# Maximum Likelihood Estimation (MLE)

1) MVUE might be difficult to _synthesize_
   → might not be _efficient_ ($\neq$ CRLB) & hence CRLB equality condition cannot be used.
   → Synthesis using "sufficient statistics" we will cover later (has its own problems)

2) MLE works always → nice _asymptotic_ performance guarantees (e.g. min variance, unbiasedness, _consistency_)
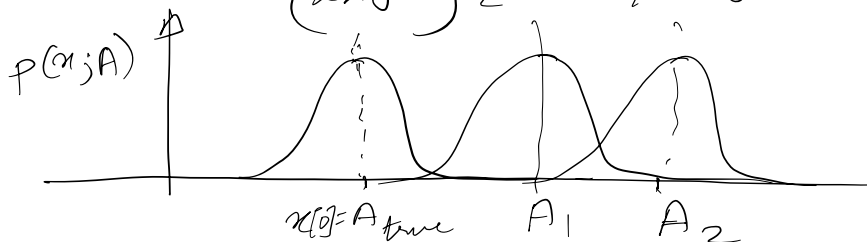
Let $p(x;\theta)$ $\left( x \in R^N ; \theta \in S \right)$ be the joint pdf / pmf of $\{x[0], \cdots, x[N-1]\}$. Then the _likelihood $f^n$_ is defined as $p(x=x_0;\theta) : S \rightarrow R$.

_Log·likelihood $f^n$_ : $\ln p(x=x_0;\theta)$

If in the Gaussian case; $x=x_0$ was observed;
$$x[0] = A + w[0] \qquad w[0] \rightsquigarrow WGN$$
$$p(x;A) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left[ -\frac{1}{2\sigma^2} (x-A)^2 \right]$$



$p(x;A)$ 

$x[0]=A_{true}$  $A_1$  $A_2$

So MLE $\hat{\theta} = \arg\max\limits_{\theta \in S} \ln p(x = x_0; \theta)$

For the example above: $\dfrac{\partial \ln p(x; A)}{\partial A} = \dfrac{1}{\sigma^2}(x[0] - A)$
$$= 0$$

$\Rightarrow \hat{A}_{MLE} = x[0]$.

Similarly for $x[n] = A + w[n]$ $\quad n = 0, 1, \cdots, N-1$
$\hat{A}_{MLE} = \dfrac{1}{N} \sum\limits_{n=0}^{N-1} x[n]$ $\qquad \hookrightarrow WGN$

\# In the above examples $\hat{A}_{MVUE} = \hat{A}_{MLE}$ can be verified — since $\hat{A}_{MVUE}$ could be calculated from CRLB.

\# However MLE is easy to find even if MVUE is hard.

Ex:
$\quad x[n] = A + w[n] \quad n = 0, \cdots, N-1 \quad \& \quad \underline{A > 0}$
$$\& \quad w[n] \sim N(0, A)$$

$p(x; A) = \dfrac{1}{(2\pi A)^{N/2}} \exp\left[ -\dfrac{1}{2A} \sum\limits_{n=0}^{N-1} (x[n] - A)^2 \right]$

$\dfrac{\partial \ln p(x; A)}{\partial A} = -\dfrac{N}{2A} + \dfrac{1}{A} \sum\limits_{n=0}^{N-1}(x[n] - A) + \dfrac{1}{2A^2} \sum\limits_{n=0}^{N-1}(x[0] - A)^2$

Check "reg" condition holds so CRLB applicable but it is difficult to see whether
$\dfrac{\partial \ln p(x, A)}{\partial A} = I(A)(\hat{A} - A)$ ??

In any case, CRLB can be computed

$$\text{Var}(\hat{A}) \geq \frac{A^2}{N(A+\frac{1}{2})} \quad \therefore \text{But no MVU estimator}$$

Now consider MLE:

$$\frac{\partial \ln p(x;A)}{\partial A} = 0 \implies \hat{A}^2 + \hat{A} - \frac{1}{N}\sum_{n=0}^{N-1} x^2[n] = 0$$

Solving & using $\hat{A} > 0 \implies \hat{A}_{MLE} = -\frac{1}{2} + \sqrt{\frac{1}{N}\sum_{n=0}^{N-1} x^2[n] + \frac{1}{4}}$

Q. So how good is this MLE? Is it MVUE? unbiased?
Is it efficient? We did not (could not) even compute var($\hat{A}_{MLE}$)!!

