

Lecture 6:

Recall that our setting was:

Data: $X_1, \dots, X_n \stackrel{iid}{\sim} f(x|\theta_0)$

Here, θ_0 denotes the true parameter that generates the data.

Input and Goal: We only have access to a random sample (X_1, \dots, X_n) , and we know that the data is generated according to a pdf $f(x|\theta)$. However, we don't know the true value of θ (i.e. we don't know θ_0) and our goal is to estimate θ_0 from the data (X_1, \dots, X_n) . Formally speaking

we'd like to design an estimator $\hat{\theta}(x_1, \dots, x_n)$ that gives us a good approximation for θ_0 .

In the previous lecture, we introduced the Maximum Likelihood Estimator (MLE): Given the sample data ($x_1 = x_1, x_2 = x_2, \dots, x_n = x_n$) we define the likelihood function

$$\text{lik}(\theta) = f(x_1, x_2, \dots, x_n | \theta)$$

and the MLE estimator, $\hat{\theta}_{\text{MLE}}$, is given

as :

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \text{lik}(\theta),$$

that is, $\hat{\theta}_{\text{MLE}}$ is the value of θ that maximizes the likelihood, or equivalently, makes the observed data sample "most probable" or "most likely".

Note that when X_i 's are iid, we have

$$\begin{aligned} f(x_1, \dots, x_n | \theta) &= f(x_1 | \theta) \times f(x_2 | \theta) \times \dots \times f(x_n | \theta) \\ &= \prod_{i=1}^n f(x_i | \theta) \end{aligned}$$

So $\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n f(x_i | \theta)$

Note also that since $\log(\cdot)$ is an increasing function, we can equivalently write:

$$\begin{aligned} \hat{\theta}_{\text{MLE}} &= \underset{\theta}{\operatorname{argmax}} \log \left(\prod_{i=1}^n f(x_i | \theta) \right) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log f(x_i | \theta) \end{aligned}$$

In summary, if we define

$\ell(\theta)$ is called
the log-likelihood
function.

$$\ell(\theta) = \sum_{i=1}^n \log f(x_i | \theta)$$

we have

$$\hat{\theta}_{\text{MLE}} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

Example : Assume X_1, \dots, X_n are generated iid according to exponential distribution with parameter λ_0 . Derive the mle estimator for λ_0 ?

So in this example, we have:

$$X_1, \dots, X_n \stackrel{iid}{\sim} f(x|\lambda) \\ = \lambda_0 \exp(-\lambda_0 x) \mathbb{1}_{\{x \geq 0\}}$$

We don't know the true value λ_0 , and we'd like to estimate it from data. We know, however, that data is generated according to $f(x|\lambda) = \lambda \exp(-\lambda x) \mathbb{1}_{\{x \geq 0\}}$

for some (unknown) λ .

The log-likelihood function is:

$$\begin{aligned} l(\lambda) &= \sum_{i=1}^n \log f(x_i|\lambda) \\ &= \sum_{i=1}^n \log(\lambda e^{-\lambda x_i}) = n \log \lambda - \lambda \sum_{i=1}^n x_i \end{aligned}$$

In order to find the maximizer of $\ell(\lambda)$
we need to solve the equation:

$$\ell'(\lambda) = 0$$

and the solution will be λ_{mle} .

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

$$\Rightarrow \ell'(\lambda) = 0 \rightarrow \lambda_{\text{mle}} = \frac{n}{\sum_{i=1}^n x_i} = \frac{1}{\bar{x}}$$

$$\text{where } \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

Note here that λ_{mle} is a **biased**
estimate of λ_0 because

$$E[\lambda_{\text{mle}}] = E\left[\frac{1}{\bar{x}}\right] \neq \frac{1}{\lambda_0}.$$

From the previous example we learned at least two facts:

1) To find the mle estimate

$$\hat{\theta}_{\text{mle}} = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

we need to solve the equation

$$\ell'(\theta) = 0,$$

and $\hat{\theta}_{\text{mle}}$ is a solution to this equation.

