

第二十七章 偏最小二乘回归分析

在实际问题中,经常遇到需要研究两组多重相关变量间的相互依赖关系,并研究用一组变量(常称为自变量或预测变量)去预测另一组变量(常称为因变量或响应变量),除了最小二乘准则下的经典多元线性回归分析(MLR),提取自变量组主成分的主成分回归分析(PCR)等方法外,还有近年发展起来的偏最小二乘(PLS)回归方法。

偏最小二乘回归提供一种多对多线性回归建模的方法,特别当两组变量的个数很多,且都存在多重相关性,而观测数据的数量(样本量)又较少时,用偏最小二乘回归建立的模型具有传统的经典回归分析方法所没有的优点。

偏最小二乘回归分析在建模过程中集中了主成分分析,典型相关分析和线性回归分析方法的特点,因此在分析结果中,除了可以提供一个更为合理的回归模型外,还可以同时完成一些类似于主成分分析和典型相关分析的研究内容,提供更丰富、深入的一些信息。

本章介绍偏最小二乘回归分析的建模方法;通过例子从预测角度对所建立的回归模型进行比较。

§1 偏最小二乘回归分析

考虑 p 个因变量 y_1, y_2, \dots, y_p 与 m 个自变量 x_1, x_2, \dots, x_m 的建模问题。偏最小二乘回归的基本作法是首先在自变量集中提出第一成分 t_1 (t_1 是 x_1, \dots, x_m 的线性组合,且尽可能多地提取原自变量集中的变异信息);同时在因变量集中也提取第一成分 u_1 ,并要求 t_1 与 u_1 相关程度达到最大。然后建立因变量 y_1, \dots, y_p 与 t_1 的回归,如果回归方程已达到满意的精度,则算法中止。否则继续第二对成分的提取,直到能达到满意的精度为止。若最终对自变量集提取 r 个成分 t_1, t_2, \dots, t_r , 偏最小二乘回归将通过建立 y_1, \dots, y_p 与 t_1, t_2, \dots, t_r 的回归式,然后再表示为 y_1, \dots, y_p 与原自变量的回归方程式,即偏最小二乘回归方程式。

为了方便起见,不妨假定 p 个因变量 y_1, \dots, y_p 与 m 个自变量 x_1, \dots, x_m 均为标准化变量。因变量组和自变量组的 n 次标准化观测数据阵分别记为

$$F_0 = \begin{bmatrix} y_{11} & \cdots & y_{1p} \\ \vdots & & \vdots \\ y_{n1} & \cdots & y_{np} \end{bmatrix}, \quad E_0 = \begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nm} \end{bmatrix}$$

偏最小二乘回归分析建模的具体步骤如下:

(1) 分别提取两变量组的第一对成分,并使之相关性达最大。

假设从两组变量分别提出第一对成分为 t_1 和 u_1 , t_1 是自变量集 $X = (x_1, \dots, x_m)^T$ 的线性组合: $t_1 = w_{11}x_1 + \cdots + w_{1m}x_m = w_1^T X$, u_1 是因变量集 $Y = (y_1, \dots, y_p)^T$ 的线性组合: $u_1 = v_{11}y_1 + \cdots + v_{1p}y_p = v_1^T Y$ 。为了回归分析的需要,要求:

- ① t_1 和 u_1 各自尽可能多地提取所在变量组的变异信息;
- ② t_1 和 u_1 的相关程度达到最大。

由两组变量集的标准化观测数据阵 E_0 和 F_0 , 可以计算第一对成分的得分向量,记

为 \hat{t}_1 和 \hat{u}_1 :

$$\hat{t}_1 = E_0 w_1 = \begin{bmatrix} x_{11} & \cdots & x_{1m} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{nm} \end{bmatrix} \begin{bmatrix} w_{11} \\ \vdots \\ w_{1m} \end{bmatrix} = \begin{bmatrix} t_{11} \\ \vdots \\ t_{n1} \end{bmatrix}$$

$$\hat{u}_1 = F_0 v_1 = \begin{bmatrix} y_{11} & \cdots & y_{1p} \\ \vdots & & \vdots \\ y_{n1} & \cdots & y_{np} \end{bmatrix} \begin{bmatrix} v_{11} \\ \vdots \\ v_{1p} \end{bmatrix} = \begin{bmatrix} u_{11} \\ \vdots \\ u_{n1} \end{bmatrix}$$

