{L}_n: \Theta \rightarrow [0, \infty)$.
- ▶ A maximum likelihood estimate (MLE):

$$\hat{\theta}_n \equiv \hat{\theta}(X^n) = \arg \max_{\theta} \ell_n(\theta) = \arg \max_{\theta} \log \mathcal{L}_n(\theta).$$

Example

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$. Find MLE of the parameter p .

Example

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$. Find MLE of the parameter p .

Probability mass function: $f(x; p) = p^x(1 - p)^{1-x}$ where $x = 0, 1$.

Example

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$. Find MLE of the parameter p .

Probability mass function: $f(x; p) = p^x(1 - p)^{1-x}$ where $x = 0, 1$.

Likelihood function:

$$\mathcal{L}_n(p) = \prod_{i=1}^n f(X_i; p) = \prod_{i=1}^n p^{X_i} (1 - p)^{1-X_i} = p^S (1 - p)^{n-S},$$

where $S = \sum_{i=1}^n X_i$.

Example

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$. Find MLE of the parameter p .

Probability mass function: $f(x; p) = p^x(1 - p)^{1-x}$ where $x = 0, 1$.

Likelihood function:

$$\mathcal{L}_n(p) = \prod_{i=1}^n f(X_i; p) = \prod_{i=1}^n p^{X_i}(1 - p)^{1-X_i} = p^S(1 - p)^{n-S},$$

where $S = \sum_{i=1}^n X_i$.

Log-likelihood:

$$\ell_n(p) = S \log p + (n - S) \log(1 - p).$$

Example

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$. Find MLE of the parameter p .

Probability mass function: $f(x; p) = p^x(1 - p)^{1-x}$ where $x = 0, 1$.

Likelihood function:

$$\mathcal{L}_n(p) = \prod_{i=1}^n f(X_i; p) = \prod_{i=1}^n p^{X_i} (1 - p)^{1-X_i} = p^S (1 - p)^{n-S},$$

where $S = \sum_{i=1}^n X_i$.

Log-likelihood:

$$\ell_n(p) = S \log p + (n - S) \log(1 - p).$$

From this we get that MLE estimate is $\hat{p}_n = S/n$.

One More Example

Example

Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$.

$\theta = (\mu, \sigma)$ is our parameter of interest.

One More Example

Example

Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$.

$\theta = (\mu, \sigma)$ is our parameter of interest.

Likelihood function has the following form (up to some constants):

$$\begin{aligned}\mathcal{L}_n(\mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sigma} \exp \left\{ -\frac{1}{2\sigma^2} (X_i - \mu)^2 \right\} = \\ &= \frac{1}{\sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 \right\} = \\ &= \frac{1}{\sigma^n} \exp \left\{ -\frac{nS^2}{2\sigma^2} \right\} \exp \left\{ -\frac{n(\bar{X} - \mu)^2}{2\sigma^2} \right\},\end{aligned}$$

where $\bar{X} = n^{-1} \sum_{i=1}^n X_i$, $S^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

(continued on the next slide)

One More Example

Example (continued)

The last statement follows from:

$$\sum_{i=1}^n (X_i - \mu)^2 = nS^2 + n(\bar{X} - \mu)^2,$$

which can be seen from:

$$\sum_{i=1}^n (X_i - \mu)^2 = \sum_{i=1}^n (X_i - \bar{X} + \bar{X} - \mu)^2.$$

Log-likelihood becomes:

$$\ell_n(\mu, \sigma) = -n \log \sigma - \frac{nS^2}{2\sigma^2} - \frac{n(\bar{X} - \mu)^2}{2\sigma^2}.$$

Let $\frac{\partial \ell(\mu, \sigma)}{\partial \mu} = 0$ and $\frac{\partial \ell(\mu, \sigma)}{\partial \sigma} = 0$, then MLE gives $\hat{\mu}_n = \bar{X}$ and $\hat{\sigma}_n = S$.

Properties of MLEs

- ▶ consistent, that is $\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta_*$, where θ_* is true value of the parameter θ ;
- ▶ *equivariant*: if $\hat{\theta}_n$ is MLE for θ then $g(\hat{\theta}_n)$ is MLE for $g(\theta)$;
- ▶ asymptotically normal: $(\hat{\theta}_n - \theta_*)/\widehat{se} \rightsquigarrow \mathcal{N}(0, 1)$;
- ▶ asymptotically optimal or efficient (for a sufficient sample size it has lower variance).

