

复杂系统建模——高维特征空间变量法

韩建勋 饶欣

(天津大学化工系)

摘 要

本文运用模式识别技术提出了一种分析和选择复杂系统变量的方法。运用该方法使得系统的建模和控制简化。本文以复杂碳化过程作为实例，用模式识别技术对其内在规律进行分析，找出了温度分布与结晶质量的定量关系。在此基础上，建立了碳化过程的动态模型。实践表明，该模型与实际过程基本吻合，并已取得明显经济效益。

关键词：模式识别，系统识别，数学模型，碳化过程。

一、复杂系统建模中高维中间特征变量的提取

1. 复杂系统建模中的问题

任何一个受控系统都有一个控制集 $U = \{u_1, \dots, u_m\}$ 和一个目标集 $S = \{s_1, \dots, s_n\}$ 。此外，由系统的其它变量及其各种变换和组合

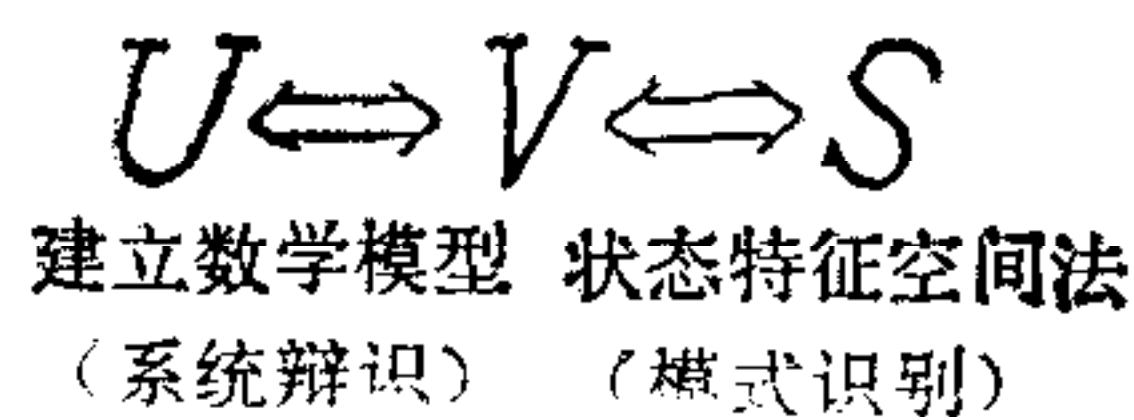


图 1 中间状态特征的引入

组成一个状态变量集 $X = \{x_1, \dots, x_k\}$ 。要实现受控系统的优化控制，需要建立 U 与 S 之间相应的数学模型，但由于实际问题的复杂性，直接建立 U 与 S 间的数学模型往往困难很大。若能找到一组中间状态特征变量 $V = (v_1, \dots$

$\dots, v_l)$ ($v_i \in X, i = 1, 2, \dots, l$)，使它与 U 之间有比较明显的函数关系，则相对 U 与 S 来说就比较容易建立数学模型。同时，它又是目标集 S 的状态特征，因此能与 S 建立一种对应关系，如图 1 所示。这样不仅可以抓住受控系统的本质特征，而且可以使建模工作简化。其关键问题是如何找出这样一组中间状态特征变量，以及如何建立 V 与 S 之间的对应关系。关键是模式识别中特征空间的提取与选择。

2. 高维中间特征变量的提取与选择

根据具体问题对目标集 S 的要求，将目标空间 S 中的点分为 p 个互不相交的子集 S_1, \dots, S_p ，并且满足

$$S = \bigcup_{i=1}^p S_i,$$

$$S_i \cap S_j = \phi, \quad i \neq j, \quad i, j = 1, \dots, p.$$

其中 ϕ 代表空集.

