

文章编号: 1000-6893(2006)04-0600-05

偏最小二乘回归在军用飞机价格预测中的应用

李寿安, 张恒喜, 童中翔, 郭 风, 董小龙

(1. 空军工程大学 工程学院, 陕西 西安 710038)

(2. 西北工业大学 航空学院, 陕西 西安 710072)

Application of Partial Least Squares Regression to the Acquisition Cost Estimating of Military Aircrafts

LI Shou-an, ZHANG Heng-xi, Dong Zhong-xiang, GUO Feng, Dong Xiaolong
(Engineering College, Air Force Engineering University, Xi'an 710038, China)

摘 要: 由于军用飞机性能要求的不断提高, 影响飞机采购价格的费用驱动因子繁多, 使得原有的价格预测模型已经不适用于现代军用飞机。分析了军用飞机采购价格样本数据少、费用驱动因子多的特点, 考虑到偏最小二乘回归(PLSR)方法在处理小样本多元数据方面的优势, 应用 PLSR 对军用飞机采购价格进行预测。PLSR 首先提取第 1、第 2 主成分对采购价格样本的特异点进行剔除; 然后进行变量投影重要度分析以筛选费用驱动因子; 最后, PLSR 对费用驱动因子进行回归建立军用飞机价格预测模型。实例表明, 在军用飞机价格预测方面, 与原有的预测模型和逐步多元回归模型相比, 应用 PLSR 预测的精度更高, 更能体现采购价格与飞机性能参数之间的关系。

关键词: 军用飞机; 采购价格; 偏最小二乘回归; 费用驱动因子

中图分类号: V37; O241.5 文献标识码: A

Abstract: With the continual promotion of the performances of military aircrafts, various Cost Drive Factors make the former Acquisition Cost estimating model not be applicable for modern military aircrafts. Analyzing small samples and various Cost Drive Factors of military aircrafts' Acquisition Cost and considering the advantages of Partial Least Squares Regression (PLSR) in analyzing multivariate data with small samples, an application of PLSR to military aircraft's Acquisition Cost estimating is proposed. Firstly, PLSR extracts the first and second principal components of military aircrafts' Acquisition Cost data to identify the outliers in samples; secondly, PLSR selects Cost Drive Factors with Variable Importance in Projection from the valid samples; finally, PLSR regresses to Cost Drive Factors and establishes the military aircraft's Acquisition Cost estimating model. The calculated results in the example show that, compared with the former estimating model and Stepwise Regression model, PLSR model has higher accuracy in military aircraft's Acquisition Cost estimating, and can reflect better the correlation between Acquisition Cost and performance parameters of military aircrafts.

Key words: military aircraft; Acquisition Cost; Partial Least Squares Regression; Cost Drive Factor

随着科学技术的迅猛发展, 军用飞机的系统复杂性和性能不断提高, 使得影响其采购价格的费用驱动因子繁多, 其采购价格急剧上涨^[1]。为节省国防资源和合理定价, 采购价格的预测在军用飞机采购活动中起着重要的作用。由于影响采购价格的费用驱动因子的增多, 原有的价格预测模型已经不适用于现代军用飞机, 采购价格的预测迫切需要新的方法和模型^[2,3]。针对军用飞机样本数据少和费用驱动因子多的特点, 本文提出应

用偏最小二乘回归(PLSR)对其采购价格进行预测。

1 军用飞机采购价格预测分析

1.1 采购价格的定义

军用飞机采购价格由其定价成本和按定价成本的 5% 的利润率计算的利润两部分组成^[3]。军用飞机的定价成本含制造成本和期间费用两部分。制造成本包括直接材料、辅助材料、备品配件、外购半成品、燃料、动力、包装物以及其它直接费用; 期间费用包括管理费用和财务费用两部分^[1,3]。产品的研制、生产费用一直被认为高度

