

1.3种 testing

① Fisher null hypothesis test

必须先找到我们想验证事件的相反情况 (H_0)，然后证明 H_0 错误。(因为证明一件事正确很难，但证明它错误容易)

One-tailed test：形式为 $H_0: \theta \leq \theta_0$ 或 $H_0: \theta \geq \theta_0$

P-value：评判标准，越小可信度越高 (level of significance)

若 $P[D|H_0] \leq P\text{-value}$ (D 为实验数据)，则可以 reject H_0 (reject H_0 at the [P-value]

level of significance)

关于 D 的选取：对测试样本取平均值 (记为 \bar{X} ，样本数量 n)，

若 H_0 为 $\theta \leq \theta_0$ ，则 D 为 n 个样本在 H_0 中平均值大于 \bar{X} 的根元率
----- ≥ ----- < -----

若 H_0 为 $\theta = \theta_0$ ，先按照 \bar{X} 倾向的一边计算，然后 χ^2 (two-tailed)
具体的根元率计算，依赖正态分布的累积值，需参考 mid 资料第 9 页①部分

推荐的 P-value：小于 0.05

② Neyman-Pearson decision theory

考虑 H_0 与 H_1 ， H_1 即为我们想要证明的情况 (与 H_0 相反) (这里的 H_0 一般为 $\theta = \theta_0$)

Type I error: $\alpha = P[\text{reject } H_0 | H_0 \text{ true}]$

Type II error: $\beta = P[\text{accept } H_0 | H_1 \text{ true}]$

Power: $1 - \beta = P[\text{accept } H_1 | H_1 \text{ true}]$ (high power: rejection of H_0 is likely if H_1 true)

Critical region: for α , $Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}}$, $\frac{|\bar{X} - \mu_0|}{\sigma/\sqrt{n}} > Z_{\alpha/2}$ (\bar{X} 的范围为 critical region)

若实际处于 critical region 内，则 reject H_0
需要注意是两边，不是中间

假设 H_1 为 $\mu = \mu_0 + \delta$

$$\beta = \frac{1}{\sqrt{n}} \int_{-\infty}^{Z_{\alpha/2} - \delta/\sqrt{n}} e^{-\frac{t^2}{2}} dt = \int_{-\infty}^{-Z_\beta} e^{-\frac{t^2}{2}} dt$$

$$-Z_\beta \approx Z_{\alpha/2} - \delta/\sqrt{n}$$

$$n \approx \frac{(Z_{\alpha/2} + Z_\beta)^2 \sigma^2}{\delta^2}$$

也就是说，Neyman-Pearson 是教你怎样选取样本数量与判定标准的 (根据目标的 α 与 β)

Operating characteristic curves (OC curves): 在非正态分布下, 不能拒绝 H_0 的概率与实得平均数的关系, (大和小提前确定) (接受 H_0)

怎样知道 β : 计算出 $d(\bar{x})$ 在横坐标上找到, 再找到数量为 n 的曲线, 横坐标为 d 处的纵坐标值即为 β (课件 416~417 页)

OC curve 的最大值 (即 $P(M_0)$ 处) 为 $1-\alpha$ $d = \frac{M-M_0}{\sigma}$

注意, OC curve 为选择方法, 即在获得平均值后, 根据曲线上的概率进行选择, 选择 H_1 的概率即为曲线上 β 值。与真实情况是 H_0 还是 H_1 无关。

③ NHST, T-test

④ T-test (原始方差未知)

$$T_{n-1} = \frac{\bar{X} - M_0}{S/\sqrt{n}} \quad (S \text{ 为样本标准差}) \quad (\text{注意: 是样本标准差, 前面的都是本身已知的标准差})$$

Reject at significance level α :

$$H_0: M = M_0 \text{ if } |T_{n-1}| > t_{\alpha/2, n-1}$$

$$H_0: M \leq M_0 \text{ if } T_{n-1} > t_{\alpha, n-1}$$

$$H_0: M \geq M_0 \text{ if } T_{n-1} < -t_{\alpha, n-1}$$

OC curve: 课件 442~443 页 $d = \frac{|M - M_0|}{\sigma}$

⑤ Chi-squared test (用来测试标准差)

