

# 双变量线性回归

教师 黄光辉  
hgh@cqu.edu.cn

## 1 第六章双变量回归:区间估计和假设检验

- 统计学预备知识
- 回归参数的区间估计
- 回归系数的假设检验
- 回归模型的方差分析

# 检验,一把双刃剑.

警惕过度地检验假设.

# 检验,一把双刃剑.

警惕过度地检验假设.

你对数据愈苛求,数据就会愈多的向你供认;

# 检验,一把双刃剑.

警惕过度地检验假设.

你对数据愈苛求,数据就会愈多的向你供认;

但是在威逼下得到的供词,在科学询查的法庭上是不容许的.

## 1 第六章双变量回归:区间估计和假设检验

- 统计学预备知识
- 回归参数的区间估计
- 回归系数的假设检验
- 回归模型的方差分析

# 双变量模型的分布假定

$$Y = \beta_0 + \beta_1 X + U$$

# 双变量模型的分布假定

$$Y = \beta_0 + \beta_1 X + U$$

独立重复观测 $(Y, X)$ ,得到观测数据

$$(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n).$$

干扰项无法直接观测,假设分布为:

$U_i \sim N(0, \sigma)$ ,  $\{U_i\}$ 序列独立同分布.



# 双变量模型的分布假定

$$Y = \beta_0 + \beta_1 X + U$$

独立重复观测 $(Y, X)$ , 得到观测数据

$$(Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n).$$

干扰项无法直接观测, 假设分布为:

$U_i \sim N(0, \sigma)$ ,  $\{U_i\}$ 序列独立同分布.

采用最小二乘法, 从上述数据得到模型参数的点估计为:

$$\begin{aligned}\widehat{\beta}_0 &= \bar{Y} - \widehat{\beta}_1 \bar{X} \\ \widehat{\beta}_1 &= \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \\ &= \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \\ &\text{记号 } x_i = X_i - \bar{X}\end{aligned}$$

# 参数的分布

$$Y = \beta_0 + \beta_1 X + U$$

$$\widehat{\beta}_0 \sim N\left(\beta_0, \sigma^2 \left(\frac{1}{n} + \frac{\bar{X}}{\sum (X_i - \bar{X})^2}\right)\right)$$

$$\widehat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{\sum (X_i - \bar{X})^2}\right)$$

$$Z_1 = \frac{\widehat{\beta}_0 - \beta_0}{\sigma_{\widehat{\beta}_0}} \sim N(0, 1), \quad Z_2 = \frac{\widehat{\beta}_1 - \beta_1}{\sigma_{\widehat{\beta}_1}} \sim N(0, 1)$$

$$\widehat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \widehat{u}_i^2 \sim \chi^2(n-2).$$

## 1 第六章双变量回归:区间估计和假设检验

- 统计学预备知识
- 回归参数的区间估计
- 回归系数的假设检验
- 回归模型的方差分析

## $\beta_1$ 的置信区间估计: t统计量

$$\widehat{\beta}_1 \sim N(\beta_1, \sigma^2 \widehat{\beta}_1)$$

## $\beta_1$ 的置信区间估计: t统计量

$$\widehat{\beta}_1 \sim N(\beta_1, \sigma_{\widehat{\beta}_1}^2)$$

$$Z = \frac{\widehat{\beta}_1 - \beta_1}{\sigma_{\widehat{\beta}_1}} = \frac{(\widehat{\beta}_1 - \beta_1)\sqrt{\sum x_i^2}}{\sigma} \sim N(0, 1)$$

## $\beta_1$ 的置信区间估计: t统计量

$$\widehat{\beta}_1 \sim N(\beta_1, \sigma_{\widehat{\beta}_1}^2)$$

$$Z = \frac{\widehat{\beta}_1 - \beta_1}{\sigma_{\widehat{\beta}_1}} = \frac{(\widehat{\beta}_1 - \beta_1)\sqrt{\sum x_i^2}}{\sigma} \sim N(0, 1)$$

