

回归分析 (01714601)

课程主页: <http://staff.ustc.edu.cn/~ynyang/lm2020>

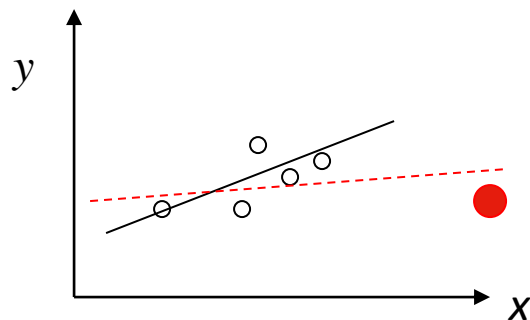
# 第26讲 影响分析（续）

2020.5.22

Jackknife



## 2. 高杠杆点: $\mathbf{x}$ 异常, $h$ 过大



所有自变量为 $\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n$ ,  $Z = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$ , 设计阵 $X = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} = (\mathbf{1}, Z) = \begin{pmatrix} 1 & \tilde{\mathbf{x}}_1^\top \\ \vdots & \vdots \\ 1 & \tilde{\mathbf{x}}_n^\top \end{pmatrix}$ ,

如果第 $i$ 个自变量 $\tilde{\mathbf{x}}_i$ 与中心 $\bar{\mathbf{x}}/n$ 的马氏距离 $d_s(\tilde{\mathbf{x}}_i, \bar{\mathbf{x}})$ 较大,则认为自变量 $\tilde{\mathbf{x}}_i$ 异常。

$$d_s(\tilde{\mathbf{x}}_i, \bar{\mathbf{x}}) = \sqrt{(\tilde{\mathbf{x}}_i - \bar{\mathbf{x}})^\top S^{-1} (\tilde{\mathbf{x}}_i - \bar{\mathbf{x}})}$$

自变量的样本均值:  $\bar{\mathbf{x}} = (\tilde{\mathbf{x}}_1 + \dots + \tilde{\mathbf{x}}_n)/n = Z^\top \mathbf{1}/n$ ,

样本方差:  $S = \sum_{k=1}^n (\tilde{\mathbf{x}}_k - \bar{\mathbf{x}})(\tilde{\mathbf{x}}_k - \bar{\mathbf{x}})^\top / (n-1) = Z^\perp{}^\top Z^\perp / (n-1)$

帽子矩阵/投影矩阵:  $H = X(X^\top X)^{-1}X^\top = (h_{ij})$ ,  $h_{ij} = \mathbf{x}_i^\top (X^\top X)^{-1} \mathbf{x}_j = (1, \tilde{\mathbf{x}}_i^\top) \begin{pmatrix} n & \mathbf{1}^\top Z \\ Z^\top \mathbf{1} & Z^\top Z \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ \tilde{\mathbf{x}}_j \end{pmatrix}$

命题1.  $H$ 的第 $i$ 个对角元  $h_{ii} = \frac{1}{n} + d_S^2(\tilde{\mathbf{x}}_i, \bar{\mathbf{x}})/(n-1)$ , 且  $\frac{1}{n} \leq h_{ii} \leq 1$ . (事实上,  $h_{ii} < 1$ ).

证明:  $X = (\mathbf{1}, Z)$ ,  $Z$ 与 $\mathbf{1}$ 正交化得 $Z$ 的中心化矩阵:  $Z^\perp = Z - \mathbf{1}\bar{\mathbf{x}}^\top = \begin{pmatrix} (\tilde{\mathbf{x}}_1 - \bar{\mathbf{x}})^\top \\ \vdots \\ (\tilde{\mathbf{x}}_n - \bar{\mathbf{x}})^\top \end{pmatrix}$ .