第一对成分 t_1 和 u_1 的协方差 $\text{Cov}(t_1, u_1)$ 可用第一对成分的得分向量 \hat{t}_1 和 \hat{u}_1 的内积来计算。故而以上两个要求可化为数学上的条件极值问题:

$$\begin{cases} \langle \hat{t}_1, \hat{u}_1 \rangle = \langle E_0 w_1, F_0 v_1 \rangle = w_1^T E_0^T F_0 v_1 \Rightarrow \max \\ w_1^T w_1 = \|w_1\|^2 = 1, \quad v_1^T v_1 = \|v_1\|^2 = 1 \end{cases}$$

利用Lagrange乘数法, 问题化为求单位向量 w_1 和 v_1 , 使 $\theta_1 = w_1^T E_0^T F_0 v_1 \Rightarrow$ 最大。问题的求解只须通过计算 $m \times m$ 矩阵 $M = E_0^T F_0 F_0^T E_0$ 的特征值和特征向量, 且 M 的最大特征值为 θ_1^2 , 相应的单位特征向量就是所求的解 w_1 , 而 v_1 可由 w_1 计算得到 $v_1 = \frac{1}{\theta_1} F_0^T E_0 w_1$ 。

(2) 建立 y_1, \dots, y_p 对 t_1 的回归及 x_1, \dots, x_m 对 t_1 的回归。

假定回归模型为

$$\begin{cases} E_0 = \hat{t}_1 \alpha_1^T + E_1 \\ F_0 = \hat{t}_1 \beta_1^T + F_1 \end{cases}$$

其中 $\alpha_1 = (\alpha_{11}, \dots, \alpha_{1m})^T$, $\beta_1 = (\beta_{11}, \dots, \beta_{1p})^T$ 分别是多对一的回归模型中的参数向量, E_1 和 F_1 是残差阵。回归系数向量 α_1, β_1 的最小二乘估计为

$$\begin{cases} \alpha_1 = E_0^T \hat{t}_1 / \|\hat{t}_1\|^2 \\ \beta_1 = F_0^T \hat{t}_1 / \|\hat{t}_1\|^2 \end{cases},$$

称 α_1, β_1 为模型效应负荷量。

(3) 用残差阵 E_1 和 F_1 代替 E_0 和 F_0 重复以上步骤。

记 $\hat{E}_0 = \hat{t}_1 \alpha_1^T$, $\hat{F}_0 = \hat{t}_1 \beta_1^T$, 则残差阵 $E_1 = E_0 - \hat{E}_0$, $F_1 = F_0 - \hat{F}_0$ 。如果残差阵 F_1 中元素的绝对值近似为 0, 则认为用第一个成分建立的回归式精度已满足需要了, 可以停止抽取成分。否则用残差阵 E_1 和 F_1 代替 E_0 和 F_0 重复以上步骤即得:

$$w_2 = (w_{21}, \dots, w_{2m})^T; \quad v_2 = (v_{21}, \dots, v_{2p})^T$$

而 $\hat{t}_2 = E_1 w_2$, $\hat{u}_2 = F_1 v_2$ 为第二对成分的得分向量。

$$\alpha_2 = \frac{E_1^T \hat{t}_2}{\|\hat{t}_2\|^2}, \quad \beta_2 = \frac{F_1^T \hat{t}_2}{\|\hat{t}_2\|^2}$$

分别为 X, Y 的第二对成分的负荷量。这时有

$$\begin{cases} E_0 = \hat{t}_1 \alpha_1^T + \hat{t}_2 \alpha_2^T + E_2 \\ F_0 = \hat{t}_1 \beta_1^T + \hat{t}_2 \beta_2^T + F_2 \end{cases}$$

