

Date: 5/7/2012

File Name: D:\A File\Research\Likelihood 12.7.23\Likelihood.16.tex

## Parametric Likelihood Inference

Xuan Yao

Maximum likelihood principle is one of the milestones in statistical literature in the past century. Here we give a brief review of the parametric likelihood inference. Throughout, we consider the following random sample from a known p.d.f. with unknown parameter  $\theta_0$ :

$$X_1, \dots, X_n \stackrel{i.i.d.}{\sim} f(x; \theta_0) \quad (1)$$

with the actual observations (realizations)

$$x_1, \dots, x_n. \quad (2)$$

## 1 Likelihood Function

*Likelihood* is the probability of observing the data we observed. Thus, for random sample (1) - (2) the likelihood is given by

$$P\{X_1 = x_1, \dots, X_n = x_n\} = \prod_{i=1}^n P\{X_i = x_i\}. \quad (3)$$

As follows, we discuss (3) for discrete and continuous p.d.f., respectively.

**Case 1:** If  $f(x; \theta_0)$  in (1) is discrete, we have  $P\{X = x\} = f(x; \theta_0)$ ; in turn, equation (3) becomes

$$P\{X_1 = x_1, \dots, X_n = x_n\} = \prod_{i=1}^n f(x_i; \theta_0). \quad (4)$$

**Case 2:** If  $f(x; \theta_0)$  in (1) is continuous, for a small constant  $\delta > 0$ , we have  $P\{X = x\} \approx P\{x - \delta < X < x + \delta\}$ ; in turn, equation (3) becomes

$$\begin{aligned} & P\{X_1 = x_1, \dots, X_n = x_n\} \\ & \approx \prod_{i=1}^n P\{x_i - \delta < X_i < x_i + \delta\} = \prod_{i=1}^n [F_X(x_i + \delta; \theta_0) - F_X(x_i - \delta; \theta_0)] \\ & = \prod_{i=1}^n [2\delta f(\xi_i; \theta_0)] = (2\delta)^n \prod_{i=1}^n f(\xi_i; \theta_0) \\ & \approx (2\delta)^n \prod_{i=1}^n f(x_i; \theta_0), \end{aligned} \quad (5)$$

where  $F_X(x; \theta_0)$  is the d.f. corresponding to  $f(x; \theta_0)$ ,  $\xi_i$  is between  $(x_i - \delta)$  and  $(x_i + \delta)$  and we assume  $f(x; \theta)$  is continuous in  $x$ . Thus, equation (5) shows that likelihood (3) is approximately proportional to  $\prod_{i=1}^n f(x_i; \theta_0)$ .

Based on (4) and (5), the *likelihood function* for  $\theta_0$  with random sample (1)-(2) is defined as

$$L(\theta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \theta), \text{ for } \theta \in \Theta, \quad (6)$$

where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $\Theta$  is the parameter space for  $\theta_0$  in (1). Note that for discrete or continuous p.d.f.  $f(x; \theta_0)$ , maximizing likelihood (3) and maximizing likelihood function (6) with respect to  $\theta$  are equivalent.

## 2 Maximum Likelihood Estimator

For random sample (1)-(2), *maximum likelihood estimator (MLE)* for  $\theta_0$  is given by

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta; \mathbf{x}), \quad (7)$$

where “arg max” is the value of argument at which the given function attains its maximum value.

Mathematically,  $\hat{\theta}$  is the value at which the likelihood function  $L(\theta; \mathbf{x})$  attains its maximum value. Recalling the relation between (3) and (6), statistically  $\hat{\theta}$  is the value of  $\theta$  in  $\Theta$  that makes the observed data has the greatest probability to be observed.

For many applications involving likelihood functions, it is more convenient to work in terms of natural logarithm of the likelihood function, called *log-likelihood*, than in terms of the likelihood function itself. Because the logarithm is a monotonically increasing function, the logarithm of a function achieves its maximum value at the same points as the function itself, and hence the log-likelihood can be used in place of the likelihood in maximum likelihood estimator and related techniques and we can write the MLE as

$$\hat{\theta} = \arg \max_{\theta \in \Theta} l(\theta; \mathbf{x}) = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \ln f(x_i; \theta), \quad (8)$$

where  $l(\theta; \mathbf{x}) = \ln L(\theta; \mathbf{x})$ .

If  $\Theta$  is open,  $l(\theta; \mathbf{x})$  is differentiable in  $\Theta$  and  $\hat{\theta}$  exists, then  $\hat{\theta}$  must satisfy the estimating equation

$$\nabla_{\theta} l(\theta; \mathbf{x}) = 0. \quad (9)$$

This is known as the *likelihood estimating equation*. So for the random sample (1)-(2), the likelihood estimating equation is given by

$$\sum_{i=1}^n \nabla_{\theta} \ln f(x_i; \theta) = 0. \quad (10)$$

Evidently, some solutions of (10) may not be the maxima or only the local maxima, thus we need to refer to other properties of the likelihood function. In the next two examples, we demonstrate how to find the MLE.

**Example:** Suppose the p.d.f. in (1) is given by  $f(x; \mu_0) = \exp\{-(x - \mu_0)^2/2\}/\sqrt{2\pi}$ , where  $x \in \mathcal{R}$  and  $\theta_0 = \mu_0$ . We find the MLE of  $\mu_0$  as follows. Using (6), we have

$$l(\mu; \mathbf{x}) = \sum_{i=1}^n \ln f(x_i; \mu) = -\frac{n}{2} \ln(2\pi) - \frac{1}{2} \sum_{i=1}^n (x_i - \mu)^2. \quad (11)$$

To maximize the log-likelihood, we differentiate (11) w.r.t.  $\mu$  and set the derivative to be zero,

$$\frac{\partial l(\mu; \mathbf{x})}{\partial \mu} = \sum_{i=1}^n (x_i - \mu) = 0. \quad (12)$$

The solution for (12) is  $\hat{\mu} = \sum_{i=1}^n x_i/n = \bar{\mathbf{x}}$ . Now, let us take the second derivative of (11),

$$\frac{\partial^2 l(\mu; \mathbf{x})}{\partial \mu^2} = -n < 0. \quad (13)$$

Thus we know that the first derivative of log-likelihood function is a decreasing. Since it attains 0 if and only if  $\hat{\mu} = \bar{\mathbf{X}}$ , the first derivative will be positive on  $(-\infty, \hat{\mu})$  and negative on  $(\hat{\mu}, \infty)$ . This means that the log-likelihood function is increasing on  $(-\infty, \hat{\mu})$  whereas decreasing on  $(\hat{\mu}, \infty)$ , thus the likelihood function attains its maximum value at  $\mu = \hat{\mu}$ . Hence the MLE for  $\mu$  is given by  $\hat{\mu} = \bar{\mathbf{X}}$ .

