

回归饱和设计值得商榷的问题

肖俊璋¹, 肖俊光²

(1 西北农林科技大学资源环境学院, 陕西杨陵 712100; 2 陕西三原县职教中心, 陕西三原 713800)

摘要: 从理论和实践上论证了回归饱和设计不能用最小二乘法估计回归模型中的参数; 同时也论证了回归饱和设计的试验方案, 不能用 D—优良性判断其优劣, 回归饱和 D—最优设计是不存在的。

关键词: 回归饱和设计; 最小二乘法; D—优良性

中图分类号: S11⁺⁵ 文献标识码: A 文章编号: 1008-505X(2011)05-1274-04

Discussions on some problems for the regression saturation design

XIAO Jun-zhang¹, XIAO Jun-guang²

(1 College of Resource and Environment, Northwest A & F University, Yangling, Shaanxi 712100, China;

2 Vocational Education Centre of Sanyuan County, Sanyuan, Shaanxi 713800, China)

Abstract: It has been proved by theory and practice that the parameters of regression model can not be estimated by the least square method in the design of regression saturation and that the experiment plan designed by regression saturation can not be evaluated by the D-optimality. As a result, the D-optimal design of regression saturation is not existed.

Key words: least square method; design of regression saturation; D-optimality

回归饱和设计就是试验的处理数等于所要确定的未知参数的个数, 也就是处理数等于所要建立的回归方程的回归系数的个数。由于没有剩余自由度, 无法检验回归方程的显著性, 所以过去认为回归饱和设计不是好的设计。自从回归饱和 D—最优设计出现以后, 一般认为其处理数少, 信息量大, 参数估计精度高, 因而颇受试验者的欢迎, 回归饱和设计的身价也大为提高。但是通过我们的实践和研究, 回归饱和设计有以下问题值得商榷:

1 关于回归模型中参数的估计问题

在回归分析中, 要达到对回归模型中参数的

无偏估计, 一般都用最小二乘法。对于回归饱和设计, 是否能用最小二乘法估计其模型中的参数, 过去是无人怀疑的。但是通过对以下试验结果的分析, 可以看出, 在回归饱和设计中, 用最小二乘法估计其模型中的参数是有问题的。在研究小麦氮、磷用量及配比试验中, 我们采用了两种试验方案, 方案一为 Box 1971 年提出的二因子二次回归饱和 D—最优设计^[1]; 方案二为一般的二因子二次回归饱和设计, 具体的试验方案及试验结果见表 1 及表 2。

根据表 1、表 2 数据, 用最小二乘法估计回归模型中的参数。经计算表 1 得编码值回归方程(1), 表 2 得编码值回归方程(2)。

$$\hat{y}_1 = 2596.2435 + 520.299x_1 + 759.549x_2 - 168.198x_1^2 - 453.477x_2^2 + 740.799x_1x_2 \quad (1)$$

$$\hat{y}_2 = 2800.767 + 498.158x_1 + 779.141x_2 - 411.296x_1^2 - 791.283x_2^2 + 1098.134x_1x_2 \quad (2)$$

表1 方案一及试验结果(回归饱和 D-最优设计)
Table 1 Experimental plan one and result (D-optimal design of regression saturation)

处理号 Treatment No.	$x_1(Z_1)$	$x_2(Z_2)$	$y(\text{kg}/\text{hm}^2)$	$\hat{y}(\text{kg}/\text{hm}^2)$
1	-1(0)	-1(0)	1435.5	1435.5
2	1(187.5)	-1(0)	994.5	994.5
3	-1(0)	1(150)	1473.0	1473.0
4	-0.1315(81.42)	-0.1315(65.15)	2430.0	2430.0
5	1(187.5)	0.3944(140.58)	3469.5	3469.5
6	0.3944(130.73)	1(150)	3373.5	3373.5

