

# Max-Likelihood estimator Lecture Notes

Gleb Pantileev

March 2025

## Contents

<b>1</b>	<b>Preface</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
<b>3</b>	<b>Basic of estimators</b>	<b>3</b>
3.1	Definitions . . . . .	3
3.2	Basic Properties . . . . .	3
3.2.1	Bias . . . . .	3
3.2.2	Efficiency . . . . .	3
3.2.3	Consistency . . . . .	3
3.2.4	Sufficient conditions fo consistency . . . . .	4
<b>4</b>	<b>Maximum-Likelihood Estimator</b>	<b>5</b>
4.1	Way to obtain . . . . .	5
4.2	Properties of MLE . . . . .	6
4.2.1	Consistency . . . . .	6
4.2.2	Invariance . . . . .	8
4.2.3	Asymptotically Normal . . . . .	9
4.2.4	Asymptotically efficient . . . . .	11
<b>5</b>	<b>Appendix</b>	<b>12</b>

# 1 Preface

The origin of this article was quite usual. During my regular studies, I have faced a problem with one of the properties of MLE. I asked several people, but no one could give me an answer, so together with my teacher, Yaroslav Alexandrovich, we decided to make this project.

## 2 Introduction

In this work I want to understand why Max-Likelihood estimators are widely used and prove their properties. In order to do so I will use (and firstly prove) Cramer-Rao lower bound, Kullback-Leibler distance and Cauchy-Swartz inequality (without proof).

## 3 Basic of estimators

### 3.1 Definitions

First of all, it is important to understand what is meant by population parameters, estimators, etc.

Population parameters- parameters which define the distribution of the random variable:  $X \sim f(x, \vec{\theta})$ . For Normal distribution, these parameters are  $\mu, \sigma^2$ , for Binomial:  $n$  and  $p$ .

Statistic- any function of the sample  $(X_1, \dots, X_n)$ . For example,  $\bar{X}, s^2, X_i, 0, g(X_1, \dots, X_n)$

Estimator- statistic used as an approximation for the population parameter  $\theta$ , based on  $(X_1, \dots, X_n)$ - random sample. Whereas estimate- estimator based on sample obtained  $(x_1, \dots, x_n)$ .

### 3.2 Basic Properties

#### 3.2.1 Bias

Estimator is called unbiased if  $E(\hat{\theta}) = \theta$ . If not, then bias of the statistics can be found as  $Bias(\hat{\theta}) = E(\hat{\theta}) - \theta$ .

#### 3.2.2 Efficiency

Efficiency of the estimator is connected with its Mean Squared Error (MSE).

$$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = Var(\hat{\theta}) + Bias^2(\hat{\theta})$$

Consider two estimators:  $\hat{\theta}$  and  $\tilde{\theta}$ . Then, if  $MSE(\hat{\theta}) < MSE(\tilde{\theta})$ ,  $\hat{\theta}$  is more efficient estimator of  $\theta$ . According to Cramer-Rao inequality theorem:

$$Var(\hat{\theta}_n) \geq \frac{(1 + bias(\hat{\theta}))^2}{I_n(\theta)}, \text{ where } I_n(\theta) = Var\left(\frac{dL(\vec{X}, \theta)}{d\theta}\right)$$

But about this later... (proof in the appendix (3))

#### 3.2.3 Consistency

Consistency is an asymptotic property which means that estimator tends to the value of the true parameter in probability as  $n$  (size of the sample) tends to infinity.

$\hat{\theta}$  is consistent estimator if  $\theta$  if:

$$P(|\hat{\theta} - \theta| > \epsilon) \rightarrow 0 \quad \forall \epsilon > 0 \iff \hat{\theta} \xrightarrow{P} \theta$$

### 3.2.4 Sufficient conditions fo consistency

Theorem states that if  $MSE(\hat{\theta}) \rightarrow 0$  as  $n \rightarrow \infty$  then  $\hat{\theta}$  is a consistent estimator of  $\theta$ .

$$MSE(\hat{\theta}) \rightarrow 0 \text{ as } n \rightarrow \infty \iff \begin{cases} E(\hat{\theta}) \rightarrow \theta \\ Var(\hat{\theta}) \rightarrow 0 \end{cases}$$

In order to prove this property, Chebychev's inequality should be used :  $P(|\hat{\theta} - \theta| \geq \epsilon) \leq \frac{MSE(\hat{\theta})}{\epsilon^2}$  (proof is in the appendix(1)). So as  $MSE(\hat{\theta}) \rightarrow 0$ ,  $P(|\hat{\theta} - \theta| \geq \epsilon) \rightarrow 0$  too.