FACT: If an efficient estimator exists, the MLE will produce it. (Exercise)

Proof: $\frac{\partial \ln p(x;\theta)}{\partial \theta} = 0$ gives MLE.

But by CRLB an efficient estimator will satisfy $\frac{\partial \ln p(x;\theta)}{\partial \theta} = I^{-1}(\theta)[\hat{\theta} - \theta]$ ——①

clearly ① $= 0$ will produce $\hat{\theta} = \theta$.

# Asymptotic guarantees later

MLE works easily for vector parameters.
$$\theta = [\theta_1 \cdots \theta_p] \; ; \; \left[\frac{\partial \ln p(x;\theta)}{\partial \theta_1} \cdots \cdots \frac{\partial \ln p(x;\theta)}{\partial \theta_p}\right] = 0_{1\times p}$$

Ex: $x[n] = A + w[n]$    $n = 0, \cdots, N-1$

$\hookrightarrow N(0, \sigma^2)$ & $\sigma^2$ unknown.

$\theta = [A \; \sigma^2]^T$,    $\dfrac{\partial \ln p(x; \theta)}{\partial A} = \dfrac{1}{\sigma^2} \sum\limits_{n=0}^{N-1} (x[n] + A)$   ①

$\dfrac{\partial \ln p(x; \theta)}{\partial \sigma^2} = -\dfrac{N}{2\sigma^2} + \dfrac{1}{2\sigma^2} \sum\limits_{n=0}^{N-1} (x[n] - A)^2$   ②
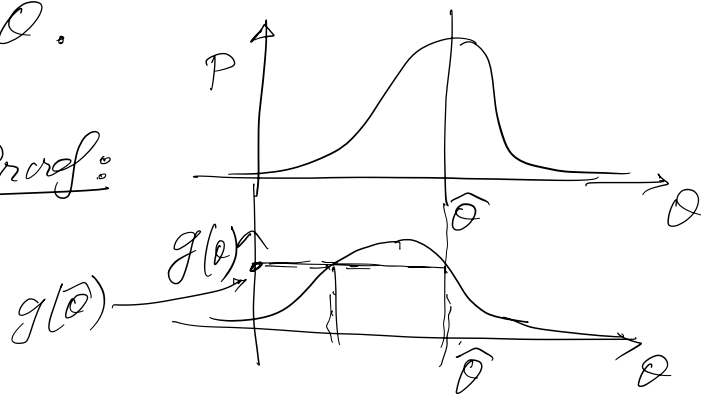
For ① $= 0$,   $\widehat{A}_{MLE} = \dfrac{1}{N} \sum\limits_{n=0}^{N-1} x[n]$

For ② $= 0$,   $\widehat{\sigma}^2 = \dfrac{1}{N} \sum\limits_{n=0}^{N-1} (x[n] - \widehat{A}_{MLE})^2$

## Invariance Property of MLE

FACT: The MLE of the parameter $\alpha = g(\theta)$ where pdf is $p(x; \theta)$ (parameterized by $\theta$) is given by $\widehat{\alpha} = g(\widehat{\theta})$ where $\widehat{\theta}$ is MLE of $\theta$.

Sketch of Proof:



$\left. \begin{array}{l} \max\limits_{\alpha} p(x; \theta) \\ \text{where } \alpha = g(\theta) \end{array} \right] \equiv \max\limits_{\alpha, \theta} \left\{ p(x; \theta) + \lambda \overbrace{[\alpha - g(\theta)]}^{H(\alpha, \theta)} \right\}$

$\dfrac{\partial H}{\partial \alpha} = 0 \; ; \; \dfrac{\partial H}{\partial \theta} = 0$

yields $\lambda = 0$,   $\dfrac{\partial p}{\partial \theta} = 0 \to \widehat{\theta}, g(\widehat{\theta})$ are the solutions.

