第四期

Journal of Southwest Agricultural University

Nov. 1986

# 重复的二次回归正交旋转设计试验 数 学 模 型 的 选 择

#### 农学系 罗泽伟

【提要】由赤池统计模型推断的信息量准则(AIC)构造出作为二次回归正交 旋转设计试验指导的回归模型的AIC统计量,其形式为.

$$AIC = N^{1}n\delta^{2} + \frac{1}{\delta^{2}}(y - x\hat{\beta})^{T}(y - x\hat{\beta}) + 2K$$

利用以上准则评判具有重复的二次回归正交旋转设计试验多个数学模型的选择。

多因子二次回归正交旋转设计是将试验设计、数学模型的建立和模型的稳定性与精确性综合考虑的多因子试验设计<sup>[1]</sup>。 近年来,我国农业科技工作者又将其应用于农作物栽培规范化农艺措施的研究中,已获得很大的社会和经济效益。对于大面积大幅度提高作物产量起到巨大的推动作用。

二次回归设计(具有正交性或旋转性是将试验点均匀地分别在P维因子空间中半径分别为  $\rho=0$ 、 $\rho=+1$ 、 $\rho=+\gamma$ 的三个球面上的,因而将试验点分为三类。分布在原点的零水平试验  $点m_o$ 个;分布在单位球面上的试验点 $m_c$ 个;分布在星号位球面上的试验点 $m_c$ 个,总的试验 次数为 $N=m_o+m_c+m_r$ 。试验不设置重复。

农业田间试验由于外因素的干扰,误差较大。根据同一试验不同重复(地点、时间)所获得的数据建立的数学模型往往有较大的差异。因此,为了提高试验的灵敏度、降低试验误差对模型的干扰、确保所得的概括大范围作物栽培生产过程的数学模型的代表性,试验需要重复。本文应用赤池信息量准则构造出二次回归正交旋转设计不同重复试验数学模型的选择标准。

## 一、基本原理

P因子二次回归正交旋转设计试验指导的数学模型其一般形式为:

$$Y = \beta_0 + \sum_{i=1}^{p} \beta_i Z_i + \sum_{1 \le i \le j \le p} \beta_{ij} Z_i Z_j + \sum_{j=1}^{p} \beta_{ij} Z_j^2 \dots$$

若记,

$$Z = (1 Z_1 \cdots Z_p Z_1 Z_2 \cdots Z_{p-1} Z_p Z_1^2 \cdots Z_p)^T$$

$$\beta = \beta_0 \beta_1 \cdots \beta_p \beta_{12} \cdots \beta_{p-1p} \beta_{11} \cdots \beta_{pp})^T$$

则①式可以写成向量形式: Y=Z<sup>T</sup>β ·······②

试验设计的结构矩阵及目标观察值矩阵见表1。

表1 P因子二次回归正交旋转设计的结构矩阵

	试		4							-		重		复	
	验号	Z <sub>0</sub>	$z_1$	$z_2 \cdots z_p$	$z_1 z_1$	2 Z <sub>1</sub>	z <sub>3</sub> … z	P - 1	$z_p z_1^2$	Z	<b>z.···z</b> <sub>p.</sub> "	I	I	r	,
	1	1	1	1 1	1	1		1	1.	1.	1*	y <sub>11</sub>	y <sub>12</sub>	y <sub>1 }</sub>	
m ç	2	1	1	1 · · · -1	1	1	***	-1	1 *	1 •	···1*	у21	у 2 2	y <sub>2</sub> r	
	:	:	:	: :	:	:		:	:	:	:		:	:	
	2 P	1	-1	-1 ···-1	1	1	•••	1	1.	1 *	···1*				
	2 <sup>p</sup> +1	1	·	0 0	n	0	•••	0	r• 2	0.	•				
	}	_							_						
m	2 <sup>P</sup> +2	1		-	0	0	•••	0	r · 3	0*	0•		į	:	
m,	:	:	:	: :	:	:		:	:	:	:				
	$2^{P}+2p$	1	0	0 ···-r	0	0	•••	0	0.	0*	0*				
	$\int 2^{p} + 2p + 1$	1	0	0 0	0	0	•••	0	0*	0.	0*				
m,	١	:	•	: :	:	:		:	:	:	:	:	;	:	
	N	1	0	0 0	0	0	•••	0	0.	0.	···0•	y <sub>N 1</sub>	<b>У</b> N 2	y <sub>N</sub> ,	

