

Stentor的天空

Stentor的技术随笔

多元线性回归

能用office07发布简直是太好了，这下子省了很多事。

1、多元线性回归模型

假定被解释变量 Y 与多个解释变量 X_1, X_2, \dots, X_k 之间具有线性关系，是解释变量的多元线性函数，称为多元线性回归模型。即

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \mu \tag{1.1}$$

其中 Y 为被解释变量， $X_j (j = 1, 2, \dots, k)$ 为 k 个解释变量， $\beta_j (j = 0, 1, 2, \dots, k)$ 为 $k+1$ 个未知参数， μ 为随机误差项。

被解释变量 Y 的期望值与解释变量 X_1, X_2, \dots, X_k 的线性方程为：

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k \tag{1.2}$$

称为多元总体线性回归方程，简称总体回归方程。

对于 n 组观测值 $Y_i, X_{1i}, X_{2i}, \dots, X_{ki} (i = 1, 2, \dots, n)$ ，其方程组形式为：

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \mu_i, (i = 1, 2, \dots, n) \tag{1.3}$$

即

$$\begin{cases} Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{21} + \dots + \beta_k X_{k1} + \mu_1 \\ Y_2 = \beta_0 + \beta_1 X_{12} + \beta_2 X_{22} + \dots + \beta_k X_{k2} + \mu_2 \\ \dots\dots\dots \\ Y_n = \beta_0 + \beta_1 X_{1n} + \beta_2 X_{2n} + \dots + \beta_k X_{kn} + \mu_n \end{cases}$$

其矩阵形式为

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

即

$$Y = X\beta + \mu \tag{1.4}$$

其中

$$Y_{n \times 1} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \text{为被解释变量的观测值向量；} \quad X_{n \times (k+1)} = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{bmatrix} \text{为解释变量的观测值矩阵；}$$

$$\beta_{(k+1) \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} \text{为总体回归参数向量；} \quad \mu_{n \times 1} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \text{为随机误差项向量。}$$

公告

昵称: [zgw21cn](#)
园龄: [11年6个月](#)
粉丝: [68](#)
关注: [1](#)
[+加关注](#)

导航

[博客园](#)
[首页](#)
[新随笔](#)
[联系](#)
[订阅](#) [RSS](#)
[管理](#)

< 2008年12月 >						
日	一	二	三	四	五	六
30	1	2	3	4	5	6
7	8	9	10	11	12	13
14	15	16	17	18	19	20
21	22	23	24	25	26	27
28	29	30	31	1	2	3
4	5	6	7	8	9	10

统计

随笔 - 54
文章 - 0
评论 - 64
引用 - 0

搜索

找找看

谷歌搜索

常用链接

[我的随笔](#)
[我的评论](#)
[我的参与](#)
[最新评论](#)
[我的标签](#)

我的标签

[SAS](#)(21)
[learn](#)(11)
[SQL](#)(4)
[R SVM](#)(3)
[ROC Excel](#)(1)
[SAS 变量名 宏变量](#)(1)
[SQL 最小二乘法](#)(1)
[TSR](#)(1)
[VBA](#)(1)
[关联规则 Apriori](#)(1)
[更多](#)

随笔分类

总体回归方程表示为：

$$E(Y) = X\beta \quad (1.5)$$

多元线性回归模型包含多个解释变量，多个解释变量同时对被解释变量 Y 发生作用，若要考察其中一个解释变量对 Y 的影响就必须假设其它解释变量保持不变来进行分析。因此多元线性回归模型中的回归系数为偏回归系数，即反映了当模型中的其它变量不变时，其中一个解释变量对因变量 Y 的均值的影响。

由于参数 $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ 都是未知的,可以利用样本观测值 $(X_{1i}, X_{2i}, \dots, X_{ki}, Y_i)$ 对它们进行估计。若计算得到的参数估计值为 $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$ ，用参数估计值替代总体回归函数的未知参数 $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ ，则得多元线性样本回归方程：

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki} \quad (1.6)$$

其中 $\hat{\beta}_j(j = 0, 1, 2, \dots, k)$ 为参数估计值， $\hat{Y}_i(i = 1, 2, \dots, n)$ 为 Y_i 的样本回归值或样本拟合值、样本估计值。

其矩阵表达形式为：

$$\hat{Y} = X\hat{\beta} \quad (1.7)$$

其中 $\hat{Y}_{n \times 1} = \begin{bmatrix} \hat{Y}_1 \\ \hat{Y}_2 \\ \vdots \\ \hat{Y}_n \end{bmatrix}$ 为被解释变量样本观测值向量 Y 的 $n \times 1$ 阶拟合值列向量； $X_{n \times (k+1)} = \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{bmatrix}$ 为解释变量 X 的 $n \times (k+1)$ 阶样本观测矩阵； $\hat{\beta}_{(k+1) \times 1} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$ 为未知参数向量 β 的 $(k+1) \times 1$ 阶估计值列向量。

样本回归方程得到的被解释变量估计值 \hat{Y}_i 与实际观测值 Y_i 之间的偏差称为残差 e_i 。

$$e_i = Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \dots + \hat{\beta}_k X_{ki}) \quad (1.8)$$

2、多元线性回归模型的假定

与一元线性回归模型相同，多元线性回归模型利用普通最小二乘法(OLS)对参数进行估计时，有如下假定：

假定1 零均值假定： $E(\mu_i) = 0, i = 1, 2, \dots, n$ ，即

$$E(\mu) = E \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \begin{bmatrix} E(\mu_1) \\ E(\mu_2) \\ \vdots \\ E(\mu_n) \end{bmatrix} = 0 \quad (2.1)$$

假定2 同方差假定(μ_i 的方差为同一常数)：

$$Var(\mu_i) = E(\mu_i^2) = \sigma^2, (i = 1, 2, \dots, n) \quad (2.2)$$

假定3 无自相关性：

$$Cov(\mu_i, \mu_j) = E(\mu_i \mu_j) = 0, (i \neq j, i, j = 1, 2, \dots, n)$$

Algorithm for Data Mining(9)
Data Structure and Algorithm
Database System(SQL Server,Oracle)(2)
Dynamic Language (1)
Numerical Computation(1)
Office Software(VBA etc)(4)
Software
Design(Language,Disign, Test etc)(2)
Statistic Analysis(SAS,R etc) (34)

随笔档案

- 2011年1月(1)
- 2010年5月(1)
- 2009年11月(2)
- 2009年9月(1)
- 2009年5月(2)
- 2009年4月(2)
- 2009年3月(3)
- 2009年2月(3)
- 2009年1月(1)
- 2008年12月(2)
- 2008年11月(3)
- 2008年10月(5)
- 2008年9月(4)
- 2008年8月(24)

Statistic

Agri521'Blog
我的英语技术博客

最新评论

- 1. Re:逻辑回归模型
能发一份word吗?
2602196069@qq.com
---厘米的阳光
- 2. Re:逻辑回归模型
能给一份word吗?
953022837@qq.com
--katty_fat
- 3. Re:多元线性回归
学神!
--知好奇行
- 4. Re:多元线性回归
质量很高，很多没看懂。。。
--AutumnLight
- 5. Re:多元线性回归
数学太好了！！我看的一头雾水。。。
--Lansing999

