

**Example 6.12** (i.i.d restricted normal distribution)

Suppose  $X_1, X_2, \dots, X_n$  are i.i.d. from  $N(\mu, 1)$  with  $0 \leq \mu < \infty$ . The log likelihood is

From the  $l(\mu, \sigma)$  in LNp 20

$$l(\mu) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n (X_i - \mu)^2$$

$$= -\frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})^2 - \frac{n}{2} (\bar{X} - \mu)^2$$

statistical modeling

sampling distribution=?  
Note:  $\bar{X} \sim N(\mu, 1/n)$

Hence the MLE of  $\mu$  is

always falls in  $(0, \infty)$  (Why?)

$$\hat{\mu} = \begin{cases} \bar{X}, & \text{if } \bar{X} \geq 0 \\ 0, & \text{if } \bar{X} < 0 \end{cases}$$

cf. moment estimator  $\hat{\mu} = \bar{X}$  (Ec)

reasonable?

cf. of  $N(\mu, 1/n)$

**Example 6.13** (i.i.d uniform(0,  $\theta$ ) distribution)

Suppose  $X_1, X_2, \dots, X_n$  are i.i.d.  $U(0, \theta)$ , where  $\theta > 0$ . Then the likelihood of  $\theta$  is

statistical modeling

$$\mathcal{L}(\theta) = \begin{cases} \theta^{-n}, & \text{if } 0 \leq X_i \leq \theta, i = 1, \dots, n \\ 0, & \text{otherwise} \end{cases}$$

$$= \begin{cases} \theta^{-n}, & \text{if } \theta \geq \max_{1 \leq i \leq n} X_i = X_{(n)} \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

Q: Are  $X_{(n)}, 2\bar{X}$  reasonable estimators?

sampling distribution=?

c.f. moment estimator  $\hat{\theta} = 2\bar{X}$  (Ec)  
 $\Rightarrow$  unbiased  
 $\Rightarrow$  but  $2\bar{X}$  might be  $< X_{(n)}$

alternative estimator  $\frac{n+1}{n} X_{(n)} = (1 + \frac{1}{n}) X_{(n)}$

Because  $\mathcal{L}(\theta)$  decreases when  $\theta$  increases, the MLE of  $\theta$  is  $X_{(n)}$ .

①  $E(X_{(n)}) = \frac{n}{n+1} \theta \leftarrow$  biased (Ec)  $\leftarrow$  always underestimate

**Example 6.14** (multinomial distribution, TBp. 272)

Suppose  $X_1, X_2, \dots, X_m$  are counts in cells  $1, 2, \dots, m$  and follow a multinomial distribution with total count  $n$  and cell probabilities  $p_1, p_2, \dots, p_m$  ( $p_i \geq 0$  for  $i = 1, 2, \dots, m$ , and  $p_1 + p_2 + \dots + p_m = 1$ ). The joint pmf of  $X_1, X_2, \dots, X_m$  is

statistical modeling

$$f(x_1, x_2, \dots, x_m | p_1, p_2, \dots, p_m) = \frac{n!}{\prod_{i=1}^m x_i!} \prod_{i=1}^m p_i^{x_i}$$

$\Rightarrow X_i \sim \text{binomial}(n, p_i)$ , but  $X_1, \dots, X_m$  not independent

dimension of parameter space =  $m-1$

where  $x_1 + x_2 + \dots + x_m = n$ . For  $n$  given, the log likelihood is

$$l(p_1, p_2, \dots, p_m) = \log n! - \sum_{i=1}^m \log x_i! + \sum_{i=1}^m X_i \log p_i$$

**MLE**  $\Rightarrow$  maximize  $l(p_1, p_2, \dots, p_m)$  subject to  $\sum_{i=1}^m p_i = 1$ .

Introduce a Lagrange multiplier  $\lambda$ , and maximize

$$l(p_1, p_2, \dots, p_m, \lambda) \equiv \log n! - \sum_{i=1}^m \log x_i! + \sum_{i=1}^m X_i \log p_i + \lambda \left( \sum_{i=1}^m p_i - 1 \right)$$

$\leftarrow l(p_1, \dots, p_m)$

Setting  $\frac{\partial l}{\partial p_i} = \frac{X_i}{p_i} + \lambda = 0$ ,  $i = 1, 2, \dots, m$  gives

$$\hat{p}_j = -\frac{X_j}{\lambda}, \quad j = 1, 2, \dots, m.$$

$$\frac{\partial l}{\partial \lambda} = \sum_{i=1}^m p_i - 1 = 0$$

Since

$$1 = \sum_{i=1}^m \hat{p}_i = \sum_{i=1}^m -\frac{X_i}{\lambda}, \quad 1 = -\frac{n}{\lambda}$$

Sampling distribution = ?

we have  $\lambda = -n$ . Hence,  $\hat{p}_i = X_i/n$ ,  $i = 1, 2, \dots, m$ .

reasonable?

