

$$Z_n = \frac{S_n - n\mu}{\delta \sqrt{n}}$$

$$M_{Z_n(t)} = E[\exp(t \cdot \frac{S_n - n\mu}{\delta \sqrt{n}})] \\ = E\left[\prod_{i=1}^n \exp\left(\frac{X_i - \mu}{\delta \sqrt{n}} \cdot t\right)\right]$$

$X_1 \dots X_n$  are iid

$$\Rightarrow = \prod_{i=1}^n E\left[\exp\left(\frac{X_i - \mu}{\delta \sqrt{n}} \cdot t\right)\right] \\ = \left\{E\left[\exp\left(t \cdot \frac{(X - \mu)}{\delta \sqrt{n}}\right)\right]\right\}^n \\ = \left\{E\left[1 + \frac{t}{\delta \sqrt{n}}(X - \mu) + \frac{t^2}{2\delta^2 n}(X - \mu)^2 + \text{higher order}\right]\right\}^n$$

$$\lim_{n \rightarrow \infty} M_{Z_n(t)} = \lim_{n \rightarrow \infty} \left(1 + \frac{t^2}{2n} + R(t)\right)^n \\ = e^{t^2/2}$$

MGF of  $N(0, 1)$

$$\hat{\mu} = \frac{1}{n} \sum X_i$$

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

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

$$\text{Var } S^2 = \frac{E[(X - \mu)^4]}{n} - \frac{\delta^4(n-3)}{n(n-1)}$$

Statistic for mean and var

$$\hat{F}(x) = \begin{cases} 0 & \text{if } x < X_{(1)} \\ k/n & \text{if } X_{(k)} \leq x \leq X_{(k+1)} \\ 1 & \text{if } x \geq X_{(n)} \end{cases}$$

quartiles vs  
order statistic

$$\hat{C}_q = X_{(k)} \quad \text{if } \frac{k-1}{n} < q < \frac{k}{n}$$

$$\hat{C}_q = \frac{1}{2} (X_{(q \cdot n)} + X_{(q \cdot n + 1)})$$

$$f_{X_{(1)}, \dots, X_{(n)}}(x_{(1)}, \dots, x_{(n)}) = n! f_{X_{(1)}}(x_{(1)}) \cdots f_{X_{(n)}}(x_{(n)})$$

$$f_{X_{(k)},(u)}(u) = k \binom{n}{k} F_{X_{(k)}}(u)^{k-1} (1 - F_{X_{(k)}}(u))^{n-k} f_{X_{(k)}}(u) \\ = n \binom{n-1}{k-1} \cdots$$

$$\hat{C}_q \stackrel{d}{=} N(C_q, \frac{q(1-q)}{nf(q)^2})$$

$X_1 \dots X_n$ , iid,  $X \sim N(0, 1)$  | Statistic dist

$$\textcircled{1} U = \sum_{i=1}^n X_i^2 \sim \chi_n^2 \quad M_{U(t)} = (1-2t)^{-\frac{n}{2}}$$

$$\frac{(n-1)S^2}{\delta^2} = \frac{\sum X_i^2 - n\bar{X}^2}{\delta^2} = \sum_{i=1}^n \frac{(X_i - \mu)^2}{\delta^2} - \frac{(\bar{X} - \mu)^2}{\delta^2 n} \sim \chi_{n-1}^2$$

$$\textcircled{2} U = \sum_{i=1}^n X_i, \quad U \sim N(0, \frac{1}{n})$$

Central limit theory

$$X_1 \dots X_n, \text{iid}, \quad X \sim N(\mu, \delta^2) \quad | \text{ Statistic dist}$$

$$\textcircled{1} \text{ for true } \delta^2 \quad \frac{\bar{X} - \mu}{\delta / \sqrt{n}} \sim N(0, 1)$$

\textcircled{2} don't know  $\delta^2$ , approx  $\delta^2$  by  $s^2$

$$\frac{\bar{X} - \mu}{s / \sqrt{n}} \stackrel{d}{\sim} t_{n-1}; \quad \frac{Z}{\sqrt{V}} \stackrel{d}{=} t_{V-1}, \quad V \sim \chi_V^2$$

$$\text{for prediction} \quad \frac{(\bar{X} - \mu) / \delta / \sqrt{n}}{\sqrt{(n-1)s^2 / n - 1}} = \frac{\bar{X} - \mu}{s / \sqrt{n}}$$

$$\frac{X^* - \bar{X}}{s / \sqrt{1 + \frac{1}{n}}} \stackrel{d}{\sim} t_{n-1}, \quad \bar{X} = \bar{x} \pm t_{n-1} \cdot s \cdot \sqrt{1 + \frac{1}{n}}$$

$$\textcircled{3} U \stackrel{d}{=} \chi_m^2, \quad V \stackrel{d}{=} \chi_n^2$$

$$\frac{U/m}{V/n} \stackrel{d}{=} F_{m,n}$$

$$\textcircled{4} 1' \quad \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\delta \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{d}{\sim} N(0, 1)$$

$$2' \quad \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \stackrel{d}{\sim} t_{n_1+n_2-2}$$

assume equal var

$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$$

$$3' \quad \text{know ratio of } \delta_1 = k \delta_2$$

$$\frac{[(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)] / \sqrt{\frac{\delta_1^2}{n_1} + \frac{\delta_2^2}{n_2}}}{\sqrt{U/n_1+n_2-2}} \stackrel{d}{\sim} t_{n_1+n_2-2}$$

$$U = \frac{(n_1-1)S_1^2}{\delta_1^2} + \frac{(n_2-1)S_2^2}{\delta_2^2} \stackrel{d}{=} \chi_{n_1+n_2-2}^2$$

$$E(x) = \mathbb{E} \bar{X} = \int x \cdot f_{\bar{X}}(x) dx \quad | \text{ method of moment}$$

$$= \theta^2 + b$$

$$\bar{\theta}^2 = \bar{X} - b \Rightarrow \bar{\theta} = \pm \sqrt{\bar{X} - b}$$