In some cases, the differentiating method is not applicable. This often happens when the domain of random sample depends on parameter.

**Example:** Suppose the p.d.f. in (1) given by uniform distribution on  $(0, \theta_0)$ . We find the MLE for  $\theta$  as follows. Since in this case, the p.d.f. is given by

$$f(x; \theta_0) = \frac{1}{\theta_0} I\{0 < x < \theta_0\}, \quad (14)$$

where  $I\{x \in A\}$  is the indicator function, i.e., for a given set  $A$ ,  $I\{x \in A\} = 1$  if  $x \in A$ ;  $I\{x \in A\} = 0$  otherwise. Thus from (6), the likelihood function is given by

$$\begin{aligned} L(\theta; \mathbf{x}) &= \prod_{i=1}^n f(x_i; \theta) = \prod_{i=1}^n \frac{1}{\theta} I\{0 < x_i < \theta\} \\ &= \frac{1}{\theta^n} \prod_{i=1}^n I\{0 < x_i < \theta\} = \frac{1}{\theta^n} I\{0 < x_1 < \theta, 0 < x_2 < \theta, \dots, 0 < x_n < \theta\} \\ &= \frac{1}{\theta^n} I\{0 < x_{(1)} < x_{(n)} < \theta\}, \end{aligned} \quad (15)$$

where  $X_{(i)}$  is the order statistic and  $x_{(i)}$  is the corresponding realization. Note that the support of  $L(\theta; \mathbf{x})$  is  $[x_{(n)}, \infty)$  and that on its support,  $L(\theta; \mathbf{x}) = 1/\theta^n$  is decreasing in  $\theta$ . Therefore  $L(\theta; \mathbf{x})$  attains its maximum value at  $\theta = X_{(n)}$ . Thus the MLE for  $\theta_0$  is given by  $\hat{\theta} = X_{(n)}$ .

### 3 Properties of MLE

Let us start this section with a convenient computational property for MLE, namely, plug-in property. Then we will present the asymptotic distribution, consistency and efficiency of MLE. At the end this section, we will discuss several disadvantages of MLE.

#### 3.1 Invariance Property

MLE holds a nice *invariance property*, which means that MLE is unaffected by re-parametrization, i.e., MLE is equivariant under one-to-one transformations.

**Theorem 3.1.** Let  $\hat{\theta}$  denote the MLE of  $\theta_0$  in random sample (1)-(2). Suppose that  $h(\cdot)$  is a one-to-one function from  $\Theta$  onto  $h(\Theta)$ . Define  $\eta \equiv h(\theta)$ . Then the MLE of  $\eta_0 = h(\theta_0)$  is  $h(\hat{\theta})$ .

**Proof:** Since  $h(\cdot)$  is onto and one-to-one,  $h^{-1}(\cdot)$  exists. Since  $\eta = h(\theta)$ , we have  $\theta = h^{-1}(\eta)$ . Hence

$$f(x; \theta_0) = f(x; h^{-1}(\eta_0)) \equiv f^*(x; \eta_0),$$

where  $f^*(x; \eta_0)$  is the p.d.f. of the random sample (1)-(2) with parameter  $\eta_0$ . Then by (6), the likelihood function for  $\eta_0$  is

$$L^*(\eta; \mathbf{x}) = \prod_{i=1}^n f^*(x_i; \eta) = \prod_{i=1}^n f(x_i; h^{-1}(\eta)) = L(h^{-1}(\eta); \mathbf{x}) = L(\theta; \mathbf{x}). \quad (16)$$

Let  $\hat{\eta}$  be the MLE for  $\eta_0$ , then we have

$$L^*(\hat{\eta}; \mathbf{x}) = \max_{\eta \in h(\Theta)} L^*(\eta; \mathbf{x}) \stackrel{(16)}{=} \max_{\eta = h(\theta) \in h(\Theta)} L(\theta; \mathbf{x}) = \max_{\theta \in \Theta} L(\theta; \mathbf{x}) = L(\hat{\theta}; \mathbf{x}), \quad (17)$$

Thus we have  $\hat{\eta} = h(\hat{\theta})$ . □

### 3.2 Consistency of MLE

A sequence of estimators  $W_n = W_n(X_1, \dots, X_n)$  is a *consistent sequence of estimators* (Casella and Berger, P468) of the parameter  $\theta$  if for every  $\epsilon > 0$  and every  $\theta \in \Theta$ ,

$$\lim_{n \rightarrow \infty} P_\theta\{|W_n - \theta| \geq \epsilon\} = 0. \quad (18)$$

The following theorem will show that MLE is consistent.

**Theorem 3.2.** *For random sample (1)-(2), let  $\hat{\theta}$  denote the MLE of  $\theta_0$ . Suppose that  $f(\theta; \mathbf{x})$  satisfies the following assumptions,*

A1  $f(x; \theta)$  is *identifiable*, i.e., if  $\theta \neq \theta'$ , then  $f(x; \theta) \neq f(x; \theta')$ .

A2 The densities  $f(x; \theta)$  have common support, and  $f(x; \theta)$  is differentiable in  $\theta$ .

A3 The parameter space  $\Omega$  contains an open set  $\omega$  of which the true parameter value  $\theta_0$  is an interior point.

