

Asymptotic Normality of Maximum Likelihood Estimators and A Discussion on their Finite Normal Approximations

Naitong Chen

December 7, 2020

1 Body

First we set up the notations and assumptions that are used throughout the report.

1.1 Notations and Assumptions

Let X_1, \dots, X_n be i.i.d. continuous random variables in \mathbb{R} with probability density function $f(x; \theta_0)$, where $\theta_0 \in \Theta$ is a unknown parameter.

Write the likelihood function

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta), \quad (1)$$

and subsequently the log-likelihood function

$$l(\theta) = \sum_{i=1}^n \log f(x_i; \theta). \quad (2)$$

Then the maximum likelihood estimator can be computed as

$$\hat{\theta} = \arg \max_{\theta \in \Theta} l(\theta) = \arg \max_{\theta \in \Theta} L(\theta). \quad (3)$$

Note then since **conditons**, a necessary consequence is that

$$\left. \frac{\partial l(\theta)}{\partial \theta} \right|_{\theta=\theta_0} = 0. \quad (4)$$

We use $\mathbb{E}_\theta[X]$ to denote expectation with respect to X under $f(x; \theta)$. Using this notation we define the Fisher information

$$I(\theta) = \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f(X; \theta) \right)^2 \right]. \quad (5)$$

Throughout the report, we assume the following regularity conditions.

- (R1) $f(x; \theta)$ is identifiable: $\theta_1 \neq \theta_2 \implies f(x; \theta_1) \neq f(x; \theta_2)$.
- (R2) $f(x; \theta)$ has common support for all $\theta \in \Theta$.
- (R3) θ_0 is an interior point in Θ .
- (R4) $f(x; \theta)$ is twice differentiable in θ .
- (R5) The integral $\int f(x; \theta)$ can be different twice in θ under the integral sign.
- (R6) $\hat{\theta}$ is the unique solution to $\frac{\partial l(\theta)}{\partial \theta} = 0$.

1.2 Intermediate Results

We prove two lemmas that are essential to developing the main results of this report.

Lemma 1: Under regularity conditions, (**define** P_{θ_0}),

$$\lim_{n \rightarrow \infty} P_{\theta_0} (L(\theta_0 | X) > L(\theta | X)) = 1, \forall \theta \neq \theta_0. \quad (6)$$

Proof. We begin by taking the log on both sides of the inequality on the LHS and rearrange.

$$L(\theta_0 | X) > L(\theta | X) \implies \sum_{i=1}^n \log f(X_i; \theta_0) > \sum_{i=1}^n \log f(X_i; \theta) \quad (7)$$

$$\implies \sum_{i=1}^n (\log f(X_i; \theta_0) - \log f(X_i; \theta)) < 0 \quad (8)$$

$$\implies \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < 0. \quad (9)$$

Since X_i 's are i.i.d. , the summands are independent, and so by the Weak Law of Large Numbers,

$$\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) \xrightarrow{p} \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right]. \quad (10)$$

Note that $-\log(\cdot)$ is strictly convex, then by Jensen's inequality, we can establish, on the RHS of the above equation,

$$-\mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] = \mathbb{E}_{\theta_0} \left[-\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] > -\log \mathbb{E}_{\theta_0} \left[\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right]. \quad (11)$$

Now note

$$\log \mathbb{E}_{\theta_0} \left[\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right] = \log \int_{\mathbb{R}} \frac{f(x; \theta)}{f(x; \theta_0)} f(x; \theta_0) dx = \log \int_{\mathbb{R}} \frac{f(x; \theta)}{d} x = \log 1 = 0. \quad (12)$$

Together, we have

$$\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) \xrightarrow{p} \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] < 0. \quad (13)$$

To show the desired equation, it is equivalent to show, by Eq. (9),

$$\lim_{n \rightarrow \infty} P_{\theta_0} \left(\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < 0 \right) = 1. \quad (14)$$

By Eq. (13), we know that for all $\epsilon > 0$,

$$\lim_{n \rightarrow \infty} P_{\theta_0} \left(\left| \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) - \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] \right| < \epsilon \right) = 1. \quad (15)$$

Again by rearranging the inequality inside, we get

$$\mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] - \epsilon < \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] + \epsilon. \quad (16)$$

Note that the probability of event Eq. (16) is less than or equal to that of

$$\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] + \epsilon. \quad (17)$$