S 为样本标准差, σ_0 为假设的, σ 为实际标准差

$$\chi^2_{n-1} = \frac{(n-1)S^2}{\sigma_0^2}$$

Reject at significance level α :

$$H_0: \sigma = \sigma_0 \text{ if } \chi^2_{n-1} > \chi^2_{\alpha/2, n-1} \text{ or } \chi^2_{n-1} < \chi^2_{1-\alpha/2, n-1}$$

$$H_0: \sigma \leq \sigma_0 \text{ if } \chi^2_{n-1} > \chi^2_{\alpha, n-1}$$

$$H_0: \sigma \geq \sigma_0 \text{ if } \chi^2_{n-1} < \chi^2_{1-\alpha, n-1}$$

OC curve: 课件 448~449 页 $\lambda = \frac{\sigma}{\sigma_0} = \frac{\text{true}}{\text{hyphesis}}$

OC curve: 只取决于分布本身, α, n, σ , 不取决于具体的测试

描述的是 H_0 为假时无法否决 H_0 的概率 (即 β)

横坐标为真实分布的平均值

例如, 若真实为 $M_0 + \tau$, 则测试结果会不同, 可能在拒绝 H_0 的区间内, 也可能在该区间外, 这些不同情况即对应纵坐标值 (端点概率)

横坐标为实际值

在假设情况中的分布 (转换, 类似于标准化), 纵坐标为对应的真情况下原假设被接受的概率

接受 H_0 的概率

2. Non-parameter statistics

non-parametric statistics: 不依赖任何参数

distribution-free statistics: 不服从常规的分布

通常，这两种合起来称为 non-parametric methods

用中位数取代平均值，用 q_1, q_2, q_3 取代方差

Fisher test

想要验证 H_0 ， $Q_+ = \#\{X_k : X_k > M_0\}$ (比 M_0 大的数量)

$$P[Q_- \leq k | M=M_0] = \sum_{x=0}^k \binom{n}{x} \frac{1}{2^n}$$

Reject at significance level α :

$$H_0: M \leq M_0 \text{ if } P[Q_- \leq k | M=M_0] < \alpha$$

$$H_0: M \geq M_0 \text{ if } P[Q_+ \leq k | M=M_0] < \alpha$$

$$H_0: M=M_0 \text{ if } P[\min(Q_-, Q_+) \leq k | M=M_0] < \alpha/2$$

Wilcoxon signed rank test (课件 46 (~462页)) (要求关于某个数对称才可使用)

将与 M_0 差的绝对值 从小到大排序，比 M_0 小的乘 -1

W_+ 为比 M_0 大的序号之和， W_- 为比 M_0 小的序号之和 (计算序号时有重复要取对应)

$$E[W] = \frac{n(n+1)}{4} \quad \text{Var}[W] = \frac{n(n+1)(2n+1)}{24} - \sum \frac{a^3 - a}{48} \quad \text{区间平均值}$$

① 可直接比较 W 和 critical value (查表)

Reject at significance level α :

$$H_0: M \leq M_0 \text{ if } |W| < \text{critical value for } \alpha$$

$$H_0: M \geq M_0 \text{ if } W_+ < \text{critical value for } \alpha$$

$$H_0: M = M_0 \text{ if } W = \min(W_+, |W_-|) < \text{critical value for } \alpha/2$$

↓ (重复项，每组都要减， a 为一组中的数量)

② 计算概率 (与 ① 是两种方法，不是共同使用)

$$M = E[W] \quad \sigma^2 = \text{Var}[W]$$

$$Z = \frac{|W| - M}{\sigma}$$

calculate $P[Z < z]$ 与 P-Value 比较

3. Proportion

考虑一个群体中，每个个体只有两种情况(0和1)

$$X=0 \text{ or } 1$$

$$p = \frac{1}{N} \sum_i X_i$$

$$\text{平均数 } \hat{p} = \bar{X} = \frac{1}{n} \sum_i X_i$$

$$\frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \text{ 为正态分布}$$

$$100(1-\alpha)\% \text{ 置信区间: } \hat{p} \pm z_{\alpha/2} \sqrt{\hat{p}(1-\hat{p})/n}$$

$$\text{为使可信度大于 } 1-\alpha, \text{ 样本数量 } n \geq \frac{z_{\alpha/2}^2 \hat{p}(1-\hat{p})}{d^2}$$

$$z = \frac{\hat{p} - p_0}{\sqrt{p_0(1-p_0)/n}}$$

Reject at significance level α :

