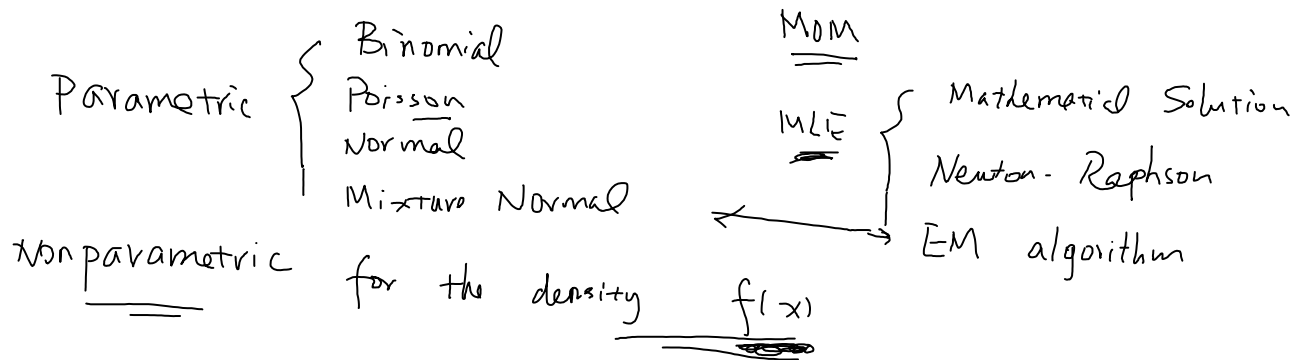


# Estimation



Kernel density estimation

Estimation is not good on the boundary.

cumulative distribution function  $F(x)$   $\rightarrow$  cdf

empirical cdf

$$\underline{F(x)} = P(X \leq x)$$

$$\underline{F_n(x)} = \frac{\# \text{ number of } X_i \leq x}{n} = \frac{1}{n} \sum 1(\underline{X_i \leq x})$$

$$\text{Let } Y_i = 1(X_i \leq x) \quad Y_i \stackrel{iid}{\sim} \text{Bernoulli}(F(x))$$

$$\underline{F_n(x)} = \frac{1}{n} \sum Y_i = \frac{\text{Bin}(n, F(x))}{n}$$

$$\mathbb{E} \underline{F_n(x)} = \frac{n F(x)}{n} = \underline{F(x)}$$

Theorem: DKW Inequality:

$$P\left(\sqrt{n} \cdot \sup_t |\underline{F_n(t)} - \underline{F(t)}| > z\right) \leq 2 \exp(-2z^2)$$

$$z = \log n$$

$$P\left(\sup_t |\underline{F_n(t)} - \underline{F(t)}| > \frac{\log n}{\sqrt{n}}\right) \rightarrow 0$$

Application ; ① Test the model assumption. Kolmogorov-Smirnov.

② Multiple testing : False Discovery Rate (FDR)

Hypothesis Testing : Yes/No question

H<sub>0</sub> : Model assumption

Ronald A. Fisher

Elements of Hypothesis Testing

Hypotheses : H<sub>0</sub> null hypothesis (conventional knowledge)

H<sub>a</sub> : Alternative hypothesis. (we want to demonstrate)

Example Clinical Trial  $\bar{r} = P(\text{Success of Surgical Procedure})$

H<sub>0</sub> :  $\bar{r} \leq 0.2$

H<sub>a</sub> :  $\bar{r} > 0.2$

Hypothesis should never depends on the data

Composite Hypothesis

Example :  $\bar{r} = \text{number of West Nile Virus}$

H<sub>0</sub> :  $\bar{r} = 55$       H<sub>a</sub> :  $\bar{r} \neq 55$

Single Hypothesis

Test Statistic : a function of the data

Rejection Region : sets of values that a researcher wants to

Rejection Region: sets of values that a researcher wants to reject  $H_0$ .  $R$

Types Of Errors:

	Reject $H_0$	Fail-to-reject $H_0$
$H_0$	Type I error	✓
$H_a$	✓	Type II

False Discovery  
False Positive = Type I

False Negative = Type II

Type I error, " $\alpha$ " significance level of  $\alpha$

Simple:  $\alpha = P(R | H_0)$

$$\alpha = \sup_{P_0 \in H_0} P(R | P_0 \text{ is true})$$

Type II error:  $\beta$

$$\beta = P(R^c | \theta \in H_a)$$

power:  $= 1 - \beta = P(R | \theta \in H_a)$

Control. Type I error  $\alpha$  at 0.1, 0.05, 0.01,  
minimize the type II error.