其中。①N= $m_c+m_r+m_o$ ,②r是设计的臂长参数,③"\*"表示经中心化处理的数据。

则, ②中未知参数向量β的极大似然估计为:

$$\hat{\beta} = (X^TX)^{-1}X^TY \dots 3$$

于是模型②的理论预测值向量为:

假定随机向量:  $\epsilon \stackrel{\wedge}{=} Y - \stackrel{\wedge}{Y} = X\beta - X\stackrel{\wedge}{\beta} \sim N(0, \sigma^2 I_N)$ 

根据赤池信息准则,模型的信息量推断统计量是:

AIC=
$$-2\ln L(\hat{\beta}, \sigma^2, Z) + 2K$$
 ......

其中,  $L(\hat{\beta}, \sigma^2, Z)$  是模型的极大似然函数, K 是模型中的有效参数个数。

在此,由于 
$$\varepsilon = Y - X\hat{\beta} \sim N(0, \sigma^2 I_N)$$
则,  $L(\hat{\beta}, \sigma^2; Z) = \frac{N}{11} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(\overline{Y} - X\hat{\beta})^T(\overline{Y} - X\hat{\beta})\right\}$ 

$$L_nL(\hat{\beta}, \sigma^2; Z) = \sum_{i=1}^{N} \ln \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(Y - X\hat{\beta})^T(Y - X\hat{\beta})\right\}$$

$$= -\frac{N}{2} \ln(2\pi) - \left\{ \frac{N}{2} \ln\sigma^2 + \frac{1}{2\sigma^2} \left( \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} \right)^T \left( -\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} \right) \right\}$$

在上式中,由于 $-\frac{N}{2}\ln\left(2\pi\right)$ 是一个常数因子,可将其略去得到。

$$\ln L(\hat{\beta}, \sigma^2; Z) = -\frac{N}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

将上式代入④,得到:

AIC=Nln
$$\sigma^2 + \frac{1}{\sigma^2} (Y - X \hat{\beta})^T (Y - X \hat{\beta})^T + 2K$$
 .....

⑤式的含意是,它可以作为某一个统计模型,其建立依据的数据资料所包含的信息量以及模型中的估计参数对这些信息的概括能力和稳定性能的模型推断综合指标。利用AIC对多个模型进行选择时,是遵循最小值原理,即是相应的AIC为最小的模型被认为是可取的。

在表 1 中,利用不同重复的目标观察值为依据可以获得不同的数学模型。例如第 <sup>t</sup> 个 重 复相应的模型为:

$$\hat{Y}t = b_0^{(i)} + \sum_{i=1}^{p} b_i^{(i)} Z_i + \sum_{1 \le i \le j \le p} b_{ij}^{(i)} Z_i Z_j + \sum_{j=1}^{p} b_{jj}^{(i)} Z_j^2$$

相应的模型推断统计量为: AIC, (t=1, 2, ..., r)

则AICio对应的模型:

$$\hat{Y}_{i0} = b^{(i0)} + \sum_{i=1}^{\rho} b_{i}^{(i0)} Z_{i} + \sum_{1 \le i \le j \le p} b_{ij}^{(j0)} Z_{i} Z_{j} + \sum_{j=1}^{\rho} b_{jj}^{(j0)} Z_{j}^{2}$$

是在上述意义下的最佳模型。

表2 四次重复试验各回归模型 的回归系数表

1	四四万水数水												
重 复	1	2	3	4									
b <sub>o</sub>	121.99	120.77	120.83	120.69									
$b_1$	-4.42	-4.04	-3.29	-3.29									
b <sub>2</sub>	0.50	0.71	1.38	1.21									
b <sub>3</sub>	-0.33	-0.13	-0.29	-0.38									
b <sub>4</sub>	0.50	-0.13	0.38	0.21									
b <sub>5</sub>	0.42	0.04	0.04	0.54									
b <sub>12</sub>	0.27	1.22	0.63	0.84									
b <sub>1 3</sub>	0.02	-0.66	-0.36	-0.41									
b <sub>1 4</sub>	0.27	-0.28	-0.48	-0.28									
b <sub>15</sub>	0.02	-0.28	-0.49	-0.28									
b 2 3	-0.11	-0.28	-0.61	-0.78									
b 2 4	0.63	0.31	0.19	0.56									
b 2 5	-0.38	-0.06	0.06	0.31									
b 3 4	-0.63	-0.06	-0.19	-0.06									
b 3 5	-0.63	-0.06	0.19	-0.06									
b <sub>45</sub>	0.50	-0.06	0.31	-0.19									
b <sub>1 1</sub>	0.75	-0.06	0.31	0.19									
b 2 2	0.50	-0.06	0.19	-0.06									
b 3 3	-0.75	0.06	0.44	-0.06									
b44	-0.50	0.06	-0.19	-0.06									
b 5 5	-0.25	0.06	-0.19	0.06									

