

Study { empirical studies { sample survey 抽样  
 observational 观察 1.1-1.2  
 experimental 实验 (有干预)  
 numerical summaries & graphical summaries 1.3-1.5  
 model fit → (plot) 2.6

PPDAL { target population: 目标群体  
 study population: 实验可能 cover 到的所有 study error  
 sample protocol: sample size Y sample randomly chosen from study population sample error 3.

point estimate  $\hat{\theta} = g(y_1, \dots, y_n)$  y: data 具体数值  
 • m.l.e:  $L(\theta)$  去 constant  $\rightarrow l(\theta) \rightarrow \frac{d}{d\theta} l(\theta) = 0$

point estimator  $\hat{\theta} = g(Y_1, \dots, Y_n)$  Y: variable function (有 distribution)

• sample distribution

$Y \sim G(\mu, \sigma)$   $Y \sim G(\mu, \frac{\sigma}{\sqrt{n}})$   
 $Y \sim P_0(\theta)$   $Y \sim G(\theta, \sqrt{\frac{\theta}{n}})$   
 $Y \sim Bi(n, \theta)$   $Y \sim G(\theta, \sqrt{\frac{\theta(1-\theta)}{n}})$   
 $Y \sim Exp(\theta)$   $Y \sim G(\theta, \frac{\theta}{\sqrt{n}})$

pivotal quantity  $\frac{Y - \mu}{\frac{\sigma}{\sqrt{n}}} \sim G(0, 1)$

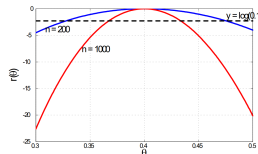
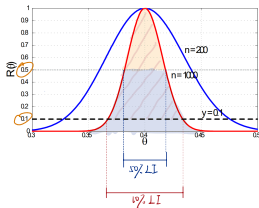
4.1-4.2

estimation

function: likelihood  $L(\theta) = \prod f(y_i; \theta)$  relative  $R(\theta) = \frac{L(\theta)}{L(\hat{\theta})}$  2.2-2.5  
 log  $l(\theta) = \log L(\theta)$   $r(\theta) = l(\theta) - l(\hat{\theta})$

interval: likelihood interval LI 4.3  
 $\{\theta: R(\theta) \geq q\}$

confidence interval CI 4.4  
 $\{\theta: P(L(Y) \leq \theta \leq U(Y)) = p\}$



$n \uparrow \rightarrow$  narrower

`> uniroot(function(x) (R(x) - q), lower, upper)$root`  
 公式

`> pnorm(1.96) - pnorm(-1.96)`  
 0.95

`> qnorm((1+0.95)/2)`  
 1.96

likelihood level (q) ← 4.6 → confidence level (p)

找 c  $p = P(W \leq c) = 2P(Z \leq \sqrt{c}) - 1$   $W \sim \chi^2_1$   
 $q = e^{-\frac{c}{2}}$

$p = P(L(Y) \leq -2 \log q) = 2P(Z \leq \sqrt{-2 \log q}) - 1$   
 $L(Y) \sim \chi^2_k$

100p% CI for  $G(\mu, \sigma)$  4.7

$\mu$  unknown,  $\sigma$  unknown  $\rightarrow$  CI (\*)

4.5

• Chi-squared  $\chi^2_k$

$W_1, W_2, \dots, W_n$  indep random variable,  $W_i \sim \chi^2_k$

$S = \sum_{i=1}^n W_i \sim \chi^2_{nk}$

$Z \sim G(0, 1) \Rightarrow Z^2 = W \sim \chi^2_1$   
 $Z_1, Z_2, \dots, Z_n \sim G(0, 1)$   $S = \sum_{i=1}^n Z_i^2 \sim \chi^2_n$

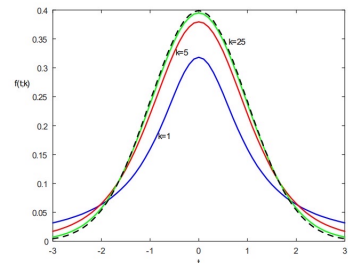
$W \sim \chi^2_k \Rightarrow P(W \geq w) = 2(1 - P(Z \leq \sqrt{w}))$   $Z \sim G(0, 1)$