## 4 Maximum-Likelihood Estimator

### 4.1 Way to obtain

The method of finding the Max-Likelihood estimator is next:

Knowing the distribution of  $X_i$ :  $X \sim f(x, \vec{\theta})$ , with unknown  $\vec{\theta}$  and getting the sample of  $n$  observations  $(X_1, \dots, X_n)$  find such  $\hat{\theta}_{ml}$  which maximises the probability of such sample.

Consider Likelihood function  $L(\theta)$ :

$$L(\vec{\theta}) = \log f(\vec{x}, \vec{\theta})$$

However, life is not such simple and there could be two possible cases:

- $\hat{\theta}_{ml}$  can be obtained by equating the derivative of Likelihood Function with respect to parameter to 0.
- But if it impossible to do,  $\hat{\theta}_{ml}$  should be obtained analytically.

Consider two examples:

First:

$$\begin{aligned} X_i \sim N(\mu, \sigma^2) \text{ } \sigma \text{ is known} &\Rightarrow f(\vec{X}, \mu) = \prod_{i=1}^n f(X_i, \mu) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{\frac{-(X_i - \mu)^2}{2\sigma^2}} \\ L(\mu) = \log(f(\vec{X}, \mu)) &= \frac{-n}{2} \cdot \log(2\pi\sigma^2) - \frac{\sum_{i=1}^n (X_i - \mu)^2}{2\sigma^2} \\ \frac{dL}{d\mu} = \frac{\sum_{i=1}^n X_i - n \cdot \hat{\mu}}{\sigma^2} = 0 &\Rightarrow \hat{\mu}_{ml} = \frac{\sum_{i=1}^n X_i}{n} = \bar{X} \end{aligned}$$

Second:

$$\begin{aligned} Y_i \sim U[0, \theta] \text{ need to estimate } \theta & f(y, \theta) = \begin{cases} \frac{y}{\theta}, & \text{if } 0 \leq y \leq \theta \\ 0, & \text{otherwise} \end{cases} \\ f(\vec{Y}, \theta) = \prod_{i=1}^n f(Y_i, \theta) &= \prod_{i=1}^n \frac{Y_i}{\theta} \\ L(\theta) = \sum_{i=1}^n \log(Y_i) - n \log \theta & \\ \frac{dL}{d\theta} = -\frac{n}{\theta} & \end{aligned}$$

This case is more inconvenient since  $\theta$  is non-negative and thus derivative of likelihood function is negative ( $L(\theta)$  constantly decreases with respect to  $\theta$ ), thus, in order to be maximized,  $\theta$  should be minimized.

$$\hat{\theta}_{ml} \rightarrow \min_{s.t. \theta \geq X_i} \Rightarrow \hat{\theta}_{ml} = X_{max}$$

## 4.2 Properties of MLE

### 4.2.1 Consistency

When it is said that estimator is consistent, it means that as size of the sample is large enough, it converges to the value of true parameter in distribution.

$$\hat{\theta}_n \xrightarrow{P} \theta$$

In order to prove this property I would use Kullback-Leibler distance between f and g, where f,g-pdfs:

$$\begin{aligned} D(f, g) &= \int f(x) \cdot \log\left(\frac{f(x)}{g(x)}\right) dx \\ D(f, g) &\geq 0 \text{ (proof in the appendix(2))} \\ D(\theta, \phi) &\equiv D(f(x, \theta), f(x, \phi)) \quad \forall \theta, \phi \in \Theta \end{aligned}$$