收稿日期: 2005-01-05; 修订日期: 2005-06-06
基金项目: 国防预研基金(98J19.3.2. JB 3201)、部委级重点课题资助项目

机密的商业信息, 一般不予公开, 因此很难采用工程估算法对采购价格进行预测。

1.2 采购价格预测分析

一般采用参数法对军用飞机采购价格进行预测。在费用估算模型中, 影响费用的特征参数为费用驱动因子, 对费用有影响的飞机特征参数很多, 不可能也没有必要全部选择, 仅需选取重要的特征参数作为费用驱动因子^[1]。费用驱动因子的量值在设计和研制初期应该易于确定。

根据参考文献[1~3], 军用飞机采购价格一般与其性能参数成对数线性关系。原有的价格预测一般采用空重、最佳高度的最大平飞速度和满油航程 3 个变量, 通过多元回归方法建立战斗机和攻击机的出厂价格预测模型。

$$y = ax^1 x^2 x^3 \quad (1)$$

式中: y 为采购价格; x_1, x_2 和 x_3 分别为飞机空重、最佳高度的最大平飞速度和满油航程; a_0, a_1, a_2 与 a_3 为相关系数。

现代军用飞机要求的日益提高, 使得影响采购价格的因素繁多, 原有的模型已经不能准确地对价格进行预测, 必须根据实际情况重新选择费用驱动因子以建立新的采购价格预测模型。

并且对于我国实际情况来说, 完成的研制项目有限, 样本很少; 而军用飞机对新技术、高性能的要求的不断提高, 使得影响采购价格的因素多而复杂, 费用驱动因子难于选择^[4]。针对样本数据少和费用驱动因子多的特点, 应当选取恰当的方法建立一种适合我国国情的价格预测模型。而 PLSR 在处理小样本多元数据方面具有独特的优势, 因此可以采用 PLSR 方法对军用飞机价格进行预测。

2 PLSR 方法

PLSR 方法的一个突出特点是它将多元线性回归分析、变量的主成分分析和变量间的典型相关分析有机地结合起来, 在一个算法下, 同时实现了回归建模、数据结构简化和两组变量间的相关分析, 给多元数据分析带来了极大的便利^[5,6]。PLSR 能够在自变量存在严重多重相关性的条件下进行回归建模, 可以比最小二乘回归更简捷地进行自变量的筛选, PLSR 方法丰富的辅助分析技术可以在建模的同时实现对自变量的筛选^[6]。

2.1 PLSR 建模步骤

设已知因变量 y 和 k 个说明性变量 $x_1, x_2,$

\dots, x_k , 样本数为 n , 构成数据表 $y = y_{n \times 1}$ 和 $X = [x_1 \ x_2 \ \dots \ x_k]_{n \times k}$ 。

(1) 将 X 与 y 进行标准化处理, 得到标准化的自变量矩阵 E_0 和因变量矩阵 F_0 。标准化处理目的是为了公式表达上的方便和减少运算误差^[2,5]。

$$x_{ij}^* = \frac{x_{ij} - \bar{x}_j}{s_j}, \quad i = 1, \dots, n; \quad j = 1, \dots, k \quad (2)$$

$$E_0 = (x_{ij}^*)_{n \times k} \quad F_0 = \begin{pmatrix} y_1 - \bar{y} \\ \vdots \\ y_n - \bar{y} \end{pmatrix}_{n \times 1} \quad (3)$$

式中: \bar{x} 是 x_j 的均值, s_j 是 x_j 的标准差; \bar{y} 是 y 的均值, s_y 是 y 的标准差。

(2) 从 E_0 中抽取一个成分 $t_1 = E_0 w_1$, 其中

$$w_1 = \frac{E_0^T F_0}{\|E_0^T F_0\|} \quad (4)$$

实施 E_0 和 F_0 在 t_1 上的回归

$$E_0 = t_1 p_1 + E_1 \quad F_0 = t_1 r_1 + F_1 \quad (5)$$

式中: p_1, r_1 是回归系数(r_1 是标量), 即

$$p_1 = \frac{E_0^T t_1}{\|t_1\|^2}, \quad r_1 = \frac{F_0^T t_1}{\|t_1\|^2} \quad (6)$$