所以  $H = P_X = P_1 + P_{Z^\perp} = \frac{\mathbf{1}\mathbf{1}^\top}{n} + Z^\perp (Z^{\perp\top} Z^\perp)^{-1} Z^{\perp\top}$ , 而  $Z^{\perp\top} Z^\perp = (n-1)S$

所以 $H$ 的第 $i$ 个对角元素:  $h_{ii} = \frac{1}{n} + (\tilde{\mathbf{x}}_i - \bar{\mathbf{x}})^\top (Z^{\perp\top} Z^\perp)^{-1} (\tilde{\mathbf{x}}_i - \bar{\mathbf{x}})$

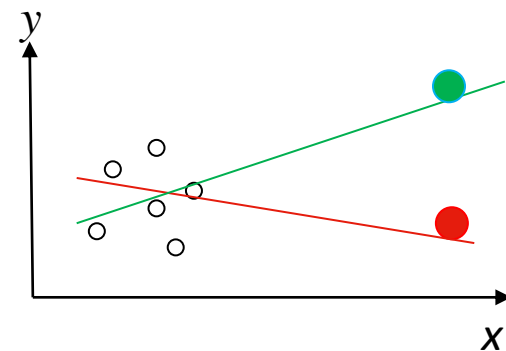
$= \frac{1}{n} + (\tilde{\mathbf{x}}_i - \bar{\mathbf{x}})^\top ((n-1)S)^{-1} (\tilde{\mathbf{x}}_i - \bar{\mathbf{x}}) = \frac{1}{n} + d_S^2(\tilde{\mathbf{x}}_i, \bar{\mathbf{x}})/(n-1)$ .

显然  $h_{ii} \geq \frac{1}{n}$ , 等号成立当且仅当  $\tilde{\mathbf{x}}_i = \bar{\mathbf{x}}$ . 由  $I_n - H \geq 0 \Rightarrow h_{ii} \leq 1$  .

定义杠杆(leverage):  $H = X(X^T X)^{-1} X^T = (h_{ij})_{1 \leq i, j \leq n}$  的第 $i$ 个对角元素 $h_{ii}$ 称为杠杆(leverage), 它度量了第 $i$ 个自变量远离中心的异常程度。

因为  $\sum_{i=1}^n h_{ii} = \text{tr}(H) = p$ , 平均来看  $h_{ii} \sim \frac{p}{n}$ , 严重高于该值即认为 $\tilde{\mathbf{x}}_i$ 是高影响的,

- 若  $h_{ii} \approx 1$ ,  $\tilde{\mathbf{x}}_i$  远离  $\bar{\mathbf{x}}$ ,  $\text{var}(e_i) = (1 - h_{ii})\sigma^2 \approx 0 \Rightarrow e_i \approx 0$ ,  $\hat{y}_i \approx y_i$ ,  $(\tilde{\mathbf{x}}_i, y_i)$  靠近回归直线而且主导回归直线的方向(右图)



- 若  $h_{ii} \approx 1/n$ , 则  $\tilde{\mathbf{x}}_i \approx \bar{\mathbf{x}}$ , 影响最小, 此时

$$\hat{y}_i = \hat{\beta}_0 + \tilde{\mathbf{x}}_i^T \hat{\boldsymbol{\gamma}} = (\bar{y} - \bar{\mathbf{x}}^T \hat{\boldsymbol{\gamma}}) + \tilde{\mathbf{x}}_i^T \hat{\boldsymbol{\gamma}} \approx \bar{y},$$

