Statistics                                           Feb 24, 2020
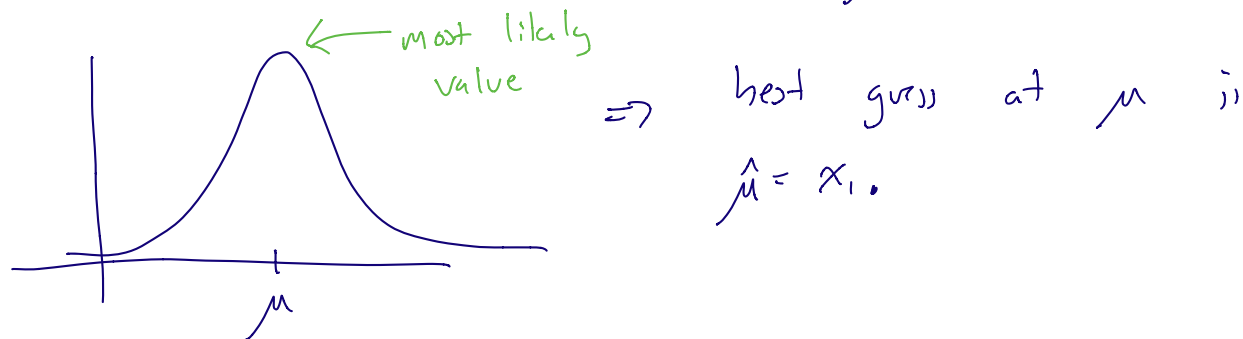
Example   Imagine you have   a   single piece of

data $x_1$, which you believe is an   observation

from a   normal   distribution:   $N(\mu, \sigma^2)$.



← most likely
value
$\Rightarrow$   best guess at $\mu$ is
$\hat{\mu} = x_1$.

I.e, we chose the value of $\mu$ that maximized
$$f(x_1; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_1-\mu)^2/2\sigma^2} . \quad \Big\} \text{ PDF of } N(\mu, \sigma^2)$$

evaluated at $x_1$.

$\Rightarrow$ This is in essence the method of Maximum Likelihood.

If $X_1, ..., X_n$ are a collection of random variables with

joint PDF $f = f(x_1, ..., x_n; \theta)$, then the **Likelihood**

function is:

$$\mathcal{L}(\theta) = f(x_1, ..., x_n; \theta)$$

$$= f(x_1; \theta) f(x_2; \theta) \cdots f(x_n; \theta) \quad \left( \text{if } X_1, ..., X_n \atop \text{are IID} \right)$$

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

Log-likelihood:   $\log \mathcal{L}(\theta) = \ell(\theta)$
$$= \sum_{i=1}^{n} \log f(x_i; \theta) \qquad \text{if IID.}$$

①

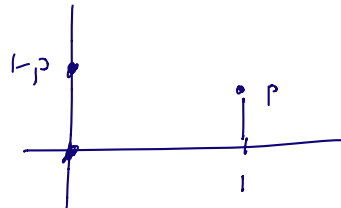The Maximum Likelihood estimator $\hat{\theta}$ is the value of $\theta$ which maximizes $L(\theta)$ or $\ell(\theta)$.

Notationally : $L(\theta) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x_i - \mu)/2\sigma^2}$

$$\sim \frac{1}{\sigma} e^{-(x_i - \mu)/2\sigma^2}$$

Example : $X_1, .., X_n \sim$ Bernoulli$(p)$    IID    r.v.'s

$X_i = 1$ with prob $p$

$X_i = 0$ with prob $1-p$.

$\Rightarrow$ mass $p$ at $X_i = 1$

mass $1-p$ at $X_i = 0$



$\Rightarrow f(x;p) = p^x (1-p)^{1-x}$

$\Rightarrow L(p) = \prod_{i=1}^{n} p^{X_i} (1-p)^{1-X_i}$

$$= p^{\Sigma X_i} (1-p)^{n - \Sigma X_i}$$

$\ell(p) = (\Sigma X_i) \log p + (n - \Sigma X_i) \log (1-p)$.

$\Rightarrow$ Solve $\ell'(p) = 0$

$\ell'(p) = \frac{1}{p} \Sigma X_i - \frac{1}{1-p} (n - \Sigma X_i) = 0$.

$\Rightarrow \boxed{\hat{p} = \frac{1}{n} \Sigma X_i}$.