(4) 设 $n \times m$ 数据阵 E_0 的秩为 $r \leq \min(n-1, m)$, 则存在 r 个成分 t_1, t_2, \dots, t_r , 使得

$$\begin{cases} E_0 = \hat{t}_1 \alpha_1^T + \dots + \hat{t}_r \alpha_r^T + E_r \\ F_0 = \hat{t}_1 \beta_1^T + \dots + \hat{t}_r \beta_r^T + F_r \end{cases}$$

把 $t_k = w_{k1}^* x_1 + \dots + w_{km}^* x_m$ ($k = 1, 2, \dots, r$), 代入 $Y = t_1 \beta_1 + \dots + t_r \beta_r$, 即得 p 个因变量的偏最小二乘回归方程式

$$y_j = a_{j1} x_1 + \dots + a_{jm} x_m, \quad (j = 1, 2, \dots, p)$$

这里 $w_h^* = (w_{h1}^*, \dots, w_{hm}^*)^T$ 满足 $\hat{t}_h = E_0 w_h^*$, $w_h^* = \prod_{j=1}^{h-1} (I - w_j \alpha_j^T) w_h$ 。

(5) 交叉有效性检验。

一般情况下, 偏最小二乘法并不需要选用存在的 r 个成分 t_1, t_2, \dots, t_r 来建立回归式, 而像主成分分析一样, 只选用前 l 个成分 ($l \leq r$), 即可得到预测能力较好的回归模型。对于建模所需提取的成分个数 l , 可以通过交叉有效性检验来确定。

每次舍去第 i 个观测 ($i = 1, 2, \dots, n$), 对余下的 $n-1$ 个观测值用偏最小二乘回归方法建模, 并考虑抽取 h 个成分后拟合的回归式, 然后把舍去的第 i 个观测点代入所拟合的回归方程式, 得到 $y_j (j = 1, 2, \dots, p)$ 在第 i 个观测点上的预测值 $\hat{y}_{(i)j}(h)$ 。对 $i = 1, 2, \dots, n$ 重复以上的验证, 即得抽取 h 个成分时第 j 个因变量 $y_j (j = 1, 2, \dots, p)$ 的预测误差平方和为

$$\text{PRESS}_j(h) = \sum_{i=1}^n (y_{ij} - \hat{y}_{(i)j}(h))^2 \quad (j = 1, 2, \dots, p)$$

$Y = (y_1, \dots, y_p)^T$ 的预测误差平方和为

$$\text{PRESS}(h) = \sum_{j=1}^p \text{PRESS}_j(h)。$$

另外, 再采用所有的样本点, 拟合含 h 个成分的回归方程。这时, 记第 i 个样本点的预测值为 $\hat{y}_{ij}(h)$, 则可以定义 y_j 的误差平方和为

$$\text{SS}_j(h) = \sum_{i=1}^n (y_{ij} - \hat{y}_{ij}(h))^2$$

定义 Y 的误差平方和为

$$\text{SS}(h) = \sum_{j=1}^p \text{SS}_j(h)$$

当 $\text{PRESS}(h)$ 达到最小值时, 对应的 h 即为所求的成分个数。通常, 总有 $\text{PRESS}(h)$ 大于 $\text{SS}(h)$, 而 $\text{SS}(h)$ 则小于 $\text{SS}(h-1)$ 。因此, 在提取成分时, 总希望比值 $\text{PRESS}(h)/\text{SS}(h-1)$ 越小越好; 一般可设定限制值为 0.05, 即当

$$\text{PRESS}(h)/\text{SS}(h-1) \leq (1-0.05)^2 = 0.95^2$$

时, 增加成分 t_h 有利于模型精度的提高。或者反过来说, 当

$$\text{PRESS}(h)/\text{SS}(h-1) > 0.95^2$$

时, 就认为增加新的成分 t_h , 对减少方程的预测误差无明显的改善作用。

为此, 定义交叉有效性为 $Q_h^2 = 1 - \text{PRESS}(h)/\text{SS}(h-1)$, 这样, 在建模的每一步计算结束前, 均进行交叉有效性检验, 如果在第 h 步有 $Q_h^2 < 1 - 0.95^2 = 0.0985$, 则模型达到精度要求, 可停止提取成分; 若 $Q_h^2 \geq 0.0975$, 表示第 h 步提取的 t_h 成分的边际贡献显著, 应继续第 $h+1$ 步计算。

