

if, want to name one

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

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} f(x|\theta)$

if distribution is $\begin{cases} \text{continuous } f(x|\theta) \\ \text{discrete } p(x|\theta) / P(X=x|\theta) \end{cases}$

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} N(\theta, 1) \quad f(x|\theta) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\theta)^2}{2}}$

$x_1, x_2, \dots, x_n \stackrel{iid}{\sim} \text{pois}(\theta) \quad P(X=x|\theta) = \frac{\theta^x \cdot e^{-\theta}}{x!}$

Joint density $f(x_1, x_2, \dots, x_n; \theta) \stackrel{iid}{=} f(x_1|\theta) f(x_2|\theta) \dots f(x_n|\theta)$
 the "chance" of observing the data given one value of θ .

Likelihood of θ : defined as:

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

one argument

treat x_1, x_2, \dots, x_n as fixed.

$\rightarrow n$ arguments, x_1, x_2, \dots, x_n (given θ)

MLE of θ : $\hat{\theta}_{MLE} = \underset{\theta}{\text{argmax}} L(\theta)$
 : argument value that max fun $L(\theta)$

`mu=1, n=100 x<-rnorm(n, mean=mu)`

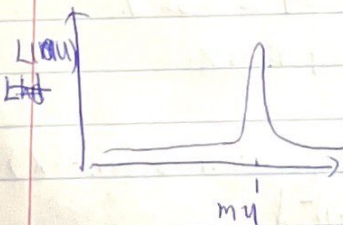
`mu-list <- seq(-2, 2, 0.01) [my guesses of mu]`

`lhd-list <- sapply(mu-list, FUN=function(mu-cand) {`

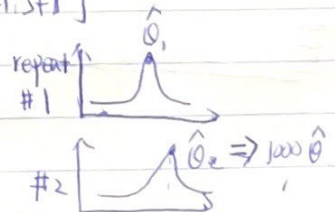
`prod(dnorm(x, mean=mu-cand)) # likelihood under μ / over μ -list`
`} # likelihood is the product of n densities / "prob of x_1, \dots, x_n around μ "`

`plot(x=mu-list, y=lhd-list); mu-list[which.max(lhd-list)]`

repeat 1000 times.

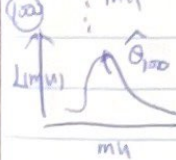
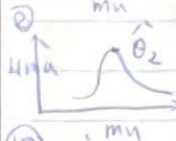
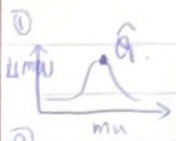


\Rightarrow `sapply(1:1000, FUN=function(i) {`
`x<-rnorm(n, mean=mu)`



`mu-list[which.max(lhd-list)]`

[each time, x_1, x_2, \dots, x_n are changing, so $L(\theta)$ changing]



distribution of $\hat{\theta}$.

$\hat{\theta}_{MLE}$ will change if data (x_1, x_2, \dots, x_n) changes
So $\hat{\theta}$ is a RV as X_1, \dots, X_n are RVs. *

[Likelihood function changes with RV: X_5]

For one sample, $\hat{\theta}$ may be away from the true θ , but stats is talking abt overall distribution (not one particular trial). So, on average, $\hat{\theta}$ is close to true θ which is 1 in this example.

intuitively, MLE is a Normal distribution

Properties of $\hat{\theta}$

- ① bias: $\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$ [$\hat{\theta}$ is a RV. we look at the expectation?]
- ② variance $\text{var}(\hat{\theta}) = E[(\hat{\theta} - E(\hat{\theta}))^2]$ by def
- ③ MSE (Mean squared error) $\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$

Cramer-Rao Inequality (Theorem): the minimum variance of unbiased $\hat{\theta}$.

Let $x_1, x_2, \dots, x_n \stackrel{iid}{\sim}$ pdf $f(\cdot; \theta)$;

Let $\hat{\theta} = g(x_1, x_2, \dots, x_n)$ be an unbiased estimator of θ .

Then under smoothness conditions of $f(\cdot; \theta)$ [differentiable]

$$\text{var}(\hat{\theta}) \geq \frac{1}{n I(\theta)}$$

where $I(\theta)$ is Fisher information

Fisher information: based on the log-likelihood of θ w/ 1 obs. x_i

$$L(\theta) = \log f(x_i; \theta) \quad f \text{ is the pdf/pmf of r.v. } x_i$$

$$I(\theta) = E\left[\left(\frac{d}{d\theta} L(\theta)\right)^2\right] \quad (\theta \text{ is a fixed, expectation is w.r.t. } X)$$

$$= E_X\left[\left(\frac{d}{d\theta} \log f(x, \theta)\right)^2\right]$$

evaluate how much information does the density provide abt θ

It's a property of the distribution alone, not depending on real data (Expectation)

一阶导数: 函数值变化大.