$$Z_3 = (n - 2) \frac{\widehat{\sigma}^2}{\sigma^2} \sim \chi(n - 2)$$

## $\beta_1$ 的置信区间估计: t统计量

$$\widehat{\beta}_1 \sim N(\beta_1, \sigma_{\widehat{\beta}_1}^2)$$

$$Z = \frac{\widehat{\beta}_1 - \beta_1}{\sigma_{\widehat{\beta}_1}} = \frac{(\widehat{\beta}_1 - \beta_1)\sqrt{\sum x_i^2}}{\sigma} \sim N(0, 1)$$

$$Z_3 = (n-2)\frac{\widehat{\sigma}^2}{\sigma^2} \sim \chi(n-2)$$

$$t = \frac{Z}{\sqrt{Z_3/(n-2)}}$$

## $\beta_1$ 的置信区间估计: t统计量

$$\widehat{\beta}_1 \sim N(\beta_1, \sigma_{\widehat{\beta}_1}^2)$$

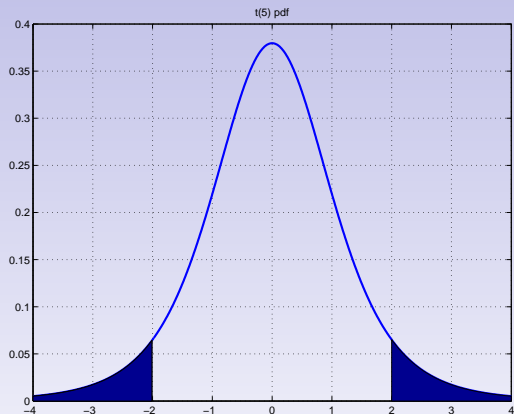
$$Z = \frac{\widehat{\beta}_1 - \beta_1}{\sigma_{\widehat{\beta}_1}} = \frac{(\widehat{\beta}_1 - \beta_1)\sqrt{\sum x_i^2}}{\sigma} \sim N(0, 1)$$

$$Z_3 = (n-2)\frac{\widehat{\sigma}^2}{\sigma^2} \sim \chi(n-2)$$

$$t = \frac{Z}{\sqrt{Z_3/(n-2)}} = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$



# $t$ 分布密度函数, $1 - \alpha$ 置信区间



$1 - \alpha$ 置信区间

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq t \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

# $\beta_1$ 的置信区间估计

$$t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n - 2)$$

## $\beta_1$ 的置信区间估计

$$t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

## $\beta_1$ 的置信区间估计

$$t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$$P \left\{ \widehat{\beta}_1 - se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2) \leq \beta_1 \leq \widehat{\beta}_1 + se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

## $\beta_1$ 的置信区间估计

$$t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$$P \left\{ \widehat{\beta}_1 - se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2) \leq \beta_1 \leq \widehat{\beta}_1 + se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$\beta_1$ 的置信区间为:

$$\left( \widehat{\beta}_1 - se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2), \widehat{\beta}_1 + se(\widehat{\beta}_1)t_{1-\frac{\alpha}{2}}(n-2) \right)$$

# $\beta_0$ 的置信区间估计

$$t = \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \sim t(n - 2)$$

## $\beta_0$ 的置信区间估计

$$t = \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

## $\beta_0$ 的置信区间估计

$$t = \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$$P \left\{ \widehat{\beta}_0 - se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2) \leq \beta_0 \leq \widehat{\beta}_0 + se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$



## $\beta_0$ 的置信区间估计

$$t = \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \sim t(n-2)$$

$$P \left\{ -t_{1-\frac{\alpha}{2}}(n-2) \leq \frac{\widehat{\beta}_0 - \beta_0}{se(\widehat{\beta}_0)} \leq t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$$P \left\{ \widehat{\beta}_0 - se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2) \leq \beta_0 \leq \widehat{\beta}_0 + se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2) \right\} = 1 - \alpha$$

$\beta_0$ 的置信区间为:

$$\left( \widehat{\beta}_0 - se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2), \widehat{\beta}_0 + se(\widehat{\beta}_0)t_{1-\frac{\alpha}{2}}(n-2) \right)$$

# 干扰项方差 $\sigma^2$ 的区间估计

$$\chi^2 = (n-2) \frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2(n-2)$$

