

Chapter 3 Point Estimate

BY YUEJIAN MO

May 21, 2018

Here two method to estimate parameters of pdf/pmf of r.v.: maximum likelihood estimation, method of moments and Bayesian estimation.

1 Maximum Likelihood Estimator 最大似然估计

1.1 Point estimator and point estimate

Definition 3.1 (A statistic 统计量). A function of one or more r.v's that does not depend on the unknown parameter vector is called statistic.

对于统计量，一旦样本确定，统计量的值也就定下来，但其分布往往包含未知参数。Following show the mean of similar words.

- Estimation 估计(方法)
- Estimator 估计量
- Estimate 估计值
- population 母体, $\mathbf{X}(r, v)$
- sample 组 x_1, x_2, \dots, x_n (iid)

1.2 Joint density and likelihood function

In formal contexts, “likelihood” is often used as a synonym “probability”. In mathematical statistics, the two terms have different meaning. *Probability* in this technical context describes the plausibility of a future outcome, given a model parameter value, without reference to any observed data. *Likelihood* describes the plausibility of a model parameter value, given specific observed data.

Since \mathbf{x} has been observed and its components are therefore fixed real numbers, we regard $f(\mathbf{x}; \boldsymbol{\theta})$ as a function of $\boldsymbol{\theta}$, and define

$$L(\boldsymbol{\theta}) = L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i; \boldsymbol{\theta}), \boldsymbol{\theta} \in \boldsymbol{\Theta},$$

as the *likelihood function* of the random sample \mathbf{x} . It also can be called: $L(\boldsymbol{\theta})$ is the likelihood function of $\boldsymbol{\theta}$.

For avoid the operation of \prod , we has *log-likelihood*

$$l(\boldsymbol{\theta}) \triangleq \log\{L(\boldsymbol{\theta})\} = \sum_{i=1}^n \log\{f(x_i; \boldsymbol{\theta})\} \text{ for } \boldsymbol{\theta} \in \boldsymbol{\Theta}$$

There is no loss of information in using $l(\boldsymbol{\theta})$ instead of $L(\boldsymbol{\theta})$ because $\log(\cdot)$ is a monotonic increasing function.

1.3 Maximum likelihood estimator and maximum likelihood estimate

To get reasonable θ , we suppose that a statistic

$$\hat{\theta} = \begin{pmatrix} \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_n \end{pmatrix} = \begin{pmatrix} u_1(\mathbf{x}) \\ \vdots \\ u_n(\mathbf{x}) \end{pmatrix} \doteq \mathbf{u}(\mathbf{x})$$

satisfies

$$L(\hat{\theta}) = \max_{\theta \in \Theta} L(\theta).$$

We call $\hat{\theta} = \mathbf{u}(\mathbf{x})$ the *maximum likelihood estimator* (MLE) of θ and call $u(\mathbf{x})$ a *maximum likelihood estimate* (mle) of θ . There is no guarantee that the MLE exists or if it does whether it is unique.

1.4 The invariance property of MLE

Theorem 3.1: (Invariance of MLE). Let $\hat{\theta} = u(X_1, \dots, X_N)$ be the MLE of $\theta_{p \times 1} \in \Theta$. If $\eta_{p \times 1} = (h_1(\theta), \dots, h_p(\theta))^T$ is a one-to-one transformation between θ and η , then $\hat{\eta} = h(\hat{\theta})$ is the MLE of η .

Theorem 3.2 (Extension of Theorem 3.1): Let $\hat{\theta}$ be the MLE of $\theta = (\theta_1, \dots, \theta_p)^T \in \Theta$. If $\eta_{r \times 1} = h(\theta) = (h_1(\theta), \dots, h_r(\theta))^T$ for $1 \leq r \leq p$ is a many-to-few transformation between θ and η , then $\hat{\eta} = h(\hat{\theta}) = (h_1(\hat{\theta}), \dots, h_r(\hat{\theta}))^T$ is the MLE of η .

This property of invariance of MLEs allows us in our discussion of maximum likelihood estimation to consider estimating $(\theta_1, \dots, \theta_p)^T$ rather than the more general $h_1(\theta_1, \dots, \theta_p), \dots, h_r(\theta_1, \dots, \theta_p)$.

2 Moment Estimator

Method of moments is proposed by the great British statistician Karl Pearson near the turn of the twentieth century. If H_0 is rejected, one way is to guess another population distribution. Alternatively, we can estimate the first and second moments of the unknown population distribution $F(\cdot)$ by using the *method of moments*.

The method of moments can be applied to both *parametric* and *nonparametric* statistics.

3 Bayesian Estimator

4 Properties of Estimators

4.1 Unbiasedness

Definition 3.2 (Unbiased estimator and bias). An estimator $\varphi(\mathbf{x})$ is an *unbiased estimator* of the parameter θ if $E\{\varphi(\mathbf{x})\} = \theta$ for $\theta \in \Theta$. Otherwise, the estimator is biased and the bias is defined by

$$b(\theta) = E\{\varphi(\mathbf{x})\} - \theta$$

where $\mathbf{x} = (X_1, \dots, X_n)^T$.