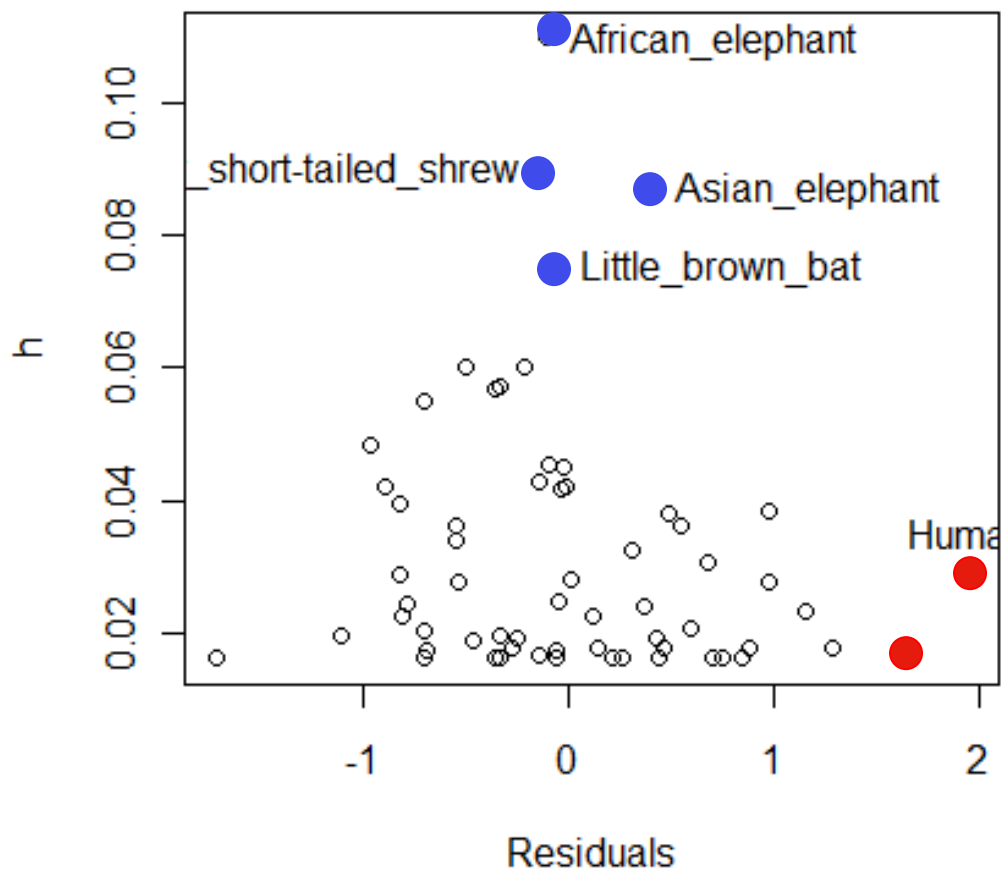
事实上, 若  $h_{ii} = 1/n$ ,  $(\tilde{\mathbf{x}}_i, \hat{y}_i) = (\bar{\mathbf{x}}, \bar{y})$  对回归结果没有任何影响。

```
> hatvalues(lm.out) #杠杆值
```

```
> influence.measures(lm.out) #杠杆值及其它影响度量
```

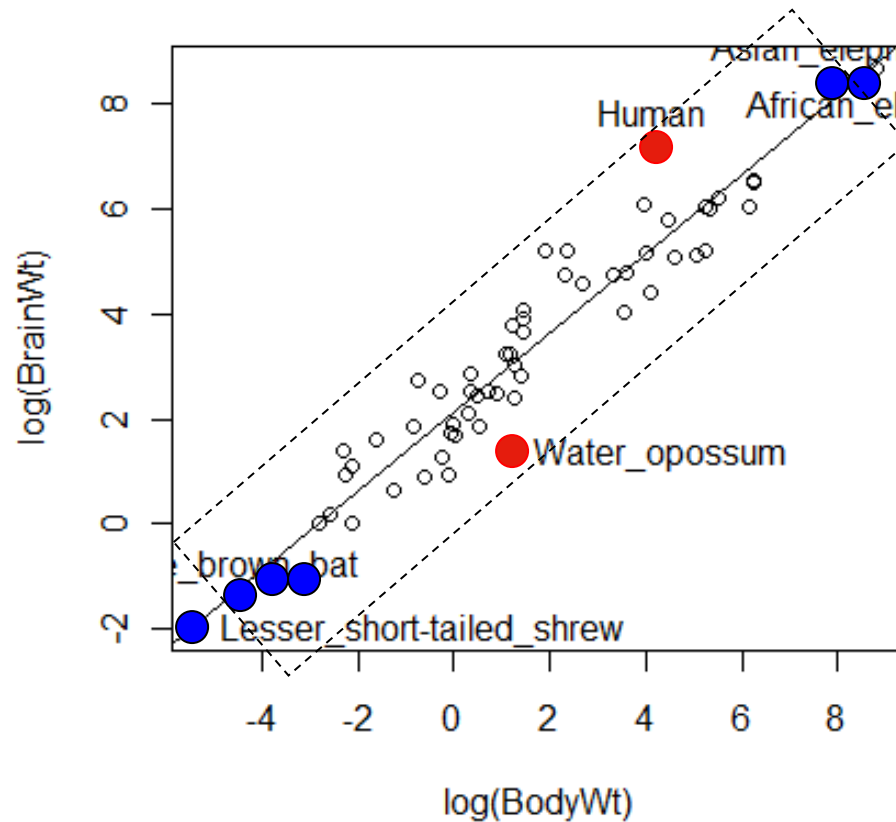
## 例1(续): 动物脑重量与体重的关系

```
a=lm(log(BrainWt)~log(BodyWt), data=brains)
```