阅读排行榜

- 1. 多元线性回归(17073)
- 2. 逻辑回归模型(16897)
- 3. 关联规则(apriori algorithm)(15010)
- 4. Excel中如何作特大值数据的条形图(11882)
- 5. 在R中使用支持向量机(SVM)(2.Kernlab包)(11207)

评论排行榜

- 1. 逻辑回归模型(22)
- 2. 关联规则(apriori algorithm)(18)
- 3. 多元线性回归(7)

$$\begin{aligned}
 E(\boldsymbol{\mu}\boldsymbol{\mu}^T) &= E \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} (\mu_1, \mu_2, \dots, \mu_n) = E \begin{bmatrix} \mu_1^2 & \mu_1\mu_2 & \cdots & \mu_1\mu_n \\ \mu_2\mu_1 & \mu_2^2 & \cdots & \mu_2\mu_n \\ \vdots & \vdots & \ddots & \vdots \\ \mu_n\mu_1 & \mu_n\mu_2 & \cdots & \mu_n^2 \end{bmatrix} \\
 &= \begin{bmatrix} E(\mu_1^2) & E(\mu_1\mu_2) & \cdots & E(\mu_1\mu_n) \\ E(\mu_2\mu_1) & E(\mu_2^2) & \cdots & E(\mu_2\mu_n) \\ \vdots & \vdots & \ddots & \vdots \\ E(\mu_n\mu_1) & E(\mu_n\mu_2) & \cdots & E(\mu_n^2) \end{bmatrix} \\
 &= \begin{bmatrix} \sigma_\mu^2 & 0 & \cdots & 0 \\ 0 & \sigma_\mu^2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_\mu^2 \end{bmatrix} = \sigma_\mu^2 \mathbf{I}_n \quad (2.3)
 \end{aligned}$$

假定4 随机误差项 μ 与解释变量 X 不相关(这个假定自动成立):

$$\text{Cov}(X_{\bar{j}}, \mu_i) = 0, (j = 1, 2, \dots, k, i = 1, 2, \dots, n) \quad (2.4)$$

假定5 随机误差项 μ 服从均值为零, 方差为 σ^2 的正态分布:

$$\mu_i \sim N(0, \sigma_\mu^2 \mathbf{I}_n) \quad (2.5)$$

假定6 解释变量之间不存在多重共线性:

$$\text{rank}(\mathbf{X}) = k+1 \leq n$$

即各解释变量的样本观测值之间线性无关, 解释变量的样本观测值矩阵 \mathbf{X} 的秩为参数个数 $k+1$, 从而保证参数 $\beta_0, \beta_1, \beta_2, \dots, \beta_k$ 的估计值唯一。

3、多元线性回归模型的参数估计

3.1 回归参数的最小二乘估计

对于含有 k 个解释变量的多元线性回归模型

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + \mu_i \quad (i = 1, 2, \dots, n)$$

设 $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ 分别作为参数 $\beta_0, \beta_1, \dots, \beta_k$ 的估计量, 得样本回归方程为:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}$$

观测值 Y_i 与回归值 \hat{Y}_i 的残差 e_i 为:

$$e_i = Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki})$$

由最小二乘法可知 $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ 应使全部观测值 Y_i 与回归值 \hat{Y}_i 的残差 e_i 的平方和最小, 即使

$$\begin{aligned}
 Q(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k) &= \sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 \\
 &= \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \cdots - \hat{\beta}_k X_{ki})^2 \quad (3.1)
 \end{aligned}$$

取得最小值。根据多元函数的极值原理, Q 分别对 $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_k$ 求一阶偏导, 并令其等于零, 即

$$\frac{\partial Q}{\partial \hat{\beta}_j} = 0, (j = 1, 2, \dots, k) \quad (3.2)$$

即

4. Excel中如何作特大值数据的条形图(6)

5. 遗传编程算法(4)

推荐排行榜

1. 逻辑回归模型(4)

2. 多元线性回归(4)

3. 通过示例学SAS (1)--从外部文件读入数据(3)

4. 在R中使用支持向量机(SVM) (1) (2)

5. 遗传编程算法(2)

$$\begin{cases} \frac{\partial Q}{\partial \hat{\beta}_0} = 2 \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki})(-1) = 0 \\ \frac{\partial Q}{\partial \hat{\beta}_1} = 2 \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki})(-X_{1i}) = 0 \\ \dots\dots \\ \frac{\partial Q}{\partial \hat{\beta}_k} = \sum (Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \dots - \hat{\beta}_k X_{ki})(-X_{ki}) = 0 \end{cases}$$

化简得下列方程组

$$\begin{cases} n\hat{\beta}_0 + \hat{\beta}_1 \sum X_{1i} + \hat{\beta}_2 \sum X_{2i} + \dots + \hat{\beta}_k \sum X_{ki} = \sum Y_i \\ \hat{\beta}_0 \sum X_{1i} + \hat{\beta}_1 \sum X_{1i}^2 + \hat{\beta}_2 \sum X_{2i}X_{1i} + \dots + \hat{\beta}_k \sum X_{ki}X_{1i} = \sum X_{1i}Y_i \\ \dots\dots \\ \hat{\beta}_0 \sum X_{ki} + \hat{\beta}_1 \sum X_{1i}X_{ki} + \hat{\beta}_2 \sum X_{2i}X_{ki} + \dots + \hat{\beta}_k \sum X_{ki}^2 = \sum X_{ki}Y_i \end{cases} \quad (3.3)$$

上述 $(k+1)$ 个方程称为正规方程，其矩阵形式为

$$\begin{bmatrix} n & \sum X_{1i} & \sum X_{2i} & \dots & \sum X_{ki} \\ \sum X_{1i} & \sum X_{1i}^2 & \sum X_{2i}X_{1i} & \dots & \sum X_{ki}X_{1i} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum X_{ki} & \sum X_{1i}X_{ki} & \sum X_{2i}X_{ki} & \dots & \sum X_{ki}^2 \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix} = \begin{bmatrix} \sum Y_i \\ \sum X_{1i}Y_i \\ \vdots \\ \sum X_{ki}Y_i \end{bmatrix} \quad (3.4)$$

因为

$$\begin{bmatrix} n & \sum X_{1i} & \sum X_{2i} & \dots & \sum X_{ki} \\ \sum X_{1i} & \sum X_{1i}^2 & \sum X_{2i}X_{1i} & \dots & \sum X_{ki}X_{1i} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \sum X_{ki} & \sum X_{1i}X_{ki} & \sum X_{2i}X_{ki} & \dots & \sum X_{ki}^2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ X_{k1} & X_{k2} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} 1 & X_{11} & X_{21} & \dots & X_{k1} \\ 1 & X_{12} & X_{22} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{1n} & X_{2n} & \dots & X_{kn} \end{bmatrix} = \mathbf{X}'\mathbf{X}$$

$$\begin{bmatrix} \sum Y_i \\ \sum X_{1i}Y_i \\ \vdots \\ \sum X_{ki}Y_i \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ X_{k1} & X_{k2} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \mathbf{X}'\mathbf{Y}$$