*Then the MLE is consistent.*

**Proof:**

By A3 we know that there exist  $\delta > 0$  such that for any  $0 < \epsilon < \delta$ ,  $(\theta_0 - \epsilon, \theta_0 + \epsilon) \subset \Theta$ . Then by S.L.L.N., we have

$$\begin{aligned} & \frac{1}{n}(l(\theta_0 - \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) \\ &= \frac{1}{n} \sum_{i=1}^n (\ln f(x_i; \theta_0 - \epsilon) - \ln f(x_i; \theta_0)) \xrightarrow{a.s.} E_0\{\ln f(x_i; \theta_0 - \epsilon)\} - E_0\{\ln f(x_i; \theta_0)\} \\ &= E_0 \left\{ \ln \frac{f(x_i; \theta_0 - \epsilon)}{f(x_i; \theta_0)} \right\} \end{aligned} \quad (19)$$

and

$$\begin{aligned} & \frac{1}{n}(l(\theta_0 + \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) \\ &= \frac{1}{n} \sum_{i=1}^n (\ln f(x_i; \theta_0 + \epsilon) - \ln f(x_i; \theta_0)) \xrightarrow{a.s.} E_0\{\ln f(x_i; \theta_0 + \epsilon)\} - E_0\{\ln f(x_i; \theta_0)\} \\ &= E_0 \left\{ \ln \frac{f(x_i; \theta_0 + \epsilon)}{f(x_i; \theta_0)} \right\}. \end{aligned} \quad (20)$$

Apply Jensen's Inequality, we get

$$E_\theta \left\{ \ln \frac{f(x; \theta')}{f(x; \theta)} \right\} < \ln E_\theta \left\{ \frac{f(x; \theta')}{f(x; \theta)} \right\} = \ln \int \frac{f(x; \theta')}{f(x; \theta)} \cdot f(x; \theta) dx = \ln \int f(x; \theta') dx = 0 \quad (21)$$

The inequality in (21) is strict due to A1, that is  $f(x; \theta)$  being identifiable, i.e. for any  $\theta' \neq \theta$ , we have  $f(x; \theta') \neq f(x; \theta)$ . From (19)-(21), we know that

$$P \left\{ \lim_{n \rightarrow \infty} \frac{1}{n}(l(\theta_0 - \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0 \right\} = P \left\{ \lim_{n \rightarrow \infty} \frac{1}{n}(l(\theta_0 + \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0 \right\} = 1 \quad (22)$$

So  $\exists N$  such that

$$P \left\{ \frac{1}{n}(l(\theta_0 - \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0 \right\} = P \left\{ \frac{1}{n}(l(\theta_0 + \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0 \right\} = 1, \text{ for all } n > N. \quad (23)$$

Hence

$$P\{(l(\theta_0 - \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0\} = P\{(l(\theta_0 + \epsilon; \mathbf{x}) - l(\theta_0; \mathbf{x})) < 0\} = 1, \text{ for all } n > N, \quad (24)$$

Note that by A2,  $f(x; \theta)$  is differentiable on  $[\theta_0 - \epsilon, \theta_0 + \epsilon]$ . Recall that  $l(\theta_0; \mathbf{x}) = \sum_{i=1}^n \ln f(x_i, \theta_0)$ .  $l(\theta_0; \mathbf{x})$  is also differentiable and continuous on  $[\theta_0 - \epsilon, \theta_0 + \epsilon]$ . Hence there exist  $\hat{\theta}$  such that  $l(\hat{\theta}; \mathbf{x}) \geq l(\theta; \mathbf{x})$  for all  $\theta \in [\theta_0 - \epsilon, \theta_0 + \epsilon]$ . However by (24),  $\hat{\theta}$  is not equal to  $\theta_0 - \epsilon$  or  $\theta_0 + \epsilon$ , therefore we have

$$\left. \frac{\partial l(\theta; \mathbf{x})}{\partial \theta} \right|_{\hat{\theta}} = 0 \text{ and } \hat{\theta} \in (\theta_0 - \epsilon, \theta_0 + \epsilon) \text{ for all } n > N. \quad (25)$$

Consequently,  $\hat{\theta}$  is the MLE and for any  $\epsilon$ ,  $\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta_0| < \epsilon) = 1$ . Hence MLE is consistent.  $\square$

### 3.3 Asymptotic Distribution

A nice asymptotic distribution will simplify computation for large sample. The following theorem will show that MLE is asymptotically normal when sample size is sufficiently large.

**Theorem 3.3.** *Suppose  $\hat{\theta}_n$  is the MLE of true value  $\theta_0$  in random sample (1)-(2). Let  $I(\theta_0)$  denote the Fisher Information in  $X_1$ . Suppose that  $f(x; \theta)$ ,  $\theta \in \Theta$  satisfies the following three assumptions,*

A4 for any  $x$ ,  $f(x; \theta)$  is three times differentiable with respect to  $\theta$  in a small neighbourhood of the true value  $\theta_0$ .

A5 For  $\theta$  in a small neighbourhood of  $\theta_0$ ,  $|\partial^3 \ln f(x; \theta) / \partial \theta^3| \leq H(x)$  and  $E\{H(X)\} < \infty$ .

A6  $E_0\{\partial \ln f(X, \theta_0) / \partial \theta|_{\theta_0}\} = 0$ ;  $E_0\{\partial^2 \ln f(X, \theta_0) / \partial \theta^2|_{\theta_0}\} = -I(\theta_0)$ ;  $I(\theta_0) > 0$ .