**Example 6.15** (Hardy-Weinberg Equilibrium, TBp. 273)

**Hardy-Weinberg law:** if gene frequencies are in equilibrium, the genotypes AA, Aa, and aa occur in a population with frequencies  $(1 - \theta)^2$ ,  $2\theta(1 - \theta)$ , and  $\theta^2$ .  $\leftarrow \theta \in (0, 1) \text{ \& } (1-\theta)^2 + 2\theta(1-\theta) + \theta^2 = 1$

prob. of getting A  
 $1(1-\theta^2) + \frac{1}{2}(2\theta(1-\theta))$   
 $+ 0 \cdot \theta^2 = 1 - \theta$   
 prob. of getting a  
 $0(1-\theta^2) + \frac{1}{2}(2\theta(1-\theta))$   
 $+ 1 \cdot \theta^2 = \theta$

		Mother	
		A [1-θ]	a [θ]
Father	A [1-θ]	AA [(1-θ) <sup>2</sup> ]	Aa [θ(1-θ)]
	a [θ]	Aa [θ(1-θ)]	aa [θ <sup>2</sup> ]

**Question:** If we sample  $n$  (a fixed number) persons from the popu-  
lation, and let  $X_1, X_2$ , and  $X_3$  (random variables) denote the counts  
 in the three cells (AA, Aa, aa), what is a suitable statistical model  
 (i.e., joint distribution) for  $(X_1, X_2, X_3)$ ?

Notice that  $n = X_1 + X_2 + X_3$ .

large population,  
 without replacement  
 $\approx$  with replacement

Statistical modeling

$$(X_1, X_2, X_3) \sim \text{multinomial}(n, p_1, p_2, p_3)$$

$X_i \sim \text{binomial}(n, p_i)$   
 but  $X_1, X_2, X_3$  not independent

$$p_1 + p_2 + p_3 = 1$$

dim = 2

$\therefore$  genetics knowledge

Log likelihood of  $\theta$  is

$$\begin{aligned} l(\theta) &= \log n! - \sum_{i=1}^3 \log X_i! \\ &\quad + X_1 \log (1 - \theta)^2 + X_2 \log [2\theta(1 - \theta)] + X_3 \log \theta^2 \\ &= \log n! - \sum_{i=1}^3 \log X_i! \\ &\quad + (2X_1 + X_2) \log (1 - \theta) + (2X_3 + X_2) \log \theta + X_2 \log 2. \end{aligned}$$

Setting  $\frac{d}{d\theta} l(\theta) = 0$ ,

$$-\frac{2X_1 + X_2}{1 - \theta} + \frac{2X_3 + X_2}{\theta} = 0$$

Sampling distribution = ?

yields the MLE of  $\theta$

can be written as a function of  $X_1$  and  $X_2$

$$\hat{\theta} = \frac{2X_3 + X_2}{2X_1 + 2X_2 + 2X_3} = \frac{2X_3 + X_2}{2n} = \hat{p}_3 + \frac{1}{2}\hat{p}_2$$

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

**Question:** What is the difference between Ex. 6.14 and Ex. 6.15?

**Chinese population data of Hong Kong in 1937:** ( $\underline{M}$ ,  $\underline{N}$  are erythrocyte antigens)

Assume Hardy-Weinberg Equilibrium?

statistical modeling = ?

Blood Type	$\underline{M}$	$\underline{MN}$	$\underline{N}$	Total
Frequency	342	500	187	1029
(by Ex.6.14)	0.333	0.486	0.182	

$$342/1029$$

dim = 2

cf.

dim = 1

① MLE is  $\hat{\theta} = 0.4247 \Rightarrow (\hat{p}_1, \hat{p}_2, \hat{p}_3) = (0.331, 0.489, 0.180)$ .

2. Approximate the sampling distribution of  $\hat{\theta}$  by bootstrap:

simulation method

cf.