记残差矩阵

$$E_1 = E_0 - t_1 p_1, \quad F_1 = F_0 - t_1 r_1 \quad (7)$$

检查收敛性, 若 y 对 t_1 的回归方程已达到满意的精度(可用交叉有效性^[2,5]确定), 则进行下一步; 否则, 令: $E_0 = E_1, F_0 = F_1$, 回到第(2)步, 对残差矩阵进行新一轮的成分提取和回归分析。

(3) 在第 h 次主成分提取与回归($h = 1, 2, \dots, m$), 回归方程满足精度要求, 这时得到 h 个成分 t_1, t_2, \dots, t_m , 实施 F_0 在 t_1, t_2, \dots, t_m 上的回归, 得

$$F_0 = r_1 t_1 + r_2 t_2 + \dots + r_m t_m \quad (8)$$

由于 t_1, t_2, \dots, t_m 均是 E_0 的线性组合, 因此, 可写成 E_0 的线性组合形式, 即

$$F_0 = r_1 E_0 w_1^* + \dots + r_m E_0 w_m^* \quad (9)$$

式中: $w_h^* = \prod_{j=1}^{h-1} (I - w_j p_j^T) w_h$, I 为单位矩阵。

最后, 就有

$$\hat{y}^* = \alpha_1 x_1^* + \dots + \alpha_p x_p^* \quad (10)$$

x_j^* 的回归系数为

$$\alpha_i = \sum_{h=1}^m r_h w_{ij}^* \quad (11)$$

式中: w_{ij}^* 是 w_h^* 的第 j 个分量。

(4) 按照标准化的逆过程, 将 $F_0(\hat{y}^*)$ 的回归方程还原为 y 对 X 的回归方程。

2.2 辅助分析方法

(1) 特异点的剔除 样本中特异点的存在会

对回归产生明显的拉动作用,使模型偏离原统计规律。PLSR 可以利用自变量中提取的主成分在二维平面图上对样本点分布结构作直观观察,并判别样本特异点^[5]。

定义第 i 个样本点对第 h 主成分 t_h 的贡献率为

$$T_{hi,2} = \frac{t_{hi,2}}{(n-1)s_{h,2}} \quad (12)$$

式中: $s_{h,2}$ 为主成分 t_h 的方差; t_{hi} 为 t_h 的第 i 个分量。则样本点 i 对主成分 t_1, t_2, \dots, t_m 的累计贡献率为

$$T_{i,2} = \frac{1}{(n-1)} \sum_{h=1}^m \frac{t_{hi,2}}{s_{h,2}} \quad (13)$$

$T_{i,2}$ 可用来判别样本中的特异点。如果某一样本点的 $T_{i,2}$ 值过大,则说明这一样本对主成分构成的贡献过大,成为一个特异点。Tracy 给出了一个统计量用以检验^[5,6]

$$\frac{n^2(n-m)}{m(n^2-1)} T_i^2 \sim F(m, n-m) \quad (14)$$

当

$$T_i^2 \geq \frac{m(n^2-1)}{n^2(n-m)} F_{\alpha}(m, n-m) \quad (15)$$

可以认为在 $1-\alpha$ 的检验水平上,样本点 i 为特异点。

当 $m=2$ 时,这个判别条件为

$$\left(\frac{t_{1i}^2}{s_1^2} + \frac{t_{2i}^2}{s_2^2} \right) \geq \frac{2(n-1)(n^2-1)}{n^2(n-2)} F_{\alpha}(2, n-2) \quad (16)$$

这是一个椭圆。一般而言,提取 2 个主成分即可包含变量系统中绝大部分变异信息,所以在 t_1/t_2 平面图上作出这个椭圆图,如果所有的样本点都落在椭圆内,则认为没有特异点;反之落在椭圆外的点就认为是特异点。

(2) 变量投影重要性分析 在 PLSR 分析中,第 j ($j=1, 2, \dots, k$) 个自变量对因变量的解释能力是以前变量投影重要性指标 VIP_j 来测度的。其定义为^[2,5]

$$VIP_j = \sqrt{\frac{k}{\text{Rd}(y)} \sum_{h=1}^m \text{Rd}(y; t_h) w_{hj}^2} \quad (17)$$