*As  $n$  goes to infinity,  $\sqrt{n}(\hat{\theta}_n - \theta_0)$  goes to  $N(0, I^{-1}(\theta_0))$  in distribution.*

**Proof:** Let us make Taylor expansion of  $\partial l(\theta; \mathbf{x}) / \partial \theta|_{\hat{\theta}_n}$  at  $\theta = \theta_0$ ,

$$0 = \left. \frac{\partial l(\theta; \mathbf{x})}{\partial \theta} \right|_{\hat{\theta}_n} = \left. \frac{\partial l(\theta; \mathbf{x})}{\partial \theta} \right|_{\theta_0} + (\hat{\theta}_n - \theta_0) \left. \frac{\partial^2 l(\theta; \mathbf{x})}{\partial \theta^2} \right|_{\theta_0} + \frac{1}{2} (\hat{\theta}_n - \theta_0)^2 \left. \frac{\partial^3 l(\theta; \mathbf{x})}{\partial \theta^3} \right|_{\theta_1}, \quad (26)$$

where  $\theta_1$  is between  $\hat{\theta}_n$  and  $\theta_0$ . So

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = \left\{ -\sqrt{n} \cdot \left. \frac{1}{n} \frac{\partial l(\theta; \mathbf{x})}{\partial \theta} \right|_{\theta_0} \right\} / \left\{ \left. \frac{1}{n} \left[ \frac{\partial^2 l(\theta; \mathbf{x})}{\partial \theta^2} \right]_{\theta_0} + \frac{1}{2} (\hat{\theta} - \theta_0)^2 \frac{\partial^3 l(\theta; \mathbf{x})}{\partial \theta^3} \right]_{\theta_1} \right\}. \quad (27)$$

From A6, we can see that

$$E_0 \left\{ \left. \frac{\partial \ln f(X_1; \theta)}{\partial \theta} \right|_{\theta_0} \right\} = 0 \text{ and } \text{Var}_0 \left\{ \left. \frac{\partial \ln f(X_1; \theta)}{\partial \theta} \right|_{\theta_0} \right\} = I(\theta_0) > 0. \quad (28)$$

Recall that  $l(\theta; \mathbf{x}) = \sum_{i=1}^n \ln f(x_i; \theta)$  and  $\ln f(X_1; \theta), \ln f(X_2; \theta), \dots, \ln f(X_n; \theta)$  are i.i.d. random variables. Hence by C.L.T.,

$$-\sqrt{n} \cdot \left. \frac{1}{n} \frac{\partial l(\theta; \mathbf{x})}{\partial \theta} \right|_{\theta_0} \xrightarrow{D} N(0, I(\theta_0)). \quad (29)$$

Apply S.L.L.N, we obtain

$$\left. \frac{1}{n} \frac{\partial^2 l(\theta; \mathbf{x})}{\partial \theta^2} \right|_{\theta_0} \xrightarrow{a.s.} E_0 \left\{ \left. \frac{\partial^2 \ln f(X_1; \theta)}{\partial \theta^2} \right|_{\theta_0} \right\} \stackrel{A6}{=} -I(\theta_0) \quad (30)$$

Then note that  $\hat{\theta}_n \xrightarrow{P} \theta_0$  as  $n$  goes to infinity. Therefore

$$\theta_1 \xrightarrow{P} \theta_0 \text{ and } \theta_1 - \theta_0 \xrightarrow{P} 0 \text{ as } n \text{ goes to infinity} \quad (31)$$

Hence as  $n \rightarrow \infty$ ,  $\theta_1$  is in arbitrary small neighbourhood of  $\theta_0$ . Therefore by A2, we have

$$\frac{1}{n} \frac{\partial^3 l(\theta; \mathbf{x})}{\partial \theta^3} \Big|_{\theta_1} = O_p(1) \quad \text{and} \quad \frac{1}{2n} (\hat{\theta}_n - \theta_0) \frac{\partial^3 l(\theta; \mathbf{x})}{\partial \theta^3} \Big|_{\theta_1} \xrightarrow{P} 0. \quad (32)$$

From (30) and (32) we know that the denominator of (27) goes to  $-I(\theta_0)$  in probability while (29) shows that the numerator of (27) goes to  $N(0, I(\theta_0))$ . Consequently, by Slutsky's Theorem, we obtain

$$\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{D} N(0, I(\theta_0))/I(\theta_0) \stackrel{D}{=} N(0, I^{-1}(\theta_0)). \quad (33)$$

□

### 3.4 Asymptotic Efficiency

A sequence of estimators  $W_n$  is *asymptotically efficient* for a parameter  $\tau(\theta_0)$  if  $\sqrt{n}[W_n - \tau(\theta_0)] \rightarrow N(0, v(\theta_0))$  in distribution and

$$v(\theta_0) = \frac{\tau'^2(\theta_0)}{E_0 [(\partial \ln f(X; \theta_0)/\partial \theta)^2]}; \quad (34)$$

that is, the asymptotic variance of  $W_n$  achieves the Cramer-Rao Lower Bound.

**Theorem 3.4.** *Maximum likelihood estimator is asymptotically efficient.*

**Proof:** Note that in our case,  $W_n = \hat{\theta}_n$  and  $\tau(\theta_0) = \theta_0$ . Followed by Theorem 3.2 and (34), we get the conclusion. □

### 3.5 Disadvantages of MLE

Although MLE does hold some convenient mathematical properties (plug-in) and good asymptotic behaviour (asymptotic normal, consistency and efficiency), it also has some disadvantages.

1. All the good statistical behaviour are based on sufficiently large sample size. Actually, for small sample, MLE may be significantly biased. We may also lose efficiency when sample size is small.
2. We need to assume the distribution of random sample according to prior experience or knowledge. All the calculation, no matter for large sample or small sample, is based on the assumed p.d.f.  $f(x; \theta)$ . However, in practice, it is quite possible that the  $f(x; \theta)$  we propose is not close to the real distribution, which will cause a vital damage to the whole process.
3. To derive a convenient way to calculate MLE, we assumed independence among  $X_1, \dots, X_n$ . This assumption may also be violated in practise.
4. In some cases, maximum likelihood estimator does not necessary exist. Even it does exist and can be calculated by differentiating the likelihood function, the calculation might be very complex and will not lead to a explicit answer.
5. Sometimes we apply Newton-Raphson, EM and etc. to give a numerical solution to MLE. This calls for more regulation on parameter space and p.d.f.. These methods may also be sensitive to the initial point for iteration

## 4 Likelihood Ratio Test

A *likelihood ratio test (LRT)* is used to compare the fitness of two models, one of which (the null model) is a special case of the other (the alternative model). Suppose the parameter  $\theta_0$  in (1) belongs to a set  $\Theta$ , then LRT can be defined as follows.