2) In general $\hat{\theta}_{\text{mle}}$ is not an unbiased estimator for the true parameter θ_0 .

The next natural question that we should ask is how good is the mle estimator? Why should we trust that the mle estimator would give us a good estimate? We'll find an answer to these questions in the next section.

Large Sample Theory For mle :

In this section we'd like find out how good the mle estimator becomes as we increase the number of data samples n .

The first thing we could expect is that $\hat{\theta}_{\text{mle}}$ will become close to the true value θ_0 as the number of samples n increases. Let's first verify this.

Lemma : $\lim_{n \rightarrow \infty} \hat{\theta}_{\text{mle}} = \theta_0$

The full proof of this lemma is beyond the scope of our class. Here, I'll provide the intuition. Recall that $\hat{\theta}_{\text{mle}}$ is the solution of the following:

$$l'(\hat{\theta}_{\text{mle}}) = \frac{d}{d\theta} l(\theta) \Big|_{\theta=\hat{\theta}_{\text{mle}}} = 0.$$

So if $\hat{\theta}_{\text{mle}}$ approaches θ_0 as n increases, we should expect that $l'(\theta_0)$ becomes very small (converges to zero).

as n increases. Instead of $\ell(\theta)$, let's look at a normalized version of it, $\frac{1}{n} \ell(\theta)$, and see how $\frac{1}{n} \ell(\theta_0)$ behaves.

$$\frac{1}{n} \ell(\theta) = \frac{1}{n} \sum_{i=1}^n \log f(x_i | \theta).$$

$$\Rightarrow \frac{1}{n} \ell'(\theta) = \frac{1}{n} \sum_{i=1}^n \frac{d}{d\theta} \log f(x_i | \theta)$$

$$\text{Let } y_i = \frac{\frac{d}{d\theta} f(x_i | \theta)}{f(x_i | \theta)}$$

then by the Law of Large Numbers (LLN)
we have:

$$\frac{1}{n} \sum_{i=1}^n y_i \rightarrow E[y]$$

$$n \rightarrow \infty$$

$$\sum_{i=1}^n \left(\frac{\frac{d}{d\theta} f(x_i | \theta)}{f(x_i | \theta)} \right) \rightarrow E \left[\frac{\frac{d}{d\theta} f(x | \theta)}{f(x | \theta)} \right]$$

Now, recall that data is generated according to $x_i \sim f(x_i | \theta_0)$.

We have :

$$\frac{1}{n} \ell'(\theta_0) \xrightarrow{n \rightarrow \infty} E \left[\frac{\frac{d}{d\theta} f(x|\theta)}{f(x|\theta_0)} \right]$$

$$= \int \frac{\frac{d}{d\theta} f(x|\theta_0)}{f(x|\theta_0)} f(x|\theta_0) dx$$

$$= \int \frac{d}{d\theta} f(x|\theta_0) dx$$

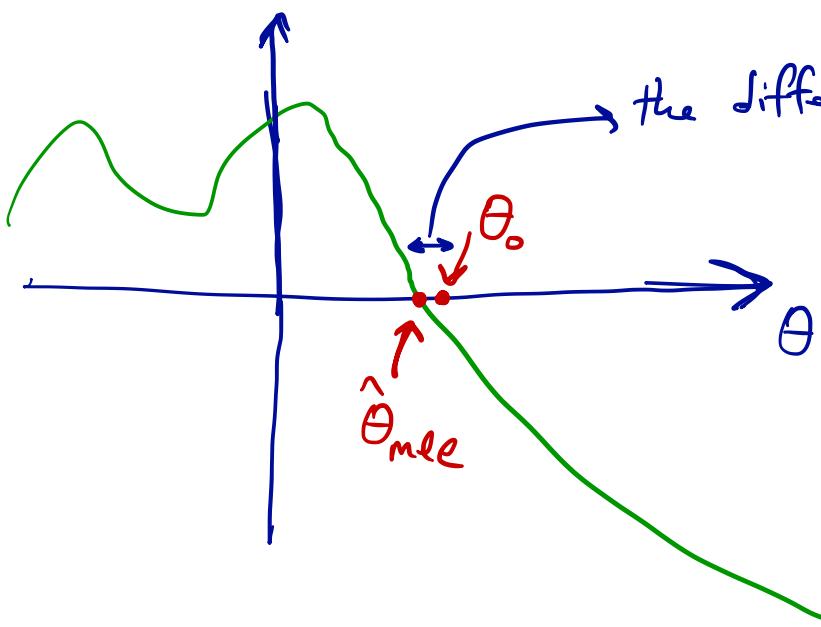
$$= \frac{d}{d\theta} \left(\int f(x|\theta_0) dx \right)$$

