

# Lecture 14

## §1 Review

### 1. Statistics 的定义

#### Definition

- Given samples:  $X_1, X_2, \dots, X_n$
- Goal: estimate an unknown parameter  $\theta$  of the true model, (we use samples).

#### The hat notation

The point estimator used to estimate a parameter  $\theta$  is usually denoted as  $\hat{\theta}$ .

## 2. MLE

### Formal Definition

- Given a model with an unknown parameter  $\theta$
- Given samples:  $X_1, X_2, \dots, X_n$
- The probability that the models generates the samples is called **likelihood**.  $L(\theta) = P(X_1, X_2, \dots, X_n | \theta)$
- To determine the best  $\theta$ , we choose  $\hat{\theta}$  such that  $L(\theta)$  (or equivalently, the **log-likelihood**  $l(\theta) = \log L(\theta)$ ) is maximize at  $\theta = \hat{\theta}$ . **Maximum likelihood estimate (MLE)**



### Likelihood function

- Given a model with an unknown parameter  $\theta$
- Given samples:  $X_1, X_2, \dots, X_n$

#### Continuous model:

- Likelihood:  $L(\theta) = \prod_i f(X_i | \theta)$  f: the model's probability density function (PDF)
- Log-Likelihood:  $l(\theta) = \sum_i \log(f(X_i | \theta))$

#### Discrete model:

- Likelihood:  $L(\theta) = \prod_i P(X_i | \theta)$  P: the model's probability mass function (PMF)
- Log-Likelihood:  $l(\theta) = \sum_i \log(P(X_i | \theta))$

### 例: Example continued

- For drug experiments,
  - N experiments  $N_1 + N_2 = N$
  - M successes  $M_1 + M_2 = M$
- Possible models: Bernoulli distribution with probability  $p$ .
- Given one model, the probability of generating such samples:  
$$L(p) = p^M (1-p)^{N-M}$$
- Maximize the probability  $L(p)$

## Example continued

- Maximize the probability

$$L(p) = p^M (1-p)^{N-M}$$



We often maximize the **logarithm** of the probability (usually easier to maximize) of generating such samples

$$l(p) = \log L(p) = M * \log p + (N-M) * \log (1-p)$$

## Example continued

We often maximize the **logarithm** of the probability (usually easier to maximize) of generating such samples

$$l(p) = \log L(p) = M * \log p + (N-M) * \log (1-p)$$

- For continuous function  $l(p)$ 
  - Find the point  $l'(p) = 0$  and  $l''(p) < 0$
  - Before we discuss optimization, all exercises with continuous cases just require you to show the point  $l'(p) = 0$ .

$$\bullet l'(p) = \frac{M}{p} - \frac{N-M}{1-p} = 0 \quad \rightarrow \hat{p} = \frac{M}{N} \text{ (sample cure rate)}$$



## Example continued

Suppose we have two coins with head probability being  $p_1$  and  $p_2$ , respectively.

$N_1$  experiments with  $M_1$  heads using coin 1 ;  $\hat{P}_1 = \frac{M_1}{N_1}$   
 $N_2$  experiments with  $M_2$  heads using coin 2.  $\hat{P}_2 = \frac{M_2}{N_2}$

Given  $p$ , the probability of generating such samples (likelihood function):

$$L(p) = p_1^{M_1} (1-p_1)^{N_1-M_1} \underbrace{(p+0.3)^{M_2}}_{P_2} \underbrace{(0.7-p)^{N_2-M_2}}_{1-P_2}$$

Loglikelihood:

$$l(p) = M_1 \log p + (N_1 - M_1) \log(1-p) + M_2 \log(p+0.3) + (N_2 - M_2) \log(0.7-p)$$

$$l'(p) = \frac{M_1}{p} + \frac{N_1 - M_1}{1-p} + \frac{M_2}{p+0.3} + \frac{N_2 - M_2}{0.7-p} \quad \rightarrow \hat{p}$$



## Example: uniform

Let  $X$  be a Uniform random variable on the interval  $[0, \theta]$

$$f(x; \theta) = \begin{cases} \frac{1}{\theta}, & \text{for } 0 \leq x \leq \theta, \\ 0, & \text{otherwise,} \end{cases} = \frac{1}{\theta} \mathbf{1}_{\{0 \leq x \leq \theta\}}$$

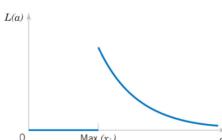
indicator function  $\mathbf{1}_A(x)$

$$\mathbf{1}_A(x) = \begin{cases} 1, & \text{if } x \in A, \\ 0, & \text{otherwise} \end{cases}$$

The likelihood function of a random sample of size  $n$  is:

