

Therefore,  $X_i \sim B(n, p_i) \Rightarrow \text{Var}(X_i) = np_i(1-p_i)$  [large when  $p=1/2$   
Small when  $p \approx 0$  or 1]

$$\Lambda = \prod_{i=1}^m \left( \frac{np_{i0}}{x_i} \right)^{x_i} \quad \text{and} \quad -2 \log \Lambda = 2 \sum_{i=1}^m x_i \log \left( \frac{x_i}{np_{i0}} \right).$$

–  $H_0$  is rejected if

$$-2 \log \Lambda(X_1, \dots, X_m) \geq c.$$

– The constant  $c$  can be determined by the fact that under  $H_0$ ,

$$\text{test statistic} \rightarrow -2 \log \Lambda \xrightarrow{D} \chi_{m-1}^2 \leftarrow \text{asymptotic null distribution}$$

(i.e.,  $-2 \log \Lambda$  is asymptotically  $\chi_{m-1}^2$  distributed when  $n \rightarrow \infty$ )

because

$$\dim(\Omega) - \dim(\Omega_0) = (m-1) - 0 = m-1.$$

Note: Not  $m \rightarrow \infty$   
(LN, CH8, p.86)

### Question 7.14 (examine distribution assumption in statistical modeling $\Rightarrow$ goodness of fit)

• In statistical modeling, we often see the statement:

$X_1, \dots, X_n$  are i.i.d. from a distribution with cdf  $F(\cdot|\Theta)$ .

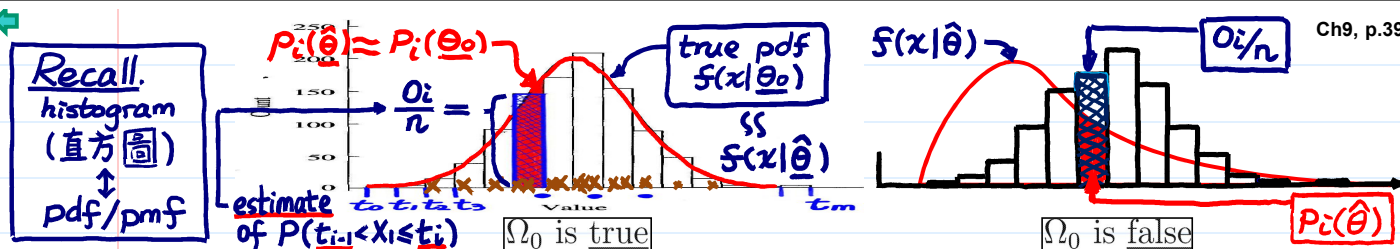
• Suppose that the independent and identical assumptions are true, can we examine (using data) whether the distribution assumption (i.e.,  $F(\cdot|\Theta)$ ) is reasonable?

original  
statistical  
model

$$\Omega = \{F(\cdot|\Theta)\} \xrightarrow{\text{enlarge}} \Omega = \{F(\cdot|\Theta)\} \cup \{\text{other distributions}\}$$

determine useful & useless  
information.  
Recall: sufficient  
statistics

e.g.  $\dim(\Omega_0) = \infty$  in LN p.5



• Observe data  $X_1, \dots, X_n$ .

• Let  $O_i, i = 1, \dots, m$ , be the number of  $X_1, \dots, X_n$  that fall in  $(t_{i-1}, t_i]$ .

• Then,

$$O_1, \dots, O_m \sim \text{Multinomial}(n, p_1, \dots, p_m),$$

where  $p_i$  is the “red area” in the graph when  $\Omega_0$  is true.

### Example 7.17 (GLR tests for goodness-of-fit, TBp.341-342)

•  $\Omega$ : the vector of cell probabilities  $\mathbf{p} = (p_1, \dots, p_m)$  that are free except for the constraints that  $p_i \geq 0, \sum_{i=1}^m p_i = 1$ .

– Note that  $\dim(\Omega) = m-1$ . Recall: Hardy-Weinberg Equilibrium (LN, CH8, p.24)

•  $H_0: \Omega_0 \subset \Omega$ ,

$$\Omega_0 = \{\mathbf{p}(\Theta) : \mathbf{p}(\Theta) = (p_1(\Theta), \dots, p_m(\Theta))\}$$

is the vector of cell probabilities from  $F(\cdot|\Theta)$ , where  $\Theta$  is a  $k$ -dimensional unknown parameter, i.e.,  $p_i(\Theta) = F(t_i|\Theta) - F(t_{i-1}|\Theta)$ .

