

Engineering Notes

ENGINEERING NOTES are short manuscripts describing new developments or important results of a preliminary nature. These Notes should not exceed 2500 words (where a figure or table counts as 200 words). Following informal review by the Editors, they may be published within a few months of the date of receipt. Style requirements are the same as for regular contributions (see inside back cover).

Method for Multivariate Analysis with Small Sample in Aircraft Cost Estimation

Shou'an Li* and Bifeng Song†

Northwestern Polytechnical University,
710072 Xi'an, Shaanxi, People's Republic of China
and

Hengxi Zhang‡

Air Force Engineering University,
710038 Xi'an, Shaanxi, People's Republic of China

DOI: 10.2514/1.28174

Nomenclature

E_0	=	standardized independent variable matrix
F_0	=	standardized dependent variable matrix
h	=	number of principal components
k	=	number of independent variables
n	=	number of samples
$Rd(X)$	=	variation accuracy that X is explained by t_1, t_2, \dots, t_h
$Rd(y)$	=	variation accuracy that y is explained by t_1, t_2, \dots, t_h
$Rd(y; t_i)$	=	variation accuracy that y is explained by t_i
$r(x_j, t_i)$	=	correlation coefficient of x_j and t_i
$r(y, t_i)$	=	correlation coefficient of y and t_i
t_1, t_2, \dots, t_h	=	extracted principal components
X	=	independent variable matrix
x_1, x_2, \dots, x_k	=	independent variables
Y	=	dependent variable matrix
y	=	dependent variable
α	=	confidence level
$\alpha_1, \alpha_2, \dots, \alpha_k$	=	regression coefficients

I. Introduction

WITH the development of science and technology, the efficiency and complexity of modern aircraft are enhanced continually, and aircraft cost increases rapidly. Life cycle cost has been a decisive factor for modern aircraft design. To control the continual rapid increased life cycle cost, aircraft cost estimation is indispensable early in the design phase.

The cost data during the development and manufacturing phase are strictly confidential for aircraft companies, so aircraft cost is hard

to be estimated by the engineering method. The parametric method for estimating aircraft cost has been used extensively in modern aircraft design. The parametric cost models are designed to be used when little is known about an aircraft design or when a readily applied validity and consistency check of detailed cost estimates is necessary [1].

Since 1966, Rand Corporation has developed a series of parametric cost models for military aircraft using multiple least-squares regression method [1]. For small homogeneous sample, only the weight and speed are selected as the cost-related explanatory variables in these models generally [1,2]. Aircraft cost-estimating relationships in the models are generally expressed as

$$y = a_0 x_1^{a_1} x_2^{a_2} \quad (1)$$

where y is the cost, x_1 is the weight, and x_2 is the speed; a_0 , a_1 , and a_2 are the related coefficients.

The parametric cost models are based on aircraft sample. The homogeneity and size of the sample are important for model accuracy [1–3]. Multiple least-squares regression method requires a larger sample relatively to the number of the explanatory variables. Some explanatory variables that have good relationships with aircraft cost have to be eliminated, and only two or three very representative variables are selected into the cost-estimating model using multiple least-squares regression due to multivariate data with small sample [1,2]. Various cost-related explanatory variables, small sample and the presence of multicollinearity that multiple least-squares regression is hard to overcome often plague aircraft cost estimation. Partial least-squares regression (PLSR) is very appropriate to cope with these problems, so PLSR is applied to multivariate analysis and parametric modeling for estimating aircraft cost in the paper.

II. PLSR Method

PLSR method is a statistical tool that has been specially designed to deal with multiple regression where the number of observations is limited, missing data are numerous, and the correlations between variables are high [4]. It is a recent technique that generalizes and combines features from principal component analysis and multiple least-squares regression.

A. PLSR Characteristics

1) High multicollinearity between variables can be overcome by PLSR. The data can be decomposed and reconstructed, the synthetic variables with the best explanation to the dependent variables can be extracted from the independent variables, and the information and noises in the data can be identified by PLSR.

2) Independent variable selection using PLSR is more convenient than using multiple least-squares regression [4,5]. The aided analysis techniques of PLSR can be helpful to select independent variables during regressing.

