

第二章：一元线性回归

马学俊 (主讲) 晁越 (助教)

苏州大学

数学科学学院

<https://xuejunma.github.io/>



- 1 引言
- 2 一元线性回归模型
- 3 参数 β_0 , β_1 的估计
- 4 最小二乘估计的性质
- 5 回归方程的显著性检验
- 6 残差分析
- 7 回归系数的区间估计
- 8 预测和控制

引言

- $F(x) = P(X \leq x); -\infty < x < \infty$
- $f_X(x) = f(x) = F'(x)$
-

$$\mathbf{E}(X) = \sum_x xp(x)$$

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} xf(x)dx$$

- $E(x)$ is the "center" of a distribution (or its r.v.) in the sense that

$$\min_b \mathbf{E}(X - b)^2 = \mathbf{E}[X - \mathbf{E}X]^2$$

Conditional Expectation

- Conditional Expectation of X when Y is given as y is that
 - $\mathbf{E}(X|Y = y) = \sum_x xp_{X|Y}(X|Y)$ for discrete r.v.
 - $\mathbf{E}(X|Y = y) = \int_{-\infty}^{\infty} xf_{X|Y}(x|y)dx$ for cont. r.v.
 - Interpretation: Note that $X|Y = y$ is a new r.v., $\mathbf{E}(X|Y = y)$ is the expectation on this r.v.
- Law of Total Expectation ¹

$$\mathbf{E}[\mathbf{E}(X|Y)] = \mathbf{E}(X)$$

- Law of Total Variance

$$\text{Var}(X) = \text{Var}[\mathbf{E}(X|Y)] + \mathbf{E}[\text{Var}(X|Y)]$$

¹Statistical Inference 2nd Edition by George Casella Roger L. Berger

if X and Y are any two r.v.s, then

$$\mathbf{E}(X) = \mathbf{E}[\mathbf{E}(X|Y)]$$

Proof:

$$\begin{aligned} \mathbf{E}X &= \int \int xf(x, y) dx dy \\ &= \int \left[\int xf(x|y) dx \right] f_Y(y) dy \\ &= \int \mathbf{E}(X|y) f_Y(y) dy = \mathbf{E}[\mathbf{E}(X|Y)] \end{aligned}$$

In general, the conditional expectation $\mathbf{E}[X|Y]$ can be defined as a r.v. $g(Y)$ such that

$$\mathbf{E}[(X - g(Y))^2] = \inf_{\text{among all reasonable function } h} \mathbf{E}[(X - h(X))^2]$$

or $\mathbf{E}[X|Y]$ is the function of Y which is "closest" to X in terms of mean square error.

For any two random variables X and Y

$$\text{Var}(X) = \text{Var}[\mathbf{E}(X|Y)] + \mathbf{E}[\text{Var}(X|Y)]$$

provided that the expectation exist. Proof:

$$\begin{aligned} \text{Var}(X) &= \mathbf{E}X \left\{ [X - \mathbf{E}(X|Y) + \mathbf{E}(X|Y) - \mathbf{E}]^2 \right\} \\ &= \mathbf{E} \left\{ [X - \mathbf{E}(X|Y)]^2 + [\mathbf{E}(X|Y) - \mathbf{E}X]^2 \right. \\ &\quad \left. + 2[X - \mathbf{E}(X|Y)][\mathbf{E}(X|Y) - \mathbf{E}X] \right\} \\ &= \mathbf{E}\{[X - \mathbf{E}(X|Y)]^2\} + \mathbf{E}\{[\mathbf{E}(X|Y) - \mathbf{E}X]^2\} \\ &= \mathbf{E}[\text{Var}(X|Y)] + \text{Var}[\mathbf{E}(X|Y)] \end{aligned}$$