Unbiasedness:  $E(T_{(n)}) = \theta$

Consistency:  $E(\bar{T}_{(n)}) \rightarrow \theta$

$$\text{Var}(\bar{T}_{(n)}) \rightarrow 0$$

Efficiency:  $\text{Var}(T_{(n)}) \geq \frac{1}{12\theta} \cdot \text{eff}(T_{(n)}) = \frac{\text{eff}(T_{(n)})}{\text{Var}(T_{(n)})}$

$$D_1 = \frac{\partial}{\partial \theta} \log L, \quad \int \frac{\partial}{\partial \theta} \log L'(x) \cdot L(x) dx = 0$$

$$= E(D_1)$$

$$= \int \frac{\partial^2}{\partial \theta^2} \log L(x) \cdot L(x) dx + \int \frac{\partial}{\partial \theta} \log L(x) \frac{\partial}{\partial \theta} L(x) dx$$

MLE estimator

$$\begin{aligned}
 &= E(D_2) + E(D_1^2) \quad E(T) = \theta \\
 &= 0 \\
 \text{Cov}(T, D_1) &\gg E(TD_1) - E(T) \cdot E(D_1) \\
 &= \int T(x) \frac{\partial \log L(x)}{\partial \theta} L(x) dx \\
 &= 1
 \end{aligned}$$

$$\frac{\text{Cov}(T, D_1)}{\text{Var } T \cdot \text{Var } D_1} = \frac{1}{\text{Var } T \cdot \text{Var } D_1} \leq 1$$

$$\Rightarrow \text{Var } T \gtrapprox \frac{1}{\text{Var } D_1} = \frac{1}{I(\theta)}$$

$$\boxed{\text{Var } T \text{ MLE} \approx \frac{1}{I(\theta)}} \rightarrow \boxed{\begin{array}{c|c} & I(\theta) \\ \hline N(0, \theta) & n/2\theta^2 \\ N(\theta, \delta^2) & n/\delta^2 \end{array}}$$

$$I(\theta) = -E(D_2) = E(D_1^2) = \text{Var}(D_1)$$

$$\frac{L_1}{L_0} = \exp\left(\frac{1}{2}n(\bar{x} - \theta_0)^2\right) \quad \boxed{\text{log ratio test}}$$

$$\log \frac{L_1}{L_0} = \frac{1}{2}n(\bar{x} - \theta_0)^2 = \frac{1}{2}\frac{(\bar{x} - \theta_0)^2}{(1/\sqrt{n})} \stackrel{d}{=} \frac{1}{2}X_1^2$$

$$\text{Reject } H_0 \text{ if } \log \frac{L_1}{L_0} > \frac{1}{2} C_{0.95} (X_1^2)$$

Decision making: the ultimate goal of using statistics is to make decisions.  
 The framework we use for this, is hypothesis testing

$$\sum_{i=1}^n N(0, 1)^2 \sim \chi_n^2$$

$$\frac{N(\mu, 1)}{\sqrt{\frac{\chi_{n-1}^2}{n-1}}} \sim t_{n-1} \Leftrightarrow \frac{\bar{X}-\mu}{S/\sqrt{n}}$$

$$\frac{\chi_{n-1}^2 / (n-1)}{\chi_{n-1}^2 / (n-1)} \sim F_{n-1, n-1}$$

MM & MLE

$$MM: E(X) = \mu = \frac{1}{n} \sum X_i$$

$$E(X^2) = \delta^2 + \mu^2 = \frac{1}{n} \sum X_i^2$$

$$\Rightarrow \mu = \bar{x}, \delta^2 = \bar{x}^2 - \bar{x}$$

$$MLE: \frac{\partial L}{\partial \mu} = 0, \frac{\partial L}{\partial \delta^2} = 0$$

$$\Rightarrow \mu \approx \bar{x}, \delta^2 \approx \bar{x}^2$$

$$MVB = \text{Var } T_{MLE} = \frac{1}{I(\theta)} = -\frac{1}{E(D_2)}$$

likelihood ratio test

test statistic

$L(\theta_1)$  or  $\log L(\theta_1) - \log L(\theta_0)$  need to assume a dist first

Distribution free method  
derive test statistic without assume dist

$$(X^2 \text{ test}): \chi_{k-p-1}^2 \quad \boxed{\frac{\sum_{i=1}^k (\text{obs} - np_i)^2}{n p_i \exp}} = U \text{ test stat}$$

$\chi_{(n-1)(c-1)}^2$  independent test

median inf test: Sign test:  $H_0: m = m_0$

$$T.S.: Z = \text{freq}(X \leq m_0)$$

$$Z|H_0 \sim N(0, 1)$$

$$\text{Pvalue} = 2 \times P(Z \geq \# \text{obs})$$

(MC) Inverse-transform method

$$\begin{aligned} X &= F^{-1}(U) \sim F(x) \\ \Rightarrow P(X < x) &= P(F^{-1}(U) \leq x) \\ &= P(U \leq F(x)) \\ &= F(x) \end{aligned}$$