注(Note): Z_1 —氮 N kg/hm²; Z_2 —磷 P₂O₅ kg/hm²

表2 方案二及试验结果
Table 2 Experimental plan two and result

处理号 Treatment No.	$x_1(Z_1)$	$x_2(Z_2)$	$y(\text{kg}/\text{hm}^2)$	$\hat{y}(\text{kg}/\text{hm}^2)$
1	-1(0)	-1(0)	1419.0	1419.0
2	0.25(117.19)	1(150)	3162.0	3162.0
3	1(187.5)	0.25(93.75)	3307.5	3307.5
4	-0.5(46.88)	-0.5(37.5)	2136.0	2136.0
5	1(187.5)	0.75(131.25)	3850.5	3850.5
6	0.75(164.06)	1(150)	3754.5	3754.5

注(Note): Z_1 —氮 N kg/hm²; Z_2 —磷 P₂O₅ kg/hm²

将方案一和方案二中每个处理各因子的编码值代入相应的回归方程,求出每个处理的回归值,见表1、表2 的 \hat{y} 列。

从表1、表2 可以看出,方案一与方案二的回归值 \hat{y} 与观察值(y)完全吻合。这里就产生一个问题,回归饱和 D-最优设计与一般的回归饱和设计的效果为什么完全相同? 经过研究,原因在于对回归模型中参数的估计上。对于一次或者二次回归饱和设计,用最小二乘法估计其模型中的参数时其矩阵形式为:

$$Ab = B \quad (3)$$

$$\text{或}, (X'X)b = X'Y \quad (4)$$

$$\text{所以}, b = A^{-1}B = (X'X)^{-1}X'Y \quad (5)$$

其中: b 在 p 元一次回归中为 $b' = (b_1, b_2, \dots, b_p)$, 在 P 元二次回归中为 $b' = (b_1, b_2, \dots, b_p, b_{11}, b_{22})$,

$$\hat{y}_1 = 2596.2357 + 520.29x_1 + 759.54x_2 - 168.218x_1^2 - 453.48x_2^2 + 740.79x_1x_2 \quad (8)$$

$$\hat{y}_2 = 2800.772 + 498.158x_1 + 779.159x_2 - 411.299x_1^2 - 791.298x_2^2 + 1098.141x_1x_2 \quad (9)$$

可以看出,回归方程(8)、(9)与回归方程(1)、(2)是一样的,回归系数很小的差异是由舍入误差

$\dots, b_{pp}, b_{12}, b_{13}, \dots, b_{p-1,p}$; X 为结构矩阵; A 为正规方程组的系数矩阵, $A = X'X$; B 为正规方程组的常数项矩阵, $B = X'Y$; $Y' = (y_1, y_2, \dots, y_N)$, N = 试验处理数。

在回归饱和设计中,结构矩阵 X 为 N 阶方阵,当 X 为非奇异矩阵时,则 X 有逆矩阵 X^{-1} , X' 有逆矩阵 $(X')^{-1}$ 。对于式(4)的等号两边各乘 $(X')^{-1}$ 得:

$$Xb = Y \quad (6)$$

$$\text{则}, b = X^{-1}Y \quad (7)$$

显然,在结构矩阵 X 为非奇异矩阵的回归饱和设计中,式(6)与式(3)等价,式(7)与式(5)等价。

回归方程(1)、(2)中的回归系数是根据方案一、二及其试验结果用式(5)计算出的,现用式(7)进行计算,得:

$$\hat{y}_1 = 2596.2357 + 520.29x_1 + 759.54x_2 - 168.218x_1^2 - 453.48x_2^2 + 740.79x_1x_2 \quad (8)$$

$$\hat{y}_2 = 2800.772 + 498.158x_1 + 779.159x_2 - 411.299x_1^2 - 791.298x_2^2 + 1098.141x_1x_2 \quad (9)$$

造成的,用式(8)、(9)求出每个处理的回归值(\hat{y})与其观察值(y)也是完全吻合的。

现在对式(7)作进一步分析,式(7)是由式(6)来的,应从式(6)谈起。式(6)就是把观察值(y_α)根

$$y_\alpha = \beta_0 + \beta_1 x_{\alpha 1} + \beta_2 x_{\alpha 2} \beta_{11} x_{\alpha 1}^2 + \beta_{22} x_{\alpha 2}^2 + \beta_{12} x_{\alpha 1} x_{\alpha 2} + \varepsilon_\alpha \quad (10)$$