3) Through principal component analysis and synthetic variable extraction, PLSR can be used to deal with multivariate data where the number of observations is fewer than that of variables, but multiple least-squares regression cannot do this [4,5].

4) The regression equation derived by PLSR can contain all the independent variables and most of the data information for

Received 4 October 2006; revision received 4 January 2007; accepted for publication 31 January 2007. Copyright © 2007 by the American Institute of Aeronautics and Astronautics, Inc. All rights reserved. Copies of this paper may be made for personal or internal use, on condition that the copier pay the \$10.00 per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923; include the code 0021-8669/07 \$10.00 in correspondence with the CCC.

*Postdoctoral Research Associate, Department of Aircraft Design, Postbox 120; ananana2002@163.com.

†Professor, Department of Aircraft Design.

multivariate analysis with small sample, so it is more effective than that derived by multiple least-squares regression.

With these proprietary characteristics, PLSR can be applied to analyze aircraft cost-estimating relationships better than multiple least-squares regression.

B. Modeling Steps of PLSR [6,7]

The dependent variable y of n observations is described by an $n \times 1$ matrix denoted Y , the independent variables x_1, x_2, \dots, x_k of n observations are described by an $n \times k$ matrix denoted X :

$$X = [x_{ij}]_{n \times k} \quad Y = [y_i]_{n \times 1} \quad i = 1, \dots, n; \quad j = 1, \dots, k \quad (2)$$

where x_{ij} is the number j independent variable of the number i observation, and y_i is the dependent variable of the number i observation.

Step 1: Standardize X and Y . The aim of standardization is to make the data dimensionless. E_0 and F_0 are the standardized forms of X and Y :

$$E_0 = \left[\frac{x_{ij} - \bar{x}_j}{s_j} \right]_{n \times k} = [x_{ij}^*]_{n \times k} \quad F_0 = \left[\frac{y_i - \bar{y}}{s_y} \right]_{n \times 1} = [y_i^*]_{n \times 1} \quad (3)$$

In Eq. (3) \bar{x}_j and \bar{y} are the mean values of x_j and y , respectively, and s_j and s_y are the standard deviations of x_j and y , respectively; $x_j = [x_{1j}, x_{2j}, \dots, x_{nj}]^T$, $y = [y_1, y_2, \dots, y_n]^T$.

Step 2: Regression analysis. The first principal component t_1 is extracted from E_0 . $t_1 = E_0 w_1$, and

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

E_0 and F_0 are regressed on t_1 :

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

where p_1 and r_1 are the regression coefficients (r_1 is a scalar quantity), and

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

The matrix remnants are described as

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

Step 3: Accuracy analysis. If the regression equation of y on t_1 achieves the required accuracy, next step continues; else, $E_0 = E_1$, $F_0 = F_1$, and step 1, 2 are repeated to extract a new principal component from the matrix remnants.

Step 4: If the number h principal component is extracted ($h = 1, 2, \dots, k$), the accuracy requirement is achieved. The regression equation of F_0 on t_1, t_2, \dots, t_h can be derived by PLSR:

$$\hat{F}_0 = r_1 t_1 + r_2 t_2 + \dots + r_h t_h \quad (8)$$

Because t_1, t_2, \dots, t_h are the linear combination of E_0 , Eq. (8) can also be expressed as Eq. (9):

$$\hat{F}_0 = r_1 E_0 w_1^* + \dots + r_h E_0 w_h^* \quad (9)$$

where $w_h^* = [\prod_{j=1}^{h-1} (I - w_j p_j')] w_i$ ($i = 1, 2, \dots, h$), I is a unit matrix.

Finally, Eq. (9) can be expressed as

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

where $x_j^* = [x_{1j}^*, x_{2j}^*, \dots, x_{nj}^*]^T$, $y^* = [y_1^*, y_2^*, \dots, y_n^*]^T$, and the regression coefficient α_j of x_j^* is

$$\alpha_j = \sum_{i=1}^n r_i w_{ij}^* \quad (11)$$

In Eq. (11), w_{ij}^* is the number j element of w_i^* .