For any two random variables X and Y

$$\text{Var}(X) = \text{Var}[\mathbf{E}(X|Y)] + \mathbf{E}[\text{Var}(X|Y)]$$

provided that the expectation exist. Proof:

$$\begin{aligned} \text{Var}(X) &= \mathbf{E}X\left\{[X - \mathbf{E}(X|Y) + \mathbf{E}(X|Y) - \mathbf{E}X]^2\right\} \\ &= \mathbf{E}\left\{[X - \mathbf{E}(X|Y)]^2 + [\mathbf{E}(X|Y) - \mathbf{E}X]^2\right. \\ &\quad \left.+ 2[X - \mathbf{E}(X|Y)][\mathbf{E}(X|Y) - \mathbf{E}X]\right\} \\ &= \mathbf{E}\{[X - \mathbf{E}(X|Y)]^2\} + \mathbf{E}\{[\mathbf{E}(X|Y) - \mathbf{E}X]^2\} \\ &= \mathbf{E}[\text{Var}(X|Y)] + \text{Var}[\mathbf{E}(X|Y)] \end{aligned}$$

$$\mathbf{E}\left\{2[X - \mathbf{E}(X|Y)][\mathbf{E}(X|Y) - \mathbf{E}X]\right\} = \mathbf{E}[\mathbf{E}(Z|Y)]$$

Real data

- The Current Population Survey (CPS) is a monthly survey of about 57,000 U.S. households conducted by the Bureau of the Census of the Bureau of Labor Statistics [cps09mar]².
- female: 1 if female, 0 otherwise
- earnings: total annual wage and salary earnings
- hours: number of hours worked per week
- week: number of weeks worked per year

cps09mar

```
1 > rm(list=ls())
2 > dat20 <- read.table("cps09mar.txt", head=TRUE, fileEncoding="utf8")
3 > head(dat20)
4 age female hisp education earnings hours week union uncov region race
5 1 52      0    0        12    146000    45   52     0     0        1    1
6 2 38      0    0        18     50000    45   52     0     0        1    1
7 3 38      0    0        14     32000    40   51     0     0        1    1
8 4 41      1    0        13     47000    40   52     0     0        1    1
9 5 42      0    0        13    161525    50   52     1     0        1    1
10 6 66      1    0        13     33000    40   52     0     0        1    1
```

²<https://www.ssc.wisc.edu/bhansen/econometrics/>

- $X = \frac{\text{earnings}}{\text{hours*week}}$
- $Y = \text{female}$
- Homework: **Write R code to check the following equations.**

$$\mathbf{E}[\mathbf{E}(X|Y)] = \mathbf{E}(X)$$

$$\text{Var}(X) = \text{Var}[\mathbf{E}(X|Y)] + \mathbf{E}[\text{Var}(X|Y)]$$

```

1 y <- dat20$earnings/(dat20$hours*dat20$week)
2 logy <- log(y)
3 plot(density(log(y)))
4 index_men <- which(dat20$female==0)
5 index_women <- which(dat20$female==1)
6 logy_mem <- log(y[index_men])
7 logy_women <- log(y[-index_men])
8 plot(density(logy_mem), ylim=c(0, 0.8), pch=1)
9 points(density(logy_women), col="red", lty=2, pch=3)
10 legend(-8, 0.6, c("man", "woman"), col = c(1,2), lty = c(1, 2),
11        pch = c(1, 3))

```

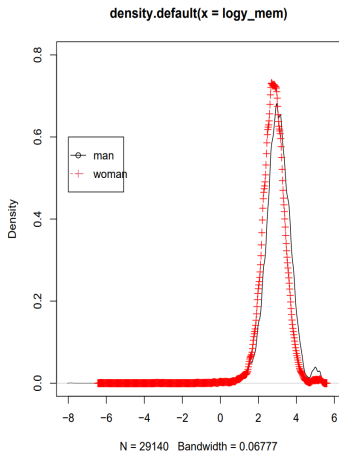


Figure: The density of $Y|female$

cps09mar

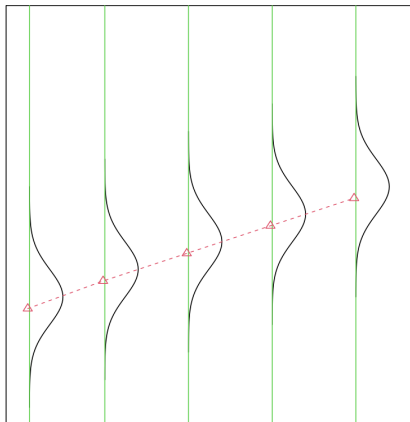
```

1 > mean(logy)
2 [1] 2.946185
3 > var(logy)
4 [1] 0.456827
5 > pp <- c((1-mean(dat20$female)), mean(dat20$female))
6 > mean_menwomean <- c( mean(logy_mem), mean(logy_women))
7 > var_menwomean <- c( var(logy_mem), var(logy_women))
8 > resu <- data.frame(pp, mean_menwomean, var_menwomean)
9 > rownames(resu) <- c("men", "women")
10 > resu
11           pp  mean_menwomean  var_menwomean
12 men    0.5742777         3.045938      0.4956187
13 women  0.4257223         2.811624      0.3729886

