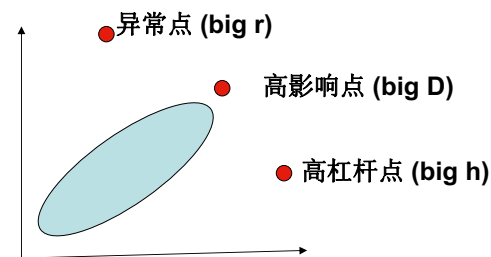


## 第二十二讲. 置信域和同时置信区间

1

回顾影响分析



定义:  $DFITS_i = \frac{\hat{Y}_i - \hat{Y}_{(-i)}}{\sqrt{h_{ii} \hat{\sigma}_{(-i)}^2}}$ , 其中  $\hat{Y}_{(-i)} = \tilde{x}_i' \hat{\beta}_{(-i)}$  = 删除第*i*行后对 $Y_i$ 的预测

证明:  $DFITS_i = r_i^* \sqrt{\frac{h_{ii}}{1-h_{ii}}}$ , 其中  $r_i^* = \frac{e_i}{\hat{\sigma}_{(-i)} \sqrt{1-h_{ii}}}$

定义: Cook 距离  $D_i = \frac{\|\hat{Y} - \hat{Y}_{(-i)}\|^2}{p \hat{\sigma}^2} = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})' X' X (\hat{\beta} - \hat{\beta}_{(-i)})}{p \hat{\sigma}^2}$

证明:  $D_i = \frac{1}{p} \times \frac{h_{ii}}{1-h_{ii}} r_i^2$ , 其中  $r_i = \frac{e_i}{\sqrt{1-h_{ii}} \hat{\sigma}}$  为标准化残差。

2

证明:  $D_i = \frac{1}{p} \times \frac{h_{ii}}{1-h_{ii}} r_i^2$ , 其中  $r_i = \frac{e_i}{\sqrt{1-h_{ii}} \hat{\sigma}}$ 。

证: 删除样本点*i*, 拟合模型  $Y_{(-i)} = X_{(-i)} \beta + \varepsilon_{(-i)}$ , 得  $\hat{\beta}_{(-i)} = (X_{(-i)}' X_{(-i)})^{-1} X_{(-i)}' Y_{(-i)}$

记  $X = \begin{pmatrix} \tilde{x}_1' \\ \vdots \\ \tilde{x}_n' \end{pmatrix}$ , 则  $X'X = \sum_{i=1}^n \tilde{x}_i \tilde{x}_i'$ ,  $X'Y = \sum_{i=1}^n \tilde{x}_i y_i$

$\Rightarrow X'X = X_{(-i)}' X_{(-i)} + \tilde{x}_i \tilde{x}_i'$ ,  $X'Y = \sum_{j=1}^n \tilde{x}_j y_j = X_{(-i)}' Y_{(-i)} + \tilde{x}_i y_i$

利用事实: 设  $A_{p \times p}$  对称,  $x, y$  为  $p \times 1$  向量, 则  $(A + xy')^{-1} = A^{-1} - \frac{A^{-1}xy'A^{-1}}{1 + x'A^{-1}y}$

$\Rightarrow \hat{\beta}_{(-i)} = (X_{(-i)}' X_{(-i)})^{-1} X_{(-i)}' Y_{(-i)} = (X'X - \tilde{x}_i \tilde{x}_i')^{-1} (X'Y - \tilde{x}_i y_i)$

$= \left( (X'X)^{-1} + \frac{(X'X)^{-1} \tilde{x}_i \tilde{x}_i' (X'X)^{-1}}{1 - \tilde{x}_i' (X'X)^{-1} \tilde{x}_i} \right) (X'Y - \tilde{x}_i y_i)$