First: work out  $F(x) = u$  (sometimes impossible)

$$\text{Second: } X = F^{-1}(u)$$

$$\text{ex: } F_x(x) = 1 - e^{-x} = u \\ x = -\log(1-u) \quad u \in (0, 1)$$

The AR algorithm

① Generate  $Y$  and  $U$ .  $U \in [0, 1]$

② if  $U < \frac{f(Y)}{mg(Y)}$  accept, or  $X$

③ generate a new  $Y$  (random walk)

$$M = \max_y \frac{f(y)}{g(y)}$$

prob of being accepted

$$P(U \leq \frac{f(y)}{mg(y)})$$

$$\begin{aligned} &= \int P(U \leq \frac{f(y)}{mg(y)}) \cdot g(y) dy \\ &= \int \frac{f(y)}{mg(y)} \cdot g(y) dy \\ &= 1/M \end{aligned}$$

wilcoxon rank test,  $H_0: X_{\text{median}} = X^*$   
 $\text{rank}(x < x^*) - \text{rank}(x > x^*)$

$$\sim N(0, \frac{n(n+1)(2n+1)}{6})$$

median equation test,  $H_0: \text{Median } X = Y$

Contingency table,  $\chi^2$  test

$$H_0: M_x = M_y \Rightarrow m^* \text{ overall}$$

$$x \leq m^* > m^* \quad n_i \text{ exp to be equal}$$

$$Y \quad n_i$$

(Sufficient statistics)

Fisher-Neyman factorisation theorem,

$$L(\vec{x}) = K_1 \left( \frac{t(\vec{x})}{S.S.} \mid \theta \right) \cdot K_2(\vec{x})$$

Exponential:  $f(x|\theta) = \exp(p(\theta) \cdot k(x) + h(x) + g(\theta))$

$$Y = \sum K(x_i) \cdot S.S.$$

$$EY = -n \frac{q'(\theta)}{p'(\theta)} \quad \text{if } EY = n \theta \Rightarrow Y_n \rightarrow MVB$$

$$\text{Var } \hat{\theta}$$

$$\text{Var } Y = \frac{n}{p'(\theta)^2} [p''(\theta)q'(\theta) - q''(\theta)p'(\theta)]$$

Completeness: A S.S. is complete if it comes from a family of dist such that if for all  $\theta$ ,  $E(U|y) = 0 \Rightarrow U = 0$

(Joint test equality of all means.)

Assume:  $k$  group:  $\bar{X}_k, S_k^2, \frac{S_k^2 \cdot (n-k)}{B}$  one way ( $\bar{X}_{kj} = \mu_j + \alpha_j$ )  $\sum \alpha_j = 0$   $\sim \chi_{n-k}^2$

$$W = \sum_{i=1}^{n-1} (n-i) S_k^2 \sim \delta \cdot I_{N-k} \quad (\text{always})$$

$$(S_B^2 = \frac{1}{k-1} \sum (\bar{X}_i - \bar{\bar{X}})^2) \quad H_0: \bar{\bar{X}} \text{ for all group}$$

$$B = n(k-1) S_B^2 \sim \delta^2 \chi_{k-1}^2 \quad \text{if } H_0 \text{ is true}$$

$$T = W + B$$

one way ANOVA:  $\frac{B/k-1}{W/n-k} \sim F_{k-1, n-k}$ . Reject if  $P$  is  $< 0.05$

rank test for (not for me) comparative inf

$$W_x \quad (\text{Sum of ranks from } X)$$

$$\sim N\left(\frac{1}{2}n(n+1), \frac{1}{12}n_1n_2(n_1+n_2+1)\right)$$

(to compare if  $X, Y$  are different)

$$H_0: X = Y$$

Summary of distribution

Drank test more powerful

$$\begin{matrix} x_1 & x_2 & \dots & x_m \\ \text{Rank} & 1 & 2 & \dots & n \end{matrix}$$

ANOVA test, for size  $\alpha$

Bonferroni method

$$\alpha_E = 1 - (1 - \alpha)^m$$

Tukey method

$$\bar{x}_i - \bar{x}_j \pm C_{0.95}(Q_{k, N-k}) \frac{S}{\sqrt{n}}$$

Two way ANOVA & Additive model

$$Y_{ij} = \mu + \alpha_i + \beta_j + \epsilon_{ij}$$

$$\text{where } \sum_{i=1}^a \alpha_i = 0 \quad \sum \beta_j = 0$$

$$T = R + C + I + W = B + W$$

Between the row

$$B = R + C + I$$

if we only one obs,  $\Rightarrow C=1$   
if I not exist

$$R = bc \sum_{\text{Row}} \sum_{\text{Col}} \frac{a}{a_i} (X_{i..} - \bar{X}_{..})^2$$

$$C = \sum_{\text{Col}} \frac{b}{b_j} (X_{.j} - \bar{X}_{..})^2$$

$$I = \sum_{\text{Row}} \sum_{\text{Col}} \frac{ab}{a_i b_j} (X_{ij..} - \bar{X}_{i..} - \bar{X}_{.j} + \bar{X}_{..})^2$$

$$T = R + C + W$$

$$B = R + C$$

$$\text{and } W = \bar{X}_{..} \Rightarrow \bar{X}$$

$$W = \sum_{\text{Row}} \sum_{\text{Col}} \sum_{\text{Inte}} (X_{ijk..} - \bar{X}_{ij..} - \bar{X}_{.jk} + \bar{X}_{..})^2$$

$$R/a-1 \quad C/b-1 \quad I/ab-a-b+1$$

$$W/ab(c-1) \quad \text{test Col} \quad W/ab(c-1)$$

$$\text{test Row} \quad \text{test Inte}$$

$$\alpha = \alpha_0 - \beta \bar{X}$$

Linear Regress

$$\text{model: } y = \alpha_0 + \beta(x - \bar{x}) = \alpha_0 + \beta u = \hat{\alpha}_0 + \hat{\beta} x \quad (\text{LS linear})$$

$$\hat{\alpha}_0 = \bar{y}, \quad \hat{\alpha}_0, \hat{\beta} \text{ for } \hat{\alpha}_0 \text{ and } \hat{\beta}$$

$$\hat{\beta} = \frac{\sum u_i y_i}{\sum u_i^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

$$\text{cov}(\hat{\alpha}_0, \hat{\beta}) = \alpha_0^2 \left( \frac{1}{n} \sum y_i^2 - \frac{1}{n^2} \sum u_i^2 \sum y_i^2 \right)^{-1/2}$$

$$\text{cov}(\hat{\alpha}_0, \hat{\beta}) = \frac{1}{n} \sum u_i^2 \sum u_i \text{cov}(y_i, y_i) = 0$$