Since $\mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] < 0$, by fixing $\epsilon = -\mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] > 0$, we have

$$\lim_{n \rightarrow \infty} P_{\theta_0} \left(\left| \frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) - \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] \right| < \epsilon \right) \quad (18)$$

$$\leq \lim_{n \rightarrow \infty} P_{\theta_0} \left(\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < \mathbb{E}_{\theta_0} \left[\log \left(\frac{f(X_1; \theta)}{f(X_1; \theta_0)} \right) \right] + \epsilon \right) \quad (19)$$

$$= \lim_{n \rightarrow \infty} P_{\theta_0} \left(\frac{1}{n} \sum_{i=1}^n \log \left(\frac{f(X_i; \theta)}{f(X_i; \theta_0)} \right) < 0 \right) = 1. \quad (20)$$

Therefore, we conclude that

$$\lim_{n \rightarrow \infty} P_{\theta_0} (L(\theta_0 | X) > L(\theta | X)) = 1, \forall \theta \neq \theta_0. \quad (21)$$

□

Lemma 2: Under regularity conditions,

$$\mathbb{E} \left(\frac{\partial \log f(X; \theta)}{\partial \theta} \right) = 0, \text{ and } I(\theta) = -\mathbb{E}_{\theta} \left[\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right] = \text{Var} \left(\frac{\partial \log f(X; \theta)}{\partial \theta} \right). \quad (22)$$

Proof. We begin with differentiating both sides of $1 = \int_{\mathbb{R}} f(x; \theta) dx$.

$$0 = \frac{\partial}{\partial \theta} \int_{\mathbb{R}} f(x; \theta) dx \quad (23)$$

$$= \int_{\mathbb{R}} \frac{\partial}{\partial \theta} f(x; \theta) dx \quad (24)$$

$$= \int_{\mathbb{R}} \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx \quad (25)$$

$$= \int_{\mathbb{R}} \frac{\partial \log f(x; \theta)}{\partial \theta} f(x; \theta) dx \quad (26)$$

$$= \mathbb{E}_{\theta} \left[\frac{\partial \log f(x; \theta)}{\partial \theta} \right]. \quad (27)$$

Differentiating with respect to θ again, we get

$$0 = \frac{\partial}{\partial \theta} \int_{\mathbb{R}} \frac{\partial \log f(x; \theta)}{\partial \theta} f(x; \theta) dx \quad (28)$$

$$= \int_{\mathbb{R}} \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} f(x; \theta) dx + \int_{\mathbb{R}} \frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial f(x; \theta)}{\partial \theta} dx \quad (29)$$

$$= \int_{\mathbb{R}} \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} f(x; \theta) dx + \int_{\mathbb{R}} \frac{\partial \log f(x; \theta)}{\partial \theta} \frac{\partial \log f(x; \theta)}{\partial \theta} f(x; \theta) dx \quad (30)$$

$$= \mathbb{E}_{\theta} \left[\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right] + \mathbb{E}_{\theta} \left[\left(\frac{\partial \log f(x; \theta)}{\partial \theta} \right)^2 \right], \quad (31)$$

where Eq. (30) used the same trick as Eq. (26). By Eq. (5), we have

$$I(\theta) = \mathbb{E}_{\theta} \left[\left(\frac{\partial \log f(x; \theta)}{\partial \theta} \right)^2 \right] = -\mathbb{E}_{\theta} \left[\frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right]. \quad (32)$$

Finally, since $\mathbb{E}_\theta \left[\frac{\partial \log f(x; \theta)}{\partial \theta} \right] = 0$,

$$\text{Var} \left(\frac{\partial \log f(x; \theta)}{\partial \theta} \right) = \mathbb{E}_\theta \left[\left(\frac{\partial \log f(x; \theta)}{\partial \theta} \right)^2 \right] - \left(\mathbb{E}_\theta \left[\frac{\partial \log f(x; \theta)}{\partial \theta} \right] \right)^2 = I(\theta). \quad (33)$$

Together, we conclude that

$$\mathbb{E} \left(\frac{\partial \log f(X; \theta)}{\partial \theta} \right) = 0, \text{ and } I(\theta) = -\mathbb{E}_\theta \left[\frac{\partial^2 \log f(X; \theta)}{\partial \theta^2} \right] = \text{Var} \left(\frac{\partial \log f(X; \theta)}{\partial \theta} \right). \quad (34)$$

□

1.3 Asymptotic Normality of MLE

We first establish the consistency of MLE.

Lemma 3: Under regularity conditions, $\hat{\theta}$ is a consistent estimator of θ_0 .

Proof. We begin by showing that the equation $\frac{\partial l(\theta|x)}{\partial \theta} = 0$ has a solution $\hat{\theta}_n$ that converges in probability to θ_0 .

