

Definition 6.11 (score equation, score function)

- The log likelihood of an i.i.d. sample of size n from a pdf/pmf $f(x|\theta)$ is

$$l(\theta) = \sum_{i=1}^n \log f(X_i|\theta).$$

score function of n obs. = sum of score function of i th obs
(Exercise) Check the previous MLE examples and identify the score equation & score function.
- The MLE maximizes $l(\theta)$ and is usually obtained by solving the score equation

$$0 = \frac{\partial}{\partial \theta} l(\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i|\theta).$$

sum of scores = 0 at $\hat{\theta}_{MLE}$
Score of i th observation X_i
a function of observations X_i 's and parameter θ
- $\frac{\partial}{\partial \theta} l(\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i|\theta)$ is called score function.

Definition 6.12 (Fisher information for one-dimensional parameter, TBp.263)

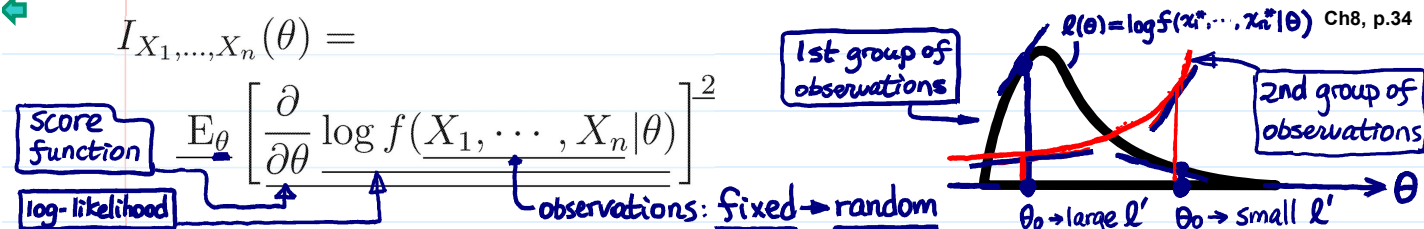
3/18

Let X_1, \dots, X_n be a sample of size n with a joint pdf/pmf f . Define

$$I_{X_1, \dots, X_n}(\theta) = E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1, \dots, X_n|\theta) \right]^2$$

This is a function of unknown parameter θ
can allow not to be i.i.d.
average (w.r.t. the parameter θ)
% $\log f(X|\theta)$ is a r.v., but not statistic
weighted average of score² → always ≥ 0
treated as random variable
 which is called the (Fisher) information of θ contained in X_1, \dots, X_n .

Question: What information does $I_{X_1, \dots, X_n}(\theta)$ offer?

**Theorem 6.4** (TBp. 276)

Let X_1, \dots, X_n be an i.i.d. sample of size n from a pdf/pmf $f(x|\theta)$.

$$I_{X_1, \dots, X_n}(\theta) = E_{\theta} \left(\frac{\partial}{\partial \theta} \log \left[\prod_{i=1}^n f(X_i|\theta) \right] \right)^2 = E_{\theta} \left[\sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i|\theta) \right]^2$$

Y₁, ..., Y_n indep.
score function of i th observation
Y_i

$$= \sum_{i=1}^n E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_i|\theta) \right]^2 + 2 \sum_{i < j} E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_i|\theta) \right] E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_j|\theta) \right]$$

∴ independent
∴ identical

$$= n E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right]^2 \equiv n \cdot I_{X_1}(\theta)$$
= 0, by (Δ) in LNp.35.
E(Y_iY_j) = E(Y_i)E(Y_j)

- $I_{X_1}(\theta)$ is the Fisher information contained in a sample of size one.
larger sample size, more information.
- $n I_{X_1}(\theta)$: interpreted as the information of θ contained in a sample of size n from $f(\cdot|\theta)$.
X, Y indep ⇒ I_{X,Y}(θ) = I_X(θ) + I_Y(θ)
more information
cor(X, Y) = 0
Y = aX + b
less information
cor(X, Y) = ±1

The Fisher informations of independent samples are additive.

