

Lesson 9



## Statistics II

---

Today's topics

- maximum likelihood (最尤推定)

来嶋 秀治 (Shuji Kijima)

Dept. Informatics,  
Graduate School of ISEE

## Remind: Probability III

### i.i.d. (multivariate distribution)

Distribution of random variables  $X$  and  $Y$  of  $(\Omega, \mathcal{F}, P)$ .

Ex1. two dice.

$\Omega = \{(1,1), (1,2), \dots, (6,5), (6,6)\}$

$X =$  sum of casts

$Y =$  product of casts

例2. poker

choose five cards,

$X =$  # of A's

$Y =$  # of spades

Prop.

$$\Pr[(X_1 = x_1) \& (X_2 = x_2) \& \cdots \& (X_n = x_n)] \\ = \Pr[X_1 = x_1] \Pr[X_2 = x_2] \cdots \Pr[X_n = x_n]$$

Prop.

Suppose **discrete** r.v.  $X_1, X_2, \dots, X_n$  are i.i.d. w/pmf  $f$ .

Then  $f(x_1, x_2, \dots, x_n) = f(x_1)f(x_2) \cdots f(x_n)$ .

Prop.

Suppose **continuous** r.v.  $X_1, X_2, \dots, X_n$  are i.i.d. w/pdf  $f$ .

Then  $f(x_1, x_2, \dots, x_n) = f(x_1)f(x_2) \cdots f(x_n)$ .

Prop.

$$\Pr[(X \leq x) \& (Y \leq y)] = \Pr[X \leq x] \Pr[Y \leq y]$$

i.e.,  $F_{XY}(x, y) = F_X(x)F_Y(y)$

4

Prop.

Suppose *continuous* r.v.  $X$  and  $Y$  are **independent**.

Then,  $f_{XY}(x, y) = f_X(x)f_Y(y)$ .

Prop.

$$\Pr[(X \leq x) \& (Y \leq y)] = \Pr[X \leq x] \Pr[Y \leq y]$$

i.e.,  $F_{XY}(x, y) = F_X(x)F_Y(y)$

Prop.

Suppose **continuous** r.v.  $X$  and  $Y$  are **independent**.

Then,  $f_{XY}(x, y) = f_X(x)f_Y(y)$ .

Proof.

$$\begin{aligned} f(x, y) &:= \frac{\partial^2}{\partial x \partial y} F_{XY}(x, y) \\ &= \frac{\partial^2}{\partial x \partial y} (F_X(x)F_Y(y)) \\ &= \frac{\partial}{\partial x} \left( \left( \frac{\partial}{\partial y} F_X(x) \right) F_Y(y) \right) + \frac{\partial}{\partial x} \left( F_X(x) \left( \frac{\partial}{\partial y} F_Y(y) \right) \right) \\ &= 0 + \frac{\partial}{\partial x} (F_X(x)f_Y(y)) \\ &= \frac{\partial}{\partial x} F_X(x)f_Y(y) + F_X(x) \frac{\partial}{\partial x} f_Y(y) \\ &= f_X(x)f_Y(y) \end{aligned}$$

## Joint distribution of i.i.d. r.v.

Prop.

Suppose  $X_1, \dots, X_n$  are **i.i.d.** with density function  $f(x)$ .

Then, the pdf of the **joint distribution** of  $(X_1, \dots, X_n)$  is

$$f(x_1, \dots, x_n) := f(x_1) \cdots f(x_n).$$



# Estimating the population parameters

---

母パラメータの推定

## Maximum likelihood (最尤推定)

### maximum likelihood

Suppose samples  $x_1, \dots, x_n$  are i.i.d.

with a prob. func.  $f(x; \theta)$

where  $\theta$  is an **unknown parameter**.

Find  $\theta$  to maximize  $\Pr[X = \mathbf{x} \mid \theta]$ .



## Statistical inference: maximum likelihood

### Example 1

The number of defective products per 10,000 products.

lot	1	2	3	4	5	6	7	8	9	10
#defective	0	2	0	0	1	1	0	3	1	0

How often do defectives appear?

Let  $X$  be a r.v. denoting #defectives,

Then  $X \sim \text{Po}(\lambda)$ , i.e.,

$$\Pr[X = x] =: f(x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

for unknown parameter  $\lambda$ .