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$$

设 $\begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{bmatrix}$ 为估计值向量

样本回归模型 $\mathbf{Y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}$ 两边同乘样本观测值矩阵 \mathbf{X} 的转置矩阵 \mathbf{X}' ，则有

$$\mathbf{X}'\mathbf{Y} = \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{X}'\mathbf{e}$$

得正规方程组：

$$\mathbf{X}'\mathbf{Y} = \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \quad (3.5)$$

由假定(6)， $R(\mathbf{X}) = k+1$ ， $\mathbf{X}'\mathbf{X}$ 为 $(k+1)$ 阶方阵，所以 $\mathbf{X}'\mathbf{X}$ 满秩， $\mathbf{X}'\mathbf{X}$ 的逆矩阵 $(\mathbf{X}'\mathbf{X})^{-1}$ 存在。因而

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \quad (3.6)$$

则为向量 $\boldsymbol{\beta}$ 的OLS估计量。

以二元线性回归模型为例，导出二元线性回归模型的OLS估计量的表达式。由(1.3)式得二元线性回归模型为

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \mu_i$$

为了计算的方便，先将模型中心化。

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ji}, x_{ji} = X_{ji} - \bar{X}_j, (j=1,2)$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, y_i = Y_i - \bar{Y}$$

$$L_{pq} = \sum x_{pi} x_{qi}, (p, q = 1, 2)$$

$$L_{jY} = \sum x_{ji} y_i, (j = 1, 2)$$

$$L_{YY} = \sum y_i^2$$

设 $\alpha_0 = \beta_0 + \beta_1 \bar{X}_1 + \beta_2 \bar{X}_2$ ，则二元回归模型改写为中心化模型。

$$Y_i = \alpha_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \mu_i \quad (3.7)$$

记

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{21} \\ 1 & x_{12} & x_{22} \end{bmatrix}, \boldsymbol{\beta} = \begin{bmatrix} \alpha_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} n & 0 & 0 \\ 0 & \sum x_{1i}^2 & \sum x_{1i} x_{2i} \\ 0 & \sum x_{2i} x_{1i} & \sum x_{2i}^2 \end{bmatrix}, \mathbf{X}'\mathbf{Y} = \begin{bmatrix} \sum Y_i \\ \sum x_{1i} Y_i \\ \sum x_{2i} Y_i \end{bmatrix} \quad (3.8)$$

将 $L_{pq} = \sum x_{pi} x_{qi}, (p, q = 1, 2)$ 代入得

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} n & 0 & 0 \\ 0 & L_{11} & L_{12} \\ 0 & L_{21} & L_{22} \end{bmatrix} \quad (3.9)$$

因为

$$\begin{aligned} \sum_{i=1}^n x_{ji} Y_i &= \sum_{i=1}^n x_{ji} (y_i + \bar{Y}) = \sum_{i=1}^n x_{ji} y_i + \bar{Y} \sum_{i=1}^n x_{ji} \\ &= \sum_{i=1}^n x_{ji} y_i = L_{jY}, (j=1,2) \end{aligned} \quad (3.10)$$

则

$$\mathbf{X}'\mathbf{Y} = \begin{bmatrix} \sum Y_i \\ L_{1Y} \\ L_{2Y} \end{bmatrix}$$

由(3.6)式得

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = \begin{bmatrix} \frac{1}{n} & 0 \\ 0 & \mathbf{L}^{-1} \end{bmatrix} \begin{bmatrix} \sum Y_i \\ L_{1Y} \\ L_{2Y} \end{bmatrix} \quad (3.11)$$

其中

$$\mathbf{L}^{-1} = \begin{bmatrix} L_{11} & L_{12} \\ L_{12} & L_{22} \end{bmatrix}^{-1} = \frac{1}{L_{11}L_{22} - L_{12}^2} \begin{bmatrix} L_{22} & -L_{12} \\ -L_{12} & L_{11} \end{bmatrix}$$

由(3.11)式可知

$$\hat{\alpha}_0 = \bar{Y}$$

$$\begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \mathbf{L}^{-1} \begin{bmatrix} L_{1Y} \\ L_{2Y} \end{bmatrix} = \frac{1}{L_{11}L_{22} - L_{12}^2} \begin{bmatrix} L_{22} & -L_{12} \\ -L_{12} & L_{11} \end{bmatrix} \begin{bmatrix} L_{1Y} \\ L_{2Y} \end{bmatrix}$$

得

$$\hat{\beta}_1 = \frac{L_{1Y}L_{22} - L_{2Y}L_{12}}{L_{11}L_{22} - L_{12}^2} \quad (3.12)$$

$$\hat{\beta}_2 = \frac{L_{2Y}L_{11} - L_{1Y}L_{12}}{L_{11}L_{22} - L_{12}^2} \quad (3.13)$$

$$\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1\bar{X}_1 - \hat{\beta}_2\bar{X}_2 \quad (3.14)$$

3.2 随机误差项 μ 的方差 σ_μ^2 的估计量

样本回归方程得到的被解释变量估计值 \hat{Y}_i 与实际观测值 Y_i 之间的偏差称为残差 e_i

$$e_i = Y_i - \hat{Y}_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki})$$

则

$$\begin{aligned} \mathbf{e} &= \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}} = (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}) - \mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}] \\ &= (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}) - \mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu})] \\ &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu} - \mathbf{X}[\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu}] \\ &= \boldsymbol{\mu} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu} \\ &= [\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\boldsymbol{\mu} \end{aligned}$$

设 $\mathbf{P} = \mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ ，可以得出 \mathbf{P} 是 n 阶对称幂等矩阵， $\mathbf{P} = \mathbf{P}'$ ， $\mathbf{P}^2 = \mathbf{P}$ 。于是

$$\mathbf{e} = \mathbf{P}\boldsymbol{\mu}$$

而残差的平方和为

$$\begin{aligned} \sum e_i^2 &= \mathbf{e}'\mathbf{e} = (\mathbf{P}\boldsymbol{\mu})'(\mathbf{P}\boldsymbol{\mu}) = \boldsymbol{\mu}'\mathbf{P}'\mathbf{P}\boldsymbol{\mu} = \boldsymbol{\mu}'\mathbf{P}\boldsymbol{\mu} \\ &= \boldsymbol{\mu}'[\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\boldsymbol{\mu} \\ E(\mathbf{e}'\mathbf{e}) &= E(\boldsymbol{\mu}'[\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\boldsymbol{\mu}) \\ &= \sigma_\mu^2 \text{tr}[\mathbf{I}_n - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \sigma_\mu^2 [\text{tr}\mathbf{I}_n - \text{tr}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \\ &= \sigma_\mu^2 [n - (k+1)] \end{aligned}$$

