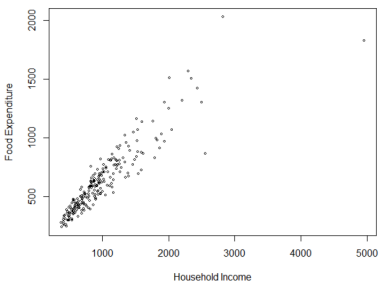




### 7 异常值与强影响点







R

```
1 > rm(list=ls())
2 > library(quantreg)
3 > data(engel)
4 > head(engel)
5     income  foodexp
6 1 420.1577 255.8394
7 2 541.4117 310.9587
8 > attach(engel)
9 > plot(income, foodexp, xlab="Household Income",
10        ylab="Food Expenditure", cex=.5)
```



R

```
1 rm(list=ls())
2 > library(MASS)
3 > p <- 3; n <- 50
4 > mu <- rep(0, p)
5 > sigmu <- diag(rep(1, p))
6 > x <- MASS::mvrnorm(n=n, mu=mu, Sigma = sigmu)
7 > sigma <- rep(1:10, length=n)
8 > e <- rnorm(n=n, mean=0, sd=sigma)
9 > beta0 <- c(1.5, 2, 1.5)
10 > y <- x %*% beta0 + x[, 1]* e
11 > plot(x[, 1], y)
12 > fit <- lm(y~-1)
13 > summary(fit)
14 Coefficients:
15 Estimate Std. Error t value Pr(>|t|)
16 x1 -0.2041      0.8627  -0.237   0.8140
17 x2  1.8931      1.0092   1.876   0.0669
18 x3  1.0701      0.9453   1.132   0.2634
19 Residual standard error: 6.405 on 47 degrees of freedom
20 Multiple R-squared: 0.08127, Adjusted R-squared: 0.02262
21 F-statistic: 1.386 on 3 and 47 DF, p-value: 0.2587
```

# 异方差性的检验

## 残差图分析法

图: 无异方差

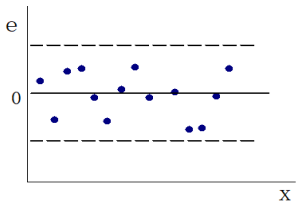
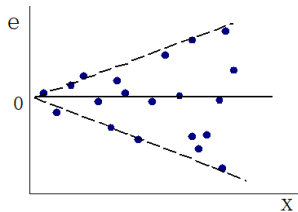


图: 存在异方差









# 例 4.3R 代码

```

1 > rm(list=ls())
2 > ex43 <- read.table("ex43.txt", head=TRUE, fileEncoding="utf8")
3 > attach(ex43)
4 > head(ex43)
5     y      x
6 1 5081 25669
7 2 2724 17885
8 > fit43 <- lm(y~x)
9 > summary(fit43)
10
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) 1.274e+02  2.744e+02  0.464      0.646
13 x           1.068e-01  8.573e-03 12.454 3.66e-13 ***
14
15 > cor.test(x=x, y =abs(fit43$residuals), method = "spearman")
16 S = 2100, p-value = 0.0008485
17 alternative hypothesis: true rho is not equal to 0
18 sample estimates:
19 rho
20 0.5766129
    
```

R

```

1 rm(list=ls())
2 library(MASS)
3 p <- 3; n <- 50
4 mu <- rep(0, p)
5 sigmu <- diag(rep(1, p))
6 x <- MASS::mvrnorm(n=n, mu=mu, Sigma = sigmu)
7 sigma <- rep(1:10, length=n)
8 e <- rnorm(n=n, mean=0, sd=sigma)
9 beta0 <- c(1.5, 2, 1.5)
10 y <- x %*% beta0 + x[, 1]* e
11 plot(x[, 1], y)
12 fit <- lm(y~1)
13 summary(fit)
14 plot(fit$residuals)
15 cor.test(x=x[, 1], y =abs(fit$residuals), method = "spearman")
16 cor.test(x=x[, 2], y =abs(fit$residuals), method = "spearman")
17 cor.test(x=x[, 3], y =abs(fit$residuals), method = "spearman")
    
```