$$H_0: p = p_0 \text{ if } |z| > z_{\alpha/2}$$

$$H_0: p \leq p_0 \text{ if } z > z_\alpha$$

$$H_0: p \geq p_0 \text{ if } z < -z_\alpha$$

考虑两个群体, p_1 和 p_2 , 要研究 $p_1 - p_2$

$$\widehat{p_1 - p_2} = \hat{p}_1 - \hat{p}_2$$

$$100(1-\alpha)\% \text{ 置信区间: } \hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \quad (\text{pooled})$$

$$\text{Reject at significance level } \alpha: (z = \frac{\hat{p}_1 - \hat{p}_2 - (p_1 - p_2)_0}{\sqrt{\hat{p}_1(1-\hat{p}_1)/n_1 + \hat{p}_2(1-\hat{p}_2)/n_2}})$$

$$H_0: p_1 - p_2 = (p_1 - p_2)_0 \text{ if } |z| > z_{\alpha/2}$$

$$\cdots \cdots \leq \cdots \cdots \quad z > z_\alpha$$

$$\cdots \cdots \geq \cdots \cdots \quad z < -z_\alpha$$

Pooled estimator for proportion (比较两个大小)

$$H_0: p_1 = p_2$$

$$\hat{p} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$$

$$Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\hat{p}(1-\hat{p})} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

Reject at significance level α :

$$H_0: p_1 = p_2 \text{ if } |Z| > z_{\alpha/2}$$

$$H_0: p_1 \leq p_2 \text{ if } Z > z_\alpha$$

$$H_0: p_1 \geq p_2 \text{ if } Z < -z_\alpha$$

4. Comparisons

① Variances

$$X_1 \sim N(\mu_1, \sigma_1^2) \quad X_2 \sim N(\mu_2, \sigma_2^2) \quad A$$

$$\frac{(n_1-1)s_1^2}{\sigma_1^2} \sim \chi_{n_1-1}^2 \quad \frac{(n_2-1)s_2^2}{\sigma_2^2} \sim \chi_{n_2-1}^2 \quad A$$

$$F\text{-distribution: } F_{\gamma_1, \gamma_2} = \frac{\chi_{\gamma_1}^2 / \gamma_1}{\chi_{\gamma_2}^2 / \gamma_2} \quad (\text{F-distribution with } \gamma_1 \text{ and } \gamma_2 \text{ degrees of freedom})$$

$$P[F_{\gamma_1, \gamma_2} < x] = 1 - P[F_{\gamma_2, \gamma_1} < \frac{1}{x}]$$

$$\text{density function: } f_{\gamma_1, \gamma_2}(x) = \gamma_1^{\gamma_1/2} \gamma_2^{\gamma_2/2} \frac{\Gamma(\frac{\gamma_1+\gamma_2}{2})}{\Gamma(\frac{\gamma_1}{2}) \Gamma(\frac{\gamma_2}{2})} \frac{x^{\gamma_1/2-1}}{(x_1 x + x_2)^{(\gamma_1+\gamma_2)/2}}$$

$$f_{\alpha, \gamma_1, \gamma_2} : P[F_{\gamma_1, \gamma_2} > f_{\alpha, \gamma_1, \gamma_2}] = \alpha$$

$$f_{1-\alpha, \gamma_1, \gamma_2} = \frac{1}{f_{\alpha, \gamma_1, \gamma_2}}$$

考慮 A 及 $\sigma_1^2 = \sigma_2^2$, 則 $\frac{s_1^2}{s_2^2}$ 服从 n_1-1 和 n_2-1 的 F 分布

$$F_{n_1-1, n_2} = \frac{s_1^2}{s_2^2}$$

F -test, reject at significance level α :

$$H_0: \sigma_1 \leq \sigma_2 \text{ if } \frac{s_1^2}{s_2^2} > f_{\alpha, n_1-1, n_2-1}$$

$$\therefore Z \dots \frac{s_2^2}{s_1^2} > f_{\alpha, n_2-1, n_1-1}$$

$$H_0: \sigma_1 = \sigma_2 \text{ if } \frac{s_1^2}{s_2^2} > f_{\alpha/2, n_1-1, n_2-1} \text{ or } \frac{s_2^2}{s_1^2} > f_{\alpha/2, n_2-1, n_1-1}$$