式中: w_{hj} 为 w_h 的第 j 个分量,用于衡量 x_j 对构造 t_h 主成分的贡献大小; $\text{Rd}(y; t_h)$ 和 $\text{Rd}(y)$ 分别为 y 由 t_h 和 t_1, t_2, \dots, t_h 所解释的变异精度,分别代表了 t_h 对 y 的解释能力和 t_1, t_2, \dots, t_m 对 y 的累计解释能力,且有

$$\text{Rd}(y; t_h) = r^2(y, t_h)$$

$$\text{Rd}(y) = \sum_{h=1}^m \text{Rd}(y; t_h)$$

$$\text{Rd}(X) = \frac{1}{k} \sum_{h=1}^m \sum_{j=1}^k r^2(x_j, t_h) \quad (18)$$

式中: $r(y, t_h)$ 和 $r(x_j, t_h)$ 分别为因变量 y 和自变量 x_j 与主成分 t_h 的相关系数。

VIP_j 定义式的意义是基于这样一个事实:由于 x_j 对 y 的解释是通过 t_h 来传递的,如果 t_h 对 y 的解释能力很强,而 x_j 在构造 t_h 时又起到了相当重要的作用,则 x_j 对 y 的解释能力就被视为很大。也就是说,如果在 $\text{Rd}(y; t_h)$ 值很大的 t_h 成分上, w_{hj} 取很大的值,则 x_j 对解释 y 就有很重要的作用。

3 PLSR 在采购价格估算中的应用

3.1 应用分析

军用飞机采购价格一般与其性能参数成对数线性关系,而 PLSR 仅能得到线性回归模型。因此,在价格预测过程中,首先应对样本数据对数化,然后对对数化的数据进行线性回归处理,再对得到的线性回归方程反对数化,可以得到采购价格与性能参数的对数线性关系式。应用 PLSR 预测军用飞机采购价格的步骤如下。

(1) 样本数据对数化。对数化的目的是将对数线性关系的数据转化为线性关系的数据。

(2) 剔除样本特异点。PLSR 对对数化数据标准化,提取第 1、2 主成分 t_1, t_2 , 在 t_1/t_2 平面上作出判断椭圆图以剔除样本特异点。

(3) 选择费用驱动因子。对自变量进行变量投影重要性分析,根据变量投影重要性指标的大小来提取费用驱动因子。

(4) PLSR 对费用驱动因子进行主成分提取,在一定精度控制下进行回归处理并得到回归方程。

(5) 对回归方程进行反对数化,可以得到采购价格与费用驱动因子之间的对数线性关系式,即价格预测模型。

3.2 实例分析

现以军用涡扇运输机采购价格预测模型的建立为例予以分析。表征涡扇运输机性能的特征参数很多,实例中取 8 个典型的特征参数,分别为最大起飞重量 x_1 、机身长 x_2 、机高 x_3 、起飞距离 x_4 、满油航程 x_5 、最佳高度的最大平飞速度 x_6 、飞机空重 x_7 和最大载油量 x_8 ; 价格用 y 表示,价格的

基准年度为 2004 年。表 1 中列出了 9 型涡扇运输机样本的性能参数与采购价格, 为了对建立的

模型进行误差分析和预测检验, 选取表中前 8 个子样为训练样本, I 机型为检验样本。

表 1 运输机样本的性能数据与价格

Table 1 Performance parameters and Acquisition Costs of transporters samples

机型	x_1/kg	x_2/m	x_3/m	x_4/m	x_5/km	$x_6/\text{m} \cdot \text{s}^{-1}$	x_7/kg	x_8/kg	$y/\text{万元}$
A	13494	23.5	8.43	867	4262	425	6597	5683	6666.7
B	6849	14.39	4.57	987	3701	746	3655	2640	3624.3
C	9979	16.9	5.12	1581	4679	874	5357	3350	6569.9
D	5670	13.34	4.57	536	3641	536	3656	1653	5586.23
E	63503	39.75	9.30	1859	6764	925	33183	21273	27768.8
F	22000	29.87	6.75	1200	2870	907	34360	5500	17575.2
G	21500	27.17	7.65	1050	2000	580	12200	5000	18137.6
H	70310	29.79	11.66	1091	7876	602	36300	36300	50476
I	21000	24.615	7.3	1300	3100	819.2	11700	6000	14250