设 V 是由系统中间状态特征变量组成的空间. 再依据 S 空间的分类, 对 V 进行相应的分类. 对任意 $s(k) \in S_i$ (k 采样时刻), 则标记 $v(k)$, $v(k) \in V_i$, 即将 $v(k)$ 归为 V_i 类. 这样 V 也被分为 V_1, \dots, V_p 类. 如果 V_1, \dots, V_p 各类有较好的聚类, 彼此之间又有较好的分类, 那么, 就能建立 V 与 S 之间一种区域对应关系. 这种通过 V 达到对 S 识别的可识别性问题其关键是 V 中特征变量 v_1, \dots, v_l 的提取与选择, 即要使 V_1, \dots, V_p 在特征空间中达到最佳的分类效果. 因此需要定义一种可分性判据 $J_{ij}(v_1, \dots, v_l)$. 它描述了特征向量 $V = (v_1, \dots, v_l)$ 对 V_i 与 V_j 两类的可分性, 还要满足以下四条性质:

1) 可分性判据 J_{ij} 与 V_i 和 V_j 两类的误识率 $P(i|j)$ 和 $P(j|i)$ (或其上、下界) 有单调关系;

2) 当 v_1, \dots, v_l 彼此独立时, 可分性判据 J_{ij} 具有可加性

$$J_{ij}(v_1, \dots, v_l) = \sum_{k=1}^l J_{ij}(v_k);$$

3) 可分性判据具有度量性

①非负性 $J_{ij} \geq 0$ 当且仅当 $i = j$ 时, $J_{ij} = 0$.

②对称性 $J_{ij} = J_{ji}$.

4) 单调性

$$J_{ij}(v_1, \dots, v_l) \leq J_{ij}(v_1, \dots, v_l, v_{l+1}).$$

为了描述 V_1, \dots, V_p 各类在特征空间 V 中总体可分性, 还需要定义一个总体目标函数 $J(J_{11}, J_{12}, \dots, J_{1p}, \dots, J_{(p-1)p})$, 一般可选 $J = \sum_{i,j=1}^p J_{ij}$. 显然 J 也满足上述四条性质.

因此对特征变量 v_1, \dots, v_l 的提取与选择就要使相应的总体目标函数 J 取极值

$$J^* = \max_{(v_1, \dots, v_l) \in X} (\min) J.$$

一种比较直观的可分性判据是根据样本在特征空间中距离概念——类内距 S_w 和类间距 S_b 定义的.

$$\text{其中 } S_w = \sum_{i=1}^p p_i E[(v - m_i)(v - m_i)^T | V_i],$$

$$S_b = \sum_{i=1}^p p_i (m_i - m_0)(m_i - m_0)^T.$$

$v \in V_i$, p_i , m_i 及 m_0 分别是 V_i 类的先验概率, 均值向量和总体均值向量. 根据实际问题的不同要求可构造和选择不同的总体目标函数如:

$$J_0 = |S_w + S_b|,$$

$$J_1 = \text{tr}(S_w + S_b),$$

$$J_2 = \ln \left(\frac{|S_b|}{|S_w|} \right),$$

$$J_3 = \text{tr}(S_w^{-1} S_b),$$

$$J_4 = \frac{\text{tr} S_b}{\text{tr} S_w},$$

$$J_5 = \frac{|S_w + S_b|}{|S_w|}.$$

由于多类问题特征变量的提取和选择与样本概率分布有关,因此在样本的概率分布不是正态或是未知情况下, $K-L$ 正交变换是一种比较实用的方法,它可以在不直接考虑样本概率分布情况下,进行特征提取与选择.