§ 2 一种更简洁的计算方法

上节介绍的算法原则和推导过程的思路在目前的文献中是最为常见的。然而, 还有一种更为简洁的计算方法, 即直接在 E_0, \dots, E_{r-1} 矩阵中提取成分 $t_1, \dots, t_r (r \leq m)$ 。要求 t_h 能尽可能多地携带 X 中的信息, 同时, t_h 对因变量系统 F_0 有最大的解释能力。注意, 无需在 F_0 中提取成分 u_h , 这可以使计算过程大为简化, 并且对算法结论的解释也更为方便。

偏最小二乘法的简洁算法的步骤如下:

(1) 求矩阵 $E_0^T F_0 F_0^T E_0$ 最大特征值所对应的特征向量 w_1 , 求得成分得分向量 $\hat{t}_1 = E_0 w_1$, 和残差矩阵 $E_1 = E_0 - \hat{t}_1 \alpha_1^T$, 其中 $\alpha_1 = E_0^T \hat{t}_1 / \|\hat{t}_1\|^2$ 。

(2) 求矩阵 $E_1^T F_0 F_0^T E_1$ 最大特征值所对应的特征向量 w_2 , 求得成分得分向量 $\hat{t}_2 = E_1 w_2$, 和残差矩阵 $E_2 = E_1 - \hat{t}_2 \alpha_2^T$, 其中 $\alpha_2 = E_1^T \hat{t}_2 / \|\hat{t}_2\|^2$ 。

⋮

(r) 至第 r 步, 求矩阵 $E_{r-1}^T F_0 F_0^T E_{r-1}$ 最大特征值所对应的特征向量 w_r , 求得成分得分向量 $\hat{t}_r = E_{r-1} w_r$ 。

如果根据交叉有效性, 确定共抽取 r 个成分 t_1, \dots, t_r 可以得到一个满意的预测模型, 则求 F_0 在 $\hat{t}_1, \dots, \hat{t}_r$ 上的普通最小二乘回归方程为

$$F_0 = \hat{t}_1 \beta_1^T + \dots + \hat{t}_r \beta_r^T + F_r$$

把 $t_k = w_{k1}^* x_1 + \dots + w_{km}^* x_m$ ($k=1, 2, \dots, r$), 代入 $Y = t_1 \beta_1 + \dots + t_r \beta_r$, 即得 p 个因变量的偏最小二乘回归方程式

$$y_j = a_{j1} x_1 + \dots + a_{jm} x_m, \quad (j=1, 2, \dots, p)。$$

这里 $w_h^* = (w_{h1}^*, \dots, w_{hm}^*)^T$ 满足 $\hat{t}_h = E_0 w_h^*$, $w_h^* = \prod_{j=1}^{h-1} (I - w_j \alpha_j^T) w_h$ 。

§ 3 案例分析

本节采用兰纳胡德 (Linnerud) 给出的关于体能训练的数据进行偏最小二乘回归建模。在这个数据系统中被测的样本点, 是某健身俱乐部的 20 位中年男子。被测变量分为两组。第一组是身体特征指标 X , 包括: 体重、腰围、脉搏。第二组变量是训练结

果指标 Y ，包括：单杠、弯曲、跳高。原始数据见表 1。

表 2 给出了这 6 个变量的简单相关系数矩阵。从相关系数矩阵可以看出，体重与腰围是正相关的；体重、腰围与脉搏负相关；而在单杠、弯曲与跳高之间是正相关的。从两组变量间的关系看，单杠、弯曲和跳高的训练成绩与体重、腰围负相关，与脉搏正相关。