Remark: the properties of MLEs stated above hold when $f(x; \theta)$ is sufficiently regular. In “difficult” cases MLEs “lose” these properties.

Consistency of MLE

- ▶ Maximizing $\ell_n(\theta)$ is equivalent to maximizing $M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; \theta)}{f(X_i; \theta_*)}$,
- ▶ because $M_n(\theta) = n^{-1}(\ell_n(\theta) - \ell_n(\theta_*))$ and $\ell_n(\theta_*)$ is a constant.
- ▶ Then

$$\begin{aligned}\mathbb{E}_{\theta_*} \left(\log \frac{f(x; \theta)}{f(x; \theta_*)} \right) &= \int \log \left(\frac{f(x; \theta)}{f(x; \theta_*)} \right) f(x; \theta_*) dx = \\ &= - \int \log \left(\frac{f(x; \theta_*)}{f(x; \theta)} \right) f(x; \theta_*) dx = -K(\theta_*, \theta).\end{aligned}$$

- ▶ Hence, $M_n(\theta) \approx -K(\theta_*, \theta)$ attains its maximum at θ_* since $-K(\theta_*, \theta_*) = 0$ and $-K(\theta_*, \theta) < 0$ for $\theta \neq \theta_*$.
- ▶ We need to show that MLE estimate converges *in probability* to the true value of the parameter.

Consistency of MLE

Theorem

- ▶ Let θ_* denote the true value of the parameter θ . Define

$$M_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log \frac{f(X_i; \theta)}{f(X_i; \theta_*)}$$

and $M(\theta) = -K(\theta_*, \theta)$.

- ▶ Suppose that $\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{\mathbb{P}} 0$
- ▶ and for every $\epsilon > 0$ $\sup_{\theta: |\theta - \theta_*| \geq \epsilon} M(\theta) < M(\theta_*)$.
- ▶ Let $\hat{\theta}_n$ denote the MLE, then

$$\hat{\theta}_n \xrightarrow{\mathbb{P}} \theta_*.$$

Consistency of MLE

Proof: since $\hat{\theta}_n$ maximizes $M_n(\theta)$, we have $M_n(\hat{\theta}_n) \geq M_n(\theta_*)$. It follows that,

$$\begin{aligned} M(\theta_*) - M(\hat{\theta}_n) &= M_n(\theta_*) - M(\hat{\theta}_n) + M(\theta_*) - M_n(\theta_*) \\ &\leq M_n(\hat{\theta}_n) - M(\hat{\theta}_n) + M(\theta_*) - M_n(\theta_*) \\ &\leq \sup_{\theta} |M_n(\theta) - M(\theta)| + M(\theta_*) - M_n(\theta_*) \xrightarrow{\mathbb{P}} 0. \end{aligned}$$

Then for any $\delta > 0$:

$$\mathbb{P} \left(M(\hat{\theta}_n) < M(\theta_*) - \delta \right) \rightarrow 0.$$

Consistency of MLE

- ▶ Take arbitrary $\epsilon > 0$.
- ▶ By assumption there exists $\delta > 0$, for which the inequality $|\theta - \theta_*| \geq \epsilon$ implies that $M(\theta) < M(\theta_*) - \delta$.
- ▶ Hence,

$$\mathbb{P}\left(|\hat{\theta}_n - \theta_*| > \epsilon\right) \leq \mathbb{P}\left(M(\hat{\theta}_n) < M(\theta_*) - \delta\right) \rightarrow 0.$$

Equivariance

Theorem

Let $\tau = g(\theta)$ be a function of the parameter θ and let $\hat{\theta}_n$ be its MLE estimate.

Then $\hat{\tau}_n = g(\hat{\theta}_n)$ is MLE for $\tau = g(\theta)$.

Proof:

- ▶ Assume that g is one to one mapping. i.e. an inverse function $h = g^{-1}$ exists.
- ▶ For any τ and $\theta = h(\tau)$:

$$\mathcal{L}_n(\theta) = \mathcal{L}_n(h(\tau)).$$

- ▶ Due to the mapping being one to one the maximizers coincide:

$$\hat{\theta}_n = h(\hat{\tau}_n)$$

and consequently

$$\hat{\tau}_n = g(\hat{\theta}_n).$$

Example

- ▶ Let $X_1, \dots, X_n \sim \mathcal{N}(\theta, 1)$.
- ▶ MLE of θ equals $\hat{\theta}_n = \overline{X}_n$.
- ▶ Let $\tau = e^\theta$.
- ▶ Then MLE for τ equals $\hat{\tau}_n = e^{\hat{\theta}_n} = e^{\overline{X}}$.