# 干扰项方差 $\sigma^2$ 的区间估计

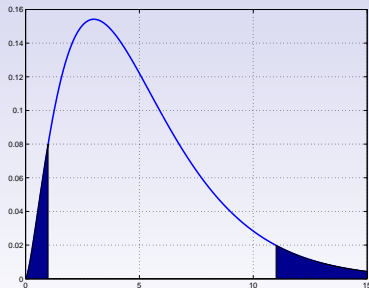
$$\chi^2 = (n-2) \frac{\widehat{\sigma^2}}{\sigma^2} \sim \chi^2(n-2)$$

$$\widehat{\sigma^2} = \frac{1}{n-2} \sum_{i=1}^n \widehat{u}_i^2$$

# 干扰项方差 $\sigma^2$ 的区间估计

$$\chi^2 = (n-2) \frac{\widehat{\sigma}^2}{\sigma^2} \sim \chi^2(n-2)$$

$$\widehat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \widehat{u}_i^2$$



## $\sigma^2$ 的置信区间

$$P \left\{ \chi_{\frac{\alpha}{2}}^2 \leq \frac{n-2}{\sigma^2} \widehat{\sigma}^2 \leq \chi_{1-\frac{\alpha}{2}}^2 \right\} = 1 - \alpha$$

## $\sigma^2$ 的置信区间

$$P \left\{ \chi_{\frac{\alpha}{2}}^2 \leq \frac{n-2}{\sigma^2} \widehat{\sigma}^2 \leq \chi_{1-\frac{\alpha}{2}}^2 \right\} = 1 - \alpha$$
$$\left( \frac{(n-2)\widehat{\sigma}^2}{\chi_{1-\frac{\alpha}{2}}^2}, \frac{(n-2)\widehat{\sigma}^2}{\chi_{\frac{\alpha}{2}}^2} \right)$$

## 1 第六章双变量回归:区间估计和假设检验

- 统计学预备知识
- 回归参数的区间估计
- 回归系数的假设检验
- 回归模型的方差分析

# 统计量的观测值

$\beta_0$ 和 $\beta_1$ 的估计值为 $\widehat{\beta}_0$ 和 $\widehat{\beta}_1$ 。



# 统计量的观测值

$\beta_0$ 和 $\beta_1$ 的估计值为 $\widehat{\beta}_0$ 和 $\widehat{\beta}_1$ 。  
统计量的观测值通常出现在什么地方呢？

# 统计量的观测值

$\beta_0$ 和 $\beta_1$ 的估计值为 $\widehat{\beta}_0$ 和 $\widehat{\beta}_1$ 。

统计量的观测值通常出现在什么地方呢？

由大数定律，统计量会出现在自己中心附近区域，越远离中心，出现的可能性越小。

从而，我们可以断定

$\widehat{\beta}_0$ 总是出现在真实的 $\beta_0$ 附近；

$\widehat{\beta}_1$ 总是出现在真实的 $\beta_1$ 附近。

# 统计量的观测值

$\beta_0$ 和 $\beta_1$ 的估计值为 $\widehat{\beta}_0$ 和 $\widehat{\beta}_1$ 。

统计量的观测值通常出现在什么地方呢？

由大数定律，统计量会出现在自己中心附近区域，越远离中心，出现的可能性越小。

从而，我们可以断定

$\widehat{\beta}_0$ 总是出现在真实的 $\beta_0$ 附近；

$\widehat{\beta}_1$ 总是出现在真实的 $\beta_1$ 附近。

由此，我们可以回答下列问题

一元线性回归的截距项是否为一个指定的常数 $\beta_0$

一元线性回归的斜率项是否为一个指定的常数 $\beta_1$ 。

## 参数异常观测的刻画：拒绝域和接受域

$$y = \beta_0 + \beta_1 x + u$$

$$\widehat{\beta}_1 \sim N\left(\beta_1, \widehat{\sigma}_{\widehat{\beta}_1}^2\right), \quad t = \frac{\widehat{\beta}_1 - \beta_1}{\text{se}(\widehat{\beta}_1)} \sim t(n-2)$$