### 二、应用实例

本段引用"重庆市杂交水稻规范化栽培技术研究"课题的五因子、播种期( $Z_1$ )、移栽叶龄( $Z_2$ )、栽植密度( $Z_3$ )、施氮水平( $Z_4$ )、施磷钾水平( $Z_3$ )、 全实施的二次回归正交旋转设计重复 4 次试验的播种—齐穗天数数据资料为例,引证本文介绍的方法。

根据Robert R. Sokal建立的各重复 相 应 二次回归模型的回归系数列入表2; 显著 性 检验见表3:

表3 各次重复回归模型的显著性检验

重复 变因	1	2	3	4		
 SS总	563,64	538.75	837.89	422.75		
SS回	544.06	475.91	365.10	364.23		
SS误	2.90	4.50	2.40	2.50		
SS失	16.67	58.34	70.39	52.93		
SS离	19.57	62.84	72.78	55.43		
F 值	11.26	2.94	1.91	2,54		
显著性	0.01	0.05	ns	0.05		

由表3可知,除第3重复外其余各重复试验均达到0.05以上的显著性水平。为了进一步判断模型的稳定性,在AIC准则意义下,即模型中的有效(显著)参数个数,得到各重复相应的回归模型中各偏回归系数的显著性分析(表4)。

表4表明,第 I 重复相应的模型中显著 ( P < 0.05 ) 的参数个数为16与另外三个模型相比达到最大,第 II 重复相应的模型中显著 ( P < 0.05 ) 的参数个数为7与另外三个模型比较为最小。由于模型的稳定性是随模型中所包含的参数增加而下降的。因此,稳定性指标以模型 I 为最差,模型 I 最优。但是,从表3的分析结果表明,由于模型 I 的显著水平最高,因而 拟合性最好。这样看来模型的拟合性和稳定性之间有时会出现矛盾。赤池弘次认为,构造 的模型推断准则是综合了矛盾两方面协调作用的统计量。

根据⑤式,以N=36次 试验, $\sigma^2$ 为试验误差的方差估计量,( $Y-X\beta$ ) $^T$ ( $Y-X\beta$ )为 离回归变异平方和,K是模型中所含的有效参数个数。求出各试验重复模型的AIC统计量。

聚4 各重复回归模型中偏回归系数的显著性分析

次 项	$X_3 \mid X_4 \mid X_5 \mid X_1 \times_2 \mid X_1$	1 1 1 1	2.67 6 4.17 2.35 0.	18.62     8.28     18.62     12.93     7.28     0	0.01 0.05 0.01 0.01 0.05	1 1 1 1 1	.01 0.38 0.38 0.04 47.53 13.	784.0824.08 0.75 0.75 0.08 95.36 27.	0.01 ns ns 0.01 0.	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	38 2.04 3.38 0.04 12.92 4.	5170.15 7.66 12.66 0.19 48.45 15.	0.01 0.05 0.01 ns 0.01 0.	1 1 1 1	260.0435.04 3.38 1.04 7.04 22.78 5.	26,1512,15 3,75 25,35 82,01 19.	0.01 0.01 ns 0.01 0.01 0.
赵	1X 1 X 1 X 1 X 1 X 5	1 1 1	0.01 2.35 0.01	0.04 7.28 0.04	ns   0.05   ns	1	78 2.53 2.53	56 5.06 5.06	01 0.05 0.05		25 7.67 7.67 1	95 28,76 28,76 4	.01 0.01 0.01	1 1 1	28 2.53 2.53 1	01 9.11 9.11 7	01. 0.05 0.05
中	X2X, X2X, x2x	1 1	0.01 6.25 2.	0.04 19.39 6.	ns 0.01 0.	1 1	2.53 1.56 0.	5.06 3.13 0.13	0.05 ns n	-	12.09 0.56 0.	5.33 2.11 0	0.01 ns n	1 1	9.53 5.06 1.	0.31 18.23 5.	0.01 0.01 0.
	$\mathbf{x_g} \left  \mathbf{x_3} \mathbf{x_4} \right  \mathbf{X_5} \left  \mathbf{X_4} \right  \mathbf{X}$	1 1 1	.25 9.256.25 4	98 19,39 19,39 12,41	05 0.010.010.010	1 1 1 1	0.06 0.060.060.060	3 0,130,130,130	str str s	1 2 2 2	0.06 0.560.561.571	.23 2.112.115.855	s ns 0.050	1 1 1	56 0.060.060.56	63 0.23 0.23 2.03	.05 ns ns ns
11 次 员	(2   X2   X3   X4   x 5	1 1 1 1 1	9 4 9 4 1	27.93 12.41 27.93 14.41 3.10	.010.010.010.01 ns	1 1 1 1 1	0.060.060.060.060.06	.13.0130.130.130.13	ns ns ns ns	1 1 1 1 1	.570.563.060,560.56	.852,1111.482,112,11	.05 ns 0.01 ns ns	1 1 1 1 1	0.560.060.060.060.06	2,030,230,230,230,23	su su su su su