Step 5: Reversing the process of standardization, Eqs. (9) or (10) can be transformed to the regression equation of y on x_1, x_2, \dots, x_k .

C. Supplementary Analysis Methods [7,8]

1. Finding out Outliers

The outliers in a sample may make the regression equation deviate from the original statistic rules. The outliers in observations can be identified through the principal component analysis on a 2-D coordinate plane using PLSR.

The contribution rate of the number j observation to t_i is defined as

$$T_{ij}^2 = t_{ij}^2 / (n-1) s_i^2 \quad (12)$$

where s_i is the standard deviation of t_i , t_{ij} is the number j element of t_i . The accumulated contribution rate of the number j observation to t_1, t_2, \dots, t_h is

$$T_j^2 = \frac{1}{(n-1)} \sum_{i=1}^h \frac{t_{ij}^2}{s_i^2} \quad (13)$$

Equation (13) can be applied to identify the outliers in observations. If T_j^2 is great, the observation is an outlier. A test statistic is offered as

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

If

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

The observation j is considered as an outlier under the confidence level of α .

If $m = 2$, Eq. (15) can be described as

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

It is an ellipse equation. Generally, two principal components t_1, t_2 can contain most of the variation of the variables. Plotting the ellipse on the t_1, t_2 coordinate plane, if all the observations are in the ellipse, there is no outlier; else, there are outliers outside the ellipse.

2. Variable Importance in Projection Analysis

The explanation ability of x_j ($j = 1, \dots, k$) to y can be measured by variable importance in projection (VIP) using PLSR. The VIP of x_j denoted VIP_j is defined as

$$VIP_j = \sqrt{\frac{k}{Rd(y)} \sum_{i=1}^h Rd(y; t_i) w_{ij}^2} \quad (17)$$

In Eq. (17), w_{ij} , which can measure the contribution rate of x_j in the construction of t_i , is the number j element of w_i ; $Rd(y; t_i)$ and $Rd(y)$ are, respectively, the variation accuracies that y is explained by t_i and t_1, t_2, \dots, t_h , and can express the explanation abilities of t_i and t_1, t_2, \dots, t_h to y . $Rd(X)$ is the variation accuracy that X is explained by t_1, t_2, \dots, t_h , and can express the explanation ability of t_1, t_2, \dots, t_h to X .

$$\begin{aligned} \text{Rd}(y; t_i) &= r^2(y; t_i) \quad \text{Rd}(y) = \sum_{i=1}^h \text{Rd}(y; t_i) \\ \text{Rd}(X) &= \frac{1}{k} \sum_{i=1}^h \sum_{j=1}^k r^2(x_j, t_i) \end{aligned} \quad (18)$$

where $r(y, t_i)$ is the correlation coefficient of y and t_i , and $r(x_j, t_i)$ is that of x_j and t_i .

If the explanation ability of t_i to y is strong and x_j is important in the construction of t_i , then the explanation ability of x_j to y is strong, and the VIP_j is great. So VIP analysis can be applied to the independent variable selection.

III. Application of PLSR to Aircraft Cost Estimation

A. Application Analysis

Generally, aircraft cost has good log-linear relationships with the cost-related explanatory variables [1–3]. PLSR can be used to deal with the linear relationships between variables. In the parametric cost modeling, first, it is necessary to transform the cost data to the logarithmic data; second, the transformed data are standardized, and two principal components t_1, t_2 are extracted to find out outliers from the observations using PLSR; third, under the control of cross validation analysis [7,8], PLSR is used to extract principal components from the standardized data, and select representative cost-related explanatory variables through VIP analysis; finally, principal components are extracted from the representative explanatory variables and the cost-estimating relationships are derived by PLSR.

B. Case Study

The application process is illustrated with estimating airframe development cost for eight fighter observations. There are a lot of cost-related parameters for fighters. Seven representative parameters are selected as the explanatory variables in the case, and they are empty weight x_1 , maximum speed x_2 , development period x_3 , combat radius x_4 , takeoff distance x_5 , climb rate x_6 , and gross payload x_7 . The airframe development cost is presented with y . A, B, C, D, E, F, G, and H in Table 1 are the fighter observations. To analyze errors and test the validity of results in the estimation, A, B, C, D, E, F, and G are selected as seven training observations, and H is selected as a test observation.