OC curve: $\lambda = \frac{\sigma_1}{\sigma_2}$, 课件 502~503 页

② Means

$$\widehat{M_1 - M_2} = \widehat{M_1} - \widehat{M_2}$$

i) 差已知: $\bar{X}_1 \sim N(M_1, \sigma_1^2/n_1)$ $\bar{X}_2 \sim N(M_2, \sigma_2^2/n_2)$

$$\frac{\bar{X}_1 - \bar{X}_2 - (M_1 - M_2)}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}} \text{ 为标准正态分布}$$

$$Z = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}}$$

Reject at level α :

$$H_0: M_1 \leq M_2 \text{ if } Z > z_\alpha$$

$$H_0: M_1 \geq M_2 \text{ if } Z < -z_\alpha$$

$$H_0: M_1 = M_2 \text{ if } Z < -z_{\alpha/2} \text{ or } Z > z_{\alpha/2}$$

OC curve: $d = \frac{|M_1 - M_2|}{\sqrt{\sigma_1^2 + \sigma_2^2}}$ $n = \frac{\sigma_1^2 + \sigma_2^2}{\sigma_1^2/n_1 + \sigma_2^2/n_2}$ (equivalent sample size) (课件 50~51)

ii) 差未知但相等 ($\sigma_1^2 = \sigma_2^2 = \sigma^2$)

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (M_1 - M_2)}{\sqrt{\sigma^2(1/n_1 + 1/n_2)}}$$

$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2} \quad (\text{用来预测方差, pooled estimator for variance})$$

$$\chi^2_{n_1+n_2-2} = \frac{(n_1-1)S_1^2}{\sigma^2} + \frac{(n_2-1)S_2^2}{\sigma^2} \quad \text{为 } n_1+n_2-2 \text{ 自由度的 chi-squared 分布}$$

$$T_{n_1+n_2-2} = \frac{Z}{\sqrt{\chi^2_{n_1+n_2-2}/(n_1+n_2-2)}} = \frac{(\bar{X}_1 - \bar{X}_2) - (M_1 - M_2)}{\sqrt{S_p^2(1/n_1 + 1/n_2)}} \quad \text{为 } n_1+n_2-2 \text{ 自由度的 T 分布}$$

$M_1 - M_2$ 的 $100(1-\alpha)\%$ 置信区间: $(\bar{X}_1 - \bar{X}_2) \pm t_{\alpha/2, n_1+n_2-2} \sqrt{S_p^2(1/n_1 + 1/n_2)}$

Student's T test: $T_{n_1+n_2-2} = \frac{(\bar{X}_1 - \bar{X}_2) - (M_1 - M_2)}{\sqrt{S_p^2(1/n_1 + 1/n_2)}}$

~~Reject at α :~~

Reject at α :

$$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0 \text{ if } |\bar{T}_{n_1+n_2-2}| > t_{\alpha/2, n_1+n_2-2}$$

$$H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0 \text{ if } \bar{T}_{n_1+n_2-2} > t_{\alpha/2, n_1+n_2-2}$$
$$\quad \quad < -t_{\alpha/2, n_1+n_2-2}$$

OC curves: $d = \frac{|\mu_1 - \mu_2|}{\sigma}$ ($n^* = 2n - 1$), 课件 520-521

iii) 方差未知且不同

$$\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2} \right) = \frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}$$

Welch-Satterthwaite Approximation:

$$\gamma := \frac{(\lambda_1 s_1^2 + \dots + \lambda_k s_k^2)^2}{\sum_{i=1}^k \frac{(\lambda_i s_i^2)^2}{n_i - 1}}, \quad (\text{a}) \quad \gamma \cdot \frac{\lambda_1 s_1^2 + \dots + \lambda_k s_k^2}{\lambda_1 \sigma_1^2 + \dots + \lambda_k \sigma_k^2} \text{ 近似为 } \gamma \text{ 自由度的卡方分布}$$

Let $k=2, \lambda_1 = \frac{1}{n_1}, \lambda_2 = \frac{1}{n_2}$

(a) $\gamma = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{(s_1^2/n_1)^2/(n_1-1) + (s_2^2/n_2)^2/(n_2-1)}$

r. $\frac{s_1^2/n_1 + s_2^2/n_2}{\sigma_1^2/n_1 + \sigma_2^2/n_2}$ 为 γ 自由度的 chi-squared 分布

$T_\gamma = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)_0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$ 为 γ 自由度的 T 分布 (正常需将 γ 向下取整)