## Maximum likelihood

$\theta$ : parameters  
e.g.  $N(\mu, \sigma)$ ,  $E(\lambda)$

### Preparation

Let  $X_1, \dots, X_n$  be **i.i.d.** with density function  $f(x; \theta)$ , and

let  $f(x_1, \dots, x_n; \theta) := f(x_1; \theta) \cdots f(x_n; \theta)$ ,

i.e., density function of a **joint distribution** of  $(X_1, \dots, X_n)$ .

remark

$X_1, \dots, X_n$  are independent

### Maximum likelihood estimation

Given **sample values**  $X_1 = a_1, \dots, X_n = a_n$ ,

let  $L(\theta | \mathbf{a}) := f(a_1, \dots, a_n; \theta)$ , called **likelihood function**, and

$\theta^* = \arg \max_{\theta} L(\theta)$  is called **maximum likelihood estimator**.

## Ex. Maximum likelihood

Poisson distribution  $Po(\lambda)$

$$f(x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

### Ex. Poisson distribution

Let  $X_1, \dots, X_n$  be independent r.v.s according to  $Po(\lambda)$ , and  $a_1, \dots, a_n$  are sample values.

$$L(\lambda; a_1, \dots, a_n) = e^{-\lambda} \frac{\lambda^{a_1}}{a_1!} e^{-\lambda} \frac{\lambda^{a_2}}{a_2!} \dots e^{-\lambda} \frac{\lambda^{a_n}}{a_n!}$$

max. likelihood estimator of  $\lambda = \operatorname{argmax}_{\lambda} L(\lambda)$

$$\frac{\partial}{\partial \lambda} L(\lambda; a_1, \dots, a_n) = \dots$$

## Ex. Maximum likelihood

Poisson distribution  $Po(\lambda)$

$$f(x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

### Ex. Poisson distribution

Let  $X_1, \dots, X_n$  be independent r.v.s according to  $Po(\lambda)$ , and  $a_1, \dots, a_n$  are sample values.

$$\begin{aligned} \log(L(\lambda)) &= \sum_{i=1}^n \log(f(a_i; \lambda)) = \sum_{i=1}^n (-\lambda + a_i \log \lambda - \log(a_i!)) \\ &= n(\bar{a} \log \lambda - \lambda) - \sum_{i=1}^n \log(a_i!) \end{aligned}$$

$$\frac{\partial}{\partial \lambda} \log(L(\lambda)) = n \left( \bar{a} \frac{1}{\lambda} - 1 \right)$$

$\bar{a}$  is the maximum likelihood estimator of  $\lambda$

## Statistical inference: maximum likelihood

### Example 2

The scores of examination.

student	1	2	3	4	5	6	7	8	9	10
score	72	89	64	52	96	64	70	83	56	70

How much ratio do they understand?

Let  $X$  be a r.v. denoting scores,

Then  $X \sim N(\mu, \sigma^2)$ , i.e.,

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

for unknown parameters  $\mu$  and  $\sigma$ .

## Ex. Maximum likelihood

### Ex. Normal distribution

Let  $X_1, \dots, X_n$  be independent r.v.s according to  $N(\mu, \sigma^2)$ , and  $a_1, \dots, a_n$  are sample values.

$$\begin{aligned}\ln(L(\mu, \sigma; a)) &= \sum_{i=1}^n \ln(f(a_i; \mu, \sigma)) \\ &= \sum_{i=1}^n \ln\left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(a_i - \mu)^2}{2\sigma^2}\right)\right) \\ &= \sum_{i=1}^n \left(-\frac{1}{2} \ln(2\pi) - \frac{1}{2} \ln \sigma^2 - \frac{(a_i - \mu)^2}{2\sigma^2}\right) \\ &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln \sigma^2 - \sum_{i=1}^n \frac{(a_i - \mu)^2}{2\sigma^2}\end{aligned}$$

## Ex. Maximum likelihood

### Ex. Normal distribution

Let  $X_1, \dots, X_n$  be independent r.v.s according to  $N(\mu, \sigma^2)$ , and  $a_1, \dots, a_n$  are sample values.

$$\frac{\partial}{\partial \mu} \ln(L(\mu, \sigma; a)) = -\frac{\partial}{\partial \mu} \frac{\sum_{i=1}^n (a_i - \mu)^2}{2\sigma^2} = -\frac{\sum_{i=1}^n (a_i - \mu)}{\sigma^2}$$

$\bar{a}$  is the maximum likelihood estimator of  $\mu$

$$\begin{aligned} & \frac{\partial}{\partial \sigma} \ln(L(\mu, \sigma; a)) \\ &= -\frac{\partial}{\partial \sigma} \frac{n}{2} \ln \sigma^2 - \frac{\partial}{\partial \sigma} \frac{\sum_{i=1}^n (a_i - \mu)^2}{2\sigma^2} \\ &= -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (a_i - \mu)^2}{\sigma^3} \end{aligned}$$