其中“ tr ”表示矩阵的迹，即矩阵主对角线元素的和。于是

$$\sigma_\mu^2 = \frac{E(\mathbf{e}'\mathbf{e})}{n - (k+1)} = E\left(\frac{\mathbf{e}'\mathbf{e}}{n - (k+1)}\right)$$

随机误差项 μ 的方差 σ_μ^2 的无偏估计量，记作 S_e^2 ，即 $E(S_e^2) = \sigma_\mu^2$ ， $S_e^2 = \hat{\sigma}_\mu^2$ ， S_e 为残差的标准差(或回归标准差)。

因此

$$S_e^2 = \frac{\sum e_i^2}{n-k-1} = \frac{\mathbf{e}'\mathbf{e}}{n-k-1} \quad (3.15)$$

其中

$$\begin{aligned} \sum e_i^2 &= \mathbf{e}'\mathbf{e} = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\ &= \mathbf{Y}'\mathbf{Y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} \\ &= \mathbf{Y}'\mathbf{Y} - 2\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} \\ &= \mathbf{Y}'\mathbf{Y} - \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{Y} \quad (3.16) \end{aligned}$$

例如,对于二元线性回归模型($k=2$)

$$S_e^2 = \frac{\mathbf{e}'\mathbf{e}}{n-3} = \frac{\sum e_i^2}{n-3} \quad (3.17)$$

$$\begin{aligned} \sum e_i^2 &= \mathbf{e}'\mathbf{e} = L_{YY} - \hat{\beta}_1 L_{1Y} - \hat{\beta}_2 L_{2Y} \\ &= \sum Y_i^2 - \hat{\beta}_1 \sum X_{1i} Y_i - \hat{\beta}_2 \sum X_{2i} Y_i \quad (3.18) \end{aligned}$$

3.3、估计参数的统计性质

1、线性性

指最小二乘估计量 $\hat{\boldsymbol{\beta}}$ 是被解释变量的观测值 Y_1, Y_2, \dots, Y_n 的线性函数。

由于

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$$

设 $\mathbf{P} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, 则矩阵 \mathbf{P} 为一非随机的 $(k+1) \times n$ 阶常数矩阵。所以

$$\hat{\boldsymbol{\beta}} = \mathbf{P}\mathbf{Y} \quad (3.19)$$

显然最小二乘估计量 $\hat{\boldsymbol{\beta}}$ 是被解释变量的观测值 Y_1, Y_2, \dots, Y_n 的线性函数。

2、无偏性

将 $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}$ 代入(3-16)式得

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu} \\ &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu} \quad (3.20) \end{aligned}$$

则

$$\begin{aligned} E(\hat{\boldsymbol{\beta}}) &= \boldsymbol{\beta} + E[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\mu}] \\ &= \boldsymbol{\beta} + (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'E(\boldsymbol{\mu}) \\ &= \boldsymbol{\beta} \end{aligned}$$

所以 $\hat{\boldsymbol{\beta}}$ 是 $\boldsymbol{\beta}$ 的无偏估计量。

3.最小方差性

设 \mathbf{P} 为 $n \times P$ 阶数值矩阵, \mathbf{X} 为 $P \times n$ 阶随机矩阵(随机变量为元素的矩阵), \mathbf{Q} 为 $n \times n$ 阶数值矩阵, 则

$$E(\mathbf{P}\mathbf{X}\mathbf{Q}) = \mathbf{P}(E(\mathbf{X}))\mathbf{Q}$$

下面推导 $\hat{\boldsymbol{\beta}}$ 的方差、协方差矩阵。

$$Var(\hat{\boldsymbol{\beta}}) = E\left[(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\right]$$

定义:

$$= E \left[\begin{bmatrix} \hat{\beta}_0 - \beta_0 \\ \hat{\beta}_1 - \beta_1 \\ \vdots \\ \hat{\beta}_k - \beta_k \end{bmatrix} (\hat{\beta}_0 - \beta_0, \hat{\beta}_1 - \beta_1, \dots, \hat{\beta}_k - \beta_k) \right]$$

$$= \begin{bmatrix} \text{Var}(\hat{\beta}_0) & \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) & \dots & \text{Cov}(\hat{\beta}_0, \hat{\beta}_k) \\ \text{Cov}(\hat{\beta}_1, \hat{\beta}_0) & \text{Var}(\hat{\beta}_1) & \dots & \text{Cov}(\hat{\beta}_1, \hat{\beta}_k) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(\hat{\beta}_k, \hat{\beta}_0) & \text{Cov}(\hat{\beta}_k, \hat{\beta}_1) & \dots & \text{Var}(\hat{\beta}_k) \end{bmatrix}$$

由(3.20)式得

$$\hat{\beta} - \beta = (X'X)^{-1} X' \mu$$

$$(\hat{\beta} - \beta)' = [(X'X)^{-1} X' \mu]' = \mu' X (X'X)^{-1}$$

所以

$$\text{Var}(\hat{\beta}) = E[(\hat{\beta} - \beta)(\hat{\beta} - \beta)']$$

$$= E[(X'X)^{-1} X' \mu \mu' X (X'X)^{-1}]$$

$$= (X'X)^{-1} X' E(\mu \mu') X (X'X)^{-1}$$

$$= (X'X)^{-1} X' \sigma_{\mu}^2 I_n X (X'X)^{-1}$$

$$= \sigma_{\mu}^2 (X'X)^{-1} \quad (3.21)$$

这个矩阵主对角线上的元素表示 $\hat{\beta}$ 的方差，非主对角线上的元素表示 $\hat{\beta}$ 的协方差。例如 $\text{Var}(\hat{\beta}_i)$ 是位于 $\sigma_{\mu}^2 (X'X)^{-1}$ 的第 i 行与第 i 列交叉处的元素(主对角线上的元素)； $\text{Cov}(\hat{\beta}_i, \hat{\beta}_j)$ 是位于 $\sigma_{\mu}^2 (X'X)^{-1}$ 的第 i 行与第 j 列交叉处的元素(非主对角线上的元素)

在应用上，我们关心的 $\hat{\beta}$ 的方差，而忽略协方差，因此把(3.21)式记作

$$\text{Var}(\hat{\beta}) = \sigma_{\mu}^2 (X'X)^{-1}_{ii} \quad (3.22)$$

记 $S^{-1} = (X'X)^{-1} = (C_{ij})$, ($i, j = 0, 1, 2, \dots, k$)，则 $\text{Var}(\hat{\beta}_i) = \sigma_{\mu}^2 C_{ii}$ ，所以 $\hat{\beta}$ 是 β 的最小方差线性无偏估计。这说明，在(1.1)式系数的无偏估计量中，OLS估计量的方差比用其它估计方法所得的无偏估计量的方差都要小，这正是OLS的优越性所在。