**Definition 4.1.** *The likelihood ratio test statistic for testing  $H_0 : \theta_0 \in \Theta_0, \Theta_0 \subset \Theta$  vs  $H_1 : \theta_0 \in \Theta_0^c$  is*

$$\lambda(\mathbf{x}) = \frac{\sup_{\theta \in \Theta_0} L(\theta; \mathbf{x})}{\sup_{\theta \in \Theta} L(\theta; \mathbf{x})}. \quad (35)$$

A *likelihood ratio test (LRT)* is any test that has a rejection region of the form  $\{\mathbf{x} : \lambda(\mathbf{x}) \leq c_\alpha\}$ , where  $c$  is any number satisfying  $0 \leq c_\alpha \leq 1$ .

**Remark 4.1.** .

1. Since the numerator of (35) is maximized over a smaller region compared to the denominator, we can conclude that likelihood ratio is always smaller than one.
2. An optimized case is when null hypothesis is true. Recall that if we have a large sample the MLE is approximately equal to the true value. Hence the likelihood ratio will be close to one. Otherwise, likelihood ratio will be close to zero

The constant  $c_\alpha$  in Definition 4.1 is decided by the level of the test. For a test of level  $\alpha$ ,

$$\alpha = P(\text{reject } H_0 | H_0) = P_{\theta_0 \in \Theta_0}(\lambda(\mathbf{x}) < c_\alpha), \quad (36)$$

and the rejection region is  $(0, c_\alpha)$ , which means that if the likelihood ratio is smaller than  $c_\alpha$ , we should reject the null hypothesis with probability  $1 - \alpha$ .

A special case for (35) is testing  $H_0 : \theta_0 = \theta_0^*$  vs  $H_1 : \theta_0 \neq \theta_0^*$ . Further, let us suppose the MLE exists. Since we have only one candidate under null hypothesis,  $\lambda(\mathbf{x})$  becomes

$$\lambda(\mathbf{x}) = \frac{L(\theta_0^*; \mathbf{x})}{\sup_{\theta \in \Theta} L(\theta; \mathbf{x})} = \frac{L(\theta_0^*; \mathbf{x})}{L(\hat{\theta}; \mathbf{x})}. \quad (37)$$

The calculation of  $c_\alpha$  calls for an explicit distribution of  $\lambda(\mathbf{x})$ . LRT has a nice  $\chi_\nu^2$  distribution when we have a large sample size. Here and throughout this note, we use  $\chi_\nu^2$  to denote the chi square distribution with  $\nu$  degrees of freedom. Let us present this property starting with the simple  $H_0 : \theta_0 = \theta_0^*$  vs  $H_1 : \theta_0 \neq \theta_0^*$ . The following theorems are cited from Theorem 10.3.1 and Theorem 10.3.3 in [1].

**Theorem 4.1.** *Suppose  $\theta_0 \in \Theta \subset \mathcal{R}$ . For testing  $H_0 : \theta_0 = \theta_0^*$  vs  $H_1 : \theta_0 \neq \theta_0^*$ , with random samples (1)-(2). Then under  $H_0$ , as  $n \rightarrow \infty$ ,  $-2 \ln \lambda(\mathbf{x}) \rightarrow \chi_1^2$  in distribution.*

**Proof:** First expand  $\ln L(\theta; \mathbf{x}) = l(\theta; \mathbf{x})$  in a Taylor series around the MLE  $\hat{\theta}$ , giving,

$$l(\theta; \mathbf{x}) = l(\hat{\theta}; \mathbf{x}) + l'(\hat{\theta}; \mathbf{x})(\theta - \hat{\theta}) + \frac{1}{2}l''(\hat{\theta}; \mathbf{x})(\theta - \hat{\theta})^2 + \dots \quad (38)$$

Now substitute the expansion for  $l(\theta_0^*; \mathbf{x})$  in  $-2 \ln \lambda(\mathbf{x}) = -2l(\theta_0^*; \mathbf{x}) + 2 \ln(\hat{\theta}; \mathbf{x})$ , and get

$$-2 \ln \lambda(\mathbf{x}) \approx -l''(\hat{\theta}; \mathbf{x})(\theta_0^* - \hat{\theta}), \quad (39)$$

where we use the fact that  $l'(\hat{\theta}; \mathbf{x}) = 0$ . Since  $l''(\hat{\theta}; \mathbf{x})$  is the observed fisher information  $\hat{I}_n(\hat{\theta})$  and  $\hat{I}_n(\hat{\theta})/n \rightarrow I(\theta_0^*) = I(\theta_0)$  under  $H_0$ . It follows from Theorem 3.3 and Slutsky's Theorem that  $-2 \ln \lambda(\mathbf{x}) \rightarrow \chi_1^2$  in distribution.  $\square$

Theorem 4.1 can be extended to the cases where the null hypothesis concerns a vector of parameters. The following generalization, which we state without proof, allows us to ensure Theorem 4.1 is true for large samples.

**Theorem 4.2. (Wilk's Theorem)** For testing  $H_0 : \theta_0 \in \Theta_0$  vs  $H_1 : \theta_0 \in \Theta_0^c$ , suppose random samples are (1)-(2) and  $\theta_0 \in \Theta$ . Then under  $H_0$ , as  $n \rightarrow \infty$ ,  $-2 \ln \lambda(\mathbf{x}) \rightarrow \chi_\nu^2$  in distribution, where  $\nu$  equals to the difference between the number of free parameters specified by  $\theta \in \Theta_0$  and the number of free parameters specified by  $\theta \in \Theta$ .