Since  $\mu = \bar{a}$  maximize  $L(\mu, \sigma)$  (independent of  $\sigma$ ),

$$(\sigma^*)^2 = \frac{\sum (a_i - \bar{a})^2}{n}$$

## Ex. maximum likelihood

### Example 2 (Hardy Weinberg)

Blood type	A	B	AB	O
samples	43	12	6	39
possible type	aa, ao	bb,bo	ab	oo

Estimate the ratio  $\theta_a, \theta_b, \theta_o$  of  $a, b, o$

Model

$$\begin{aligned}
 p_A &= \theta_a^2 + 2\theta_a\theta_o \\
 p_B &= \theta_b^2 + 2\theta_b\theta_o \\
 p_{AB} &= 2\theta_a\theta_b \\
 p_O &= \theta_o^2
 \end{aligned}$$

Hardy Weinberg equilibrium

where

$$\begin{aligned}
 p_A + p_B + p_{AB} + p_O &= 1 \\
 \theta_a + \theta_b + \theta_o &= 1
 \end{aligned}$$



## Ex. maximum likelihood

### Example 2 (Hardy Weinberg)

Blood type	A	B	AB	O
samples	43	12	6	39
possible type	aa, ao	bb,bo	ab	oo

Estimate the ratio  $\theta_a, \theta_b, \theta_o$  of  $a, b, o$

$$\begin{aligned}
 & \log L(\theta) \\
 &= \log \left( \frac{n!}{x_A! x_B! x_{AB}! x_O!} p_A^{x_A} p_B^{x_B} p_{AB}^{x_{AB}} p_O^{x_O} \right) \\
 &= x_A \log p_A + x_B \log p_B + x_{AB} \log p_{AB} + x_O \log p_O + \log \left( \frac{n!}{x_A! x_B! x_{AB}! x_O!} \right) \\
 &= x_A \log(\theta_a^2 + 2\theta_a\theta_o) + x_B \log(\theta_b^2 + 2\theta_b\theta_o) + x_{AB} \log(2\theta_a\theta_b) + x_O \log(\theta_o^2) \\
 &+ \log \left( \frac{n!}{x_A! x_B! x_{AB}! x_O!} \right)
 \end{aligned}$$

## Ex. maximum likelihood

### Example 2 (Hardy Weinberg)

Blood type	A	B	AB	O
samples	43	12	6	39
possible type	aa, ao	bb,bo	ab	oo

Estimate the ratio  $\theta_a, \theta_b, \theta_o$  of  $a, b, o$

solution:  $\hat{\theta}_a = 0.285, \hat{\theta}_b = 0.094, \hat{\theta}_o = 0.621$

Blood type	A	B	AB	O
samples	43	12	6	39
	$\theta_a^2 + 2\theta_a\theta_o$	$\theta_b^2 + 2\theta_b\theta_o$	$2\theta_a\theta_b$	$\theta_o^2$
estimation	43.5	12.6	5.4	38.5



## Topics Related to max. likelihood

---

### Today's topics

- **KL divergence**
- **Consistency of max. likelihood**
- **Fisher Information (Fisher情報量)**
  - sharpness of max. likelihood
- **Cramer-Rao inequality**
  - a kind of central limit theorem

## Kullback-Leibler divergence

### Def. (KL divergence)

Let  $f, g$  be density function.

$$\text{KL}(g, f) = E_g \left[ \log \left( \frac{g(X)}{f(X)} \right) \right] = \int g(x) \log \left( \frac{g(x)}{f(x)} \right) dx$$

### Prop.

- $\text{KL}(g, f) \geq 0$
- $\text{KL}(g, f) = 0$  iff  $g \equiv f$

note that

✓  $\text{KL}(g, f) \neq \text{KL}(f, g)$ , in general

✓ *triangle inequality* does not hold, in general

i.e.  $\text{KL}(g, f) + \text{KL}(f, h) > \text{KL}(g, h)$  may hold.