Definition 3.3 (MSE). Given an estimator $Y = \varphi(\mathbf{x})$ of θ , the *mean square error* (MSE) of the estimator is defined by

$$\text{MSE} = E\{\varphi(\mathbf{x}) - \theta\}^2 = \text{Var}\{\varphi(\mathbf{x})\} + b^2(\theta)$$

If the estimator $\varphi(\mathbf{x})$ is unbiased, then $\text{MSE} = \text{Var}(\varphi(\mathbf{x}))$.

4.2 Efficiency

Maybe two estimators share the same bias for the same unknown parameter, so we have a notion of efficiency to choose the unbiased estimator with the *smaller* variance.

Definition 3.4 (Relative efficiency). Let $\hat{\theta}_1$ and $\hat{\theta}_2$ be two unbiased estimators for a parameter θ . If

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2),$$

we say that $\hat{\theta}_1$ is *more efficient* than $\hat{\theta}_2$. The *relative efficiency* of $\hat{\theta}_1$ to $\hat{\theta}_2$ is defined by the ratio

$$\frac{\text{Eff}_{\hat{\theta}_1}}{\text{Eff}_{\hat{\theta}_2}} = \frac{\text{Var}(\hat{\theta}_2)}{\text{Var}(\hat{\theta}_1)}$$

Let $u = \{\hat{\theta} : E(\hat{\theta}) = \theta\}$ denote the family of unbiased estimators of θ . We try to find the $\hat{\theta}^* \in u$ with the smallest variance. Here are found by Cramer and Rao, if we could find a constant c_0 satisfying

$$\text{Var}(\hat{\theta}) \geq c_0, \forall \hat{\theta} \in u,$$

thus, the $\hat{\theta}^*$ is equivalent to finding the lower bound c_0 .

Theorem 3.3 (The general CR inequality). Let $\tau(\theta)$ be an arbitrary function of the unknown θ . If (i) $\theta = T(\mathbf{x})$ is an unbiased estimator of $\tau(\theta)$, and (ii) the support of the population density $f(x; \theta)$ does not depend on the parameter θ , then

$$\text{Var}(\hat{\theta}) \geq \frac{\{\tau'(\theta)\}^2}{I_n(\theta)},$$

where $I_n(\theta)$ is the Fisher information. The right hand side is called the *Cramer–Rao lower bound*.

Theorem 3.4 (Alternative expression). Let $I_n(\theta)$ denote the information, If $E\{S(\theta)\} = 0$, then

$$I_n(\theta) = E\left\{-\frac{d^2 \log L(\theta; \mathbf{x})}{d\theta^2}\right\} = nI(\theta),$$

where

$$I(\theta) = E\left[\left\{\frac{d \log f(X; \theta)}{d\theta}\right\}^2\right] = E\left\{-\frac{d^2 \log f(X; \theta)}{d\theta^2}\right\}$$

denote the Fisher information for a single sample.

Definition 3.5 (UMVUE). An estimator θ^* is called a UMVUE of θ if it is unbiased and has the smallest variance among all unbiased estimators.

Definition 3.6 (Efficient estimator). If an unbiased estimator $\theta = T(\mathbf{x})$ for $\tau(\theta)$ has variance equal to the Cramer-Rao lower bound, then θ is called an *efficient estimator* for $\tau(\theta)$.

Chi-square distribution

Notation: $X \sim \chi^2(n)$

4.3 Sufficiency

Definition 3.7(Sufficient statistic). A statistic $T(\mathbf{x})$ is said to be a sufficient statistic of θ if the conditional distribution of \mathbf{x} , given $T(\mathbf{x})=t$, does not depend on θ for any value of t . In discrete case, this means that

$$\Pr\{X_1 = x_1, \dots, X_n = x_n; \theta | T(\mathbf{x}) = t\} = h(\mathbf{x})$$

Thm 3.5 (Factorization theorem) A statistic $T(\mathbf{x})$ is a sufficient statistic of the unknown parameter θ iff the joint pdf(or pmf) can be written in the form

$$f(x_1, \dots, x_n; \theta) = f(\mathbf{x}; \theta) = g(T(\mathbf{x}); \theta) \times h(\mathbf{x}),$$

Definition 3.8 (Joint sufficient statistics). Let $X_1, \dots, X_n \sim \text{iid } f(x; \theta)$. The statistics $T_1(\mathbf{x}), \dots, T_r(\mathbf{x})$ are said to be jointly sufficient if the conditional distribution of \mathbf{x} , given

4.4 Completeness

Definition 3.9 (Completeness). Let X_1, \dots, X_n denote a random sample from the pdf (or pmf) $f(x; \theta)$ with parameter space and let

Theorem 3.7 (Lehmann-Scheffe Theorem). Let $T(\mathbf{x})$ is a complete sufficient statistic for θ . If $g(T)$ is an unbiased estimator of $\tau(\theta)$, then $g(T)$ is the unique UMVUE for $\tau(\theta)$.

Fisher 信息量

5 Reference

- https://en.wikipedia.org/wiki/Likelihood_function
- <https://www.zhihu.com/question/33567579>