3. 特征提取的 $K-L$ 正交变换法^[2]

设特征空间中有一组正交归一的向量 $\Phi = (\phi_1, \dots, \phi_l)^T$, 它将特征空间 $V^{(l)}$ 中的特征向量 $v = (v_1, \dots, v_l)^T$ 变换成一组新的特征向量 $w = (w_1, \dots, w_l)^T$, $w = \Phi v$. 由 Φ 的正交归一性有 $\Phi\Phi^T = \Phi^T\Phi = I_{l \times l}$, 所以有

$$v = \Phi^T w. \quad (1)$$

特征向量组成的自相关阵 R 为

$$R = \sum_{i=1}^p P(v_i) E\{v_i v_i^T\}. \quad (2)$$

其中 $P(v_i)$, v_i 分别是 V_i 类的先验概率, V_i 类所有样本组成的样本矩阵将 (1) 式代入 (2) 式, 则

$$R = \sum_{i=1}^p P(v_i) E\{\Phi^T w_i w_i^T \Phi\} = \Phi^T \left[\sum_{i=1}^p P(V_i) E\{w_i w_i^T\} \right] \Phi,$$

其中 w_i 是 v_i 经 Φ 变换以后新的特征向量 w 组成的样本矩阵. 由于 w_i 具有正交性质, 所以

$$w_i^T w_j = \begin{cases} 1, & \text{如 } i = j, \\ 0, & \text{如 } i \neq j. \end{cases}$$

由上可知

$$\left[\sum_{i,j=1}^p p(v_i) E\{w_i w_j\} \right] = D_\lambda.$$

D_λ 为一对角阵

$$D_\lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots \\ & & & \lambda_l \end{pmatrix},$$

所以 $R = \Phi^T D_\lambda \Phi$, 即

$$R\Phi^T = \Phi^T D_\lambda. \quad (3)$$

由 (3) 式可以看出, 对于 p 类样本组成的特征向量自相关阵 R , D_λ 的对角元素是 R 的特征值 λ_j 组成的. $\lambda_j (j = 1, \dots, l)$ 相应的各特征向量 ϕ_j , 构成 Φ 阵.

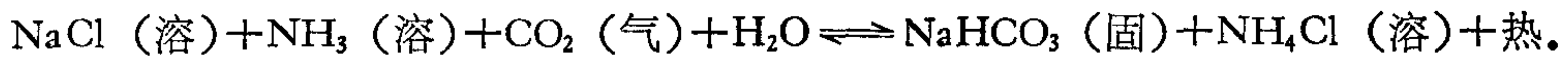
类同, 式 (3) 中的 R 也可换为 J_0, J_1, \dots, J_5 等形式. 现在可以得到特征提取与选择的原则, 就是欲达到对 $V_i, (i = 1, \dots, p)$ 各类的最佳分类, 应通过对总体目标函数特征值 λ 的选取找出相应的正交特征向量 ϕ_j , 并组成 Φ 阵, 使变换后新的样本的类内距 S'_w 尽可能的小, 类间距 S'_b 尽可能大.

在完成特征变量的选择后, 可以建立 V 与 S 间的一种区域对应关系, 再针对所选择的特征变量, 建立 U 与 V 之间的数学模型.

二、实例: 碳化复杂过程的建模

1. 碳化过程特征变量——温度分布

纯碱 (Na_2CO_3) 的生产是一个复杂过程, 其中氨盐水碳酸化又是整个生产的中心环节, 其化学反应方程式^[3]为



整个过程在碳化塔内进行。这是一个气、液、固三相物系的吸收、反应、传热和结晶同时进行的复杂过程。由于“制碱”和“清洗”工况的切换, 塔内气液相分布不均和多塔耦合等现象, 导致系统流体力学的非平稳严重, 使系统具有严重的非线性和时变性。因此, 它是一个多变量、多干扰、多不确定因素, 既连续又间歇, 多目标的分布参数受控系统。