Theorem 6.4 (TBp. 276)Under appropriate smoothness conditions on f , pdf/pmf

$$I_{X_1}(\theta) \equiv E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right]^2 = \text{Var}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right] = -E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X_1|\theta) \right]$$

$\xleftarrow{\text{cf. same r.v. } Y_1}$

Proof (for pdf case): Since $\int f(x|\theta) dx = 1$ for all θ ,

$$\begin{aligned} 0 &= \frac{\partial}{\partial \theta} \int f(x|\theta) dx = \int \frac{\partial}{\partial \theta} f(x|\theta) dx = \int \left[\frac{\frac{\partial}{\partial \theta} f(x|\theta)}{f(x|\theta)} \right] f(x|\theta) dx \\ &= \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] f(x|\theta) dx = E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right] \quad \dots (\Delta) \\ &\Rightarrow \text{Var}_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right] \quad \textcircled{Y_1} \quad \text{score function of one observation} \\ &= E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right]^2 - \left\{ E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right] \right\}^2 \quad (\Rightarrow \text{average of scores} = 0) \\ &= E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right]^2. \quad \textcircled{Y_1} \end{aligned}$$

$$\begin{aligned} 0 &= \frac{\partial^2}{\partial \theta^2} \int f(x|\theta) dx = \frac{\partial}{\partial \theta} \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right] f(x|\theta) dx \\ &= \int \left[\frac{\partial^2}{\partial \theta^2} \log f(x|\theta) \right] f(x|\theta) dx + \int \left[\frac{\partial}{\partial \theta} \log f(x|\theta) \right]^2 f(x|\theta) dx \\ &= E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X_1|\theta) \right] + E_{\theta} \left[\frac{\partial}{\partial \theta} \log f(X_1|\theta) \right]^2. \quad \text{I}_{X_1}(\theta) \end{aligned}$$

(need smoothness of f for interchanging integration and differentiation.)**Note (TBp. 278, 279).**For i.i.d. case, ℓ : log-likelihood of n observations

$$\begin{aligned} E_{\theta} [l'(\theta)^2] &= I_{X_1, \dots, X_n}(\theta) = n \cdot I_{X_1}(\theta) = -n E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X_1|\theta) \right] \quad \text{use joint distribution of } X_1, \dots, X_n \\ &\stackrel{\text{same } \theta}{=} - \sum_{i=1}^n E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \log f(X_i|\theta) \right] = -E_{\theta} \left[\frac{\partial^2}{\partial \theta^2} \sum_{i=1}^n \log f(X_i|\theta) \right] = -E_{\theta} [l''(\theta)] \\ &\quad \because \text{identical} \end{aligned}$$

- **interpretation:** when $|E_{\theta} [l''(\theta)]|$ is large at $\theta = \theta_0$, $l(\theta)$ is, on average, changing rapidly in a vicinity of θ_0 . log-likelihood of X_1, \dots, X_n

Example 6.18 (Fisher information of i.i.d. Bernoulli $B(\theta)$)

Let X_1, \dots, X_n be i.i.d. from Bernoulli distribution $B(\theta)$ (i.e., the pmf of X_i is, statistical modeling (e.g., toss a coin) $\theta^{x_i}(1-\theta)^{1-x_i}$, for $x_i \in \{0, 1\}$), joint pmf $\propto \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$

then $E(X_i) = \theta$ and $\text{Var}(X_i) = \theta(1-\theta)$.

- For a single observatoin X_i , the first and second derivatives of its log likelihood are:

$$\begin{aligned} \log f(x|\theta) &= x \log \theta + (1-x) \log(1-\theta), \quad \text{score} \\ \frac{\partial}{\partial \theta} \log f(x|\theta) &= x/\theta - (1-x)/(1-\theta) = (x-\theta)/(\theta(1-\theta)), \\ \frac{\partial^2}{\partial \theta^2} \log f(x|\theta) &= -x/\theta^2 - (1-x)/(1-\theta)^2. \quad \begin{matrix} 1/\theta, \text{ when } x=1 \\ -1/(1-\theta), \text{ when } x=0 \end{matrix} \end{aligned}$$

x : fixed

➡ The Fisher information of a single observation, say X_1 , is

X : random

by definition 6.12
LNp.33

$$I_{X_1}(\theta) = \frac{E_{\theta} \left[\frac{X_1 - \theta}{\theta(1-\theta)} \right]^2}{\theta^2(1-\theta)^2} = \frac{E_{\theta}[(X_1 - \theta)^2]}{\theta^2(1-\theta)^2}$$

$$= \frac{\text{Var}_{\theta}(X_1)}{\theta^2(1-\theta)^2} = \frac{\theta(1-\theta)}{\theta^2(1-\theta)^2} = \frac{1}{\theta(1-\theta)}$$

$$I_{X_1}(\theta) = \text{Var}_{\theta} \left[\frac{X_1 - \theta}{\theta(1-\theta)} \right]$$