Let  $\theta^*$ - true value of  $\theta$

Consider new function  $M(\theta)$ :

$$\begin{aligned} M_n(\theta) &= \frac{1}{n} \cdot \sum_{i=1}^n \log \frac{f(X_i, \theta)}{f(X_i, \theta^*)} = \frac{1}{n} (L_n(\theta) - L_n(\theta^*)) \\ \text{thus if } M_n(\theta) &\rightarrow \max \Leftrightarrow L_n(\theta) \rightarrow \max \end{aligned}$$

By Law of Large Numbers:

$$\begin{aligned} M_n(\theta) &\rightarrow M(\theta) = E_{\theta^*} \left( \log \frac{f(X_i, \theta)}{f(X_i, \theta^*)} \right) = \\ \int \log \frac{f(x, \theta)}{f(x, \theta^*)} f(x, \theta^*) dx &= - \int \log \frac{f(x, \theta^*)}{f(x, \theta)} f(x, \theta^*) dx = -D(\theta^*, \theta) \\ M(\theta) &= -D(\theta^*, \theta) \leq 0 \Rightarrow M(\theta) \text{ max at } \theta = \theta^* \\ M_n(\theta) \text{ max at } \theta &= \hat{\theta}_{ml}, \quad M_n(\theta) \rightarrow M(\theta), \text{ thus } \hat{\theta}_{ml} \rightarrow \theta^* \end{aligned}$$

But this proof is dirty, now consider more formal proof! Stay hard) Uniform convergence should be proved (for all values of  $\theta \in \Theta$ ) estimator converges to the value of true parameter.

$$\begin{aligned} M(\theta) &= -D(\theta^*, \theta) \leq 0 \\ M(\theta^*) &= -D(\theta^*, \theta^*) = 0 \end{aligned}$$

$$M(\theta) = E_{\theta^*}(M_n) = E_{\theta^*}(M_i)$$

$$\begin{aligned} M_n(\theta) &= \frac{\sum_{i=1}^n M_i}{n} = \frac{\sum_{i=1}^n \log \frac{f(X_i, \theta)}{f(X_i, \theta^*)}}{n} = \frac{\log f(\vec{x}, \theta) - \log f(\vec{x}, \theta^*)}{n} \\ M_n(\theta^*) &= 0 \end{aligned}$$

Assumptions made:

1)  $M_n(\theta)$  converges to  $M(\theta) \forall \theta \in \Theta$ :

$$\sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0$$

2) For all  $\theta \neq \theta^*$   $M(\theta)$  is less than the value of  $M(\theta^*)$ :

$$\sup_{\theta: |\theta - \theta^*| \geq \epsilon} M(\theta) < M(\theta^*) \quad \forall \epsilon > 0$$

3) For the finite samples  $\hat{\theta}_{ml}$  determine the distribution better (due to a "bad sample" with outliers or due to a small size).

$$M_n(\hat{\theta}_{ml}) \geq M_n(\theta^*) \quad (\text{for finite samples only})$$

Proof:

$$\begin{aligned} M(\theta^*) - M(\hat{\theta}_{ml}) &= M_n(\theta^*) - M(\hat{\theta}_{ml}) + M(\theta^*) - M_n(\theta^*) \leq M_n(\hat{\theta}_{ml}) - M(\hat{\theta}_{ml}) + M(\theta^*) - M_n(\theta^*) = \\ &= M_n(\hat{\theta}_{ml}) - M(\hat{\theta}_{ml}) \leq \sup_{\theta \in \Theta} |M_n(\theta) - M(\theta)| \xrightarrow{P} 0 \end{aligned}$$

$$\text{Thus, } M(\theta^*) - M(\hat{\theta}_{ml}) \xrightarrow{P} 0 \Leftrightarrow P(|M(\theta^*) - M(\hat{\theta}_{ml})| > \delta) \rightarrow 0 \quad \forall \delta > 0$$

Now consider two events:

$$\begin{aligned} A &= \left\{ |\hat{\theta}_{ml} - \theta^*| > \epsilon \right\} \quad \forall \epsilon > 0 \\ B &= \left\{ M(\theta^*) - M(\hat{\theta}_{ml}) > 0 \right\} \end{aligned}$$

If event A happens, then  $\hat{\theta}_{ml} \neq \theta^*$  so under the second assumption it is true that  $M(\hat{\theta}_{ml}) < M(\theta^*)$   
(B happens too)  $\Rightarrow$  A- subset of B  $\Rightarrow P(A) \leq P(B)$

$$P(|\hat{\theta}_{ml} - \theta^*| > \epsilon) \leq P(M(\theta^*) - M(\hat{\theta}_{ml}) > 0) \rightarrow 0$$

Thus  $P(|\hat{\theta}_{ml} - \theta^*| > \epsilon) \rightarrow 0$  - proved.