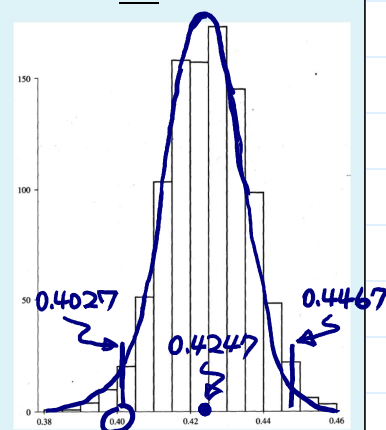
exact method  
asymptotic method

- Generate 1000 random counts from multinomial with  $n = 1029$  and cell probabilities 0.331, 0.489, and 0.180.
- From each of the 1000 experiments, a MLE value  $\hat{\theta}^*$  was determined.

From the histogram of the 1000 estimates (Figure 8.7 in Textbook), which should approximate the sampling distribution of  $\hat{\theta}$ .

- Looks like Normal.
- Standard deviation of the 1000 values gives estimated standard error of  $\hat{\theta}$ :

$$s_{\hat{\theta}} = 0.011. \quad 0.4247 \pm 2 \times (0.011) = [0.4027, 0.4467]$$



### Example 6.16 (Muon Decay, TBp. 266 & 271)

- Let  $\Theta$  be the angle at which electrons are emitted in muon decay.
- Let  $X = \cos(\Theta)$ . It has a distribution with pdf

data

$$f(x|\alpha) = \frac{1}{2}(1 + \alpha x), \quad -1 \leq x \leq 1, -1 \leq \alpha \leq 1.$$

a pdf?

parameter

- The mean of  $X$  is

1st moment

$$\mu = \alpha/3 \Rightarrow \alpha = 3\mu.$$

Ec

- The moments estimator of  $\alpha$  based on a sample  $X_1, \dots, X_n$  is  $\hat{\alpha} = 3\bar{X}$ .
- The log likelihood of  $\alpha$  is

$$l(\alpha) = \sum_{i=1}^n \log(1 + \alpha X_i) - n \log 2.$$

i.i.d.  $\sim f(x|\alpha)$

unbiased

Setting the derivative equal to zero, the MLE of  $\alpha$  satisfies the nonlinear equation

asymptotic method (LNp.39)

sampling distribution = ?

simulation method

$$0 = \frac{d}{d\alpha} l(\alpha) = \sum_{i=1}^n \frac{X_i}{1 + \alpha X_i}.$$

Dahlquist & Bjorck (1974)  
Chapter 6

- The MLE of  $\alpha$  has no easy close-form solution.

$\Rightarrow$  can use an iterative method to numerically solve for MLE.

$\Rightarrow$  method of moments estimate could be used as a starting value.



**Example 6.17** (i.i.d. Gamma distribution, TBp. 270)

- Suppose  $X_1, X_2, \dots, X_n$  are i.i.d.  $\Gamma(\alpha, \lambda)$ . The joint pdf is

*Statistical modeling*  $\rightarrow$  
$$f(x_1, x_2, \dots, x_n | \alpha, \lambda) = \prod_{i=1}^n \frac{1}{\Gamma(\alpha)} \lambda^\alpha x_i^{\alpha-1} e^{-\lambda x_i}$$

- The log likelihood is

$$\begin{aligned} l(\alpha, \lambda) &= \sum_{i=1}^n [\alpha \log \lambda + (\alpha - 1) \log X_i - \lambda X_i - \log \Gamma(\alpha)] \\ &= n\alpha \log \lambda + (\alpha - 1) \sum_{i=1}^n \log X_i - \lambda \sum_{i=1}^n X_i - n \log \Gamma(\alpha). \end{aligned}$$

- Setting 
$$\begin{cases} 0 = \frac{\partial l}{\partial \alpha} = n \log \lambda + \sum_{i=1}^n \log X_i - n \frac{\Gamma'(\alpha)}{\Gamma(\alpha)} \\ 0 = \frac{\partial l}{\partial \lambda} = \frac{n\alpha}{\lambda} - \sum_{i=1}^n X_i \end{cases}$$

- The MLE then satisfies

$(\hat{\alpha}, \hat{\lambda}) \rightarrow \begin{cases} \hat{\lambda} = \hat{\alpha} / \bar{X} \\ n \log \hat{\alpha} - n \log \bar{X} + \sum_{i=1}^n \log X_i - n \frac{\Gamma'(\hat{\alpha})}{\Gamma(\hat{\alpha})} = 0 \end{cases}$

- 2nd part is a nonlinear equation  $\Rightarrow$  **no easy closed-form solution.**  
 $\Rightarrow$  can use iterative method to find (approximate) the solution  
 $\Rightarrow$  method of moments estimates can be used as initial value.