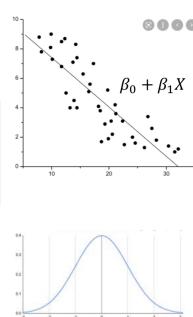
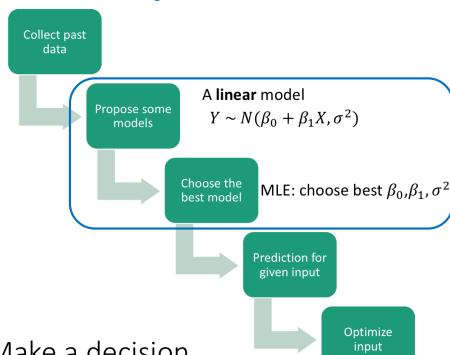
$$L(\theta) = \prod_{i=1}^n f(x_i; \theta) = \frac{1}{\theta^n} \prod_{i=1}^n \mathbf{1}_{\{0 \leq x_i \leq \theta\}} = \begin{cases} \frac{1}{\theta^n}, & \text{if } \theta \geq \max\{x_1, x_2, \dots, x_n\} \\ 0, & \text{if } \theta < \max\{x_1, x_2, \dots, x_n\} \end{cases}$$

$$\hat{\theta} = \max\{x_1, x_2, \dots, x_n\}$$



Calculus methods don't work here because  $L(\theta)$  is maximized at the **discontinuity**. Clearly,  $\theta$  cannot be smaller than  $\max(x_i)$ , thus the MLE is  $\max\{x_1, x_2, \dots, x_n\}$ .

## 3. linear regression



- PDF for normal  $\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$ .
- Samples:  $(X_1, Y_1), \dots, (X_N, Y_N)$
- For the model with  $\beta_0, \beta_1, \sigma^2$ , the likelihood is

$$\frac{1}{(\sqrt{2\pi})^n \sigma^n} \exp\left[-\frac{1}{2} \frac{\sum_i (Y_i - \beta_1 X_i - \beta_0)^2}{\sigma^2}\right]$$

- Given  $\sigma^2$ , to maximize the likelihood, we only need to minimize

$$\sum_i (Y_i - \beta_1 X_i - \beta_0)^2$$

- Taking derivative over  $\beta_0$  and  $\beta_1$ , we have

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) X_i = 0$$

$$\sum_i (Y_i - \beta_1 X_i - \beta_0) = 0$$

MLE:  $\left\{ \begin{array}{l} \widehat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} \\ \widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{X} \end{array} \right.$

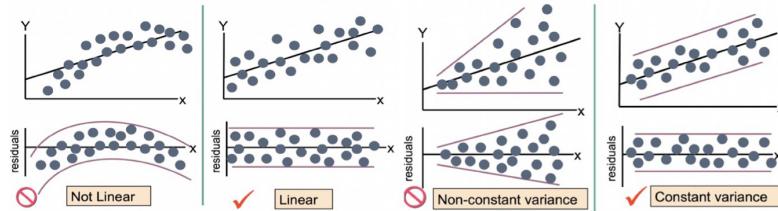
## Residual Analysis: check assumptions

Linear regression assumes that...

1. The relationship between  $X$  and  $Y$  is **linear**
2. The variance of  $Y - \beta_0 - \beta_1 X$  at every value of  $X$  is the **same** (homogeneity)

$$\text{Residual: } e_i = Y_i - \widehat{\beta}_0 - \widehat{\beta}_1 X_i$$

MLE:  $\left\{ \begin{array}{l} \widehat{\beta}_1 = \frac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2} \\ \widehat{\beta}_0 = \bar{Y} - \widehat{\beta}_1 \bar{X} \end{array} \right.$



## 4. Interval estimation

### Recap: interval estimation

Whether a drug can cure a disease:  $\hat{p} = \frac{\sum_i X_i}{n}$

- Drug 1:  $\hat{p}_1 = 90\%$ . **10 experiments**.
- Drug 2:  $\hat{p}_2 = 80\%$ . **10000 experiments**.

Which drug do you think is more effective?

With **more** data, we believe the estimator is **closer** to the true parameter.

### Main technique: Central limit theorem

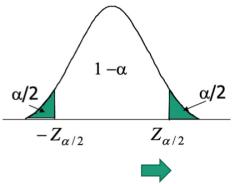


No matter what the true distribution is, the **sample mean** will be very close to the **normal distribution**, as long as the sample size is **large**.

# Interval Estimation for mean

- We have data  $X_1, X_2, \dots, X_n$  that are sampled from some distribution
- Their mean is  $\mu$ , which we want to estimate
- We can easily give a point estimate:  $\bar{X}$  (sample mean)
- How to get an interval estimate??
  - Use Central Limit Theorem!

## Interval Estimation - example



$$Z = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0,1)$$

$$\rightarrow P\left(-z_{\alpha/2} \leq \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \leq z_{\alpha/2}\right) = 1 - \alpha$$

$$\rightarrow P\left(\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}\right) = 1 - \alpha$$

With probability  $1 - \alpha$ ,  $\mu$  is within the interval  $[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$

You can also write as  $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

## Some observations

**1- $\alpha$  Confidence interval:**

$$[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}}]$$

We typically call  $1 - \alpha$  as the **confidence level**.

