

1 **Part 3: Properties of MLEs**

2 **Summary**

3 In this Part we explore some of the important properties possessed by
4 maximum likelihood estimators.

5 Briefly stated, we will see that, asymptotically, these estimators cannot be
6 beat, in the sense that they have minimal MSE, and also are approximately
7 normally distributed with variance that is easy to calculate. This makes
8 quantifying the error in the estimators straightforward.

1 Invariance of the MLE

2 The invariance result can be stated as follows: If $\hat{\theta}_{\text{MLE}}$ is the MLE for θ , then
3 $g(\hat{\theta}_{\text{MLE}})$ is the MLE for $g(\theta)$.

4 This is a useful, important result because it is often the case that we are
5 truly interested in estimating some function of the parameter θ , not θ itself.

6 Using this, we can say, for example, that if X_1, X_2, \dots, X_n are i.i.d.
7 normal(μ, σ^2), then the MLE for σ is

8
$$\sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

9 and the MLE for the **signal to noise ratio** μ/σ is

10
$$\bar{X} / \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$$

- 1 One consequence of this property is that estimation is not affected by how
2 one chooses to **parameterize** a model.
- 3 For example, some choose to represent the Exponential distribution using
4 the “mean parameterization” which leads to $E(X) = \mu$, whereas we have
5 used the “rate parameterization” which leads to $E(X) = 1/\lambda$. Because
6 of invariance, these two approaches will lead to the same estimates, i.e.,
7 $\hat{\mu} = 1/\hat{\lambda}$.

1 **Exercise:** Does unbiasedness have this invariance property?

2 No. If $\hat{\theta}$ is an unbiased estimator for θ , no guarantee
3 that $g(\hat{\theta})$ is unbiased for $g(\theta)$. For example,
4 s is a biased estimator for σ , even when s^2 is
5 unbiased for σ^2 .

6 **Exercise:** Suppose that X_1, X_2, \dots, X_n are i.i.d. $\text{Exponential}(\lambda)$. What is the
7 MLE of the **tail probability** $P(X_i > a)$? assuming $a > 0$.

8 In exponential case, $P(X_i > a) = 1 - F_{X_i}(a) = e^{-\lambda a}$ (if $a > 0$)

9 The MLE for λ is $\hat{\lambda} = 1/\bar{x}$. So MLE for
10 $e^{-\lambda a}$ is $e^{-a/\bar{x}}$.

1 Consistency of the MLE

2 An estimator $\hat{\theta}$ is **consistent for θ** if as the sample size n increases, the
3 estimator converges to the true value of the parameter.

4 Formally, we would write: For any $\epsilon > 0$,
5 $\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta_0| > \epsilon) = 0$ true value of parameter θ

"convergence in prob."

6 This is a fairly low standard placed on an estimator, and, fortunately, under
7 certain "regularity conditions," the MLE is consistent.

Asymptotic Normality of the MLE

Another great property of the MLE: The MLE is asymptotically normal, with variance that is not difficult to calculate, under certain “regularity conditions.” (See page 516 from Casella and Berger, Second Edition.)

This will allow an easy route to:

1. Calculate standard errors for MLEs

2. Constructing confidence intervals for parameters

3. Constructing confidence intervals for general functions of parameters (via the Delta method)

4. Constructing hypothesis tests for parameters

From Casella and Berger, Second Edition

516

ASYMPTOTIC EVALUATIONS

Section 10.6

Although estimators such as d_n , called *superefficient*, can be constructed in some generality, they are more of a theoretical oddity than a practical concern. This is because the values of θ for which the variance goes below the bound are a set of Lebesgue measure 0. However, the existence of superefficient estimators serves to remind us to always be careful in our examination of assumptions for establishing properties of estimators (and to be careful in general!).

10.6.2 Suitable Regularity Conditions

The phrase “under suitable regularity conditions” is a somewhat abused phrase, as with enough assumptions we can probably prove whatever we want. However, “regularity conditions” are typically very technical, rather boring, and usually satisfied in most reasonable problems. But they are a necessary evil, so we should deal with them. To be complete, we present a set of regularity conditions that suffice to rigorously establish Theorems 10.1.6 and 10.1.12. These are not the most general conditions but are sufficiently general for many applications (with a notable exception being if the MLE is on the boundary of the parameter space). Be forewarned, the following is not for the fainthearted and can be skipped without sacrificing much in the way of understanding.

These conditions mainly relate to differentiability of the density and the ability to interchange differentiation and integration (as in the conditions for Theorem 7.3.9). For more details and generality, see Stuart, Ord, and Arnold (1999, Chapter 18), Ferguson (1996, Part 4), or Lehmann and Casella (1998, Section 6.3).

