

Module 5: Statistical inference (II)

Siyue Yang

05/21/2022

Outline

This module we will review

- Basics of parametric inference
- Methods for generating parametric estimators
- Maximum likelihood estimators
- Delta method
- Optimization method for finding MLE in R (Newton-Raphson, EM algorithm)

Parametric inference

Definition (Parametric models)

$$\mathfrak{F} = \{f(x; \theta) : \theta \in \Theta\}$$

where the $\Theta \subset \mathbb{R}^k$ is the parameter space and $\theta = (\theta_1, \dots, \theta_k)$ is the parameter.

Goal of parametric inference

- estimate the parametric θ (assume we known the form of the density).

Parameter of interest and nuisance parameter

Often, we are interested in estimating some function $T(\theta)$.

For example, if $X \sim N(\mu, \sigma^2)$, then

- Parameters: $\theta = (\mu, \sigma)$
- Parameter space: $\Theta = \{(\mu, \sigma) : \mu \in \mathbb{R}, \sigma > 0\}$

If the goal is to estimate the μ then

- Parameter of interest: $T(\theta) = \mu$
- Nuisance parameter: σ

Methods for generating parametric estimators

- ① Method of moments
- ② Maximum likelihood

Method of moments

Suppose that the parameter $\theta = (\theta_1, \dots, \theta_k)$ has k components.

- For $1 \leq j \leq k$, define the j^{th} moment

$$\alpha_j \equiv \alpha_j(\theta) = \mathbb{E}_\theta(X^j) = \int x^j dF_\theta(x)$$

- The j^{th} sample moment

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

- The method of moments estimator $\hat{\theta}_n$

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

$$\vdots$$

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

Maximum likelihood

- Parametric model: $f(x; \theta)$, X_1, \dots, X_n iid
- Likelihood function

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

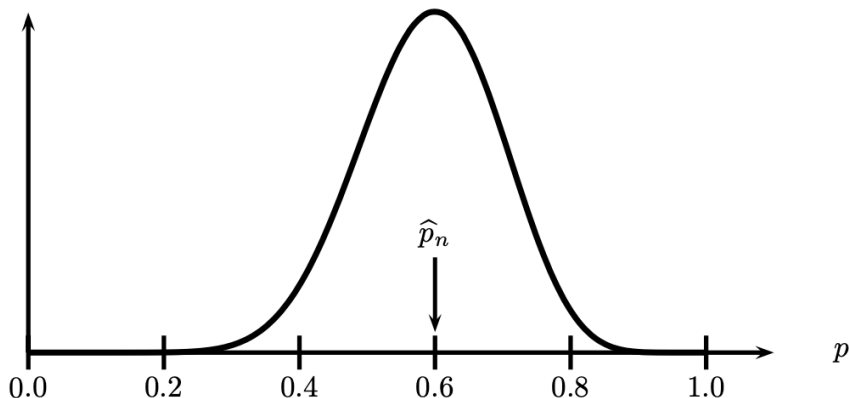
- The log-likelihood function

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

- The maximum likelihood estimator (MLE)

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \mathcal{L}(\theta)$$

An example of MLE



Likelihood function for Bernoulli with $n = 20$ and $\sum_{i=1}^n X_i = 12$. The MLE is $\hat{p}_n = 12/20 = 0.6$.

Why is maximum likelihood estimation so popular?

- A unified framework for estimation.
- Under mild regularity conditions, MLEs are
 - ① **consistent** \rightarrow converge to the true value in probability as $n \rightarrow \infty$, i.e.

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| \leq \epsilon) = 1 \quad \forall \epsilon > 0$$

- ② **asymptotically normal** $\rightarrow \sqrt{n}(\hat{\theta} - \theta) \sim N(0, \sigma^2)$ for large n
- ③ **asymptotically efficient** \rightarrow achieve the lowest variance for large n
- ④ **equivariant** \rightarrow if $\hat{\theta}$ is the MLE for θ then $g(\hat{\theta})$ is the MLE for $g(\theta)$

Steps to find the MLE

- 1 Write out the likelihood

$$\mathcal{L}(\theta) = f(X_1, \dots, X_n; \theta)$$

- 2 Simplify the log likelihood

$$\ell(\theta) = \log \mathcal{L}(\theta)$$

- 3 Take the derivative of $\ell(\theta)$ with respect to the parameter of interest, θ Set = 0
- 4 Solve for θ (get $\hat{\theta}_{MLE}$)
- 5 Check that $\hat{\theta}_{MLE}$ is a maximum ($\frac{\partial^2}{\partial \theta^2} \ell(\theta) < 0$)

Exercise

Suppose we have an iid sample $\{X_1, \dots, X_n\}$ with $X_i \sim \text{Bernoulli}(p)$. Find the MLE for p .