## Consistency of Maximum likelihood

Thm. (consistency of maximum likelihood estimator)

If  $f$  is **strictly positive** (i.e.,  $f(x) > 0$  for any  $x$ ) and **continuous**, then the **maximum likelihood estimator**  $\theta^*$  of  $\theta_0$  is **consistent**.

i.e.,  $\forall \varepsilon > 0, \lim_{n \rightarrow \infty} \Pr[|\theta^* - \theta_0| < \varepsilon] = 1$

Rem.

$$\text{KL}(g, f) = E_g[\log g(X)] - E_g[\log f(X)]$$

$$\begin{aligned} \theta^* &= \arg \min_{\theta} \text{KL}(f_{\theta^*}, f_{\theta}) \\ &= \arg \min_{\theta} (E_{f_{\theta^*}}[\log f_{\theta^*}(X)] - E_{f_{\theta^*}}[\log f_{\theta}(X)]) \\ &= \arg \max_{\theta} E_{f_{\theta^*}}[\log f_{\theta}(X)] \end{aligned}$$

From the law of large number, intuitively

$$\arg \max_{\theta} \frac{1}{n} \sum_i^n \log f_{\theta}(x_i) \rightarrow \arg \max_{\theta} E_{f_{\theta^*}}[\log f_{\theta}(X)]$$

## Fisher Information

Proposition.

For a density function  $f(x; \theta)$ ,

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

# Fisher Information

## Proposition.

For a density function  $f(x; \theta)$ ,

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

$$\begin{aligned} & \mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) \right] \\ &= \mathbb{E} \left[ \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} \right] - \mathbb{E} \left[ \frac{\left( \frac{\partial}{\partial \theta} f(x; \theta) \right)^2}{(f(x; \theta))^2} \right] \\ &= \int_{-\infty}^{\infty} \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} f(x; \theta) dx - \mathbb{E} \left[ \frac{\left( \frac{\partial}{\partial \theta} f(x; \theta) \right)^2}{(f(x; \theta))^2} \right] \\ &= \int_{-\infty}^{\infty} \frac{\partial^2}{\partial \theta^2} f(x; \theta) dx - \mathbb{E} \left[ \frac{\left( \frac{\partial}{\partial \theta} f(x; \theta) \right)^2}{(f(x; \theta))^2} \right] \\ &= -\mathbb{E} \left[ \frac{\left( \frac{\partial}{\partial \theta} f(x; \theta) \right)^2}{(f(x; \theta))^2} \right] \end{aligned}$$

$$\begin{aligned} & \frac{\partial^2}{\partial \theta^2} \log f(x; \theta) \\ &= \frac{\partial}{\partial \theta} \left( \frac{\partial}{\partial \theta} \log f(x; \theta) \right) \\ &= \frac{\partial}{\partial \theta} \left( \frac{\frac{\partial}{\partial \theta} f(x; \theta)}{f(x; \theta)} \right) \\ &= \frac{\frac{\partial^2}{\partial \theta^2} f(x; \theta)}{f(x; \theta)} - \frac{\left( \frac{\partial}{\partial \theta} f(x; \theta) \right)^2}{(f(x; \theta))^2} \end{aligned}$$

$$\begin{aligned} & \int_{-\infty}^{\infty} f(x; \theta) dx = 1 \\ & \int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} f(x; \theta) dx = 0 \\ & \int_{-\infty}^{\infty} \frac{\partial^2}{\partial \theta^2} f(x; \theta) dx = 0 \end{aligned}$$

## Fisher Information

### Proposition.

For a density function  $f(x; \theta)$ ,

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

### Def.

For a density function  $f(x; \theta)$ ,

the **Fisher Information**  $I(\theta)$  is defined by

$$I(\theta) := \mathbb{E} \left[ \left( \frac{\partial}{\partial \theta} \log f(X; \theta) \right)^2 \right] = -\mathbb{E} \left[ \frac{\partial^2}{\partial \theta^2} \log f(X; \theta) \right]$$



## A type of central limit theorem

Thm.

Let  $f(x;\theta)$  be positive and twice differentiable, and let  $\theta^*$  be a max likelihood estimator of  $\theta_0$ .

Then  $\sqrt{n}(\theta^* - \theta_0) \sim N\left(0, \frac{1}{I(\theta)}\right)$  when  $n \rightarrow +\infty$ .

... meaning that

$$E[\theta^* | \theta_0] \simeq \theta_0$$

$$\text{Var}[\theta^* | \theta_0] \simeq \frac{1}{nI(\theta)}$$

for  $n \rightarrow +\infty$ .

## A type of central limit theorem

Thm. (Cramer-Rao's bound)

Let  $X_1, \dots, X_n$  be i.i.d. with density function  $f(x; \theta)$ , where  $f(x; \theta)$  is positive and twice differentiable, and let  $I(\theta)$  be the Fisher Information of  $f(x; \theta)$ .

If  $\hat{\theta}$  is an unbiased estimator, then

$$\text{Var}[\hat{\theta}] \geq \frac{1}{nI(\theta)}$$