AIC<sub>1</sub>: 104.48 AIC<sub>2</sub>: 231.41 AIC<sub>3</sub>: 560.89 AIC<sub>4</sub>: 365.98

于是可以推断,在模型的拟合性和稳定性下,四个模型的优良性顺序是: I、I、I、I、I。在本例中,模型的拟合性是最佳模式推断的主要因素。事实上,尽管模型 I 的稳定性最差,但由于它的拟合性比其它模型强得多(显著性水平的数量级差别),因而弥补了稳定性方面的不足,最终被选为最佳模型。第 I 模型的AIC估计量与第 IV 模型相应的AIC值小134.57。这反映出了这两个几乎相同拟合性的模型在稳定性方面表现出的巨大差别,即在综合试验资料提供的信息量能力方面,模型 II 比之模型 II 优良得多。因此,在只有 II、IV 两模型供选择时,按AIC准则,应该毫不迟疑地选择前者作为生产过程的预测和控制的数学模型。笔者认为,从AIC准则和统计模型的优良性推断,能尽可能将模型的拟合性集中到少数几个参数上的模型是好的模型。

利用二次回归正交旋转设计试验建立刻划某一生产过程的数学模型,其目的的还主要在于研讨该过程中涉及的主要动态因子的数量关系,进而制订生产过程的最佳 决策 方案。因此,数学模型的代表性是必须的基础,而拟合性及稳定性则又是统计模型代表性 的主要 内容。可见本文构造出的二次回归模型的优良性推断准则,其实际意义是显然的。

#### 参考文献

- ①茆诗松等: 回归分析及其试验设计, 华东师大出版社, 1981
- ②朱伟勇、回归设计及其应用,数学的实践与认识(3),1978
- ③朱伟勇,回归设计及其应用,数学的实践与认识(4),1978
- ④刘璋温:选择回归模型的几个准则,数学的实践与认识(4),1983
- ⑤刘璋温. 赤池信息量准则AIC及其意义,数学的实践与认识(3),1980
- ⑥重庆市《杂交中稻高产栽培技术规范研究》协作组.杂交中稻高产栽培技术规范研究试验、示范总结、西农科技(3),1986
- @Robert, R. Sokal, Biometry, W. H. Freeman and Company, 1969

# CHOICE OF MATHEMATICAL MODEL FROM REGRESSIONAL ORTHOGONAL AND VOTATIONAL DESIGN EXPERIMENT AT TWO EXPONENTS (ROVDETE) WITH REPLICATION

Luo Zewei

(Department of Agronomy, Southwest Agricultural University)

#### Abstract

The statistics, a H. Akaike information criterion dealt with ROVDETE was constructed according to information criterion concerned about statistical model inference presented by H. Akaike. It had the following general form:

AIC=N\*
$$\ln \sigma^2$$
+ (Y-X\* $\hat{\beta}$ )<sup>T</sup>\* (Y-X\* $\hat{\beta}$ )/ $\sigma^2$ +2\*K

The formula above could be used to inference the choice of multi-mathematical model from the replicated ROVDETE.