## 参数异常观测的刻画：拒绝域和接受域

$$y = \beta_0 + \beta_1 x + u$$

$$\widehat{\beta}_1 \sim N\left(\beta_1, \widehat{\sigma}_{\widehat{\beta}_1}^2\right), \quad t = \frac{\widehat{\beta}_1 - \beta_1}{\widehat{se}(\widehat{\beta}_1)} \sim t(n-2)$$

$$P\left\{-t_{1-\frac{\alpha}{2}} \leq t \leq t_{1-\frac{\alpha}{2}}\right\} = P\left\{-t_{1-\frac{\alpha}{2}} \leq \frac{\widehat{\beta}_1 - \beta_1}{\widehat{se}(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}\right\} = 1 - \alpha$$

## 参数异常观测的刻画：拒绝域和接受域

$$y = \beta_0 + \beta_1 x + u$$

$$\widehat{\beta}_1 \sim N\left(\beta_1, \widehat{\sigma}_{\widehat{\beta}_1}^2\right), \quad t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$

$$P\left\{-t_{1-\frac{\alpha}{2}} \leq t \leq t_{1-\frac{\alpha}{2}}\right\} = P\left\{-t_{1-\frac{\alpha}{2}} \leq \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}\right\} = 1 - \alpha$$

$$P\left\{\widehat{\beta}_1 - t_{1-\frac{\alpha}{2}} se(\widehat{\beta}_1) \leq \beta_1 \leq \widehat{\beta}_1 + t_{1-\frac{\alpha}{2}} se(\widehat{\beta}_1)\right\} = 1 - \alpha$$

## 参数异常观测的刻画：拒绝域和接受域

$$y = \beta_0 + \beta_1 x + u$$

$$\widehat{\beta}_1 \sim N\left(\beta_1, \widehat{\sigma}_{\widehat{\beta}_1}^2\right), \quad t = \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \sim t(n-2)$$

$$P\left\{-t_{1-\frac{\alpha}{2}} \leq t \leq t_{1-\frac{\alpha}{2}}\right\} = P\left\{-t_{1-\frac{\alpha}{2}} \leq \frac{\widehat{\beta}_1 - \beta_1}{se(\widehat{\beta}_1)} \leq t_{1-\frac{\alpha}{2}}\right\} = 1 - \alpha$$

$$P\left\{\widehat{\beta}_1 - t_{1-\frac{\alpha}{2}} se\left(\widehat{\beta}_1\right) \leq \beta_1 \leq \widehat{\beta}_1 + t_{1-\frac{\alpha}{2}} se\left(\widehat{\beta}_1\right)\right\} = 1 - \alpha$$

$H_0$ : 真实的斜率是 $\beta_1$ ;  $H_1$ : 真实的斜率不是 $\beta_1$ . 接受域为:

$$\left(\widehat{\beta}_1 - t_{1-\frac{\alpha}{2}} se\left(\widehat{\beta}_1\right), \widehat{\beta}_1 + t_{1-\frac{\alpha}{2}} se\left(\widehat{\beta}_1\right)\right)$$

# 干扰项方差的假设检验

单侧检验:

$$H_0 : \sigma^2 = \sigma_0^2; H_1 : \sigma^2 > \sigma_0^2$$



# 干扰项方差的假设检验

单侧检验:

$$H_0 : \sigma^2 = \sigma_0^2; H_1 : \sigma^2 > \sigma_0^2$$

如果原假设成立, 那么

$$\chi^2 = (n - 2) \frac{\hat{\sigma}^2}{\sigma_0^2} \sim \chi^2(n - 2).$$

# 干扰项方差的假设检验

单侧检验:

$$H_0 : \sigma^2 = \sigma_0^2; H_1 : \sigma^2 > \sigma_0^2$$

如果原假设成立, 那么

$$\chi^2 = (n - 2) \frac{\hat{\sigma}^2}{\sigma_0^2} \sim \chi^2(n - 2).$$

支持 $H_1$ 的证据为

$$\frac{\hat{\sigma}^2}{\sigma_0^2} > c,$$

# 干扰项方差的假设检验

单侧检验:

$$H_0 : \sigma^2 = \sigma_0^2; H_1 : \sigma^2 > \sigma_0^2$$

如果原假设成立, 那么

$$\chi^2 = (n - 2) \frac{\hat{\sigma}^2}{\sigma_0^2} \sim \chi^2(n - 2).$$

支持 $H_1$ 的证据为

$$\frac{\hat{\sigma}^2}{\sigma_0^2} > c,$$

这个事件出现的概率为 $1 - \alpha$ , 从而我们有拒绝域

$$\left( \hat{\sigma}^2 : (n - 2) \frac{\hat{\sigma}^2}{\sigma_0^2} > \chi_{1-\alpha}^2 \right)$$