#### 4.2.2 Invariance

The invariance property means that if  $\hat{\theta}_{ml}$ - max-likelihood estimator of  $\theta$ , then  $g(\hat{\theta}_{ml}) = \hat{g}_{ml}$ - max-likelihood estimator of  $g(\theta)$ .

Proof:

- Let  $L(\hat{\theta})$ - is maximum of  $L(\theta)$ :  $\frac{dL}{d\theta}(\hat{\theta}) = 0$
- Let  $\theta = g(\eta)$ , so in order to get the exact value of  $\eta$ , function  $g$  should be invertible:  $\eta = g^{-1}(\theta)$ .

$$\text{Let } \frac{dL}{d\eta} = \frac{\partial L}{\partial g(\eta)} \cdot \frac{dg(\eta)}{d\eta} = 0 \Rightarrow$$

$$\begin{cases} \frac{dL}{dg(\eta)} = 0 \\ \frac{dg(\eta)}{d\eta} = 0 \end{cases}$$

Consider  $g(x)$  such that its derivative does not equal to 0 at any point (so there would be no lie points). Then  $\frac{dL}{dg(\eta)} = 0$  at some certain point  $\hat{\eta}$ . So  $\hat{\eta}$  is a maximiser of likelihood function and there is only one such point, then  $\hat{\theta} = g(\hat{\eta})$ . Taking the inverse function ( $g^{-1}(x)$ ) of both sides, we obtain  $\hat{\eta} = g^{-1}(\hat{\theta})$ . Overall, limitations of the function:  $g(\epsilon)$  should be invertible and strictly monotonic (its derivative should not be equal to 0 at any point, otherwise it would be necessary to check the obtained point analytically).

It is a very useful property because usually it is necessary to estimate probability of event, which is a function of  $\theta$  also. For example, Poisson distribution:  
 $X \sim Pois(\lambda)$ ,  $\hat{\lambda}_{ml} = \bar{X}$ , and  $P(X = 5) = \frac{\lambda^5}{5!} e^{-\lambda}$ , so  $\hat{\lambda}_{ml} = \bar{X}$  can be used to estimate this probability:

$$\hat{p} = \frac{\hat{\lambda}_{ml}^5}{5!} e^{-\hat{\lambda}_{ml}} = \frac{\bar{X}^5}{5!} e^{-\bar{X}}$$



### 4.2.3 Asymptotically Normal

Refresh previous results:

$$S_n(\theta) \equiv S(\vec{X}, \theta)$$

$$E(S_n) = 0$$

$$I_n(\theta) \equiv \text{Var}(S_n) = E(-S'_n) \text{ (see appendix(2))}$$

Now let's approximate the  $S_n(\theta)$  function by Taylor series:

$$S_n(\theta) = S_n(\theta_0) + S'_n(\theta_0)(\theta - \theta_0) + o(\theta - \theta_0)$$

Let  $\theta = \hat{\theta}_n$  (approximate the value of  $S_n(\theta)$  at point  $\hat{\theta}_n$ )

Let  $\theta_0 = \theta$ - true value of parameter.

$$S_n(\hat{\theta}_n) = S_n(\theta) + S'_n(\theta)(\hat{\theta}_n - \theta) + o(\hat{\theta}_n - \theta)$$

Now consider the case in which  $\hat{\theta}_{ml}$  is obtained by nullifying the  $S_n(\theta)$ . In this case  $S_n(\hat{\theta}_{ml})=0$

$$0 = S_n(\theta) + S'_n(\theta)(\hat{\theta}_n - \theta) + o(\hat{\theta}_n - \theta)$$