式中: $\alpha = 1, 2, \dots, N$ 为处理号; y_α 为第 α 处理的观察值; $\beta_1, \beta_2, \dots, \beta_{12}$ 为6个待定参数; ε_α 是第 α 处

$$\begin{cases} y_1 = b_0 - b_1 - b_2 + b_{11} + b_{22} + b_{12} \\ y_2 = b_0 + b_1 - b_2 + b_{11} + b_{22} - b_{12} \\ y_3 = b_0 - b_1 + b_2 + b_{11} + b_{22} - b_{12} \\ y_4 = b_0 - 0.1315b_1 - 0.1315b_2 + 0.017292b_{11} + 0.017292b_{22} + 0.017292b_{12} \\ y_5 = b_0 + b_1 + 0.3944b_2 + b_{11} + 0.155551b_{22} + 0.3944b_{12} \\ y_6 = b_0 + 0.3944b_1 + b_2 + 0.155551b_{11} + b_{22} + 0.3944b_{12} \end{cases} \quad (11)$$

将试验结果 $y_1 = 1435.5, y_2 = 994.5, y_3 = 1473.0, y_4 = 2430.0, y_5 = 3469.5, y_6 = 3373.5$ 代入联立方程组(11),用式(7)解出式(11)中的6个回归系数,就得回归方程(8)。显然,这样所配置的回归方程,其回归值(\hat{y}_α)必然与观察值(y_α)完全吻合,这就回答了前面所提出的问题。

把联立方程组(11)与回归模型(10)比较,联立方程组(11)的每个方程比回归模型(10)少误差项 ε_α ,这就是说,在联立方程组(11)中把观察值 y_α 作为没有试验误差的一般变量,但实际上 y_α 是带有试验误差 ε_α 的随机变量。因此所求出的回归系数可靠性较差。

综上所述,在回归饱和设计中,要估计回归模型中的参数,最小二乘法是用不上的,式(3)、(5)是用最小二乘法建立的,但它们与式(6)、(7)等价,这就是说,表面上是用最小二乘法估计回归系数,实际上等于把观察值(y_α)当作没有试验误差的一般变量,根据模型代入试验方案中的每个处理组成联立方程组求解回归系数,建立回归方程。这样所配置的回归方程可靠性较差。

2 用D-优良性评价回归饱和设计的问题

对于给定的回归模型和给定的因子空间某一区域上,可以编制多种试验方案,对于这些试验方案的好坏可以从不同角度进行评价。D-优良性是从对回归模型中的参数的估计好坏来评价的。对于回归饱和设计的试验方案用D-优良性评价其优劣,长期以来被人们认为是正确的,所以出现了不少的回归饱和D-最优设计的试验方案,并且长期被应用。

据回归模型代入方案中的每个处理组成求解回归系数的联立方程组。以二元二次回归为例,其模型为:

理的试验误差。如果按方案一进行试验,设 $\beta_0, \beta_1, \dots, \beta_{12}$ 的估计值为 b_0, b_1, \dots, b_{12} ,则式(6)为:

但是通过以下分析,完全证明回归饱和设计的试验方案不能用D-优良性来评价。

反映试验方案D-优良性好坏的数量指标是回归系数的密集椭球体的体积(该体积的平方称为广义方差),该体积小,则回归系数对回归模型中参数的估计精度高,即D-优良性好,反之则反。回归模型中的参数在不同的试验方案中,可以得到不同的最小二乘估计,也可得到不同的回归系数的密集椭球体的体积。在给定的回归模型下和给定的因子空间某一区域上,可以编制多种试验方案,回归系数的密集椭球体的体积愈小,试验方案的D-优良性愈好。如果在给定的因子空间某一区域上,能使回归系数的密集椭球体的体积达到最小的那个试验方案,称为该因子区域上的D-最优方案。回归系数的密集椭球体的体积计算公式如下^[1]:

$$V(\varepsilon) = \frac{(m+2)^{m/2} \pi^{m/2}}{\Gamma\left(\frac{m}{2} + 1\right) \sqrt{|A(\varepsilon)|}} \quad (12)$$