For general linear model (easier to solve)

$$y = A\theta + e \quad (\mu=0) \quad y = \mu + A'\theta + e, \quad e = y_i - \hat{y}_i$$

$$A^T A \hat{\theta} = A^T y \Rightarrow \hat{\theta} = (A^T A)^{-1} A^T y$$

$$\hat{\theta} \stackrel{d}{=} N_p(\theta, \sigma^2 (A^T A)^{-1})$$

$$S^2 = \frac{1}{n-p} (y^T y - \hat{\theta}^T A^T y) \text{ for } S^2$$

$$\text{Var}(\alpha + 2\beta) = \text{Var}(\alpha) + 4\text{Var}(\beta) + 4\text{cov}(\alpha, \beta)$$

$$E(y_i) = A \cdot \theta, \quad \text{Var}(y_i) = \text{Var}(A \cdot \theta), \text{ ex.}$$

$$T = R + C + I + W = B + W$$

$$\text{set if } c \geq 1 \text{ and } B=R+C$$

$$\text{if } c \leq 1, \text{ we use } T = R + C + W$$

$$\text{Then, } W = I + W$$

$$24-6$$

$$I/ab-a-b+1 \sim F_{m,n}$$

$$W/ab(c-1) \sim \chi^2_{n-2}$$

$$\text{test R.C.I.} = 0$$

$$y \stackrel{d}{=} N(\alpha_0 + \beta u, \frac{\sigma^2}{n} + \frac{(x_0 - \bar{x})^2 \cdot \sigma^2}{\sum u_i^2})$$

$$= R^2 \sum (y_i - \bar{y})^2 + S^2(n-2)$$

$$\sum (y_i - \bar{y})^2 = \frac{1}{n} \sum u_i^2 + S^2(n-2)$$

$$\text{ANOVA test } \beta=0: \frac{B/1}{W/n-2} \sim F_{1, n-2}$$

$$\frac{\hat{\alpha}_0 - \alpha_0}{S/\sqrt{n}} \stackrel{d}{=} t_{n-2}, \quad \frac{\hat{\beta} - \beta}{S/\sqrt{\sum u_i^2}} \stackrel{d}{=} t_{n-2}$$

$$= \frac{R^2(n-2)}{1-R^2}$$

$$\Rightarrow \frac{R\sqrt{n-2}}{\sqrt{1-R^2}} \sim t_{n-2}$$

# Point estimation ①

## Interval estimation

If you have estimator, for a value you know its distribution, Then you can infer a Interval estimation, it is a reasonable estimation.

Unbiased estimator      does Unbiased estimator must be sufficient Statistics?

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \stackrel{d}{=} N(0, 1)$$

$$\frac{\bar{X} - \mu}{\sqrt{s^2/n}} \stackrel{d}{=} t_{n-1}$$

$$MLE$$

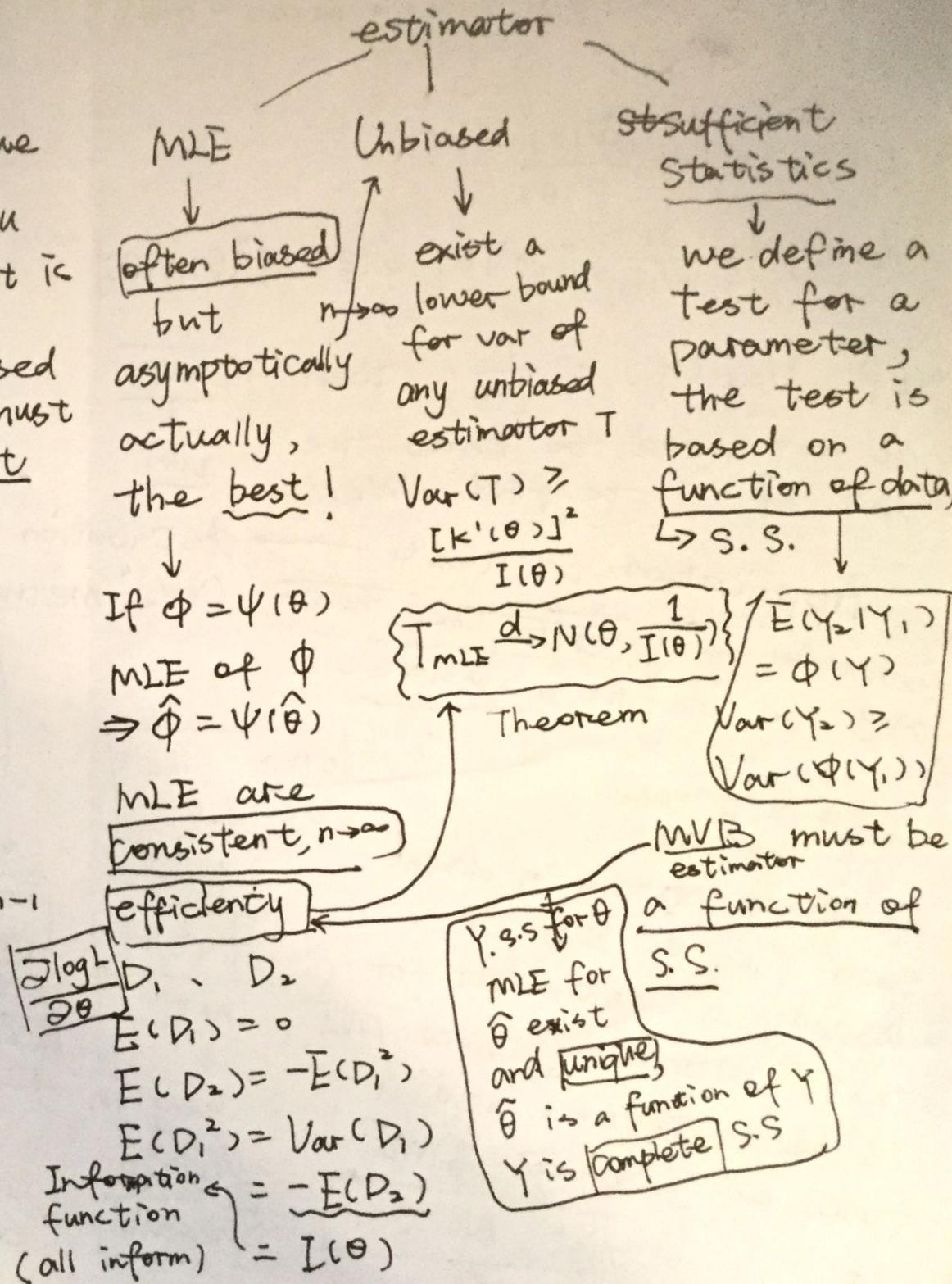
$$\begin{aligned} U &\stackrel{d}{=} \chi_m^2 \\ V &\stackrel{d}{=} \chi_n^2 \\ Z &= \frac{U/m}{V/(m+n)} \end{aligned}$$