Asymptotic Normality

► Let $s(X; \theta) = \frac{\partial \log f(X; \theta)}{\partial \theta}$ be a *score function*.

► Then Fisher information is defined as:

$$I_n(\theta) = \mathbb{V}_\theta \left(\sum_{i=1}^n s(X_i; \theta) \right) = \sum_{i=1}^n \mathbb{V}_\theta (s(X_i; \theta)).$$

Lemma

Let $f(x; \theta)$ be continuously differentiable in θ . Then it holds

$$\mathbb{E}_\theta(s(X; \theta)) = 0 \quad \text{and} \quad \mathbb{V}_\theta(s(X; \theta)) = \mathbb{E}_\theta(s^2(X; \theta)).$$

Proof: We know that $1 = \int f(x; \theta) dx$.

Then

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int f(x; \theta) dx = \int \frac{\partial}{\partial \theta} f(x; \theta) dx = \int \frac{\partial f(x; \theta)}{\partial \theta} \cdot \frac{1}{f(x; \theta)} f(x; \theta) dx \\ &= \int \frac{\partial \log f(x; \theta)}{\partial \theta} f(x; \theta) dx = \int s(x; \theta) f(x; \theta) dx = \mathbb{E}_\theta[s(X; \theta)]. \end{aligned}$$

Theorem

The following equality holds: $I_n(\theta) = nI(\theta)$.

Also,

$$I(\theta) = -\mathbb{E}_\theta \left(\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right) = - \int \left(\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right) f(x; \theta) dx.$$

Theorem

Let $se = \sqrt{\mathbb{V}(\hat{\theta}_n)}$. Under appropriate regularity conditions, the following hold:

1. $se \approx \sqrt{1/I_n(\theta_*)}$ and $\frac{\hat{\theta}_n - \theta_*}{se} \rightsquigarrow \mathcal{N}(0, 1)$,
2. Let $\widehat{se} = \sqrt{1/I_n(\hat{\theta}_n)}$. Then $\frac{\hat{\theta}_n - \theta_*}{\widehat{se}} \rightsquigarrow \mathcal{N}(0, 1)$.

Asymptotic Normality

Proof:

Let $\ell_n(\theta) = \log \mathcal{L}_n(\theta)$.

Then $0 = \ell'_n(\hat{\theta}_n) \approx \ell'_n(\theta_*) + (\hat{\theta}_n - \theta_*)\ell''_n(\theta_*)$.

We obtain

$$\begin{aligned}\hat{\theta}_n - \theta_* &\approx -\ell'_n(\theta_*)/\ell''_n(\theta_*), \\ \sqrt{n}(\hat{\theta}_n - \theta_*) &\approx \frac{\frac{1}{\sqrt{n}}\ell'_n(\theta_*)}{-\frac{1}{n}\ell''_n(\theta_*)}.\end{aligned}\tag{1}$$

Asymptotic Normality

Let $Y_i = \frac{\partial \log f(X_i; \theta)}{\partial \theta}$.

From the Lemma on the Slide 32 it follows that $\mathbb{E}(Y_i) = 0$ and $\mathbb{V}(Y_i) = I(\theta_*)$.

Then, according to CLT, for the numerator in (1) it holds that

$$n^{-1/2} \sum_{i=1}^n Y_i = \sqrt{n} \bar{Y} = \sqrt{n}(\bar{Y} - 0) \rightsquigarrow W \sim \mathcal{N}(0, I(\theta_*)).$$

Asymptotic Normality

Define $A_i = -\partial^2 \log f(X_i; \theta) / \partial \theta^2$. Then $\mathbb{E}(A_i) = I(\theta_*)$ and for the denominator in (1) it holds that $\overline{A} \xrightarrow{\mathbb{P}} I(\theta_*)$.