$$= \frac{d}{d\theta} (1) = 0$$

So, we have

$$\frac{1}{n} \ell'(\theta_0) \xrightarrow{n \rightarrow \infty} 0$$

$$\frac{1}{n} \ell'(\theta)$$



the difference of θ_0 and $\hat{\theta}_{mle}$ would
vanish by n .

So we now that θ_0 and $\hat{\theta}_{mle}$ will
be come closer and closer as n increases.

We now want to refine this picture
a little further and see how small
is the difference $\hat{\theta}_{mle} - \theta_0$ in terms
of n .

Theorem (Central-limit type relation for $\hat{\theta}_{MLE} - \theta_0$):

For n sufficiently large, we have:

$$\hat{\theta}_{MLE} - \theta_0 \sim \mathcal{N}\left(0, \frac{1}{\sqrt{n I(\theta)}}\right),$$

↗ normal distribution
 ↗ zero-mean
 ↗ variance is
 of order $1/n$

where $I(\theta)$, is the so-called Fisher Information, given as follows: Assuming $x \sim f(x|\theta_0)$:

$$I(\theta) = E \left[\left(\frac{d}{d\theta} (\ln f(x|\theta_0)) \right)^2 \right]$$

$$= \int \left(\frac{\frac{d}{d\theta} f(x|\theta_0)}{f(x|\theta_0)} \right)^2 f(x|\theta_0) dx.$$

proof: Again, the full proof is beyond the scope of this class and needs some extra smoothness assumptions on f . In the following, I will provide an intuitive sketch.

Recall that $\ell'(\hat{\theta}_{\text{MLE}}) = 0$. Using Taylor Series, we can write:

$$0 = \ell'(\hat{\theta}_{\text{MLE}}) \approx \ell'(\theta_0) + \ell''(\theta_0)(\hat{\theta}_{\text{MLE}} - \theta_0)$$

As a result, we have:

$$\begin{aligned} \hat{\theta}_{\text{MLE}} - \theta_0 &= - \frac{\ell'(\theta_0)}{\ell''(\theta_0)} \\ &= - \frac{\frac{1}{n} \ell'(\theta_0)}{\frac{1}{n} \ell''(\theta_0)} \end{aligned} \quad (*)$$

let's try to compute both terms $\frac{1}{n} \ell'(\theta_0)$ and $\frac{1}{n} \ell''(\theta_0)$. We have:

$$\begin{aligned} \frac{1}{n} \ell'(\theta_0) &= \frac{1}{n} \sum_{i=1}^n \frac{d}{d\theta} f'(x_i | \theta_0) \\ &\xrightarrow{\text{CLT}} \text{Normal}(\text{mean}, \sqrt{\frac{\text{var}}{n}}) \end{aligned}$$

where:

$$\text{mean} = E \left[\frac{d}{d\theta} \log f(x|\theta_0) \right]$$

$$\begin{aligned}
 &= \int \frac{\frac{d}{d\theta} f(x|\theta_0)}{f(x|\theta_0)} f(x|\theta_0) dx = \int \frac{d}{d\theta} f(x|\theta_0) dx \\
 &= \frac{d}{d\theta} \left(\int f(x|\theta_0) dx \right) \\
 &= \frac{d}{d\theta} (1) = 0
 \end{aligned}$$

$$\Rightarrow \text{mean} = 0$$

$$\Rightarrow \text{Var} = E \left[\left(\frac{d}{d\theta} (\log f(x|\theta)) \right)^2 \right] = I(\theta_0)$$