Choose Rejection region  $P(R | H_0) \leq \alpha$

distribution of  $(T | H_0)$

Example 6.2.3:

### Example 6.2.3

$$H_0: \tau = 0.2$$

$$H_a: \tau = 0.5$$

$n = 25$  subjects.  $Y$  be the number of success

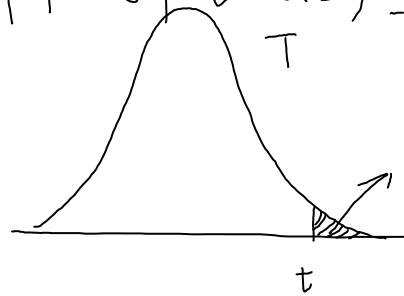
$$Y \sim \text{Bin}(25, \tau)$$

$$R = \{Y > 8\}. \quad \text{Type I and II?}$$

$$\text{Type I} = P(R | H_0) = P(Y > 8 | \tau = 0.2) = 0.047$$

$$\text{Type II} = P(R^c | H_a) = P(Y \leq 8 | \tau = 0.5) = 0.054$$

P-value.  $P(\quad | H_0)$



P-value is the smallest " $\alpha$ " that we can reject  $H_0$

If  $\alpha < \text{P-value}$ , fail-to-reject  $H_0$

$\alpha > \text{P-value}$ , reject  $H_0$

Ex  $n = 25$ ,  $Y = 10$ ,  $\alpha = 0.01$

$$P(Y \geq 10 | \tau = 0.2) = 0.005 < \alpha \rightarrow \text{reject } H_0$$

① Hypotheses ② " $\alpha$ "

③ data  $\rightarrow$  Statistic

④ Find the distribution of  $T$  under  $H_0$

⑤ Rejection region.  $\{T, \text{rejecting}\} \rightarrow$  make a decision

④' p-value,  $\{p\text{-value}, \alpha\} \rightarrow \text{make a decision}$

### 6.3 Generalized Likelihood Ratio Test (GLRT)

$$\underline{H}_0: \theta \in \Theta_0$$

$$H_a: \theta \in \Theta_1$$

$$\text{Test Statistic: } LR = \frac{\max_{\theta \in \Theta_1} L(\theta)}{\max_{\theta \in \Theta_0} L(\theta)}$$

$$\Lambda = \frac{\max_{\theta \in \Theta_0 \cup \Theta_1} L(\theta)}{\max_{\theta \in \Theta_0} L(\theta)}$$

$$\text{Reject region: } \{ \Lambda > k^* \}$$

Test for one population mean.

$$\underline{H}_0: \mu = 17.60 = \mu_0$$

$$H_a: \mu \neq 17.60$$

$$\underline{\underline{\sigma = 1.3}}$$

$$X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$$

$$L(\mu) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{\sum (x_i - \mu)^2}{2\sigma^2} \right\}$$

$$\underline{H}_0: \max_{\theta \in \Theta_0} L(\mu) = L(17.60) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{\sum (x_i - 17.6)^2}{2\sigma^2} \right\}$$

$$H_0 \cup H_a: \max L(\mu) = L(\bar{x}) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^n \exp \left\{ -\frac{\sum (x_i - \bar{x})^2}{2\sigma^2} \right\}$$

$$\Lambda = \frac{\exp \left\{ -\frac{\sum (x_i - \bar{x})^2}{2\sigma^2} \right\}}{\exp \left\{ -\frac{\sum (x_i - \mu_0)^2}{2\sigma^2} \right\}} = \exp \left\{ \frac{\sum (x_i - \mu_0)^2 - \sum (x_i - \bar{x})^2}{2\sigma^2} \right\}$$

$$= c \cdot \exp \left\{ \frac{n (\bar{x} - \mu_0)^2}{2 \sigma^2} \right\}$$

Reject  $H_0$  Regm  $R = \{ \Lambda > k^* \} = \left\{ c \exp \left( \frac{n (\bar{x} - \mu_0)^2}{2 \sigma^2} \right) > k^* \right\}$

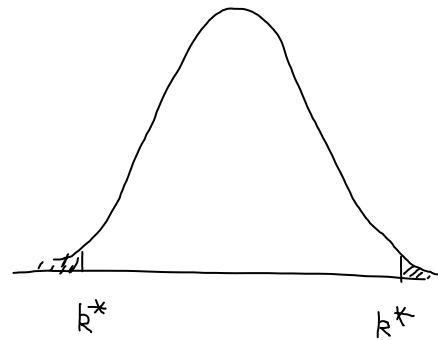
$$= \left\{ \frac{n (\bar{x} - \mu_0)^2}{2 \sigma^2} > k^* \right\} = \left\{ (\bar{x} - \mu_0)^2 > k^* \right\}$$

$$= \left\{ |\bar{x} - \mu_0| > k^* \right\} = \left\{ \frac{|\bar{x} - \mu_0|}{\sigma/\sqrt{n}} > k^* \right\}$$