若要直接建立控制集 $U = \{\text{中段气量, 下段气量, 进卤量, \dots}\}$ 与目标集 $S = \{\text{NaHCO}_3 \text{ 结晶质量, NaCl 转化率, 制碱周期, \dots}\}$ 之间的数学模型有较大困难。因此需要寻找中间状态变量。由机理分析得知 NaHCO_3 结晶过程是关键, 它不仅比 NaCl 转化率更重要, 而且它的改善有利于延长制碱周期, 此外, 塔内两个主要状态变量 CO_2 浓度分布和碳化卤的温度分布都与 NaHCO_3 结晶过程有密切关系。但检测 CO_2 浓度分布没有在线手段, 甚至查定时也难用人工分析获取, 为此选择温度为中间状态变量。由于塔内吸收和结晶是连续变化过程, 若只根据塔上某点温度控制整个碳化塔, 将有严重缺陷。因此, 必须沿塔高的温度分布进行控制, 通过 5 个实测点的温度值, 利用有二阶导数的三次样条插值函数, 得到各采样时刻温度沿塔高分布的插值函数 $S(x)$ 。

$$\begin{aligned} S(x) = & (x - x_j)(x - x_{j+1})(x - 2x_j + x_{j+1})M_{j+1}/6(x_{j+1} - x_j) \\ & - (x - x_j)(x - x_{j+1})(x - x_j - 2x_{j+1})M_j/6(x_{j+1} + x_j) \\ & + [(x - x_j)f_{j+1} - (x - x_{j+1})f_j]/(x_{j+1} - x_j). \end{aligned}$$

其中 x_j , f_j 和 M_j 分别是第 j 检测点的塔高, 温度的实测值和 $S(x)$ 在 x_j 处的二阶导数值。

2. 特征变量的提取与选择

利用前面的方法可找出了碳化过程重要目标 (NaHCO_3 结晶质量) 与特征变量 (温度分布) 的对应关系, 并依此选出了特征分量。这些参数为建模提供了中间变量, 也为控制系统的设计提供依据。

结晶质量用 NaHCO_3 固体沉降时间 z 表示, 并由它组成目标集 Z 。根据实际要求将 Z 分为好、中、差三个不相交的子集。

$$Z_1 = \{z(k) | z(k) \leq 200 \text{ 秒}, k = 1, \dots, n_1\} \text{——好子集}$$

$$Z_2 = \{z(k) | 200 < z(k) \leq 230 \text{ 秒}, k = 1, \dots, n_2\} \text{——中子集}$$

$$Z_3 = \{z(k) | z(k) > 230 \text{ 秒}, k = 1, \dots, n_3\} \text{——差子集}$$

特征向量由温度插值函数 $S(x)$ 上的不同点的温度组成, 并组成特征空间 V , 可分性判据的总体目标函数为

$$J_0 = |S_w + S_b|.$$

运用 $K-L$ 正交变换法, 最后选择了温度分布上的 6 个点组成特征向量

$$v_k = [S_{29}(k), S_{25}(k), S_{20}(k), S_{17}(k), S_{10}(k), S_1(k)]^T.$$

其中 $S_i(k)$ 代表 k 采样时刻温度分布插值 $S(x)$ 上的第 i 点的值, $i = 29, 25, 20, 17, 10, 1$. 它构成一个 6 维特征空间 $V^{(6)}$. 把对温度分布函数的研究转化为对 $V^{(6)}$ 的研究, 依据 Z 的分类对 $V^{(6)}$ 分类, 若 $z(k) \in Z_i$, 则令 $v(k) \in V_i$ 类. 利用 70 组数据进行计算, 进一步将 $V^{(6)}$ 中点按最佳可分性, 投影到直观的二维空间(平面)上, 如图 2(a) 所示. 图中 1, 2, 3 代表的好, 中, 差三类有明显的聚类, 各类之间又有明显的分类. 这说明特征变量温度分布确实与结晶质量有对应关系.

由于对 6 点的温度控制难以实现, 进一步用 $K-L$ 特征提取, 得 3 个主要特征, 从而将原 6 维特征空间 $V^{(6)}$ 降至 3 维, 其特征向量为

$$v'_k = [S_{25}(k), S_{20}(k), S_{17}(k)]^T.$$

相似地将 $V^{(3)}$ 投影到直观的二维空间上, 从图 2(b) 可看出它基本保持了原 $V^{(6)}$ 的分布结构.