Due to a consistency of  $\hat{\theta}_n$  the first order precision is enough since  $|\hat{\theta}_n - \theta| \xrightarrow{P} 0$

$$\begin{aligned} \hat{\theta}_n - \theta &\approx \frac{S_n(\theta)}{-S'_n(\theta)} \\ \sqrt{n}(\hat{\theta}_n - \theta) &\approx \frac{\frac{1}{\sqrt{n}}S_n(\theta)}{-\frac{1}{n}S'_n(\theta)} \end{aligned}$$

$\frac{1}{\sqrt{n}}S_n(\theta) = \frac{1}{\sqrt{n}} \sum_{i=1}^n S_i(\theta) = \sqrt{n} \cdot \bar{S}$ , then by CLT:

$$\frac{\bar{S} - E(\bar{S})}{\text{Var}(\bar{S})} \sim N(0, 1)$$

$$E(\bar{S}) = E(S_i) = 0$$

$$\text{Var}(\bar{S}) = \frac{1}{n^2} \text{Var}(S_n) = \frac{1}{n^2} \cdot I_n(\theta) = \frac{1}{n} I(\theta),$$

Thus,

$$\bar{S} \sim N(0, \frac{1}{n} I(\theta)) \Rightarrow \sqrt{n} \cdot \bar{S} \sim N(0, I(\theta))$$

$-\frac{1}{n}S'_n(\theta) = -\frac{1}{n}\sum_{i=1}^n S'_i(\theta) = -\overline{S'}$  by CLT is distributed

$$-\overline{S'} \sim N(E(-\frac{1}{n}S'_n(\theta)), Var(-\frac{1}{n}S'_n(\theta))) = N(E(-\overline{S'}), Var(-\overline{S'}))$$

$$E(-\overline{S'}) = E(-S'_i) = I(\theta)$$

$$Var(-\overline{S'}) = \frac{1}{n} \cdot Var(S'_i(\theta)) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

$$\begin{aligned} \text{Finally, } \frac{1}{\sqrt{n}}S_n(\theta) &\sim N(0, I(\theta)) \\ -\frac{1}{n}S'_n(\theta) &\xrightarrow{P} I(\theta) \end{aligned}$$

Then, by theorem from probability theory and statistics if  $X_n \rightsquigarrow X$  and  $Y_n \rightsquigarrow const$ , then  $X_n \cdot Y_n \rightsquigarrow cX$

$$\sqrt{n}(\hat{\theta}_n - \theta) \approx \frac{\frac{1}{\sqrt{n}}S_n(\theta)}{-\frac{1}{n}S'_n(\theta)} \rightsquigarrow \frac{N(0, I(\theta))}{I(\theta)} = N(0, \frac{1}{I(\theta)})$$

Meaning that  $\hat{\theta}_n \sim N(\theta, I_n^{-1}(\theta))$

So, the unicorn has been proved, but this is not useful because in the real world, true value of  $\theta$  cannot be obtained and used. So the question arises, having only estimator for theta, how can we use it in calculating its variance?

Def  $se$ - standard error of the estimator:  $se = \sqrt{\frac{1}{I_n(\theta)}}$  and  $\hat{se}$ - Max-likelihood estimator of  $se$ :  $\hat{se} = \sqrt{\frac{1}{I_n(\hat{\theta}_n)}}$  (That is my guess, maybe will be proved later)). Then, let's find the distribution of the estimator using  $\hat{se}$ :

$$\hat{\theta}_n \sim N(\theta, se^2)$$

$$\frac{\hat{\theta}_n - \theta}{\hat{se}} = (\hat{\theta} - \theta) \cdot \sqrt{n} \cdot \sqrt{I(\hat{\theta}_n)} = \sqrt{n} \cdot (\hat{\theta} - \theta) \cdot \sqrt{I(\theta)} \cdot \frac{\sqrt{I(\hat{\theta}_n)}}{\sqrt{I(\theta)}} =$$