The computation of  $\nu$  is usually straight forward. Most often,  $\Theta$  can be represented as a subset of  $q$ -dimensional Euclidean space that contains an open subset in  $\mathcal{R}^q$ , and  $\Theta_0$  can be represented as a subset of  $p$ -dimensional Euclidean space that contains an open subset in  $\mathcal{R}^p$ , where  $p < q$ . Then  $\nu = q - p$  is the degrees of freedom for the test statistic.

Rejection of  $H_0$  for small values of  $\lambda(\mathbf{x})$  is equivalent to rejection for large values of  $-2 \ln \lambda(\mathbf{x})$ . Thus,

$$H_0 \text{ is rejected if and only if } -2 \ln \lambda(\mathbf{x}) \geq \chi_{\nu, \alpha}^2, \quad (40)$$

and the asymptotic rejection region is  $(\chi_{\nu, \alpha}^2, \infty)$ . Here and through out this note  $\chi_{\nu, \alpha}^2$  is the constant such that  $P\{\chi_\nu^2 > \chi_{\nu, \alpha}^2\} = \alpha$ .

Next let us present an example using Wilk's Theorem.

**Example 4.1.** Let  $\theta = (p_1, p_2, p_3, p_4, p_5)$ , where  $p_j$ 's are non-negative and sum to 1. For (1)-(2),  $f(j; \theta) = p_j, j = 1, \dots, 5$ . Find the LRT test statistic for testing  $H_0 : p_1 = p_2 = p_3$  and  $p_4 = p_5$  vs  $H_1 : H_0$  is not true, and the asymptotic rejection region.

**Sol 4.1.** The likelihood function under  $\Theta$  is

$$L(\theta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \theta) = \prod_{i=1}^n p_i^{y_i}, \text{ where } y_j = \text{number of } x_1, \dots, x_n \text{ equal to } j. \quad (41)$$

The full parameter space,  $\Theta$ , with four free parameters, is really a four-dimensional set since  $p_5 = 1 - p_1 - p_2 - p_3 - p_4$ . The parameter set is defined by

$$\sum_{j=1}^4 p_j \leq 1 \text{ and } p_j \geq 0, \quad j = 1, \dots, 4, \quad (42)$$

a subset of  $\mathcal{R}^4$  containing an open subset of  $\mathcal{R}^4$ . Thus  $q = 4$ . There is only one free parameter in the set specified by  $H_0$  because once  $p_1$  is fixed,  $p_2 = p_3$  must equal to  $p_1$  and  $p_4 = p_5$  must equal  $(1 - 3p_1)/2$ . Thus  $p = 1$  and the degrees of freedom is  $\nu = 4 - 1 = 3$ .

To calculate  $\lambda(\mathbf{x})$ , the MLE of  $\theta$  under both  $\Theta_0$  and  $\Theta$  must be determined. By setting

$$\frac{\partial}{\partial p_j} l(\theta; \mathbf{x}) = 0, \text{ for each of } j = 1, \dots, 4 \quad (43)$$

and using the facts that  $p_5 = 1 - p_1 - p_2 - p_3 - p_4$  and  $y_5 = n - y_1 - y_2 - y_3 - y_4$ , we can verify that the MLE of  $p_j$  under  $\Theta$  is  $\hat{p}_j = y_j/n$ . Under  $H_0$ , the likelihood function reduces to

$$L(\theta; \mathbf{x}) = \prod_{i=1}^n f(x_i; \theta) = p_1^{y_1 + y_2 + y_3} \left( \frac{1 - 3p_1}{2} \right)^{y_4 + y_5}. \quad (44)$$

Using the same method as (43), the MLE's under  $H_0$  are  $p_{10} = p_{20} = p_{30} = (y_1 + y_2 + y_3)/(3n)$  and  $p_{40} = p_{50} = (1 - 3\hat{p}_{10})/2$ . Substituting these values and the  $\hat{p}_j$  values into  $L(\theta; \mathbf{x})$  and combining terms with the same exponent yield

$$\lambda(\mathbf{x}) = \left( \frac{y_1 + y_2 + y_3}{3y_1} \right)^{y_1} \left( \frac{y_1 + y_2 + y_3}{3y_2} \right)^{y_2} \left( \frac{y_1 + y_2 + y_3}{3y_3} \right)^{y_3} \left( \frac{y_4 + y_5}{2y_4} \right)^{y_4} \left( \frac{y_4 + y_5}{2y_5} \right)^{y_5}. \quad (45)$$



Thus the test statistic is

$$-2 \ln \lambda(\mathbf{x}) = 2 \sum_{i=1}^5 y_i \ln \left( \frac{y_i}{m_i} \right), \quad (46)$$

where  $m_1 = m_2 = m_3 = (y_1 + y_2 + y_3)/3$  and  $m_4 = m_5 = (y_4 + y_5)/2$ . The asymptotic size  $\alpha$  test rejects  $H_0$  if  $-2 \ln \lambda(\mathbf{x}) \geq \chi_{3,\alpha}^2$ .

Although likelihood ratio test is not necessarily unbiased, we can approach the unbiasedness by increasing sample size. In other words, likelihood ratio test is consistent.

**Theorem 4.3.** *The likelihood ratio test is consistent.*

**Proof:** We need to show that if true value  $\theta_0 \neq \theta_0^*$ , we reject  $H_0$  with probability one as  $n$  goes to infinity. We reject the null hypothesis if  $\lambda(\mathbf{x}) < c$ , or equivalently, if

$$-\ln \lambda(\mathbf{x}) = \sum_{i=1}^n \ln f(x_i; \hat{\theta}_n) - \sum_{i=1}^n \ln f(x_i; \theta_0^*) > c. \quad (47)$$