表 1 体能训练数据

No	体重(x_1)	腰围(x_2)	脉搏(x_3)	单杠(y_1)	弯曲(y_2)	跳高(y_3)
1	191	36	50	5	162	60
2	189	37	52	2	110	60
3	193	38	58	12	101	101
4	162	35	62	12	105	37
5	189	35	46	13	155	58
6	182	36	56	4	101	42
7	211	38	56	8	101	38
8	167	34	60	6	125	40
9	176	31	74	15	200	40
10	154	33	56	17	251	250
11	169	34	50	17	120	38
12	166	33	52	13	210	115
13	154	34	64	14	215	105
14	247	46	50	1	50	50
15	193	36	46	6	70	31
16	202	37	62	12	210	120
17	176	37	54	4	60	25
18	157	32	52	11	230	80
19	156	33	54	15	225	73
20	138	33	68	2	110	43
均值	178.6	35.4	56.1	9.45	145.55	70.3
标准差	24.6905	3.202	7.2104	5.2863	62.5666	51.2775

表 2 相关系数矩阵

	1	0.8702	-0.3658	-0.3897	-0.4931	-0.2263
	0.8702	1	-0.3529	-0.5522	-0.6456	-0.1915
	-0.3658	-0.3529	1	0.1506	0.225	0.0349
	-0.3897	-0.5522	0.1506	1	0.6957	0.4958
	-0.4931	-0.6456	0.225	0.6957	1	0.6692
	-0.2263	-0.1915	0.0349	0.4958	0.6692	1

利用如下的 MATLAB 程序：

```

clc,clear
load pz.txt %原始数据存放在纯文本文件 pz.txt 中
mu=mean(pz);sig=std(pz); %求均值和标准差
rr=corrcoef(pz); %求相关系数矩阵
data=zscore(pz); %数据标准化,变量记做 X*和 Y*
n=3;m=3; %n 是自变量的个数,m 是因变量的个数
x0=pz(:,1:n);y0=pz(:,n+1:end); %原始的自变量和因变量数据
e0=data(:,1:n);f0=data(:,n+1:end); %标准化后的自变量和因变量数据

```

```

num=size(e0,1);%求样本点的个数
chg=eye(n); %w 到 w*变换矩阵的初始化
for i=1:n
%以下计算 w, w*和 t 的得分向量,
    matrix=e0'*f0*f0'*e0;
    [vec,val]=eig(matrix); %求特征值和特征向量
    val=diag(val); %提出对角线元素,即提出特征值
    [val,ind]=sort(val,'descend');
    w(:,i)=vec(:,ind(1)); %提出最大特征值对应的特征向量
    w_star(:,i)=chg*w(:,i); %计算 w*的取值
    t(:,i)=e0*w(:,i); %计算成分 ti 的得分
    alpha=e0*t(:,i)/(t(:,i)'*t(:,i)); %计算 alphai
    chg=chg*(eye(n)-w(:,i)*alpha'); %计算 w 到 w*的变换矩阵
    e=e0-t(:,i)*alpha'; %计算残差矩阵
    e0=e;
%以下计算 ss(i)的值
    beta=t\f0; %求回归方程的系数,数据标准化,没有常数项
    cancha=f0-t*beta; %求残差矩阵
    ss(i)=sum(sum(cancha.^2)); %求误差平方和
%以下计算 press(i)
    for j=1:num
        t1=t(:,1:i);f1=f0;
        she_t=t1(j,:);she_f=f1(j,:); %把舍去的第 j 个样本点保存起来
        t1(j,:)=[];f1(j,:)=[]; %删除第 j 个观测值
        beta1=[t1,ones(num-1,1)]\f1; %求回归分析的系数,这里带有常数项
        cancha=she_f-she_t*beta1(1:end-1,:)-beta1(end,:); %求残差向量
        press_i(j)=sum(cancha.^2); %求误差平方和
    end
    press(i)=sum(press_i);
    Q_h2(1)=1;
    if i>1, Q_h2(i)=1-press(i)/ss(i-1); end
    if Q_h2(i)<0.0975
        fprintf('提出的成分个数 r=%d',i); break
    end
end
end
beta_z=t\f0; %求 Y*关于 t 的回归系数
xishu=w_star*beta_z; %求 Y*关于 X*的回归系数,每一列是一个回归方程
mu_x=mu(1:n);mu_y=mu(n+1:end); %提出自变量和因变量的均值
sig_x=sig(1:n);sig_y=sig(n+1:end); %提出自变量和因变量的标准差
ch0=mu_y-(mu_x./sig_x*xishu).*sig_y; %计算原始数据回归方程的常数项
for i=1:m
    xish(:,i)=xishu(:,i)./sig_x'*sig_y(i); %计算原始数据回归方程的系数
end
sol=[ch0;xish] %显示回归方程的系数,每一列是一个方程,每一列的第一个数是常数项
save mydata x0 y0 num xishu ch0 xish