Since θ_0 is an interior point of Θ , we can find $a > 0$ such that $\theta_0 \in (\theta_0 - a, \theta_0 + a) \subset \Theta$. Then define the event

$$S_n = \{x : l(\theta_0 | x) > l(\theta_0 - a | x) \cap l(\theta_0 | x) > l(\theta_0 + a | x)\}. \quad (35)$$

Lemma 1 says, under P_{θ_0} , when n approaches infinity, θ_0 is the unique maximizer to $L(\theta | x)$, which implies that it is also the unique maximizer to $l(\theta | x)$. Then clearly $\lim_{n \rightarrow \infty} P_{\theta_0}(S_n) = 1$.

Note since $l(\theta_0 | x) > l(\theta_0 - a | x)$ and $l(\theta_0 | x) > l(\theta_0 + a | x)$, with $f(\theta; x)$ being continuous and differentiable, there must exist, for any $x \in S_n$, a local maximum in $(\theta_0 - a, \theta_0 + a)$. Denote this value $\hat{\theta}_n$, then $\frac{\partial l(\theta|x)}{\partial \theta} \Big|_{\theta=\hat{\theta}_n} = 0$.

Then for all $a > 0$ small enough, we can find a sequence of $\hat{\theta}$ such that

$$\lim_{n \rightarrow \infty} P_{\theta_0}(|\hat{\theta}_n - \theta_0| < a) = 1. \quad (36)$$

By choosing $\hat{\theta}_n$ to be the one closest to θ_0 , denoted θ_n^* , we have identified a sequence $(\theta_n^*)_{n \geq 1}$, independent of a , such that

$$\lim_{n \rightarrow \infty} P_{\theta_0}(|\hat{\theta}_n - \theta_0| < a) = 1, \forall a > 0. \quad (37)$$

This precisely means that $\hat{\theta}_n \xrightarrow{P} \theta_0$.

Now note that under regularity conditions, $\hat{\theta}_n$ is the unique solution to $\frac{\partial l(\theta|x)}{\partial \theta} = 0$. Therefore, we conclude that $\hat{\theta}$ is a consistent estimator of θ_0 . □

We now use this fact to prove the asymptotic normality of MLE.

Theorem 1: Under regularity conditions,

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \mathcal{N}\left(0, [I(\theta_0)]^{-1}\right), \quad (38)$$

where $I(\theta_0) = \mathbb{E}_{\theta_0} \left[\left(\frac{\partial \log f(X; \theta)}{\partial \theta} \Big|_{\theta=\theta_0} \right)^2 \right]$.

Proof. By the Mean Value Theorem, for $f : \mathbb{R} \rightarrow \mathbb{R}$ continuous on $[a, b]$ and differentiable on (a, b) , for all $c \in (a, b)$,

$$\frac{f(a) - f(b)}{a - b} = f'(c). \quad (39)$$

Let $f(\theta) = l'(\theta) = \frac{\partial l(\theta)}{\partial \theta}$, $a = \hat{\theta}$, $b = \theta_0$, $c = \theta_1 \in (\theta_0, \hat{\theta})$. Then with $l''(\theta) = \frac{\partial^2 l(\theta)}{\partial \theta^2}$, the above equation becomes

$$\frac{l'(\hat{\theta}) - l'(\theta_0)}{\hat{\theta} - \theta_0} = l''(\theta_1). \quad (40)$$

We know $l'(\hat{\theta}) = 0$, then the equation above becomes

$$0 = l'(\theta_0) + (\hat{\theta} - \theta_0)l''(\theta_1) \implies \sqrt{n}(\hat{\theta} - \theta_0) = -\frac{\sqrt{n}l'(\theta_0)}{l''(\theta_1)}. \quad (41)$$

We first look at the denominator of Eq. (41). By Lemma 3, $\hat{\theta} \xrightarrow{P} \theta_0$. Then since $\theta_1 \in (\theta_0, \hat{\theta})$, we must have $\theta_1 \xrightarrow{P} \theta_0$. Then by Proposition 10.7 from the lecture notes,

$$l''(\theta_1) \xrightarrow{P} l''(\theta_0). \quad (42)$$

Now by the Weak Law of Large Numbers and Lemma 2,

$$l''(\theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log f(x_i; \theta) \Big|_{\theta=\theta_0} \xrightarrow{P} \mathbb{E}_{\theta_0} \left[\frac{\partial^2}{\partial \theta^2} \log f(X_1; \theta) \Big|_{\theta=\theta_0} \right] = -I(\theta_0). \quad (43)$$