```

讨论

- 分析 female 对 Y 的影响
- 如果 female 是连续变量



R

```

1 y <- seq(-4,4,0.1)
2 x <- dnorm(y)
3 plot((x+0.1),y,xlim=c(0,4.5), ylim=c(-4,10),xlab="",
4 type="l", ylab="", xaxt="n", yaxt="n")
5 abline(v=min(x+0.1), col=3)
6 points((x+1),(y+1),type="l")
7 abline(v=min(x+1), col=3)
8 points((x+2),(y+2),type="l")
9 abline(v=min(x+2), col=3)
10 points((x+3),(y+3),type="l")
11 abline(v=min(x+3), col=3)
12 points((x+4),(y+4),type="l")
13 abline(v=min(x+4), col=3)
14 index <- which(x==max(x))
15 meanx <- x[index] + c(0.1, seq(1:4))
16 meany <- y[index] + c(0, seq(1:4))
17 points(meanx-0.42 ,meany-0.42, pch=2, col=2, type="o", lty=2)

```

一元线性回归模型

一元线性回归模型的数学形式为：

$$y = \beta_0 + \beta_1 x + \epsilon \quad (1)$$

通常假定：

$$\begin{cases} E(\epsilon) = 0 \\ \text{var}(\epsilon) = \sigma^2 \end{cases}$$

对式 (1) 两端求条件期望，得到回归方程：

$$E(y|x) = \beta_0 + \beta_1 x$$

一元线性回归模型经验方程

如果获得 n 组样本观测值 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, 样本模型:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, 2, \dots, n \quad (2)$$

满足:

$$\begin{cases} E(\epsilon_i) = 0 \\ \text{var}(\epsilon_i) = \sigma^2 \end{cases} \quad i = 1, 2, \dots, n$$

对 (2) 两端分别求期望和方差, 得到:

$$E(y_i) = \beta_0 + \beta_1 x_i, \quad \text{var}(y_i) = \sigma^2, \quad i = 1, 2, \dots, n$$

$E(y_i) = \beta_0 + \beta_1 x_i$ 从平均意义上表达了变量 y 与 x 的统计规律性。

用 $\hat{\beta}_0, \hat{\beta}_1$ 分别表示 β_0, β_1 的估计值, 获得 y 关于 x 的一元线性经验回归方程

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$$

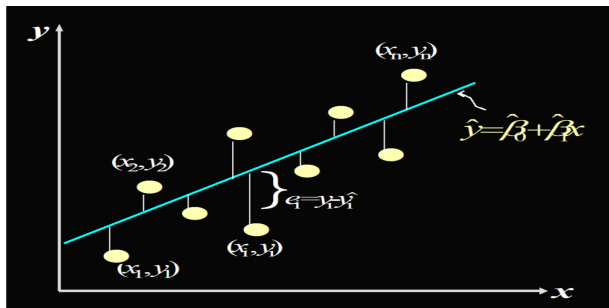
普通最小二乘估计

普通最小二乘估计 (Ordinary Least Square Estimation, 简记为 OLSE) 就是寻找参数 β_0, β_1 的估计值使离差平方和达到极小

$$\begin{aligned} Q(\hat{\beta}_0, \hat{\beta}_1) &= \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 \\ &= \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 \end{aligned}$$

- $\hat{\beta}_0, \hat{\beta}_1$ 称为 β_0, β_1 的最小二乘估计
- $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ 为 y_i 的**回归拟合值**，简称回归值或拟合值
- $e_i = y_i - \hat{y}_i$ 为 y_i 的残差。

普通最小二乘估计



- 残差平方和 $\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$
- 整体上刻画了 n 个样本观测点 $(x_i, y_i), i = 1, \dots, n$ 到回归直线 $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ 距离的长短。

普通最小二乘估计

$$\begin{cases} \frac{\partial Q}{\partial \beta_0} \big|_{\beta_0=\hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0 \\ \frac{\partial Q}{\partial \beta_1} \big|_{\beta_1=\hat{\beta}_1} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) x_i = 0 \end{cases}$$