The following four assumptions are sufficient to prove Theorem 10.1.6, consistency of MLEs:

- (A1) We observe X_1, \dots, X_n , where $X_i \sim f(x|\theta)$ are iid.
- (A2) The parameter is *identifiable*; that is, if $\theta \neq \theta'$, then $f(x|\theta) \neq f(x|\theta')$.
- (A3) The densities $f(x|\theta)$ have common support, and $f(x|\theta)$ is differentiable in θ .
- (A4) The parameter space Ω contains an open set ω of which the true parameter value θ_0 is an interior point.

True value not on boundary of parameter space

The next two assumptions, together with (A1)–(A4) are sufficient to prove Theorem 10.1.12, asymptotic normality and efficiency of MLEs.

- (A5) For every $x \in \mathcal{X}$, the density $f(x|\theta)$ is three times differentiable with respect to θ , the third derivative is continuous in θ , and $\int f(x|\theta) dx$ can be differentiated three times under the integral sign.
- (A6) For any $\theta_0 \in \Omega$, there exists a positive number c and a function $M(x)$ (both of which may depend on θ_0) such that

$$\left| \frac{\partial^3}{\partial \theta^3} \log f(x|\theta) \right| \leq M(x) \quad \text{for all } x \in \mathcal{X}, \quad \theta_0 - c < \theta < \theta_0 + c.$$

with $E_{\theta_0}[M(X)] < \infty$.

1 To start, we focus on the case where the parameter θ is one-dimensional.

2 The formal result can be stated as follows:

3 If X_1, X_2, \dots, X_n are i.i.d. $f(x; \theta_0)$, then, under suitable regularity condi-
4 tions, as n increases,

$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow N(0, I^{-1}(\theta_0))$$

6 where

$$I(\theta) = E_{\theta} \left(- \frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right)$$

8 is the **Fisher Information**.

9 In the definition of $I(\theta)$, X is random, and the subscript on E_{θ} is meant to
10 emphasize that X has distribution specified by $f(x; \theta)$.

1 **Comment:** It would be technically more correct to define the Fisher infor-
2 mation as

$$3 \quad I(\theta) = E_{\theta} \left(- \frac{\partial^2 \log f(X; t)}{\partial t^2} \Big|_{t=\theta} \right)$$

4 since this makes a distinction between the θ which is serving as the input
5 into the function $I(\theta)$ and the θ used in the two derivatives of $\log f(X; \theta)$.

6 You do not typically see it written this way (even in advanced texts) because
7 it just makes things more messy without much conceptual gain.

1 **Exercise:** Explain the practical value of the above result.

2
$$\sqrt{n}(\hat{\theta} - \theta_0) \rightarrow N(0, I^{-1}(\theta_0)) \text{ as } n \rightarrow \infty$$

3
$$\sqrt{n}(\hat{\theta} - \theta_0) \approx N(0, I^{-1}(\theta_0)) \text{ for } n \text{ "large"}$$

4
$$\hat{\theta}_{MLE} \approx N\left(\theta_0, \frac{1}{nI(\theta_0)}\right)$$

6
$$\text{So, } E(\hat{\theta}_{MLE}) \approx \theta_0, \text{ i.e. approx. unbiased for } \theta$$

7 and
$$V(\hat{\theta}_{MLE}) \approx \frac{1}{nI(\theta_0)} \approx \frac{1}{nI(\hat{\theta})}$$

8
$$\text{SE of } \hat{\theta}_{MLE} \approx \frac{1}{\sqrt{nI(\hat{\theta})}}$$

11 Also, $\hat{\theta}_{MLE} \approx \text{Normal}$.

1 **Exercise:** Show that

$$nI(\theta) = E_{\theta} \left[-\frac{\partial^2}{\partial \theta^2} \log L(\theta) \right]$$