# 干扰项方差的假设检验

单侧检验:

$$H_0 : \sigma^2 = \sigma_0^2; H_1 : \sigma^2 > \sigma_0^2$$

如果原假设成立, 那么

$$\chi^2 = (n-2) \frac{\hat{\sigma}^2}{\sigma_0^2} \sim \chi^2(n-2).$$

支持 $H_1$ 的证据为

$$\frac{\hat{\sigma}^2}{\sigma_0^2} > c,$$

这个事件出现的概率为 $1 - \alpha$ , 从而我们有拒绝域

$$\left( \hat{\sigma}^2 : (n-2) \frac{\hat{\sigma}^2}{\sigma_0^2} > \chi_{1-\alpha}^2 \right)$$

也就是

$$\hat{\sigma}^2 > \frac{\sigma_0^2 \chi_{1-\alpha}^2}{(n-2)}$$

## 精确的显著性水平: p-value

为了判断一个假设是否显著, 选择统计量 $T(X_1, X_2, \dots, X_n)$ 。  
观测样本, 得到数据 $x_1, x_2, \dots, x_n$ , 带入统计量, 得到观测值  
 $t(x_1, x_2, \dots, x_n)$ 。

## 精确的显著性水平: p-value

为了判断一个假设是否显著, 选择统计量 $T(X_1, X_2, \dots, X_n)$ 。  
观测样本, 得到数据 $x_1, x_2, \dots, x_n$ , 带入统计量, 得到观测值  
 $t(x_1, x_2, \dots, x_n)$ 。

按照 $T$ 的分布, 计算出概率

$$P\{T > t\} = p$$

$p$ 就称为统计量 $T$ 的p值, 它表示在原假设成立时,  $T$ 统计量超过实际观测值的概率。

# 精确的显著性水平: p-value

为了判断一个假设是否显著, 选择统计量 $T(X_1, X_2, \dots, X_n)$ 。  
观测样本, 得到数据 $x_1, x_2, \dots, x_n$ , 带入统计量, 得到观测值  
 $t(x_1, x_2, \dots, x_n)$ 。

按照 $T$ 的分布, 计算出概率

$$P\{T > t\} = p$$

$p$ 就称为统计量 $T$ 的 $p$ 值, 它表示在原假设成立时,  $T$ 统计量超过实际观测值的概率。

$p$ 值大, 超过 $t$ 的事件容易发生, 所以原假设成立;

$p$ 值小, 超过 $t$ 的事件不容易发生, 所以原假设应该被拒绝。

# 精确的显著性水平: p-value

为了判断一个假设是否显著, 选择统计量 $T(X_1, X_2, \dots, X_n)$ 。  
观测样本, 得到数据 $x_1, x_2, \dots, x_n$ , 带入统计量, 得到观测值  
 $t(x_1, x_2, \dots, x_n)$ 。

按照 $T$ 的分布, 计算出概率

$$P\{T > t\} = p$$

$p$ 就称为统计量 $T$ 的 $p$ 值, 它表示在原假设成立时,  $T$ 统计量超过实际观测值的概率。

$p$ 值大, 超过 $t$ 的事件容易发生, 所以原假设成立;