Hence,

$$\sqrt{n}(\hat{\theta}_n - \theta_*) \rightsquigarrow \frac{W}{I(\theta_*)} \stackrel{d}{=} \mathcal{N}\left(0, \frac{1}{I(\theta_*)}\right).$$

Assume that $I(\theta)$ is a continuous function of its argument, then $I(\hat{\theta}_n) \xrightarrow{\mathbb{P}} I(\theta_*)$ and

$$\frac{\hat{\theta}_n - \theta_*}{\widehat{se}} = \sqrt{n}I^{1/2}(\hat{\theta}_n)(\hat{\theta}_n - \theta_*) = \left\{ \sqrt{n}I^{1/2}(\theta_*)(\hat{\theta}_n - \theta_*) \right\} \sqrt{\frac{I(\hat{\theta}_n)}{I(\theta_*)}}.$$

The first factor converges in distribution to $\mathcal{N}(0, 1)$, the second – to 1. \square

Cramér–Rao bound

Theorem

Let $\hat{\theta}_n = g(X_1, \dots, X_n)$ be an unbiased estimate of θ_* . Then

$$\mathbb{V}(\hat{\theta}_n) \geq \frac{1}{I_n(\theta_*)}.$$

Proof: Since estimate $\hat{\theta}_n$ is unbiased for $\theta = \theta_*$, we have:

$$\begin{aligned} 1 &= \int g(x_1, \dots, x_n) \frac{\partial f(x_1, \dots, x_n; \theta)}{\partial \theta} dx_1 \dots dx_n \\ &= \int g(x_1, \dots, x_n) \frac{\partial \log f(x_1, \dots, x_n; \theta)}{\partial \theta} f(x_1, \dots, x_n; \theta) dx_1 \dots dx_n \\ &= \mathbb{E}_\theta \hat{\theta}_n s(X_1, \dots, X_n; \theta) = \mathbb{E}_\theta \hat{\theta}_n \left[\sum_{i=1}^n s(X_i; \theta) \right]. \end{aligned}$$

Cramér–Rao bound

Proof (continued): Hence we arrive at:

$$\mathbb{E}_\theta \hat{\theta}_n \left[\sum_{i=1}^n s(X_i; \theta) \right] = 1.$$

But, since $\mathbb{E}_\theta(s(X; \theta)) = 0$:

$$\mathbb{E}_\theta \hat{\theta}_n \left[\sum_{i=1}^n s(X_i; \theta) \right] = \mathbb{E}_\theta (\hat{\theta}_n - \theta) \left[\sum_{i=1}^n s(X_i; \theta) \right] = \text{Cov}_\theta(\hat{\theta}_n, \sum_{i=1}^n s(X_i; \theta)).$$

Using $I_n(\theta) = \mathbb{V}_\theta(\sum_{i=1}^n s(X_i; \theta))$ and $\text{Cov}_\theta(X, Y) \leq \sqrt{\mathbb{V}_\theta(X)\mathbb{V}_\theta(Y)}$ we get

$$1 \leq \mathbb{V}_\theta(\hat{\theta}_n) I_n(\theta). \quad \square$$

- ▶ Let $X_1, \dots, X_n \sim \mathcal{N}(\theta, \sigma^2)$.
- ▶ MLE for θ is $\hat{\theta}_n = \bar{X}_n$.
- ▶ Denote the *sample median* by $\tilde{\theta}_n$, it can also be used to estimate θ .

$$\begin{aligned}\sqrt{n}(\hat{\theta}_n - \theta_*) &\rightsquigarrow \mathcal{N}(0, \sigma^2); \\ \sqrt{n}(\tilde{\theta}_n - \theta_*) &\rightsquigarrow \mathcal{N}\left(0, \sigma^2 \frac{\pi}{2}\right).\end{aligned}$$

- ▶ More generally, consider two estimators T_n and U_n and suppose that

$$\begin{aligned}\sqrt{n}(T_n - \theta_*) &\rightsquigarrow \mathcal{N}(0, t^2); \\ \sqrt{n}(U_n - \theta_*) &\rightsquigarrow \mathcal{N}(0, u^2).\end{aligned}$$

Optimality

We define the *asymptotic relative efficiency* of U_n to T_n by $ARE(U, T) = t^2/u^2$.

In the Normal example above, $ARE(\tilde{\theta}_n, \hat{\theta}_n) = 2/\pi = 0.63$.

Asymptotic relative efficiency can be interpreted as a fraction of the data that is “effectively” used for the estimation.

Theorem

Let $\hat{\theta}_n$ be the MLE and $\tilde{\theta}_n$ is any other estimator. Then, assuming appropriate regularity conditions (e.g. that the assumed parametric model is true), this

$$ARE(\tilde{\theta}_n, \hat{\theta}_n) \leq 1.$$

Thus, MLE has the smallest (asymptotic) variance and we say that the MLE is efficient or asymptotically optimal.