可以得到残差的性质：

$$\begin{cases} \sum_{i=1}^n e_i = 0 \\ \sum_{i=1}^n x_i e_i = 0 \end{cases}$$

整理后得到正规方程组

$$\begin{cases} n\hat{\beta}_0 + \left(\sum_{i=1}^n x_i\right)\hat{\beta}_1 = \sum_{i=1}^n y_i \\ \left(\sum_{i=1}^n x_i\right)\hat{\beta}_0 + \left(\sum_{i=1}^n x_i^2\right)\hat{\beta}_1 = \sum_{i=1}^n x_i y_i \end{cases}$$

普通最小二乘估计

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{cases}$$

其中 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ 。

记：

$$L_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n(\bar{x})^2$$

$$L_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n(\bar{x}\bar{y})$$

则：

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{L_{xy}}{L_{xx}} \end{cases}$$

例 2 Homework R code

$$\bar{x} = \frac{49.2}{15} = 3.28, \quad \bar{y} = \frac{396.2}{15} = 26.413$$

$$L_{xx} = \sum_{i=1}^n x_i^2 - n(\bar{x})^2 = 196.16 - 15 \times (3.28)^2 = 34.784$$

$$L_{xy} = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} = 1470.65 - 1299.536 = 171.114$$

得到：

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 26.413 - 4.919 \times 3.28 = 10.279 \\ \hat{\beta}_1 = L_{xy}/L_{xx} = 171.114/34.784 = 4.919 \end{cases}$$

于是得到回归方程：

$$\hat{y} = 10.275 + 4.919x$$

最大似然估计

最大似然估计 (maximum likelihood estimation, 简记为 MLE) 是利用总体的分布密度或概率分布的表达式及其样本所提供的信息求未知参数估计量的一种方法。似然函数并不局限于独立同分布的样本。

- 连续型随机变量：似然函数是样本的联合密度函数
- 离散型随机变量：似然函数是样本的联合概率函数

对于一元线性回归模型参数的最大似然估计，如果已经得到样本观测值 $(x_i, y_i), i = 1, \dots, n$ ，那么在假设 $\epsilon_i \sim N(0, \sigma^2)$ 时， y_i 服从如下正态分布：

$$y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

y_i 的分布密度为：

$$f_i(y_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2\sigma^2} [y_i - (\beta_0 + \beta_1 x_i)]^2 \right\}$$

y_1, y_2, \dots, y_n 的似然函数为：

$$\begin{aligned} L(\beta_0, \beta_1, \sigma^2) &= \prod_{i=1}^n f_i(y_i) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 \right\} \end{aligned}$$

取对数似然函数为：

$$\ln(L) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2$$

至此与最小二乘原理相同。

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 \text{(有偏估计)}$$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n-2} \sum_{i=1}^n [y_i - (\beta_0 + \beta_1 x_i)]^2 \text{(无偏估计)}$$

最小二乘估计的性质

- 线性

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} y_i$$

- 无偏性

$$E(Y_i) = \beta_0 + \beta_1 x_i$$

$$\begin{aligned} E(\hat{\beta}_1) &= \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} E(y_i) \\ &= \sum_{i=1}^n \frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} (\beta_0 + \beta_1 x_i) = \beta_1 \end{aligned}$$

其中用到 $\sum (x_i - \bar{x}) = 0$, $\sum (x_i - \bar{x}) x_i = \sum (x_i - \bar{x})^2$

最小二乘估计的性质

- β_0, β_1 的方差

$$\text{var}(\hat{\beta}_1) = \sum_{i=1}^n \left[\frac{x_i - \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^2 \text{var}(y_i) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{var}(\hat{\beta}_0) = \left[\frac{1}{n} + \frac{(\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \sigma^2$$

$$\text{cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{\bar{x}}{L_{xx}} \sigma^2$$

$$\hat{\beta}_0 \sim N\left(\beta_0, \left(\frac{1}{n} + \frac{(\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) \sigma^2\right)$$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{L_{xx}}\right)$$