Hence,

$$\frac{1}{n} \ell'(\theta_0) \sim N(0, \sqrt{\frac{I(\theta_0)}{n}}) . \quad (**)$$

Let's compute $\frac{1}{n} \ell''(\theta_0)$. We have :

$$\frac{1}{n} \ell''(\theta_0) = \frac{1}{n} \sum_{i=1}^n \frac{\frac{\partial^2}{\partial \theta^2} \log f(x_i|\theta_0)}{\frac{\partial^2}{\partial \theta^2}}$$

law of large numbers

$$\longrightarrow E \left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta_0) \right]$$

we will prove later on that :

$$E \left[\frac{d^2 \ln f(x|\theta_0)}{d\theta^2} \right] = -I(\theta_0)$$

(to-be-proven)

Hence :

$$\frac{1}{n} \ell''(\theta_0) \rightarrow -I(\theta_0) \quad (***)$$

Let us now put together the three relations
(*) , (**) , and (***), to obtain:

$$\hat{\theta}_{mle} - \theta_0 \sim N(0, \frac{1}{I(\theta_0)} \sqrt{\frac{I(\theta_0)}{n}})$$

$$\sim N(0, \frac{1}{\sqrt{n I(\theta_0)}}).$$

And we get the result of the theorem.

Let us now show the relation (to-be-proven):

Recall that for every θ we have

$$\int f(x|\theta) dx = 1,$$

as $f(x|\theta)$ are probability density functions.

Let us define the function

$$g(\theta) = \int \frac{d}{d\theta} (\log f(x|\theta)) \cdot f(x|\theta) dx$$

note that

$$\begin{aligned} g(\theta) &= \int \frac{\frac{d}{d\theta} f(x|\theta)}{f(x|\theta)} \cdot f(x|\theta) dx \\ &= \int \frac{d}{d\theta} f(x|\theta) = \frac{d}{d\theta} \int f(x|\theta) dx \\ &= \frac{d}{d\theta} (1) = 0 \end{aligned}$$

Hence, $g(\theta) = 0$ for every θ . As a result, we can also say that,

$\frac{d}{d\theta} g(\theta) = 0$, which gives:

$$0 = \frac{d}{d\theta} \left[\int \left(\frac{d}{d\theta} \log f(x|\theta) \right) \cdot f(x|\theta) dx \right]$$

$$= \int \frac{d}{d\theta} \left[\left(\frac{d}{d\theta} \log f(x|\theta) \right) \cdot f(x|\theta) \right] dx$$

and hence, using the relation $(h_1 h_2)' = h_1' h_2 + h_2' h_1$,
we have:

$$0 = \int \left(\frac{d^2}{d\theta^2} \log f(x|\theta) \right) f(x|\theta) dx \\ + \int \left[\frac{d}{d\theta} \log f(x|\theta) \right]^2 f(x|\theta) dx$$

and hence :

$$\int \left(\frac{d^2}{d\theta^2} \log f(x|\theta) \right) f(x|\theta) dx = - \int \left(\frac{d}{d\theta} \log f(x|\theta) \right)^2 f(x|\theta) dx$$

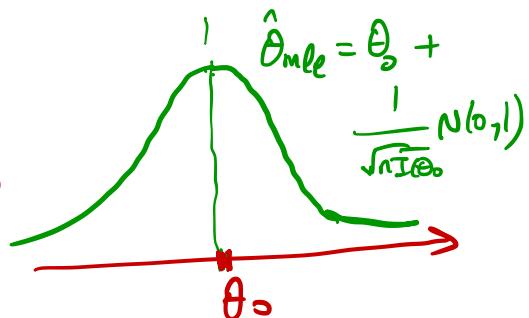
letting $\theta = \theta_0$ we get :

$$E \left[\frac{d^2}{d\theta^2} \log f(x|\theta_0) \right] = - I(\theta_0)$$

Confidence intervals for mle :