用 S_{ϵ}^2 代替 σ_{μ}^2 则得 $\hat{\beta}_i$ 的标准估计量的估计值，乃称为标准差。

$$S(\hat{\beta}_i) = \sqrt{C_{ii} S_{\epsilon}^2} \quad (3.23)$$

其中

$$S_{\epsilon}^2 = \frac{e'e}{n-k-1}$$

对于二元回归模型($k=2$)，求估计量 $\hat{\beta}_1, \hat{\beta}_2$ 的方差，由(3.22)式得

$$\text{Var}(\hat{\beta}) = \sigma_{\mu}^2 (X'X)^{-1}_{ii} = \sigma_{\mu}^2 \begin{bmatrix} \frac{1}{n} & 0 \\ 0 & L^{-1} \end{bmatrix}_{ii}$$

其中

$$L = \begin{bmatrix} L_{11} & L_{12} \\ L_{12} & L_{22} \end{bmatrix}$$

于是

$$\text{Var}\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \sigma_{\mu}^2 L_{ii}^{-1} \frac{\sigma_{\mu}^2}{L_{11}L_{22} - L_{12}^2} \begin{bmatrix} L_{22} & -L_{12} \\ -L_{12} & L_{11} \end{bmatrix}$$

所以

$$\text{Var}(\hat{\beta}_1) = \sigma^2(\hat{\beta}_1) = \frac{L_{22}}{L_{11}L_{22} - L_{12}^2} \sigma_{\mu}^2 \quad (3.24)$$

$$\text{Var}(\hat{\beta}_2) = \sigma^2(\hat{\beta}_2) = \frac{L_{11}}{L_{11}L_{22} - L_{12}^2} \sigma_{\mu}^2 \quad (3.25)$$

$$S(\hat{\beta}_1) = \sqrt{\frac{L_{22}}{L_{11}L_{22} - L_{12}^2} S_e^2} \quad (3.26)$$

$$S(\hat{\beta}_2) = \sqrt{\frac{L_{11}}{L_{11}L_{22} - L_{12}^2} S_e^2} \quad (3.27)$$

其中

$$S_e^2 = \frac{\mathbf{e}'\mathbf{e}}{n-3}$$

4. 显著性检验

4.1 拟合优度检验

4.1.1 总离差平方和分解

设具有 k 个解释变量的回归模型为

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + \mu_i$$

其回归方程为

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}$$

离差分解：

$$Y_i - \bar{Y} = (Y_i - \hat{Y}_i) + (\hat{Y}_i - \bar{Y})$$

总离差平方和分解式为：

$$\sum (Y_i - \bar{Y})^2 = \sum (\hat{Y}_i - \bar{Y})^2 + \sum (Y_i - \hat{Y}_i)^2 \quad (4.1)$$

即

$TSS = ESS + RSS$ (4.2) 总离差平方和分解为回归平方和与残差平方和两部分。 $\sum (Y_i - \bar{Y})^2$ 体现了观测值 y_1, y_2, \dots, y_n 总波动大小，称为总偏差平方和，记作TSS。 $\sum (\hat{Y}_i - \bar{Y})^2$ 体现了 n 个估计值 $\hat{y}_1, \hat{y}_2, \dots, \hat{y}_n$ 的波动大小，它是由于 Y 与自变量 x_1, x_2, \dots, x_k 的变化而引起，被称作为回归平方和，记为ESS (Explained Sum of Squares) 或U； $\sum (Y_i - \hat{Y}_i)^2$ 称为残差平方和，记为RSS (Residual Sum of Squares) 或Q。

4.1.2 样本决定系数

对于多元回归方程，其样本决定系数为复决定系数或多重决定系数。

$$R_{YX}^2, (i = 1, 2, \dots, k), \text{ 简记为 } R^2。$$

$$R^2 = \frac{ESS}{TSS} \quad (4.3)$$

根据式(4.2)

$$R^2 = 1 - \frac{RSS}{TSS} \quad (4.4)$$

因为

$$TSS = \sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2$$

由(3.16)式知

$$RSS = Y'Y - \hat{\beta}'X'Y$$

所以

$$ESS = TSS - RSS = \hat{\beta}'X'Y - n\bar{Y}^2$$

$$R^2 = \frac{\hat{\beta}'X'Y - n\bar{Y}^2}{Y'Y - n\bar{Y}^2} \quad (4.5)$$

R^2 作为检验回归方程与样本值拟合优度的指标： $R^2 (0 \leq R^2 \leq 1)$ 越大，表示回归方程与样本拟合的越好；反之，回归方程与样本值拟合较差。

具体的，当 $k = 2$ 时，求样本决定系数

$$R^2 = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = \frac{\sum y_i^2 - \sum e_i^2}{\sum y_i^2}$$

由(3.8)式，得 $\sum e_i^2 = L_{YY} - \hat{\beta}_1 L_{1Y} - \hat{\beta}_2 L_{2Y}$ ，因此有

$$R^2 = \frac{\hat{\beta}_1 L_{1Y} + \hat{\beta}_2 L_{2Y}}{L_{YY}} \quad (4.6)$$

4.1.3 调整后的样本决定系数

在使用 R^2 时，容易发现 R^2 的大小与模型中的解释变量的数目有关。如果模型中增加一个新解释变量，总离差 TSS 不会改变，但总离差中由解释变量解释的部分，即回归平方和 ESS 将会增加，这就是说 R^2 与模型中解释变量个数有关。但通过增加模型中解释变量的数目而使 R^2 增大是错误的，显然这样 R^2 来检验被回归方程与样本值拟合优度是不合适的，需要对 R^2 进行调整，使它不但能说明已被解释离差与总离差的关系，而且又能说明自由度的数目。

以 \bar{R}^2 表示调整样本决定系数，

$$\bar{R}^2 = 1 - \frac{S_e^2}{S_y^2} \quad (4.7)$$

其中

$$S_e^2 = \frac{\sum e_i^2}{n-k-1}, S_y^2 = \frac{\sum (Y_i - \bar{Y})^2}{n-1}$$

这里 $n-k-1$ 是残差平方和的自由度， $n-1$ 是总离差平方和的自由度。

由(4.7)式得

$$\bar{R}^2 = 1 - \frac{\sum e_i^2}{\sum (Y_i - \bar{Y})^2} \times \frac{n-1}{n-k-1} = 1 - (1 - R^2) \frac{n-1}{n-k-1}$$

其中， n 是样本观测值的个数， k 是解释变量的个数。从式中可以看出，当增加一个解释变量时，由前面分析可知

R^2 会增加，引起 $(1 - R^2)$ 减少，而 $\frac{n-1}{n-k-1}$ 增加，因而 \bar{R}^2 不会增加。这样用 \bar{R}^2 判定回归方程拟合优度，就消除了 R^2 对解释变量个数的依赖。

R^2 或 \bar{R}^2 只能说明在给定的样本条件下回归方程与样本观测值拟合优度，并不能做出对总体模型的推测，因此不能单凭 R^2 或 \bar{R}^2 来选择模型，必须对回归方程和模型中各参数的估计量做显著性检验。

4.2 方程显著性检验

由离差平方和分解(4.2)式可知, 总离差平方和 TSS 的自由度为 $n-1$, 回归平方和 ESS 是由 k 个解释变量 X_1, X_2, \dots, X_k 对 Y 的线性影响决定的。因此它的自由度为 k 。所以, 残差平方和的自由度由总离差平方和的自由度减去回归平方和的自由度, 即为 $n-k-1$ 。