二阶导数: (函数的) 斜率变化大, 曲率↑

Alternative def: $I(\theta) = -E\left[\frac{d^2}{d\theta^2} L(\theta)\right]$

so when θ is given, $L(\theta)$ is known. $L(\theta)$ may or may not depend on θ .

proof.

$$\int f(x; \theta) dx = 1$$

$$[\text{def } L(\theta) = \log f(x; \theta)]$$

$$\frac{d}{d\theta} \int f(x; \theta) dx = \frac{d}{d\theta} 1 = 0.$$

because of the smoothness condition of $f(\cdot; \theta)$

swap and \int

$$\frac{d}{d\theta} \int \frac{d}{d\theta} f(x; \theta) dx = \frac{d}{d\theta} 1 = 0. \quad (1)$$

$$\frac{d}{d\theta} \log f(x; \theta) = \frac{\frac{d}{d\theta} f(x; \theta)}{f(x; \theta)} \quad [\text{chain rule}]$$

$$\therefore \frac{d}{d\theta} f(x; \theta) = \left(\frac{d}{d\theta} \log f(x; \theta) \right) \cdot f(x; \theta)$$

apply $\int \cdot dx$ to both sides:

$$\int \frac{d}{d\theta} f(x; \theta) dx = \int \left(\frac{d}{d\theta} \log f(x; \theta) \right) f(x; \theta) dx.$$

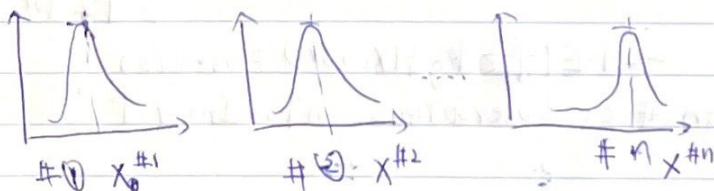
only value a RV takes

$$\overset{(1)}{\downarrow} 0 = E\left(\frac{d}{d\theta} \log f(X; \theta) \right) \quad [\text{by def of expectation}]$$

→ a fun of X → RV

Intuitively, FI indicates "how fast is the likelihood changing at current value of θ ." (curvature)

FI is the average of first derivative over all X s.



$I(\theta) \uparrow \rightarrow \frac{dL(\theta)}{d\theta} \uparrow \rightarrow$ change θ a little bit, $L(\theta)$ changes a lot/the most.
→ θ_{true} is a sensitive spot of Likelihood function \rightarrow differentiate θ from other θ s. $I(\theta_{\text{true}})$ is the max \uparrow

$$\frac{d}{d\theta} f(x; \theta) = \left(\frac{d}{d\theta} \log f(x; \theta) \right) \cdot f(x; \theta) \quad \textcircled{2}$$

(a·b)' = a'b + b'a

take derivative on ②

$$\frac{d^2}{d\theta^2} f(x; \theta) = \left(\frac{d^2}{d\theta^2} \log f(x; \theta) \right) \cdot f(x; \theta) + \frac{d}{d\theta} \log f(x; \theta) \cdot \left(\frac{d}{d\theta} f(x; \theta) \right)$$

apply integration

$$\int \frac{d^2}{d\theta^2} f(x; \theta) dx \stackrel{\text{smoothness}}{=} \frac{d^2}{d\theta^2} \int f(x; \theta) dx = \frac{d^2}{d\theta^2} 1 = 0 \quad [\text{LHS}]$$

$$[\text{RHS}] = \int \frac{d^2}{d\theta^2} \log f(x; \theta) \cdot f(x; \theta) dx + \int \frac{d}{d\theta} \log f(x; \theta) \cdot \left(\frac{d}{d\theta} \log f(x; \theta) \right) \cdot f(x; \theta) dx$$

$$= E\left(\frac{d^2}{d\theta^2} \log f(X; \theta)\right) + E\left[\left(\frac{d}{d\theta} \log f(X; \theta)\right)^2\right]$$

$$\text{LHS} = 0 = \text{RHS} = E\left(\frac{d^2}{d\theta^2} \log f(X; \theta)\right) + E\left[\left(\frac{d}{d\theta} \log f(X; \theta)\right)^2\right] \quad \square$$