- 加权最小二乘法, Box-Cox 变换法,(参考文献 [1])
- 方差稳定性变换法

一元线性回归普通最小二乘法的残差平方和为:

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - E(y_i))^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

一元线性回归的加权最小二乘的离差平方和为:

$$Q(\beta_0, \beta_1) = \sum_{i=1}^n \omega_i (y_i - E(y_i))^2 = \sum_{i=1}^n \omega_i (y_i - \beta_0 - \beta_1 x_i)^2$$

加权最小二乘估计为

$$\begin{cases} \hat{\beta}_{0\omega} = \bar{y}_\omega - \hat{\beta}_{1\omega} \bar{x}_\omega \\ \hat{\beta}_{1\omega} = \frac{\sum_{i=1}^n \omega_i (x_i - \bar{x}_\omega)(y_i - \bar{y}_\omega)}{\sum_{i=1}^n \omega_i (x_i - \bar{x}_\omega)^2} \end{cases}$$

其中  $\bar{x}_\omega = \frac{1}{\sum \omega_i} \omega_i x_i$  为自变量的加权平均,  $\bar{y}_\omega = \frac{1}{\sum \omega_i} \omega_i y_i$  为自变量的加权平均。

## 一元加权最小二乘估计

为了消除异方差的影响, 观测值的权数应该是观测值误差项方差的倒数, 即

$$\omega_i = \frac{1}{\sigma_i^2}$$

$\sigma_i^2$  为第  $i$  个观测值误差项方差。误差项方差较大的观测值接受较小的权数；误差项方差较小的观测值接受较大的方差。

在社会经济研究中，经常会遇到误差项方差与  $x$  的幂函数  $x^m$  成比例，其中， $m$  为待定未知参数

$$\omega_i = \frac{1}{x_i^m}$$





当误差项  $\epsilon_i$  存在异方差

$$\partial = \sum_{i=1}^n \omega_i (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_k x_{ik})^2$$
$$\mathbf{W} = \begin{pmatrix} \omega_1 & & & \vdots \\ & \omega_2 & & \\ & & \ddots & \\ \vdots & & & \omega_n \end{pmatrix}$$
$$\hat{\beta}_w = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-1}\mathbf{X}'\mathbf{W}\mathbf{y}$$

# 多元加权最小二乘法：权函数的确定方法

通常取权函数  $W$  为某个自变量  $x_j, j = 1, 2, \dots, p$  的幂函数，即  $W = x_j^m$

★ 在  $x_1, x_2, \dots, x_p$  这  $p$  个自变量中取哪一个？

这只需计算每个自变量  $x_j$  与普通残差的等级相关系数，选取等级相关系数最大的自变量构造权函数。例 4.4



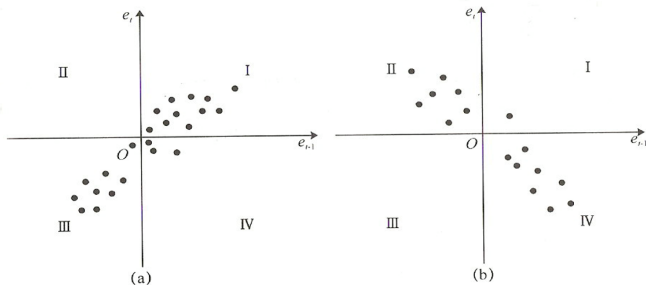
# 自相关性带来的问题

- ① 参数的估计值不再具有最小方差线性无偏性。
- ② 均方误差 MSE 可能严重低估误差项的方差。
- ③ 容易导致对 t 值评价过高, 常用的 F 检验和 t 检验失效。如果忽视这一点, 可能导致得出回归参数统计检验为显著, 但实际上并不显著的严重错误结论。
- ④ 当存在序列相关时, 仍然是 的无偏估计量, 但在任一特定的样本中, 可能严重歪曲 的真实情况, 即最小二乘估计量对抽样波动变得非常敏感
- ⑤ 如果不加处理地运用普通最小二乘法估计模型参数, 用此模型进行预测和结构分析将会带来较大的方差甚至错误的解释。

