

## 1 *The Likelihood Function*

In the previous chapter, we introduced statistical models  $\{P_\theta : \theta \in \Omega\}$  which describe the probability models that could generate the observed data. In this chapter we will develop inferences that depend only on the model  $\{P_\theta : \theta \in \Omega\}$  and the data  $s$ .

The simplest setting has the statistical model resulting in discrete distribution.

Suppose that we observe the data  $s$  and that the pmf is  $p_\theta$ . The **likelihood function**  $L(\cdot|s)$  is defined on the parameter space  $\Omega$  by

$$L(\theta|s) = p_\theta(s), \quad \theta \in \Omega.$$

For the observed data  $s$ , the likelihood is the probability of observing  $s$  when the true value of the parameter is  $\theta$ . This induces an ordering on  $\Omega$  in that we believe that  $\theta_1$  is more plausible as the true value of the parameter than  $\theta_2$  if

$$p_{\theta_1}(s) > p_{\theta_2}(s) \quad \text{or} \quad L(\theta_1|s) > L(\theta_2|s).$$

STOP

End

10/22/21

**Remarks:**

- We note that  $L(\theta|s)$  is the probability of the value  $s$  given that the true value of the parameter is  $\theta$ .
- The likelihood  $L(\theta|s)$  **is not** the probability of  $\theta$  given that we have observed  $s$ .
- In many cases the value of  $L(\theta|s)$  is small for all  $\theta$ . Thus, we are interested in the relative value of the likelihood.
- The above implies that we should consider **likelihood ratios**

$$\frac{L(\theta_1|s)}{L(\theta_2|s)} \text{ or } \log(L(\theta_1|s)) - \log(L(\theta_2|s))$$

in determining inferences for  $\theta$  based on the likelihood function.

# Statistics 630

---

We now define the **likelihood function** for a random sample  $X_1, \dots, X_n$  from a distribution with pmf or pdf  $f_\theta(\cdot)$ . The joint pmf or pdf  $X_1, \dots, X_n$  is

$$f_\theta(x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i).$$

Now, suppose that  $x_1, \dots, x_n$  are the observed values of  $X_1, \dots, X_n$ . Then the **likelihood function** is

$$L(\theta|x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i).$$

The notation  $L(\theta|x_1, \dots, x_n)$  indicates that we regard  $L(\cdot|x_1, \dots, x_n)$  as a function of  $\theta$ . After the data are observed,  $x_1, \dots, x_n$  are viewed as constants.

**Remark:** In the case where we have several parameters  $\theta_1, \dots, \theta_k$ , we define the **likelihood function** for a random sample  $X_1, \dots, X_n$  from a distribution with pmf or pdf as  $f_{\theta_1, \dots, \theta_k}(\cdot)$

---

$$L(\theta_1, \dots, \theta_k|x_1, \dots, x_n) = \prod_{i=1}^n f_{\theta_1, \dots, \theta_k}(x_i).$$

# Statistics 630

---

For a random sample from a discrete distribution with one parameter  $\theta$ ,

$$L(\theta|x_1, \dots, x_n) = P_\theta[X_1 = x_1, \dots, X_n = x_n].$$

We compare the likelihood function at two parameter values,  $\theta_1$  and  $\theta_2$ . If

$$L(\theta_1|x_1, \dots, x_n) > L(\theta_2|x_1, \dots, x_n),$$

then the observed data  $X_1 = x_1, \dots, X_n = x_n$  are more likely to have occurred if  $\theta = \theta_1$  than if  $\theta = \theta_2$ .

For a continuous rv  $X$  with pdf  $f_\theta(x)$ , we can write that

$$P_\theta[x - \epsilon < X < x + \epsilon] = \int_{x-\epsilon}^{x+\epsilon} f_\theta(x) dx \approx 2\epsilon f_\theta(x) = 2\epsilon L(\theta|x).$$

Thus, if

$$\frac{2\epsilon L(\theta_1|x)}{2\epsilon L(\theta_2|x)} = \frac{L(\theta_1|x)}{L(\theta_2|x)} > 1,$$

we feel that  $X$  is more likely to be near  $x$  when  $\theta = \theta_1$ .

# Statistics 630

---

Example 44: In a sample of 50 adult Americans, only 14 correctly described the Bill of Rights as the first ten amendments to the U. S. Constitution. Estimate the proportion of Americans that can give a correct description of the Bill of Rights.

Let  $Y_i = 1$  if the person is correct and  $Y_i = 0$  if the person is wrong. Then the pmf of a single observation is

$$p_\theta(y_i) = \theta^{y_i}(1 - \theta)^{1-y_i}$$

and the likelihood function is

$$L(\theta|y_1, \dots, y_n) = \prod_{i=1}^n \theta^{y_i}(1 - \theta)^{1-y_i} = \theta^y(1 - \theta)^{n-y}.$$

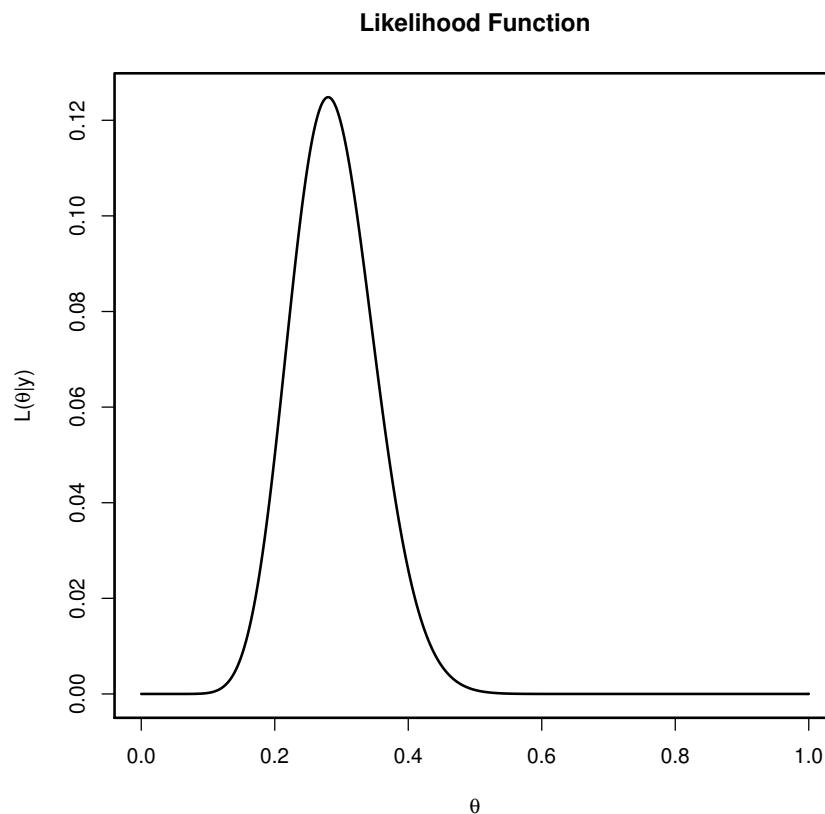
where  $y = \sum_{i=1}^n y_i$ .

$$\text{where } y = \sum_{i=1}^n y_i \sim \text{Bin}(n, \theta)$$

For our data,

$$L(\theta|y) = \theta^y(1-\theta)^{n-y} = \theta^{14}(1-\theta)^{50-14}.$$

The following is a plot of the likelihood function:



## 1.1 Sufficient Statistics

*(What summary data is needed to express the relative likelihood values.)*

In the previous example, the value of the likelihood function depended on the observed sample  $\{y_1, \dots, y_n\}$  only through  $y = \sum_{i=1}^n y_i$ . Such a simplification of the likelihood function will occur for many of the models that we use in this course.

A statistic  $T(s)$  is said to be a **sufficient statistic** for the model  $\{P_\theta : \theta \in \Omega\}$  (or simply for  $\theta$ ) if, whenever  $T(s_1) = T(s_2)$ , then

$$L(\theta|s_1) = c(s_1, s_2)L(\theta|s_2).$$

We will typically use the Factorization Theorem to check for a sufficient statistic:

~~Factorization Theorem~~ Suppose that the density (or pmf) for the model  $P_\theta$  is given by  $f_\theta(s)$ . Then  $T$  is a sufficient statistic for the model if the density factors

$$f_\theta(s) = h(s)g_\theta(T(s)),$$

where  $g_\theta$  and  $h$  are nonnegative and  $h$  does not depend on  $\theta$ .

NOTE:  $S$  is always sufficient if  $T$  factors themselves as a sufficient. The question

### Example 44:

The pmf of a single observation is  $p_\theta(y_i) = \theta^{y_i}(1-\theta)^{1-y_i}$ , and the **likelihood function** (or joint pmf) is

$$f_\theta(y_1, \dots, y_n) = \prod_{i=1}^n \theta^{y_i}(1-\theta)^{1-y_i} = \theta^{T(y_1, \dots, y_n)}(1-\theta)^{n-T(y_1, \dots, y_n)} \times 1.$$

where  $T(y_1, \dots, y_n) = \sum_{i=1}^n y_i$ . Set  $g_\theta(T) = \theta^T(1-\theta)^{n-T}$  and  $h(y_1, \dots, y_n) = 1$ . By the Factorization Theorem,  $T$  is sufficient for  $\theta$ .

## Example: Two-Parameter Weibull Distribution

Let  $X_1, \dots, X_n$  be a random sample from the Weibull distribution with pdf

$$f_{\alpha,\beta}(x) = \alpha\beta x^{\alpha-1} e^{-\beta x^\alpha}, \quad x > 0, \alpha > 0, \beta > 0.$$

1. Suppose  $\alpha = 2$ . Then the likelihood is

$$f_\beta(x) = (2\beta)^n (\prod_i x_i) e^{-\beta \sum_i x_i^2}$$

$X_i^2 \sim \text{Exp}$   
 $\Rightarrow \sum_i X_i^2 \sim \text{Gamma}$

We use the Factorization Theorem to show that  $\sum_{i=1}^n X_i^2$  is sufficient for  $\beta$ .

2. Suppose that  $\beta = 1$ . Then the likelihood is

$$f_\alpha(x) = \alpha^n (\prod_i x_i^{\alpha-1}) e^{-\sum_i x_i^\alpha}$$

We cannot reduce the data using the Factorization Theorem to anything simpler than  $T(X_1, \dots, X_n) = (X_1, \dots, X_n)$ . Thus, the entire set of  $n$  observations forms the sufficient statistic.

## 1.2 Maximum likelihood estimates

We now use the likelihood  $L(\theta|s)$  for the data  $s$  to define a point estimate of  $\theta$ .

For the observed data  $s$ , the *maximum likelihood estimates* are values  $\hat{\theta}$  such that

$$L(\hat{\theta}(s)|s) \geq L(\theta|s).$$

for all  $\theta$  in the parameter space  $\Omega$ .

The intuition of maximum likelihood is that  $\hat{\theta}$  is the parameter value such that the observed data  $s$  *is the most probable*. In a certain sense, then, the maximum likelihood estimates are those parameter values that *are the most consistent with the observed data*.

Often we use the **log-likelihood function** for inference:

$$\ell(\theta|s) = \log(L(\theta|s))$$

Since  $\log(x)$  is a one-to-one and increasing function of  $x$ ,  $\ell(\theta_1|s) > \ell(\theta_2|s)$  iff  $L(\theta_1|s) > L(\theta_2|s)$ . Thus, we could have defined the mle as the value of the parameter(s) that maximizes the log-likelihood.

## Log-likelihood for a Random Sample:

Let  $X_1, \dots, X_n$  be a random sample from a distribution with pmf or pdf  $f_\theta(\cdot)$ .

Then the **likelihood function** is

$$L(\theta|x_1, \dots, x_n) = \prod_{i=1}^n f_\theta(x_i),$$

and the log-likelihood function is

$$\ell(\theta|x_1, \dots, x_n) = \sum_{i=1}^n \log(f_\theta(x_i)).$$

**Multiparameter Case:** For the observed data  $s$ , the *maximum likelihood estimates* are values  $\hat{\theta}_1(s), \dots, \hat{\theta}_k(s)$  such that

$$L(\hat{\theta}_1(s), \dots, \hat{\theta}_k(s)|s) \geq L(\theta_1, \dots, \theta_k|s).$$

for all  $(\theta_1, \dots, \theta_k)$  in the parameter space  $\Omega$ . As in the single parameter case, we can use *log-likelihood function* for inference:

$$\ell(\theta_1, \dots, \theta_k|s) = \log(L(\theta_1, \dots, \theta_k|s)).$$

Let  $X_1, \dots, X_n$  be a random sample from a distribution with pmf or pdf  $f_{\theta_1, \dots, \theta_k}(\cdot)$ . Then the *likelihood function* is

$$L(\theta_1, \dots, \theta_k|x_1, \dots, x_n) = \prod_{i=1}^n f_{\theta_1, \dots, \theta_k}(x_i),$$

and the log-likelihood function is

$$\ell(\theta_1, \dots, \theta_k|x_1, \dots, x_n) = \sum_{i=1}^n \log(f_{\theta_1, \dots, \theta_k}(x_i)).$$

## 1.3 Computation of the MLE

We now consider the case where  $\theta$  is a single parameter ( $k = 1$ ) and  $\Omega$  is contained in the real line. The mle is the value  $\hat{\theta}(s)$  that maximizes the likelihood  $L(\theta|s)$  or equivalently, the log-likelihood  $\ell(\theta|s) = \log(L(\theta|s))$ .