$\hookrightarrow$  LNp.14

- Rainfall amount data** (Ex 6.7-6.8, LNp.14-16):

- Take the initial value as the method of moments estimates

$$\hat{\alpha} = 0.375, \quad \hat{\lambda} = 1.674.$$

By an iterative procedure, the MLE's are computed:

$$\hat{\alpha} = 0.441, \quad \hat{\lambda} = 1.96$$

**Q:** Which estimate is closer to the true values of parameters?

**C.f.**

**asymptotic method**  
(LNp.43)

$\Rightarrow$  of little practical difference from the moment estimates.

$\hookrightarrow$  check histograms in LNp.16

- Exact sampling distribution of the MLE is intractable  $\Rightarrow$  can use simulation to approximate:

$\therefore$  no close form

- Generate many, say 1000, samples of size 227 from Gamma with  $\alpha = 0.441, \lambda = 1.96$ .
- Form MLE of  $\alpha, \lambda$  for each sample.
- Construct histogram of the 1000 MLE's.

- From the histograms of the simulated MLEs:

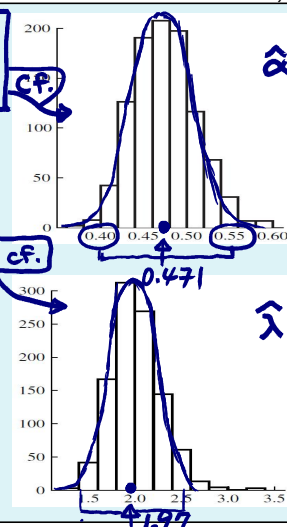
- The histograms look like normal.



$\hat{\alpha}$  might be biased

The histograms are centered at  $\hat{\alpha} = 0.471$  and  $\hat{\lambda} = 1.97$ .

The histograms in LNp.16



moment estimator  
 $S_{\hat{\alpha}} = 0.06, S_{\hat{\lambda}} = 0.34$   
(LNp.16)

cf.  $s_{\hat{\alpha}} = 0.04, s_{\hat{\lambda}} = 0.28$ .

Q: Which estimator will you prefer?

Sampling distribution of MLE's are less dispersed than those of the method of moments estimates.  
 $0.471 \pm 2 \times (0.04) = [0.391, 0.551]$   
 $1.97 \pm 2 \times (0.28) = [1.41, 2.53]$

### Summary (advantages of MLE)

1. easy to interpret  $\leftarrow$  LNp.18, definition
2. widely applicable
3. the range of the MLE coincides with the range of the parameter  $\leftarrow$  check Ex.6.12, 6.13 (LNp.22)
4. invariance under reparameterizations  $\leftarrow$  Thm 6.1 (LNp.19)
5. nice theoretical properties  $\leftarrow$  e.g. asymptotic properties.  
(Thm 6.5 ~ 6.9, LNp 38 ~ 43)

❖ Reading: textbook, 8.5, 8.5.1

### Large sample (asymptotic) theory for method of moment estimator and MLE

more data  
more information  
less random

Recall: 1. Law of Large Number  
2. Central Limit Theorem  $\rightarrow$  for sum/average estimator  $\hat{\theta}$

#### Definition 6.10 (consistent, TBp. 266)

Let  $\hat{\theta}_n$  be an estimator of a parameter  $\theta$  based on a sample of size  $n$ . Then  $\hat{\theta}_n$  is called **consistent in probability** if  $\hat{\theta}_n$  converges in probability to  $\theta$  as  $n$  tends to infinity, i.e. for any  $\epsilon > 0$ ,

sample size is large

LLN  
cf.  
 $\rightarrow$  for  $\bar{X}_n$

$$P(|\hat{\theta}_n - \theta| > \epsilon) \rightarrow 0 \text{ as } n \rightarrow \infty.$$

$$\begin{aligned} \hat{\theta}_n &\xrightarrow{P} \theta \\ \hat{\theta}_n &\xrightarrow{d} \theta \end{aligned}$$

$\because \theta$  is a constant

➤ method of moment estimator

#### Theorem 6.2 (consistency of method of moment estimator, TBp. 266)

The weak law of large numbers implies that

the  $g_i$ 's in LNp.9

$$\frac{1}{n} \sum_{i=1}^n X_i^k \equiv \hat{\mu}_k \xrightarrow{P} \mu_k \text{ in probability as } n \rightarrow \infty.$$

FYI,  $\frac{\sqrt{n}(\hat{\mu}_k - \mu_k)}{\sqrt{\mu_{2k} - \mu_k^2}} \xrightarrow{d} N(0,1)$

If the function relating  $\mu_k$  and  $\theta_j$  are continuous, method of moments estimators are consistent.