(1) 原有价格预测模型 原有的价格预测一般采用空重、最佳高度的最大平飞速度和满油航程 3 个自变量, 通过多元回归方法建立采购价格预测模型为

$$y = 34.86x_1^{0.18} x_2^{-0.59} x_7^{0.88} \quad (19)$$

(2) 逐步多元回归模型 逐步多元回归是一种通用的变量筛选方法, 它是向前变量选择法和向后变量排除法的结合^[2]。逐步多元回归对样本数据进行回归, 得到回归方程为

$$y = 1.42x_1^{0.92} \quad (20)$$

方程式表明采购价格仅与最大起飞重量 x_1 有关, 因此不能真实地反映采购价格与飞机性能之间的关系。

(3) PLSR 预测模型

① 剔除样本异点 以 x_1, x_2, \dots, x_8 为自变量, y 为因变量, PLSR 提取两个主成分 t_1, t_2 。计算 t_1, t_2 的方差, 取置信度为 95%, 根据式 (16), 在 t_1/t_2 平面上作出椭圆图, 所有样本均在椭圆内, 如图 1 所示, 因此样本没有特异点。

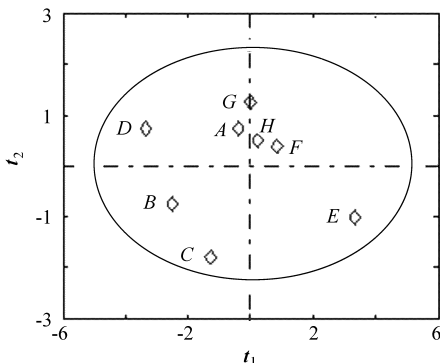


图 1 样本筛选椭圆图

Fig 1 Judgment ellipse of samples

控制下, PLSR 提取所有自变量中主成分进行回归, 得到 y 与 x_1, x_2, \dots, x_8 的回归方程, 并计算了每个自变量的 VIP, $Rd(X)$ 和 $Rd(y)$ 。结果如下

$$y = 16.48x_1^{0.39} x_2^{-0.08} x_3^{0.26} x_4^{-0.67} x_5^{-0.15} x_6^{0.45} x_7^{0.40} x_8^{0.21}$$

$$\text{VIP: } 1.21 \quad 1.12 \quad 1.10 \quad 0.79 \quad 0.68$$

$$0.54 \quad 1.21 \quad 1.12$$

$$Rd(X) = 86.5\%, Rd(y) = 64.8\% \quad (21)$$

$Rd(X)$ 和 $Rd(y)$ 计算值表明, 提取的主成分仅包含自变量中 86.5% 的变异信息和反映因变量中 64.8% 的变异信息。并且从每个自变量 VIP 值来看, x_4, x_5, x_6 的 VIP 值都小于 0.8, 因此他们对 y 的解释能力较小, 在回归分析中可以将它们排除。

采用 PLSR 对 x_1, x_2, x_3, x_7, x_8 与 y 进行回归, 获得了 x_1, x_2, x_3, x_7, x_8 与 y 之间回归方程。

$$y = 2.20x_1x_2^{-1.94} x_3^1 x_7^{0.79} x_8^{-0.66}$$

$$\text{VIP: } 1.02 \quad 0.95 \quad 0.94 \quad 1.04 \quad 1.01$$

$$Rd(X) = 99.6\%, Rd(y) = 97.1\% \quad (22)$$

提取的主成分包含了自变量中 99.6% 的变异信息和反映了因变量中 97.1% 的变异信息, 并且各个变量 VIP 值相对平均, 因此回归的方程是令人满意的。

(4) 精度分析 根据式 (19)、(20) 和 (22) 可以计算训练样本值、检验样本值, 将样本值与真实值相比较可以得到计算误差。计算的平均训练误差、检验样本值及其误差如表 2 所示。