The **score function** is defined to be the first partial derivative of the log-likelihood function with respect to  $\theta$  :

$$S(\theta|s) = \frac{\partial \ell(\theta|s)}{\partial \theta}.$$

*Assume  $\ell$  is differentiable  
in  $\theta$  at its max.*

To obtain the MLE, we solve the **score equation**:

$$S(\theta|s) = \frac{\partial \ell(\theta|s)}{\partial \theta} = 0.$$

We also need to check that our solution is an global maximum. To check for a local maximum, we can check that

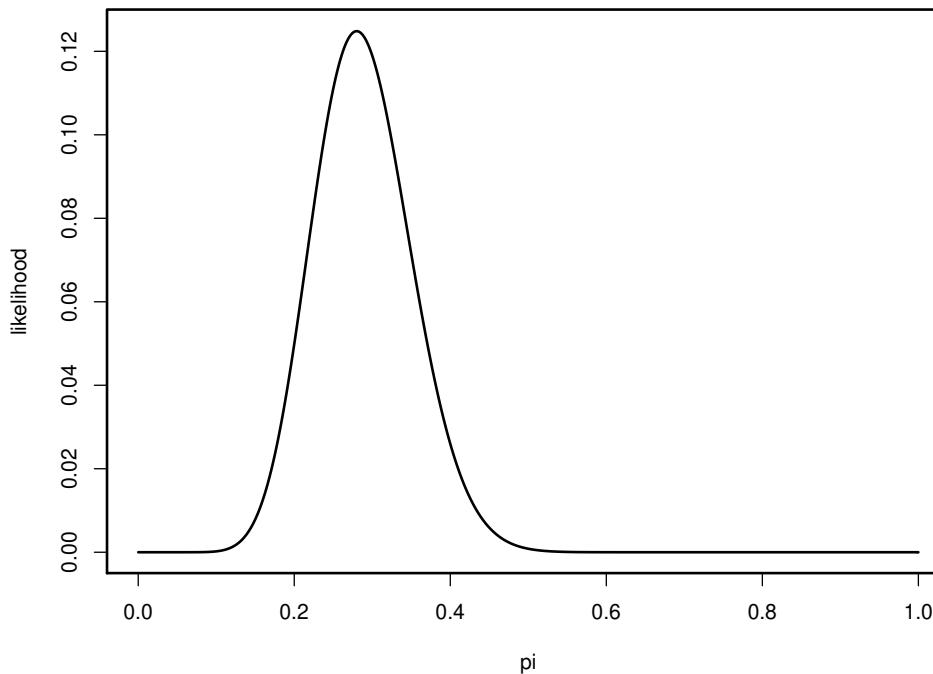
$$\left. \frac{\partial S(\theta|s)}{\partial \theta} \right|_{\theta=\hat{\theta}(s)} = \left. \frac{\partial^2 \ell(\theta|s)}{\partial \theta^2} \right|_{\theta=\hat{\theta}(s)} < 0.$$

*ST R: Monday 10/25/21*

Example 44: The likelihood function is

$$L(\theta|y_1, \dots, y_n) = \prod_{i=1}^n \theta^{y_i} (1-\theta)^{1-y_i} = \theta^y (1-\theta)^{n-y} = \theta^{14} (1-\theta)^{50-14}.$$

The following is a plot of the likelihood function:



It is an easier calculus problem to maximize the **log-likelihood**:

$$\begin{aligned}\ell(\theta|y) &= \log(L(\theta|y)) = \log [\theta^y(1-\theta)^{n-y}] \\ &= y\log(\theta) + (n-y)\log(1-\theta)\end{aligned}$$

To maximize, take the derivative and set = 0 to obtain the score equation:

$$\frac{\partial \ell(\theta|y)}{\partial \theta} = \frac{y}{\theta} - \frac{n-y}{1-\theta} = 0$$

We obtain the mle:

$$\hat{\theta}(y) = \frac{y}{n} = \frac{\sum_i y_i}{n} = \frac{14}{50} = 0.28$$

---

---

Example 45 Let  $X_1, \dots, X_n$  be a random sample from the **normal distribution**  $f_{\mu, \sigma}(\cdot)$ . The **likelihood function** is

$$L(\mu, \sigma | x_1, \dots, x_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \left(\frac{x_i - \mu}{\sigma}\right)^2\right].$$

We need to maximize the **log-likelihood**,  $\ell = \log(L)$ .

$$\begin{aligned}\ell(\mu, \sigma) &= \log(L(\mu, \sigma | x_1, \dots, x_n)) \\ &= -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2.\end{aligned}$$

To find the maximizer of  $\ell(\mu, \sigma)$ , take partial derivatives with respect to  $\mu$  and  $\sigma$ . Then set them equal to 0.

$$\frac{\partial \ell(\mu, \sigma)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \mu).$$

Setting this equal to 0 yields  $\mu = \bar{x}$ . Now

$$\frac{\partial \ell(\mu, \sigma)}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\implies \sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2.$$

Substituting in  $\mu = \bar{x}$ , we obtain

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

*We expect sample variance to not reflect population variance if n is large.*

*Note that  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$  is called unbiased sample variance.*

By computing second partials, one may verify that, in fact,

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

maximize the likelihood. The quantity  $\hat{\sigma}^2$  is a version of the *sample variance*.

Since we know that  $\mu$  and  $\sigma^2$  are the mean and variance of the normal distribution, it makes sense that we would use the sample mean and variance to estimate them.

We call the quantities  $\bar{x}$  and  $\hat{\sigma}^2$  **maximum likelihood estimates**. They are fixed values computed from the observed values  $x_1, \dots, x_n$ .

If we replace  $x_1, \dots, x_n$  in the above expression by the rvs  $X_1, \dots, X_n$ , we obtain **maximum likelihood estimators** which are random variables and have a probability distribution called a **sampling distribution**. We derived the sampling distributions of  $\bar{X}$  and  $\hat{\sigma}^2$  in Chapter 4.

Example 46 Let  $X_1, \dots, X_n$  be a random sample from the exponential( $\lambda$ ) distribution. Find the mle of  $\lambda$ .

$$E[X_i] = \frac{1}{\lambda}$$

The log-likelihood is

$$\ell(\lambda) = \log(L(\lambda|x_1, \dots, x_n)) = \log\left(\lambda^n e^{-\lambda \sum_{i=1}^n x_i}\right) = n \log(\lambda) - \lambda \sum_{i=1}^n x_i.$$

We differentiate the log-likelihood and set = 0 to get the score equation:

$$\frac{\partial \ell(\lambda)}{\partial \lambda} = \frac{n}{\lambda} - \sum_{i=1}^n x_i = 0.$$

Now solve to get the mle of  $\lambda$ :

$$\hat{\lambda} = \frac{n}{\sum_{i=1}^n x_i} = \frac{1}{\bar{x}}$$

MLE  
or  $= \frac{1}{\bar{x}^2}$

$\Rightarrow \bar{x}$  is the mle of mean  $(\frac{1}{\bar{x}})$

Example 47 Suppose  $X_1, \dots, X_n$  is a random sample from the [gamma distribution](#)

$$f_{\alpha, \lambda}(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} I_{(0, \infty)}(x).$$

The [log-likelihood](#) is

$$\begin{aligned}\ell(\alpha, \lambda) &= \log(L(\alpha, \lambda | x_1, \dots, x_n)) = \log \left( \prod_{i=1}^n \frac{\lambda^\alpha}{\Gamma(\alpha)} x_i^{\alpha-1} e^{-\lambda x_i} \right) \\ &= n\alpha \log(\lambda) + (\alpha - 1) \sum_{i=1}^n \log(x_i) - \lambda \sum_{i=1}^n x_i - n \log(\Gamma(\alpha))\end{aligned}$$

The [score equations \(also called likelihood equations\)](#) are

$$\frac{\partial \ell(\alpha, \lambda)}{\partial \alpha} = n \log \lambda + \sum_{i=1}^n \log(x_i) - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} = 0$$

and

$$\frac{\partial \ell(\alpha, \lambda)}{\partial \lambda} = \frac{n\alpha}{\lambda} - \sum_{i=1}^n x_i = 0.$$

We solve the second equation for  $\lambda$ :

$$\hat{\lambda} = \frac{n\hat{\alpha}}{\sum_{i=1}^n x_i} = \frac{\hat{\alpha}}{\bar{x}}.$$

We substitute this into the first equation to obtain a nonlinear equation for  $\hat{\alpha}$ :

$$n \log(\hat{\alpha}) - n \log(\bar{x}) + \sum_{i=1}^n \log(x_i) - n \frac{\Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})} = 0.$$

This equation cannot be solved in closed form. An iterative method of solving the equation must be used. This requires starting values such as the moment estimates which we will cover later.

Example 48 Let  $X_1, \dots, X_n$  be a random sample from the uniform distribution on the interval  $[0, \theta]$ ,  $\theta > 0$ . Find the mle of  $\theta$ .

+ Score function  
work on this example

The likelihood is

$$L(\theta|x_1, \dots, x_n) = \begin{cases} \frac{1}{\theta^n}, & 0 \leq x_i \leq \theta \\ 0 & \text{otherwise} \end{cases}$$

This is a decreasing function of  $\theta$  for  $\theta > 0$ . (Why?)

Also, the constraints on the  $x_i$ 's imply that  $x_i \leq \theta$ , all  $i$ . Equivalently,

$$x_{(n)} = \max\{x_1, \dots, x_n\} \leq \theta.$$

To maximize a decreasing function, we take the smallest allowable value of  $\theta$  to be the mle:

$$\hat{\theta} = \max\{X_1, \dots, X_n\} = X_{(n)}.$$

## 1.4 Invariance Property of MLEs

A basic property of the MLE is **invariance**. Suppose we let  $X_1, \dots, X_n$  form a random sample from a distribution with pmf or pdf  $f_\theta(x)$  and let  $\hat{\theta}$  be the mle of  $\theta$ .

Consider the alternative parameterization using  $\tau = \psi(\theta)$  where  $\psi$  is a one-to-one function of  $\theta$ . The *plug-in estimate* of  $\tau$  is given by  $\hat{\tau} = \psi(\hat{\theta})$ .

Alternatively, we could find the mle of  $\tau$  using the new parameterization.

The **invariance property of the mle** says that it makes no difference which parameterization we use for the finding the mle:

- If  $\hat{\theta}$  is the mle of  $\theta$  and  $\psi(\theta)$  is a one-to-one function, then  $\psi(\hat{\theta})$  is the mle of  $\psi(\theta)$ .
- Example 6.2.7 in the text shows that a plug-in estimate can behave badly when  $\psi$  is not one-to-one.