```

计算得只要提出两个成分 t_1, t_2 即可，交叉有效性 $Q_2^2 = -0.3263$ 。 w_h 与 w_h^* 的取值见表 3，成分 t_h 的得分 \hat{t}_h 见表 4。

表 3 w_h 与 w_h^* 的取值

自变量	w_1	w_2	w_1^*	w_2^*
x_1	-0.5899	-0.4688	-0.5899	-0.3679
x_2	-0.7713	0.5680	-0.7713	0.6999
x_3	0.2389	0.6765	0.2389	0.6356

表 4 成分 t_h 的得分 \hat{t}_h

No	1	2	3	4	5	6	7	8	9	10
\hat{t}_1	-0.6429	-0.7697	-0.9074	0.6884	-0.4867	-0.2291	-1.4037	0.7436	1.7151	1.1626
\hat{t}_2	-0.5914	-0.1667	0.5212	0.68	-1.1328	0.0717	0.0767	0.2106	0.6549	-0.1668
No	11	12	13	14	15	16	17	18	19	20
\hat{t}_1	0.3645	0.7433	1.1867	-4.3898	-0.8232	-0.749	-0.3929	1.1993	1.0485	1.9424
\hat{t}_2	-0.7007	-0.6983	0.757	0.76	-0.9738	0.5211	0.2034	-0.7827	-0.3729	1.1294

标准化变量 \tilde{y}_k 关于成分 t_1, t_2 的回归模型如下

$$\tilde{y}_k = r_{1k}t_1 + r_{2k}t_2, \quad k = 1, 2, 3$$

由于成分 t_h 可以写成原变量的标准化变量 \tilde{x}_j 的函数，即有

$$t_h = w_{1h}^*\tilde{x}_1 + w_{2h}^*\tilde{x}_2 + w_{3h}^*\tilde{x}_3, \quad h = 1, 2$$

由此可得由成分 t_1, t_2 所建立的偏最小二乘回归模型为

$$\begin{aligned} \tilde{y}_k &= r_{1k}(w_{11}^*\tilde{x}_1 + w_{21}^*\tilde{x}_2 + w_{31}^*\tilde{x}_3) + r_{2k}(w_{12}^*\tilde{x}_1 + w_{22}^*\tilde{x}_2 + w_{32}^*\tilde{x}_3) \\ &= (r_{1k}w_{11}^* + r_{2k}w_{12}^*)\tilde{x}_1 + (r_{1k}w_{21}^* + r_{2k}w_{22}^*)\tilde{x}_2 + (r_{1k}w_{31}^* + r_{2k}w_{32}^*)\tilde{x}_3 \end{aligned}$$

有关 $r_h = (r_{h1}, r_{h2}, r_{h3})$ 的计算结果见表 5。

表 5 回归系数 r_h

k	1	2	3
r_1	0.3416	0.4161	0.1430
r_2	-0.3364	-0.2908	-0.0652

所以，有

$$\tilde{y}_1 = -0.0778\tilde{x}_1 - 0.4989\tilde{x}_2 - 0.1322\tilde{x}_3$$

$$\tilde{y}_2 = -0.1385\tilde{x}_1 - 0.5244\tilde{x}_2 - 0.0854\tilde{x}_3$$

$$\tilde{y}_3 = -0.0604\tilde{x}_1 - 0.1559\tilde{x}_2 - 0.0073\tilde{x}_3$$

将标准化变量 $\tilde{y}_k, \tilde{x}_k (k = 1, 2, 3)$ 分别还原成原始变量 $y_k, x_k (k = 1, 2, 3)$ ，则回归方程为

$$y_1 = 47.0197 - 0.0167x_1 - 0.8237x_2 - 0.0969x_3$$

$$y_2 = 612.5671 - 0.3509x_1 - 10.2477x_2 - 0.7412x_3$$

$$y_3 = 183.9849 - 0.1253x_1 - 2.4969x_2 - 0.0518x_3$$

为了更直观、迅速地观察各个自变量在解释 $y_k (k=1,2,3)$ 时的边际作用，可以绘制回归系数图，见图 1。这个图是针对标准化数据的回归方程的。

从回归系数图中可以立刻观察到，腰围变量在解释三个回归方程时起到了极为重要的作用。然而，与单杠及弯曲相比，跳高成绩的回归方程显然不够理想，三个自变量对它的解释能力均很低。