这是 Welch's test, Reject at α :

$$H_0: \mu_1 - \mu_2 = (\mu_1 - \mu_2)_0 \text{ if } |T_\gamma| > t_{\alpha/2, \gamma}$$

$$H_1: \mu_1 - \mu_2 \neq (\mu_1 - \mu_2)_0 \text{ if } T_\gamma > t_{\alpha, \gamma}$$

$$H_1: \mu_1 - \mu_2 \geq (\mu_1 - \mu_2)_0 \text{ if } T_\gamma < -t_{\alpha, \gamma}$$

5. Non-parametric comparisons

有 X、Y 两个分布, 比较 $X > Y$ 的概率

$X_1 \sim X_m$ 和 $Y_1 \sim Y_n$ ($m \leq n$); 从大到小排序, 记为 R_i ($i=1 \sim m+n$), 重复的记为序号
令 $W_m = X_1 \sim X_m$ 的序号之和

$$E[W_m] = \frac{m(m+n+1)}{2} \quad \text{Var}[W_m] = \frac{mn(m+n+1)}{12}$$

若有重复项，则 $\text{Var}[W_m] = \frac{mn(m+n+1)}{12 - \sum \frac{t^2 - t}{12}}$ (t 为重复的项数) (#如果为负数不要减)

$$H_0: P[X > Y] = \frac{1}{2}$$

$$Z = \frac{W_m - E[W_m]}{\sqrt{\text{Var}[W_m]}}$$

, 若 $Z < Z_{\alpha/2}$ 则 reject , 大于 $Z_{\alpha/2}$ 则 reject $H_0: P[X < Y] > \frac{1}{2}$

Paired test: 有 X, Y 两个分布，研究 $D = X - Y$ (要求 X, Y 正态分布)

$$H_0: \mu_D = 0$$

$$\bar{T}_{n-1} = \frac{\bar{D} - \mu_D}{\sqrt{s_D^2/n}} \quad (\text{paired } T \text{ test}), \text{ reject 规则和前面都一样, 不再叙述}$$

pooled vs paired: pooled (池子, 所有放一起), paired (一对一对比较)

$$T_{\text{pooled}} = \frac{\bar{X} - \bar{Y}}{\sqrt{2s_p^2/n}} \quad \text{critical value} = t_{\alpha/2, 2n-2}$$

$$T_{\text{paired}} = \frac{\bar{X} - \bar{Y}}{\sqrt{s_D^2/n}} \quad \text{critical value} = t_{\alpha/2, n-1}$$

Pooled test 可以更容易 Reject H_0 , 更好

$$\frac{\sigma_D^2}{n} = \frac{2\sigma^2}{n}(1 - \rho_{XY})$$

若 $\rho_{XY} > 0$, paired 更好, 否则 pooled 更好

Correlation

$$\text{Var}[X] = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

$$\text{Cov}[X, Y] = \frac{1}{n-1} \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\text{correlation coefficient (相当于 } \rho_{XY}): R = \hat{\rho} = \frac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$$

更多内容: 课件 558~563页

b. Categorical data

Categorical random variable: X can take n values link with probabilities

$$p_1 \sim p_k \quad (p_1 + \dots + p_k = 1)$$

Multinomial trial (with parameters $p_1 \sim p_k$): result in exactly one from k possible outcomes, possibility is p_i respectively

Multinomial random variables: n 次 M... trial 中出现结果 i 的次数

Multinomial distribution: $f_{X_1 X_2 \dots X_k}(x_1, x_2, \dots, x_k) = \frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}$
 n 次 M ... dis...: $E[X_i] = np_i$, $\text{Var}[X_i] = np_i(1-p_i)$
 $\text{Cov}[X_i, X_j] = -np_i p_j$

Pearson statistic: $\sum_{i=1}^k \frac{(X_i - np_i)^2}{np_i}$ 近似为 $k-1$ 自由度 Chi-squared 分布

Cochran's rule: Pearson statistic 要求 n 足够大, 即 $E[X_i] \geq 1$ 对所有 i , 且 $E[X_i] \geq 5$ 对 80% 以上的 i

Pearson's Chi-squared Goodness-of-fit Test

$$H_0: p_i = p_{i0} \quad (i=1, \dots, k)$$

$$\chi^2_{k-1} = \sum_{i=1}^k \frac{(X_i - np_{i0})^2}{np_{i0}} \quad (\text{reject } H_0 \text{ at } \alpha \text{ if } \chi^2_{k-1} > \chi^2_{\alpha, k-1})$$