# Statistics 630

---

Example 44 again Consider a sequence of Bernoulli trials. Let  $Y_i = 1$  if the person is correct and  $Y_i = 0$  if the person is wrong. Then the pmf of a single observation is

$$f_\theta(y_i) = \theta^{y_i} (1 - \theta)^{1-y_i},$$

and the likelihood function is

$$L(\theta|y) = \prod_{i=1}^n \theta^{y_i} (1 - \theta)^{1-y_i} = \theta^y (1 - \theta)^{n-y},$$

where  $y = \sum_{i=1}^n y_i$ . The mle of  $\theta$  was found to be

$$\hat{\theta} = \frac{Y}{n}.$$

Consider the new parameter and its mle using **invariance**:

$$\xi = \frac{\theta}{1 - \theta}, \quad \hat{\xi} = \frac{\hat{\theta}}{1 - \hat{\theta}} = \frac{\frac{Y}{n}}{1 - \frac{Y}{n}} = \frac{Y}{n - Y}.$$

STOR 10(27 | z) WEDNESDAY

The likelihood in terms of  $\xi$  is

$$\text{lik}(\xi) = \prod_{i=1}^n \frac{\xi^{y_i}}{\xi + 1} = \frac{\xi^y}{(\xi + 1)^n}.$$

The likelihood equation is

$$\frac{\partial \ell(\xi)}{\partial \xi} = \frac{\partial}{\partial \xi} [y \log(\xi) - n \log(\xi + 1)] = \frac{y}{\xi} - \frac{n}{\xi + 1}$$

Set this equal to zero and solve to get the mle of  $\xi$ :

$$\hat{\xi} = \frac{Y}{n - Y}.$$

## 2 Method of Moments Estimators

We now introduce another approach to estimation of a parameter  $\theta$  when sampling from a distribution with pmf/pdf  $f_\theta(x)$  where  $\theta \in \Omega$ . Recall that the first population moment of a rv  $X$  is

$$\mu_X = E_\theta(X) \stackrel{\text{choose } \theta \text{ s.t.}}{=} m(\theta) \approx \bar{X}$$

**Note:**  $E_\theta(X)$  denotes the expectation of  $X$  using the distribution with pmf/pdf  $f_\theta(x)$ . This expectation typically is a function of  $\theta$ .

We now introduce **sample moments**. Given a random sample  $X_1, \dots, X_n$ , the first sample moment is  $m_1 = \bar{X}$ , the sample mean.

### Method of moments:

Express the population mean as a function of the unknown parameter  $\theta$ , solve for  $\theta$ , and substitute the sample mean for the population mean.

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the exponential( $\lambda$ ) distribution. Find the moment estimator of  $\lambda$ .

The exponential( $\lambda$ ) distribution has pdf

$$f_\lambda(x) = \lambda e^{-\lambda x}, \quad x > 0.$$

From the properties of the gamma distribution, the mean of an exponential rv is

$$E_\lambda(X) = \frac{1}{\lambda}.$$

Solve for  $\lambda$ :

$$\lambda = \frac{1}{E_\lambda(X)}.$$

Substitute  $\bar{X}$  for  $E_\lambda(X)$ :

$$\hat{\lambda} = \frac{1}{\bar{X}}.$$

We can extend the method of moments to estimation of multiple parameters  $\theta_1, \dots, \theta_k$  when sampling from a distribution with pmf/pdf  $f_{\theta_1, \dots, \theta_k}(x)$  where  $(\theta_1, \dots, \theta_k) \in \Omega$ . Recall that the  $j^{th}$  population moment of a rv  $X$  is

$$\mu_j = E_{\theta_1, \dots, \theta_k}(X^j).$$

**Note:**  $E_{\theta_1, \dots, \theta_k}(X^j)$  denotes the expectation of  $X^j$  using the distribution with pmf/pdf  $f_{\theta_1, \dots, \theta_k}(x)$ .

We now introduce **sample moments**. Given a random sample  $X_1, \dots, X_n$ , the  $j$ th sample moment is

$$m_j = \frac{1}{n} \sum_{i=1}^n X_i^j, \quad j = 1, 2, \dots.$$

## Multiparameter method of moments:

The population moments are functions of the unknown parameters. Express unknown parameters as functions of the population moments, and then substitute sample moments for population moments.

Example 47 again Suppose  $X_1, \dots, X_n$  is a random sample from the **gamma distribution**

$$f_{\alpha, \lambda}(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} I_{(0, \infty)}(x).$$

Recall that the first two moments of the gamma distribution are

$$E_{\alpha, \lambda}(X) = \frac{\alpha}{\lambda} \quad \text{and} \quad E_{\alpha, \lambda}(X^2) = \frac{\alpha(\alpha + 1)}{\lambda^2}.$$

Now, solve these two equations for  $\alpha$  and  $\lambda$ . This gives

$$\alpha = \frac{[E_{\alpha, \lambda}(X)]^2}{E_{\alpha, \lambda}(X^2) - [E_{\alpha, \lambda}(X)]^2} \Rightarrow \hat{\alpha} = \frac{\bar{x}^2}{\bar{x}^2 - s_x^2}$$

and

$$\lambda = \frac{E_{\alpha, \lambda}(X)}{E_{\alpha, \lambda}(X^2) - [E_{\alpha, \lambda}(X)]^2}. \Rightarrow \hat{\lambda} = \frac{\bar{x}}{\bar{x}^2 - s_x^2}$$

To obtain the method of moments estimators, we simply replace  $E(X)$  by  $\bar{X}$  and  $E(X^2)$  by  $m_2$ . This gives

$$\hat{\alpha} = \frac{\bar{X}^2}{m_2 - \bar{X}^2} \rightsquigarrow \hat{\alpha} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

and

$$\hat{\lambda} = \frac{\bar{X}}{m_2 - \bar{X}^2}.$$

We know that both  $\alpha$  and  $\lambda$  must be positive. Are  $\hat{\alpha}$  and  $\hat{\lambda}$  positive?

Since each data value  $X_i$  is positive,  $\bar{X}$  must be positive. This means the two estimates are positive if and only if  $m_2 - \bar{X}^2 > 0$ . *\* we problem ~ )*

It's easy to show that

$$m_2 - \bar{X}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2, \quad \text{possible to get nonsensical solutions.}$$

which obviously is positive.

## 3 Properties of Estimators

We will use the MLE  $\hat{\theta}$  as an estimate of the true value of  $\theta$ . If we are interested in estimating the characteristic of  $\psi(\theta)$ , a natural estimate is the plug-in estimate  $\psi(\hat{\theta})$ . More generally we can consider properties of a general estimator  $T$  of  $\psi(\theta)$ . Some important questions about estimators include the following:

- How close is the expected value of the estimator to the parameter it estimates?
- How close can we expect an estimate to be to the parameter it estimates?
- How is the behavior of an estimator related to sample size?
- How does the estimator compare to other estimators?

## 3.1 Bias

Let  $T$  be an estimator of a parametric quantity  $\psi(\theta)$ . A basic property of the estimator is its **bias**. The **bias** of an estimator  $T$  is defined as

$$\text{Bias}_\theta(T) = E_\theta(T) - \psi(\theta).$$

If  $\text{Bias}_\theta(T) = 0$  for all  $\theta \in \Omega$ , we say that  $T$  is an **unbiased estimator** of  $\psi(\theta)$ .

---

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the  $\text{exponential}(\lambda)$  distribution. Check the bias of the maximum likelihood estimator as an estimator of  $\lambda$ :

$$\hat{\lambda} = \frac{1}{\bar{X}} = \frac{n}{\sum_{i=1}^n X_i}.$$

We will use the fact the  $Y = \sum_{i=1}^n X_i \sim \text{Gamma}(n, \lambda)$ .

---

$$\begin{aligned}
 E_\lambda\left(\frac{1}{Y}\right) &= \int_0^\infty \frac{1}{y} \frac{\lambda^n}{\Gamma(n)} y^{n-1} e^{-\lambda y} dy \\
 &= \frac{\lambda^n}{\Gamma(n)} \frac{\Gamma(n-1)}{\lambda^{n-1}} \int_0^\infty \frac{\lambda^{n-1}}{\Gamma(n-1)} y^{n-2} e^{-\lambda y} dy = \frac{\lambda}{n-1}
 \end{aligned}$$

Thus,

$$E_\lambda(\hat{\lambda}) = E_\lambda\left(\frac{n}{Y}\right) = \frac{n}{n-1} \lambda \quad \text{and} \quad \text{Bias}_\lambda(\hat{\lambda}) = \frac{\lambda}{n-1}.$$


---

**Remark:** Sometimes one can adjust a biased estimator to create an unbiased estimator. We will illustrate that here, but will reserve judgment on which estimator is better. In the above example, let

$$\tilde{\lambda} = \frac{n-1}{n} \hat{\lambda} = \frac{n-1}{Y}.$$

Then  $E_\lambda(\tilde{\lambda}) = \lambda$ .

Example 45 again Let  $X_1, \dots, X_n$  be a random sample from the [normal distribution](#)  $f_{\mu, \sigma}(\cdot)$ . The maximum likelihood estimators were found to be

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

Check the bias of these estimators.

We know from Chapter 3 that  $E(\bar{X}) = \mu$  for any distribution. Thus,  $\bar{X}$  is an unbiased estimator of  $\mu$ .

From Chapter 4, we know that  $U = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \sim \chi^2(n-1)$ . Thus,

$$E(\hat{\sigma}^2) = E\left(\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right) = \frac{\sigma^2}{n} E(U) = \frac{(n-1)\sigma^2}{n}.$$

In fact, the previous result holds for any distribution with mean  $\mu$  and variance  $\sigma^2$ :

$$\begin{aligned}
 E[\hat{\sigma}^2] &= E\left[\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2\right] = \frac{1}{n} E\left[\sum_{i=1}^n X_i^2 - n(\bar{X})^2\right] \\
 &= \frac{1}{n} \left\{ \sum_{i=1}^n E(X_i^2) - nE[(\bar{X})^2] \right\} \\
 &= \frac{1}{n} \left\{ \sum_{i=1}^n (\sigma^2 + \mu^2) - n \left[ \frac{\sigma^2}{n} + \mu^2 \right] \right\} \\
 &= \frac{1}{n} \left\{ n\sigma^2 + n\mu^2 - (\sigma^2 + n\mu^2) \right\} = \frac{1}{n}(n\sigma^2 - \sigma^2) = \frac{n-1}{n}\sigma^2
 \end{aligned}$$

Thus,  $\hat{\sigma}^2$  has bias:

$$\text{Bias}(\hat{\sigma}^2) = \frac{n-1}{n}\sigma^2 - \sigma^2 = -\frac{1}{n}\sigma^2.$$

*Lec 11*  $S^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$  is unbiased for  $\sigma^2$

## 3.2 Standard Error

To measure variation in the estimator  $T$  of  $\psi(\theta)$ , one can evaluate either the variance or the standard deviation of its sampling distribution. The most commonly reported quantity is the standard error:

$$SE_\theta(T) = \sqrt{\text{Var}_\theta(T)}.$$

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the exponential( $\lambda$ ) distribution. Obtain the standard error of the maximum likelihood estimator,

$$\hat{\lambda} = \frac{1}{\bar{X}} = \frac{n}{\sum_{i=1}^n X_i}.$$

As in finding the bias, let  $Y = \sum_{i=1}^n X_i$ . Then (need  $\sim \mathcal{Z}$ )

$$E_\lambda \left( \frac{1}{Y^2} \right) = \int_0^\infty \frac{1}{y^2} \frac{\lambda^n}{\Gamma(n)} y^{n-1} e^{-\lambda y} dy = \frac{\lambda^2}{(n-1)(n-2)}.$$

Thus,

$$\begin{aligned}\text{Var}_\lambda(\hat{\lambda}) &= n^2 \text{Var}_\lambda\left(\frac{1}{Y}\right) = n^2 \lambda^2 \left[ \frac{1}{(n-1)(n-2)} - \left(\frac{1}{n-1}\right)^2 \right] \\ &= \lambda^2 \frac{n^2}{(n-2)(n-1)^2}\end{aligned}$$

and

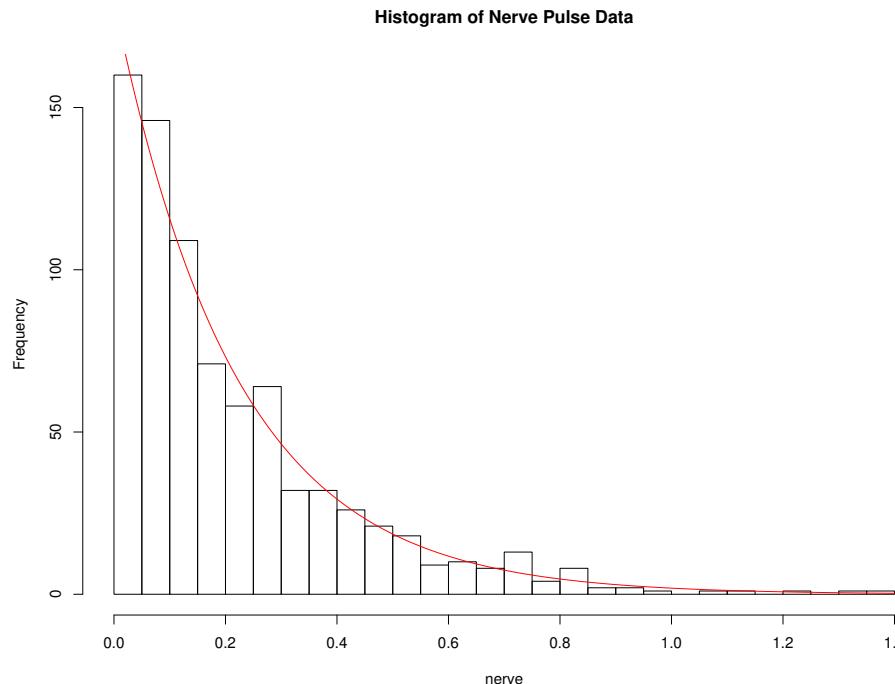
$$SE_\lambda(\hat{\lambda}) = \frac{\lambda n}{(n-1)\sqrt{n-2}}.$$