# 自相关性的诊断

## 图示检验法

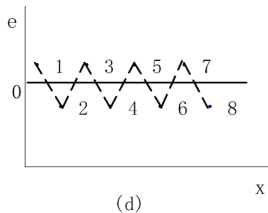
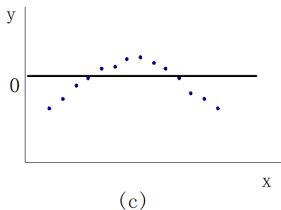
① 绘制  $(e_t, e_{t-1})$  的散点图。



# 自相关性的诊断

## 图示检验法

### ② 按照时间顺序绘制回归残差项 $e_t$ 的图形





# 自相关性的诊断

## ● DW 检验

DW 检验是 J.Durbin 和 G.S.Watson 于 1951 年提出的一种适用于小样本的一种检验方法。

DW 检验只能用于检验随机扰动项具有一阶自回归形式的序列相关问题。

这种检验方法是建立计量经济学模型中最常用的方法，一般的计算机软件都可自动产生出 D.W 值。

随机误差项的一阶自回归形式为

$$\epsilon_t = \rho\epsilon_{t-1} + \mu_t$$

为了检验序列的相关性，构造的假设是

$$H_0: \rho = 0$$



# 自相关性的诊断

## ● DW 检验

定义 DW 统计量为：

$$DW = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=2}^n e_t^2}$$

如果认为  $\sum_{t=2}^n e_t^2$  与  $\sum_{t=2}^n e_{t-1}^2$  近似相等

$$DW = \frac{\sum_{t=2}^n e_t^2 + \sum_{t=2}^n e_{t-1}^2 - 2 \sum_{t=2}^n e_t e_{t-1}}{\sum_{t=2}^n e_t^2} \approx 2 \left[ 1 - \frac{\sum_{t=2}^n e_t e_{t-1}}{\sum_{t=2}^n e_t^2} \right]$$

$$\hat{\rho} = \frac{\sum_{t=2}^n e_t e_{t-1}}{\sqrt{\sum_{t=2}^n e_t^2} \sqrt{\sum_{t=2}^n e_{t-1}^2}} \approx \frac{\sum_{t=2}^n e_t e_{t-1}}{\sum_{t=2}^n e_t^2}$$

因此

$$DW \approx 2(1 - \hat{\rho})$$

# DW 检验

图: DW 值与  $\hat{\rho}$  对应关系

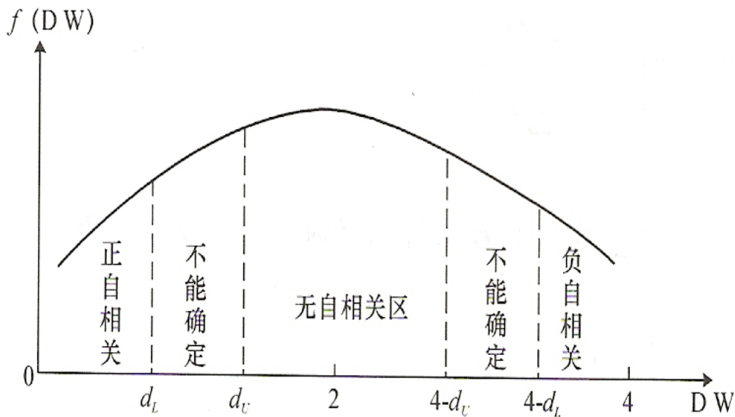
$\hat{\rho}$	D. W	误差项的自相关性
-1	4	完全负自相关
$(-1, 0)$	$(2, 4)$	负自相关
0	2	无自相关
$(0, 1)$	$(0, 2)$	正自相关
1	0	完全正自相关