为了考察这三个回归方程的模型精度，我们以 (\hat{y}_{ik}, y_{ik}) 为坐标值，对所有的样本点绘制预测图。 \hat{y}_{ik} 是第 k 个变量，第 i 个样本点 (y_{ik}) 的预测值。在这个预测图上，如果所有点都能在图的对角线附近均匀分布，则方程的拟合值与原值差异很小，这个方程的拟合效果就是满意的。体能训练的预测图见图 2。

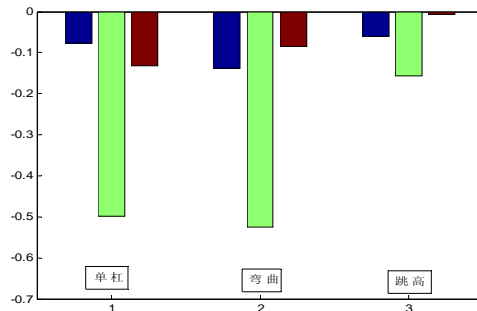


图 1 回归系数的直方图

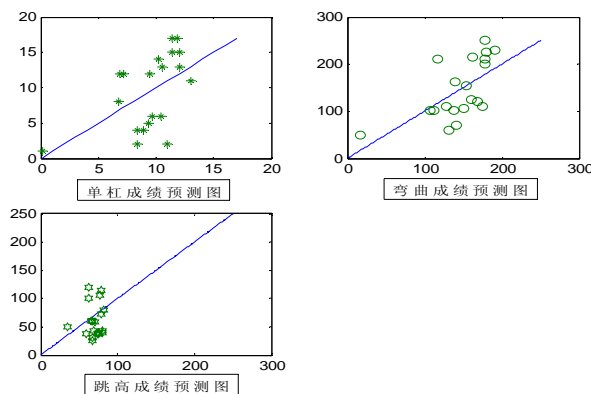


图 2 体能训练预测图

画直方图的 MATLAB 程序为：bar(xishu')

画体能训练的预测图的 MATLAB 程序如下：

```
load mydata
ch0=repmat(ch0,num,1);
yhat=ch0+x0*xish; %计算 Y 的预测值
y1max=max(yhat); %求预测值的最大值
y2max=max(y0); %求观测值的最大值
ymax=max([y1max;y2max]) %求预测值和观测值的最大值
cancha=yhat-y0; %计算残差
subplot(2,2,1)
```



```
%画直线 y=x,并画预测图
plot(0:ymax(1),0:ymax(1),yhat(:,1),y0(:,1),'*')
subplot(2,2,2)
plot(0:ymax(2),0:ymax(2),yhat(:,2),y0(:,2),'O')
subplot(2,2,3)
plot(0:ymax(3),0:ymax(3),yhat(:,3),y0(:,3),'H')
```

习题二十七

1. (化工试验例子) 考察的指标 (因变量) y 表示原辛烷值, 自变量 x_1 表示直接蒸馏成分, x_2 表示重整汽油, x_3 表示原油热裂化油, x_4 表示原油催化裂化油, x_5 表示聚合物, x_6 表示烷基化物, x_7 表示天然香精。7 个变量表示 7 个成分含量的比例 (满足 $x_1 + x_2 + \cdots + x_7 = 1$)。表 6 给出 12 种混合物中 7 种成分和 y 的数据。试用偏最小二乘方法建立 y 与 x_1, x_2, \cdots, x_7 的回归方程, 用于确定 7 种构成元素 x_1, x_2, \cdots, x_7 对 y 的影响。