We often wish to estimate the standard error since it depends on the unknown parameter. Here, the estimated standard error of  $\hat{\lambda}$  is

$$s_{\hat{\lambda}} = \frac{\hat{\lambda} n}{(n-1)\sqrt{n-2}}.$$

$\hat{\lambda}$   $\sim$   $\frac{\lambda n}{\sqrt{n-2}}$  for large  $n$

Example 43 again Cox and Lewis (1966) reported 799 waiting times between successive pulses along a nerve fiber. The data appear in the following histogram:



The mean is  $\bar{x} = 0.2186$  and  $\hat{\lambda} = 1/0.2186 = 4.575$ . The exponential distribution with  $\lambda = 4.575$  is superimposed on the histogram. The estimated standard error of  $\hat{\lambda}$  is

$$s_{\hat{\lambda}} = \frac{\hat{\lambda}n}{(n-1)\sqrt{n-2}} = \frac{4.575(799)}{798\sqrt{797}} = 0.1623.$$

### 3.3 Mean Squared Error

A means of judging how well  $\hat{\theta}$  estimates  $\theta$  is to use the **mean squared error**. The mean squared error of an estimator  $T$  of a parametric quantity  $\psi(\theta)$  is

$$\text{MSE}_\theta(T) = E_\theta[(T - \psi(\theta))^2].$$

We may express the mean squared error of  $T$  as

$$\begin{aligned}\text{MSE}_\theta(T) &= E_\theta \left[ (T - E_\theta(T) + E_\theta(T) - \psi(\theta))^2 \right] \\ &= \text{Var}_\theta(T) + [\text{Bias}_\theta(T)]^2.\end{aligned}$$

The mean squared error is particularly useful for comparing two or more estimators. If we can show that

$$\text{MSE}_\theta(T_1) \leq \text{MSE}_\theta(T_2)$$

no matter what the value of  $\theta$  is, then this is a good reason to prefer  $T_1$  to  $T_2$ .

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the exponential( $\lambda$ ) distribution. Obtain the mean squared error of the moment estimator and of the bias adjusted estimator.

Earlier we found the bias and variance of  $\hat{\lambda}$ :  $\frac{1}{\bar{X}}$

$$\text{Bias}_{\lambda}(\hat{\lambda}) = \frac{\lambda}{n-1} \quad (\text{bias} \rightarrow 0 \text{ as } n \rightarrow \infty)$$

$$\text{Var}_{\lambda}(\hat{\lambda}) = \lambda^2 \frac{n^2}{(n-2)(n-1)^2}$$

Thus, the mean squared error is

$$\begin{aligned} \text{MSE}_{\lambda}(\hat{\lambda}) &= \text{Var}_{\lambda}(\hat{\lambda}) + \text{Bias}^2(\hat{\lambda}) \\ &= \lambda^2 \frac{n^2}{(n-2)(n-1)^2} + \left( \frac{\lambda}{n-1} \right)^2 \\ &= \lambda^2 \frac{n+2}{(n-1)(n-2)} \end{aligned}$$

The bias-adjusted estimator is

$$\tilde{\lambda} = \frac{n-1}{Y} = \frac{n-1}{n} \hat{\lambda}.$$

$$Y = \sum x_i$$

where  $\hat{\lambda} = \frac{1}{Y}$

We then have

$$\text{Bias}_\lambda(\tilde{\lambda}) = 0$$

$$\begin{aligned}\text{Var}_\lambda(\tilde{\lambda}) &= \left(\frac{n-1}{n}\right)^2 \text{Var}_\lambda(\hat{\lambda}) \\ &= \left(\frac{n-1}{n}\right)^2 \lambda^2 \frac{n^2}{(n-2)(n-1)^2} = \frac{\lambda^2}{n-2}\end{aligned}$$

Thus,

$$\text{MSE}_\lambda(\tilde{\lambda}) = \text{Var}_\lambda(\tilde{\lambda}) = \frac{\lambda^2}{n-2} < \lambda^2 \frac{n+2}{(n-1)(n-2)} = \text{MSE}_\lambda(\hat{\lambda}).$$

Example 45 again Let  $X_1, \dots, X_n$  be a random sample from the **normal distribution**  $f_{\mu, \sigma}(\cdot)$ . The mles were found to be

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

These were also found to be the moment estimators. We earlier obtained the expectations of these estimators and determined that the following properties:

- $\bar{X}$  is an unbiased estimator of  $\mu$ .
- $\text{Bias}_{\mu, \sigma}(\hat{\sigma}^2) = \frac{n-1}{n}\sigma^2 - \sigma^2 = -\frac{1}{n}\sigma^2$ .

Does this mean  $\hat{\sigma}^2$  is a bad estimator?

We could correct the bias in estimating  $\sigma^2$  by using the sample variance of Chapter 5:

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \approx \frac{\sqrt{n-1}}{\hat{\sigma}^2}$$

We next obtain the variances and mean squared errors of the two estimators of  $\sigma^2$ ,  $\hat{\sigma}^2$  and  $S^2$ .

$$\begin{aligned}\text{Var}_{\mu,\sigma}(\hat{\sigma}^2) &= \text{Var}_{\mu,\sigma} \left( \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 \right) \\ &= \frac{\sigma^4}{n^2} \text{Var} \left( \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \right) = \frac{2(n-1)\sigma^4}{n^2}\end{aligned}$$

The MSE of  $\hat{\sigma}^2$  becomes

$$MSE(\hat{\sigma}^2) = \text{Var}(\hat{\sigma}^2) + \text{Bias}^2(\hat{\sigma}^2) = \frac{2(n-1)\sigma^4}{n^2} + \left( \frac{-\sigma^2}{n} \right)^2 = \frac{(2n-1)\sigma^4}{n^2}$$

We next consider the sample variance,  $S^2$ .

$$\begin{aligned}\text{Var}_{\mu,\sigma}(S^2) &= \text{Var}_{\mu,\sigma} \left( \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right) \\ &= \frac{\sigma^4}{(n-1)^2} \text{Var} \left( \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sigma^2} \right) = \frac{2(n-1)\sigma^4}{(n-1)^2} = \frac{2\sigma^4}{n-1}.\end{aligned}$$

If we compare the MSEs of the two estimators, we see that the MSE of the biased estimator, the MLE  $\hat{\sigma}^2$ , is smaller than that of the unbiased estimator, the sample variance  $S^2$ :

$$\frac{MSE(S^2)}{MSE(\hat{\sigma}^2)} = \frac{\frac{2\sigma^4}{n-1}}{\frac{(2n-1)\sigma^4}{n^2}} = \frac{2n^2}{(2n-1)(n-1)} > 1.$$

---

---

Example 48 again Let  $X_1, \dots, X_n$  be a random sample from the uniform distribution on the interval  $[0, \theta]$ ,  $\theta > 0$ . Check the bias and mse of the mle of  $\theta$ . Compare the mle to the moment estimator.

First we need to find the pdf of  $\hat{\theta} = X_{(n)} = \max\{X_1, \dots, X_n\}$ . It is most easily derived by using the cdf. For  $0 \leq x \leq \theta$ ,

$$F_\theta(x) = P[X_{(n)} \leq x] = P[X_1 \leq x, \dots, X_n \leq x] = \prod_{i=1}^n P[X_i \leq x] = \left(\frac{x}{\theta}\right)^n.$$

$\underset{\theta}{\textcircled{X}_{(n)}} \sim \text{Beta}(n, 1)$

The pdf of  $X_{(n)}$  is then

$$f(x) = \frac{nx^{n-1}}{\theta^n}, \quad 0 \leq x \leq \theta.$$

The first two moments are

$$E(\hat{\theta}) = \frac{1}{\theta^n} \int_0^\theta x n x^{n-1} dx = \frac{n}{\theta^n} \int_0^\theta x^n dx = \frac{n}{\theta^n} \frac{x^{n+1}}{n+1} \Big|_0^\theta = \frac{n}{n+1} \theta,$$

$$E(\hat{\theta}^2) = \frac{1}{\theta^n} \int_0^\theta x^2 n x^{n-1} dx = \frac{n}{\theta^n} \int_0^\theta x^{n+1} dx = \frac{n}{\theta^n} \frac{x^{n+2}}{n+2} \Big|_0^\theta = \frac{n}{n+2} \theta^2,$$

$$\begin{aligned} \text{Var}(\hat{\theta}) &= E(\hat{\theta}^2) - (E(\hat{\theta}))^2 = \frac{n}{n+2} \theta^2 - \left( \frac{n}{n+1} \right)^2 \theta^2 \\ &= \frac{n}{(n+2)(n+1)^2} \theta^2. \end{aligned}$$

The resulting bias and MSE of  $\hat{\theta}$  are

$$\text{Bias}(\hat{\theta}) = E(\hat{\theta}) - \theta = -\frac{\theta}{n+1},$$

$$MSE(\hat{\theta}) = \frac{n}{(n+2)(n+1)^2} \theta^2 + \left( -\frac{\theta}{n+1} \right)^2 = \frac{2\theta^2}{(n+1)(n+2)}.$$

# Statistics 630

---

We could adjust the mle to remove the bias:

$$\check{\theta} = \frac{n+1}{n} \hat{\theta}.$$

The resulting variance and MSE of  $\check{\theta}$  are

$$\begin{aligned} MSE(\check{\theta}) &= \text{Var}(\check{\theta}) = \left( \frac{n+1}{n} \right)^2 \text{Var}(\hat{\theta}) \\ &= \left( \frac{n+1}{n} \right)^2 \frac{n}{(n+2)(n+1)^2} \theta^2 \\ &= \frac{\theta^2}{n(n+2)}, \end{aligned}$$

$$\frac{MSE(\hat{\theta})}{MSE(\check{\theta})} = \frac{\frac{2\theta^2}{(n+1)(n+2)}}{\frac{\theta^2}{n(n+2)}} = \frac{2n}{n+1} > 1.$$

Hence, the bias-adjusted estimator has smaller MSE for  $n > 1$ .

# Statistics 630

---

We now find the method of moments estimator and examine its properties.

$$E(X) = \frac{\theta}{2} \implies \text{the moment estimator is } \tilde{\theta} = 2\bar{X}.$$

The resulting bias and MSE are

$$\begin{aligned}\text{Bias}(\tilde{\theta}) &= E(\tilde{\theta}) - \theta = 2\frac{\theta}{2} - \theta = 0, \\ \text{MSE}(\tilde{\theta}) &= \text{Var}(\tilde{\theta}) = 4\text{Var}(\bar{X}) = 4\frac{\theta^2}{12n} = \frac{\theta^2}{3n}.\end{aligned}$$

We see that the MSE goes to zero as a function of  $\frac{1}{n}$  for the moment estimator and as a function of  $\frac{1}{n^2}$  for the other two estimators.

## 3.4 Consistency of Estimators

*Definition:* Let  $T_n$  be an estimator of  $\psi(\theta)$  based on a sample  $X_1, \dots, X_n$  from  $f_\theta(x)$ . Then the sequence of estimators  $\{T_n, n = 1, 2, \dots\}$  is said to be **consistent for  $\psi(\theta)$**  if

$$T_n \xrightarrow{P} \psi(\theta).$$

By the Weak Law of Large Numbers, the sample moments converge in probability to the population moments:

$$\frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{P} E(X^k). \quad \overbrace{\sum_{i=1}^n g(x_i)}^{\sim} \rightarrow \bar{E}(g(x))$$

If the functions that relate the estimates to the sample moments are continuous, then moment estimators will converge in probability to the parameters.

