

S&DS 242/542: Theory of Statistics

Lecture 14: Consistency and asymptotic normality of the MLE

Recap: Consistency and asymptotic normality

Given data $X_1, \dots, X_n \stackrel{iid}{\sim} f(x | \theta)$ from a parametric model, an estimator $\hat{\theta}$ for θ is *consistent* if

$$\hat{\theta} \rightarrow \theta$$

in probability as $n \rightarrow \infty$.

It is *asymptotically normal* if, furthermore, for some asymptotic variance $v(\theta)$,

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}(0, v(\theta))$$

in distribution as $n \rightarrow \infty$.

We showed last lecture why a method of moments estimator $\hat{\theta}$ for a parameter $\theta \in \mathbb{R}$ is usually asymptotically normal, and used the delta method to derive its asymptotic variance.

Consistency and asymptotic normality of the MLE

Theorem

Let $f(x | \theta)$ be a parametric model, with a single parameter $\theta \in \mathbb{R}$. Let $X_1, \dots, X_n \stackrel{iid}{\sim} f(x | \theta)$, and let $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ be the MLE. Under regularity conditions for $f(x | \theta)$, as $n \rightarrow \infty$,¹

(a) $\hat{\theta}$ is consistent.

(b) $\hat{\theta}$ is asymptotically normal, and $\sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}(0, \frac{1}{I(\theta)})$.

The function $I(\theta)$ in this asymptotic variance has the two equivalent forms

$$I(\theta) = \text{Var}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X | \theta) \right] = -\mathbb{E}_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X | \theta) \right]$$

where \mathbb{E}_{θ} and Var_{θ} denote expectation/variance over $X \sim f(x | \theta)$.

¹In this course, we won't discuss the exact conditions, which are technical. Three of the conditions are that θ is not on the boundary of the parameter space, $\theta \mapsto \log f(x | \theta)$ is twice differentiable, and $I(\theta)$ is non-zero.

Fisher Information

The function

$$I(\theta) = \text{Var}_\theta \left[\frac{\partial}{\partial \theta} \log f(X | \theta) \right] = -\mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f(X | \theta) \right]$$

is called the **Fisher information**. The quantity $\frac{\partial}{\partial \theta} \log f(X | \theta)$ is called the **score**.

The first expression for $I(\theta)$ states that the Fisher information is the variance of the score. We will see that the score has mean zero under the true parameter θ :

$$\mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \log f(X | \theta) \right] = 0$$

So this first expression may also be written as

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

Example of the Poisson model

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Poisson}(\lambda)$. The MLE is $\hat{\lambda} = \bar{X}$.

$$f(x|\lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \Rightarrow \log f(x|\lambda) = -\lambda + x \log \lambda - \log(x!)$$

$$\text{The score is: } \frac{\partial}{\partial \lambda} \log f(x|\lambda) = -1 + \frac{x}{\lambda}$$

$$E_{\lambda} \left[\frac{\partial}{\partial \lambda} \log f(x|\lambda) \right] = E_{\lambda} \left[-1 + \frac{X}{\lambda} \right] = -1 + \frac{E_{\lambda}[X]}{\lambda} = -1 + \frac{\lambda}{\lambda} = 0$$

$$I(\lambda) = \text{Var}_{\lambda} \left[\frac{\partial}{\partial \lambda} \log f(x|\lambda) \right] = \text{Var}_{\lambda} \left[-1 + \frac{X}{\lambda} \right] = \frac{1}{\lambda^2} \text{Var}_{\lambda}[X] = \frac{1}{\lambda}$$

$$\text{Alternatively: } \frac{\partial^2}{\partial \lambda^2} \log f(x|\lambda) = -\frac{x}{\lambda^2}$$

$$I(\lambda) = -E_{\lambda} \left[\frac{\partial^2}{\partial \lambda^2} \log f(x|\lambda) \right] = E_{\lambda} \left[\frac{X}{\lambda^2} \right] = \frac{1}{\lambda^2} E_{\lambda}[X] = \frac{1}{\lambda}$$

$$\text{So } \sqrt{n}(\hat{\lambda} - \lambda) \rightarrow \mathcal{N}(0, \frac{1}{I(\lambda)}) = \mathcal{N}(0, \lambda)$$

Example of the Pareto model

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Pareto}(\theta, 1)$. The MLE is $\hat{\theta} = \frac{n}{\sum_{i=1}^n \log X_i}$.