表 6 化工试验的原始数据

序号	x_1	x_2	x_3	x_4	x_5	x_6	x_7	y
1	0	0.23	0	0	0	0.74	0.03	98.7
2	0	0.1	0	0	0.12	0.74	0.04	97.8
3	0	0	0	0.1	0.12	0.74	0.04	96.6
4	0	0.49	0	0	0.12	0.37	0.02	92.0
5	0	0	0	0.62	0.12	0.18	0.08	86.6
6	0	0.62	0	0	0	0.37	0.01	91.2
7	0.17	0.27	0.1	0.38	0	0	0.08	81.9
8	0.17	0.19	0.1	0.38	0.02	0.06	0.08	83.1
9	0.17	0.21	0.1	0.38	0	0.06	0.08	82.4
10	0.17	0.15	0.1	0.38	0.02	0.1	0.08	83.2
11	0.21	0.36	0.12	0.25	0	0	0.06	81.4
12	0	0	0	0.55	0	0.37	0.08	88.1

2. 试对表 7 的 38 名学生的体质和运动能力数据, 用偏最小二乘法建立 5 个运动能力指标与 7 个体质变量的回归方程。

表 7 学生体质与运动能力数据

序号	体 质 情 况							运 动 能 力				
	x_1	x_2	x_3	x_4	x_5	x_6	x_7	y_1	y_2	y_3	y_4	y_5
1	46	55	126	51	75.0	25	72	6.8	489	27	8	360
2	52	55	95	42	81.2	18	50	7.2	464	30	5	348
3	46	69	107	38	98.0	18	74	6.8	430	32	9	386
4	49	50	105	48	97.6	16	60	6.8	362	26	6	331
5	42	55	90	46	66.5	2	68	7.2	453	23	11	391
6	48	61	106	43	78.0	25	58	7.0	405	29	7	389
7	49	60	100	49	90.6	15	60	7.0	420	21	10	379
8	48	63	122	52	56.0	17	68	7.0	466	28	2	362
9	45	55	105	48	76.0	15	61	6.8	415	24	6	386
10	48	64	120	38	60.2	20	62	7.0	413	28	7	398
11	49	52	100	42	53.4	6	42	7.4	404	23	6	400

12	47	62	100	34	61.2	10	62	7.2	427	25	7	407
13	41	51	101	53	62.4	5	60	8.0	372	25	3	409
14	52	55	125	43	86.3	5	62	6.8	496	30	10	350
15	45	52	94	50	51.4	20	65	7.6	394	24	3	399
16	49	57	110	47	72.3	19	45	7.0	446	30	11	337
17	53	65	112	47	90.4	15	75	6.6	420	30	12	357
18	47	57	95	47	72.3	9	64	6.6	447	25	4	447
19	48	60	120	47	86.4	12	62	6.8	398	28	11	381
20	49	55	113	41	84.1	15	60	7.0	398	27	4	387
21	48	69	128	42	47.9	20	63	7.0	485	30	7	350
22	42	57	122	46	54.2	15	63	7.2	400	28	6	388
23	54	64	155	51	71.4	19	61	6.9	511	33	12	298
24	53	63	120	42	56.6	8	53	7.5	430	29	4	353
25	42	71	138	44	65.2	17	55	7.0	487	29	9	370
26	46	66	120	45	62.2	22	68	7.4	470	28	7	360
27	45	56	91	29	66.2	18	51	7.9	380	26	5	358
28	50	60	120	42	56.6	8	57	6.8	460	32	5	348
29	42	51	126	50	50.0	13	57	7.7	398	27	2	383
30	48	50	115	41	52.9	6	39	7.4	415	28	6	314
31	42	52	140	48	56.3	15	60	6.9	470	27	11	348
32	48	67	105	39	69.2	23	60	7.6	450	28	10	326
33	49	74	151	49	54.2	20	58	7.0	500	30	12	330
34	47	55	113	40	71.4	19	64	7.6	410	29	7	331
35	49	74	120	53	54.5	22	59	6.9	500	33	21	348
36	44	52	110	37	54.9	14	57	7.5	400	29	2	421
37	52	66	130	47	45.9	14	45	6.8	505	28	11	355
38	48	68	100	45	53.6	23	70	7.2	522	28	9	352