If  $g$  is continuous

$$g\left(\frac{1}{n} \sum_{i=1}^n X_i^k\right) \rightarrow g(E(X^k))$$

# Statistics 630

---

**Remark:** A convenient way to show that an estimator is consistent is to show that its MSE goes to zero as  $n$  tends to infinity. This result can be proved in the same manner as the Weak Law of Large Numbers:

$$\begin{aligned} 0 &\leq \lim_{n \rightarrow \infty} P_\theta[|T_n - \psi(\theta)| \geq \varepsilon] \\ &\leq \lim_{n \rightarrow \infty} \frac{E_\theta[(T_n - \psi(\theta))^2]}{\varepsilon^2} \\ &= \lim_{n \rightarrow \infty} \frac{\text{MSE}_\theta(T_n)}{\varepsilon^2} = 0. \end{aligned}$$

Since  $\text{MSE}_\theta(T_n) = \text{Var}_\theta(T_n) + [\text{Bias}_\theta(T_n)]^2$ , the sequence of estimators  $\{T_n\}$  is consistent if

$$\text{Var}_\theta(T_n) \rightarrow 0, \quad \text{and} \quad \text{Bias}_\theta(T_n) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

*Any time  $\text{MSE} \rightarrow 0$ , we have a consistent estimator.*

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the exponential ( $\lambda$ ) distribution. It is easy to show the consistency of the moment estimator,

$$\hat{\lambda}_n = \frac{1}{\bar{X}} = \frac{n}{\sum_{i=1}^n X_i}.$$

Since

$$\lim_{n \rightarrow \infty} \text{MSE}(\hat{\lambda}_n) = \lim_{n \rightarrow \infty} \lambda^2 \frac{n+2}{(n-1)(n-2)} = 0, \quad \hat{\lambda}_n \xrightarrow{P} \lambda.$$

Example 43 again Let  $X_1, \dots, X_n$  be a random sample from the uniform distribution on the interval  $(0, \theta)$ ,  $\theta > 0$ . Check the consistency the mle of  $\theta$ .

Since

$$\lim_{n \rightarrow \infty} \text{MSE}(\hat{\theta}) = \lim_{n \rightarrow \infty} \frac{2\theta^2}{(n+1)(n+2)} = 0, \quad \hat{\theta}_n \xrightarrow{P} \theta.$$

STOP 11/11/21 *Malay*:

NOTE (STATED BY) W. Waller (s.2)

### 3.5 Asymptotic Distribution of Estimators

We need to obtain the sampling distribution of point estimators in order to derive the properties of the estimators and also confidence intervals and hypothesis tests based on the estimators. However, it is often not practical to obtain the exact sampling distribution of an estimator  $T_n$  calculated from a random sample  $X_1, \dots, X_n$ . In such cases, we can obtain the large sample distribution of the statistic and use it to approximate the sampling distribution for finite sample size  $n$ . A basic result that is useful for statistics that are sample means is the Central Limit Theorem:

Let  $X_1, \dots, X_n$  be a random sample from a distribution having variance  $\sigma^2$  (with  $0 < \sigma^2 < \infty$ ) and mean  $\mu$ . Let  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ . Then, for each real number  $z$ ,

$$\lim_{n \rightarrow \infty} P\left(\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \leq z\right) = \Phi(z).$$

Thus,  $\frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} \xrightarrow{D} Z$ , where  $Z$  is a standard normal rv.

## Another Method of Finding Asymptotic Distributions

In many situations, such as in using method of moments estimators, one is interested in obtaining an asymptotic distribution of an estimator that is a function of a statistic that we know is asymptotically normal. The so-called delta method (or propagation of errors) enables us to obtain the large sample distribution of the estimator of interest.

Suppose that we know that

$$\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, V(\theta)),$$

and we are interested in the asymptotic distribution of the estimator  $g(T_n)$  of  $g(\theta)$ .

We assume that  $g(t)$  is a continuous function that has a sufficient number of derivatives. Then we can apply Taylor's theorem to  $g$ :

$$g(t) = g(\theta) + (t - \theta)g'(\theta) + \text{higher order terms.}$$

We now apply the Taylor expansion to  $T_n$  and ignore the higher order terms:

$$g(T_n) = g(\theta) + (T_n - \theta)g'(\theta).$$

We rearrange the terms to get

$$\sqrt{n}(g(T_n) - g(\theta)) = g'(\theta)\sqrt{n}(T_n - \theta).$$

Since  $\sqrt{n}(T_n - \theta) \xrightarrow{D} N(0, V(\theta))$  by assumption, we obtain the result that

$$\sqrt{n}(g(T_n) - g(\theta)) \xrightarrow{D} N(0, [g'(\theta)]^2 V(\theta)).$$

So  $g(\bar{X}_n) \approx N(g(\theta), [g'(\theta)]^2 V(\theta)/n)$

**Remark:** We can apply this result to most method of moment estimators. The Central Limit Theorem implies that  $\bar{X}_n$  is asymptotically normal. Since the MOM estimator is a smooth function of  $\bar{X}_n$ , we can use this result to obtain the asymptotic distribution of the MOM estimator.

Example 46 again Let  $X_1, \dots, X_n$  be a random sample from the exponential( $\lambda$ ) distribution. The moment estimator of  $\lambda$  was found to be  $\hat{\lambda} = 1/\bar{X}_n$ . From the properties of the gamma distribution, the mean and variance of an exponential rv are

$$E_\lambda(X) = \frac{1}{\lambda} \quad \text{and} \quad \text{Var}_\lambda(X) = \frac{1}{\lambda^2}. \quad \square^2$$

By the Central Limit Theorem,  $\sqrt{n}(\bar{X}_n - 1/\lambda) \xrightarrow{D} N(0, 1/\lambda^2)$ .

Now  $g(t) = 1/t$  and  $g'(t) = -1/t^2$ . Thus, we apply the delta method to  $\hat{\lambda} = g(\bar{X}_n) = 1/\bar{X}_n$ :

$$\sqrt{n}(\hat{\lambda} - \lambda) = \sqrt{n}(g(\bar{X}_n) - g(1/\lambda)) \xrightarrow{D} N(0, (1/\lambda^2)(g'(1/\lambda))^2).$$

Since

$$\frac{1}{\lambda^2} \left( g' \left( \frac{1}{\lambda} \right) \right)^2 = \left( \frac{1}{\lambda^2} \right) \left( \frac{-1}{(1/\lambda)^2} \right)^2 = \lambda^2,$$

$N(x, \frac{1}{\lambda^2})$

we get the result

$$\sqrt{n}(\hat{\lambda} - \lambda) \xrightarrow{D} N(0, \lambda^2).$$

End C6

Wednesday 11/3/21

### 3.5.1 Another Approach to Estimating the Standard Error

Here we could easily find an exact expression for the standard error of our estimator. For other estimators, this may not be the case. We can approximate the standard error computationally using a technique known as the **bootstrap**. We will first discuss the **parametric bootstrap**.

*known distribution for the data*

Suppose that we take a random sample  $X_1, \dots, X_n$  from a distribution with parameter  $\theta$ . Suppose that  $\hat{\theta} = S(X_1, \dots, X_n)$  for some statistic  $S$  whose sampling distribution is difficult to derive.

- If we knew the true value of the parameter, say  $\theta_0$ , we could generate *new data*  $x_1, \dots, x_n$  from  $f_{\theta_0}(x)$  and compute  $s_1 = s(x_1, \dots, x_n)$ .  
*using the parameter value for the data.*
- If we repeat this a large number  $B$  of times, we obtain a random sample of values of  $S$ ,  $(s_1, \dots, s_B)$ .
- We use this random sample to estimate the sampling distribution of  $S$ .
- Since we do not know  $\theta_0$ , we will replace it by  $\hat{\theta}$  and carry out this procedure.

*Example:* We will use the parametric bootstrap to estimate the standard error of the maximum likelihood estimator of the rate parameter for the exponential distribution.

Here  $\hat{\lambda} = 4.575$  and  $n = 799$ .

- We generate  $B = 1000$  samples of size  $n = 799$  from the exponential distribution with  $\lambda = 4.575$ .
- Compute  $\hat{\lambda}$  for each sample.
- Find the standard deviation of the 1000 values of  $\hat{\lambda}$ .

# Statistics 630

---

Following is the R code for carrying out these calculations.

```
> mean(nerve)
[1] 0.2185732
> 1/mean(nerve)
[1] 4.575126
> temp=rep(0,1000)
> for (i in 1:1000)temp[i]= 1/mean(rexp(799,rate=1/mean(nerve)))
> sd(temp)
[1] 0.1629491
```

Here we were curiously coincident with our previous estimate of the standard error of  $\hat{\lambda}$ ,  $s_{\hat{\lambda}} = 0.1623$ .

Other bootstrap samples of size  $B = 1000$  resulted in estimated standard errors of

0.1640392, 0.1633224, 0.1664940, 0.1648791, and 0.1573375.

Another approach does not use an assumed form for the pdf  $f_\theta(x)$ . The nonparametric bootstrap involves resampling from the observed data. We will suppose that we want to approximate the sampling distribution of a statistic,  $S(X_1, \dots, X_n)$ .

- Take a sample  $x_1^*, \dots, x_n^*$  of size  $n$  with replacement from the observed data,  $x_1, \dots, x_n$ .
- Compute the observed value of the statistic,  $s_1^* = s(x_1^*, \dots, x_n^*)$ .
- If we repeat this a large number  $B$  of times, we obtain  $B$  bootstrap values of  $S$ ,  $(s_1^*, \dots, s_B^*)$ .
- We use these bootstrap values to estimate the sampling distribution of  $S$ .

# Statistics 630

---

Following is the R code to generate 1000 bootstrap estimates of  $\lambda$  and compute the estimated standard error of  $\hat{\lambda}$ :

```
> temp=rep(0,1000)
> for(i in 1:1000)temp[i]=1/mean(sample(nerve,replace=TRUE))
> sd(temp)
[1] 0.1594686
```

Other bootstrap estimates of the standard error of  $\hat{\lambda}$  tended to be larger than those obtained using the parametric bootstrap.

## 3.6 Large Sample Theory for Maximum Likelihood Estimators

The MLE has excellent large sample properties under certain regularity conditions. We suppose that  $X_1, \dots, X_n$  is a random sample from a population with pdf or pmf  $f_\theta(x)$ .

- The density  $f_\theta(x)$  is a smooth function of  $\theta$ .
- The support of the distribution,  $\{x : f_\theta(x) > 0\}$  does not depend on the parameter,  $\theta$ .  
*→ rule out uniform parameter specification endpoint.*  
*for example.*
- The parameter space  $\Omega$  satisfies certain conditions.
- $\text{Var} \left( \frac{\partial \log(f_\theta(X))}{\partial \theta} \right)$  is finite. ★

We denote the “true” value of  $\theta$  by  $\theta_0$ .

Important properties of the MLE include the following:

1.  $\hat{\theta}_n$  is **consistent**. That is

$$\hat{\theta}_n \xrightarrow{P} \theta_0 \quad \text{as } n \rightarrow \infty$$

2. The MLE is asymptotically normal:

$\rightarrow$  proof  $\rightarrow$  essentially the  $\ell(\hat{\theta}|x)$   
 $\rightarrow$  a sum and the CLT  
 applies to sums.

$$\frac{\hat{\theta}_n - \theta_0}{\sqrt{V_n(\theta_0)}} \xrightarrow{D} N(0, 1) \quad \text{as } n \rightarrow \infty$$

We usually interpret this to mean

$$\hat{\theta}_n \underset{\text{approx}}{\sim} N(\theta_0, \hat{V}_n)$$

where

$$\hat{V}_n = V_n(\hat{\theta}_n)$$

We need to define one other quantity to obtain  $V_n(\theta_0)$ :

Fisher's information in  $\theta$  based on the observation  $X$  is defined as

$$I(\theta) = E_\theta \left[ \left( \frac{\partial}{\partial \theta} \log f_\theta(X) \right)^2 \right]$$

*equivalent  
to C of the log  
presence of the*

For computing  $I(\theta)$ , we often use the result that

$$I(\theta) = -E_\theta \left[ \left( \frac{\partial^2}{\partial \theta^2} \log f_\theta(X) \right) \right]$$

*This is fisher  
information based  
on 1 obs. we want  
it for n obs.*

This quantity can be estimated in several ways:

$I(\hat{\theta}_n)$  – Plug in

~~$\hat{I}(\hat{\theta}_n)$~~   $\hat{I}(\hat{\theta}_n) = -\frac{1}{n} \frac{\partial^2 \ell(\hat{\theta}_n)}{\partial \theta^2}$  – Hessian or observed information

*Take expectation to get  $I_n(\hat{\theta}_n)$*

~~to first order to do it the way to the following model.~~

$$E\left[\left(\sum_{i=1}^n \log(f_\theta(x_i))\right)^2\right] = \text{Var}(\bar{S}(\theta|X))$$