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

Delta-method

- ▶ Let $\tau = g(\theta)$ where g is a smooth function.
- ▶ MLE for τ is $\hat{\tau} = g(\hat{\theta})$.
- ▶ What is the distribution¹ of $\hat{\tau}$?

Theorem

If $\tau = g(\theta)$ where g is differentiable and $g'(\theta) \neq 0$, then

$$\frac{\hat{\tau}_n - \tau_*}{\widehat{se}(\hat{\tau}_n)} \rightsquigarrow \mathcal{N}(0, 1),$$

where $\hat{\tau}_n = g(\hat{\theta}_n)$ and $\widehat{se}(\hat{\tau}_n) = |g'(\hat{\theta}_n)| \widehat{se}(\hat{\theta}_n)$.

¹Distribution of an estimator is also called its sampling distribution

Theorem (continued)

Proof:

$$\begin{aligned}\hat{\tau}_n &= g(\hat{\theta}_n) \approx g(\theta_*) + (\hat{\theta}_n - \theta_*)g'(\theta_*) = \tau_* + (\hat{\theta}_n - \theta_*)g'(\theta_*), \\ \sqrt{n}(\hat{\tau}_n - \tau_*) &\approx \sqrt{n}(\hat{\theta}_n - \theta_*)g'(\theta_*), \\ \frac{\sqrt{nI(\theta_*)}(\hat{\tau}_n - \tau_*)}{g'(\theta_*)} &\approx \sqrt{nI(\theta_*)}(\hat{\theta}_n - \theta_*).\end{aligned}$$

(finished on the next slide)

Theorem (continued)

We know that $\sqrt{nI(\theta_*)}(\hat{\theta}_n - \theta_*)$ converges in distribution to $\mathcal{N}(0, 1)$. Then

$$\frac{\sqrt{nI(\theta_*)}(\hat{\tau}_n - \tau_*)}{g'(\theta_*)} \rightsquigarrow \mathcal{N}(0, 1).$$

From this we can conclude that

$$\hat{\tau}_n \approx \mathcal{N}(\tau_*, se^2(\hat{\tau}_n)), \quad se^2(\hat{\tau}_n) = \frac{(g'(\theta_*))^2}{nI(\theta_*)}.$$

If we replace θ_* with $\hat{\theta}_n$ the result will still hold.

Example

Example

► Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$.

► Statistic to estimate:

$$\psi = g(p) = \log \frac{p}{1-p}.$$

► Fisher information is

$$I(p) = \frac{1}{p(1-p)}.$$

► Estimate of the standard error is

$$\widehat{se} = \sqrt{\frac{\widehat{p}_n(1 - \widehat{p}_n)}{n}}.$$

Example

Example (continued)

MLE of ψ is

$$\hat{\psi}_n = \log \frac{\hat{p}_n}{1 - \hat{p}_n}.$$

Since $g'(p) = 1/(p(1 - p))$, according to the delta method we obtain

$$\widehat{se}(\hat{\psi}_n) = |g'(\hat{p}_n)| \widehat{se}(\hat{p}_n) = \frac{1}{\sqrt{n\hat{p}_n(1 - \hat{p}_n)}}.$$

Example

- ▶ Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$.
- ▶ Suppose that μ is known and σ is unknown.
- ▶ We want to estimate $\psi = \log \sigma$.
- ▶ The log-likelihood is

$$\ell(\sigma) = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2.$$

(continued on the next slide)

Example

Example (continued)

Taking the derivative and equating it to zero we get:

$$\hat{\sigma}_n = \sqrt{\frac{\sum_{i=1}^n (X_i - \mu)^2}{n}}.$$

To get the standard error we need the Fisher information. First,

$$\begin{aligned}\log f(X; \sigma) &= -\log \sigma - \frac{(X - \mu)^2}{2\sigma^2}, \\ \frac{\partial^2 \log f(X; \sigma)}{\partial \sigma^2} &= \frac{1}{\sigma^2} - \frac{3(X - \mu)^2}{\sigma^4}, \\ I(\sigma) &= -\frac{1}{\sigma^2} + \frac{3\sigma^2}{\sigma^4} = \frac{2}{\sigma^2}.\end{aligned}$$

(continued on the next slide)

Example

Example (continued)