式中 $V(\varepsilon)$ 是试验方案 ε 的密集椭球体体积; m 是回归系数的个数; $\Gamma(x)$ 是 Γ 函数; $|A(\varepsilon)|$ 是试验方案 ε 的信息矩阵 $A(\varepsilon)$ 的行列式。

由式(12)可知,当回归模型确定后, m 也就确定了,密集椭球体的体积 $V(\varepsilon)$ 就决定于相应的信息矩阵行列式 $|A(\varepsilon)|$ 的值。密集椭球体的体积 $V(\varepsilon)$ 与相应的信息矩阵行列式 $|A(\varepsilon)|$ 的平方根成反比,因此, $|A(\varepsilon_1)| > |A(\varepsilon_2)|$ 与 $V(\varepsilon_1) < V(\varepsilon_2)$ 等价,这就是说,信息矩阵行列式 $|A(\varepsilon)|$ 值的大小,反映了试验方案D-优良性的好坏。

信息矩阵 $A(\varepsilon)$ 就是正规方程组的系数矩阵,它是在用最小二乘法估计回归模型中的参数时建立

的。如果一个回归试验方案,通过试验后,不能用最小二乘法估计回归模型中的参数,当然也就不能建立信息矩阵,所以该试验方案就不能用D-优良性判断其优劣。从第一个问题的讨论中可知,回归饱和设计是不能用最小二乘法估计回归模型中的参数,因而回归饱和设计的试验方案不能用D-优良性判断其优劣。当然,回归饱和D-最优设计是不存在的。回归饱和设计虽然处理数少,但从对回归模型中的参数的估计来看是不好的设计。

3 讨论

1) 回归饱和D-最优设计有一次回归饱和D-最优设计和二次回归饱和D-最优设计。由于因子数 $P \geq 4$ 的二次回归饱和D-最优方案未曾找到,所以又根据 $P=2,3$ 的二次回归饱和D-最优方案的谱点结构,提出了 $P \geq 4$ 较好的二次回归饱和设计的方案^[1]。所有这些设计从未有人提出怀疑和意见,一般认为该设计处理数少,参数估计精度高,因而颇受欢迎,有一段时间许多书刊都在介绍这种设计,各地也纷纷举办学习班学习这种设计,这种设计也被写入高等院校本科生和研究生的教材^[2-3],因而该设计在土壤肥料试验中被广泛应用。为什么过去没有人对这种设计提出怀疑和意见呢?很大可能是忽略了回归饱和设计不能用最小乘法估计回归模型中参数这一事实。目前还有人坚持原来的看法,还是用最小二乘法估计回归模型中参数时的信

息矩阵A的行列式 $|A|$ 值的大小来评价回归饱和设计的优劣就是例证。

2) 前面的论证只是证明回归饱和设计不能用D-优良性评价其试验方案的优劣,当然就没有D-最优方案。回归饱和设计是个不好的设计,它没有剩余自由度,不能检验回归方程的显著性;它对回归模型中参数的估计,精确度低。回归饱和设计适于探索性试验,由于条件限制不得已采用了这种设计,其结果可作进一步试验的参考。如何提高回归饱和设计对回归模型中参数估计的精度,就是设置重复,随着重复次数的增加,其精度不断增强。

参 考 文 献:

- [1] 茅诗松,丁元,贾纪芗,等. 回归分析及其试验设计[M]. 上海:华东师范大学出版社,1981. 259-268,299-302.
Mao S S, Ding Y, Jia J X et al. Regression analysis and experiment design [M]. Shanghai: East China Normal University Press, 1981. 259-268, 299-302.
- [2] 毛达如,申建波. 植物营养研究方法[M]. 北京:中国农业大学出版社, 2005. 266-269.
Mao D R, Shen J B. Plant nutrition research methods [M]. Beijing: China Agricultural University Press, 2005. 266-296.
- [3] 王兴仁,张福锁,等. 现代肥料试验设计[M]. 中国农业出版社,1996. 81-90, 103-105.
Wang X R, Zhang F S et al. Modern fertilizer experimental design [M]. Beijing: China Agricultural Press, 1996. 91-90.