$p$ 值小, 超过 $t$ 的事件不容易发生, 所以原假设应该被拒绝。

一般取0.001, 0.0001等较小的数作为 $p$ 值的临界值。



## 1 第六章双变量回归:区间估计和假设检验

- 统计学预备知识
- 回归参数的区间估计
- 回归系数的假设检验
- 回归模型的方差分析

# 回归模型的离差

用相对于中点的离差刻画数据的变异：  
因变量的总离差

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

回归直线的离差

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

残差的离差

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

# 回归模型的离差

用相对于中点的离差刻画数据的变异：  
因变量的总离差

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

回归直线的离差

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

残差的离差

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

离差平方和分解公式：

$$TSS = ESS + RSS$$

# 回归模型拟合优度判定：判定系数

总的离差平方和 = 回归直线的离差平方和 + 残差的离差平方和

# 回归模型拟合优度判定：判定系数

总的离差平方和 = 回归直线的离差平方和 + 残差的离差平方和

$$1 = \frac{ESS}{TSS} + \frac{RSS}{TSS}$$

# 回归模型拟合优度判定：判定系数

总的离差平方和 = 回归直线的离差平方和 + 残差的离差平方和

$$1 = \frac{ESS}{TSS} + \frac{RSS}{TSS}$$

$$r^2 = \frac{ESS}{TSS}$$

称为回归模型的判定系数。

样本相关系数：

$$r = \pm\sqrt{r^2}$$

# 回归模型显著性判定：方差分析

$$y = \beta_0 + \beta_1 x + u$$

问题：变量 $x$ 和 $y$ 之间确实存在线性关系么？

# 回归模型显著性判定：方差分析

$$y = \beta_0 + \beta_1 x + u$$

问题：变量 $x$ 和 $y$ 之间确实存在线性关系么？

$$H_0 : \beta_1 = 0; \quad H_1 : \beta_1 \neq 0.$$



# 回归模型显著性判定：方差分析

$$y = \beta_0 + \beta_1 x + u$$

问题：变量 $x$ 和 $y$ 之间确实存在线性关系么？

$$H_0 : \beta_1 = 0; \quad H_1 : \beta_1 \neq 0.$$

原假设成立时有

$$F = \frac{ESS}{RSS/(n-2)} = \frac{\widehat{\beta}_1^2 \sum \widehat{x}_i^2}{\sum \widehat{u}_i^2 / (n-2)} \sim F(1, n-2)$$

# 一元线性回归方差分析表

变异来源	$SS^*$	df (自由度)	$MSS^*$ (均方和)
源于回归 (ESS)	$\sum \hat{y}_i^2$	1	$\sum \hat{y}^2$
源于剩余 (RSS)	$\sum \hat{u}_i^2$	n-2	$\frac{1}{n-2} \sum \hat{u}_i^2$
TSS	$\sum y_i$	n-1	F

# 一元线性回归方差分析表

变异来源	$SS^*$	df (自由度)	$MSS^*$ (均方和)
源于回归 (ESS)	$\sum \hat{y}_i^2$	1	$\sum \hat{y}^2$
源于剩余 (RSS)	$\sum \hat{u}_i^2$	n-2	$\frac{1}{n-2} \sum \hat{u}_i^2$
TSS	$\sum y_i$	n-1	F

一元线性回归中，F检验和T检验的作用一样；

多元线性回归中，F检验用于线性结构检验，T检验用于单个系数检验。

# 回归分析的应用：均值预测

$$Y = \beta_0 + \beta_1 X + U$$

假定  $U \sim N(0, \sigma^2)$ , 那么  $Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$

$$Y_0 = \beta_0 + \beta_1 X_0 = E[Y|X = X_0].$$

# 回归分析的应用：均值预测

$$Y = \beta_0 + \beta_1 X + U$$

假定  $U \sim N(0, \sigma^2)$ , 那么  $Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$

$$Y_0 = \beta_0 + \beta_1 X_0 = E[Y|X = X_0].$$

$\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_0$  是  $Y_0 = \beta_0 + \beta_1 X_0$  的最优线性无偏估计.

# 回归分析的应用：均值预测

$$Y = \beta_0 + \beta_1 X + U$$

假定  $U \sim N(0, \sigma^2)$ , 那么  $Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$

$$Y_0 = \beta_0 + \beta_1 X_0 = E[Y|X = X_0].$$

$\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_0$  是  $Y_0 = \beta_0 + \beta_1 X_0$  的最优线性无偏估计.  
可以证明, 统计量

$$\hat{Y}_0 \sim N\left(\beta_0 + \beta_1 x_0, \sigma^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum x_i^2} \right]\right)$$