Now we look at the numerator of Eq. (41). By Lemma 2, $\mathbb{E}_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X_1; \theta) \Big|_{\theta=\theta_0} \right] = 0$ and $I(\theta_0) = \text{Var} \left(\frac{\partial}{\partial \theta} \log f(X_1; \theta) \Big|_{\theta=\theta_0} \right)$. Then by the Central Limit Theorem,

$$\sqrt{n}l'(\theta_0) = \sqrt{n} \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i; \theta) \Big|_{\theta=\theta_0} \quad (44)$$

$$= \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(x_i; \theta) \Big|_{\theta=\theta_0} - \mathbb{E}_{\theta_0} \left[\frac{\partial}{\partial \theta} \log f(X_1; \theta) \Big|_{\theta=\theta_0} \right] \right) \quad (45)$$

$$\xrightarrow{d} \mathcal{N} \left(0, \text{Var} \left(\frac{\partial}{\partial \theta} \log f(X_1; \theta) \Big|_{\theta=\theta_0} \right) \right) \quad (46)$$

$$= \mathcal{N}(0, I(\theta_0)). \quad (47)$$

Together by Eqs. (41), (43) and (47) and Slutsky's Theorem,

$$\sqrt{n}(\hat{\theta} - \theta_0) = -\frac{\sqrt{n}l'(\theta_0)}{l''(\theta_1)} \xrightarrow{d} \frac{\mathcal{N}(0, I(\theta_0))}{I(\theta_0)} = \mathcal{N}(0, [I(\theta_0)]^{-1}). \quad (48)$$

Therefore we conclude that $\sqrt{n}(\hat{\theta} - \theta_0) \xrightarrow{d} \mathcal{N}(0, [I(\theta_0)]^{-1})$. \square

Out of all possible estimators for θ_0 , how good is the MLE? The next result gives us some insight into the quality of MLE compared to others in terms of the variance of these estimators.

Theorem 2: Let $Y = u(X_1, \dots, X_n)$ be an unbiased estimator of θ_0 such that $\mathbb{E}_{\theta_0}[Y] = \theta_0$. Then under regularity conditions, $\text{Var}(Y) \geq \frac{1}{nI(\theta_0)}$.

Proof. We first expand $E_{\theta_0}[Y]$.

$$\theta_0 = E_{\theta_0}[Y] = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} u(x_1, \dots, x_n) \prod_{i=1}^n f(x_i; \theta_0) dx_1 \cdots dx_n. \quad (49)$$

Differentiating both sides with respect to θ_0 gives

$$1 = \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} u(x_1, \dots, x_n) \left(\sum_{i=1}^n \frac{\partial f(x_i; \theta_0)}{\partial \theta_0} \frac{1}{f(x_i; \theta_0)} \right) \prod_{i=1}^n f(x_i; \theta_0) dx_1 \cdots dx_n \quad (50)$$

$$= \int_{\mathbb{R}} \cdots \int_{\mathbb{R}} u(x_1, \dots, x_n) \left(\sum_{i=1}^n \frac{\partial \log f(x_i; \theta_0)}{\partial \theta_0} \right) \prod_{i=1}^n f(x_i; \theta_0) dx_1 \cdots dx_n \quad (51)$$

By writing $Z = \sum_{i=1}^n \frac{\partial \log f(X_i; \theta_0)}{\partial \theta_0}$, we get, with ρ denoting the correlation coefficient between Y and Z ,

$$1 = \mathbb{E}_{\theta_0}[YZ] = \mathbb{E}_{\theta_0}[Y]\mathbb{E}_{\theta_0}[Z] + \rho\sqrt{\text{Var}(Y)}\sqrt{\text{Var}(Z)} \implies \rho = \frac{1}{\sqrt{\text{Var}(Y)}\sqrt{\text{Var}(Z)}}. \quad (52)$$

We note that since X_i 's are i.i.d. ,

$$\text{Var}(Z) = \text{Var}\left(\sum_{i=1}^n \frac{\partial \log f(X_i; \theta_0)}{\partial \theta_0}\right) \quad (53)$$

$$= \sum_{i=1}^n \text{Var}\left(\frac{\partial \log f(X_i; \theta_0)}{\partial \theta_0}\right) \quad (54)$$

$$= n\text{Var}\left(\frac{\partial \log f(X_1; \theta_0)}{\partial \theta_0}\right) \quad (55)$$

$$= nI(\theta_0). \quad (56)$$

Then

$$\rho = \frac{1}{\sqrt{nI(\theta_0)}\sqrt{\text{Var}(Y)}}. \quad (57)$$

By definition, $\rho^2 \leq 1$, then

$$\rho^2 = \frac{1}{nI(\theta_0)\text{Var}(Y)} \leq 1 \implies \text{Var}(Y) \geq \frac{1}{nI(\theta_0)}. \quad (58)$$