$$\frac{(n-1) \cdot s^2}{\sigma^2} \stackrel{d}{=} \chi_{n-1}$$

$$Z \stackrel{d}{=} F_{m, n}$$

$\frac{MVB}{Var(T_n)}$  then  
 $\rightarrow$   $\frac{1}{n}$  estimator  
 $=$   $\frac{1}{n}$  efficient  
 $\sim T_n$  say +  
 $\Rightarrow$  eff. say +  
 $\text{eff. } \frac{1}{n}$  eff.  
 $\text{eff. } \frac{1}{n}$

$$\frac{S_1^2/t_1^2}{S_2^2/t_2^2} \stackrel{d}{=} F_{n_1-1, n_2-1}$$



$$f(x|\theta) = \exp(p(\theta) \cdot k(x) + H(x) + q(\theta))$$

$$Y = \sum_{i=1}^n k(X_i) \quad S.S \quad \text{?}$$

$$E(Y) = -n \frac{q'(\theta)}{p'(\theta)}$$

$$\text{Var}(Y) = \frac{n}{p'(\theta)} [P''(\theta) q'(\theta) - q''(\theta) p'(\theta)]$$

If  $\text{Var}(Y) = \frac{1}{I(\theta)}$ , then  $\rightarrow$  MVB  
do not need to know the dist of  $Y$  to find MVB.

Application in test  $\rightarrow$  "distribution free" method  
fundamental thing we define S. statistic  $\downarrow$  (Central limit theorem)

• MLE case if the test is about param in distr.  
such func of S.S if usually a estimator for that parameter, Since MLE is the best, and also a func of S.S  $\Rightarrow$  MLE

dist of the

dist of that  
key in  
define threshold

the dist. some property about to test whether param of a dist not to test a  $\rightarrow$

case 1 ① test  $H_0$ : a uniform dist

case 1 ② test  $H_0$ : a poisson dist

$J_m^2$  test  $U = \sum_{i=1}^k \frac{(f_i - np_i)^2}{np_i} \xrightarrow{d} \chi^2_{k-p-1}$

case 2 ③ test independency  $\Gamma \times C$  contingency table  $\chi^2_{(r-1)(c-1)}$

case 3

① sign test (about the value of median)  
 $a' Z = \text{freq}(X < m)$   
 $Z \sim \text{bin}(n, 0.5)$

②  $b' Z \sim N\left(\frac{n}{2}, \frac{n}{4}\right)$   
 based on the confidence interval of  $Z$

Equality of medians  
 $\alpha'$  test  
 comparative inference

concerned with  $\alpha, \beta$

$W_1 \stackrel{d}{=} N\left(\frac{1}{2}n_1(n_1+n_2+1), \frac{1}{12}n_1n_2(n_1+n_2+1)\right)$

$$W = \sum_{i=1}^k (n_i - 1) S_i^2, \quad S_i^2 = \frac{1}{n_i-1} \sum (x_i - \bar{x}_i)^2 \Rightarrow W = \sum_{i=1}^k \cdot \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$

$$B = n(k-1) S_B^2, \quad S_B^2 = \frac{1}{k-1} \sum (\bar{x}_i - \bar{\bar{x}})^2 \Rightarrow B = \sum_{i=1}^k n_i (\bar{x}_i - \bar{\bar{x}})^2$$

estimator for  $S^2 = \frac{W}{(n-1) \cdot k}$

~~test  $\mu_i = \mu_j$~~

$$t = \frac{\bar{\mu}_i - \bar{\mu}_j}{\sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \text{ compare } t \sim t_{(n-1) \cdot k}$$