# 回归分析的应用：均值预测

$$Y = \beta_0 + \beta_1 X + U$$

假定  $U \sim N(0, \sigma^2)$ , 那么  $Y \sim N(\beta_0 + \beta_1 x, \sigma^2)$

$$Y_0 = \beta_0 + \beta_1 X_0 = E[Y|X = X_0].$$

$\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 X_0$  是  $Y_0 = \beta_0 + \beta_1 X_0$  的最优线性无偏估计.  
可以证明, 统计量

$$\hat{Y}_0 \sim N\left(\beta_0 + \beta_1 x_0, \sigma^2 \left[ \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum x_i^2} \right]\right)$$

我们有

$$t = \frac{\hat{Y}_0 - (\beta_0 + \beta_1 x_0)}{se(\hat{Y}_0)} \sim t(n-2)$$

由此可以得到  $Y_0$  的置信区间.

# 回归结果的评价和分析

用线性模型拟合数据,结果是好,还是不好?



# 回归结果的评价和分析

用线性模型拟合数据,结果是好,还是不好?

判断的准则:

- 1 所估计参数的符号是否与预期的一样;

# 回归结果的评价和分析

用线性模型拟合数据,结果是好,还是不好?

判断的准则:

- 1 所估计参数的符号是否与预期的一样;
- 2 统计检验的显著性是否与他人研究的结果一致;

# 回归结果的评价和分析

用线性模型拟合数据,结果是好,还是不好?

判断的准则:

- 1 所估计参数的符号是否与预期的一样;
- 2 统计检验的显著性是否与他人研究的结果一致;
- 3 从判定系数来看,对变异的解释是否足够明显.

# 回归结果的评价和分析

用线性模型拟合数据,结果是好,还是不好?

判断的准则:

- ① 所估计参数的符号是否与预期的一样;
- ② 统计检验的显著性是否与他人研究的结果一致;
- ③ 从判定系数来看,对变异的解释是否足够明显.

除了上述判断的依据外,我们还需要检验残差序列是否满足正态性假定.

# 残差序列正态性检验

检验的方法有：

- ① 残差直方图分析；
- ② 残差正态概率图分析；
- ③ 雅克-贝拉检验(Jarque-Bera Test)( $JB$ )；
- ④ 安德森-达林检验(Anderson-Darling Normality Test)( $A^2$ ).

# 残差序列正态性检验

检验的方法有:

- ① 残差直方图分析;
- ② 残差正态概率图分析;
- ③ 雅克-贝拉检验(Jarque-Bera Test)( $JB$ );
- ④ 安德森-达林检验(Anderson-Darling Normality Test)( $A^2$ ).

$JB$ 检验用的较多, 统计量为:

$$JB = n \left[ \frac{S^2}{6} + \frac{(K - 3)^2}{24} \right]$$

S: 偏态参数;

K: 峰态参数.

# 偏态和峰态的定义

偏态的定义:

$$S = \frac{E[(X - \mu)^3]}{\sigma^3}$$

# 偏态和峰态的定义

偏态的定义:

$$S = \frac{E[(X - \mu)^3]}{\sigma^3}$$

峰态的定义:

$$K = \frac{E[(X - \mu)^4]}{[E(X - \mu)^2]^2}$$



# 偏态和峰态的定义

偏态的定义:

$$S = \frac{E[(X - \mu)^3]}{\sigma^3}$$

峰态的定义:

$$K = \frac{E[(X - \mu)^4]}{[E(X - \mu)^2]^2}$$

- ①  $K = 3$ , 常态, 例如正态分布.
- ②  $K < 3$ , 扁峰, 尾部肥而短.
- ③  $K > 3$ , 尖峰, 尾部细而长.

# 偏态和峰态的定义

偏态的定义:

$$S = \frac{E[(X - \mu)^3]}{\sigma^3}$$

峰态的定义:

$$K = \frac{E[(X - \mu)^4]}{[E(X - \mu)^2]^2}$$

- ①  $K = 3$ , 常态, 例如正态分布.
- ②  $K < 3$ , 扁峰, 尾部肥而短.
- ③  $K > 3$ , 尖峰, 尾部细而长.

正态分布为常态.

$$H_0 : u_i \sim N(0, \sigma^2)$$

$$JB \xrightarrow{P} \chi^2(2)$$