检验回归方程是否显著,

第一步, 作出假设

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

备择假设 H_1 : $\beta_1, \beta_2, \dots, \beta_k$ 不同时为0

第二步, 在 H_0 成立的条件下, 计算统计量 F

$$F = \frac{ESS/k}{RSS/(n-k-1)} \sim F(k, n-k-1)$$

第三步, 查表临界值

对于假设 H_0 , 根据样本观测值计算统计量 F 给定显著水平 α , 查第一个自由度为 k , 第二个自由度为 $n-k-1$ 的 F 分布表得临界值 $F_\alpha(k, n-k-1)$ 。当 $F \geq F_\alpha(k, n-k-1)$ 时, 拒绝 H_0 , 则认为回归方程显著成立; 当 $F < F_\alpha(k, n-k-1)$ 时, 接受 H_0 , 则认为回归方程无显著意义。

4.3 参数显著性检验

回归方程显著成立, 并不意味着每个解释变量 X_1, X_2, \dots, X_k 对被解释变量 Y 的影响都是重要的。如果某个解释变量对被解释变量 Y 的影响不重要, 即可从回归模型中把它剔除掉, 重新建立回归方程, 以利于对经济问题的分析和对 Y 进行更准确的预测。为此需要对每个变量进行考查, 如果某个解释变量 X 对被解释变量 Y 的作用不显著, 那么它在多元线性回归模型中, 其前面的系数可取值为零。因此必须对 β_i 是否为零进行显著性检验。

由(3.23)式

$$S(\hat{\beta}_i) = \hat{\sigma}(\hat{\beta}_i) = \sqrt{C_{ii} S_e^2} \quad (4.8)$$

其中

$$S_e^2 = \frac{\mathbf{e}'\mathbf{e}}{n-k-1}$$

C_{ii} 为 L^{-1} 的第 i 个对角元素, 而 $L = \tilde{X}'\tilde{X}$, \tilde{X} 是中心化的数据阵。

对回归系数 $\hat{\beta}_i$ 进行显著性 t 检验, 步骤如下:

(1) 提出原假设 $H_0: \beta_i = 0$; 备择假设 $H_1: \beta_i \neq 0$ 。

(2) 构造统计量 $t = \frac{\hat{\beta}_i - \beta_i}{S(\hat{\beta}_i)}$, 当 $\beta_i = 0$ 成立时, 统计量 $t = \frac{\hat{\beta}_i}{S(\hat{\beta}_i)} \sim t(n-k-1)$ 。这里 $S(\hat{\beta}_i)$ 是 $\hat{\beta}_i$ 的标准差, k 为解释变量个数, 计算由式(4.8)给出。

(3) 给定显著性水平 α , 查自由度为 $n-k-1$ 的 t 分布表, 得临界值 $\frac{t_\alpha}{2}(n-k-1)$ 。

(4) 若 $|t| \geq \frac{t_\alpha}{2}(n-k-1)$, 则拒绝 $H_0: \beta_i = 0$, 接受 $H_1: \beta_i \neq 0$, 即认为 β_i 显著不为零。若

$|t| < \frac{t_\alpha}{2}(n-k-1)$, 则接受 $H_0: \beta_i = 0$, 即认为 β_i 显著为零。

5. 回归变量的选择与逐步回归

5.1 变量选择问题

在实际问题中, 影响因变量 Y 的因素 (自变量) 很多, 人们希望从中挑选出影响显著的自变量来建立回归关系式, 这就涉及到自变量选择的问题。

在回归方程中若漏掉对 Y 影响显著的自变量, 那么建立的回归式用于预测时将会产生较大的偏差。但回归式若包含的变量太多, 且其中有些对 Y 影响不大, 显然这样的回归式不仅使用不方便, 而且反而会影响预测的精度。因

而选择合适的变量用于建立一个"最优"的回归方程是十分重要的问题。

选择"最优"子集的变量筛选法包括逐步回归法(Stepwise),向前引入法(Forward)和向后剔除法(Backward)。

向前引入法是从回归方程仅包括常数项开始,把自变量逐个引入回归方程。具体地说,先在 m 个自变量中选择一个与因变量线性关系最密切的变量,记为 x_{i_1} ,然后在剩余的 $m-1$ 个自变量中,再选一个 x_{i_2} ,使得 $\{x_{i_1}, x_{i_2}\}$ 联合起来二元回归效果最好,第三步在剩下的 $m-2$ 个自变量中选择一个变量 x_{i_3} ,使得 $\{x_{i_1}, x_{i_2}, x_{i_3}\}$ 联合起来回归效果最好, ...如此下去,直至得到"最优"回归方程为止。

向前引入法中的终止条件为,给定显著性水平 α ,当某一个对将被引入变量的回归系数作显著性检查时,若 $p\text{-value} \geq \alpha$,则引入变量的过程结束,所得方程即为"最优"回归方程。

向前引入法有一个明显的缺点,就是由于各自变量可能存在着相互关系,因此后续变量的选入可能会使前面已选入的自变量变得不重要。这样最后得到的"最优"回归方程可包含一些对 Y 影响不大的自变量。

向后剔除法与向前引入法正好相反,首先将全部 m 个自变量引入回归方程,然后逐个剔除对因变量 Y 作用不显著的自变量。具体地说,从回归式 m 个自变量中选择一个对 Y 贡献最小的自变量,比如 x_{j_1} ,将它从回归方程中剔除;然后重新计算 Y 与剩下的 $m-1$ 个自变量回归方程,再剔除一个贡献最小的自变量,比如 x_{j_2} ,依次下去,直到得到"最优"回归方程为止。向后剔除法中终止条件与向前引入法类似。

向后剔除法的缺点在于,前面剔除的变量有可能因以后变量的剔除,变为相对重要的变量,这样最后得到的"最优"回归方程中有可能漏掉相对重要的变量。

逐步回归法是上述两个方法的综合。向前引入中被选入的变量,将一直保留在方程中。向后剔除法中被剔除的变量,将一直排除在外。这两种方程在某些情况下会得到不合理的结果。于是,可以考虑到,被选入的变量,当它的作用在新变量引入后变得微不足道时,可以将它删除;被剔除的变量,当它的作用在新变量引入情况下变得重要时,也可将它重新选入回归方程。这样一种以向前引入法为主,变量可进可出的筛选变量方法,称为逐步回归法。