$W \sim \chi^2_2 \Rightarrow W \sim Exp(1/2)$   $P(W \geq w) = e^{-w/2}$

• t distrib  $T \sim t(k)$   
 ↑ d.f

$Z \sim G(0, 1)$   $U \sim \chi^2_k \Rightarrow T = \frac{Z}{\sqrt{U/k}} \sim t_k$

k 越大, t(k) 曲线越接近  $G(0, 1)$



$p\text{-value} = P(D \geq d ; H_0)$

$p\text{-value}$	Interpretation
$p\text{-value} > 0.10$	No evidence against $H_0$ based on the observed data.
$0.05 < p\text{-value} \leq 0.10$	Weak evidence against $H_0$ based on the observed data.
$0.01 < p\text{-value} \leq 0.05$	Evidence against $H_0$ based on the observed data.
$0.001 < p\text{-value} \leq 0.01$	Strong evidence against $H_0$ based on the observed data.
$p\text{-value} \leq 0.001$	Very strong evidence against $H_0$ based on the observed data.

$H_0: \theta = \theta_0$   
 $p\text{-val} \leq 0.05 \Leftrightarrow 95\% \text{ CI for } \theta \text{ 不包含 } \theta_0$

5.1  $Bi(n, \theta)$   $D = |Y - n\theta|$   $n \text{ large} \xrightarrow{CLT} \frac{Y - n\theta}{\sqrt{n\theta(1-\theta)}} \sim G(0, 1)$

5.2  $G(\mu, \sigma)$   $H_0: \mu = \mu_0$   $p\text{-val} = 2[1 - P(T \leq \frac{|\bar{y} - \mu_0|}{\frac{\sigma}{\sqrt{n}}})]$   $T \sim t_{n-1}$   
 $H_0: \sigma^2 = \sigma_0^2$   $u = \frac{(n-1)S^2}{\sigma_0^2}$   $U \sim \chi_{n-1}^2$   $\begin{cases} P(U \leq u) > 0.5 & p\text{-val} = 2P(U \geq u) \\ P(U \leq u) < 0.5 & p\text{-val} = 2P(U \leq u) \end{cases}$

5.3 likelihood ratio test statistics  $\lambda(\theta_0) = -2 \log(L(\theta_0))$   $\lambda(\theta_0) \sim \chi_1^2$

b. Gaussian Response Model 比较 2 个 Variate  $X$  与  $Y$  在 关联性

$r = \frac{S_{xy}}{\sqrt{S_{xx} S_{yy}}}$   $S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$   $S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2$   $S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$   
 $\text{corr}(X, Y) = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{E[(X - \mu_x)(Y - \mu_y)]}{\sigma_x \sigma_y}$

• linear regression model:  $Y_i \sim G(\alpha + \beta x_i)$   $\alpha = \bar{y} - \hat{\beta} \bar{x}$   $\hat{\beta} = \frac{S_{xy}}{S_{xx}}$

- check linear reg model  $\begin{cases} \text{scatterplot \& fitted line} \\ \text{residual plot } (x_i, \hat{r}_i) \\ \text{standard residual plot } (x_i, \hat{r}_i^*) \\ \text{standard residual plot } (\hat{\mu}_i, \hat{r}_i^*) \end{cases}$

• 2 population 比较 b.4

• General  $G(\mu, \sigma)$  Response Model b.5  $Y_i \sim G(\mu(x_i), \sigma) \rightarrow \mu(x_i) = \beta_0 + \sum_{j=1}^k \beta_j x_{ij}$   
 $H_0: \beta_j = 0$   $T \sim t(n-k-1)$

7. Multinomial Model 比较 2 个 model

$\lambda(\theta_0) = 2 \sum_{j=1}^k Y_j \log(\frac{Y_j}{e_j})$   $Y_j = \text{observed}$   $e_j = \text{expected}$   $p\text{-value} = P(W \geq \lambda(\theta_0))$   $W \sim \chi^2(k-1-p)$

Pearson goodness of fit test 当有 observed 与 expected value 时直接套公式

$D = \sum_{j=1}^k \frac{(Y_j - E_j)^2}{E_j} \sim \chi_{k-1-p}^2$   $d = \sum_{j=1}^k \frac{(Y_j - e_j)^2}{e_j}$   $p\text{-value} = P(D \geq d)$   $D \sim \chi^2(k-1-p)$

Two Way Table d.f = (row-1)(col-1) ← 先根据 distribution  $f$   $e_j$

8. Cause & Effect

causation

Vassociate, x casually related 原因