高斯-马尔可夫条件

$$\begin{cases} E(\epsilon_i) = 0, & i = 1, 2, \dots, n \\ \text{cov}(\epsilon_i, \epsilon_j) = \begin{cases} \sigma^2, & i = j \\ 0, & i \neq j \end{cases} & i, j = 1, 2, \dots, n \end{cases}$$

t 检验

t 检验用于检验回归系数的显著性，检验的原假设是：

$$H_0: \beta_1 = 0$$

对立假设是：

$$H_1: \beta_1 \neq 0$$

由

$$\hat{\beta}_1 \sim N\left(\beta_0, \frac{\sigma^2}{L_{xx}}\right)$$

当原假设 $H_0: \beta_1 = 0$ 成立时，有

$$\hat{\beta}_1 \sim N\left(0, \frac{\sigma^2}{L_{xx}}\right)$$

构造 t 统计量

$$t = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2/L_{xx}}} = \frac{\hat{\beta}_1 \sqrt{L_{xx}}}{\hat{\sigma}}$$

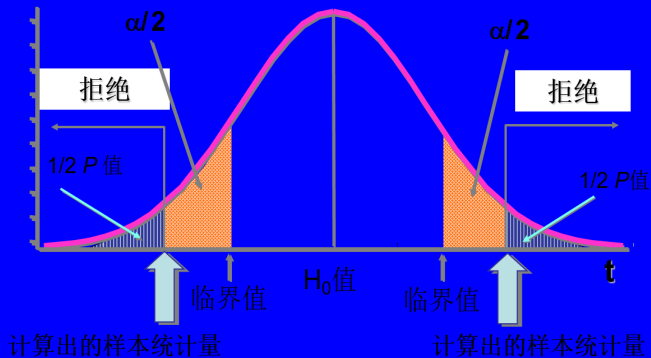
其中

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n e_i^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

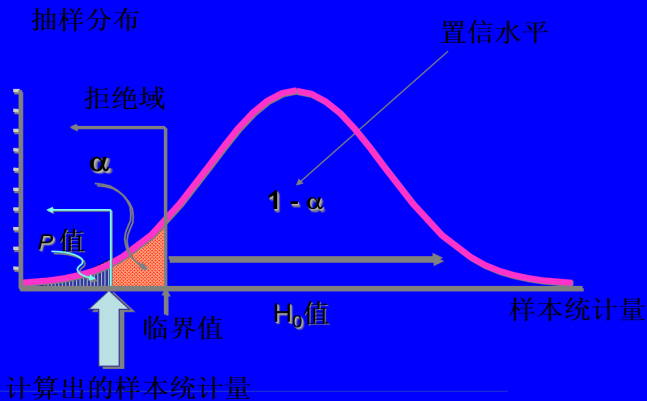
P - value

- p 值即显著性概率值 Significance Probability Value
- 是当原假设为真时得到目前的样本以及更极端样本的概率，所谓极端就是与原假设相背离
- 它是用此样本拒绝原假设所犯弃真错误的真实概率，被称为观察到的 (或实测的) 显著性水平