3

注意到  $h_{ii} = \tilde{x}_i' (X'X)^{-1} \tilde{x}_i$

$\hat{\beta}_{(-i)} = \left( (X'X)^{-1} + \frac{(X'X)^{-1} \tilde{x}_i \tilde{x}_i' (X'X)^{-1}}{1 - \tilde{x}_i' (X'X)^{-1} \tilde{x}_i} \right) (X'Y - \tilde{x}_i y_i)$   
 $= \hat{\beta} - (X'X)^{-1} \tilde{x}_i y_i + \frac{1}{1-h_{ii}} (X'X)^{-1} \tilde{x}_i \tilde{x}_i' (X'X)^{-1} X'Y - \frac{1}{1-h_{ii}} (X'X)^{-1} \tilde{x}_i \tilde{x}_i' (X'X)^{-1} \tilde{x}_i y_i$   
 $= \hat{\beta} - (X'X)^{-1} \tilde{x}_i y_i + \frac{1}{1-h_{ii}} (X'X)^{-1} \tilde{x}_i \tilde{x}_i' \hat{\beta} - \frac{h_{ii}}{1-h_{ii}} (X'X)^{-1} \tilde{x}_i y_i$   
 $= \hat{\beta} - \frac{1}{1-h_{ii}} \left\{ (X'X)^{-1} \tilde{x}_i (y_i - \tilde{x}_i' \hat{\beta}) \right\} = \hat{\beta} - \frac{1}{1-h_{ii}} \left\{ (X'X)^{-1} \tilde{x}_i e_i \right\}$

$\Rightarrow \hat{\beta} - \hat{\beta}_{(-i)} = \frac{1}{1-h_{ii}} \left\{ (X'X)^{-1} \tilde{x}_i e_i \right\}$

所以  $D_i = \frac{(\hat{\beta} - \hat{\beta}_{(-i)})' X' X (\hat{\beta} - \hat{\beta}_{(-i)})}{p \hat{\sigma}^2} = \frac{1}{p} \times \frac{h_{ii}}{1-h_{ii}} r_i^2$

注: 由  $\hat{Y}_i - \hat{Y}_{(-i)} = (X\hat{\beta} - X\hat{\beta}_{(-i)})_i = \frac{1}{1-h_{ii}} \left\{ \tilde{x}_i' (X'X)^{-1} \tilde{x}_i e_i \right\} = \frac{h_{ii} e_i}{1-h_{ii}} \Rightarrow DFITS_i = r_i^* \sqrt{\frac{h_{ii}}{1-h_{ii}}}$

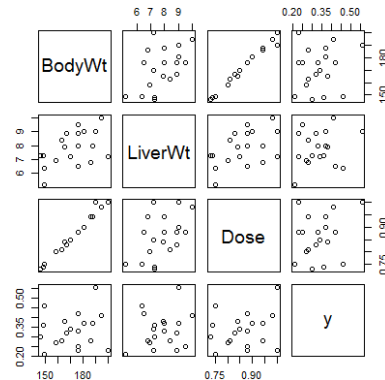
4

例2. (Dennis Cook数据)。一项试验希望研究动物肝脏对某种药物的吸收能力。为此随机选取19只小白鼠，口服该药物，剂量大小(Dose)由体重决定 (大约为 40mg/kg 乘以体重)。一段时间后，测出肝中与所含药物重量，除以服用的剂量，得到肝吸收百分比 y。我们研究 y 与体重 BodyWt, 肝重LiverWt, 剂量Dose的关系。

理论上，这种决定药物剂量的方法决定了所测得到的 y 应该与三个变量无关。

```
lm( y ~ ., data = rat)
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.266    0.195   1.4   0.19
BodyWt      -0.021    0.008  -2.7   0.02 *
LiverWt      0.014    0.017   0.8   0.42
Dose         4.178    1.523   2.7   0.02 *
---
Multiple R-squared:  0.36,
F-statistic: 2.9 on 3 and 15 DF, p-value: 0.072
```

BodyWt和Dose都显著。



5

因为BodyWt, Dose几乎完全成正比，模型中去掉BodyWt 或Dose, 变量不再显著：

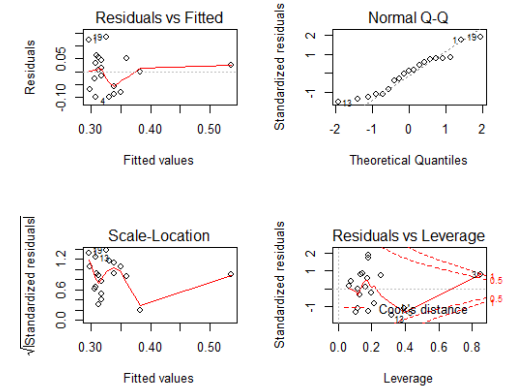
```
lm( y ~ . - Dose, data = rat)
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.17836    0.22778   0.8   0.4
BodyWt       0.00035    0.00151   0.2   0.8
LiverWt       0.01233    0.02041   0.6   0.6
```

```
lm( y ~ . - BodyWt, data = rat)
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.1174    0.2191   0.5   0.6
LiverWt      0.0087    0.0201   0.4   0.7
Dose         0.1736    0.2863   0.6   0.6
```

