

Lecture 7 成对数据的统计分析

x_1, x_2, \dots, x_n

必修 一元 频率分布 \bar{x}, S^2 样本估计总体

二元 相关性 $\left\{ \begin{array}{l} 1^\circ \text{ 数值数据之间} \\ 2^\circ \text{ 有序度量之间} \end{array} \right. \quad x \text{ 的取值分布连续区间}$

$\left\{ \begin{array}{l} 及格 \rightarrow 0 \\ 不及格 \rightarrow 1 \end{array} \right. \quad \left\{ \begin{array}{l} 男 \rightarrow 0 \\ 女 \rightarrow 1 \end{array} \right.$

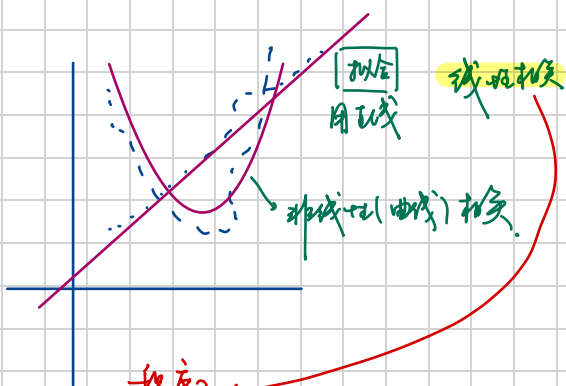
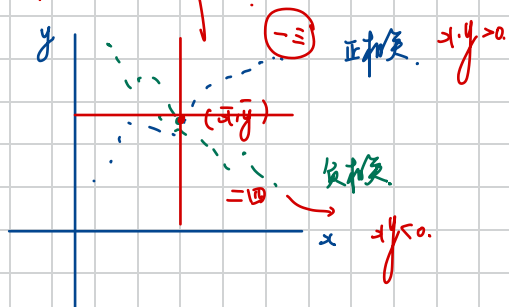
$(x_1, y_1), \dots, (x_n, y_n)$

2° 分类 2x2列联表

X \ Y	0	1
0	n_{11}	n_{12}
1	n_{21}	n_{22}

样本估计总体
误差

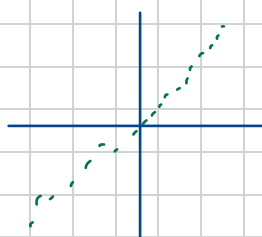
1° Pearson (线性) 相关系数



平移 $(x_i - \bar{x}, y_i - \bar{y})$

体重 vs 身高 $\textcircled{m} \rightarrow \textcircled{m}$ 同一个样本 L_{xy} 个100倍

$$L_{xy} = (x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})$$



$L_{xy} > 0 \rightsquigarrow$ 正相关 $L_{xy} < 0 \rightsquigarrow$ 负相关

相关程度 \rightarrow 数值大小

$$\vec{a} = (x_1, x_2) \quad \vec{b} = (y_1, y_2) \quad \text{定义}$$

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

$$\cos \theta = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} = \frac{x_1 y_1 + x_2 y_2}{\sqrt{x_1^2 + x_2^2} \sqrt{y_1^2 + y_2^2}}$$

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

$$= \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} \in [-1, 1]$$

$$\vec{a} = (x_1 - \bar{x}, x_2 - \bar{x}, \dots, x_n - \bar{x}) \quad \vec{b} = (y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_n - \bar{y})$$

取到 -1 和 1 说明 $\vec{a} \parallel \vec{b}$

$$\text{定义 } \vec{a} \cdot \vec{b} = (x_1 - \bar{x})(y_1 - \bar{y}) + \dots + (x_n - \bar{x})(y_n - \bar{y})$$

即 $\exists \lambda \in \mathbb{R}, \vec{a} = \lambda \vec{b}$

$$|\vec{a}| = \sqrt{\vec{a} \cdot \vec{a}}$$

$$\Rightarrow x_1 - \bar{x} = \lambda (y_1 - \bar{y}) \quad r = \frac{1}{\lambda} \quad \text{若 } \lambda = 0, r = 1, k > 0$$

命题 (Cauchy-Schwarz) 在 \mathbb{R}^n 空间中, 设非零向量 $\vec{a} \cdot \vec{b} = x_1 y_1 + \dots + x_n y_n$

$$\text{证: } (\vec{a} \cdot \vec{b})^2 \leq |\vec{a}|^2 |\vec{b}|^2$$

$$\text{证: } |\vec{a}|^2 = x_1^2 + \dots + x_n^2 \geq 0 \quad \vec{a} = \vec{0} \Leftrightarrow \vec{a} \cdot \vec{a} = 0$$