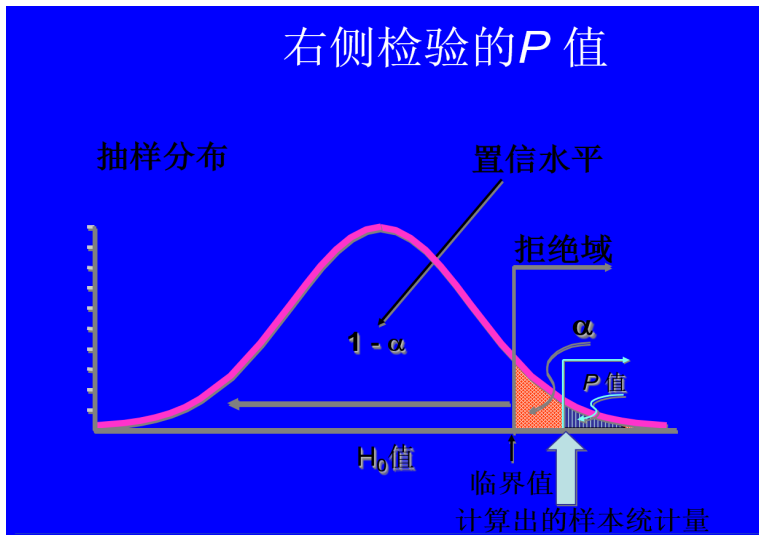
双侧检验的 P 值



左侧检验的 P 值



右侧检验的 P 值



利用 p 值进行检验的决策准则

- 若 p 值 $> \alpha$, 不能拒绝 H_0
- 若 p 值 $< \alpha$, 拒绝 H_0

双侧检验 p 值 $= 2 \times$ 单侧检验 p 值

F 检验

F 检验是根据平方和分解式，直接从回归效果检验回归方程的显著性。

平方和分解式是：

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

$$SST = SSR + SSE$$

总的离差平方和： $SST = \sum_{i=1}^n (y_i - \bar{y})^2$,

回归平方和： $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$,

残差平方和： $SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$ 。

F 检验

构造 F 统计量如下

$$F = \frac{SSR/1}{SSE/(n-2)}$$

图：一元线性回归方差分析表

方差来源	自由度	平方和	均方	F值	P值
回归	1	SSR	SSR/1	$\frac{SSR/1}{SSE/(n-2)}$	$P(F>F值)$ =P值
残差	$n-2$	SSE	SSE/ (n-2)		
总和	$n-1$	SST			

相关系数的显著性检验

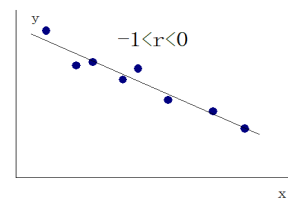
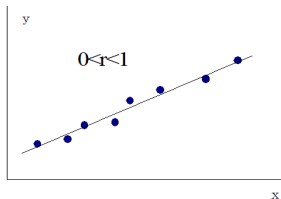
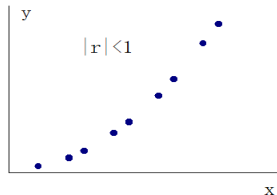
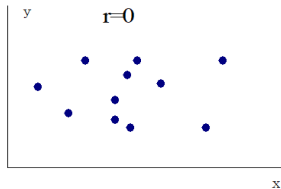
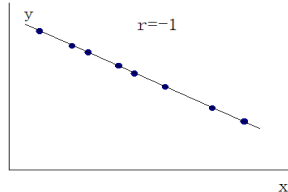
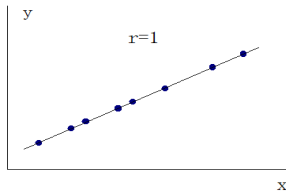
相关系数的显著性检验

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

$$= \frac{L_{xx}}{\sqrt{L_{xx} L_{yy}}} = \hat{\beta}_1 \sqrt{\frac{L_{xx}}{L_{yy}}}$$

r 为 x 与 y 的简单相关系数，简称相关系数。

相关系数直观意义图



相关系数的显著性检验

两变量间相关程度的强弱分为以下几个等级：

- 当 $|r| \geq 0.8$ 时，视为高度相关；
- 当 $0.5 \leq |r| \leq 0.8$ ，视为低度相关；
- 当 $0.3 \leq |r| \leq 0.5$ 时，视为低度相关；
- 当 $|r| < 0.3$ 时，表明两个变量之间的相关程度极弱，在实际应用中可视为不相关。

三种检验的关系

对于一元线性回归，这三种检验的结果是完全一致的。

$$H_0: \beta = 0 \quad t = \frac{\hat{\beta}_1}{\sqrt{\hat{\sigma}^2/L_{xx}}} = \frac{\hat{\beta}_1\sqrt{L_{xx}}}{\hat{\sigma}}$$

$$H_0: \rho = 0 \quad t = \frac{\sqrt{n-2}r}{\sqrt{1-r^2}}$$

$$H_0: \text{回归无效} \quad F = \frac{SSR/1}{SSE/(n-2)}$$

决定系数

回归平方和与总离差平方和之比定义为决定系数，也称为判定系数、确定系数。记为 r^2

$$r^2 = \frac{SSR}{SST} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

可以证明

$$r^2 = \frac{SSR}{SST} = \frac{L_{xy}^2}{L_{xx}L_{yy}} = (r)^2$$

决定系数 r^2 是一个反映直线与样本观测拟合优度的相对指标，是因变量的变异中能用自变量解释的比例。其数值在 0 ~ 1 之间，可以用百分数表示。

残差概念与残差图

残差: $e_i = y_i - \hat{y}_i = y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i$