– Note that  $\dim(\Omega_0) = k$ .

Assume  $m-1 > k$

$$P_{AA} = (1-\theta)^2, P_{Aa} = 2\theta(1-\theta), P_{aa} = \theta^2$$

e.g.

- Goodness-of-fit tests judge plausibility of the models in  $H_0$  relative to  $H_A: \Omega \setminus \Omega_0$  using data  $O_1, \dots, O_m$ .
- Likelihood ratio test
  - The LR statistic is

$$\Lambda = \frac{\max_{\mathbf{p} \in \Omega_0} \left[ \frac{n!}{o_1! o_2! \dots o_m!} p_1(\underline{\mathbf{p}})^{o_1} \dots p_m(\underline{\mathbf{p}})^{o_m} \right]}{\max_{\mathbf{p} \in \Omega} \left( \frac{n!}{o_1! o_2! \dots o_m!} p_1^{o_1} \dots p_m^{o_m} \right)} = \prod_{i=1}^m \left( \frac{p_i(\hat{\underline{\mathbf{p}}})}{\hat{p}_i} \right)^{o_i}$$

where  
 \*  $\hat{\underline{\mathbf{p}}} = \underline{\mathbf{p}}_i(\hat{\underline{\Theta}})$  ← restricted  
 \*  $\hat{\underline{\Theta}}$  is the MLE of  $\underline{\Theta}$ , and ← multinomial joint pmf  
 \*  $\hat{p}_1, \hat{p}_2, \dots, \hat{p}_m$  are the unrestricted MLE, i.e., ← MLE based on  $O_1, \dots, O_m$ , rather than  $X_1, \dots, X_n$

$$\hat{p}_i = o_i/n, i = 1, 2, \dots, m.$$

- Then, since  $O_i = n\hat{p}_i$ ,

$$-2 \log \Lambda = 2 \sum_{i=1}^m n\hat{p}_i \log \left( \frac{n\hat{p}_i}{np_i(\hat{\underline{\Theta}})} \right) = 2 \sum_{i=1}^m O_i \log \left( \frac{O_i}{E_i} \right),$$

where

\*  $O_i = n\hat{p}_i$  is the observed counts, and

\*  $E_i = np_i(\hat{\underline{\Theta}})$  is the expected counts.

- $H_0$  is rejected if  $-2 \log \Lambda \geq c$ .

← when  $H_0$  is true

use ratio to compare

cf.

← the test statistic in Ex.7.16 (LNp.38)

- To decide the constant  $c$ , since

$$\dim(\Omega) - \dim(\Omega_0) = m - k - 1,$$

under  $H_0$ , the large sample (i.e.,  $n$  is large) distribution of  $-2 \log \Lambda$  is  $\chi_{m-k-1}^2$ , i.e.,  $c = \chi_{m-k-1}^2(\alpha)$ . ← (1- $\alpha$ )-quantile of  $\chi_{m-k-1}^2$

### Question 7.15

Compare • the  $\Omega_0$  and  $\Omega$  in item 3, LNp.5, and ← data:  $X_1, \dots, X_n$

$\dim = k$  ←  $\dim = \infty$

• the  $\Omega_0$  and  $\Omega$  in Ex.7.17, LNp.39. ← data:  $O_1, \dots, O_m$

$\dim = k$  ←  $\dim = m-1$

Are they different? What and Why different?

### Example 7.18 (Pearson's Chi-square test, TBp.342-343)

- For the same  $H_0$  vs.  $H_A$  and data  $O_1, \dots, O_m$  in Ex.7.17, LNp.39, Pearson's chi-square test for goodness of fit:

- The test statistic is:

$$X^2 = \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i}$$

determine what observations "more extreme"

(reasonable?)

use difference to compare

Q: Why divided by  $E_i$ ?

Consider the 2 cases

①  $O_i = 2, E_i = 1$

②  $O_i = 999, E_i = 1000$  when  $n = 5000$

- $H_0$  is rejected if  $X^2 \geq c$ .

- Under  $H_0$ ,  $X^2 \xrightarrow{D} \chi_{m-k-1}^2$ .

← when  $n \rightarrow \infty$

← asymptotic null distribution