(Ec) $\rightarrow \frac{1}{\theta(1-\theta)}$

by Thm. 6.4
LNp.35

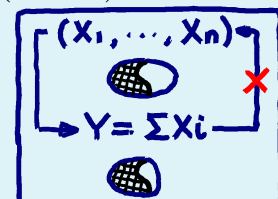
$$I_{X_1}(\theta) = -E_{\theta} \left[-\frac{X_1}{\theta^2} - \frac{1 - X_1}{(1-\theta)^2} \right]$$

$$= \frac{\theta}{\theta^2} + \frac{1-\theta}{(1-\theta)^2} = \frac{1}{\theta} + \frac{1}{1-\theta} = \frac{1}{\theta(1-\theta)}$$

• The Fisher information of observations X_1, \dots, X_n is

(Exercise) Use joint pmf to obtain

$$I_{X_1, \dots, X_n}(\theta) = n I_{X_1}(\theta) = \frac{n}{\theta(1-\theta)}$$



Notice that $I_{X_1, \dots, X_n}(\theta)$

\bar{X} estimate θ
 $\text{Var}(\bar{X}) = \frac{\theta(1-\theta)}{n}$
• large when $\theta \approx 1/2$
• small when $\theta \approx 0$ or 1

– increases when n increases.

– increases when $\theta \downarrow 0$ or $\theta \uparrow 1$.

– reaches a minimum $4n$ at $\theta = 0.5$.

larger sample size, more information

more information when θ is close to 0 or 1

minimum information when $\theta = 1/2$

• Consider a single observation $Y \sim \text{Binomial}(n, \theta)$. The pmf of Y is

$X_1, \dots, X_n \rightarrow$ non-invertible transformation ΣX_i

$$f(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}, \quad \text{for } y \in \{0, 1, \dots, n\}.$$

y : fixed

⊖ The second derivative of log likelihood is

$$\frac{\partial^2 \log f(y|\theta)}{\partial^2 \theta} = -y/\theta^2 - (n-y)/(1-\theta)^2.$$

check sufficient statistic in Ex. b.21
LNp.48

Y : random

⊖ The Fisher information of Y , is

$$I_Y(\theta) = -E_{\theta} \left[-\frac{Y}{\theta^2} - \frac{n-Y}{(1-\theta)^2} \right] = \frac{n\theta}{\theta^2} + \frac{n-n\theta}{(1-\theta)^2} = \frac{n}{\theta(1-\theta)}$$

(X_1, \dots, X_n)

& $Y = \Sigma X_i$ carry same amount of information

⊖ Note that $I_Y(\theta)$ is the same as $I_{X_1, \dots, X_n}(\theta)$. ← reasonable? ←

Theorem 6.5 (consistency of MLE, TBp. 275)

3/30

Under appropriate smoothness conditions of f , the MLE from an i.i.d. sample is consistent.

$\hat{\theta}_{MLE} \xrightarrow{P} \theta_0$ ← more general than LLN.

X_1, \dots, X_n i.i.d. $f(x|\theta_0)$

Proof (sketch, for pdf case): Denote the true value of θ by θ_0 . The MLE

maximizes $\frac{l(\theta)}{n} = \frac{1}{n} \sum_{i=1}^n \log f(X_i|\theta)$.

Y_i

Y_1, \dots, Y_n are i.i.d.

What is the difference between $\int [\log f(x|\theta)] f(x|\theta_0) dx$ and $\int [\log f(x|\theta)] f(x|\theta) dx$?

The weak law of large numbers implies

$$\forall \theta, \quad \frac{l(\theta)}{n} \xrightarrow{P} E_{\theta_0}[\log f(X|\theta)] = \int \underbrace{[\log f(x|\theta)]}_{\text{variable}} \underbrace{f(x|\theta_0)}_{\text{fixed}} dx \quad \text{as } n \rightarrow \infty.$$