如何解释这种现象？

从第一个图可以发现第3个小白鼠的拟合值过大，其杠杆值  $h = 0.85$ ；

从第4个图发现其  $D = 0.93$ ，是高影响点。



Influence measures of  
lm(formula = y ~ ., data = rat) :

	dfb.1	dfb.BdyW	dfb.LvrW	dfb.Dose	dfbet	cov.r	cook.d	hat
1	-0.038	0.315	-0.704	-0.244	0.892	0.631	0.169	0.178
2	0.143	-0.098	-0.482	0.126	-0.609	1.016	0.089	0.179
3	-0.231	-1.668	0.305	1.747	1.905	7.401	0.930	0.851
4	0.125	-0.127	-0.304	0.140	-0.494	0.860	0.057	0.108
5	0.522	-0.396	0.550	0.275	-0.909	1.524	0.203	0.392
6	0.002	0.014	0.029	-0.017	0.043	1.567	0.000	0.161
7	-0.184	0.150	-0.083	-0.118	0.310	1.289	0.025	0.137
8	-0.297	0.059	0.246	-0.040	0.426	1.520	0.047	0.254
9	-0.010	0.018	0.000	-0.017	0.043	1.402	0.000	0.067
10	-0.006	0.010	-0.003	-0.009	-0.014	1.496	0.000	0.120
11	-0.291	0.194	0.101	-0.173	-0.410	1.066	0.041	0.120
12	0.217	-0.025	0.052	-0.009	0.269	1.444	0.019	0.172
13	-0.772	0.144	0.766	-0.120	-1.099	0.972	0.273	0.316
14	-0.035	-0.046	-0.077	0.060	-0.142	1.461	0.005	0.131
15	0.019	0.041	-0.055	-0.038	0.119	1.359	0.004	0.076
16	0.123	-0.006	0.329	-0.049	-0.447	1.375	0.051	0.217
17	-0.104	0.015	-0.028	0.002	-0.126	1.607	0.004	0.195
18	-0.154	0.190	0.162	-0.189	0.352	1.270	0.032	0.149
19	0.856	-0.250	-0.295	0.171	0.995	0.517	0.200	0.178

7

删除第三行数据重新分析：

```
lm(formula = y ~ ., data = rat[-3, ])
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.311    0.205   1.52   0.15
BodyWt      -0.008    0.019  -0.42   0.68
LiverWt      0.009    0.019   0.48   0.64
Dose         1.485    3.713   0.40   0.70

Multiple R-squared:  0.0211,
F-statistic: 0.1 on 3 and 14 DF, p-value: 0.958
```

去除BodyWt, Dose之一：

```
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.3268    0.1957   1.67   0.12
BodyWt      -0.0003    0.0013  -0.25   0.81
LiverWt      0.0065    0.0170   0.38   0.71
```

```
      Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.3227    0.198   1.63   0.12
LiverWt      0.0061    0.017   0.36   0.72
Dose        -0.0553    0.253  -0.22   0.83
```

所有变量都不显著。

8

lm( y ~ ., data = rat)				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.266	0.195	1.4	0.19
BodyWt	-0.021	0.008	-2.7	0.02 *
LiverWt	0.014	0.017	0.8	0.42
Dose	4.178	1.523	2.7	0.02 *

lm(formula = y ~ ., data = rat[-3, ])				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.311	0.205	1.52	0.15
BodyWt	-0.008	0.019	-0.42	0.68
LiverWt	0.009	0.019	0.48	0.64
Dose	1.485	3.713	0.40	0.70

lm( y ~ . - Dose, data = rat)				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.17836	0.22778	0.8	0.4
BodyWt	0.00035	0.00151	0.2	0.8
LiverWt	0.01233	0.02041	0.6	0.6