# DW 检验

根据样本容量  $n$  和解释变量的数目  $k$  (这里包括常数项), 查 DW 分布表, 得临界值  $d_L$  和  $d_U$ , 然后依下列准则考察计算得到的 DW 值, 以决定模型的自相关状态:

$0 \leq D.W. \leq d_L$ ,	误差项 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ 间存在正相关;
$d_L < D.W. \leq d_U$ ,	不能判定是否有自相关;
$d_U < D.W. < 4 - d_U$ ,	误差项 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ 间无自相关;
$4 - d_U \leq D.W. < 4 - d_L$ ,	不能判定是否有自相关;
$4 - d_L \leq D.W. \leq 4$ ,	误差项 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ 间存在负相关。

# DW 检验



# DW 检验

D.W 检验尽管有着广泛的应用,但也有明显的缺点和局限性。

- DW 检验有一个不能确定的区域,一旦 DW 值落在这个区域,就无法判断。这时,只有增大样本容量或选取其他方法。
- DW 统计量的上、下界表要求  $n > 15$ , 这是因为样本如果再小,利用残差就很难对自相关的存在性作出比较正确的诊断。
- DW 检验不适应随机项具有高阶序列相关的检验。

# 自相关问题的处理方法

## • 迭代法

以一元线性回归模型为例，设一元线性回归模型的误差项存在一阶自相关

$$y_t = \beta_0 + \beta_1 x_t + \epsilon_t$$

$$\epsilon_t = \rho \epsilon_{t-1} + \mu_t$$

$$\begin{cases} E(\mu_t) = 0, & t = 1, 2, \dots, n \\ \text{cov}(\mu_t, \mu_s) = \begin{cases} \sigma^2, & t = s \\ 0, & t \neq s \end{cases} & t, s = 1, 2, \dots, n \end{cases}$$

根据一元线性回归模型  $y_t = \beta_0 + \beta_1 x_t + \epsilon_t$ ，有

$$y_{t-1} = \beta_0 + \beta_1 x_{t-1} + \epsilon_{t-1}$$

# 迭代法

变形后有

$$(y_t - \rho y_{t-1}) = (\beta_0 - \rho \beta_0) + \beta_1(x_t - \rho x_{t-1}) + (\epsilon_t - \rho \epsilon_{t-1})$$

令：

$$y'_t = y_t - \rho y_{t-1}$$

$$x'_t = x_t - \rho x_{t-1}$$

$$\beta'_0 = \beta_0(1 - \rho)$$

$$\beta'_1 = \beta_1$$

得到有随机独立误差项，满足线性回归基本假设的

$$y'_t = \beta'_0 + \beta'_1 x'_t + \mu_t \quad (*)$$

## 32 / 46



# 自相关问题的处理方法

## ● 差分法

一阶差分法通常适用于原模型存在较高程度的一阶自相关的情况。在迭代法中，当  $\rho$

$$(y_t - \rho y_{t-1}) = (\beta_0 - \rho \beta_0) + \beta_1(x_t - \rho x_{t-1}) + (\epsilon_t - \epsilon_{t-1})$$

为：

$$(y_t - \rho y_{t-1}) = \beta_1(x_t - x_{t-1}) + (\epsilon_t - \epsilon_{t-1})$$

以  $\Delta y_t = y_t - y_{t-1}$ ， $\Delta x_t = x_t - x_{t-1}$ ，得到

$$\Delta y_t = \beta_1 \Delta x_t + \mu_t$$

上式是不带有常数项的回归方程

$$\hat{\beta}_1 = \frac{\sum_{t=2}^n \Delta y_t \Delta x_t}{\sum_{t=2}^n \Delta x_t^2}$$