Exercise

Suppose we have an iid sample $\{X_1, \dots, X_n\}$ with $X_i \sim \text{Bernoulli}(p)$. Find the MLE for p .

1. The likelihood

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

where $S = \sum_i X_i$

Exercise

Suppose we have an iid sample $\{X_1, \dots, X_n\}$ with $X_i \sim \text{Bernoulli}(p)$. Find the MLE for p .

1. The likelihood

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

where $S = \sum_i X_i$

2. Log-likelihood

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

Exercise

Suppose we have an iid sample $\{X_1, \dots, X_n\}$ with $X_i \sim \text{Bernoulli}(p)$. Find the MLE for p .

1. The likelihood

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

where $S = \sum_i X_i$

2. Log-likelihood

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

3. MLE (Solved the scoring equation)

$$\ell'_n(p) = 0$$

The MLE is $\hat{p}_n = S/n$.

Score function and Fisher information

- Score function

$$s(X; \theta) = \frac{\partial \log f(X; \theta)}{\partial \theta}$$

- Fisher information

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

Asymptotic normality

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

- ① $se \approx \sqrt{1/I_n(\theta)}$ and

$$\frac{(\hat{\theta}_n - \theta)}{se} \rightsquigarrow N(0, 1).$$

- ② Let $\widehat{se} = \sqrt{1/I_n(\hat{\theta}_n)}$. Then,

$$\frac{(\hat{\theta}_n - \theta)}{\widehat{se}} \rightsquigarrow N(0, 1)$$

- Let

$$C_n = (\hat{\theta}_n - z_{\alpha/2}\widehat{se}, \hat{\theta}_n + z_{\alpha/2}\widehat{se})$$

Then, $\mathbb{P}_\theta(\theta \in C_n) \rightarrow 1 - \alpha$ as $n \rightarrow \infty$.

Elements of likelihood estimation

One random variable: Given a model for X which assumes X has a density $f(x; \theta)$, $\theta \in \Theta \subset \mathbb{R}^k$, we have the following definitions:

likelihood function

$$L(\theta; x) = c(x)f(x; \theta)$$

log-likelihood function

$$\ell(\theta; x) = \log L(\theta; x)$$

score function

$$u(\theta) = \partial \ell(\theta; x) / \partial \theta$$

observed information function

$$j(\theta) = -\partial^2 \ell(\theta; x) / \partial \theta \partial \theta^T$$

expected information (in one observation)

$$i(\theta) = E_{\theta} \left\{ U(\theta) U(\theta)^T \right\}$$

Elements of likelihood estimation (i.i.d.)

Independent observations: When we have X_i independent, identically distributed from $f(x_i; \theta)$, then, denoting the observed sample $\mathbf{x} = (x_1, \dots, x_n)$ we have:

likelihood function	$L(\theta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \theta)$
log-likelihood function	$\ell(\theta) = \ell(\theta; \mathbf{x}) = \sum_{i=1}^n \ell(\theta; x_i)$
maximum likelihood estimate	$\hat{\theta} = \hat{\theta}(\mathbf{x}) = \arg \sup_{\theta} \ell(\theta)$
score function	$U(\theta) = \ell'(\theta) = \sum U_i(\theta)$
observed information function	$j(\theta) = -\ell''(\theta) = -\ell''(\theta; \mathbf{x})$
observed (Fisher) information	$j(\hat{\theta})$
expected (Fisher) information	$i(\theta) = E_{\theta} \left\{ U(\theta) U(\theta)^T \right\} = ni_1(\theta)$

Delta method

Theorem (The Delta Method).