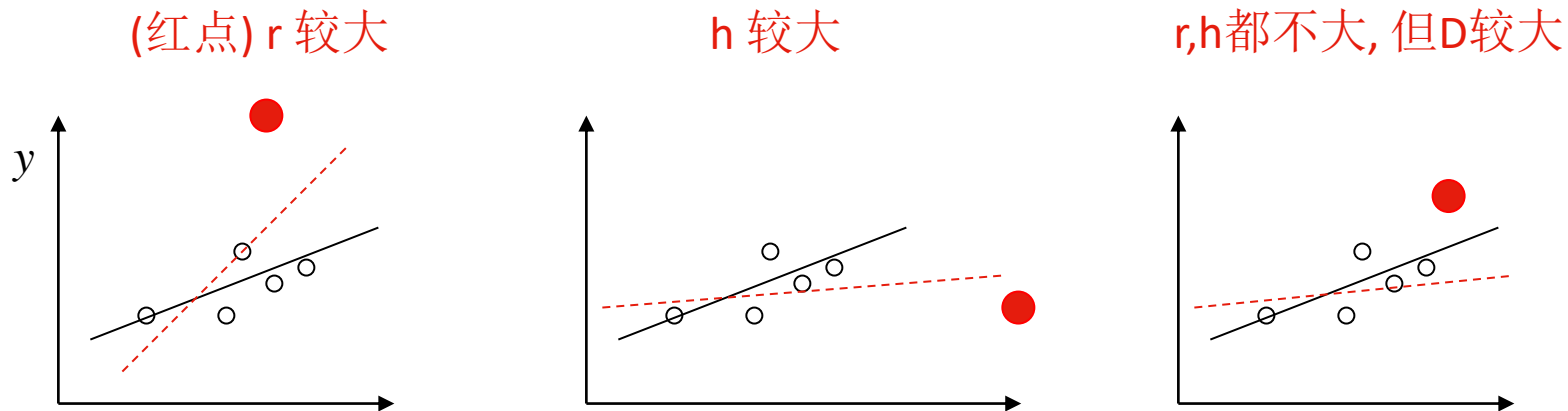
- x-异常 (高杠杆)
- y-异常 (outlier)

异常点outlier(红点Human, Water\_opossum) 远离回归线;  
高杠杆点 (蓝色) 在两端。



对回归直线走向起决定作用的是外围 (虚线) 附近的那些点 (高影响点)

### 3. 高影响点: $(x,y)$ 异常



有时, 第 $i$ 个样本点 $(\mathbf{x}_i, y_i)$ 的自变量和响应都不异常, 但它们合在一起(第三图)可能对回归分析的影响可能较大, 我们以

**Cook距离 $D$ 、DFFITS、DFBETAS**

度量这种影响, 它们都是基于“在原始模型中删除一行数据所导致的拟合的变化”得到的。其中 $D$ 应用最为广泛, 其它影响度量通常作为补充。

# Jackknife方法定义影响度量

John Tukey 用 delete-1 method (leave-one-out) 方法考察数据点的影响, 即删除一个样本点 导致模型拟合效果的变化程度。这可以称为Jackknife方法 (但习惯上Jackknife专指用于求偏差、方差的deletel-1方法)