# 差分法

一阶差分法的应用条件是自相关系数  $\rho = 1$ ，在实际应用中， $\rho$  接近 1 时我们就采用差分法而不用迭代法，这有两个原因。

第一，迭代法需要用样本估计自相关系数  $\rho$ ，对  $\rho$  的估计误差会影响迭代法的使用效率；

第二，差分法比迭代法简单，人们在建立时序数据的回归模型时，更习惯于用差分法。

但是完全的  $\rho = 1$  情况几乎是见不到的，实际应用时  $\rho$  较大就行！

## BOX-COX 变换

BOX-COX 变换是由博克斯 (BOX) 与考克斯 (COX) 在 1964 年提出的一种应用非常广泛的变换, 它是对因变量  $y$  做如下变换:

$$y^{(\lambda)} = \begin{cases} \frac{y^\lambda - 1}{\lambda}, \lambda \neq 0 \\ \ln y, \lambda = 0 \end{cases}$$

其中,  $\lambda$  为待定参数。此变换要求  $y$  的各分量都大于 0。否则可用下面推广的 BOX-COX 变换

$$y^{(\lambda)} = \begin{cases} \frac{(y + a)^\lambda - 1}{\lambda}, \lambda \neq 0 \\ \ln(y + a), \lambda = 0 \end{cases}$$

即先对  $y$  做平移, 使得  $y + a$  的各个分量都大于 0 后再做 BOX-COX 变换。

对于不同的  $\lambda$ , 所做的变换也不同, 所以这是一个变换族。它包含一些常用的变换, 如对数变换 ( $\lambda = 0$ ), 平方根变换 ( $\lambda = 1/2$ ) 和倒数变换 ( $\lambda = -1$ )。

# BOX-COX 变换

寻找合适的  $\lambda$ ，使得变换后

$$\mathbf{y}^{(\lambda)} = \begin{pmatrix} y_1^{(\lambda)} \\ y_2^{(\lambda)} \\ \vdots \\ y_n^{(\lambda)} \end{pmatrix} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

从而符合线性回归模型的各项假设：误差分量等方差、不相关等。

事实上，BOX-COX 变换不仅可以处理异方差性问题，还能处理自相关、误差非正态、回归函数非线性等情况。

# BOX-COX 变换

经过计算可得  $\lambda$  的最大似然估计 (参见参考文献 [2])

$$L_{\max}(\lambda) = (2\pi e \hat{\sigma}_{\lambda}^2)^{-\frac{n}{2}} |\mathbf{J}|$$

式中,  $\hat{\sigma}_{\lambda}^2 = \frac{1}{n} SSE(\lambda, y^{(\lambda)})$ ,  $|\mathbf{J}| = \prod_{i=1}^n \left| \frac{dy_i^{(\lambda)}}{dy_i} \right| = \prod_{i=1}^n y_i^{(\lambda-1)}$

令  $z^{(\lambda)} = \frac{y^{(\lambda)}}{|\mathbf{J}|}$ , 对  $L_{\max}$  取对数并略去与  $\lambda$  无关的常数项, 可得

$$\ln L_{\max}(\lambda) = -\frac{n}{2} \ln SSE(\lambda, z^{(\lambda)})$$

为找出  $\lambda$ , 使得  $\ln L_{\max}(\lambda)$  达到最大, 只需使  $SSE(\lambda, z^{(\lambda)})$  达到最小即可。它的解析解比较难找, 通常是给出一系列  $\lambda$  的值, 计算对应的  $SSE(\lambda, z^{(\lambda)})$ , 取使得  $SSE(\lambda, z^{(\lambda)})$  达到最小的  $\lambda$  即可。