lm( y ~ . - Dose, data = rat[-3, ])				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.3268	0.1957	1.67	0.12
BodyWt	-0.0003	0.0013	-0.25	0.81
LiverWt	0.0065	0.0170	0.38	0.71

lm( y ~ . - BodyWt, data = rat)				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.1174	0.2191	0.5	0.6
LiverWt	0.0087	0.0201	0.4	0.7
Dose	0.1736	0.2863	0.6	0.6

lm( y ~ . - BodyWt, data = rat[-3, ])				
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.3227	0.198	1.63	0.12
LiverWt	0.0061	0.017	0.36	0.72
Dose	-0.0553	0.253	-0.22	0.83

Influence measures of  
lm(formula = y ~ . - Dose, data = rat) :

	dfb.1	dfb.BdyW	dfb.LvrW	dfbet	cov.r	cook.d	hat	inf
1	0.00102	0.28867	-0.43179	0.5333	1.103	0.092246	0.1647	
2	0.12469	0.16552	-0.49159	-0.6007	1.054	0.114968	0.1717	
3	-0.78638	0.53400	0.34751	1.1370	0.375	0.296102	0.1350	*

Influence measures of  
lm(formula = y ~ . - BodyWt, data = rat) :

	dfb.1	dfb.LvrW	dfb.Dose	dfbet	cov.r	cook.d	hat	inf
1	0.03135	-0.38849	0.24283	0.4831	1.128	0.076530	0.1558	
2	0.11665	-0.48512	0.17765	-0.5862	1.081	0.110240	0.1747	
3	-1.11935	0.15633	0.95488	1.3922	0.430	0.453009	0.1986	*

这可能表明，某些高影响点，只有在有严重复共线性时才是高影响的。最初模型(所有数据、所有变量)分析的问题在于自变量之间的高度共线性(复共线性)导致了高影响点。

9

10

## 置信域

模型:  $Y = X_{n \times p} \beta + \varepsilon$

$A$  为  $q \times p$  行满秩矩阵, 则正态假设下  $A\beta$  的  $(1-\alpha)100\%$  置信域:

$$\left\{ b = A\beta: (b - A\hat{\beta})' [A(X'X)^{-1}A']^{-1} (b - A\hat{\beta}) / q\hat{\sigma}^2 \leq F_{q, n-p}(1-\alpha) \right\}$$

$$\text{由 } \hat{\beta} \sim N(\beta, \sigma^2(X'X)^{-1}) \Rightarrow A\hat{\beta} \sim N(A\beta, \sigma^2 A(X'X)^{-1}A')$$

$$\Rightarrow \frac{(A\hat{\beta} - A\beta)'(A(X'X)^{-1}A')^{-1}(A\hat{\beta} - A\beta)}{\sigma^2} \sim \chi_q^2$$

$$\Rightarrow \frac{(A\hat{\beta} - A\beta)'(A(X'X)^{-1}A')^{-1}(A\hat{\beta} - A\beta) / q\sigma^2}{RSS / (n-p)\sigma^2} \sim F_{q, n-p}$$

分子只与  $\hat{Y} = P_X Y$  有关, 分母只与  $e = (I - P_X)Y$  有关, 故独立。

$$= \frac{(A\hat{\beta} - A\beta)'(A(X'X)^{-1}A')^{-1}(A\hat{\beta} - A\beta)}{q\hat{\sigma}^2} \text{ 是枢轴变量}$$

11

特别地,  $A = \mathbf{a}'$  是行向量时,  $\mathbf{a}'\beta$  的  $(1-\alpha)100\%$  置信区间:

$$\left\{ \mathbf{a}'\beta: |\mathbf{a}'\beta - \mathbf{a}'\hat{\beta}| \leq t_{n-p}(\alpha/2) \hat{\sigma} \sqrt{\mathbf{a}'(X'X)^{-1}\mathbf{a}} \right\}$$

特别地,  $A = (0, \dots, 1, 0, \dots, 0)'$ ,  $A\beta = \beta_k$  的  $(1-\alpha)100\%$  置信区间:

$$\left\{ \beta_k: |\beta_k - \hat{\beta}_k| \leq t_{n-p}(\alpha/2) \hat{\sigma} \sqrt{(X'X)^{-1}_{kk}} = t_{n-p}(\alpha/2) \hat{\sigma} / \|\mathbf{x}_k^\perp\| \right\}$$