We found from the previous section that

$$\hat{\theta}_{\text{mle}} - \theta_0 \sim \frac{1}{\sqrt{n I(\theta_0)}} N(0, 1) \rightarrow$$



Let's now find β such that

$$\Pr \left\{ \hat{\theta}_{\text{mle}} - \beta \leq \theta_0 \leq \hat{\theta}_{\text{mle}} + \beta \right\} \approx 1 - \alpha$$

Let $\beta \triangleq \frac{\gamma}{\sqrt{n I(\theta_0)}}$, we have:

$$\Pr \left\{ \hat{\theta}_{\text{mle}} - \beta \leq \theta_0 \leq \hat{\theta}_{\text{mle}} + \beta \right\}$$

$$= \Pr \left\{ -\beta \leq \hat{\theta}_{\text{mle}} - \theta_0 \leq \beta \right\}$$

$$= \Pr \left\{ \frac{-\gamma}{\sqrt{n I(\theta_0)}} \leq \hat{\theta}_{\text{mle}} - \theta_0 \leq \frac{\gamma}{\sqrt{n I(\theta_0)}} \right\}$$

$$= \Pr \left\{ -\frac{\gamma}{\sqrt{n I(\theta_0)}} \leq \frac{1}{\sqrt{n I(\theta_0)}} N(0, 1) \leq \frac{\gamma}{\sqrt{n I(\theta_0)}} \right\}$$

$$= \Pr \{ -\gamma \leq N(0,1) \leq \gamma \}$$

we know that by letting $\gamma = z_{\alpha/2}$

we have that

$$\Pr \{ -z_{\alpha/2} \leq N(0,1) \leq +z_{\alpha/2} \} = 1 - \alpha$$

$$\Rightarrow \beta = z_{\alpha/2} \cdot \frac{1}{\sqrt{n I(\theta_0)}}$$

or

$$\Pr \left\{ \hat{\theta}_0 - \frac{z_{\alpha/2}}{\sqrt{n I(\theta_0)}} \leq \theta_0 \leq \hat{\theta}_0 + \frac{z_{\alpha/2}}{\sqrt{n I(\theta_0)}} \right\} \approx 1 - \alpha$$

Now, in reality, we don't know what θ_0 is (in fact we are estimating θ_0 !). For finding confidence intervals when we only have access to data (X_1, \dots, X_n) we can use the estimate

$\hat{\theta}_{\text{mle}}$ instead of θ_0 (in the formula for the confidence interval). That is,

$$\Pr \left\{ \hat{\theta}_{\text{mle}} - \frac{z_{\alpha/2}}{\sqrt{n I(\hat{\theta}_{\text{mle}})}} \leq \theta_0 \leq \hat{\theta}_{\text{mle}} + \frac{z_{\alpha/2}}{\sqrt{n I(\hat{\theta}_{\text{mle}})}} \right\} \approx 1 - \alpha$$

Cramer-Rao Bound :

Recall that for our estimators, we typically prefer the one which have smaller variance. In this regard, a fundamental question is what is the minimum possible variance in terms of the number of samples n . Given our setting $x_1, x_2, \dots, x_n \sim f(x|\theta_0)$, the Cramer-Rao bound

tells us that for any unbiased estimator of θ_0 , call it $T(x_1, \dots, x_n)$, the variance is bound by

$$\text{Var}(T) \geq \frac{1}{n I(\theta_0)}.$$

Recall that the mle estimator has the distribution (for n large) :

$$\hat{\theta}_{\text{mle}} - \theta_0 \sim \frac{1}{\sqrt{n I(\theta_0)}} N(0, 1)$$

So, when n is large, we have $\text{Var}(\hat{\theta}_{\text{mle}}) = \frac{1}{n I(\theta_0)}$

which matches the Cramer-Rao lower bound.