Ex: $x[n] \sim N(0, \sigma^2) \rightarrow \sigma^2$ unknown $\rightarrow$ Power in dB to be estimated.

$$P = 10 \log_{10} \sigma^2$$

$$P(x, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp\left[-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} x^2[n]\right]$$

$$\frac{\partial \ln p(x; \sigma^2)}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{n=0}^{N-1} x^2[n] = 0$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{n=0}^{N-1} x^2[n]$$

. From invariance property:

$$\hat{P} = 10 \log_{10} \hat{\sigma}^2 = 10 \log_{10} \frac{1}{N} \sum_{n=0}^{N-1} x^2[n]$$

---

Recall:

Def$^n$: A seq of r.v.s $X_1, X_2, \ldots$ converges in prob. to a r.v. $X$, if for every $\varepsilon > 0$,

$$\lim_{n \to \infty} P(|X_n - X| \geq \varepsilon) = 0$$

or

$$\lim_{n \to \infty} P(|X_n - X| < \varepsilon) = 1$$

\# They converge almost surely to a r.v. $X$ if for every $\varepsilon > 0$,

$$P\left(\lim_{n \to \infty} |X_n - X| < \varepsilon\right) = 1$$

\# Then converge in distribution to a r.v. $X$ if

$$\lim_{n \to \infty} F_{X_n}(x) = F_X(x)$$

at all points where $F_X(x)$ is continuous.

FACT: A.S.C $\Rightarrow$ Conv. in prob $\Rightarrow$ Conv. in distribution

FACT 1: Weak Law of Large Nos:
Let $X_1, X_2 \cdots$ be iid r.v. with $EX_i = \mu$ and $Var\ X_i = \sigma^2 < \infty$. Define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$. Then for every $\varepsilon > 0$

$$\lim_{n \to \infty} P\left( |\bar{X}_n - \mu| < \varepsilon \right) = 1$$

i.e. $\bar{X}_n \xrightarrow{P} \mu$.

FACT 2: Strong law of large nos:
Let $X_1, X_2 \cdots$ be iid r.v's with $EX_i = \mu$ and $Var\ X_i = \sigma^2 < \infty$, define $\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$. Then for every $\varepsilon > 0$,

$$P\left( \lim_{n \to \infty} |\bar{X}_n - \mu| < \varepsilon \right) = 1$$

ie $\bar{X}_n \xrightarrow{a.s} \mu$.

FACT 2 $\Rightarrow$ FACT 1
(FACT 1 is still stated separately mainly for historical reasons & easy proof)

Central limit Theorem: Let $X_1, X_2, \dots$
be a seq. of iid r.v.'s with
$E X_i = \mu < \infty$ and $0 < \text{Var } X_i = \sigma^2 < \infty$
Define $\overline{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$. Let $F_n(x)$ denote
the cdf of $\frac{\sqrt{n}(\overline{X}_n - \mu)}{\sigma}$. Then

for any $x \in (-\infty, \infty)$,

$$\lim_{n \to \infty} F_n(x) = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2}{2}} dy$$

i.e. $F_n(x) \xrightarrow{d} N(0,1)$

---

Def$^n$: A sequence of estimators $\widehat{\theta}_N = \theta_N(x[0], \dots, x[N-1])$
is a consistent sequence of estimators of
$\theta$ if for every $\varepsilon > 0$, & $\forall \theta$,

$$\lim_{N \to \infty} P_\theta \left( | \widehat{\theta}_N(x) - \theta | < \varepsilon \right) = 1$$

$$\equiv \lim_{N \to \infty} P_\theta \left( | \widehat{\theta}_N(x) - \theta | \geq \varepsilon \right) = 0$$

\# $\widehat{\theta}_N(x)$ converges in prob. to $\theta$ (true) for
each $\theta$.

Ex: $x[n] = A + w[n]$ ; $w[n] \sim N(0,1)$

Consider the seq. of sample means

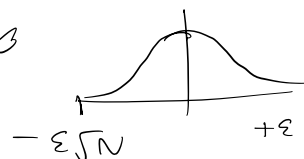$$\overline{x}_N = \frac{1}{N} \sum x[i^0] \qquad \overline{x}_N \sim N\left(A, \frac{1}{N}\right)$$