误差: $\epsilon_i = y_i - \beta_0 - \beta_1 x_i$

残差项 e_i 是误差项 ϵ 的估计值。

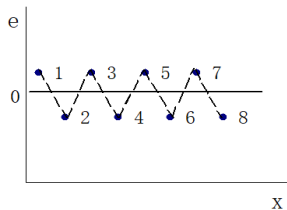
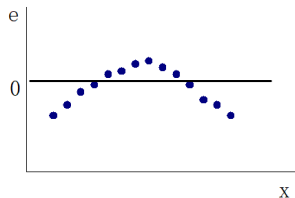
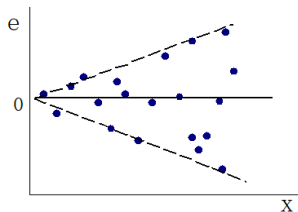
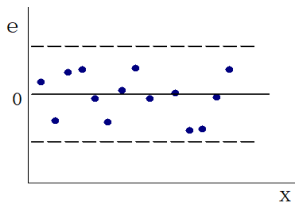


图: 残差图

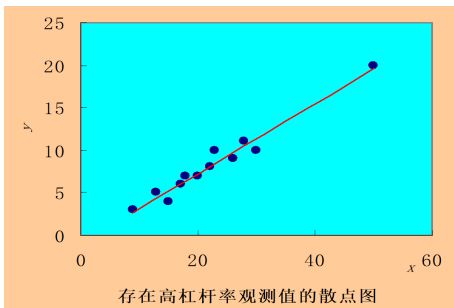
残差的性质

- 性质 1: $E(e_i) = 0$ 证明:

$$E(e_i) = E(y_i) - E(\hat{y}_i) = (\beta_0 + \beta_1 x_i) - (\hat{\beta}_0 + \hat{\beta}_1 x_i) = 0$$

- 性质 2: $var(e_i) = \left[1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{L_{xx}}\right] \sigma^2 = (1 - h_{ii}) \sigma^2$

其中 $h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{L_{xx}}$ 称为杠杆值。



残差的性质

- 性质 3: 残差满足约束条件

$$\begin{cases} \sum_{i=1}^n e_i = 0 \\ \sum_{i=1}^n x_i e_i = 0 \end{cases}$$

这表明残差 e_1, e_2, \dots, e_n 的相关的, 不是独立的。

改进的残差

标准化残差：

$$ZRE_i = \frac{e_i}{\hat{\sigma}}$$

学生化残差：

$$SRE_i = \frac{e_i}{\hat{\sigma}\sqrt{1-h_{ii}}}$$

回归系数的区间估计

由 $\hat{\beta}_1 \sim N(\beta_0, \frac{\sigma^2}{L_{xx}})$ 可得

$$t = \frac{\hat{\beta}_1 - \beta_1}{\sqrt{\hat{\sigma}^2/L_{xx}}} = \frac{(\hat{\beta}_1 - \beta_1)\sqrt{L_{xx}}}{\hat{\sigma}}$$

服从自由度为 $n-2$ 的 t 分布

$$P(|\frac{(\hat{\beta}_1 - \beta_1)\sqrt{L_{xx}}}{\hat{\sigma}}| < t_{\alpha/2}(n-2)) = 1 - \alpha$$

上式等价于

$$P(\hat{\beta}_1 - t_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{L_{xx}}} < \beta_1 < \hat{\beta}_1 + t_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{L_{xx}}}) = 1 - \alpha$$

即得到 β_1 的置信度为 $1 - \alpha$ 的置信区间为：

$$(\hat{\beta}_1 - t_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{L_{xx}}}, \hat{\beta}_1 + t_{\alpha/2} \frac{\hat{\sigma}}{\sqrt{L_{xx}}})$$

单值预测

单值预测就是用单个值作为因变量新值的预测值。

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$$E(\hat{y}_0) = E(y) = \beta_0 + \beta_1 x_0$$

区间预测

区间预测就是对于给定的显著水平 α ，找到一个区间 (T_1, T_2) ，使对应于某特定的 x_0 的实际值 y_0 以 $1 - \alpha$ 的概率被区间 (T_1, T_2) 包含，用公示表示：

$$P(T_1 < y_0 < T_2) = 1 - \alpha$$

一、因变量新值的区间预测

首先要给出估计值 $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ 的分布，在正态性假设下

$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ 服从正态分布，其期望为 $E(\hat{y}_0) = \beta_0 + \beta_1 x_0$ ，计算其方差

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0 = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_0 = \sum_{i=1}^n \left[\frac{1}{n} + \frac{(x_i - \bar{x})(x_0 - \bar{x})}{L_{xx}} \right] y_i$$

$$\text{var}(\hat{y}_0) = \sum_{i=1}^n \left[\frac{1}{n} + \frac{(x_i - \bar{x})(x_0 - \bar{x})}{L_{xx}} \right] \text{var}(y_i) = \left[\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}} \right] \sigma^2$$