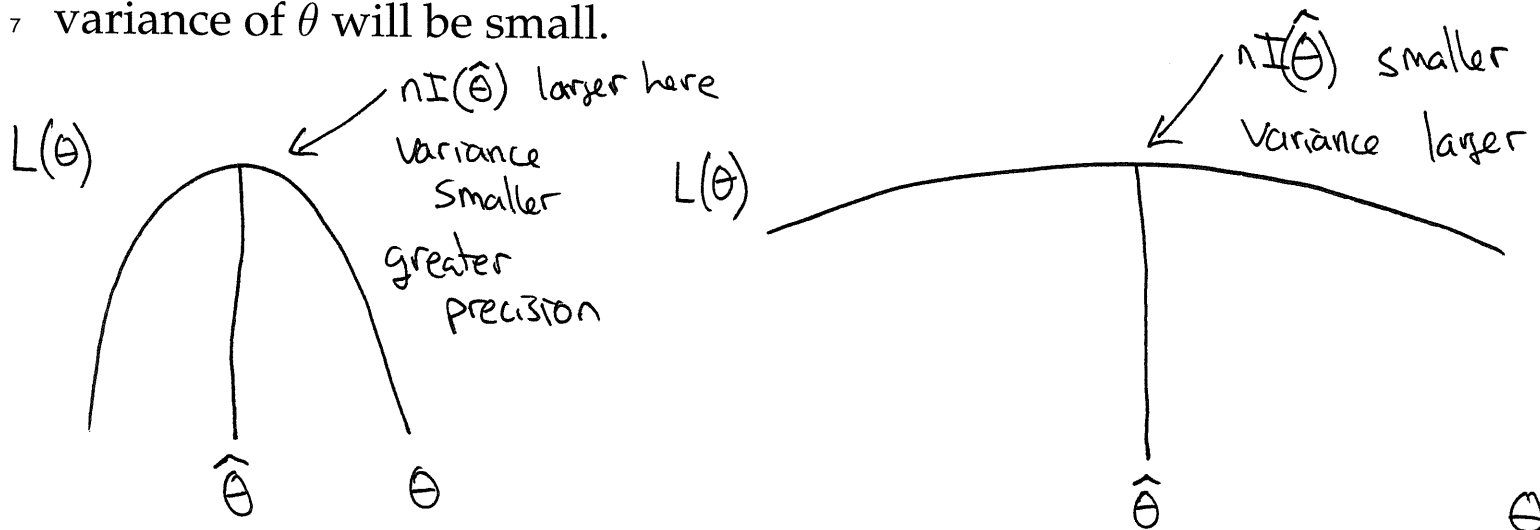
$$nI(\theta) = n E_{\theta} \left[-\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right]$$

$$= \sum_{i=1}^n E_{\theta} \left[-\frac{\partial^2 \log f(X_i; \theta)}{\partial \theta^2} \right] \quad \text{since } X_1, X_2, \dots, X_n \text{ are iid}$$

$$= E_{\theta} \left[-\frac{\partial^2}{\partial \theta^2} \sum_{i=1}^n \log f(X_i; \theta) \right]$$

$$= E_{\theta} \left[-\frac{\partial^2}{\partial \theta^2} \ell(\theta) \right]$$

- 1 The previous results show that the variance of the MLE is approximately
- 2 equal to the reciprocal of the expected value of the negative second deriva-
- 3 tive of the log likelihood at its peak.
- 4 And, “the negative of the second derivative” is the amount of curvature at
- 5 the peak of the likelihood.
- 6 If the log likelihood has a “sharp” peak, then $nI(\hat{\theta})$ will be large, and the
- 7 variance of $\hat{\theta}$ will be small.



- 1 **Exercise:** Assume X_1, X_2, \dots, X_n are i.i.d. $N(\mu, \sigma^2)$, where σ^2 is known. We
2 know that \bar{X} is the MLE for μ . Derive the Fisher Information for μ , and use
3 it to approximate the standard error and distribution of \bar{X} .

$$4 \quad \ell(\mu) = -\frac{\sum (x_i - \mu)^2}{2\sigma^2} + \text{constant} \quad (\text{see Slide 30 of previous part})$$

$$6 \quad \frac{\partial \ell(\mu)}{\partial \mu} = \frac{\sum (x_i - \mu)}{\sigma^2} = \frac{\sum x_i - n\mu}{\sigma^2}$$

$$9 \quad \frac{\partial^2 \ell(\mu)}{\partial \mu^2} = -\frac{n}{\sigma^2}$$

$$12 \quad nI(\mu) = E_{\mu} \left[-\frac{\partial^2}{\partial \mu^2} \ell(\mu) \right] = E_{\mu} \left[\frac{n}{\sigma^2} \right] = \frac{n}{\sigma^2}$$

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

The asymptotic normality of MLE tells me that \bar{X} is approx. normal with mean μ_0 and variance

$$\frac{1}{nI(\mu_0)} = \frac{1}{n/\sigma^2} = \frac{\sigma^2}{n}$$

In fact, we already knew that \bar{X} is exactly normal with mean μ_0 and variance σ^2/n

- 1 **Exercise:** Assume that X_1, X_2, \dots, X_n are i.i.d. $\text{Exponential}(\lambda)$. Derive the
2 Fisher Information for λ and the approximate distribution of its MLE.