Goodness-of-fit Test for discrete distribution

根据对应分布的自身特性, 分出几类, 得到每类的数量和期望, 对其使用上面方法 (不满足 Cochran's rule 的类要去掉), 看能否 reject (大于 χ_α)
 课件 579 ~ 583 页

For continuous distribution

先分类 (分段), $p = p_i = P[a_{i-1} \leq X \leq a_i] = \int_{a_{i-1}}^{a_i} f(x) dx$ (每段概率都相等)

注意: 是根据 H_0 分段, 不是根据测试数据

计算出每类的 O_i 和 E_i , 计算 χ^2 , 看能否 reject (大于 χ_α)

课件 584 ~ 588 页

Contingency Table

	<6h	6-9h	>9h
Row	n_{11}	n_{12}	n_{13}
avg	n_{21}	n_{22}	n_{23}
high	n_{31}	n_{32}	n_{33}

marginal row and column sums:

$$n_{i \cdot} = \sum_{j=1}^c n_{i,j} \quad n_{\cdot j} = \sum_{i=1}^r n_{i,j}$$

$H_0: p_{ij} = p_i \cdot p_j$ (完全独立)

$$E_{ij} = n \cdot \hat{p}_{ij} = n \cdot \hat{p}_i \cdot \hat{p}_j = n \cdot \frac{n_i}{n} \cdot \frac{n_j}{n} = \frac{n_i \cdot n_j}{n}$$

$$\chi^2_{(r-1)(c-1)} = \sum_{i=1}^r \sum_{j=1}^c \frac{(a_{ij} - E_{ij})^2}{E_{ij}}$$

为 chi-squared 分布，看能否 reject ($\lambda > \chi_{\alpha}$)

例子：课件 59 ~ 600 页

Comparing proportions (test for homogeneity)

每一类个体的概率设为 1，后面都一样，自由度仍为 $(r-1)(c-1)$

7. linear regression

Independent variable X , 随意的, predictor variable, regressor

Dependent variable Y , 随机的, 要研究 Y 与 X 的关系, response variable

① Simple linear regression

$$Y|X = \beta_0 + \beta_1 X + E \quad (E[E] = 0)$$

需要计算: b_0 : estimator for β_0

b_1 : estimator for β_1

会有 n 组 (X, Y) 数据, 写作 $Y_i = Y|X_i \quad (i=1 \sim n)$, 则有 $(X_1, Y_1) \sim (X_n, Y_n)$

$Y_i = b_0 + b_1 X_i + e_i$ (e_i 称为 residual), 目标是使 e_i 小

最小二乘法 (Least-Squares Estimation): $e_1^2 + e_2^2 + \dots + e_n^2 \rightarrow \text{minimum}$

Error sum of squares: $SS_E = e_1^2 + \dots + e_n^2 = \sum_n (Y_i - (b_0 + b_1 X_i))^2$

$$\frac{\partial SS_E}{\partial b_0} = -2 \sum_n (Y_i - b_0 - b_1 X_i) = 0, \quad \frac{\partial SS_E}{\partial b_1} = -2 \sum_n (Y_i - b_0 - b_1 X_i) X_i = 0$$

$$\text{Then } nb_0 + b_1 \sum_n X_i = \sum_n Y_i \quad b_0 \sum_n X_i + b_1 \sum_n X_i^2 = \sum_n X_i Y_i$$

$$\text{Then } b_1 = \frac{\sum_n X_i Y_i - \sum_n X_i \cdot \sum_n Y_i}{n \sum_n X_i^2 - (\sum_n X_i)^2} \quad b_0 = \frac{1}{n} \sum_n Y_i - \frac{b_1}{n} \sum_n X_i$$

$$S_{XX} = \sum_n (X_i - \bar{X})^2 \quad S_{YY} = \sum_n (Y_i - \bar{Y})^2 \quad S_{XY} = \sum_n (X_i - \bar{X})(Y_i - \bar{Y})$$

$$\text{P.s. } b_0 = \bar{Y} - b_1 \bar{X}, \quad b_1 = \frac{S_{XY}}{S_{XX}}, \quad SS_E = S_{YY} - b_1 S_{XY}$$

$\frac{b_1 - \beta_1}{\sigma / \sqrt{\sum (x_i - \bar{x})^2}}$ 和 $\frac{b_0 - \beta_0}{\sigma / \sqrt{\frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2}}}$ 都服从标准正态分布 (其中 σ 为 error 的标准差, 基于理论而不是计算)