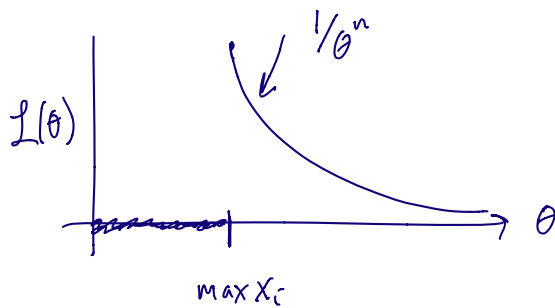Example: Let $X_1, ..., X_n \sim U(0, \theta)$ IID.

$$f(x; \theta) = \begin{cases} \frac{1}{\theta} & \text{for } x \in [0, \theta] \\ 0 & \text{otherwise} \end{cases}$$

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

$$L'(\theta) = \frac{d}{d\theta} \frac{1}{\theta^n} = -n \frac{1}{\theta^{n-1}} = 0$$

$$= \begin{cases} \frac{1}{\theta^n} & \text{if all } x_i \leq \theta \\ 0 & \text{otherwise} \end{cases}, \quad \text{i.e. if } \max x_i > \theta.$$

Plot:



$$\Rightarrow \quad \hat{\theta} = \max X_i.$$

## Properties of the Maximum Likelihood Estimator

### Consistency

true value of $\theta$

$$\Rightarrow \quad \hat{\theta} \xrightarrow{P} \theta_*$$

i.e. $P\left( |\hat{\theta} - \theta_*| > \epsilon \right) \to 0$ for any $\epsilon > 0$ as $n \to \infty$.

### Sketch of Proof

Consider first the Kullback-Leibler "distance"

$$D(f, g) = \int f(x) \log\left( \frac{f(x)}{g(x)} \right) dx.$$

"distance" between pdfs $f$ and $g$.

$$\Rightarrow \quad D(f,g) \geq 0$$

$$D(f,f) = 0$$

Write $D(\theta, \psi) = D\left(f(x;\theta), f(x;\psi)\right)$.

We say that the statistical model $\mathcal{F}$ is $\underline{\text{identifiable}}$

if $\quad \theta \neq \psi \quad \Rightarrow \quad D(\theta, \psi) > 0$.

Now, maximizing $\ell(\theta)$ is equivalent to maximizing

$$M(\theta) = \frac{1}{n} \sum_i \log \frac{f(x_i;\theta)}{f(x_i;\theta_*)}$$

$$= \frac{1}{n} \sum_i \left( \log f(x_i;\theta) - \log f(x_i;\theta_*) \right)$$

$$= \frac{1}{n} \left( \ell(\theta) - \underbrace{\ell(\theta_*)}_{} \right)$$

<span style="color:red">constant with respect to $\theta$.</span>

If $\quad \mathbb{E}\left( \log \frac{f(x_i;\theta)}{f(x_i;\theta_*)} \right)$ exists, then by the

Law of Large Numbers, as $n \to \infty$ $M(\theta)$ converges

to

$$\mathbb{E}_{\theta_*}\left( \log \frac{f(x_i;\theta)}{f(x_i;\theta_*)} \right) = \int \log \frac{f(x;\theta)}{f(x;\theta_*)} f(x;\theta_*) \, dx$$

$$= -\int \log \frac{f(x;\theta_*)}{f(x;\theta)} f(x;\theta_*) \, dx$$

$$= -D(\theta_*, \theta).$$

So for large $n$, $M(\theta) \simeq -D(\theta_*, \theta)$, which is

maximized at $\theta = \theta_*$ since $-D(\theta_*, \theta_*) = 0$ and

$-D(\theta_*, \theta) < 0$ for $\theta \neq \theta_*$.

$\Rightarrow$ The maximizer of $M(\theta)$ tends to $\theta_*$.

## Equivariance

Thm: Let $\tau = g(\theta)$ be a function of $\theta$.

Let $\hat{\theta}$ be the MLE of $\theta$. Then $\hat{\tau} = g(\hat{\theta})$

is the MLE of $\tau$.

## Asymptotic Normality

Goal show that $\hat{\theta} \to N(\theta_*, ?)$

Def: Score function

$$s(x; \theta) = \frac{\partial}{\partial \theta} \log f(x; \theta)$$

Fisher Information:

$$I(\theta) = \text{Var}\left(s(x; \theta)\right)$$

$$I_n(\theta) = \text{Var}\left(\sum_i s(x_i; \theta)\right) \qquad \text{and since } x_i$$
$$\text{are IID}$$

$$= \sum_i \text{Var}\left(s(x_i; \theta)\right).$$

Compute the expected value of score function:

$$\mathbb{E}\left(s(X;\theta)\right) = \int \frac{\partial}{\partial\theta} \log f(x;\theta) \quad f(x;\theta) \; dx$$

$$= \int \frac{1}{f(x;\theta)} \left(\frac{\partial}{\partial\theta} f(x;\theta)\right) f(x;\theta) \; dx$$

$$= \frac{\partial}{\partial\theta} \int_{-\infty}^{\infty} f(x;\theta) \; dx$$

$$= \frac{\partial}{\partial\theta} 1 \quad = 0$$

$$\Rightarrow \quad \mathrm{Var}\left(s(x;\theta)\right) = \mathbb{E}\left(s(x;\theta)^2\right) - \left(\mathbb{E}\left(s(x;\theta)\right)\right)^2$$

$$= \mathbb{E}\left(s(x;\theta)^2\right)$$

Thm: $I_n(\theta) = n \, I(\theta)$ and furthermore,

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

$$= -\mathbb{E}_\theta\left(\frac{\partial}{\partial\theta} s(x;\theta)\right).$$

$$= \mathrm{Var}\left(s(x;\theta)\right)$$

$$= \mathbb{E}\left(s(x;\theta)^2\right)$$

Thm: Let $\text{se} = \sqrt{\mathrm{Var}(\hat{\theta})}$, "under appropriate regularity conditions"