所有数据: $\mathbf{y}_{n \times 1}, X_{n \times p}$	$\mathbf{y}, X$ 删除第 <i>i</i> 行: $\mathbf{y}_{(-i)}, X_{(-i)}$ ( $n-1$ ) $\times 1$ ( $n-1$ ) $\times p$	差异( <i>DF</i> )
$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$	$\mathbf{y}_{(-i)} = X_{(-i)}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_{(-i)}$	
$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$	$\hat{\boldsymbol{\beta}}^{(-i)} = (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T \mathbf{y}_{(-i)}$	<i>DFBETAS</i> : $\propto \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}^{(-i)}$
$\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$	$\hat{\mathbf{y}}^{(-i)} = X\hat{\boldsymbol{\beta}}^{(-i)}$  假设 $X_{(-i)}^T X_{(-i)}$ 可逆	<i>DFFITS</i> : $\propto (\hat{\mathbf{y}} - \hat{\mathbf{y}}^{(-i)})_i = \mathbf{x}_i^T \hat{\boldsymbol{\beta}} - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(-i)}$  <i>Cook's D</i> : $\propto \ \hat{\mathbf{y}} - \hat{\mathbf{y}}^{(-i)}\ ^2 = \ X\hat{\boldsymbol{\beta}} - X\hat{\boldsymbol{\beta}}^{(-i)}\ ^2$

下标(-*i*)表示删除数据第*i*行:  $\mathbf{y}_{(-i)}, X_{(-i)}$ .

上标(-*i*)表示基于 $\mathbf{y}_{(-i)}, X_{(-i)}$ 得到的估计、拟合等:  $\hat{\boldsymbol{\beta}}^{(-i)}, \hat{\mathbf{y}}^{(-i)}$ .

例如,  $\hat{\boldsymbol{\beta}}^{(-i)} : p \times 1, \mathbf{y}_{(-i)} : (n-1) \times 1, \hat{\mathbf{y}}^{(-i)}_{n \times 1} = X \hat{\boldsymbol{\beta}}^{(-i)} \neq X_{(-i)} \hat{\boldsymbol{\beta}}^{(-i)}$ .



数据:  $(\tilde{\mathbf{x}}_i, y_i), i = 1, \dots, n,$

模型:  $y_i = \beta_0 + \tilde{\mathbf{x}}_i^\top \boldsymbol{\gamma} + \varepsilon_i \hat{=} \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, i = 1, \dots, n$

所有数据  $\mathbf{y}_{n \times 1} = (y_1, \dots, y_n)^\top, X_{n \times p} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top,$

模型:  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon},$

$\boldsymbol{\beta}$ 的LS估计:  $\hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y},$

拟合值:  $\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}},$

$\sigma^2$ 的LS估计:  $\hat{\sigma}^2 = \frac{RSS}{n-p} = \frac{1}{n-p} \|\mathbf{y} - X\hat{\boldsymbol{\beta}}\|^2$

删除第  $i$  个样本点  $(y_i, \tilde{\mathbf{x}}_i)$ ,

数据:  $\mathbf{y}_{(-i)} = (y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n)^\top$ ,  $X_{(-i)} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1}, \mathbf{x}_{i+1}, \dots, \mathbf{x}_n)^\top$ , 行数  $(n-1)$

模型:  $\mathbf{y}_{(-i)} = X_{(-i)}\boldsymbol{\beta} + \boldsymbol{\varepsilon}_{(-i)}$

$\boldsymbol{\beta}$  的LS估计:  $\hat{\boldsymbol{\beta}}^{(-i)} = \left( X_{(-i)}^\top X_{(-i)} \right)^{-1} X_{(-i)}^\top \mathbf{y}_{(-i)}$ ,

$\sigma^2$  的LS估计:  $\hat{\sigma}^{(-i)2} = \frac{1}{n-1-p} \| \mathbf{y}_{(-i)} - X_{(-i)}\hat{\boldsymbol{\beta}}^{(-i)} \|^2$

$\mathbf{y}$  的拟合值:  $\hat{\mathbf{y}}^{(-i)} = X\hat{\boldsymbol{\beta}}^{(-i)}$