Hence $\quad P_A\left(|\overline{x}_N - A| < \varepsilon\right)$

$$= \int_{A-\varepsilon}^{A+\varepsilon} \left(\frac{N}{2\pi}\right)^{1/2} e^{-(N/2)(\overline{x}_N - A)^2} d\overline{x}_N = \int_{-\varepsilon}^{+\varepsilon} \left(\frac{N}{2\pi}\right)^{1/2} e^{-\frac{N}{2}(y^2)} dy$$

$$\left(y = \overline{x}_N - \theta\right)$$

$$= \int_{-\varepsilon\sqrt{N}}^{+\varepsilon\sqrt{N}} \left(\frac{1}{2\pi}\right)^{1/2} e^{-\left(\frac{1}{2}\right)t^2} dt = P\left(-\varepsilon\sqrt{N} < Z < \varepsilon\sqrt{N}\right) \quad Z \sim N(0,$$

$$\left(t = y\sqrt{n}\right) \qquad \rightarrow 1 \quad \text{as} \quad N \rightarrow \infty$$



$-\varepsilon\sqrt{N} \qquad +\varepsilon$

Hence $\quad \overline{x}_N$ is consistent estimator
for $A$.

---

<u>Usually a sufficient condition is used to</u>
check:

Recall, from Chebychev's Inequality,

$$P_\theta\left(|\widehat{\theta}_N - \theta| \geq \varepsilon\right) \leq \frac{E_\theta\left[(\widehat{\theta}_N - \theta)^2\right]}{\varepsilon^2}$$

But $\quad E_\theta\left[(\widehat{\theta}_N - \theta)^2\right] = \text{Var}_\theta \widehat{\theta}_N + \left[\text{Bias}_\theta \widehat{\theta}_N\right]^2$

Hence if $\quad \text{Var} \rightarrow 0 \quad \& \quad \text{Bias} \rightarrow 0 \quad \text{as} \quad N \rightarrow \infty$
then $\quad \widehat{\theta}_N \xrightarrow{P} \theta$

__Thm:__ If $\hat{\theta}_n$ is a sequence of estimators of a
parameter $\theta$ satisfying

    1> $\displaystyle\lim_{N \to \infty}$   $Var_\theta\left(\hat{\theta}_N\right) = 0$

    2> $\displaystyle\lim_{N \to \infty}$   $Bias_\theta\left(\hat{\theta}_N\right) = 0$    $\Bigg\}$   $\underline{\underline{\forall \theta}}$

1)

then $\hat{\theta}_N$ is a consistent sequence of estimators
of $\theta$.

__Ex:__ Easy to verify for the example above

    $E_\theta \, \overline{x}_N = \theta$        $Var_\theta\left(\overline{x}_N\right) = \dfrac{1}{N}$

---

## Consistency of MLE s:

__Thm:__ Let $x[1] \cdots x[n]$ be iid $p(x[i]; \theta)$
and let $L(\theta; x) = \displaystyle\prod_{i=1}^{n} p(x[i]; \theta)$ be the
likelihood fn. Let $\hat{\theta}_n$ denote the MLE
of $\theta$. Under "some __regularity conditions__" on
$p(x[i]; \theta)$, for every $\varepsilon > 0$ & $\forall \theta$,
    $\displaystyle\lim_{n \to \infty}$   $P_\theta\left(\left|\hat{\theta}_n - \theta\right| \geq \varepsilon\right) = 0$

i.e.   $\hat{\theta}_n$ is a consistent estimator of $\theta$

## Outline of proof:

    Let $\theta_0$ be the true parameter value.

Let $\quad L_n(\theta) = \frac{1}{n} \sum\limits_{i=1}^{n} \ln p(x[i]; \theta)$ $\quad$ depends on sample

$\&\quad L(\theta) = \int\limits_X [\ln p(x; \theta)] \, p(x; \theta_0) \, dx$

$\qquad\qquad = E_{\theta_0}[\ln p(x; \theta)]$

Then $\quad$ by $\quad$ L.L.N $\quad \forall \theta$,

$$L_n(\theta) \longrightarrow E_{\theta_0}[\ln p(x; \theta)] = L(\theta)$$

$\qquad\qquad\qquad\qquad\qquad\qquad$ does
$\qquad\qquad\qquad\qquad\qquad\qquad$ not depend on
$\qquad\qquad\qquad\qquad\qquad\qquad$ sample.

Claim: $\quad L(\theta) < L(\theta_0) \quad \forall \theta$
$[\text{Infact } P(L(\theta) < L(\theta_0)) = 1]$

Proof: Exercise

Then: proof follow from:
$\quad$ 1. $\widehat{\theta_n}$ is the maximizes of $L_n(\theta)$ $(MLE)$
$\quad$ 2. $\theta_0$ is the minimizes of $L(\theta)$ (by claim
$\quad$ 3. $\forall \theta$, $L_n(\theta) \longrightarrow L(\theta)$ by LLN.