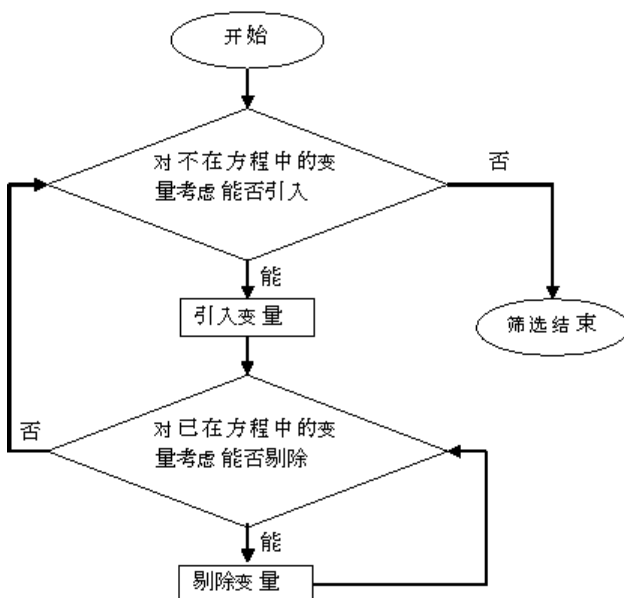
5.2 逐步回归分析

5.2.1 基本思想

逐个引入自变量。每次引入对 Y 影响最显著的自变量,并对方程中的老变量逐个进行检验,把变为不显著的变量逐个从方程中剔除掉,最终得到的方程中既不漏掉对 Y 影响显著的变量,又不包含对 Y 影响不显著的变量。

5.2.2 筛选的步骤

首先给出引入变量的显著性水平 α_m 和剔除变量的显著性水平 α_{out} ,然后按下图筛选变量。



5.2.3 逐步筛选法的基本步骤

逐步筛选变量的过程主要包括两个基本步骤：一是从回归方程中考虑剔除不显著变量的步骤；二是从不在方程中的变量考虑引入新变量的步骤。

(1) 考虑可否引入新变量的基本步骤。假设已入选 r 个变量，不在方程中的变量记为 $x_{j_1}, \dots, x_{j_{n-r}}$ 。

1. 计算不在方程中的变量 x_{j_i} 的偏回归平方和 P_{j_i} ：

$P_{j_i} = Q(i_1, \dots, i_r) - Q(i_1, \dots, i_r, j_i)$ ， Q 表示括号中这些变量的回归模型的残差平方和。并设

$P_{j_0} = \max(P_{j_1}, \dots, P_{j_{n-r}})$ ，即不在方程中的变量 x_{j_0} 是对 Y 影响最大的变量。

1. 检验变量 x_{j_0} 对 Y 的影响是否显著。对变量 x_{j_0} 作回归系数的显著性检验，即检验 $H_0: \beta_{j_0} = 0$ ，检验统计量为

$$F_{j_0} = \frac{P_{j_0}}{Q(i_1, \dots, i_r, j_0)/(n-r-2)} \text{ 及 } p = P\{F \geq F_{j_0}\}, \text{ 其中 } F \sim F(1, n-r-1).$$

若 $p < \alpha_m$ ，则引入 x_{j_0} ，并转入考虑可否剔除变量的步骤。若 $p \geq \alpha_m$ ，则逐步筛选变量的过程结束。

(2) 考虑可否剔除变量的基本步骤。假设已引入回归方程的变量为 $x_{i_1}, x_{i_2}, \dots, x_{i_r}$ 。

1. 计算已在方程中的变量 x_{i_k} 的偏回归平方和 P_{i_k} 。

$$\begin{aligned} P_{i_k} &= Q(i_1, \dots, i_{k-1}, i_{k+1}, \dots, i_r) - Q(i_1, \dots, i_r) \\ &= U(i_1, \dots, i_r) - U(i_1, \dots, i_{k-1}, i_{k+1}, \dots, i_r) \end{aligned}$$

其中 Q 表示括号中这些变量的回归模型的残差平方和， U 表示其回归平方和。设

$P_{i_0} = \min(P_{i_1}, \dots, P_{i_r})$ ，即相应的变量 x_{i_0} 是方程中对 Y 影响最小的变量。

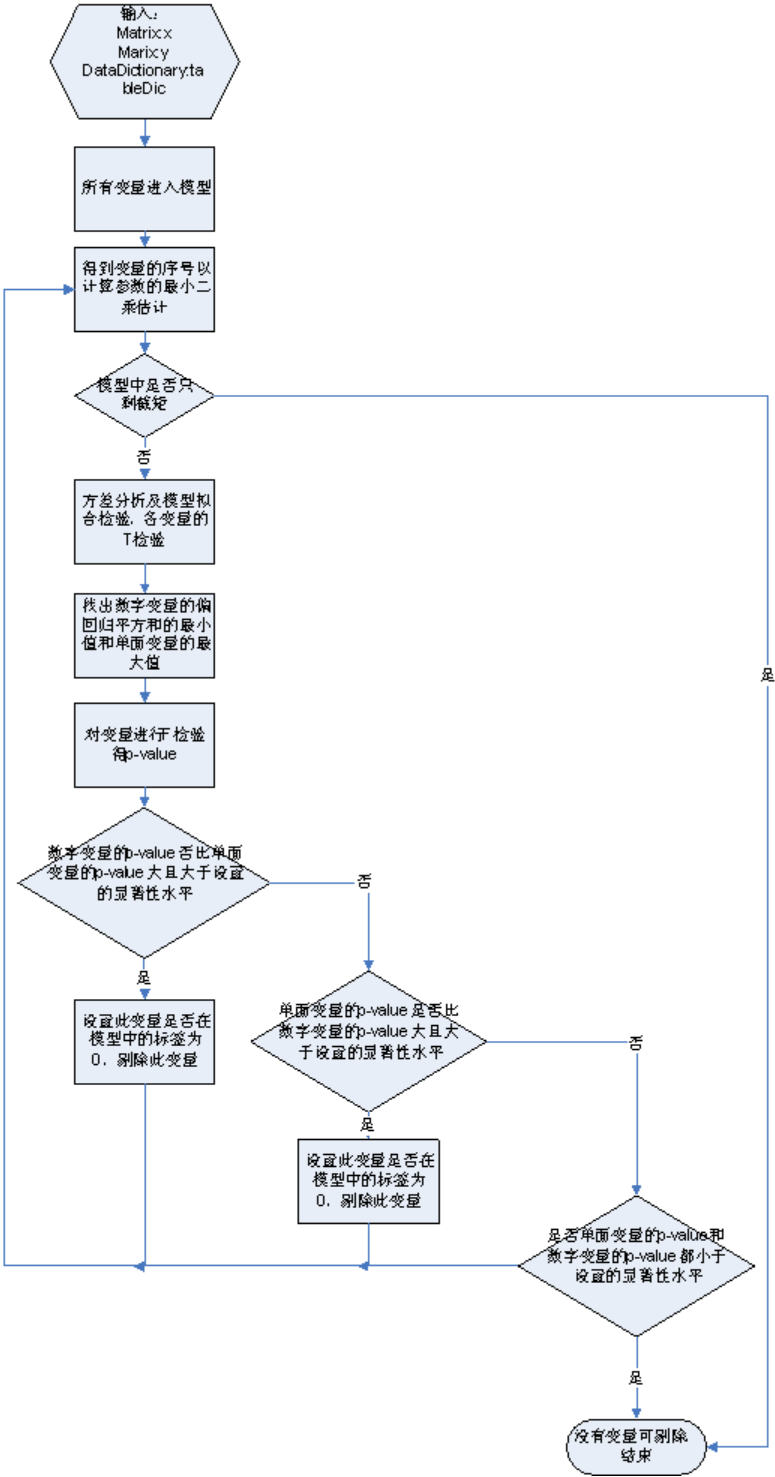
1. 检验 x_{i_0} 对 Y 的影响是否显著。对变量 x_{i_0} 进行回归系数的显著性检验，即检验 $H_0: \beta_{i_0} = 0$ ，检验统计量为

$$F_{i_0} = \frac{P_{i_0}}{Q(i_1, \dots, i_r)/(n-r-1)} \text{ 及 } p = P\{F \geq F_{i_0}\}, \text{ 其中 } F \sim F(1, n-r-1).$$