其中  $\mathbf{x}_k$  为  $X$  的第  $k$  列,  $\mathbf{x}_k^\perp = \mathbf{x}_k - P_{X_{(-k)}} \mathbf{x}_k$

$A = (\mathbf{a}_1, \dots, \mathbf{a}_q)'$ ,  $A\beta$  的置信域即  $q$  个参数  $\mathbf{a}_1'\beta, \dots, \mathbf{a}_q'\beta$  的同时置信估计 (为一椭球)。

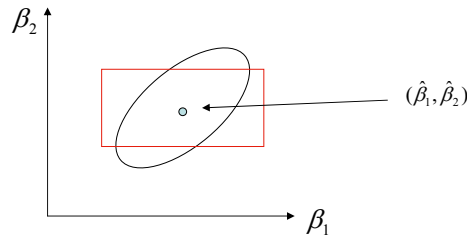
12

## 同时置信区间

为了简便实用, 对于多个参数, 能否构造长方形置信域

$$\{(\beta_1, \dots, \beta_q): L_j \leq \beta_j - \hat{\beta}_j \leq U_j, j=1, \dots, q\}$$

称为同时置信区间, simultaneous confidence interval



最为常用的是Bonferroni同时置信区间, 其它还有Scheffe同时置信区间, 对于方差分析模型, 有Tukey同时置信区间。

13

## Bonferroni 同时置信区间 :

$\mathbf{a}_1' \beta, \dots, \mathbf{a}_q' \beta$  的 Bonferroni  $(1-\alpha)100\%$  同时置信区间为

$$\bigcap_{j=1}^q \left\{ \mathbf{a}_j' \beta: |\mathbf{a}_j' \beta - \mathbf{a}_j' \hat{\beta}| \leq t_{n-p} \left( \frac{\alpha}{2q} \right) \hat{\sigma} \sqrt{\mathbf{a}_j' (X'X)^{-1} \mathbf{a}_j} \right\}$$

注意每一个  $\mathbf{a}_j' \beta$  的置信区间  $I_j$  的置信水平是  $(1-\alpha/q)$  而不是  $1-\alpha$

$$\text{证明: } P(I_j) = 1 - \frac{\alpha}{q}, \text{ 所以 } P\left(\bigcap_{j=1}^q I_j\right) = 1 - P\left(\bigcup_{j=1}^q I_j^c\right) \geq 1 - \sum_{j=1}^q P(I_j^c) = 1 - \alpha$$

即 Bonferroni 同时置信区间的置信度为  $1-\alpha$

$$\text{Bonferroni 不等式: } P\left(\bigcup_{j=1}^q E_j\right) \leq \sum_{j=1}^q P(E_j)$$

14

特别地,  $\beta_{i_1}, \dots, \beta_{i_q}$  的  $(1-\alpha)100\%$  Bonferroni 同时置信区间为

$$\bigcup_{j=i_1, \dots, i_q} \left\{ \beta_j: |\beta_j - \hat{\beta}_j| \leq t_{n-p} (\alpha/2q) \hat{\sigma} / \|\mathbf{x}_j^\perp\| \right\}$$

注1: 当  $I_j, j=1, \dots, q$  独立时 ( $\mathbf{a}_1, \dots, \mathbf{a}_q$  正交时)

$$P\left(\bigcap_{j=1}^q I_j\right) = \prod_{j=1}^q P(I_j) = (1-\alpha/q)^q \approx 1-\alpha$$