Expand the first term in (47) at true value  $\theta_0$ , we can re-write it as

$$\begin{aligned} -\ln \lambda(\mathbf{x}) &= \sum_{i=1}^n \ln f(x_i; \theta_0) + \sum_{i=1}^n \sum_{r=1}^s \frac{\partial \ln f(x_i; \theta_0)}{\partial \theta_{0,r}} (\hat{\theta}_{n,r} - \theta_{0,r}) + n o_p(\|\hat{\theta}_n - \theta_0\|) - \sum_{i=1}^n \ln f(x_i; \theta_0^*) \\ &= \sum_{i=1}^n \ln \frac{f(x_i; \theta_0)}{f(x_i; \theta_0^*)} + \sum_{i=1}^n \mathbf{J}(\theta_0; x_i)^T (\hat{\theta}_n - \theta_0) + n o_p(\|\hat{\theta}_n - \theta_0\|), \end{aligned} \quad (48)$$

where  $\mathbf{J}$  is Fisher Score. To use L.L.N. and C.L.T., we manipulate (48) into a more convenient form and split it into three parts, namely,  $nA$ ,  $\sqrt{n} \cdot B \cdot C$ , and  $\sqrt{n} o_p(\|C\|)$ ,

$$\begin{aligned} -\ln \lambda(\mathbf{x}) &= n \cdot \frac{1}{n} \sum_{i=1}^n \ln \frac{f(x_i; \theta_0)}{f(x_i; \theta_0^*)} + \sqrt{n} \left[ \frac{1}{n} \sum_{i=1}^n \mathbf{J}(x_i; \theta_0)^T - 0 \right] \cdot \sqrt{n} (\hat{\theta}_n - \theta_0) + \sqrt{n} o_p(\sqrt{n} \|\hat{\theta}_n - \theta_0\|) \\ &= nA + \sqrt{n} \cdot B \cdot C + \sqrt{n} o_p(\|C\|). \end{aligned} \quad (49)$$

By L.L.N, as  $n$  tends to infinity,  $A$  tends to

$$E_{\theta_0} \ln \frac{f(x_i; \theta_0)}{f(x_i; \theta_0^*)} = E_{\theta_0} \left( -\ln \frac{f(x_i; \theta_0^*)}{f(x_i; \theta_0)} \right)$$

with probability one. Observe that  $-\ln(\bullet)$  is a convex function, we can apply Jensen's Inequality to the limit of  $A$ ,

$$A \rightarrow E_{\theta_0} \left( -\ln \frac{f(x_i; \theta_0^*)}{f(x_i; \theta_0)} \right) > -\ln E_{\theta_0} \frac{f(x_i; \theta_0^*)}{f(x_i; \theta_0)} = -\ln \int \frac{f(x_i; \theta_0^*)}{f(x_i; \theta_0)} f(x_i; \theta_0) dx = -\ln 1 = 0. \quad (50)$$

Hence we have proved that  $A \rightarrow \text{constant} > 0$  with probability one. Consequently,  $A \rightarrow n \cdot \text{constant} = \infty$  with probability one.

By LLN, the second term in (49) is bounded. This suffice to show that  $-\ln \lambda(\mathbf{x})$  will be greater than any given constant as  $n$  goes to infinity with probability one. In other words, we reject null hypothesis with probability one.  $\square$

In the end of this section, we present an example of small sample size. In this case, we can deduce the exact distribution of  $\lambda(\mathbf{x})$  without requiring  $n \rightarrow \infty$  or applying Wilk's Theorem.

**Example 4.2.** *Suppose the p.d.f. in (1) is given by  $N(\mu_0, \sigma_0^2)$ . Find the test statistic for  $H_0: \mu_0 = \mu_0^*$  vs  $H_1: \mu_0 \neq \mu_0^*$ .*

**Sol 4.2.** By Definition 4.1, we can write

$$\lambda(\mathbf{x}) = \frac{\sup_{\sigma^2} (2\pi\sigma^2)^{-n/2} \exp\{-\sum_{i=1}^n (x_i - \mu_0^*)^2 / (2\sigma^2)\}}{\sup_{\mu, \sigma^2} (2\pi\sigma^2)^{-n/2} \exp\{-\sum_{i=1}^n (x_i - \mu)^2 / (2\sigma^2)\}}. \quad (51)$$

Note that the MLE for the numerator is

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0^*)^2; \quad (52)$$

and the MLE for the denominator are

$$\hat{\mu} = \bar{x}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (53)$$

Therefore we can calculate  $\lambda(\mathbf{x})$  by plugging (52) and (53) back to (51)

$$\ln \frac{1}{\lambda(\mathbf{x})} \propto \frac{\hat{\sigma}^2}{\hat{\sigma}^2} = \frac{\sum_{i=1}^n (x_i - \mu_0^*)^2 / n}{\sum_{i=1}^n (x_i - \bar{x})^2 / (n-1)} \quad (54)$$

To simplify our test rule further we use the following equation, which can be established by expanding  $\hat{\sigma}^2$ .

$$\hat{\sigma}^2 = \hat{\sigma}^2 + (\bar{x} - \mu_0^*)^2 \quad (55)$$

Therefore,

$$\ln \frac{1}{\lambda(\mathbf{x})} \propto 1 + (\bar{x} - \mu_0^*)^2 / \hat{\sigma}^2 \quad (56)$$

Because  $s^2 = (n-1)^{-1} \sum_{i=1}^n (x_i - \bar{x})^2 = n\hat{\sigma}^2$ ,  $\hat{\sigma}^2 / \hat{\sigma}^2$  is a monotone increasing function of  $|T_n|$  where

$$T_n = \frac{\sqrt{n}(x - \mu_0^*)}{s}. \quad (57)$$

Therefore the likelihood ratio tests reject for small values of  $\lambda(\mathbf{x})$ , or equivalently, large values of  $|T_n|$ . Because  $T_n$  has a  $T$  distribution under  $H_0$ , the size  $\alpha$  critical value is  $t_{n-1, 1-\alpha/2}$ . We should reject null hypothesis if  $|T_n| \geq t_{n-1, 1-\alpha/2}$ .

