

DOI: 10.3969/j.issn.1001-8972.2010.09.020

多元线性回归分析及其应用

林彬 湛江师范学院信科院 524048

Multiple linear regression analysis and its application

Lin Bin (Zhanjiang Normal College, Zhanjiang, 524048, China)

摘要

本文研究了多元线性回归理论及应用,探讨了多元线性回归模型中未知参数的估计及其参数的检验问题,以实例进行了验证。

关键词

多元线性回归分析; 回归模型; 检验问题。

Abstract

We study the multiple linear regression theory and application, and estimation of unknown parameters and parameter testing problem in multivariate linear regression model, and then an example is solved to verify its effectiveness.

Key words

Multiple linear regression analysis; regression model; parameter testing problem

1 概述

回归分析是一种传统的应用性较强的科学方法,是现代应用统计学的一个重要的分支,在各个科学领域都得到了广泛的应用。它不仅能够把隐藏在大规模原始数据群体中的重要信息提炼出来,把握住数据群体的主要特征,从而得到变量间相关关系的数学表达式,利用概率统计知识对此关系进行分析,以判别其有效性,还可以利用关系式,由一个或多个变量值去预测和控制另一个因变量的取值,从而知道这种预测和控制达到的程度,并进行因素分析。

2 多元线性回归数学模型

设可预测的随机变量为 y , 它受到 p 个非随机因素 $x_1, x_2, \dots, x_{p-1}, x_p$, 和不可预测的随机因素 ε 的影响。多元线性回归数学模型为

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_{p-1} x_{p-1} + \beta_p x_p + \varepsilon \quad (1)$$

$$\varepsilon \sim N(0, \sigma^2)$$

其中 $\beta_0, \beta_1, \dots, \beta_p$ 为回归系数

对 y 和 $x_1, x_2, \dots, x_{p-1}, x_p$ 分别进行 n 次独立观测, 取得 n 组数据(样本)

$$y_i, x_{i1}, x_{i2}, \dots, x_{ip-1} \quad (i = 1, 2, 3, \dots, n)$$

则有

$$\begin{cases} y_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_{p-1} x_{1p-1} + \varepsilon_1 \\ y_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_{p-1} x_{2p-1} + \varepsilon_2 \\ \vdots \\ y_n = \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_{p-1} x_{np-1} + \varepsilon_n \end{cases} \quad (2)$$

其中 $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ 相互独立, 且服从 $N(0, \sigma^2)$ 分布。

令

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np-1} \end{bmatrix}$$

则式(2)用矩阵形式表示为

$$Y = X\beta + \varepsilon \quad (3)$$

$$\varepsilon \sim N(0, \sigma^2 I_n)$$

3 模型参数 的最小二乘法估计与误差方差 σ^2 的估计

的最小二乘法估计即选择 $\hat{\beta}$ 使误差项的平方和为最小值, 这时 $\hat{\beta}$ 的值

的点估计。

$$S(\beta) = \varepsilon^T \varepsilon = (y - x\beta)^T (y - x\beta) \quad (4)$$

为了求 β ，由(4)式将 $S(\beta)$ 对 β 求导，并令其为零，得

$$\frac{dS(\beta)}{d\beta} = \frac{d[(y - x\beta)^T (y - x\beta)]}{d\beta} = \frac{d[(y^T y - \beta^T x^T y - y^T x\beta + \beta^T x^T x\beta)]}{d\beta} = 0 \quad (5)$$

由(5)式可解出 $\hat{\beta}$

$$\hat{\beta} = (x^T x)^{-1} (x^T y) \quad (6)$$

对残差向量 ε

$$\hat{\varepsilon} = y - \hat{y} = y - x\hat{\beta} = [I - x(x^T x)^{-1} x^T] y \quad (7)$$

则残差平方和

$$\hat{\varepsilon}^T \hat{\varepsilon} = \hat{\varepsilon}^T [I - x(x^T x)^{-1} x^T] y = y^T y - \hat{\beta}^T x^T y \quad (8)$$

又因为 $E(y) = x\beta$ ，因此

$$E(\hat{\varepsilon}^T \hat{\varepsilon}) = \sigma^2 (n - p)$$

$$\sigma^2 = \frac{1}{n - p} (\hat{\varepsilon}^T \hat{\varepsilon}) \quad (9)$$

4 模型检验

多元线性回归数学模型建立后，是否与实际数据有较好的拟合度，其模型线性关系的显著性如何等，还需通过数理统计进行检验。常用的统计检验有 R 检验和 F 检验。

4.1 R 检验

$$R = \sqrt{1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (10)$$

R 是复相关系数，用于测定回归模型的拟合优度，R 越大，说明 Y 与 x_1, x_2, \dots, x_{p-1} 的线性关系越显著。 \bar{y} 为 y_i 的平均值，R 取值范围为 $0 < R \leq 1$ 。

4.2 F 检验

$$F = \frac{U/m}{Q/(n-m-1)} \sim F(m, n-m-1) \quad (11)$$

式中 $Q = \sum (y_i - \bar{y})^2$ ， $U = \sum (\hat{y}_i - \bar{y})^2$ ，

m 为自变量个数，n 为数据个数。

F 服从 $F(m, n-m-1)$ 分布，取显著

性水平为 α ，如果 $F > F_{\alpha}(m, n-m-1)$ ，表明回归模型显著，可从用于预测。反之，回归模型不能用于预测。

5 应用实例

某医院为了解病人对医院工作的满意程度 Y 和病人的年龄 X_1 、病情的严重程度 X_2 和病人的忧虑程度 X_3 之间的关系，随机调查了该医院的 10 位病人，得数据如表 1 所示。

表 1 病人满意度的调查数据

年龄 (x_1)	病情程度 (x_2)	忧虑程度 (x_3)	满意程度 (y)
50	51	2.3	48
36	46	2.3	57
40	48	2.2	66
41	44	1.8	70
28	43	1.8	89
49	54	2.9	36
42	50	2.2	46
45	48	2.4	54
52	62	2.9	26
29	50	2.1	77

使用 MATLAB 语言编程并计算得下面结果：

```
RegCoff=
175.5249
-1.1713
-0.5117
-19.6453
R=0.9603
F=23.7098
FX=9.0886 0.4105 2.5260
TX=9.0224 0.8376 79.7754
```

从结果可以得出，回归模型为

$$y = 175.5249 - 1.1713x_1 - 0.5117x_2 - 19.6453x_3$$

取 $\alpha = 0.05$ 对方程和回归系数进行检验。查 F 分布表可得 $F_{0.05}(3, 6) = 4.76$ ， $F_{0.05}(1, 6) = 5.99$

本例中的方程检验值 $F = 23.7098 > 4.76$ ，说明模型的回归效果高度显著。

$F_1 = 9.0886 > 5.99$ ，说明 x_1 显著。

$F_2 = 0.4105 < 5.99$ ，说明 x_2 很不显著。

$F_3 = 2.5260 < 5.99$ ，说明 x_3 不显著。

R 为 0.9603 接近 1，表明线性相关性较强。

在实际中，由于 Y 的影响因素还有很多，使 Y 与 X 关系更为复杂，而且记录数据的准确性、可靠性、异常数据等问题，将影响 Y 的预测分析。

参考文献

- [1] 梅长林，范金城. 数据分析方法. 高等教育出版社. 2006.2.
- [2] 何晓群. 现代统计分析与应用. 中国人民大学出版社. 2007.8.

作者简介

林彬(1962-)，男，湛江人，博士，讲师，研究方向：应用数学。