即 Bonferroni 方法基本能精确地控制置信水平。

注2: 而不独立时, Bonferroni 过于严格 (置信水平可能会远大于  $1-\alpha$ ), 因而是一种过于保守的方法。

15

## Scheffe 同时置信区间:

Scheffe 同时置信区间:

$$P\left(|\mathbf{a}' \beta - \mathbf{a}' \hat{\beta}| \leq \hat{\sigma} \sqrt{\mathbf{a}' (X'X)^{-1} \mathbf{a}} \sqrt{pF_{p, n-p}(1-\alpha)}, \forall \mathbf{a} \in R^p\right) \geq 1-\alpha$$

$$\text{引理: } \sup_u \{(u'v)^2 / u' Au\} \leq v' A^{-1} v, u, v \in R^n, A > 0$$

$$\text{证: 由 Cauchy 不等式 } (u'v)^2 \leq u' Au \times v' A^{-1} v.$$

$$\text{证: } P\left(|\mathbf{a}' \beta - \mathbf{a}' \hat{\beta}| \leq \hat{\sigma} \sqrt{\mathbf{a}' (X'X)^{-1} \mathbf{a}} \sqrt{pF_{p, n-p}(1-\alpha)}, \forall \mathbf{a} \in R^p\right)$$

$$= P\left(\sup_{\mathbf{a} \in R^p} \frac{|\mathbf{a}' \beta - \mathbf{a}' \hat{\beta}|}{\sqrt{\mathbf{a}' (X'X)^{-1} \mathbf{a}}} \leq \hat{\sigma} \sqrt{pF_{p, n-p}(1-\alpha)}\right)$$

$$= P\left((\beta - \hat{\beta})' X' X (\beta - \hat{\beta}) / p \hat{\sigma}^2 \leq F_{p, n-p}(1-\alpha)\right) = 1-\alpha$$

注: Scheffe 同时置信区间含无穷多个区间, 过长而较少使用

16

## 多重检验 (multiple testing)

与建立 $\beta_1, \dots, \beta_q$ 的同时置信区间类似的一个问题是同时检验多个假设

$$H_1: \beta_1 = 0, \dots, H_q: \beta_q = 0$$

更一般地，实际问题中可能考虑同时检验多个假设  $H_1, \dots, H_q$

如果每个检验的水平都是 $\alpha$ ，而且各个检验是独立的，那么同时检验 $q$ 个假设的I型错误率为 $1 - (1 - \alpha)^q \approx q\alpha$ 。

更一般地，实际问题中可能考虑同时检验多个假设  $H_1, \dots, H_q$

如果每个检验的水平都是 $\alpha$ ，而且各个检验是独立的，那么同时即使各个检验不独立，同时检验 $q$ 个假设的I型错误率也会远大于 $\alpha$ （由Bonferroni不等式，至多为 $q\alpha$ ）。

为此需要校正（下调）单个检验的水平。

17

Bonferroni校正方法将单个检验的水平调整为 $\alpha/q$ ，由Bonferroni不等式， $q$ 个检验的I型错误率不会超过 $\alpha$ 。

换个角度看，Bonferroni方法将 $q$ 个检验的最小p值 $\min p_j$ 调整为 $q \times \min p_j$

类似于Bonferroni同时置信区间，将单个检验的水平下调为 $\alpha/q$ 对于独立情形是精确的；对于不独立情形可能太谨慎/保守了。所以Bonferroni方法主要用于检验之间独立，或者不独立当相关结构未知的情形。

18

```
> lm(log(Species)~., data=lakes)->a
Coefficients:
      Estimate Std. Error t    value Pr(>|t|)
(Intercept)  3.638e+00  5.313e-01  6.847  2.87e-07 ***
MaxDepth    -3.299e-04  2.815e-03 -0.117  0.9076
MeanDepth    6.357e-04  3.993e-03  0.159  0.8747
Cond        -2.334e-05  2.203e-04 -0.106  0.9164
Elev        -1.976e-04  8.751e-05 -2.258  0.0326 *
Lat         -2.345e-02  9.115e-03 -2.573  0.0161 *
Long        -4.422e-05  3.000e-03 -0.015  0.9884
Dist        -8.352e-02  4.569e-02 -1.828  0.0790 .
NLakes       6.471e-05  1.577e-04  0.410  0.6850
Photo       1.653e-04  2.298e-04  0.719  0.4784
Area        1.337e-07  6.051e-08  2.209  0.0362 *
---
F-statistic: 5.226 on 10 and 26 DF, p-value: 0.0003291
```

比如上述结果中如果同时检验10个自变量是否存在某个或某些是显著的，最小p值为0.0161，Bonferroni校正后的p值：0.161不显著。（与回归方程的显著性F检验相比过于保守）。

19