## 5 Likelihood Ratio Confidence Interval

In the previous section, we derived the fact that  $-2 \ln \lambda(\mathbf{x})$  has an asymptotic chi squared distribution. For fixed  $\theta_0^*$  in  $H_0 : \theta_0 = \theta_0^*$ , the acceptance region is given by

$$\{\lambda(\mathbf{x}) : -2 \ln \lambda(\mathbf{x}) \leq \chi_{1, \alpha}^2\}, \quad (58)$$

Then by inverting the LRT, we can conclude that for (1)-(2), the set

$$\left\{ \theta : -2 \ln \left( \frac{L(\theta; \mathbf{x})}{L(\hat{\theta}; \mathbf{x})} \right) \leq \chi_{1, \alpha}^2 \right\} \quad (59)$$

is an approximate  $1 - \alpha$  confidence interval.

**Example 5.1.** The p.d.f. in (1) is given by Bernoulli( $p$ ) and  $Y = \sum_{i=1}^n X_i$ . We have the approximate  $1 - \alpha$  confidence set

$$\left\{ p : -2 \ln \left( \frac{p^y (1-p)^{n-y}}{\hat{p}^y (1-\hat{p})^{n-y}} \right) \leq \chi_{1, \alpha}^2 \right\}. \quad (60)$$

For some special distributions, we can find the exact distribution of  $\lambda(\mathbf{x})$ . In this case, we can get a more accurate confidence interval by inverting the LRT of  $H_0 : \theta_0 = \theta_0^*$  vs  $H_1 : \theta_0 \neq \theta_0^*$ . The confidence interval is of form

$$\text{accept } H_0 \text{ if } \frac{L(\theta_0^*; \mathbf{x})}{L(\hat{\theta}; \mathbf{x})} \leq k(\theta_0), \quad (61)$$

with the resulting confidence region

$$\{\theta : L(\theta; \mathbf{x}) \geq k'(\mathbf{x}, \theta)\}, \quad (62)$$

for some function  $k'$  that gives  $1 - \alpha$  confidence.

**Example 5.2.** Suppose that the p.d.f. of (1) is exponential( $\lambda$ ). Find the confidence interval for  $\lambda$  by inverting a level  $\alpha$  test of  $H_0 : \lambda = \lambda_0^*$  vs  $H_1 : \lambda \neq \lambda_0^*$ .

**Sol 5.1.** For random sample (1)-(2), the LRT statistic is given by

$$\frac{\exp(-\sum x_i/\lambda_0^*)/\lambda_0^{*n}}{\sup_{\lambda>0} \exp(-\sum x_i/\lambda)/\lambda^n} = \frac{\exp(-\sum x_i/\lambda_0^*)/\lambda_0^{*n}}{e^{-n}/(\sum x_i/n)} = \left(\frac{\sum x_i}{n\lambda_0^*}\right)^n e^n e^{-\sum x_i/\lambda_0^*}. \quad (63)$$

For fixed  $\lambda_0^*$ , the acceptance region is given by

$$A(\lambda_0^*) = \left\{ \mathbf{x} : \left(\frac{\sum x_i}{\lambda_0^*}\right)^n e^{-\sum x_i/\lambda_0^*} \geq k^* \right\}, \quad (64)$$

where  $k^*$  is a constant chosen to satisfy  $P_{\lambda_0^*}(\mathbf{X} \in A(\lambda_0^*)) = 1 - \alpha$  and the constant  $e^n/n$  has been absorbed into  $k^*$ . Inverting this acceptance region gives the  $1 - \alpha$  confidence set

$$C(\mathbf{x}) = \left\{ \lambda : \left(\frac{\sum x_i}{\lambda}\right)^n e^{-\sum x_i/\lambda} \geq k^* \right\}, \quad (65)$$

which is an interval in the parameter space  $\Theta$ .

The expression defining  $C(\mathbf{x})$  depends on  $\mathbf{x}$  only through  $\sum x_i$ . So the confidence interval can be expressed in the form

$$C\left(\sum x_i\right) = \left\{ \lambda : L\left(\sum x_i\right) \leq \lambda \leq U\left(\sum x_i\right) \right\} \quad (66)$$

where  $L$  and  $U$  are functions determined by the constraints that the set (64) has probability  $1 - \alpha$  and

$$\left(\frac{\sum x_i}{L(\sum x_i)}\right)^n e^{-\sum x_i/L(\sum x_i)} = \left(\frac{\sum x_i}{U(\sum x_i)}\right)^n e^{-\sum x_i/U(\sum x_i)}. \quad (67)$$

If we set

$$\frac{\sum x_i}{L(\sum x_i)} = a \text{ and } \frac{\sum x_i}{U(\sum x_i)} = b, \text{ where } a > b \text{ are constants.} \quad (68)$$

Then (67) becomes  $a^n e^{-1} = b^n e^{-b}$ , which yields easily to numerical solution. To work out some details, let  $n = 2$  and note that  $\sum X_i \sim \Gamma(2, \lambda)$  and  $\sum X_i/\lambda \sim \Gamma(2, 1)$ . Hence from (68), the confidence interval becomes

$$\left\{ \lambda : \frac{1}{a} \sum x_i \leq \lambda \leq \frac{1}{b} \sum x_i \right\},$$

where  $a$  and  $b$  satisfy

$$P_\lambda \left( \frac{1}{a} \sum X_i \leq \lambda \leq \frac{1}{b} \sum X_i \right) = P \left( b \leq \frac{\sum X_i}{\lambda} \leq a \right) = 1 - \alpha$$

Then

$$P \left( b \leq \frac{\sum X_i}{\lambda} \leq a \right) = \int_b^a t e^{-t} dt = e^{-b}(b+1) - e^{-a}(a+1). \quad (69)$$

To get, for example, a 90% confidence interval, we must simultaneously satisfy the probability condition and the constraint. To Three decimal places, we got  $a = 5.480$ ,  $b = 0.441$ , with a confidence coefficient of 0.90006. Thus,

$$P_{\lambda} \left( \frac{1}{5.480} \sum X_i \leq \lambda \leq \frac{1}{0.441} \sum X_i \right) = 0.90006.$$

## References

- [1] George Casella and Roger L. Berger, *Statistical Inference*. Duxbury Press, 2nd Edition, 2002.