$\qquad\qquad\qquad\qquad\qquad\qquad L(\theta)$

$\qquad\qquad\qquad\qquad\qquad\qquad L_n(\theta)$

$\qquad\qquad \widehat{\theta} \qquad \theta_0 \qquad\qquad\qquad\qquad \theta$

Q. What about asymptotic variance of the estimator?

Def$^{\underline{n}}$: For an estimator $\hat{\theta}_n$ of $\theta$, if
$$k_n (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, \sigma^2), \text{ where}$$
$\{k_n\}$ is a sequence of constants, then $\sigma^2$ is called the asymptotic variance of $\hat{\theta}_n$.

Example: For $n$ iid $N(\mu, \sigma^2)$ observations $X_1, \cdots, X_n$, $\bar{X}_n = \frac{1}{n} \sum_{i=1}^{n} X_i$,
$$\sqrt{n} (\bar{X}_n - \mu) \xrightarrow{d} N(0, \sigma^2)$$

Note: We need this $\sqrt{n}$, otherwise variance goes to $0$.

Def$^{\underline{n}}$: A sequence of estimators $\hat{\theta}_n$ is asymptotically efficient for $\theta$ if $\sqrt{n} (\hat{\theta}_n - \theta) \xrightarrow{d} N(0, I^{-1}(\theta))$ where $I(\theta) = -E\left[ \frac{\partial^2}{\partial \theta^2} \ln p(x; \theta) \right]$

$\left(\begin{array}{l} \text{i.e. asymptotically } \hat{\theta}_n \text{ is unbiased} \\ \text{\& } \quad " \quad n \text{Var}(\hat{\theta}_n) \text{ achieves CRLB} \end{array}\right)$

$\hookrightarrow$ Why? Recall $\text{Var}(\hat{\theta}_n) \to 0$.

__Thm:__ Let $x[1], \cdots, x[n]$ be iid $p(x[i]; \theta_0)$, $\hat{\theta}$ be MLE for $\theta_0$. Under "reg. conditions" on $p(x; \theta_0)$

$$\sqrt{n}(\hat{\theta} - \theta_0) \to N(0, I^{-1}(\theta_0)) \quad \text{where}$$

$I^{-1}(\theta_0)$ is the CRLB. $\Rightarrow \hat{\theta}$ is consistent & asymptotically efficient estimator of $\theta_0$.

__Proof:__ Let $L(\theta) = \sum_{i=1}^{n} \ln p(x_i; \theta)$ be the likelihood function & let $\theta_0$ be the true value of $\theta$.

__Taylor-series of__ $\frac{\partial L}{\partial \theta}$ around $\theta_0$:

$$\frac{\partial L}{\partial \theta} = \frac{\partial L}{\partial \theta}\Big|_{\theta_0} + (\theta - \theta_0)\frac{\partial^2 L}{\partial \theta^2}\Big|_{\theta_0}$$

For finding the MLE: $\hat{\theta}$, we put

$$\frac{\partial L}{\partial \theta} = 0 \implies (\hat{\theta} - \theta_0) = -\frac{\frac{\partial L}{\partial \theta}\big|_{\theta_0}}{\frac{\partial^2 L}{\partial \theta^2}\big|_{\theta_0}}$$

or

$$\sqrt{n}(\hat{\theta} - \theta_0) = \frac{-\frac{1}{\sqrt{n}}\frac{\partial L}{\partial \theta}\big|_{\theta_0}}{\frac{1}{n}\frac{\partial^2 L}{\partial \theta^2}\big|_{\theta_0}}$$

Now:

$$-\frac{1}{\sqrt{n}}\frac{\partial L}{\partial \theta}\Big|_{\theta_0} \xrightarrow{d} N(0, I(\theta_0)) \left.\begin{array}{c} \\ \\ \end{array}\right\} \text{Proof skipped}$$

$$\frac{1}{n}\frac{\partial^2 L}{\partial \theta^2}\Big|_{\theta_0} \xrightarrow{P} I(\theta_0)$$

Hence $\sqrt{n}\left(\hat{\theta} - \theta_0\right) \xrightarrow{d} N\left(0, \dfrac{1}{I(\theta_0)}\right)$