$$f(x|\theta) = \frac{\theta}{x^{\theta+1}} \quad (\text{for } x > 1) \Rightarrow \log f(x|\theta) = \log \theta - (\theta+1) \log x.$$

$$\frac{\partial}{\partial \theta} \log f(x|\theta) = \frac{1}{\theta} - \log x \quad \leftarrow \text{This is the score.}$$

$$\text{Recall: } \mathbb{E}_\theta[\log X] = \frac{1}{\theta}. \text{ So } \mathbb{E}_\theta\left[\frac{\partial}{\partial \theta} \log f(x|\theta)\right] = \frac{1}{\theta} - \frac{1}{\theta} = 0.$$

$$\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) = -\frac{1}{\theta^2}$$

$$\Rightarrow I(\theta) = -\mathbb{E}_\theta\left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta)\right] = \frac{1}{\theta^2}.$$

Example of the Pareto model

Alternatively: $I(\theta) = \text{Var}_\theta \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right] = \text{Var}_\theta [\log X]$.

$$\mathbb{E}_\theta [(\log X)^2] = \int_1^\infty (\log x)^2 \cdot \frac{\theta}{x^{\theta+1}} dx \quad (u = \log x, \quad x = e^u)$$

$$= \int_0^\infty u^2 \cdot \frac{\theta}{e^{u(\theta+1)}} e^u du$$

$$= \int_0^\infty u^2 \cdot \underbrace{\theta e^{-\theta u}}_{\text{PDF of Exponential}(\theta)} du = \dots = \frac{2}{\theta^2}$$

$$\Rightarrow \text{Var}_\theta [\log X] = \mathbb{E}_\theta [(\log X)^2] - (\mathbb{E}_\theta \log X)^2 = \frac{2}{\theta^2} - \frac{1}{\theta^2} = \frac{1}{\theta^2}$$

$$\Rightarrow I(\theta) = \frac{1}{\theta^2}$$

$$\text{So } \sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}(0, \frac{1}{I(\theta)}) = \mathcal{N}(0, \theta^2).$$

Comparison with method of moments

Recall from last lecture that for the method of moments estimator, we instead had (when $\theta > 2$)

$$\sqrt{n}(\hat{\theta} - \theta) \rightarrow \mathcal{N}\left(0, \frac{\theta(\theta - 1)^2}{\theta - 2}\right)$$

For any $\theta > 2$, it holds that

$$\frac{\theta(\theta - 1)^2}{\theta - 2} > \theta^2.$$

So when n is large, the method-of-moments estimator has larger standard error than the MLE. This is a general phenomenon, which we will discuss next lecture.

Plug-in estimator for the Pareto mean

Suppose we are interested in estimating the mean $\frac{\theta}{\theta-1}$ of this Pareto distribution, instead of θ . A plug-in estimator based on the MLE $\hat{\theta}$ would be $\frac{\hat{\theta}}{\hat{\theta}-1}$. To compute its asymptotic variance:

Apply delta method w/ $g(x) = \frac{x}{x-1} = 1 + \frac{1}{x-1}$, $g'(x) = -\frac{1}{(x-1)^2}$

$$\sqrt{n} \left(\underbrace{\frac{\hat{\theta}}{\hat{\theta}-1} - \frac{\theta}{\theta-1}}_{\text{True mean of Pareto}(\theta,1)} \right) \rightarrow \mathcal{N}(0, \theta^2 \cdot g'(\theta)^2)$$

$$\text{True mean of Pareto}(\theta,1) = \mathcal{N}\left(0, \frac{\theta^2}{(\theta-1)^4}\right)$$

Comparison to the sample mean

The plug-in estimate $\frac{\hat{\theta}}{\hat{\theta}-1}$ is not the only reasonable estimate for the mean of the Pareto distribution: What if we decided to simply use the sample mean \bar{X} ?

For this estimate \bar{X} , the CLT shows

$$\sqrt{n} \left(\bar{X} - \frac{\theta}{\theta-1} \right) \rightarrow \mathcal{N} \left(0, \frac{\theta}{(\theta-1)^2(\theta-2)} \right)$$

where $\frac{\theta}{(\theta-1)^2(\theta-2)}$ is the variance of the Pareto distribution (again assuming $\theta > 2$).

It may be checked that this variance is greater than the variance $\frac{\theta^2}{(\theta-1)^4}$ for the plug-in estimate using the MLE.

Comparison to the sample mean

When n is large, the plug-in estimate $\frac{\hat{\theta}}{\hat{\theta}-1}$ using the MLE $\hat{\theta}$ is more accurate than the sample mean \bar{X} .