Therefore for any unbiased estimator Y of θ_0 , we have $\text{Var}(Y) \geq \frac{1}{nI(\theta_0)}$. \square

Note the lower bound on the variance is exactly the asymptotic variance of $\hat{\theta}$. This means that asymptotically, the MLE achieves the smallest possible variance out of all unbiased estimators of θ . While this is a nice property, it remains rather theoretical. In practice, when we work with finitely many observations, it is not clear how far away we are from asymptotic normality. And existing heuristics from undergraduate statistics courses such as calling samples larger than 25 or 30 large enough is far from satisfying. Fortunately, [AL15] have done some pioneer work aimed at answering this exact question. We turn in the next section for a brief illustration of their result on the the closeness to the asymptotic normal distribution from the MLE approximated using finitely many observations.

2 Open questions and research directions

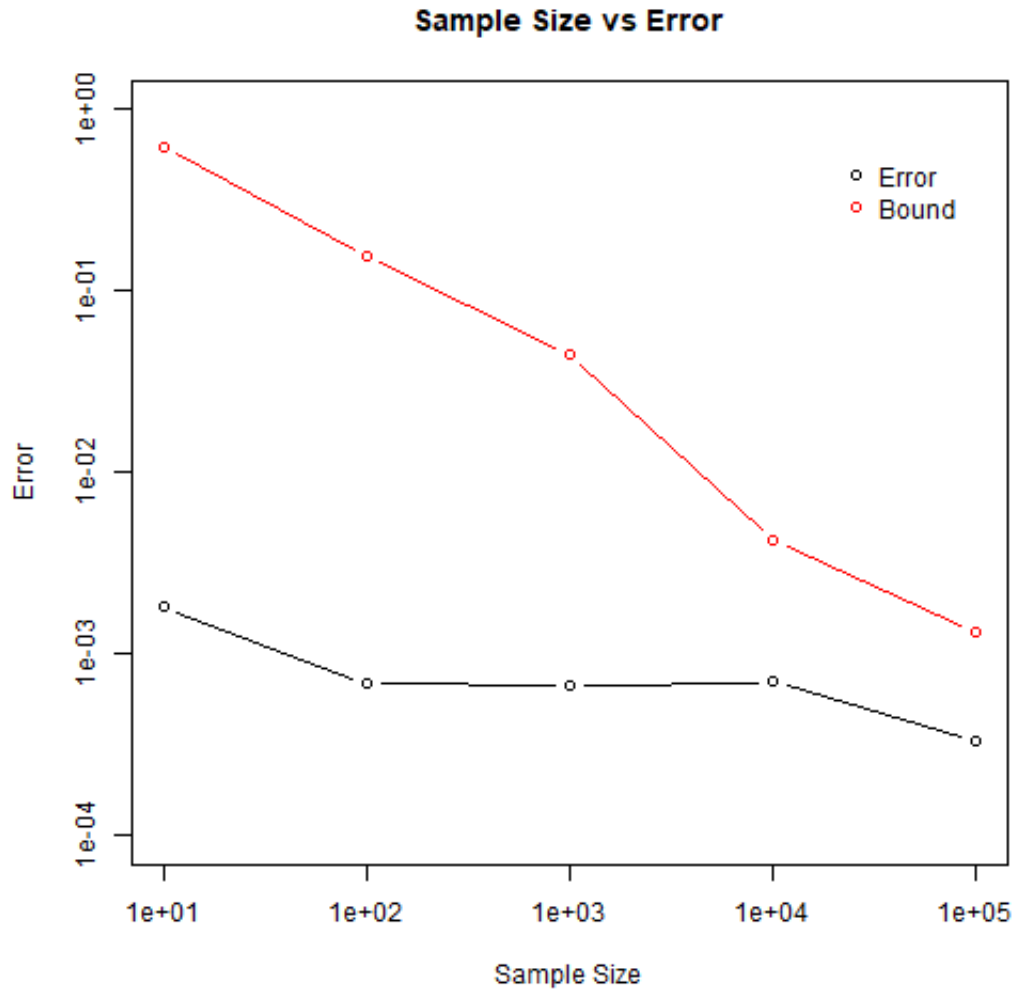


Figure 1: Simulation Gamma unknown scale

References

- [AL15] A. Anastasiou and C. Ley. “Bounds for the asymptotic normality of the maximum likelihood estimator using the Delta method”. In: *arXiv preprint arXiv:1508.04948* (2015).