~~test  $\mu_0 = \mu_1 \dots \mu_k$~~

$$\text{Test stat} = \frac{B/k-1}{W/N+k} \sim F_{k-1, N-k}$$

If  $E(u(z)) = 0$

$$\text{proof} \int_{-\infty}^{+\infty} u(z) \cdot f(z) dz = 0$$

$$\Rightarrow u(z) = 0$$

$$\text{if } f(z) = 0$$

$$\Rightarrow \theta \cdot \int_z u(z) = 0$$

if  $u(z)$  integrate to 0

then,  $u(z)$  can be non-zero.

completeness

<del>test <math>\mu_i = \mu_j</math></del>	<del>test <math>\mu_i</math> CI</del>	<del><math>\chi^2</math> test</del>	<del><math>\chi^2_{k-p-1}</math> test</del>	<del><math>\sum_{i=1}^k \frac{(obs - np_i)^2}{np_i}</math></del>
$t = \frac{\bar{\mu}_i - \bar{\mu}_j}{\sqrt{\frac{1}{n_i} + \frac{1}{n_j}}} \sim t_{(n-1) \cdot k}$	$\hat{\mu}_i \pm t_{0.95} \sqrt{\frac{s^2}{n_i}}$	$\chi^2_{(n-1) \cdot (k-1)}$	$\chi^2_{(r-1) \cdot (c-1)}$	independency test contingency table
$\sim \mathcal{N}(0, 1)$	$t_{0.975} \sqrt{\frac{s^2}{n_i}}$	$\chi^2$	$\chi^2$	Equality of median con. table
		$\text{rank test}$		Equality of median
				$W_i \stackrel{d}{=} N(\frac{1}{2}(n_1+n_2+1)n_1, \frac{1}{12}n_1n_2(n_1+n_2+1))$

~~Sign test~~

$$z = \text{freq } (x \leq m)$$

$$\text{Under } H_0: z \sim \text{bin}(n, 0.5)$$

$$z \in [0.95, N(\frac{n}{2}, \frac{n}{4})]$$

test whether true or not

$$\begin{aligned} P &= P(z \geq obs) \\ &< 0.05 \\ &\text{reject } H_0 \end{aligned}$$

$$L(\vec{x}) = k_1 \vec{x} \cdot \vec{1}(\theta) \cdot k_2(\vec{x}) \quad \text{S.S}$$

$$f(x|\theta) = \exp(p(\theta)) \cdot k(x) + H(x) + g(\theta)$$

$$Y = \sum_{i=1}^n k(x_i) \quad \text{S.S}$$

$$E(Y) = -n \frac{g'(\theta)}{p'(\theta)} = \text{if } E(Y) = n\theta \Rightarrow Y \rightarrow \text{MVB}$$

$$\text{Var}(Y) = \frac{n}{p'(\theta)^2} [p''(\theta)g'(\theta) - g''(\theta)p'(\theta)]$$

Seldom used

for complete S.S  $Y$ , MLE is a func of  $Y$

head to estimate parameter  
 $(r-1) \times (c-1)$

$$H_0: m=3$$

$$H_1: m>3$$

Let  $Z = \text{freq}(X > 3)$

Under  $H_0$ ,  $Z \stackrel{d}{=} \text{Bin}(77, 0.5)$

$Z = 45$  (observed)

$$\Rightarrow P = P(Z \geq 45)$$

If  $P < 0.05 \Rightarrow \text{reject } H_0$

### Sign test

### Comparative Hypo

- ① observations are independent
- ② Obs are normally dist with same var
- ③ mean differs by group

$$f_X(x) = e^{-x}$$

$$F_X(x) = \int_0^x e^{-y} dy$$

$$= -e^{-y} \Big|_0^x$$

$$= 1 - e^{-x}$$

**MC** ① Simulate observation from  $F_Y(y)$

$$\text{given } f_X(x) = \frac{1}{\pi \sqrt{y(1-y)}} \Rightarrow F_X(x) = \frac{2}{\pi} \arcsin(\sqrt{y}) \Rightarrow y = \sin^2 \frac{\pi}{2} u, u \in [0, 1]$$

② how to use ar-r algo

if  $U \leq \frac{f_X(Y)}{M \cdot F_Y(Y)}$  accept  $Y$ .  $U \in [0, 1]$ , if rejected, start again.

$$③ M = \max_x \frac{f_X(x)}{f_Y(x)} = \max_x e^{-x} \cdot \frac{\pi}{1-x} \sin x$$

$$\text{OR } h(x) = \frac{f_X(x)}{f_Y(x)} \Rightarrow h(x) = 0 \Rightarrow h(x) = M$$

④ Probability of being accepted

$$P(U \leq \frac{f_X(Y)}{M \cdot F_Y(Y)})$$

$$\Rightarrow \int P(U \leq \frac{f_X(y)}{M \cdot F_Y(y)}) \cdot f_Y(y) dy$$

$$= \int \frac{f_X(y)}{M \cdot F_Y(y)} \cdot f_Y(y) dy$$

$$= \frac{1}{M}$$

Simulate example

for each simulated  $Y$ ,  
allow a 50/50 chance of  
setting it to be  $-Y$  or  $Y$ .

Multiple comparisons problem

$$H_0: \mu_1 = \mu_2 = \dots = \mu_n$$

(5)

$$\text{and } \frac{\nu S^2}{\sigma^2} \stackrel{d}{=} \chi_{\nu}^2$$

$$\text{Set } Q_{k,v} = \frac{R_k}{S_v}$$

$\Rightarrow$  standardised rank distribution with  $k$  and  $v$  degrees of freedom.