① $\text{se} \approx \sqrt{1/I_n(\theta)}$ and $\dfrac{\hat{\theta} - \theta}{\text{se}} \rightsquigarrow N(0,1)$

② Let $\hat{\text{se}} \approx \sqrt{1/I_n(\hat{\theta})}$, then $\dfrac{\hat{\theta} - \theta}{\hat{\text{se}}} \rightsquigarrow N(0,1)$.

# Asymptotic Confidence Intervals

Let $C = \left( \hat{\theta} - z_{\alpha/2}\, \hat{se}, \; \hat{\theta} + z_{\alpha/2}\, \hat{se} \right)$

then $P_\theta\left(\theta \in C\right) \to 1 - \alpha$ as $n \to \infty$.

(Same exact proof as before, just rearrange terms).

# Optimality

Suppose we want to estimate $\mu$ from $X_1, .., X_n \sim N(\mu, \sigma^u)$ IID.

Let $\hat{\mu} = MLE = \frac{1}{n} \sum X_i$.

We could alternatively estimate $\mu$ using $\tilde{\mu} = \text{median}(X_1,.., X_n)$.

It can be shown that

$$\sqrt{n}\left(\hat{\mu} - \mu\right) \rightsquigarrow N\left(0, \sigma^2\right)$$

$$\sqrt{n}\left(\tilde{\mu} - \mu\right) \rightsquigarrow N\left(0, \sigma^2 \pi/2\right).$$

In general, let $T, U$ be two estimators of $\theta$, each of which is asymptotically normal:

$$\sqrt{n}\left(T - \theta\right) \rightsquigarrow N(0, t^2)$$

$$\sqrt{n}\left(U - \theta\right) \rightsquigarrow N(0, u^2)$$

$ARE(U, T)$ = asymptotic relative efficiency of $U$ to $T$

$= t^2/u^2$ = ratio of variances.

Back to our example:

$$ARE(\tilde{\mu}, \hat{\mu}) = \frac{Var(\hat{\mu})}{Var(\tilde{\mu})} = \frac{\sigma^2}{\frac{\pi}{2}\sigma^2} = \frac{2}{\pi} \approx .63.$$

<u>Thm</u>: If $\hat{\theta}$ is the MLE of $\theta$, and $\tilde{\theta}$ is any other asymptotically normal estimator, then

$$ARE(\tilde{\theta}, \hat{\theta}) \leq 1.$$

$\Rightarrow$ the MLE is <u>efficient</u>   or   <u>asymptotically optimal</u>

$\Rightarrow$ MLE has the smallest asymptotic variance.

## Multiparameter Models

Extend to models with several parameters:

Let $\theta = (\theta_1, ..., \theta_k)$, and let

$\hat{\theta} = (\hat{\theta}_1, ..., \hat{\theta}_k)$ be the MLE, i.e, the

solution to the system of equations

$$\frac{\partial}{\partial \theta_1} \ell(\theta) = 0$$

$$\frac{\partial}{\partial \theta_2} \ell(\theta) = 0$$

$$\vdots$$

$$\frac{\partial}{\partial \theta_k} \ell(\theta) = 0$$

system of $k$ equations in the $k$ unknowns $\theta_1, ..., \theta_k$.

Let us also define $H_{jk} = \frac{\partial^2 \ell}{\partial \theta_j \partial \theta_k}$

## Fisher Information Matrix:

$$I_n(\theta) = - \begin{pmatrix} \mathbb{E}(H_{11}) & \cdots\cdots & \mathbb{E}(H_{1k}) \\ \mathbb{E}(H_{21}) & & \vdots \\ \vdots & & \vdots \\ \mathbb{E}(H_{k1}) & & \mathbb{E}(H_{kk}) \end{pmatrix}$$

and $\quad J_n = I_n^{-1}$. ( Question for home: why does $J_n$ exist? ).

__Thm__ Under the same regularity conditions on $f$ as before,

$$\hat{\theta} - \theta \approx N(0, J_n).$$

$\underset{\textcolor{red}{\llcorner \text{this is a } k\text{-dimensional vector.}}}{}$

And furthermore if $\hat{\theta}_j$ is the $j^{th}$ component of $\hat{\theta}$, then

$$\frac{\hat{\theta}_j - \theta_j}{\hat{se}_j} \rightsquigarrow N(0,1)$$

where $\quad \hat{se}_j^2 = J_n(j,j)$

and $\quad \text{Cov}(\hat{\theta}_j, \hat{\theta}_k) \approx J_n(j,k).$

9