构造 $\vec{a} - t\vec{b} \quad \forall t \in \mathbb{R}$

$$(\vec{a} - t\vec{b}) \cdot (\vec{a} - t\vec{b}) \geq 0$$

$$\vec{a}^2 - 2t \vec{a} \cdot \vec{b} + t^2 \vec{b}^2 \geq 0, \quad \forall t \in \mathbb{R} \Rightarrow \Delta = 4(\vec{a} \cdot \vec{b})^2 - 4\vec{a}^2 \vec{b}^2 \leq 0$$

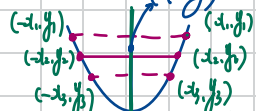
$$\Rightarrow (\vec{a} \cdot \vec{b})^2 \leq \vec{a}^2 \vec{b}^2$$

$$\Rightarrow (x_1 y_1 + \dots + x_n y_n)^2 \leq (x_1^2 + \dots + x_n^2)(y_1^2 + \dots + y_n^2)$$

$$\text{取 } t = \frac{\vec{a} \cdot \vec{b}}{\vec{b}^2} \Rightarrow \vec{a} - t\vec{b} = \vec{0} \Rightarrow \vec{a} \parallel \vec{b}$$

$$\angle_{xy} \begin{cases} > 0 \text{ 正} \\ < 0 \text{ 负} \end{cases}$$

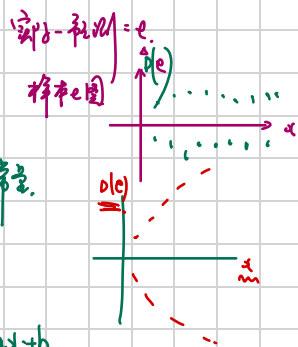
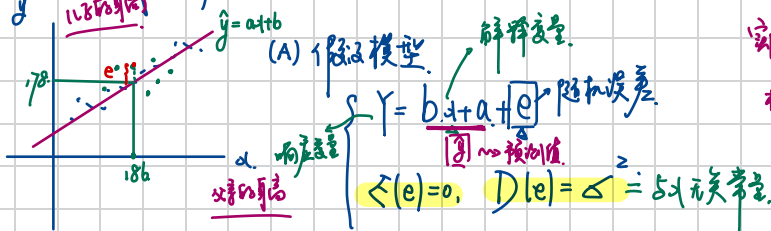
$$r = \frac{\angle_{xy}}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \in [-1, 1]$$



$|r| \rightarrow 1$ 线性强
 $|r| = 1$ 线性
 $|r| = 0$ 没有线性相关性

$$\angle_{xy} = (-x_1 + \bar{x})(y_1 - \bar{y}) + (-x_2 + \bar{x})(y_2 - \bar{y}) + \dots = 0 \Rightarrow r = \frac{\angle_{xy}}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} = 0$$

2° 一元线性回归模型



1° 表中 $x \rightarrow y$. 回归分析 $x \rightarrow$ 总体的 $y \rightarrow \bar{y}$
 期望值 = 预测值 $a + bx$

2° $\sigma^2 =$ 随机误差 e_i 的方差 = 回归分析中代入 x 后得到的总体的方差

3° σ^2 与 x 无关

(B) 用样本数据估计未知参数 a, b, σ^2

$$y = y + e = b_0 + a + e$$

$$(x_1, y_1), \dots, (x_n, y_n) \Rightarrow \sum_{i=1}^n [y_i - (bx_i + a)] = \sum_{i=1}^n e_i$$

最佳的 a, b

$$\sum_{i=1}^n |y_i - (bx_i + a)| \xrightarrow{\text{也可} \Rightarrow \text{这太难了}} \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - (bx_i + a))^2 \checkmark = Q(a, b) \min$$

$$Q(a, b) = \sum_{i=1}^n (y_i - bx_i - (\bar{y} - b\bar{x}) + (\bar{y} - b\bar{x}) - a)^2$$

$$= \sum_{i=1}^n \left[\underbrace{(y_i - \bar{y}) - b(x_i - \bar{x})}_{\triangle} + \underbrace{(\bar{y} - b\bar{x}) - a}_{\triangle} \right]^2$$

$$= \sum_{i=1}^n \left[(y_i - \bar{y}) - b(x_i - \bar{x}) \right]^2 + 2 \sum_{i=1}^n \left[(y_i - \bar{y}) - b(x_i - \bar{x}) \right] (\bar{y} - b\bar{x}) + n(\bar{y} - b\bar{x} - a)^2$$

$$= 0 + 0 + n(\bar{y} - b\bar{x} - a)^2 \geq 0 \Rightarrow a = \bar{y} - b\bar{x}$$

所证 \hat{y} 的表达式 $\hat{y} = \hat{a} + b\hat{x} = [\bar{y} - b\bar{x}] + b\hat{x}$ 在表中

$$= \sum_{i=1}^n \left[(y_i - \bar{y})^2 - 2b(x_i - \bar{x})(y_i - \bar{y}) + b^2(x_i - \bar{x})^2 \right]$$

$$= \sum_{i=1}^n (y_i - \bar{y})^2 - \left(2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right) b + b^2 \sum_{i=1}^n (x_i - \bar{x})^2$$