= |since  $\hat{\theta}_n$  is a consistent estimator of  $\theta$ , then  $\hat{\theta}_n \xrightarrow{P} \theta$ ; assume  $I(\theta)$  is a continuous function, so  $\lim_{\hat{\theta}_n \rightarrow \theta} I(\hat{\theta}_n) \xrightarrow{P} I(\theta)$ | =

$$\Rightarrow \frac{\sqrt{I(\hat{\theta}_n)}}{\sqrt{I(\theta)}} \xrightarrow{P} 1 \text{ and } \sqrt{n} \cdot (\hat{\theta} - \theta) \cdot \sqrt{I(\theta)} \rightsquigarrow N(\theta, I_n^{-1}(\theta))$$

It means that  $\sqrt{n} \cdot (\hat{\theta} - \theta) \cdot \sqrt{I(\theta)} \cdot \frac{\sqrt{I(\hat{\theta}_n)}}{\sqrt{I(\theta)}} \sim N(0, 1)$

$$\hat{\theta}_n \sim N(\theta, I_n^{-1}(\hat{\theta}_n)) = N(\theta, \hat{se}^2)$$

Example: in order to prove what we have just found, I decided to use MLE for the population mean ( $\bar{X}$ ), which distribution parameters each of us knows:

$$X_i \sim N(\mu, \sigma^2), \sigma \text{ is known} \Leftrightarrow f(x_i, (\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\text{Need } I(\vec{\theta}) : S_i(\mu) = \frac{\partial}{\partial \mu} \left( -\frac{1}{2} \log 2\pi\sigma^2 - \frac{(X_i - \mu)^2}{2\sigma^2} \right) = \frac{X_i - \mu}{\sigma^2}$$

$$-S'_i(\mu) = -\frac{\partial}{\partial \mu} \left( \frac{X_i - \mu}{\sigma^2} \right) = \frac{1}{\sigma^2}$$

$$\text{Recall the previous result: } \hat{\theta}_n \sim N\left(\theta, \frac{1}{n \cdot I(\theta)}\right) \Rightarrow \bar{X} \sim N\left(\mu, \frac{1}{n \cdot \frac{1}{\sigma^2}}\right) = N\left(\mu, \frac{\sigma^2}{n}\right)$$

$$\text{Also, since } \hat{\theta}_n \sim N(\theta, I_n^{-1}(\hat{\theta}_n)), \text{ then } \bar{X} \sim N\left(\mu, \frac{\hat{\sigma}_{ml}^2}{n}\right)$$

this is written in blue, because I'm not sure about it.=)

#### 4.2.4 Asymptotically efficient

If an estimator has this property, it means that its variance is the smallest among all other estimators (as  $n \rightarrow \infty$ )

Consider two estimators of parameter  $\theta$ : Max-Likelihood estimator ( $\hat{\theta}_{ml}$ ) and another estimator  $\tilde{\theta}$ . Asymptotical efficiency means that:

$$Var(\hat{\theta}_{ml}) \leq Var(\tilde{\theta}) \quad \forall \tilde{\theta} \in \Theta$$

In order to prove this property, it is enough to recall the Cramer-Rao lower bound inequality:

$$Var(\hat{\theta}_n) \geq \frac{(1 + bias'(\theta))^2}{I_n(\theta)}$$

Then, as it was found before, if  $S_n(\hat{\theta}) = 0$ , then  $\hat{\theta}_n \sim N(\theta, \frac{1}{I_n(\theta)})$ , meaning that this estimator is unbiased and its variance is equal to the Cramer-Rao lower bound. However, there is a piece of uncertainty, which I want to explain in my future work: consider a harder example of the  $\hat{\theta}_{ml}$  for the bound of uniform distribution(it was calculated in "Way to obtain" part of the notes):  $X_i \sim U[0, \theta]$ , then,  $\hat{\theta}_{ml} = X_{max}$ , but at this point  $S_n(\theta) \neq 0$ . So it changes the proof of the asymptotically Normal property, but I have not understood, how this should be done... YET!!!