1. PLSR Analysis

The variables x_1, x_2, \dots, x_7 of the eight observations are described by the matrix $X = [x_{ij}]_{8 \times 7}$, and the variable y of the observations is described by the matrix $Y = [y_i]_{8 \times 1}$.

First, two principal components t_1, t_2 are extracted from X , and the standard deviations of t_1, t_2 are calculated. According to Eq. (16), under the confidence level of 95%, the judgment ellipse is plotted on the t_1, t_2 coordinate plane. All the fighter observations are in the ellipse (see Fig. 1), so there is no outlier in the observations.

Second, the principal components are extracted from the independent variables of the training observations, and the regression equation of y on X is derived by PLSR under the control of cross validation analysis. The VIP values of x_1, x_2, \dots, x_7 , $\text{Rd}(X)$,

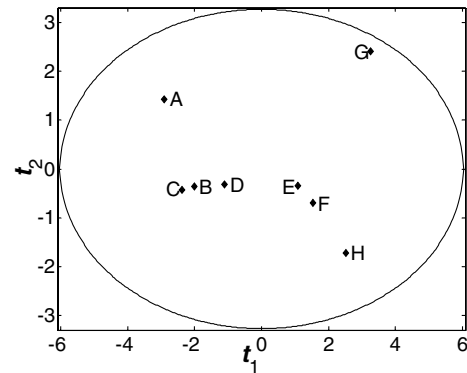


Fig. 1 The judgment ellipse of fighter observations.

and $\text{Rd}(y)$ of the equation are calculated. The regression equation and the calculated results are described as

$$\begin{aligned} y &= 2.332 \times 10^{-5} x_1^{0.449} x_2^{-0.963} x_3^{0.861} x_4^{-0.512} x_5^{0.555} x_6^{1.254} x_7^{0.617} \\ \text{VIP: } &1.21 \quad 0.48 \quad 1.25 \quad 1.28 \quad 0.65 \quad 0.55 \quad 1.17 \\ \text{Rd}(X) &= 95.3\%, \quad \text{Rd}(y) = 89.6\% \end{aligned} \quad (19)$$

The $\text{Rd}(X)$ and $\text{Rd}(y)$ show that the extracted principal components can contain 95.3% variation in the independent variables and 89.6% in the dependent variable. It is seen from the VIPs, the VIP values of the variables x_2, x_5, x_6 are obviously smaller. It shows that x_2, x_5, x_6 are not important in cost modeling relatively and should be eliminated.

The variables x_1, x_3, x_4, x_7 , are selected as the cost drive factors, and the regression equation of y on x_1, x_3, x_4, x_7 is derived as follows:

$$\begin{aligned} y &= 2.081 \times 10^{-2} x_1^{-3.273} x_3^{2.748} x_4^{1.151} x_7^{3.302} \\ \text{VIP: } &1.00 \quad 1.03 \quad 0.97 \quad 1.00 \quad \text{Rd}(X) = 97.6\% \\ \text{Rd}(y) &= 94.1\% \end{aligned} \quad (20)$$

The extracted principal components can contain 97.6% variation in the independent variables and 94.1% in the dependent variable, and the VIP values of the independent variables in Eq. (20) are relatively average. So the equation is very satisfied.