Joint Fisher Information $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} f(\cdot; \theta)$

$$I_n(\theta) = E\left[\left(\frac{d}{d\theta} \log f(X_1, X_2, \dots, X_n; \theta)\right)^2\right]$$

$$= E\left[\left(\frac{d}{d\theta} \sum_{i=1}^n \log f(X_i; \theta)\right)^2\right] \quad [\text{iid}]$$

$$I_n(\theta) = -E\left[\frac{d^2}{d\theta^2} \log f(X_1, X_2, \dots, X_n; \theta)\right]$$

$$= -E\left(\frac{d^2}{d\theta^2} \sum_{i=1}^n \log f(X_i; \theta)\right) \quad [\text{put derivative inside sum}]$$

$$= -E\left(\sum_{i=1}^n \frac{d^2}{d\theta^2} \log f(X_i; \theta)\right)$$

$$= \sum_{i=1}^n E\left(\frac{d^2}{d\theta^2} \log f(X_i; \theta)\right)$$

$$= -n \cdot E\left(\frac{d^2}{d\theta^2} \log f(X_i; \theta)\right) = n \cdot I(\theta) \quad \rightarrow \text{FI for one RV.}$$

"FI scaled with # of observations, $n \uparrow$, $I_n(\theta) \uparrow$."

Cramer-Rao lower bound: $\hat{\theta} = g(X_1, X_2, \dots, X_n)$ is unbiased, $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} f(\cdot, \theta)$

Then $\text{Var}(\hat{\theta}) \geq \frac{1}{n I(\theta)} = \frac{1}{I_n(\theta)}$

correlation

$$\text{Cor}(X, Y) = \frac{\text{COV}(X, Y)}{\sqrt{\text{Var}X} \sqrt{\text{Var}Y}} \quad [-1, 1]$$

$$\text{Cor}^2(X, Y) = \frac{\text{COV}^2(X, Y)}{\text{Var}(X)\text{Var}(Y)} \quad [0, 1] \leq 1.$$

Cauchy-Schwartz $\text{COV}^2(X, Y) \leq \text{Var}(X)\text{Var}(Y)$

proof of use some RV T, Z which satisfy:

Cramer-Rao $\text{Var}(\hat{\theta}) \geq \frac{\text{COV}(\hat{\theta}, T)}{\text{Var}(T)} \quad \text{for some } T \quad \frac{1}{nI(\theta)}$

based on previous knowledge: Let $T = \frac{d}{d\theta} \log f(x_1, x_2, \dots, x_n; \theta) = \frac{d}{d\theta} \sum_{i=1}^n \log f(x_i; \theta)$

$$\text{Var}(T) = E[T^2] - (ET)^2$$

$$= E\left[\left(\frac{d}{d\theta} \log f(x_1, \dots, x_n; \theta)\right)^2\right] - (ET)^2 = I_n(\theta) - (ET)^2 = I_n(\theta) = nI(\theta)$$

def of $I_n(\theta)$

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

$$E(T) = \int \dots \int \sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i; \theta)}{f(x_i; \theta)} f(x_1, x_2, \dots, x_n) dx_1 \dots dx_n$$

$$= \int \dots \int \left(\sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i; \theta)}{f(x_i; \theta)} \right) f(x_1; \theta) f(x_2; \theta) \dots f(x_n; \theta) dx_1 dx_2 \dots dx_n$$

$$= 0 \quad / \quad E(T) = E\left(\sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i; \theta)}{f(x_i; \theta)}\right) \stackrel{\text{iid}}{=} n E\left(\frac{\frac{d}{d\theta} f(x_1; \theta)}{f(x_1; \theta)}\right) \stackrel{\text{shown before}}{=} n \cdot 0 = 0$$

$$\text{COV}(\hat{\theta}, T) = E(\hat{\theta}T) - E(\hat{\theta}) \frac{E(T)}{\theta} = E(\hat{\theta}T)$$

$$E(\hat{\theta}T) \stackrel{\text{bedef}}{=} \int \dots \int g(\theta, x_1, x_2, \dots, x_n) \left(\sum_{i=1}^n \frac{\frac{d}{d\theta} f(x_i; \theta)}{f(x_i; \theta)} \right) \left(\prod_{i=1}^n f(x_i; \theta) \right) dx_1 \dots dx_n$$