Let  $Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$

Under  $H_0$ :  $Z \sim N(0, 1)$

$$\alpha = P(|Z| > k^* | H_0)$$



$$k^* = z_{\alpha/2}$$

$$R = \left\{ \frac{|\bar{x} - \mu_0|}{\sigma/\sqrt{n}} > z_{\alpha/2} \right\} = \{ |Z| > z_{\alpha/2} \}$$

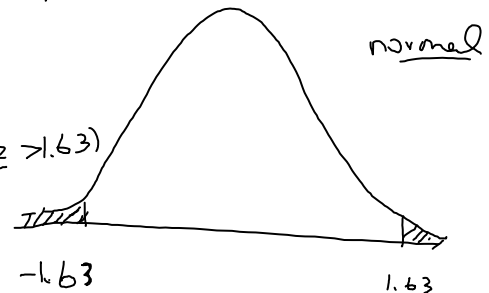
$$Z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}} = \frac{17.30 - 17.60}{1.3/\sqrt{50}} = -1.63$$

Reject  $H_0$  Regm  $\alpha = 0.10$   $R = \{ |Z| > 1.645 \}$

→ fail-to-reject  $H_0$

p-value =  $P(|Z| > 1.63 | H_0) = 2P(Z > 1.63)$

$$= 0.103$$



Some results

(1)  $H_0: \tau = \tau_0$  vs  $H_a: \tau \neq \tau_0$   $\wedge$

$$2 \log \Lambda | H_0 \sim \chi_1^2 \quad \text{when } n \rightarrow \infty$$

(2) Nested null and alternative

$$\underline{H_0}: \theta \in \underline{\underline{\Theta_{p-r}}}$$

$$H_a: \theta \in \underline{\underline{\Theta_p}}$$

$$2 \log \Lambda | H_0 \sim \chi_r^2 \quad \text{when } n \rightarrow \infty.$$

H<sub>0</sub>: rate is constant

H<sub>a</sub>: rate is seasonal

$$\tilde{Y} = (Y_1, Y_2, \dots, Y_{12})$$

Multinomial "c" category

$$\sum_{i=1}^c \pi_i = 1$$

$$(Y_1, Y_2, \dots, Y_c) \sim \text{Multinomial}(n, \underline{\underline{\pi}})$$

$$P(Y_1 = y_1, \dots, Y_c = y_c) = \frac{n!}{y_1! \dots y_c!} \prod_{i=1}^c \pi_i^{y_i}$$

$$\hat{\pi}_i = \frac{y_i}{n}$$

Let  $\theta_1, \theta_2, \dots, \theta_{12}$  be the daily rate of suicide  
 $d_1, d_2, \dots, d_{12}$  be the number of days in each month

$$\pi_i = \underline{\underline{d_i \theta_i}}$$

$$\tilde{Y} = \text{Multinomial}(n, \underline{\underline{\pi}})$$

$$\underline{H_0}: \theta_1 = \theta_2 = \dots = \theta_{12}$$

H<sub>a</sub>: At least two  $\theta$ 's are not the same

CONF

$$\frac{n}{\pi_i} \dots y_i$$

H<sub>a</sub>. At least two  $\theta$ 's are not the same

GLRT  $L(\hat{\theta}) = \frac{n!}{y_1! \dots y_k!} \prod_{i=1}^k (d_i \theta_i)^{y_i}$

H<sub>0</sub>.  $\theta = \frac{1}{365}$

$$L(\hat{\theta}) = \frac{n!}{\prod y_i!} \prod \left( \frac{d_i}{365} \right)^{y_i}$$

$1 = \sum \pi_i = \sum d_i \cdot \theta = \theta \sum d_i = \theta \cdot 365$

H<sub>0</sub> or H<sub>a</sub>.  $\max L(\hat{\theta}) = L\left(\frac{y_1}{n}, \frac{y_2}{n}, \dots, \frac{y_k}{n}\right) = \frac{n!}{\prod y_i!} \prod \left( \frac{y_i}{n} \right)^{y_i}$

$$\Lambda = \prod \frac{\left( \frac{y_i}{n} \right)^{y_i}}{\left( \frac{d_i}{365} \right)^{y_i}} = \prod \left( y_i \cdot \frac{365}{n d_i} \right)^{y_i}$$

$$2 \log \Lambda = 2 \sum_i y_i \log \left( y_i \cdot \frac{365}{n d_i} \right)$$

Under H<sub>0</sub>.  $2 \log \Lambda = \chi^2_{k-1}$

P-value:  $= P(\chi^2_{k-1} > \text{test statistic})$

Reject the H<sub>0</sub>