Suppose that

$$\frac{\sqrt{n}(Y_n - \mu)}{\sigma} \rightsquigarrow N(0, 1)$$

and that g is a differentiable function such that $g'(\mu) \neq 0$. Then

$$\frac{\sqrt{n}(g(Y_n) - g(\mu))}{|g'(\mu)|\sigma} \rightsquigarrow N(0, 1).$$

In other words,

$$Y_n \approx N\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{implies that} \quad g(Y_n) \approx N\left(g(\mu), (g'(\mu))^2 \frac{\sigma^2}{n}\right).$$

Exercise

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ and let $\psi = g(p) = \log(p/(1-p))$.

Exercise

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ and let $\psi = g(p) = \log(p/(1-p))$.

The Fisher information function is $I(p) = 1/(p(1-p))$

Exercise

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ and let $\psi = g(p) = \log(p/(1-p))$.

The Fisher information function is $I(p) = 1/(p(1-p))$

The estimated standard error of the MLE \hat{p}_n is

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

Exercise

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ and let $\psi = g(p) = \log(p/(1-p))$.

The Fisher information function is $I(p) = 1/(p(1-p))$

The estimated standard error of the MLE \hat{p}_n is

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

The MLE of ψ is $\hat{\psi} = \log \hat{p}/(1-\hat{p})$. Since, $g'(p) = 1/(p(1-p))$, according to the delta method

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

Exercise

Let $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ and let $\psi = g(p) = \log(p/(1-p))$.

The Fisher information function is $I(p) = 1/(p(1-p))$

The estimated standard error of the MLE \hat{p}_n is

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

The MLE of ψ is $\hat{\psi} = \log \hat{p}/(1-\hat{p})$. Since, $g'(p) = 1/(p(1-p))$, according to the delta method

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

An approximate 95 percent confidence interval is

$$\hat{\psi}_n \pm \frac{2}{\sqrt{n\hat{p}_n(1-\hat{p}_n)}}$$

Sometimes, there is no closed-form solution, so we need to use optimization methods to find the maximum of the log-likelihood.

- `optim()` find values of some parameters that **minimizes** some function.
- Newton-Raphson
- EM-algorithm

Newton-Raphson

Derivative of the log-likelihood around θ^3 :

$$0 = \ell'(\hat{\theta}) \approx \ell'(\theta^j) + (\hat{\theta} - \theta^j) \ell''(\theta^j)$$

Solving for $\hat{\theta}$ gives

$$\hat{\theta} \approx \theta^j - \frac{\ell'(\theta^j)}{\ell''(\theta^j)}.$$

This suggests the following iterative scheme:

$$\hat{\theta}^{j+1} = \theta^j - \frac{\ell'(\theta^j)}{\ell''(\theta^j)}$$

In the multiparameter case, the mle $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ is a vector and the method becomes

$$\hat{\theta}^{j+1} = \theta^j - H^{-1} \ell'(\theta^j)$$

where $\ell'(\theta^j)$ is the vector of first derivatives and H is the matrix of second derivatives of the log-likelihood.

Expectation-Maximization (EM) algorithm

Idea: Iterate between taking an expectation then maximizing.

Suppose we have data Y whose density $f(y; \theta)$ leads to a log-likelihood that is hard to maximize. However we can find another variable Z s.t. $f(y; \theta) = \int f(y, z; \theta) dz$ and $f(y, z; \theta)$ is easy to maximize.

- Pick a starting value θ^0 . Now for $j = 1, 2, \dots$, repeat steps E and M below:
- (The E-step): Calculate

$$J(\theta | \theta^j) = \mathbb{E}_{\theta^j} \left(\log \frac{f(Y^n, Z^n; \theta)}{f(Y^n, Z^n; \theta^j)} \mid Y^n = y^n \right).$$

The expectation is over the missing data Z^n treating θ^j and the observed data Y^n as fixed.

- (M-step) Find θ^{j+1} to maximize $J(\theta | \theta^j)$

Resources

This tutorial is based on

- Harvard Biostatistics Summer Pre Course [\[link\]](#)
- “All of Statistics” by Larry A. Wasserman [\[link\]](#)