## 5 Appendix

(1) Proof of Chebychev's inequality

Consider  $g(x)$ - non-decreasing, non-negative function

$$\begin{aligned}
 E(g(x)) &= \int_{-\infty}^{+\infty} g(x)f(x)dx = \int_{g(x) \geq \epsilon} g(x)f(x)dx + \int_{g(x) < \epsilon} g(x)f(x)dx \\
 \int_{g(x) \geq \epsilon} g(x)f(x)dx + \int_{g(x) < \epsilon} g(x)f(x)dx &\geq \epsilon \cdot \int_{g(x) \geq \epsilon} f(x)dx + \int_{g(x) < \epsilon} g(x)f(x)dx = \\
 &= |\text{since } g(x)\text{-non-negative, second integral is a positive constant}| = \\
 E(g(x)) &\geq \epsilon \cdot \int_{g(x) \geq \epsilon} f(x)dx = \epsilon \cdot P(g(x) \geq \epsilon) \Rightarrow P(g(x) \geq \epsilon) \leq \frac{E(g(x))}{\epsilon} \\
 \text{Let } g(x) &= |\hat{\theta} - \theta|, \text{ then } P(|\hat{\theta} - \theta| \geq \epsilon) = P((\hat{\theta} - \theta)^2 \geq \epsilon^2) \leq \frac{E(\hat{\theta} - \theta)^2}{\epsilon} = \frac{MSE(\hat{\theta})}{\epsilon}
 \end{aligned}$$

(2) Proof of  $D(f, g) \geq 0$

$$D(f, g) = \int \log \frac{f(x)}{g(x)} f(x)dx = \int -\log \frac{g(x)}{f(x)} f(x)dx = E(-\log \frac{g(x)}{f(x)})$$

Consider  $y = -\ln(x)$ - convex function, then, by Jensen's inequality:

$$\begin{aligned}
 E(y(x)) &\geq y(E(x)) \\
 E(-\log \frac{g(x)}{f(x)}) &\geq -\log E(\frac{g(x)}{f(x)}) \\
 D(f, g) &\geq -\log(\int \frac{g(x)}{f(x)} \cdot f(x)dx) \\
 D(f, g) &\geq -\log(1) = 0
 \end{aligned}$$

(3) Proof of Cramer-Rao lower bound inequality:

- Consider the score function:  $S_n(\theta) = S(\vec{X}, \theta) \equiv \frac{\partial}{\partial \theta} \log f(\vec{X}, \theta) = n \cdot S_i(\theta)$
- $I_n(\theta) = \text{Var}(S_n(\theta)) = n \cdot \text{Var}(S_i(\theta)) = nI(\theta)$ , where  $I(\theta) \equiv I_1(\theta)$

First of all, need to find  $E(S_n(\theta))$ :

$$\begin{aligned}
 1 &= \int f(\vec{x}, \theta) d\vec{x} \Leftrightarrow 0 = \frac{\partial}{\partial \theta} \int f(\vec{x}, \theta) d\vec{x} = |\text{under regularity conditions}| = \int \frac{\partial}{\partial \theta} f(\vec{x}, \theta) d\vec{x} = \\
 &= \left| \frac{\partial}{\partial \theta} \log f(\vec{x}, \theta) = \frac{\frac{\partial f}{\partial \theta}}{f} \Leftrightarrow f \frac{\partial}{\partial \theta} \log f(\vec{x}, \theta) = \frac{\partial f}{\partial \theta} \right| = \int \frac{\partial}{\partial \theta} \log f(\vec{x}, \theta) f(\vec{x}, \theta) d\vec{x} = E(S_n(\theta))
 \end{aligned}$$

Proof for the most common case: biased estimator of n observations  $\hat{\theta}_n(\vec{x}) \equiv \hat{\theta}_n$

According to Cauchy-Swartz inequality:  $Cov^2(a, b) \leq Var(a) \cdot Var(b)$

$$\begin{aligned} Cov^2(\hat{\theta}_n, S_n(\theta)) &= \left( \int (\hat{\theta}_n - E(\hat{\theta}_n))(S_n(\theta) - E(S_n(\theta))) d\vec{x} \right)^2 = \\ &= \left( \int (\hat{\theta}_n - E(\hat{\theta}_n))(S_n(\theta) - E(S_n(\theta))) d\vec{x} \right)^2 \leq \int (\hat{\theta}_n - E(\hat{\theta}_n))^2 d\vec{x} \cdot \int (S_n(\theta) - E(S_n(\theta)))^2 d\vec{x} = Var(\hat{\theta}_n) \cdot I_n(\theta) \end{aligned}$$