从而得

$$\hat{y}_0 \sim N(\beta_0 + \beta_1 x_0, \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}} \right) \sigma^2)$$

因变量新值的区间预测

$$\text{记 } h_{00} = \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}}$$

$$\hat{y}_0 \sim N(\beta_0 + \beta_1 x_0, h_{00} \sigma^2)$$

y_0 与 \hat{y}_0 独立

$$\text{var}(y_0 - \hat{y}_0) = \text{var}(y_0) + \text{var}(\hat{y}_0) = \sigma^2 + h_{00} \sigma^2$$

于是

$$y_0 - \hat{y}_0 \sim N(0, (1 + h_{00}) \sigma^2)$$

进而可知统计量

$$t = \frac{y_0 - \hat{y}_0}{\sqrt{1 + h_{00}} \hat{\sigma}} \sim t(n - 2)$$

因变量新值的区间预测

$$P(|\frac{y_0 - \hat{y}_0}{\sqrt{1 + h_{00}}\hat{\sigma}}| \leq t_{\alpha/2}(n-2)) = 1 - \alpha$$

由此可以求得 y_0 的置信度为 $1 - \alpha$ 的置信区间为：

$$\hat{y}_0 \pm t_{\alpha/2}(n-2)\sqrt{1 + h_{00}}\hat{\sigma}$$

当样本量 n 较大， $|x_0 - \bar{x}|$ 较小时， h_{00} 接近 0， y_0 的置信度为 95% 的置信区间近似为：

$$\hat{y}_0 \pm 2\hat{\sigma}$$

因变量新值的平均值的区间预测

由于 $E(y_0) = \beta_0 + \beta_1 x_0$ 是常数,

$$\hat{y}_0 - E(y_0) \sim N(\beta_0 + \beta_1 x_0, (\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}}) \sigma^2)$$

可得置信水平为 $1 - \alpha$ 的置信区间为

$$\hat{y}_0 \pm t_{\alpha/2}(n-2) \sqrt{h_{00}} \hat{\sigma}$$

控制问题

控制问题相当于预测问题的反问题。给定 y 的预期范围 (T_1, T_2) , 如何控制自变量 x 的值才能以 $1 - \alpha$ 的概率保证

$$P(T_1 < y < T_2) = 1 - \alpha$$

通常用近似的预测区间来确定 x

$$\begin{cases} \hat{y}(x) - 2\hat{\sigma} > T_1 \\ \hat{y}(x) + 2\hat{\sigma} < T_2 \end{cases}$$

把 $\hat{y}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$ 代入求得
当 $\hat{\beta}_1 > 0$ 时

$$\frac{T_1 + 2\hat{\sigma} - \hat{\beta}_0}{\hat{\beta}_1} < x < \frac{T_2 - 2\hat{\sigma} - \hat{\beta}_0}{\hat{\beta}_1}$$

当 $\hat{\beta}_1 < 0$ 时

$$\frac{T_2 - 2\hat{\sigma} - \hat{\beta}_0}{\hat{\beta}_1} < x < \frac{T_1 + 2\hat{\sigma} - \hat{\beta}_0}{\hat{\beta}_1}$$

思考题

Anscombe(1973)³构造了四组数据集，他们具有相同的相关系数。

```

1 > anscombe <- read.table ('anscombe.txt',header = T)
2 > head(anscombe)
3      Y1 X1    Y2 X2      Y3 X3    Y4 X4
4 1 8.04 10 9.14 10  7.46 10 6.58  8
5 2 6.95  8 8.14  8  6.77  8 5.76  8
6 3 7.58 13 8.74 13 12.74 13 7.71  8
7 4 8.81  9 8.77  9  7.11  9 8.84  8
8 5 8.33 11 9.26 11  7.81 11 8.47  8
9 6 9.96 14 8.10 14  8.84 14 7.04  8

```

- 画出散点图
- 建立线性模型
- 画出拟合图

³<https://github.com/xuejunma/ar2021/blob/main/anscombe.txt>

作业

1 小作业

- p.50 2.8, 2.11
- p.51 编程计算

2 大作业

考虑没有截距项的模型：

$$y_i = \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n$$

- β_1 的估计式
- $SST = SSR + SSE$ 是否成立，说出理由。
- R^2 一定在大于零吗，说出理由。