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

$$Q(a, b) \geq Q(\hat{a}, \hat{b})$$

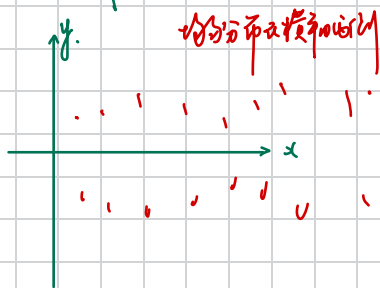
\hat{a}, \hat{b} 是 a, b 的 最小二乘估计
 $\hat{y} = \hat{b}x + \hat{a}$ 拟合的直线

拟合的直线

(c) 误差分析

$$\begin{cases} \hat{y} = \hat{a} + \hat{b}x + e \\ \text{Var}(e) = \sigma^2, \quad E(e) = 0 \end{cases}$$

残差图



R^2 决定系数

$$\begin{array}{ccc} y_i & \hat{y}_i & \bar{y} \\ \parallel & \parallel & \parallel \\ b\hat{x}_i + e_i & b\hat{x}_i + \hat{a} & \text{所有真实值的均值} \\ \uparrow & & \\ e_i \text{ 残差} & & \end{array}$$

总偏差平方和 (总差)

$$S^2 = \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S^2 = \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2$$

R^2

$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n e_i^2 \quad \text{拟合}$$

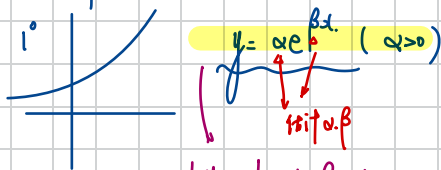
$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (\hat{a} + \hat{b}x_i - \bar{y})^2$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \in [0, 1]$$

R^2 模型拟合度

$$\begin{aligned} &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{b}x_i - \hat{a})(\hat{b}x_i + \hat{a} - \bar{y}) \\ &= \sum_{i=1}^n (y_i - \hat{b}x_i - \hat{a})(\hat{b}x_i + \hat{a} - \bar{y}) \\ &= \sum_{i=1}^n [(y_i - \bar{y}) - \hat{b}(x_i - \bar{x})] \hat{b}(x_i - \bar{x}) \\ &= \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - \hat{b} \sum_{i=1}^n (x_i - \bar{x})^2 \\ &= 0 \cdot \hat{b} = 0 \end{aligned}$$

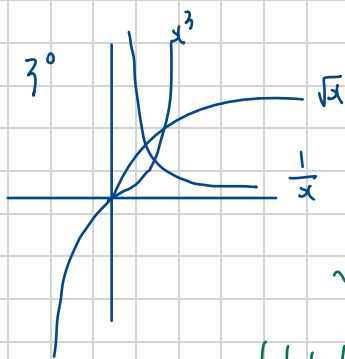
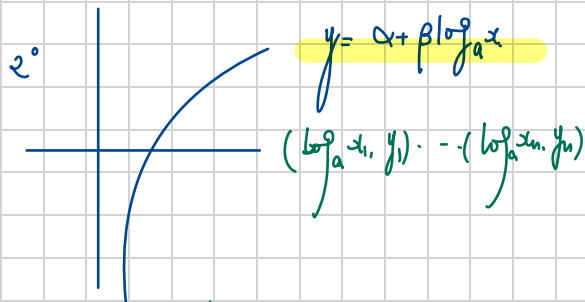
(d) 非线性回归模型



$\ln y = \ln \alpha + \beta \cdot x$

有 α, β

$(x_1, \ln y_1), (x_2, \ln y_2), \dots, (x_n, \ln y_n)$



$y = \alpha x^\beta \quad (\alpha > 0)$

$\Rightarrow \ln y = \ln \alpha + \beta \ln x$

$(\ln x_1, \ln y_1), \dots, (\ln x_n, \ln y_n)$

4°

$(\alpha \neq 0)$

$y = \frac{x}{\alpha x + \beta}$

$\frac{1}{y} = \frac{\alpha x + \beta}{x + \frac{\beta}{\alpha}}$

$= \frac{1}{\alpha} - \frac{\beta}{\alpha} \cdot \frac{1}{\alpha x + \beta}$

$\frac{1}{y} = \alpha + \beta \cdot \frac{1}{x}$

$(\frac{1}{x_1}, \frac{1}{y_1}), \dots, (\frac{1}{x_n}, \frac{1}{y_n})$

3° 分类变量的 χ^2 检验

X, Y 是否独立? \leadsto 教子成绩是否比男甲、乙校类别是否独立?