$$\begin{aligned} Cov(\hat{\theta}, S_n(\theta)) &= E(\hat{\theta}_n \cdot S_n) - E(\hat{\theta}_n) \cdot E(S_n) = |E(S_n) = 0| = E(\hat{\theta}_n \cdot S_n) \\ E(\hat{\theta}_n \cdot S_n) &= \int \hat{\theta}_n \cdot S_n f(\vec{x}) d\vec{x} = \int \hat{\theta}_n \cdot \frac{\partial}{\partial \theta} \log f(\vec{x}) \cdot f(\vec{x}) d\vec{x} = \int \hat{\theta}_n \cdot \frac{\partial f(\vec{x})}{\partial \theta} d\vec{x} \\ \text{Under regularity conditions: } &\int \hat{\theta}_n \cdot \frac{\partial f(\vec{x})}{\partial \theta} d\vec{x} = \frac{\partial}{\partial \theta} \int \hat{\theta}_n \cdot f(\vec{x}) d\vec{x} = \frac{\partial}{\partial \theta} E(\hat{\theta}_n) \end{aligned}$$

Since  $\hat{\theta}_n$  could be a biased estimator of  $\theta$ , then  $E(\hat{\theta}_n) = \theta + bias(\hat{\theta}_n)$ ,

where  $bias(\hat{\theta}_n)$  is a function of  $\theta$ , so:

$$Cov(\hat{\theta}, S_n(\theta)) = \frac{\partial}{\partial \theta} E(\hat{\theta}_n) = \frac{\partial}{\partial \theta} (\theta + bias(\hat{\theta}_n)) = 1 + bias'(\hat{\theta}_n)$$

Finally,  $Cov^2(\hat{\theta}, S_n(\theta)) = (1 + bias'(\hat{\theta}_n))^2 \leq Var(\hat{\theta}_n) \cdot I_n(\theta)$

$$Var(\hat{\theta}_n) \geq \frac{(1 + bias'(\hat{\theta}_n))^2}{I_n(\theta)}$$

Another interesting fact about the score function: Let's calculate the  $Var(S_n)$ :

$$Var(S_n) = E(S_n)^2 - E^2(S_n) = E(S_n)^2 = E\left(\frac{\partial}{\partial \theta} \log f(\vec{x})\right)^2$$

In order to obtain an interesting result, we have to make some rearrangements:

$$\begin{aligned} \frac{\partial^2}{\partial \theta^2} \log f &= \frac{\partial}{\partial \theta} \left( \frac{\frac{\partial f}{\partial \theta}}{f} \right) = \frac{\frac{\partial^2 f}{\partial \theta^2} \cdot f - \left( \frac{\partial f}{\partial \theta} \right)^2}{f^2} = \frac{\partial^2 f}{\partial \theta^2} \cdot \frac{1}{f} - \left( \frac{\frac{\partial f}{\partial \theta}}{f} \right)^2 = \frac{\partial^2 f}{\partial \theta^2} \cdot \frac{1}{f} - \left( \frac{\partial}{\partial \theta} \log f \right)^2 \Rightarrow \\ E\left(\frac{\partial}{\partial \theta} \log f\right)^2 &= \int \left( \frac{\partial}{\partial \theta} \log f \right)^2 d\vec{x} = \int \frac{\partial^2 f}{\partial \theta^2} \cdot \frac{1}{f} \cdot f d\vec{x} - \int \frac{\partial^2}{\partial \theta^2} \log f \cdot f d\vec{x} = \\ &= |\text{Under regularity conditions}| = \frac{\partial^2}{\partial \theta^2} \int f d\vec{x} + E\left(-\frac{\partial^2}{\partial \theta^2} \log f\right) = E\left(-\frac{\partial^2}{\partial \theta^2} \log f\right) \end{aligned}$$

And we obtained the unexpected result: on average the second derivative of the likelihood function is negative and thus it is a concave function so if it has a critical point, it is a point of maximum. Also, if  $S_n(\vec{x}, \theta)$ - continuous function, then the point of maximum will be the only one.