In this sense, $\hat{\theta}_{\text{mle}}$ has the minimum possible variance when n is large.

(An intuitive sketch of) the proof of Cramer-Rao:

Define the random variable

$$L = \sum_{i=1}^n \frac{d}{d\theta} (\log f(x_i | \theta_0))$$

We saw before that :

$$\begin{aligned} E[L] &= n E\left[\frac{d}{d\theta} \log f(x | \theta_0) \right] \\ &= n E\left[\frac{\frac{d}{d\theta} f(x | \theta_0)}{f(x | \theta_0)} \right] = 0 \end{aligned}$$

Also,

$$\begin{aligned} \text{Var}(L) &= \text{Var}\left(\frac{d}{d\theta} (\log f(x, \theta_0))\right) + \dots + \text{Var}\left(\frac{d}{d\theta} (\log f(x_n, \theta_0))\right) \\ &= n \text{Var}\left(\frac{d}{d\theta} (\log f(x, \theta_0))\right) \\ &= n E\left[\left(\frac{\frac{d}{d\theta} f(x, \theta_0)}{f(x, \theta_0)}\right)^2\right] \\ &= n I(\theta_0) \end{aligned}$$

(note that if y_1, \dots, y_n are iid, then)

$$\text{Var}(y_1 + y_2 + \dots + y_n) = \text{Var}(y_1) + \dots + \text{Var}(y_n) = n \text{Var}(y_1)$$

So, to summarize:

$$\begin{cases} E[L] = 0 \\ \text{Var}(L) = n I(\theta_0) \end{cases}$$

Now, take any unbiased estimator of θ_0 ; i.e.

$$E[T] = \theta_0, \text{ when } x_i \stackrel{iid}{\sim} f(x | \theta_0) \text{ for any value of } \theta_0.$$

We will first show that $E[L \cdot T] = 1$.

We have:

$$E[L \cdot T]$$

$$= \int_{(x_1, \dots, x_n)} T(x_1, \dots, x_n) \cdot L(x_1, \dots, x_n) f(x_1, \dots, x_n | \theta_0) dx_1 dx_2 \dots dx_n$$

$$= \int_{(x_1, \dots, x_n)} T(x_1, \dots, x_n) L(x_1, \dots, x_n) \prod_{i=1}^n f(x_i | \theta_0) dx_i$$

$$= \int_{(x_1, \dots, x_n)} T(x_1, \dots, x_n) \left(\sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i | \theta_0)}{f(x_i | \theta_0)} \right) \prod_{i=1}^n f(x_i | \theta_0) dx_i$$

$$= \int_{(x_1, \dots, x_n)} T(x_1, \dots, x_n) \frac{d}{d\theta} \left(\prod_{i=1}^n f(x_i | \theta_0) \right) dx_1 \dots dx_n$$

$$= \frac{d}{d\theta} \left(\underbrace{\int_{(x_1, \dots, x_n)} T(x_1, \dots, x_n)}_{\text{constant}} \cdot \prod_{i=1}^n f(x_i | \theta_0) dx_1 \dots dx_n \right)$$

$$= \left. \frac{d}{d\theta} (\theta) \right|_{\theta=\theta_0} = 1$$

(note that)

$$\left(\frac{d}{d\theta} \left(\prod_{i=1}^n f(x_i | \theta) \right) \right) = \sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i | \theta)}{f(x_i | \theta)} \prod_{j=1}^n f(x_j | \theta)$$

Now that we have $E[L \cdot T] = 1$, we can use the Cauchy-Schwarz inequality to deduce that:

$$(E[L \cdot T])^2 \leq E[L^2] \cdot E[T^2]$$

$$\underbrace{1}_{\geq} \leq \underbrace{n I(\theta_0)}_{\downarrow} \cdot \underbrace{\text{var}(T)}_{\downarrow}$$

$$\Rightarrow \text{var}(T) \geq \frac{1}{n I(\theta_0)}.$$