Rewrite this as a joint hypothesis of pairwise hypothesis

- ① individually a size of 0.05
- ② joint size of 0.05

Bonferroni's correction:

$$\underline{\alpha_E} = 1 - (1 - \underline{\alpha_c})^m$$

experimental for single comparison

reasonable approximate

$$\underline{\alpha_c} = \frac{\underline{\alpha_E}}{m}$$

$\Rightarrow$  the actual  $\alpha$  would be much less than  $\underline{\alpha_E}$

Tukey method.

1. standardised range distribution

$$z_1, \dots, z_k \sim N(\mu, \sigma^2)$$

$$R_k = z_{(k)} - z_{(1)}$$

$S^2$  is an estimator of  $\sigma^2$

2. Tukey CI

$$\bar{x}_i - \bar{x}_j \pm C_{0.95}(Q_{k,N-k}) \frac{s}{\sqrt{n}}$$

Model interpretation

$$y_{ij} = \mu_j + \epsilon_i, \epsilon_i \sim N(0, \sigma^2)$$

one-way

(fit a different mean for each diff kategorie level)

Additive models

if more than one variable

Disaggregation

additive model

each level of factor

simply adds something to its expected value

textbook contrast

$$\text{let } \alpha_i = \mu_i - \mu, \beta_j = \mu_j - \mu \quad (6)$$

$$\Rightarrow Y_{ij} = \mu + \underbrace{\alpha_i}_{\text{"simply add"}} + \underbrace{\beta_j}_{\text{}} + \epsilon_{ij}$$

where  $\sum_{i=1}^k \alpha_i = 0$      $\sum_{j=1}^b \beta_j = 0$

$$y = \begin{cases} & \text{in total } k' \cdot k \text{ group} \\ k' & T = \sum (X_{ij} - \bar{X})^2 = \sum_{i=1}^k \sum_{j=1}^b (\bar{X}_i - \bar{X}_{..})^2 + \\ & \sum_{i=1}^k \sum_{j=1}^b (\bar{X}_j - \bar{X}_{..})^2 + \\ & \sum_{i=1}^k \sum_{j=1}^b (X_{ij} - \bar{X}_i - \bar{X}_j + \bar{X}_{..})^2 \end{cases}$$

$$k' = 2, \quad k = 2$$

$$F \quad M$$

$$y = \begin{bmatrix} \text{Av}_R \\ \text{Av}_C \end{bmatrix}$$

Hypothesis test

• Does 'row factor' has effect

• - - column - - -

$$\text{Total SS} = R + C + W$$

Row Col W within the cell

Interaction term

$$T = R + C + I + w \quad (B = R + C + I)$$

Hypothesis testing: F-test

Estimation: MLE, row average

$$T = \sum \sum \sum (X_{ijk} - \bar{X}_{...})^2$$

$$= bc \sum_{i=1}^k (\bar{X}_{i..} - \bar{X}_{...})^2 + ac \sum_{j=1}^b (\bar{X}_{.j} - \bar{X}_{...})^2 + c \sum_{i=1}^k \sum_{j=1}^b (\bar{X}_{ij} - \bar{X}_{i..} - \bar{X}_{.j} + \bar{X}_{...})^2$$

The problem  $E(Y|x) = \mu(x)$

Linear regression:  $E(Y|x) = \alpha + \beta x$

assume  $\text{Var}(Y_i) = \sigma^2$

estimate what the value of  $y$  would be given different value of  $X$ . (have some connect with I (previous))

$$\sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 = \min \Rightarrow \text{MLE for } \alpha \text{ and } \beta$$

$$\hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}$$

$$\hat{\alpha}_0 = \bar{y}, E(\hat{\alpha}_0) = \bar{y}, \text{Var} = \frac{\hat{\beta}^2}{n}$$

$$E(\hat{\beta}) =$$

$$E(y_i) = \{\hat{\alpha}_0 + \hat{\beta} (x - \bar{x})\}$$

$$\alpha_0 = \alpha + \beta \bar{x}$$

$$E(Y_i) = \alpha_0 + \beta u_i$$

$\hat{A}_0$  be the estimator of  $\alpha_0$ ,  $\hat{B}$  for  $\beta$ .

$$E(\hat{A}_0) = \alpha_0$$

$$\boxed{\text{Var}(\hat{A}_0) = \frac{\sigma^2}{n} \frac{\bar{x}^2}{n}}$$

$$E(\hat{B}) = \beta$$

$$\boxed{\text{Var}(\hat{B}) = \frac{\sigma^2}{\sum_{i=1}^n u_i^2}}$$

$$\text{Cov}(\hat{A}_0, \hat{B}) = 0$$

$\hat{A}$  to be the estimator of  $\alpha$

$\hat{B}$  for  $\beta$  ( $x$  as constant  
not a variable)

$$E(\hat{A}) = \alpha$$

$$\text{Var}(\hat{A}) = \left( \frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n u_i^2} \right) \sigma^2$$

$$\text{Cov}(\hat{A}, \hat{B}) = -\frac{\bar{x}}{\sum_{i=1}^n u_i^2} \sigma^2$$

Let  $\hat{m}(x)$  be the estimator  
for  $\mu(x)$

$$\Rightarrow \hat{m}(x) = \hat{A}_0 + (\bar{x} - \bar{x}) \hat{B}$$

$$u_i = x_i - \bar{x}$$

$$E(\hat{m}(x)) = \mu(x)$$

$$\text{Var}(\hat{m}(x)) = \left( \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n u_i^2} \right) \sigma^2$$

$$\hat{d}^2 = \sum_{i=1}^n (y_i - \hat{A}_0 - \hat{B} u_i)^2$$

$$\hat{D}^2 = \frac{1}{E(D^2)} (n-2) \hat{d}^2$$

$$S^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{A}_0 - \hat{B} u_i)^2$$

unbiased estimator for  $\sigma^2$

$$\Rightarrow \sum_{i=1}^n (y_i - \bar{y})^2 = \hat{B}^2 \sum_{i=1}^n (x_i - \bar{x})^2 + d^2$$

$$\hat{A}_0 \stackrel{d}{=} N(\alpha_0, \frac{\sigma^2}{n}), \quad \hat{B} \stackrel{d}{=} N(\beta, \frac{\sigma^2}{\sum u_i^2})$$