In the Pareto model, one intuition is that the distribution is heavy-tailed, and the sample mean \bar{X} is heavily influenced by rare but large data values. In contrast, $\hat{\theta}$ estimate the shape of the Pareto distribution in a more robust way, and then estimates the mean from its relationship to the shape of the distribution.

A downside of this plug-in approach is that it is model-dependent: The estimate $\frac{\hat{\theta}}{\hat{\theta}-1}$ relies strongly on the correctness of the Pareto model, whereas \bar{X} would be a reasonable estimate of the mean of the data distribution even if the Pareto model doesn't hold true.

Proof sketch: consistency of the MLE

To explain why the MLE $\hat{\theta}$ is consistent, recall that $\hat{\theta}$ is the value of θ which maximizes

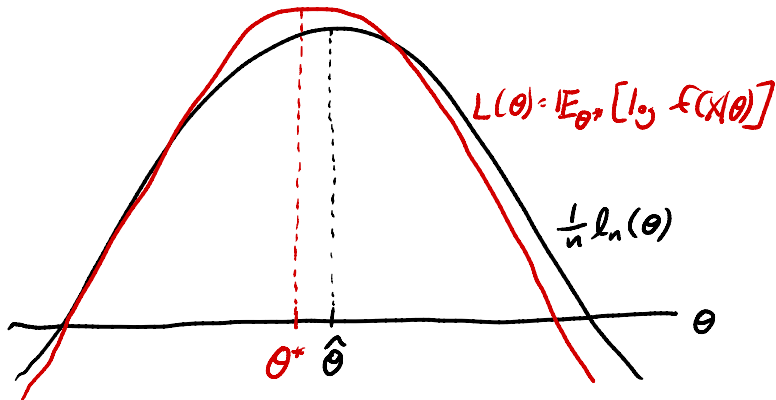
$$\frac{1}{n} \ell_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(X_i | \theta)$$

Suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f(x | \theta^*)$ with true parameter θ^* . Fixing *any* θ (not necessarily θ^*), the above is the sample average of n IID random variables, so the LLN implies

$$\frac{1}{n} \ell_n(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(X_i | \theta) \rightarrow \mathbb{E}_{\theta^*}[\log f(X | \theta)]$$

Here $\mathbb{E}_{\theta^*}[\log f(X | \theta)]$ is the expected value of each $\log f(X_i | \theta)$ when $X_i \sim f(x | \theta^*)$, where the log-likelihood is evaluated at an arbitrary parameter θ which may be different from θ^* .

Proof sketch: consistency of the MLE



Under suitable conditions, the value of θ maximizing $\frac{1}{n} \ell_n(\theta)$ (which is the MLE $\hat{\theta}$) converges in probability to the value of θ maximizing the limiting function $L(\theta) = \mathbb{E}_{\theta^*}[\log f(X | \theta)]$.

Proof sketch: consistency of the MLE

We claim that θ maximizing $L(\theta) = \mathbb{E}_{\theta^*}[\log f(X | \theta)]$ is exactly the true parameter θ^* :

$$\begin{aligned} L(\theta) - L(\theta^*) &= \mathbb{E}_{\theta^*} [\log f(X|\theta) - \log f(X|\theta^*)] \\ &= \mathbb{E}_{\theta^*} \left[\log \frac{f(X|\theta)}{f(X|\theta^*)} \right] \end{aligned}$$

The function $\log x$ is concave. So by Jensen's inequality

$$\begin{aligned} L(\theta) - L(\theta^*) &\leq \log \mathbb{E}_{\theta^*} \left[\frac{f(X|\theta)}{f(X|\theta^*)} \right] \\ &= \log \int \frac{f(X|\theta)}{f(X|\theta^*)} \cdot \cancel{f(X|\theta^*)} d\mu \end{aligned}$$

Proof sketch: consistency of the MLE

$$L(\theta) - L(\theta^*) \leq \log \int \underbrace{f(x|\theta)}_{\text{This is a PDF, so } \int f(x|\theta) dx = 1} dx = 0$$

So $L(\theta) - L(\theta^*) \leq 0$ for every θ , meaning that $L(\theta)$ is maximized at θ^* . This explains the consistency of $\hat{\theta}$.