**Statistics 630**

---

Score function

## Fisher's Information in a Random Sample

Suppose now that we have a random sample:  $X_1, \dots, X_n$  are independent with pmf or pdf  $f_\theta(x)$ . The likelihood is

$$L(\theta) = f_\theta(x_1) \times \cdots \times f_\theta(x_n).$$

Then

$$\log(L(\theta)) = \log[f_\theta(x_1) \times \cdots \times f_\theta(x_n)] = \sum_{i=1}^n \log(f_\theta(x_i))$$

and

$$\frac{\partial^2 \log(L(\theta))}{\partial \theta^2} = \sum_{i=1}^n \frac{\partial^2 \log(f_\theta(x_i))}{\partial \theta^2}$$

$\checkmark C_i(\theta) = \frac{\partial^2 \log(f_\theta(x_i))}{\partial \theta^2}$   
 $I_n(\theta) = n \sum_i C_i(\theta)$

Then the Fisher information  $I_n(\theta)$  in  $X_1, \dots, X_n$  is given by

$$I_n(\theta) = -E_\theta \left[ \frac{\partial^2 \log(L(\theta))}{\partial \theta^2} \right] = \sum_{i=1}^n -E_\theta \left[ \frac{\partial^2 \log(f_\theta(x_i))}{\partial \theta^2} \right] = n I(\theta)$$

One can show that

$$V_n(\theta_0) = \frac{1}{nI(\theta_0)}.$$
\*

Thus, we can say that

$$\underbrace{\sqrt{nI(\theta_0)}(\hat{\theta}_n - \theta_0)}_{\text{point is } \rightarrow \infty} \xrightarrow{D} N(0, 1) \quad \text{as } n \rightarrow \infty$$

$\downarrow$   
as  $\sqrt{nI(\theta)}$  SE

**Remarks:**

- We say that  $\hat{\theta}_n$  is **asymptotically unbiased**, which means that in large samples the MLE has approximately the desired mean.
- $\hat{\theta}_n$  is asymptotically efficient. This means that in large sample, it has the smallest variance among all asymptotically unbiased estimators.
- The third result says that we can use a relatively simple distribution to provide confidence intervals for  $\theta$ . In general, the actual sampling distribution of  $\hat{\theta}_n$  is very messy.

- $\sqrt{\hat{V}_n} = \left[ \sqrt{nI(\hat{\theta})} \right]^{-1}$  provides the estimated *asymptotic standard error* (ASE) for  $\hat{\theta}$ .
- $-\frac{\partial^2 \ell(\hat{\theta})}{\partial \theta^2}$  measures the *curvature* of the log-likelihood function.
- The greater the curvature, the greater the information about  $\theta$  and the smaller the ASE.

# Statistics 630

---

Example 44 again The log-likelihood for a random sample from the Bernoulli distribution is

$$\ell(\theta) = \log(L(\theta|y)) = \underbrace{y}_\text{sum} \log(\theta) + (n - y) \log(1 - \theta)$$

and the log likelihood for a single Bernoulli rv is

$$\log(f_\theta(y_i)) = y_i \log(\theta) + (1 - y_i) \log(1 - \theta)$$

The first and second derivatives are

$$\frac{\partial \log(f_\theta(y_i))}{\partial \theta} = \left( \frac{y_i}{\theta} - \frac{1 - y_i}{1 - \theta} \right)$$
$$\frac{\partial^2 \log(f_\theta(y_i))}{\partial \theta^2} = -\frac{y_i}{\theta^2} - \frac{1 - y_i}{(1 - \theta)^2}$$

— can either sum this and take expectation or take another derivative of this and take expectation and multiply by 1

Fisher's information in  $\theta$  for a single Bernoulli observation is

$$I(\theta) = -E_\theta \left[ \frac{\partial^2 \log(f_\theta(Y_i))}{\partial \theta^2} \right] = \frac{\theta}{\theta^2} + \frac{(1 - \theta)}{(1 - \theta)^2} = \frac{1}{\theta(1 - \theta)}.$$

# Statistics 630

---

Then Fisher's information in  $\theta$  for Bernoulli sample is  $nI(\theta) = \frac{n}{\theta(1-\theta)}$ ,

$$V_n(\theta_0) = \frac{1}{nI(\theta_0)} = \frac{\theta_0(1-\theta_0)}{n}, \quad \hat{\theta} = \frac{\sum y_i}{n}$$

and the estimated asymptotic standard error is

$$ASE = \sqrt{V_n(\hat{\theta})} = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}.$$

The asymptotic properties of the MLE imply that

$$\frac{\sqrt{n}(\hat{\theta} - \theta_0)}{\sqrt{\theta_0(1-\theta_0)}} \xrightarrow{D} N(0, 1) \quad \text{as } n \rightarrow \infty$$

and that

$$\hat{\theta} \stackrel{\text{approx}}{\sim} N(\theta_0, \frac{\theta_0(1-\theta_0)}{n})$$

This equivalent to our earlier normal approximation to the binomial distribution.

STOP CN

Fri day 11/5/21

Example 45 again Suppose that  $X \sim N(\mu, \sigma^2)$  where  $\sigma^2$  is known. Find the Fisher information in  $X$ .

$$f_\mu(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\log(f_\mu(x))) = -\log(\sqrt{2\pi}\sigma) - (x - \mu)^2/2\sigma^2$$

$$\frac{\partial \log(f_\mu(x))}{\partial \mu} = 2(x - \mu)/2\sigma^2 = \frac{x - \mu}{\sigma^2}$$

$$I(\mu) = E \left[ \frac{(X - \mu)^2}{\sigma^4} \right] = \frac{E[(X - \mu)^2]}{\sigma^4} = \frac{\sigma^2}{\sigma^4} = \frac{1}{\sigma^2}$$

$$I_n(\mu) = \frac{1}{\sigma^2}, \quad V_n(\mu) = \frac{\sigma^2}{n}$$

ML  
 $\hat{\mu} = \bar{x}$   
 $\text{var}(\bar{x}) = \frac{\sigma^2}{n}$

In this case  $\text{var}(\hat{\mu})$  is exactly  $\text{var}(\bar{x})$ , not always the case. In general, it's an approximation.

Example 45 continued Suppose now that  $X_1, \dots, X_n$  form a random sample from a  $N(\mu, \sigma^2)$  population where  $\sigma^2$  is known. Find the Fisher information in  $X_1, \dots, X_n$ .

From the preceding slide,

$$I(\mu) = \frac{1}{\sigma^2}.$$

Then the Fisher information  $I_n(\mu)$  in  $X_1, \dots, X_n$  is given by

$$I_n(\mu) = nI(\mu) = \frac{n}{\sigma^2}.$$

The large sample results for the MLE imply that

$$\frac{\sqrt{n}(\bar{X}_n - \mu_0)}{\sqrt{\sigma^2}} = \frac{\bar{X}_n - \mu_0}{\sigma/\sqrt{n}} \xrightarrow{D} N(0, 1) \quad \text{as } n \rightarrow \infty$$

We know that this result holds exactly for any  $n$ .

# Statistics 630

---

Example 46, again Find the asymptotic distribution of the mle.

We need to obtain Fisher's information,  $I(\lambda)$ . The likelihood and log-likelihood of a single  $X$  are  $\chi_i \sim \text{IID } \exp(\lambda) \quad \hat{\lambda} = \frac{1}{\bar{X}}$

$$f_\lambda(x) = \lambda e^{-\lambda x} \quad \text{and} \quad \log(f_\lambda(x)) = \log(\lambda) - \lambda x.$$

Then

$$\frac{\partial \log(f_\lambda(x))}{\partial \lambda} = \frac{1}{\lambda} - x,$$

and

$$\frac{\partial^2 \log(f_\lambda(x))}{\partial \lambda^2} = -\frac{1}{\lambda^2}.$$

Hence,  $I(\lambda) = \underbrace{\frac{1}{\lambda^2}}$  and  $I_n = \frac{1}{\bar{X}^2} \approx n \left( \frac{x^2}{n} \right)$

$$\frac{\sqrt{n}(\hat{\lambda} - \lambda_0)}{\lambda_0} \xrightarrow{D} N(0, 1) \quad \text{as} \quad n \rightarrow \infty$$

If we were want estimates, use delta method,  
 $\text{MLEs} \Rightarrow$  Fisher info.

## 4 Confidence Intervals

Point estimators provide a limited amount of information about the unknown parameter in that they only specify a single value. A more useful approach to estimation is to specify a range of plausible values together with some measure of plausibility. We will develop confidence intervals first for the mean of a normal population and then for general parameters.

+ my note: use CI's to do simulations to get weighted avg of stock price dependent on future anchoring

### 4.1 The $Z$ Confidence Interval for a Normal Mean

Again suppose that  $X_1, \dots, X_n$  forms a random sample from a normal distribution  $N(\mu, \sigma^2)$  where  $\sigma^2$  is known. We want to construct a random interval that will contain the mean  $\mu$  with a specified probability  $\gamma$ . Thus, we form an interval  $(l(x_1, \dots, x_n), u(x_1, \dots, x_n))$  such that

$$P_\mu[l(X_1, \dots, X_n) \leq \mu \leq u(X_1, \dots, X_n)] = \gamma.$$

We will use the sampling distribution of  $\bar{X}_n$  from Chapter 4 to derive this interval.

We let  $\phi(x)$  and  $\Phi(x)$  denote the pdf and cdf, respectively, of the standard normal distribution, and let  $c$  be a constant such that

$$P(-c < Z < c) = \int_{-c}^c \phi(x)dx = \Phi(c) - \Phi(-c) = \gamma.$$

↗ fixed  
 ↗ don't depend  
 on data-

Next

$$\begin{aligned} -c < Z < c &\iff -c < \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}} < c \\ &\iff \bar{X}_n - c\frac{\sigma}{\sqrt{n}} < \mu < \bar{X}_n + c\frac{\sigma}{\sqrt{n}} \end{aligned}$$

Thus,

$$P\left(\bar{X}_n - c\frac{\sigma}{\sqrt{n}} < \mu < \bar{X}_n + c\frac{\sigma}{\sqrt{n}}\right) = \gamma.$$

This says that with probability equal to  $\gamma$ , the unknown value of  $\mu$  falls between the random variables,

$$l(X_1, \dots, X_n) = \bar{X}_n - c\frac{\sigma}{\sqrt{n}} \text{ and } u(X_1, \dots, X_n) = \bar{X}_n + c\frac{\sigma}{\sqrt{n}}.$$

End CN

MN  $\mathcal{N}(\bar{x}, \sigma^2)$

The value of  $c$  is determined by the fact that  $P(-c < Z < c) = \gamma$ . This implies that

$$\gamma = \Phi(c) - \Phi(-c) = \Phi(c) - (1 - \Phi(c)) = 2\Phi(c) - 1.$$

Thus,  $\Phi(c) = \frac{1+\gamma}{2}$  and  $c = \Phi^{-1}((1 + \gamma)/2) = z_{(1+\gamma)/2}$ ,

where  $z_\alpha$  denotes the  $\alpha^{th}$  quantile of the standard normal distribution.

**Remark:** We can form the *likelihood interval*  $C(x_1, \dots, x_n)$  for  $\mu$  by letting

$$C(x_1, \dots, x_n) = \{\mu : L(\mu | x_1, \dots, x_n) \geq k(x_1, \dots, x_n)\}$$

for some  $k(x_1, \dots, x_n)$ . The text on page 327 shows that the likelihood interval for  $\mu$  is the same as the interval that we derived above.

## 4.2 A General Definition of Confidence Intervals

We now define a confidence interval for a function  $\psi(\theta)$  of a parameter  $\theta$  where the joint distribution of  $X_1, \dots, X_n$  has pdf/pmf  $f_\theta(x_1, \dots, x_n)$  where  $\theta \in \Omega$ .

We say that an interval  $C(X_1, \dots, X_n) = (l(X_1, \dots, X_n), u(X_1, \dots, X_n))$  is a level  $\gamma$  confidence interval for  $\psi(\theta)$  if

$$\begin{aligned} P_\theta(\psi(\theta) \in C(X_1, \dots, X_n)) &= P_\theta(l(X_1, \dots, X_n) \leq \psi(\theta) \leq u(X_1, \dots, X_n)) \\ &\geq \underbrace{\gamma}_{\text{for every } \theta \in \Omega}. \end{aligned}$$

Thus, the random interval  $C(X_1, \dots, X_n)$  has a probability of at least  $\gamma$  (usually 0.95 or larger) of containing the true value of  $\psi(\theta)$ .