F test

$$\frac{D^2}{\sigma^2} \stackrel{d}{=} X_{n-2}^2$$

$$\Rightarrow \frac{(n-2) S^2}{\sigma^2} \stackrel{d}{=} \chi_{n-2}^2$$

$$\frac{\hat{A}_0 - \alpha_0}{S/\sqrt{n}} \stackrel{d}{=} t_{n-2}$$

$$\frac{\hat{B} - \beta}{S/\sqrt{\sum u_i^2}} \stackrel{d}{=} t_{n-2}$$

$$\hat{A} = \hat{A}_0 + \bar{x} \cdot \hat{B}$$

Intuitive, the correlation  
of  $x, y$  has to do with  
 $\beta$ , if  $\beta = 0$ ,  $\Rightarrow$  independent

Sample correlation

$$R = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

(8)

$\Rightarrow$  Test based on  $\frac{R\sqrt{n-2}}{\sqrt{1-R^2}} \stackrel{d}{=} t_{n-2}$

The General Linear model

$$\underline{Y} = A\underline{\theta} + E$$

It is assumed that

$$E(E) = 0, D(E) = \sigma^2 I$$

$$\underline{Y} \stackrel{d}{=} N_n(A\underline{\theta}, \sigma^2 I)$$

multivariate normal

$$\Rightarrow Y_i = A_i \underline{\theta} + \epsilon_i$$

Goal is to estimate  $\underline{\theta}$ , method of LS

$$\min (\underline{Y} - A\underline{\theta})^T (\underline{Y} - A\underline{\theta})$$

$$\Rightarrow \hat{\underline{\theta}} = (A^T A)^{-1} A^T \underline{Y}$$

$$D(\hat{\underline{\theta}}) = \sigma^2 (A^T A)^{-1} \quad \text{find estimator of } \sigma^2$$

$$\text{denote } d^2 = (\underline{Y} - A\hat{\underline{\theta}})^T (\underline{Y} - A\hat{\underline{\theta}}) \Rightarrow \text{error}$$

$$(\underline{Y} - A\underline{\theta})^T (\underline{Y} - A\underline{\theta}) = (\underline{Y} - A\hat{\underline{\theta}})^T (\underline{Y} - A\hat{\underline{\theta}}) + (\hat{\underline{\theta}} - \underline{\theta})^T A^T A (\hat{\underline{\theta}} - \underline{\theta})$$

$$= d^2 + (\hat{\underline{\theta}} - \underline{\theta})^T A^T A (\hat{\underline{\theta}} - \underline{\theta})$$

$$E(d^2) = E[(\underline{Y} - A\hat{\underline{\theta}})^T (\underline{Y} - A\hat{\underline{\theta}})] = (n-p)\sigma^2$$

$$\tilde{s}^2 = \frac{d^2}{n-p} \quad \text{(unbiased estimator for } \sigma^2)$$

$$\beta = \rho \cdot \frac{\sigma_y}{\sigma_x}$$

constant

So, the hypotheses

$$H_0: \rho = 0 \iff H_0: \beta = 0$$

test  $\beta = 0$  by test  $\rho$  (easier)

$$\text{Since: } \sum_{i=1}^n (Y_i - \bar{Y})^2 = \hat{\beta}^2 \sum_{i=1}^n (X_i - \bar{X})^2 + d^2$$

$$\beta = \rho \cdot \frac{\sigma_y}{\sigma_x}$$

$$\Rightarrow \sum_{i=1}^n (Y_i - \bar{Y})^2 = D^2 + R^2 \sum_{i=1}^{n-1} (Y_i - \bar{Y})^2$$

$$F = \frac{\text{regression MS}}{\text{residual MS}} = \frac{R^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}{D^2 / (n-2)} = \frac{(n-2) R^2}{1 - R^2} \stackrel{d}{=} F_{1, n-2} \quad (t_{n-2}^2 = F_{1, n-2})$$

$H_0: \beta = 0$  (under  $H_0$ )

$$\Rightarrow y^T y = d^2 + \hat{\theta}^T A^T A \hat{\theta}$$

$$A^T y = A^T (A \hat{\theta} + \hat{e}) = A^T A \hat{\theta}$$

( $A^T \hat{e} = 0$  by orthogonality)

$$\Rightarrow y^T y = d^2 + \hat{\theta}^T A^T y$$

$$\begin{aligned} d^2 &= (y - A \hat{\theta})^T (y - A \hat{\theta}) \\ &= y^T y - \hat{\theta}^T A^T y \end{aligned}$$

$$E(D^2) = (n-p) \sigma^2$$

$$\left\{ S^2 = \frac{1}{n-p} (y^T y - \hat{\theta}^T A^T y) \right\} = \frac{D^2}{\text{rank of } AA}$$

$$D^2 = S^2 (n-p)$$

$$\left\{ \frac{D^2}{\sigma^2} = \frac{(n-p)S^2}{\sigma^2} \stackrel{d}{=} \chi^2_{n-p} \right\}$$

\* Application of the general linear model:

- One-way and Two way ANOVA
- Multiple linear Regression