Proof sketch: Definition of Fisher information

Next, let us check that the two definitions of $I(\theta)$ are the same, and that the score has mean zero:

$$\text{Applying again } \int f(x|\theta) dx = 1:$$

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int f(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} f(x|\theta) dx \end{aligned}$$

$$\text{Note that } \frac{\partial}{\partial \theta} \log f(x|\theta) = \frac{\frac{\partial}{\partial \theta} f(x|\theta)}{f(x|\theta)}$$

$$\begin{aligned} \Rightarrow 0 &= \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] \cdot f(x|\theta) dx \\ &= \mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right] \end{aligned}$$

Proof sketch: Definition of Fisher information

Differentiating again in θ :

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] \cdot f(x|\theta) dx \\ &= \int \frac{\partial}{\partial \theta} \left\{ \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] \cdot f(x|\theta) \right\} dx \\ &= \int \left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right] \cdot f(x|\theta) dx \\ &\quad + \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] \cdot \underbrace{\frac{\partial}{\partial \theta} f(x|\theta)}_{= \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] \cdot f(x|\theta)} dx \\ &= \mathbb{E}_\theta \left[\frac{\partial^2}{\partial \theta^2} \log f(X|\theta) \right] + \mathbb{E}_\theta \left[\left(\frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right] \end{aligned}$$

Proof sketch: Definition of Fisher information

$$\begin{aligned}\Rightarrow -\mathbb{E}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right] &= \mathbb{E}_{\theta} \left[\left(\frac{\partial}{\partial \theta} \log f(X|\theta) \right)^2 \right] \\ &= \text{Var}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right]\end{aligned}$$

1/c score has mean 0.

There are now two definitions of Fisher information, $I(\theta)$.

Proof sketch: Asymptotic normality

Finally, let us show the convergence in distribution

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

when the true parameter is θ .

The MLE $\hat{\theta}$ maximizes $l_n(\theta) = \sum_{i=1}^n \log f(X_i | \theta)$. So

$0 = l'_n(\hat{\theta})$. Since $\hat{\theta} - \theta \rightarrow 0$ in probability as $n \rightarrow \infty$,

for large n , $\hat{\theta}$ is close to θ . Taylor expand around θ :

$$0 = l'_n(\hat{\theta}) \approx l'_n(\theta) + l''_n(\theta) \cdot (\hat{\theta} - \theta)$$

$$\Rightarrow -l''_n(\theta) \cdot (\hat{\theta} - \theta) \approx l'_n(\theta)$$

$$\Rightarrow \sqrt{n}(\hat{\theta} - \theta) \approx -\frac{\sqrt{n}}{l''_n(\theta)} \cdot l'_n(\theta) = \frac{l'_n(\theta)/\sqrt{n}}{-l''_n(\theta)/n}$$

Proof sketch: Asymptotic normality

$$\sqrt{n}(\hat{\theta} - \theta) \approx \frac{l'_n(\theta)/\sqrt{n}}{-l''_n(\theta)/n}$$

$$\text{Recall } l_n(\theta) = \sum_{i=1}^n \log f(X_i | \theta)$$

For the denominator,

$$-\frac{l''_n(\theta)}{n} = \frac{1}{n} \sum_{i=1}^n \left(-\frac{\partial^2}{\partial \theta^2} \log f(X_i | \theta) \right)$$

$$\begin{aligned} &\rightarrow \mathbb{E}_{\theta} \left[-\frac{\partial^2}{\partial \theta^2} \log f(X | \theta) \right] \text{ by LLN} \\ &= I(\theta) \text{ (in probability)} \end{aligned}$$

Proof sketch: Asymptotic normality

$$\sqrt{n}(\hat{\theta} - \theta) \approx \frac{l'_n(\theta)/\sqrt{n}}{-l''_n(\theta)/n}$$

For the numerator:

$$\frac{l'_n(\theta)}{\sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \underbrace{\frac{\partial}{\partial \theta} \log f(X_i | \theta)}_{\text{score evaluated at } X_i}$$

$$\begin{aligned} &\rightarrow \mathcal{N}\left(0, \text{Var}_\theta \left[\frac{\partial}{\partial \theta} \log f(X|\theta) \right] \right) \text{ by CLT} \\ &\quad \text{(in distribution)} \\ &= \mathcal{N}(0, I(\theta)) \end{aligned}$$

By Slutsky's lemma: $\sqrt{n}(\hat{\theta} - \theta) \rightarrow \frac{1}{I(\theta)} \cdot \mathcal{N}(0, I(\theta)) = \mathcal{N}(0, \frac{1}{I(\theta)})$