In the last example we saw that

$(\bar{X}_n - z_{(1+\gamma)/2} \frac{\sigma}{\sqrt{n}} < \mu < \bar{X}_n + z_{(1+\gamma)/2} \frac{\sigma}{\sqrt{n}})$  is a level- $\gamma$  confidence interval for  $\mu$  when sampling from a normal distribution with known variance  $\sigma^2$ .

## 4.3 The $t$ Confidence Interval for a Normal Mean

We now consider the more realistic setting where  $X_1, \dots, X_n$  is a random sample from a normal  $(\mu, \sigma^2)$  distribution where

$\theta = (\mu, \sigma^2) \in \Omega = \mathcal{R} \times (0, \infty)$  is the unknown parameter. We will derive a level  $\gamma$  confidence interval for  $\psi(\theta) = \mu$ . The derivation will be similar to that where  $\sigma^2$  was known, but we will need to use the  $t$  distribution rather than the normal distribution in our derivation.

From Slide 44 of Chapter 4, we know that

$$T = \frac{\bar{X}_n - \mu}{S/\sqrt{n}} \sim t(n-1).$$

We let  $f_{n-1}(x)$  denote the pdf of the  $t$  distribution with  $n-1$  df, and let  $c$  be a constant such that

$$P(-c < T < c) = \int_{-c}^c f_{n-1}(x)dx = \gamma.$$

Next

$$-c < T < c \iff -c < \frac{\bar{X}_n - \mu}{S/\sqrt{n}} < c$$

$$\iff \bar{X}_n - c\frac{S}{\sqrt{n}} < \mu < \bar{X}_n + c\frac{S}{\sqrt{n}}$$

Typical way to get a CI  
on  $\bar{X}$  when we have a  
sample from a normal  
dist. If we don't have  
a r.s. from a normal  
 $\checkmark$  but  $n$   
is large  
enough s.t.  $\bar{X} \approx N$

$$P\left(\bar{X}_n - c\frac{S}{\sqrt{n}} < \mu < \bar{X}_n + c\frac{S}{\sqrt{n}}\right) = \gamma. \quad \text{People still use this approach for C.I.s.}$$

This says that the probability equals  $\gamma$  that the unknown value of  $\mu$  falls between the random variables,  $l = \bar{X}_n - c\frac{S}{\sqrt{n}}$  and  $u = \bar{X}_n + c\frac{S}{\sqrt{n}}$ .

Here  $c = t_{(1+\gamma)/2}(n-1)$  where  $t_\alpha(\lambda)$  denotes the  $\alpha^{th}$  quantile of the  $t(\lambda)$  distribution.

# Statistics 630

---

In applications, we observe the data  $X_1 = x_1, \dots, X_n = x_n$  and compute  
 $l(x_1, \dots, x_n) = \bar{x}_n - t_{(1+\gamma)/2}(n-1) \frac{s}{\sqrt{n}}$  and  
 $u(x_1, \dots, x_n) = \bar{x}_n + t_{(1+\gamma)/2}(n-1) \frac{s}{\sqrt{n}}.$

We say that the interval  $(l, u)$  is a confidence interval for  $\mu$  with confidence coefficient  $\gamma$ .

Often we will write the  $\gamma$  confidence interval for  $\mu$  as

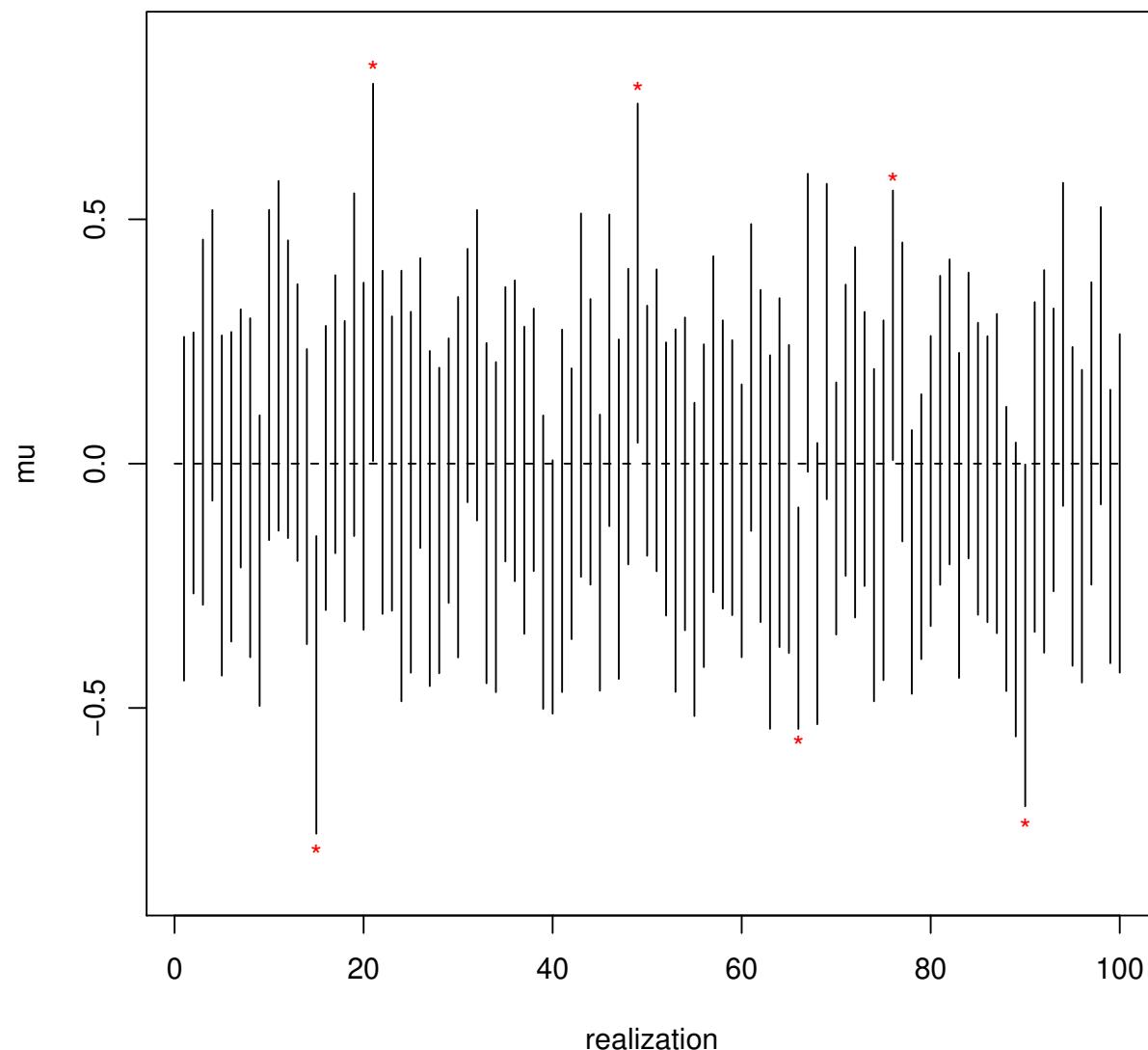
$$\bar{x} \pm t_{(1+\gamma)/2}(n-1) \frac{s}{\sqrt{n}} \quad \cancel{\text{X}}$$

We say that  $\mu$  lies in this observed interval with confidence (not probability)  $\gamma$ .

The frequentist interpretation of confidence intervals follows:

If we construct a large number of level  $\gamma$  confidence intervals for  $\mu$ , about 100( $\gamma$ )% of them will contain  $\mu$ . We illustrate this in the figure on the next slide.

## 100 Confidence Intervals for the Mean



**Remark:** The construction of the confidence interval for  $\mu$  depended on the random variable  $T$  which was a function of the data and the parameter. The important aspect of  $T$  was that it had a distribution that did not depend on the unknown parameter. We call such a random variable a **pivot**.

Example 45 again Suppose that  $X_1, \dots, X_n$  form a random sample from a  $N(\mu, \sigma^2)$  distribution. We want to construct a level  $\gamma$  confidence interval for  $\sigma^2$ .

Earlier we learned that

$$W = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{\sigma^2} \sim \chi^2(n-1).$$

We can use  $W$  as a pivot.

Let  $g_{n-1}(x)$  be the pdf of the  $\chi^2$  distribution with  $n - 1$  df, and let  $c_1$  and  $c_2$  be constants such that

$$P(c_1 < W < c_2) = \int_{c_1}^{c_2} g_{n-1}(x)dx = \gamma.$$

Then

$$\begin{aligned} c_1 < W < c_2 &\iff c_1 < \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{\sigma^2} < c_2 \\ &\iff \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{c_2} < \sigma^2 < \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{c_1}. \end{aligned}$$

Thus,

$$\left( \frac{\sum_{i=1}^n (x_i - \bar{x}_n)^2}{c_2}, \frac{\sum_{i=1}^n (x_i - \bar{x}_n)^2}{c_1} \right)$$

is a level  $\gamma$  confidence interval for  $\sigma^2$ .

## 4.4 Pivots

We summarize our approach to forming confidence intervals through the use of pivots. Suppose that  $X_1, \dots, X_n$  form a random sample from a distribution with pmf or pdf  $f_\theta(x)$ .

A pivot,  $W(X_1, \dots, X_n, \theta)$ , is a random variable whose distribution does not depend on  $\theta$ . (Often  $W(x_1, \dots, x_n, \theta)$  is a monotone function of  $\theta$ .)

Define  $a$  and  $b$  to be constants such that

$$P(a < W(X_1, \dots, X_n, \theta) < b) = \gamma.$$

When we observe  $X_1 = x_1, \dots, X_n = x_n$ , we define the level  $\gamma$  confidence interval for  $\theta$  as

$$C(x_1, \dots, x_n) = \{\theta : a < W(x_1, \dots, x_n, \theta) < b\}.$$

## 4.5 Approximate Confidence Intervals for Based on the MLE

Whenever the asymptotic distribution of an estimator  $\hat{\theta}$  is normal, i.e.,

$$\hat{\theta} \stackrel{\text{approx}}{\sim} N(\theta_0, \hat{V}_n),$$

we can form an approximate pivot:

$$\frac{\hat{\theta} - \theta}{\sqrt{\hat{V}_n}} \stackrel{\text{approx}}{\sim} N(0, 1)$$

We use this pivot to form an approximate  $100(\gamma)\%$  confidence interval for  $\theta$ :

$$\hat{\theta} \pm Z_{(1+\gamma)/2} \text{ASE}(\hat{\theta})$$

where  $Z_{(1+\gamma)/2}$  is a value such that  $P[Z \leq Z_{(1+\gamma)/2}] = (1 + \gamma)/2$  and  $Z$  is a standard normal rv. This interval is called the **Wald confidence interval for  $\theta$** .

---

Example 46 again Suppose that  $X \sim \text{Binomial}(n, \theta)$ . We can represent  $X$  as  $X_1 + \dots + X_n$  where  $X_i \sim \text{Bernoulli}(\theta)$ . Then the Fisher's information in a single  $X_i$  is

$$I(\theta) = \frac{1}{\theta(1-\theta)},$$

Then

$$V_n(\theta) = \frac{1}{nI(\theta)} = \frac{\theta(1-\theta)}{n}$$

and

$$\text{ASE}(\hat{\theta}) = \sqrt{\frac{1}{nI(\hat{\theta})}} = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}}$$

A level  $\gamma$  confidence interval for  $\theta$  is given by

$$\hat{\theta} \pm Z_{(1+\gamma)/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n}} \quad \text{where} \quad \hat{\theta} = \frac{x}{n} = \frac{\sum_{i=1}^n x_i}{n}$$

Alternatively, we could form a different pivot that is also based on the asymptotic properties of the MLE. We know that

$$\frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\theta(1-\theta)}} \xrightarrow{D} N(0, 1) \quad \text{as } n \rightarrow \infty$$

Thus,

$$P\left[-Z_{(1+\gamma)/2} < \frac{\sqrt{n}(\hat{\theta} - \theta)}{\sqrt{\theta(1-\theta)}} < Z_{(1+\gamma)/2}\right] = \gamma$$

We obtain the level  $1 - \alpha$  confidence interval:

$$\left\{ \theta : n(\hat{\theta} - \theta)^2 < \theta(1-\theta)Z_{(1+\gamma)/2}^2 \right\}$$

or

$$\frac{\hat{\theta} + \frac{Z_{(1+\gamma)/2}^2}{2n} \pm Z_{(1+\gamma)/2} \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{n} + \frac{Z_{(1+\gamma)/2}^2}{4n^2}}}{1 + Z_{(1+\gamma)/2}^2/n}$$

This interval is called the score confidence interval for  $\theta$ .