Estimator for variance

$$S^2 = \frac{SS_E}{n-2} = \frac{1}{n-2} \sum (y_i - \hat{y}_{i|x_i})^2 \text{ 为 unbiased estimator}$$

$$(n-2)S^2 / \sigma^2 = \frac{SS_E}{\sigma^2}, \text{ 为 } n-2 \text{ 自由度的 chi-squared 分布}$$

100(1-\alpha)% 置信区间: $\beta_1 = b_1 \pm t_{\alpha/2, n-2} \frac{S}{\sqrt{S_{xx}}}$, $\beta_0 = b_0 \pm t_{\alpha/2, n-2} \frac{S \sqrt{\sum x_i^2}}{\sqrt{n S_{xx}}}$

$$H_0: \beta_0 = \beta_0' \text{ and } H_0: \beta_1 = \beta_1'$$

Significant: if a regression has evidence that slope $\beta_1 \neq 0$

$$H_0: \beta_1 = 0$$

$$T_{n-2} = \frac{b_1}{S / \sqrt{S_{xx}}}, \text{ reject at } \alpha \text{ if } |T_{n-2}| > t_{\alpha/2, n-2}$$

$$\hat{y}_{i|x} = \beta_0 + \beta_1 x = \bar{Y} - \beta_1 \bar{x} + \beta_1 x = \bar{Y} + \beta_1 (x - \bar{x})$$

这里是根据新的 x 计算
对应的 y

$$\frac{\hat{y}_{i|x} - \bar{y}_{i|x}}{S \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}} \text{ 为 } n-2 \text{ 自由度的 T 分布}$$

100(1-\alpha)% 置信区间: $\hat{y}_{i|x} \pm t_{\alpha/2, n-2} S \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}$ (confidence band)

Prediction band: $T_{n-2} = \frac{\hat{y}_{i|x} - \bar{y}_{i|x}}{S \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}}}$

band is $\hat{y}_{i|x} \pm t_{\alpha/2, n-2} S \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{S_{xx}}} \quad (100(1-\alpha)%)$

区别: confidence band 针对的是拟合本身存在误差, 给的区间,
prediction band 针对数据本身具有随机性, 给的空间

$$SS_T (\text{Total Sum of squares}): SS_T = S_{yy} = \sum_n (y_i - \bar{y})^2$$

$$\text{Coefficient of determination: } R^2 = \frac{SS_T - SS_E}{SS_T} \quad (\text{越接近 } 100\%, \text{ 越好})$$

$$R^2 = \frac{S_{XY}^2}{S_{XX} S_{YY}}$$

$$\bar{T}_{n-2} = \frac{B_1}{\sqrt{S^2/S_{XX}}} = \frac{R}{\sqrt{1-R^2}} \sqrt{n-2}$$

$$F_{1,n-2} = (n-2) \frac{R^2}{1-R^2} = (n-2) \frac{SS_T - SS_E}{SS_E}$$

F Test for significance

Reject $H_0: \beta_1 = 0$ at α if $F_{1,n-2} > f_{\alpha, 1, n-2}$

Test for correlation

$H_0: \rho = 0$, reject at α if $\left| \frac{R\sqrt{n-2}}{\sqrt{1-R^2}} \right| > t_{\alpha/2, n-2}$

多次測試

Y_{ij} 为第 j 次测出 X_i 的对应 Y 值

Sample mean $\bar{Y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$

Internal sum of squares: $\sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$

Error sum of squares due to pure error: $SS_{E;pe} = \sum_{i=1}^k \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$

~~$\frac{1}{6} SS_{E;pe} = \sum_{i=1}^k \frac{1}{n_i} \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_i)^2$~~

$$= \sum_{i=1}^k \sum_{j=1}^{n_i} Y_{ij}^2 - \frac{1}{n_i} \left(\sum_{j=1}^{n_i} Y_{ij} \right)^2$$

选取 $\bar{Y} = \bar{Y}_i$, 误差最小

Error sum of squares due to lack of fit: $SS_{E;ff} = SS_E - SS_{E;pe}$

Test for lack of fit

$F_{k-2, n-k} = \frac{SS_{E;ff} / (k-2)}{SS_{E;pe} / (n-k)}$, H_0 : model appropriate

Reject at α : $F_{k-2, n-k} > f_{\alpha, k-2, n-k}$

② Multiple linear regression (太难, 跳过)