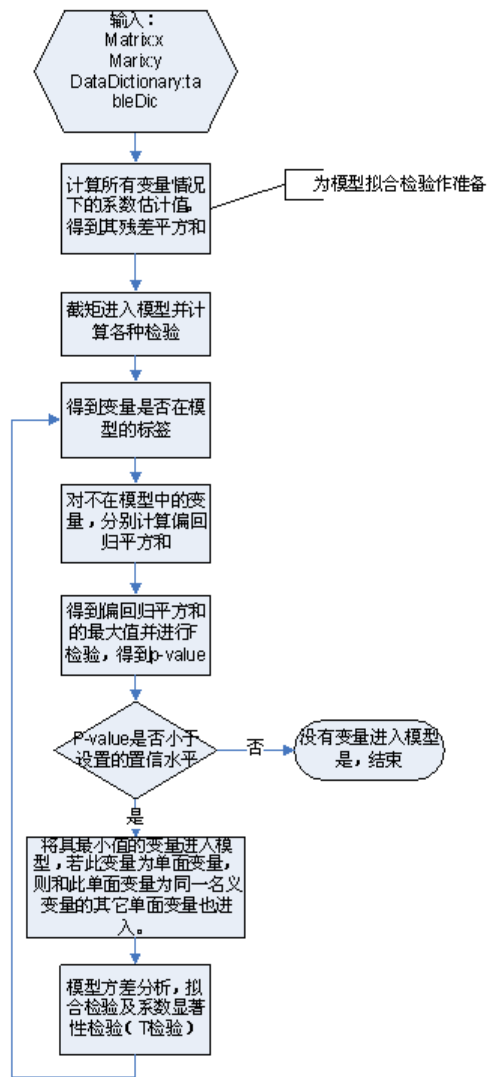
若 p 大于等于 α_{out} ，则剔除 x_{i_0} ，重新建立 Y 与其余 $r-1$ 个变量的回归方程，然后再检验方程中最不重要的变量可否删除，直到方程中没有变量可删除后，转入考虑能否引入新变量的步骤。

5.3 流程图

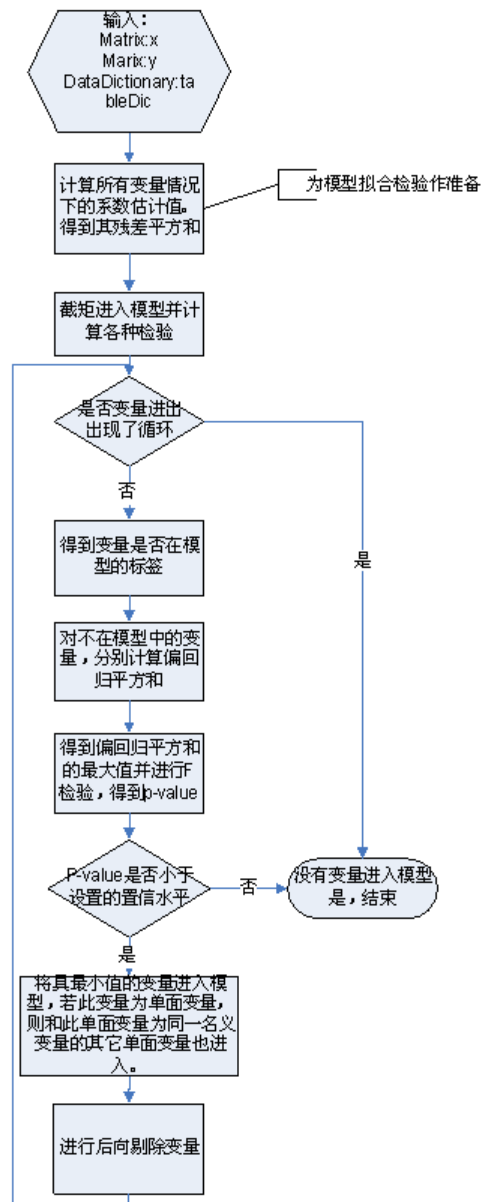
(1) 后向选择



(2) 前向引入 (Forward)



(3)逐步回归(Stepwise)



分类: [Algorithm for Data Mining](#)

好文要顶

关注我

收藏该文

[zgw21cn](#)

关注 - 1

粉丝 - 68

+ 加关注

« 上一篇: [逻辑回归Stepwise的R代码](#)

» 下一篇: [Sql server 行列互换制作交叉表格报表](#)

posted on 2008-12-24 13:14 [zgw21cn](#) 阅读(17073) 评论(7) 编辑 收藏

评论

#1楼 2009-01-07 09:43 [aliu](#)

请问怎么把多元化回归模型转化为中心化的模型？

在你上面的ols回归只是以二元回归为例，把二元回归模型改写成了中心化模型。

[支持\(0\)](#) [反对\(0\)](#)

#2楼 2009-01-07 17:49 [stentor](#)

@aliu

对观测数据进行中心化处理只是为了方便下面得出估计值B的具体的表达式而已。若需要对多元化的数据进行中心化处理，只需将各变量减去其平均值即可。

[支持\(0\)](#) [反对\(0\)](#)

#3楼 2009-04-08 16:30 [lz](#)

帅哥，有没有word版本，传一个给我，参考下，谢谢了
邮箱：leezhiwei@126.com

支持(0) 反对(0)

#4楼 2015-11-22 20:03 丁丁的大森林

有没有word哦。。。

支持(0) 反对(0)

#5楼 2016-03-01 22:39 Lansing999

数学太好了！！我看的一头雾水。。。

支持(0) 反对(0)

#6楼 2016-10-21 14:38 AutumnLight

质量很高，很多没看懂。。。

支持(0) 反对(0)

#7楼 2016-10-29 15:19 知好奇行

学神！

支持(0) 反对(0)

[刷新评论](#) [刷新页面](#) [返回顶部](#)

注册用户登录后才能发表评论，请 [登录](#) 或 [注册](#)， [访问](#) 网站首页。

- 【推荐】超50万行VC++源码：大型组态工控、电力仿真CAD与GIS源码库
- 【活动】腾讯云服务器推出云产品采购季 1核2G首年仅需99元
- 【推荐】精品问答：精品问答：Python 技术 1000 问
- 【推荐】免费下载《阿里工程师的自我修养》