2. Multiple Least-Squares Regression Analysis

The empty weight x_1 and maximum speed x_2 are selected as cost drive factors and the cost-estimating model is derived by multiple least-squares regression method [1,8]:

$$y = 8.577 \times 10^{-5} x_1^{2.23} x_2^{-0.58} \quad (21)$$

3. Stepwise Regression Analysis

Stepwise regression is a variable selection method that combines the forward variable selection method with the backward variable elimination method [7]. The regression equation is derived by stepwise regression as follows:

Table 1 The cost data of fighter observations

Observations	x_1 , kg	x_2 , km/h	x_3 , y	x_4 , km	x_5 , m	x_6 , m/s	x_7 , kg	y , million ¥
A	3,845	1812	7	400	915	135	1300	145.50
B	4,580	2280	7	420	825	150	1430	179.00
C	4,120	2170	6	450	800	143	1500	131.25
D	5,660	2105	7	500	720	154	1650	100.18
E	6,915	2290	10	770	685	178	1970	375.68
F	7,480	2310	9	810	650	193	2290	671.87
G	11,890	1890	13	860	785	200	3600	1468.94
H	7,810	2370	9	915	530	225	2500	628.64

Table 2 Calculated results and errors

Method	Average training error, %	Test sample value, million ¥	Test error, %
Multiple least-squares regression	33.33	484.02	23.00
Stepwise regression	29.84	387.85	38.30
PLSR, Eq. (19)	24.43	572.83	8.88
PLSR, Eq. (20)	18.40	674.53	7.30

$$y = 0.215x_3^{3.41} \quad (22)$$

Equation (22) shows that the airframe development cost has only to do with the development period x_3 , so it cannot reflect the correlations between the airframe development cost and the cost-related explanatory variables completely.

4. Accuracy Analysis

The values of the observations are calculated by Eq. (19–22), and the calculating errors can be obtained from the calculated values and the actual values. The average training errors, test sample values, and test errors calculated by different methods are listed in Table 2.

It is seen from the calculated results and errors in the Table 2, the average training error and test error calculated by PLSR are smaller than those calculated by multiple least-squares regression and stepwise regression, and through VIP analysis and the explanatory variable selection, Eq. (20) has higher accuracy than Eq. (19) in the airframe development cost estimation using PLSR.

IV. Conclusions

Multivariate data with small sample in aircraft cost estimation is hard to deal with by the ordinary regression methods. PLSR method can effectively be used to estimate aircraft cost with reasonable accuracy. Through finding out the outliers in observations, VIP

analysis, and the explanatory variable selection, the good relationship for estimating aircraft cost can be derived by PLSR. The application of PLSR to aircraft cost estimation can play a very important role in the cost data analysis.

PLSR has been designed to deal with the problems such as small sample data and the high multicollinearity between variables. PLSR can not only be used to estimate aircraft cost, but can also be used in other statistical studies where the observations are few, the explanatory variables are various, or the multicollinearity between variables is high.

References

- [1] Hess, R. W., and Romanoff, H. P., "Aircraft Airframe Cost Estimating Relationships," Rand Corp., Rept. R-3255-AF, Santa Monica, CA, 1987.
- [2] Zhang, H. X., *Modern Aircraft Efficiency and Cost Analysis*, Aviation Industry Press, Beijing, 2001, pp. 115–121.
- [3] Castagne, S., Curran, R., Rothwell, A., Price, M., Benard, E., and Raghunathan, S., "A Generic Tool for Cost Estimating in Aircraft Design," AIAA Paper 2004-6235, 2004, pp. 1–15.
- [4] Wold, S., Trygg, J., Berglund, A., and Antti, H., "Some Recent Developments in PLS Modeling," *Chemometrics and Intelligent Laboratory Systems*, Vol. 58, No. 2, 2001, pp. 131–150.
- [5] Izenman, A. J., "Reduced-Rank Regression for the Multivariate Linear Model," *Journal of Multivariate Analysis*, Vol. 5, No. 2, 1975, pp. 248–264.
- [6] Wold, S., Ruhe, A., Wold, H., and Dunn, W. J., III, "The Collinearity Problem in Linear Regression: The Partial Least Squares Approach to Generalized Inverses," *SIAM Journal on Scientific Computing*, Vol. 5, No. 3, 1984, pp. 735–743.
- [7] Wang, H. W., *Partial Least-Squares Regression Method and Applications*, National Defense Industry Press, Beijing, 1999, pp. 34–56.
- [8] Zhang, H. X., Zhu, J. Y., and Guo, J. L., *Introduction to Military Aircraft Type Development Engineering*, National Defense Industry Press, Beijing, 2004, pp. 171–181.