注意  $\hat{\mathbf{y}}^{(-i)} = X\hat{\boldsymbol{\beta}}^{(-i)}$  的第  $i$  行  $\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^{(-i)}$  是对  $y_i$  的预测,  
 $X_{(-i)}\hat{\boldsymbol{\beta}}^{(-i)}$  是删除第  $i$  行数据后得到的  $\mathbf{y}_{(-i)}$  的拟合

定义（影响度量）：

$$\text{DFBETAS}_i(k) = \frac{\hat{\beta}_k - \hat{\beta}_k^{(-i)}}{\hat{\sigma}^{(-i)} \sqrt{\left((X^\top X)^{-1}\right)_{kk}}}, \quad 1 \leq k \leq p, 1 \leq i \leq n.$$

$$\text{DFFITS}_i = \frac{(\hat{\mathbf{y}} - \hat{\mathbf{y}}^{(-i)})_i}{\hat{\sigma}^{(-i)} \sqrt{h_{ii}}} = \frac{\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^{(-i)}}{\hat{\sigma}^{(-i)} \sqrt{h_{ii}}}, \quad 1 \leq i \leq n.$$

$$\text{Cook 距离 } D_i = \frac{\|\hat{\mathbf{y}} - \hat{\mathbf{y}}^{(-i)}\|^2}{p \hat{\sigma}^2} = \frac{\|X \hat{\boldsymbol{\beta}} - X \hat{\boldsymbol{\beta}}^{(-i)}\|^2}{p \hat{\sigma}^2}, \quad 1 \leq i \leq n.$$

- (1) 因为  $\text{var}(\hat{\beta}_k) = \sigma^2 \left((X^\top X)^{-1}\right)_{kk}$ ,  $\text{DFBETAS}_i(k)$  是  $\hat{\beta}_k - \hat{\beta}_k^{(-i)}$  的标准化, 分母是  $\hat{\beta}_k$  的标准差,
- (2)  $\text{var}(\mathbf{x}_i^\top \hat{\boldsymbol{\beta}}) = \sigma^2 \mathbf{x}_i^\top (X^\top X)^{-1} \mathbf{x}_i = \sigma^2 h_{ii}$ ,  $\text{DFFITS}_i$  是  $\mathbf{x}_i^\top \hat{\boldsymbol{\beta}} - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^{(-i)}$  的标准化;
- (3)  $D_i$  类似于回归方程显著性的F检验。

需要注意各个度量不服从标准的 $t, F$ 分布，通常如下判别影响点：

若  $|DFBETAS_i(k)| \geq \frac{2}{\sqrt{n}}$ ，第  $i$  个样本点对于估计  $\beta_k$  可视为是高影响的。

若  $|DFFITS_i| \geq 2\sqrt{\frac{p}{n}}$ ，则样本点  $i$  视为是高影响的。

若  $D_i > 0.5$ ，第  $i$  个样本点影响较大；若  $D_i > 1$ ，第  $i$  个样本点影响很大。

R: 影响度量

```
> dfbetas(lm.out)
> dffits(lm.out)
> cooks.distance(lm.out)
> hatvalues(lm.out)
> covratio(lm.out) # 参数估计方差的比  $\det[(X_{(-i)}^T X_{(-i)})^{-1} \hat{\sigma}_{(-i)}^2] / \det[(X^T X)^{-1} \hat{\sigma}^2]$ 
> influence.measures(lm.out) # 所有上面的影响度量值。
```