## BOX-COX 变换

- BOX-COX 变换是一个幂变换族，其中当变换参数  $\lambda = 0$  时成为对数变换，而对数变换则是比幂变换应用更广泛的变换，很多场合都可以首先尝试对数据作对数变换。
- 从概率分布的角度看，当数据本身服从对数正态分布时，对数据取对数变换后就服从正态分布。对数正态分布是右偏分布，有厚重的右尾。
- 从数据看，如果数据中一些数值很大，但是小数值的数据更密集，个数也更多，大数值的数据较较疏松，个数较少，这样的数据很可能服从对数正态分布，可以尝试对数变换。
- 对回归分析问题，如果只对因变量作对数变换，就是 BOX-COX 变换  $\lambda = 0$  时的特例。也可以考虑只对自变量作对数变换，或者同时对因变量和对自变量作对数变换。



## 关于因变量 $y$ 的异常值

在残差分析中,认为超过  $\pm 3\hat{\sigma}$  的残差为异常值。标准化残差

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

## 学生化残差

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

$h_{jj}$  为  $H = X(X'X)^{-1}X'$  的主对角线元素。

当观测数据中存在关于  $y$  的异常观测值时，普通残差、标准化残差、学生化残差这三种残差都不再适用。这是由于异常值把回归线拉向自身，使异常值本身的残差减小，而其余观测值的残差增大，这时回归标准差  $\hat{\sigma}$  也会增大，因而用传统的“ $3\sigma$ ”准则不能正确分辨出异常值。解决这个问题的方法是改用删除残差。



## 关于因变量 $y$ 的异常值

删除残差的构造思想是：在计算第  $i$  个观测值的残差时，用删除掉的第  $i$  个观测值的其余  $n - 1$  个观测值拟合回归方程，计算出第  $i$  个观测值的删除拟合值  $\hat{y}_{(i)}$ ，这个删除拟合值与第  $i$  个值无关，不受第  $i$  个值是否为异常值的影响，由此定义第  $i$  个观测值的删除残差为

$$e_{(i)} = y_i - \hat{y}_{(i)}$$

可以证明

$$e_{(i)} = \frac{e_i}{1 - h_{ii}}$$

进一步，可以给出第  $i$  个观测值的删除学生残差，记为  $ser_{(i)}$ 。

$$SRE_{(i)} = SRE_i \left( \frac{n - p - 2}{n - p - 1 - SRE_i^2} \right)^{\frac{1}{2}}$$

$|SRE_{(i)}| > 3$  的观测值即判定为异常值。

# 关于自变量 $x$ 的异常值对回归的影响 (高杠杆点 High-leverage point)

- 在  $D(e_i) = (1 - h_{ii})\sigma^2$  中,  $h_{ii}$  为帽子矩阵中主对角线的第  $i$  个元素, 它是调节  $e_i$  方差大小的杠杆, 因而称  $h_{ii}$  为第  $i$  个观测值的杠杆值。类似于一元线性回归, 多元线性回归的杠杆值  $h_{ii}$  也表示自变量的第  $i$  次观测与自变量平均值之间距离的远近。
- 较大的杠杆值的残差偏小, 这是因为 **杠杆值大的观测点远离样本中心**, 能够把回归拉向自身, 因而把杠杆值大的样本点称为 **强影响点**。
- $tr(\mathbf{H}) = \sum_{i=1}^n h_{ii} = p + 1$ , 则杠杆值的平均值为

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h_{ii} = \frac{p + 1}{n}$$

一个杠杆值  $h_{ii}$  大于 2 倍或者 3 倍的  $\bar{h}$ , 就认为是大的。





## 思考题

Consider a random sample  $X_1, X_2, \dots, X_n \sim Unif(0, \theta)$ .

- i Find the estimator for  $\theta$  through MoM, denoted by  $\hat{\theta}_{MM}$ .
- ii Find the MLE  $\hat{\theta}_{MLM}$ .
- iii What are the expectation and variance of  $\hat{\theta}_{MM}$  and  $\hat{\theta}_{MLM}$ ? Which estimator is better?

## 作业

- p.124 4.4
- p.124 4.9
- p.124 4.14