表 2 计算结果与误差

Table 2 Calculated results and errors

PLSR	方法	原有模型	逐步多元回归
平均训练误差 / %	23.21	24.80	13.24
检验样本值 / 亿	10510	13680	13760
检验误差 / %	26.25	4.0	3.44

② 预测模型的建立 在交叉有效性分析地

PLSR 模型的平均训练误差和检验误差分别为 13.24% 和 3.44%，与原有模型的 23.21% 和 26.25%、逐步多元回归模型的 24.8% 和 4.0% 相比，应用 PLSR 预测军用飞机采购价格的精度最高，效果最好。

4 结论

样本少和影响价格的费用驱动因子繁多是军用飞机采购价格预测中两个难点，而 PLSR 能够利用丰富的辅助分析技术克服自变量中存在的多重相关性和选择有效费用驱动因子，在处理小样本多元数据方面具有独特的优势。因此，可采用 PLSR 预测军用飞机采购价格。

实例表明，与原有的预测模型和逐步多元回归模型相比，PLSR 在军用飞机价格预测方面具有更高的精度，并且能更好地反映采购价格与性能参数之间的关系。所以说，PLSR 在军用飞机价格预测中的应用是可行的、有效的。

参 考 文 献

- [1] 张恒喜. 现代飞机效率分析[M]. 北京: 航空工业出版社, 2001.
Zhang H X. Modern aircraft efficiency and cost analysis [M]. Beijing: Aviation Industry Press, 2001. (in Chinese)
- [2] 张恒喜, 郭基联, 朱家元. 小样本多元数据分析方法及应用[M]. 西安: 西北工业大学出版社, 2002. 22-43.
Zhang H X, Guo J L, Zhu J Y. Multivariate data analysis methods and applications with few observations[M]. Xi'an: Northwest Polytechnic University Press, 2002. 24-43. (in Chinese)
- [3] 张恒喜, 但福堂. 飞机全寿命费用分析与控制[M]. 西安: 空军工程学院, 1991.
Zhang H X, Dan F T. Aircraft life cycle cost analysis and control[M]. Xi'an: Air Force Engineering Institute, 1991.

40-43. (in Chinese)

- [4] Palen A. System cost modeling an integrated approach [R]. AIAA-92-1279, 1992.
- [5] 王惠文. 偏最小二乘回归方法及其应用[M]. 北京: 国防工业出版社, 1999.
Wang H W. Partial Least squares Regression method and applications[M]. Beijing: National Defense Industry Press, 1999. (in Chinese)
- [6] 蒋红卫, 夏结来. 偏最小二乘回归及其应用[J]. 第四军医大学学报, 2003, 24(3): 280-283.
Jiang H W, Xia J L. Partial least squares regression and its application[J]. Journal of the Fourth Military Medicine University, 2003, 24(3): 280-283. (in Chinese)

作者简介:



李寿安(1977-) 男, 江西萍乡人, 空军工程大学工程学院飞行器设计专业博士研究生。主要研究领域为飞机生存力设计, 飞机可靠性、维修性及保障性等。电话: 029-83687482; E-mail: ananan2002@163.com。



张恒喜(1937-) 男, 江苏姜堰人, 现任空军工程大学工程学院装备维修保障工程研究所所长, 教授, 博士生导师, 空军级专家, 享受政府特殊津贴, 中国数量经济学会常务理事, 军事系统工程委员会委员, 空军军标委委员。1965年毕业于西北工业大学飞机设计专业。主要研究方向为军用飞机型号发展工程。获得科研成果多项。



郭风(1979-) 男, 河南商丘人, 空军工程大学工程学院装备管理工程硕士研究生。主要研究方向为飞机效率分析。

(责任编辑: 李铁柏)