## 例1(续): 动物脑重量与体重的关系

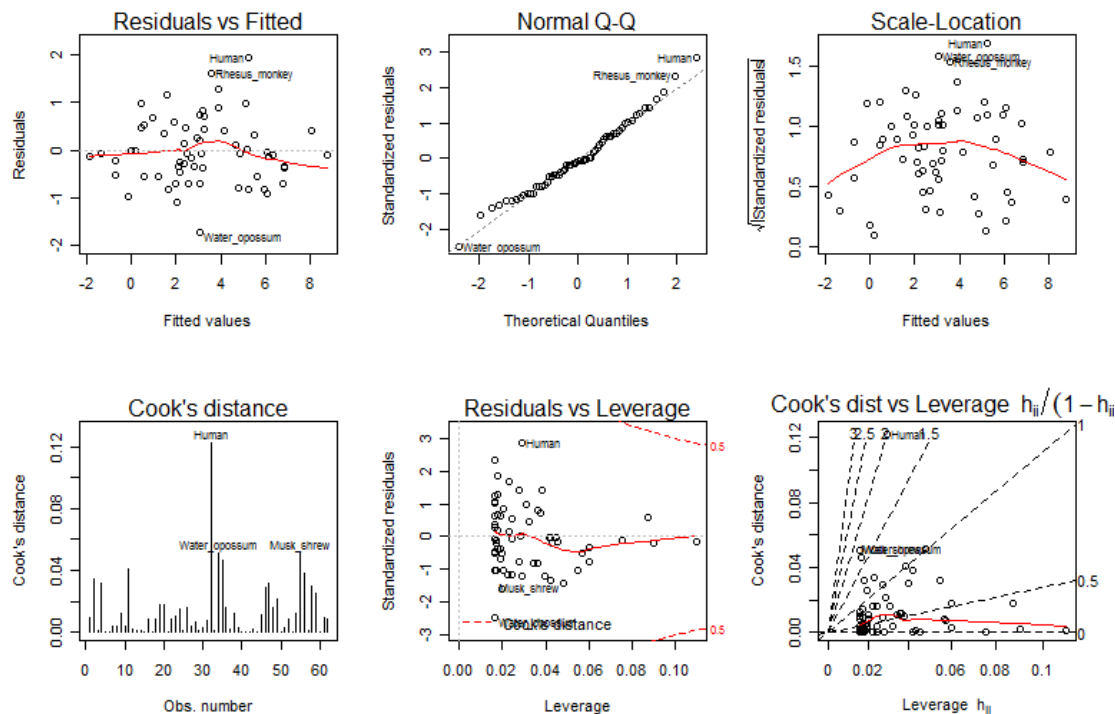
```
a=lm(log(BrainWt)~log(BodyWt), data=brains)
```

```
b=influence.measures(a)$infmtat
```

	DFBETA1	DFBETA2	DFFIT	Cov.ratio	D	hat
	dfb.1_	dfb.1(BW	dffit	cov.r	cook.d	hat
Arctic_fox	0.129	-0.005	0.139	1.011	0.01	0.016
Owl_monkey	0.261	-0.147	0.265	0.961	0.034	0.023
Beaver	-0.052	0.017	-0.052	1.048	0.001	0.018
Cow	-0.041	-0.212	-0.252	1.055	0.032	0.055
Gray_wolf	-0.006	-0.007	-0.012	1.06	0	0.025
Goat	0.014	0.014	0.025	1.057	0	0.023
Roe_deer	0.059	0.035	0.086	1.041	0.004	0.019
Guinea_pig	-0.092	0.036	-0.092	1.039	0.004	0.019
Vervet	0.144	0.005	0.159	0.999	0.012	0.016
Chinchilla	0.081	-0.048	0.083	1.05	0.004	0.024
Ground_squirrel	0.259	-0.219	0.288	1.003	0.041	0.038
Arctic_ground_squir	-0.068	0.028	-0.068	1.047	0.002	0.02
African_giant_rat	-0.05	0.02	-0.05	1.05	0.001	0.019
Human	0.22	0.353	0.527	0.797	0.122	0.029

## # R 回归诊断图（残差分析+影响分析）

```
> myfit = lm( BrainWt ~ BodyWt, data=log(brains))  
> plot(myfit, which=1:6) #default: which=c(1,2,3,5)
```



6个图分别为

- (1) 残差图：线性？等方差？
- (2) qqnorm: 误差正态？
- (3) 刻度-位置图(残差图的补充)：线性？等方差？
- (4) Cook's D: 影响分析
- (5) 残差-杠杆图: 影响分析
- (6)  $D$  vs  $h$ : 影响分析

红色实线为非线性拟合，红色虚线为D-等高线.