3 Already found $\hat{\lambda}_{\text{MLE}} = 1/\bar{X}$, $f(x; \lambda) = \lambda e^{-\lambda x}$, $x > 0$

4 $\log f(x; \lambda) = \log \lambda - \lambda x$, $x > 0$

5
$$\frac{\partial \log f(x; \lambda)}{\partial \lambda} = \frac{1}{\lambda} - x$$

6
$$\frac{\partial^2 \log f(x; \lambda)}{\partial \lambda^2} = -\frac{1}{\lambda^2}$$

7
8
9
10 So, $I(\lambda) = E_{\lambda} \left[-\frac{\partial^2 \log f(X; \lambda)}{\partial \lambda^2} \right] = E_{\lambda} \left[+\frac{1}{\lambda^2} \right] = \frac{1}{\lambda^2}$

11
12 So, $\hat{\lambda}_{\text{MLE}}$ is approx. normal with mean λ_0 and ~~var~~
variance $\frac{1}{nI(\lambda_0)} = \lambda_0^2/n$. In practice, use $\hat{\lambda}^2/n$ as
approx. variance.

- 1 **Exercise:** Assume that X_1, X_2, \dots, X_n are i.i.d. $\text{Poisson}(\lambda)$. Derive the Fisher
2 Information for λ and the approximate distribution of its MLE.

3 Approx. variance is $\hat{\lambda}^2/n$

4

5

6

7

8

9

10

11

12

1 Efficiency of the MLE

2 It is possible to show that the lower bound on the variance of any unbiased
3 estimator is $1/nI(\theta_0)$. This is called the **Cramer-Rao lower bound**.

4 Since the MLE is asymptotically unbiased, we know that we are achieving
5 the Cramer-Rao lower bound, at least asymptotically.

6 Hence, we are minimizing the MSE (asymptotically). This is objectively the
7 best rationale for utilizing maximum likelihood estimators.

1 The Delta Method

2 Often, our objective is not to estimate θ , but instead some function of θ , call
3 it $g(\theta)$. For example, often want to estimate the probability of some event

4 A. This probability is not a parameter, but it is a function of θ .

5 From the invariance property we know that the MLE of $g(\theta)$ is $g(\hat{\theta})$. But,
6 we also need to calculate a SE or construct a confidence interval for this
7 estimate.

8 One approach would be to “reparameterize” the model in terms of $g(\theta)$, and
9 then find $I(g(\theta))$, and then use $1/n I(g(\theta))$ as the variance, etc. Fortunately,
10 there is an easier way...

1 Recall the **Delta Method**:

2 If Y_n is approximately $N(\mu, \sigma_n^2)$ and $\sigma_n^2 \rightarrow 0$ as n increases, then
3 $g(Y_n)$ is approximately $N(g(\mu), (g'(\mu))^2 \sigma_n^2)$, assuming that $g'(\mu)$ ex-
4 ists and is not zero.

5 Here we will use this result with the MLE for θ in place of Y_n .

6 Then, $\mu = \theta_0$ and $\sigma_n^2 = 1/nI(\theta_0)$.

7 We can conclude that $g(\hat{\theta})$ is approximately
8
$$N\left(g(\theta_0), g'(\theta_0)^2 \left(\frac{1}{nI(\theta_0)}\right)\right)$$

9 In practice, we use the approximate variance:

10
$$g'(\hat{\theta})^2 \left(\frac{1}{nI(\hat{\theta})}\right)$$

- 1 **Exercise:** Suppose that X_1, X_2, \dots, X_n are i.i.d. Exponential(λ). We seek to
2 estimate $\tau = P(X_i > a)$ for some constant $a > 0$. Construct the MLE and
3 find its standard error. Also construct a confidence interval for τ .

4 By invariance prop, $\hat{\tau}_{MLE} = e^{-a\hat{\lambda}}$, i.e. $g(\lambda) = e^{-a\lambda}$

5 We know $I(\lambda) = 1/\lambda^2$ from Slide 14

6 Seek to use Δ -method

7 $g'(\lambda) = -ae^{-a\lambda}$

8 So, the Δ -method says that $\hat{\tau}_{MLE} = g(\hat{\lambda}_{MLE})$ is

9 approx. normal with mean $g(\lambda_0) = e^{-a\lambda_0}$ and variance

10
$$(g'(\lambda_0))^2 \left(\frac{1}{nI(\lambda_0)} \right) = a^2 e^{-2a\lambda_0} \left(\frac{\lambda_0^2}{n} \right)$$

1 In practical terms, $\hat{\tau}_{MLE}$ is approx. normal with
2 variance $\left(a^2 e^{-2a/\bar{x}}\right) \left(\frac{1}{n\bar{x}^2}\right)$
3

4 and the approx. SE is

5
$$\hat{SE} = \left(a e^{-a/\bar{x}}\right) \left(\frac{1}{\sqrt{n}\bar{x}}\right)$$

6

7
8 A $100(1-\alpha)\%$ CI for τ is

9
$$\hat{\tau} \pm z_{\alpha/2} \hat{SE}$$

10
11 Since
$$P\left(-z_{\alpha/2} \leq \frac{\hat{\tau} - \tau_0}{\hat{SE}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

12