- ▶ Thus, we obtain

$$\widehat{se} = \frac{\hat{\sigma}_n}{\sqrt{2n}}.$$

- ▶ Let $\psi = g(\sigma) = \log \sigma$, then $\hat{\psi}_n = \log \hat{\sigma}_n$.
- ▶ Since $g'(\sigma) = 1/\sigma$, we have

$$\widehat{se}(\hat{\psi}_n) = \frac{1}{\hat{\sigma}_n} \frac{\hat{\sigma}_n}{\sqrt{2n}} = \frac{1}{\sqrt{2n}}.$$

Problem Statement

Method of Moments

Maximum Likelihood Estimation

Delta Method

Multiparameter Models

We can extend the idea of the delta method to the models with several parameters. Let $\theta = (\theta_1, \dots, \theta_k)$ and let $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ be the MLE for θ . Log-likelihood will be

$$\ell_n(\theta) = \sum_{i=1}^n \log f(X_i; \theta).$$

Let

$$H_{jj} = \frac{\partial^2 \ell_n}{\partial \theta_j^2}, \quad H_{jk} = \frac{\partial^2 \ell_n}{\partial \theta_j \partial \theta_k}.$$

Define the Fisher Information Matrix by

$$I_n(\theta) = - \begin{pmatrix} \mathbb{E}_\theta(H_{11}) & \mathbb{E}_\theta(H_{12}) & \cdots & \mathbb{E}_\theta(H_{1k}) \\ \mathbb{E}_\theta(H_{21}) & \mathbb{E}_\theta(H_{22}) & \cdots & \mathbb{E}_\theta(H_{2k}) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}_\theta(H_{k1}) & \mathbb{E}_\theta(H_{k2}) & \cdots & \mathbb{E}_\theta(H_{kk}) \end{pmatrix}.$$

Define the precision matrix:

$$J_n(\theta) = I_n^{-1}(\theta).$$

Theorem

Under appropriate regularity conditions,

$$\hat{\theta} - \theta_* \rightsquigarrow \mathcal{N}(0, J_n).$$

Also, if $\hat{\theta}_j$ is the j -th component of vector $\hat{\theta}$

$$\frac{\hat{\theta}_j - \theta_{j,*}}{\widehat{se}_j} \rightsquigarrow \mathcal{N}(0, 1),$$

where $\widehat{se}_j^2 = J_n(j, j)$ is the j -th diagonal element of the matrix J_n .

The approximate covariance is

$$\text{Cov}(\hat{\theta}_j, \hat{\theta}_k) \approx J_n(j, k).$$

Let $\tau = g(\theta_1, \dots, \theta_k)$ be a function of the parameters and let its gradient be

$$\nabla g = \begin{pmatrix} \frac{\partial g}{\partial \theta_1} \\ \vdots \\ \frac{\partial g}{\partial \theta_k} \end{pmatrix}.$$

Theorem

Suppose that $\nabla g(\hat{\theta}) \neq 0$. Let $\hat{\tau} = g(\hat{\theta})$. Then

$$\frac{\hat{\tau} - \tau_*}{\widehat{se}(\hat{\tau})} \rightsquigarrow \mathcal{N}(0, 1),$$

where $\widehat{se}(\hat{\tau}) = \sqrt{(\hat{\nabla} g)^T \hat{J}_n (\hat{\nabla} g)}$, $\hat{J}_n = J_n(\hat{\theta})$, $\hat{\nabla} g = \nabla g(\hat{\theta})$.

Example

Let $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$, $\tau = g(\mu, \sigma) = \sigma/\mu$. The Fisher Information Matrix is

$$I_n(\mu, \sigma) = \begin{pmatrix} \frac{n}{\sigma^2} & 0 \\ 0 & \frac{2n}{\sigma^2} \end{pmatrix},$$

$$J_n = I_n^{-1}(\mu, \sigma) = \frac{1}{n} \begin{pmatrix} \sigma^2 & 0 \\ 0 & \frac{\sigma^2}{2} \end{pmatrix}, \quad \nabla g = \begin{pmatrix} -\frac{\sigma}{\mu^2} \\ \frac{1}{\mu} \end{pmatrix},$$

$$\widehat{se}(\hat{\tau}) = \sqrt{(\hat{\nabla} g)^T \hat{J}_n (\hat{\nabla} g)} = \frac{1}{\sqrt{n}} \sqrt{\frac{1}{\hat{\mu}^4} + \frac{\hat{\sigma}^2}{2\hat{\mu}^2}}.$$

Thank you for your attention!