By Thm.5.1, item 4, LN, Ch1-6, p.89

**Theorem 6.3** (justification for estimating standard errors, TBp. 266-267)

- Recall: In LNp.12,  $\sigma_{\hat{\theta}}(\hat{\theta}) \xrightarrow{\text{estimate}} \sigma_{\hat{\theta}}(\theta)$  **estimated standard error**  $\xleftarrow{\text{代入法}}$  **standard error (as a function of  $\theta$ ) of sampling distribution of  $\hat{\theta}$ :**
- Consider the standard error of the form:  $\sigma_{\hat{\theta}}(\theta) = \frac{1}{\sqrt{n}} \sigma^*(\theta)$  **usually  $\sigma_{\hat{\theta}}(\theta) \rightarrow 0$  as  $n \rightarrow \infty$** 
  - Ex. 6.5 (LNp.11):  $\sigma_{\hat{\lambda}}(\lambda) = \frac{1}{\sqrt{n}} \sqrt{\lambda}$  **not a linear function of  $n$**   $\lambda \sim P(\lambda)$  **irrelevant to  $n$**
  - Ex. 6.6 (LNp.13):  $\sigma_{\hat{\mu}}(\mu, \sigma) = \frac{1}{\sqrt{n}} \sigma$ , and  $\sigma_{\hat{\sigma}^2}(\mu, \sigma) \approx \frac{1}{\sqrt{n}} \sqrt{2} \sigma^2$  **i.i.d.  $\sim N(\mu, \sigma^2)$**
- Let  $\theta_0 =$  true parameter, and  $\sigma_{\hat{\theta}} \equiv \sigma_{\hat{\theta}}(\theta_0) = \frac{1}{\sqrt{n}} \sigma^*(\theta_0)$ . **True standard error of sampling distribution**
- Estimate  $\sigma_{\hat{\theta}}$  by **a fixed value**  $\rightarrow s_{\hat{\theta}} \equiv \frac{1}{\sqrt{n}} \sigma^*(\hat{\theta})$ . **a r.v.**  $\rightarrow$   **$\sigma^*(\cdot)$**
- If (1)  $\sigma^*(\theta)$  is continuous in  $\theta$ , and (2)  $\hat{\theta}$  is consistent ( $\hat{\theta} \xrightarrow{P} \theta_0$ ), then
  - By Thm 5.1, item 4, LN, Ch1-6, p. 89**  $\rightarrow \sigma^*(\hat{\theta}) \xrightarrow{P} \sigma^*(\theta_0)$ , i.e., **e.g., moment estimator, MLE.**
  - as  $n \rightarrow \infty$ ,  $\frac{s_{\hat{\theta}}}{\sigma_{\hat{\theta}}} \rightarrow 1$  in probability** ( $\Rightarrow s_{\hat{\theta}} \xrightarrow{P} \sigma_{\hat{\theta}} \rightarrow 0$ )
  - Why not say  $s_{\hat{\theta}} \xrightarrow{P} \sigma_{\hat{\theta}}$ ?**  $\frac{s_{\hat{\theta}}}{\sigma_{\hat{\theta}}} \parallel \frac{\sigma^*(\hat{\theta})/\sigma^*(\theta_0)}{\sigma^*(\hat{\theta})/\sigma^*(\theta_0)}$  **same convergence rate** **when  $n$  is large**

➤ **MLE**

**Note.** the following discussion is mainly for (1) the case of i.i.d. sample, and (2) one-dimensional parameter.

**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 **score function of  $n$  obs. = sum of score function of  $i$ th obs**

$$l(\theta) = \sum_{i=1}^n \log f(X_i|\theta).$$
**(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  $X_i$**
- $\frac{\partial}{\partial \theta} l(\theta) = \sum_{i=1}^n \frac{\partial}{\partial \theta} \log f(X_i|\theta)$  is called score function. **a function of observations  $X_i$ 's and parameter  $\theta$**

**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

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

**This is a function of unknown parameter  $\theta$**  **can allow not to be i.i.d.** **average (w.r.t. the parameter  $\theta$ )**  **$\log f(X|\theta)$  is a r.v., but not statistic** **treated as random variable** **weighted average of score<sup>2</sup>**

which is called the **(Fisher) information** of  $\theta$  contained in  $X_1, \dots, X_n$ .

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