The length of the confidence interval is affected by several factors

- As the sample size  **$n$  increases**, the length of CI **decreases**
- As the variance  **$\sigma^2$  increases**, the length of CI **increases**
- As the confidence level increases ( **$\alpha$  decreases**), the length of CI **increases**.

## Interval Estimation for proportion (mean)

Mean:  $p$   
Variance:  $p(1-p)$   
 $\hat{p} = \bar{X}$  (sample cure rate)

$$Z = \frac{\sqrt{n}(\hat{p} - p)}{\sqrt{p(1-p)}} \sim N(0,1)$$

Recall:  
 $Z = \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \sim N(0,1)$   
 Interval is  $\bar{X} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$

$$\rightarrow P\left(-z_{\alpha/2} \leq \frac{\sqrt{n}(\hat{p} - p)}{\sqrt{p(1-p)}} \leq z_{\alpha/2}\right) = 1 - \alpha$$

$$\rightarrow P\left(\hat{p} - z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}\right) = 1 - \alpha$$

With probability  $1 - \alpha$ , use the approximate interval

$$[\hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}] \quad \hat{p} \pm z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

## §2 Exercises

### 例: Additional Exercise 1:

- Suppose there are **100** white or black balls in the bag, 90 of them are one color and 10 of them are another color. I randomly pick 1 ball out of the bag and find out that it's white. So, by the maximum likelihood estimate, which is the majority color in the bag (what ~~color~~<sup>is</sup> are the 90 balls)?

case 1: 90 white 10 black  $L(\text{case 1}) = \frac{90}{100} = \text{white}$

case 2: 90 black 10 white  $L(\text{case 2}) = \frac{10}{100}$



100 balls

90 in one color  
10 in another

Draw one sample randomly



### 例: Additional Exercise 2:

- Suppose we ask people to select their favorite ice cream brand: Haagen-Dazs (type 1), Walls (type 2), Ben & Jerry's (type 3).
- We can view the large population contains 3 different types of individuals, and let  $\theta_1, \theta_2, \theta_3$  denote the proportion of individuals of type each type. Clearly,  $0 \leq \theta_i \leq 1$  and  $\theta_1 + \theta_2 + \theta_3 = 1$ .  $\theta_1, \theta_2, \theta_3$  unknown.
- Suppose we ask 100 students in CUHK(SZ), and 45 choose HD, 30 choose Walls, and 25 choose BJ. Find the MLE for  $\theta_1, \theta_2, \theta_3$ .

$$L(\theta_1, \theta_2, \theta_3) = \theta_1^{45} \cdot \theta_2^{30} \cdot \theta_3^{25} \quad l(\theta_1, \theta_2, \theta_3) = 45\ln\theta_1 + 30\ln\theta_2 + 25\ln\theta_3$$

$$\begin{cases} \frac{45}{\theta_1} - \frac{25}{1-\theta_1-\theta_2} = 0 \\ \frac{30}{\theta_2} - \frac{25}{1-\theta_1-\theta_2} = 0 \end{cases} \rightarrow \hat{\theta}_1 = 0.45 \quad \hat{\theta}_2 = 0.3 \quad \hat{\theta}_3 = 0.25$$



### 例: Additional Exercise 3:

- A geologist studied the composition of rocks in the shoreline area of Lake Michigan. He randomly selected 100 samples from the area, each containing 10 stones, and recorded the number of limestone stones in each sample. The geologist's data are as follows:

number of limestone in a sample	0	1	2	3	4	5	6	7	8	9	10
number of samples	0	1	6	7	23	26	21	12	3	1	0



### Additional Exercise 3:

number of limestone in a sample	0	1	2	3	4	5	6	7	8	9	10
number of samples	0	1	6	7	23	26	21	12	3	1	0

$$X_i = \# \text{ limestone in sample } i$$

- Assume that these 100 observations are independent of each other. Find the maximum likelihood estimate (MLE) of the proportion  $p$  of limestone in the stones in this area.

$$\begin{aligned} L(p) &= \prod_{i=1}^{100} \left[ \binom{10}{x_i} p^{x_i} (1-p)^{10-x_i} \right] \\ &= \left[ \prod_{i=1}^{100} \left( \binom{10}{x_i} \right) \right] \cdot P^{\sum x_i} \cdot (1-p)^{1000 - \sum x_i} \\ &\quad l(p) = \log \left[ \prod_{i=1}^{100} \left( \binom{10}{x_i} \right) \right] + \sum x_i \cdot \log p - (1000 - \sum x_i) \log (1-p) \end{aligned}$$

$$\hat{p} = \frac{M}{1000}$$