## Discussion of the Confidence Intervals for the Binomial Distribution

The two forms of confidence intervals for the binomial proportion will produce similar results for large sample sizes ( $n > 100$ ) as long as the observed number of successes  $x$  and observed number of failures  $n - x$  are not close to zero (at least 15). For smaller sample sizes and in the situation where  $x$  or  $n - x$  is small, the score interval has much better coverage probability.

*Example:* In a survey of 277 randomly selected adult shoppers, 69 stated that if an advertised item is unavailable they request a rain check, construct 95% CIs for  $\theta$ .

The 95% Wald interval is

$$0.249 \pm 1.96\sqrt{0.249(1 - 0.249)/277} = 0.249 \pm 0.051$$

The 95% Wald interval is (.198, .300).

# Statistics 630

---

We can also compute the 95% score confidence interval:

$$\frac{0.249 + \frac{1.96^2}{2(277)} \pm 1.96 \sqrt{\frac{0.249 \cdot 0.751}{277} + \frac{1.96^2}{4(277)^2}}}{1 + (1.96^2)/277}$$

We obtain the 95% score confidence interval for  $\theta$ :  $(0.202, 0.303)$ . This appears similar to the 95% Wald interval,  $(.198, .300)$ .

**Bill of Rights Data:** Recall that  $x = 14$  and  $n = 50$  in this example.

The 95% Wald interval is  $(0.156, 0.404)$  whereas the 95% score interval is  $(0.175, 0.417)$ . There is a greater difference between the two intervals due to the small sample size.

End CP

Friday 11/12/12

Example 46 again Suppose that  $X_1, \dots, X_n$  is a random sample from an exponential ( $\lambda$ ) distribution. Obtain exact and approximate level  $\gamma$  confidence interval for  $\lambda$ .

**Exact confidence interval:** Consider the random variable

$W = 2\lambda \sum_{i=1}^n X_i$ . Its moment generating function is

$$M_W(s) = M_{2\lambda \sum_{i=1}^n X_i}(s) = \prod_{i=1}^n M_{X_i}(2\lambda s) = \left( \frac{\lambda}{\lambda - 2\lambda s} \right)^n = \left( \frac{1}{1 - 2s} \right)^{2n/2}.$$

Thus,  $W$  has a chi-squared distribution with  $2n$  degrees of freedom. Let  $a$  and  $b$  be values such that  $G_{2n}(b) - G_{2n}(a) = \gamma$  where  $G_{2n}(x)$  is the cdf of the chi-squared distribution with  $2n$  degrees of freedom. Then

$$\begin{aligned} \gamma &= P[a \leq W \leq b] = P[a \leq 2\lambda \sum_{i=1}^n X_i \leq b] \\ &= P\left[\frac{a}{2 \sum_{i=1}^n X_i} \leq \lambda \leq \frac{b}{2 \sum_{i=1}^n X_i}\right] \end{aligned}$$

Thus, an exact level  $\gamma$  confidence interval for  $\lambda$  is

$$\left[ \frac{a}{2 \sum_{i=1}^n X_i}, \frac{b}{2 \sum_{i=1}^n X_i} \right]$$

**Approximate confidence interval I:** From slide 64, we have the result that  $\hat{\lambda} \sim AN(\lambda_0, \lambda_0^2/n)$ . For the Wald interval, we will use the pivot  $Z = \frac{\hat{\lambda} - \lambda_0}{\hat{\lambda}/\sqrt{n}}$ .

This has an approximate  $N(0, 1)$  distribution and we have

$$\begin{aligned}\gamma &= P \left[ -Z_{(1+\gamma)/2} \leq \frac{\hat{\lambda} - \lambda_0}{\hat{\lambda}/\sqrt{n}} \leq Z_{(1+\gamma)/2} \right] \\ &= P \left[ \hat{\lambda} - \frac{Z_{(1+\gamma)/2} \hat{\lambda}}{\sqrt{n}} \leq \lambda_0 \leq \hat{\lambda} + \frac{Z_{(1+\gamma)/2} \hat{\lambda}}{\sqrt{n}} \right].\end{aligned}$$

Thus,

$$\left[ \hat{\lambda} \left( 1 - \frac{Z_{(1+\gamma)/2}}{\sqrt{n}} \right), \hat{\lambda} \left( 1 + \frac{Z_{(1+\gamma)/2}}{\sqrt{n}} \right) \right]$$

is an approximate level  $\gamma$  confidence interval for  $\lambda$ .

**Approximate confidence interval II:** From slide 64, we have the result that  $\hat{\lambda} \sim AN(\lambda_0, \lambda_0^2/n)$ . For the score interval, we will use the pivot  $Z = \frac{\hat{\lambda} - \lambda_0}{\lambda_0/\sqrt{n}}$ . This has an approximate  $N(0, 1)$  distribution and we have

$$\begin{aligned}\gamma &= P \left[ -Z_{(1+\gamma)/2} \leq \frac{\hat{\lambda} - \lambda_0}{\lambda_0/\sqrt{n}} \leq Z_{(1+\gamma)/2} \right] \\ &= P \left[ \frac{-Z_{(1+\gamma)/2}}{\sqrt{n}} \leq \frac{\hat{\lambda}}{\lambda_0} - 1 \leq \frac{Z_{(1+\gamma)/2}}{\sqrt{n}} \right] \\ &= P \left[ \frac{\hat{\lambda}}{1 + \frac{Z_{(1+\gamma)/2}}{\sqrt{n}}} \leq \lambda_0 \leq \frac{\hat{\lambda}}{1 - \frac{Z_{(1+\gamma)/2}}{\sqrt{n}}} \right].\end{aligned}$$

Thus,

$$\left[ \frac{\hat{\lambda}}{1 + \frac{Z_{(1+\gamma)/2}}{\sqrt{n}}}, \frac{\hat{\lambda}}{1 - \frac{Z_{(1+\gamma)/2}}{\sqrt{n}}} \right]$$

is an approximate level  $\gamma$  confidence interval for  $\lambda$ .

Example 39 again Find a 95% confidence intervals for the rate parameter in the nerve pulse data. For this set of data,  $n = 799$  and  $\sum_{i=1}^{799} X_i = 174.64$ .

- Exact 95% confidence interval with equal tail probabilities: Let  $G$  denote the cdf of the  $\chi^2(1598)$  distribution. Then the exact 95% ci is

$$\left[ \frac{G^{-1}(0.025)}{2 \sum_{i=1}^{799} X_i}, \frac{G^{-1}(0.975)}{2 \sum_{i=1}^{799} X_i} \right] = \left[ \frac{1489.1}{2(174.64)}, \frac{1710.7}{2(174.64)} \right] = (4.262, 4.898) \quad \text{X}$$

- 95% Wald interval based on the mle:

$$\hat{\lambda} \pm Z_{0.975} \hat{\lambda} / \sqrt{n} = 4.575 \pm 1.96 \times 4.575 / \sqrt{799} = 4.575 \pm 0.317$$

Thus, the approximate 95% confidence interval is  $(4.258, 4.892) \quad \text{X}$

all 3 are  
very close  
for large  
 $n$ .

- 95% Score interval based on the mle:

$$\left[ \frac{\hat{\lambda}}{1 + \frac{Z_{0.975}}{\sqrt{n}}}, \frac{\hat{\lambda}}{1 - \frac{Z_{0.975}}{\sqrt{n}}} \right] = \left[ \frac{4.575}{1 + \frac{1.96}{\sqrt{799}}}, \frac{4.575}{1 - \frac{1.96}{\sqrt{799}}} \right] = (4.278, 4.916) \quad \text{X}$$

## 4.6 Bootstrap Confidence Intervals

The bootstrap can be used to obtain approximate confidence intervals for a function of a population parameter,  $\psi(\theta) = T(F_\theta)$  where  $T$  is a function of the distribution  $F_\theta$ . For example,  $T$  could refer to a moment of  $F_\theta$  or a quantile of  $F_\theta$ . We can write the estimate based on the original data as  $\hat{\psi} = T(F_n)$  where  $F_n$  is the empirical cdf,

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(x_i).$$

This is a step function with jumps of size  $1/n$  at each single value of  $x_i$ . We note that  $F_n$  is a cdf that puts weight  $1/n$  at each observed value of  $x_i$ .

To construct a confidence interval for  $\psi(\theta) = T(F_\theta)$ , we need to know the distribution of  $\hat{\psi} = T(F_n)$ . The bootstrap provides a computational method for approximating the distribution of  $\hat{\psi}$ .

Since  $F_\theta$  is not known, we can use the nonparametric bootstrap to approximate the distribution of  $\hat{\psi} = T(F_n)$ :

- Generate  $B$  bootstrap samples of size  $n$  from the distribution with cdf  $F_n$ .  
That is, take a sample of size  $n$  with replacement from  $\{x_1, \dots, x_n\}$ .
- For each bootstrap sample  $\{x_1^*, \dots, x_n^*\}$ , compute the estimate of  $\psi(\theta)$ , say  $\hat{\psi}_j^*$ ,  $j = 1, \dots, B$ .
- Compute the *bootstrap percentile confidence interval* for  $\psi(\theta)$  by computing

$$(\hat{\psi}_{(1-\gamma)/2}, \hat{\psi}_{(1+\gamma)/2}),$$

where  $\hat{\psi}_p$  refers to the  $p^{th}$  quantile of the bootstrap estimates,  $\{\hat{\psi}_1^*, \dots, \hat{\psi}_B^*\}$ .

**Remark:** If the form of  $F_\theta$  is known, the parametric bootstrap replaces the first step by taking a random sample of size  $n$  from a distribution with cdf  $F_{\hat{\theta}}$  where  $\hat{\theta}$  is the mle of  $\theta$ .

**$t$  Confidence Intervals** An alternative approach to forming a bootstrap confidence interval is to form a  $t$  interval centered at the estimate for the observed data and then using the bootstrap estimate of the standard error of the estimator. This results in the  $t$  confidence interval for  $\psi(\theta)$ :

$$\hat{\psi} \pm t_{(1+\gamma)/2} \sqrt{\widehat{\text{Var}}_{\hat{F}}(\hat{\psi})},$$

where  $\widehat{\text{Var}}_{\hat{F}}(\hat{\psi})$  is the sample variance of the  $B$  bootstrap estimates of  $\psi$ .

**Remark:** The above methods for forming confidence intervals are very basic applications of the bootstrap. There are other approaches using the bootstrap for confidence intervals that result in better performance, especially when the bootstrap distribution of the estimator is skewed. See the text by Efron and Tibshirani for more details.

Example 39 again Find a 95% confidence interval for the rate parameter in the nerve pulse data.

- Wald interval based on the mle: From Slide 64,  $I(\lambda) = 1/\lambda^2$ . Thus, the approximate 95% confidence interval is  $(4.258, 4.892)$ :

$$\hat{\lambda} \pm Z(0.025)\hat{\lambda}/\sqrt{n} = 4.575 \pm 1.96 \times 4.575/\sqrt{799} = 4.575 \pm 0.317$$

- Bootstrap confidence intervals:

```
> temp=rep(0,10000)
> for (i in 1:10000)temp[i]=1/mean(sample(nerve,replace=TRUE))
> quantile(temp,c(0.025,0.975))
  2.5%    97.5%
4.280029 4.896433
> 1/mean(nerve)-1.96*sd(temp)
[1] 4.26618
> 1/mean(nerve)+1.96*sd(temp)
[1] 4.884072
```

The resulting confidence intervals are  $(4.280, 4.896)$  and  $(4.266, 4.884)$ .

- Parametric bootstrap confidence intervals:

When  $F_\theta$  has a known form, we can use the parametric bootstrap to approximate the distribution of  $\hat{\psi} = T(F_n)$ :

```
> temp=rep(0,10000)
> for (i in 1:10000)temp[i]=1/mean(rexp(799,rate=1/mean(nerve)))
> quantile(temp,c(0.025,0.975))
  2.5%    97.5%
4.275859 4.904623
> 1/mean(nerve)-1.96*sd(temp)
[1] 4.258208
> 1/mean(nerve)+1.96*sd(temp)
[1] 4.892044
```

The resulting confidence intervals are  $(4.276, 4.905)$  and  $(4.258, 4.892)$ .

# Statistics 630

---

We summarize the various 95% confidence intervals:

Method	Interval
Exact	(4.262, 4.898)
Wald	(4.258, 4.892)
Score	(4.278, 4.916)
Nonparametric bootstrap quantile	(4.280, 4.896)
Nonparametric $t$	(4.266, 4.884)
Parametric bootstrap quantile	(4.276, 4.905)
Parametric $t$	(4.258, 4.892)