Y \ X	X	
	0	1
0	n_{11}	n_{12}
1	n_{21}	n_{22}

样本数据

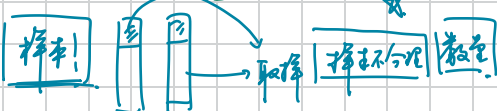
学校	数学成绩		合计
	不优秀 (Y=0)	优秀 (Y=1)	
甲校 (X=0)	33	10	43
乙校 (X=1)	38	7	45
合计	71	17	88



1° $\frac{n_{11}n_{22}}{n_{12}n_{21}} = \frac{10}{43} = 0.23$

$\frac{n_{12}n_{21}}{n_{11}n_{22}} = \frac{7}{45} = 0.16$

独立



误差=? 犯错误p=?

χ^2 (卡方) 检验

1. 零假设 H_0 : 假设 X 与 Y 无关 \Rightarrow 独立.

$$P\{X=0 | Y=0\} = P\{X=0 | Y=1\} \Rightarrow P(B|A) = P(B|\bar{A})$$

$\Leftrightarrow \{X=0\}$ 与 $\{Y=0\}$ 独立.

$\{A, B\}, \{A, \bar{B}\}, \{\bar{A}, B\}, \{\bar{A}, \bar{B}\}$ 独立.

\Rightarrow 随机变量 X, Y 独立.

2.

	Y=0	Y=1	
X=0	a	c	a+c
X=1	b	d	b+d
	a+b	c+d	n = a+b+c+d

独立概率

$$P\{Y=0 | X=0\} = \frac{a}{a+c}$$

$$P\{Y=0 | X=1\} = \frac{b}{b+d}$$

$$\left| \frac{a}{a+c} - \frac{b}{b+d} \right|$$

$$P\{Y=0 | X=0\} \neq P\{Y=0 | X=1\}$$

$$P\{X=0, Y=0\} = P\{X=0\} P\{Y=0\}$$

第1个假设

$$\frac{a}{a+b+c+d} - \frac{(a+c)(b+d)}{[a+b+c+d]^2} = \frac{a}{n} - \frac{(a+c)(b+d)}{n^2}$$

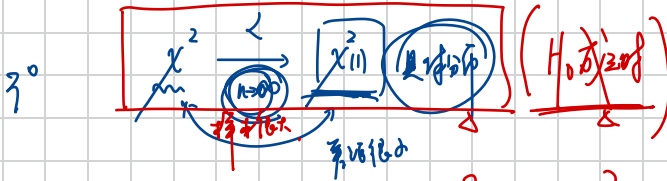
$$= \frac{1}{n} \left[a - \frac{(a+c)(b+d)}{n} \right]$$

1. 期望差值

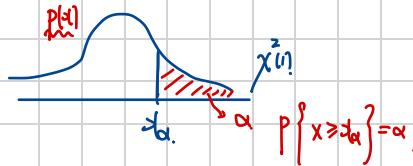
$$\chi^2 = \frac{\left(a - \frac{(a+c)(b+d)}{n} \right)^2}{\frac{(a+c)(b+d)}{n}} + \frac{\left(b - \frac{(a+b)(b+d)}{n} \right)^2}{\frac{(a+b)(b+d)}{n}} + \frac{\left(c - \frac{(a+c)(c+d)}{n} \right)^2}{\frac{(a+c)(c+d)}{n}} + \frac{\left(d - \frac{(c+d)(b+d)}{n} \right)^2}{\frac{(c+d)(b+d)}{n}}$$

$$= \frac{n(ad-bc)^2}{(a+c)(a+b)(b+d)(b+c)}$$

尽可能在 H_0 成立情况下
不能太大



记 $P\{\chi^2 \geq \chi_\alpha\} \approx P\{\chi^2(1) \geq \chi_\alpha\} = \alpha$.



(b) χ^2 独立性检验中可能用到的小概率值 α 与对应的临界值 x_α 如下, 即 $P\{\chi^2 \geq x_\alpha\} = \alpha$:

α	0.1	0.05	0.01	0.005	0.001
x_α	2.706	3.841	6.635	7.879	10.828

$\chi^2 = 12 > \chi_{0.001}$ $P\{\chi^2 > \chi_{0.001}\} = 0.001$

小概率事件 但发生了

1° 没有问题! \Rightarrow 拒绝 H_0

\Rightarrow x, y 有关, 犯错误概率不超过 0.001

即 99.99% 的把握 x, y 有关 \checkmark

$\chi^2 = 2.6 < 2.706 = \chi_{0.1}$ $\alpha = 0.1$

没有充分证据说明 H_0 为错误的。

接受 H_0 , 即认为 x, y 独立 (暂定)

注: 此处不可用 $0.1 = \alpha$ 来证明 x, y 独立的结论
的正确性!