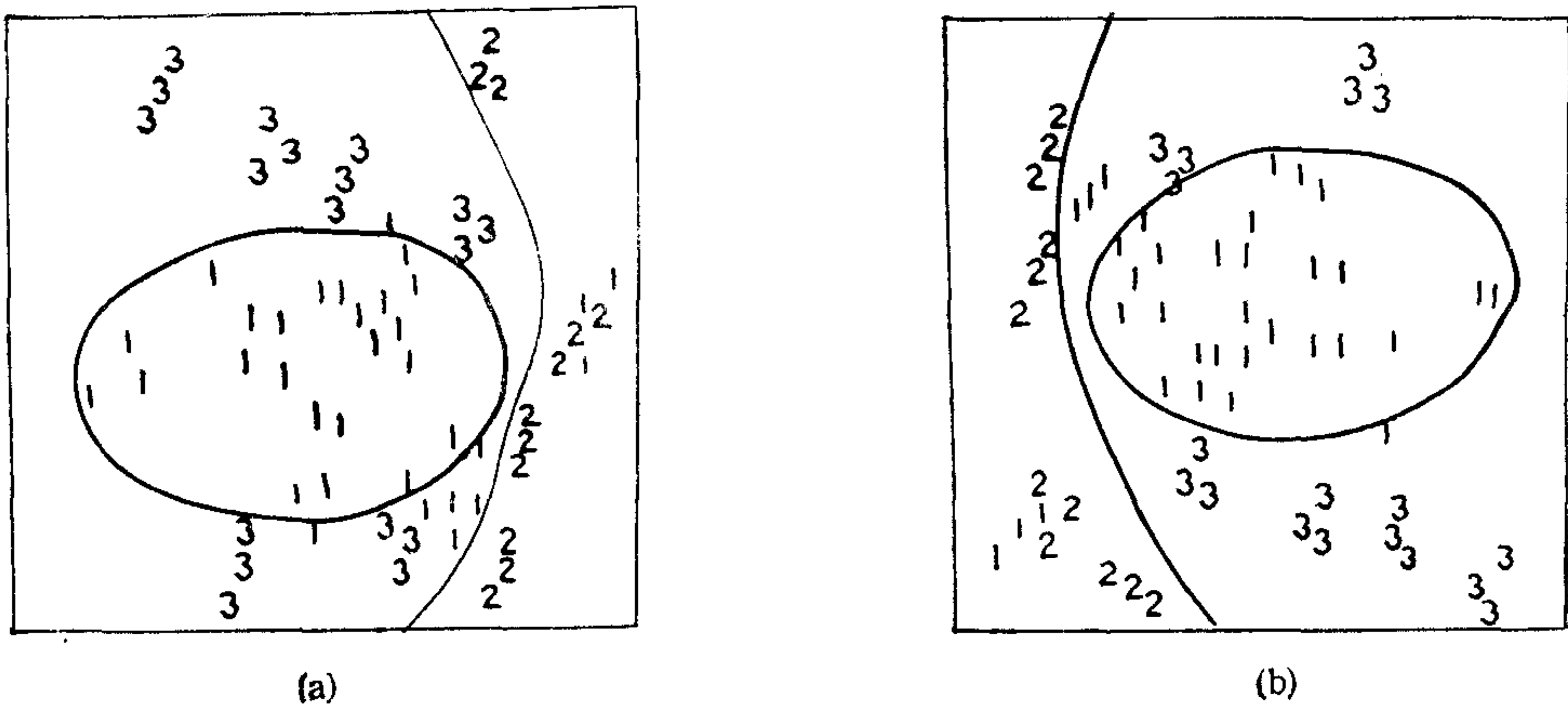


图 2 高维特征向量在二维空间上的分布

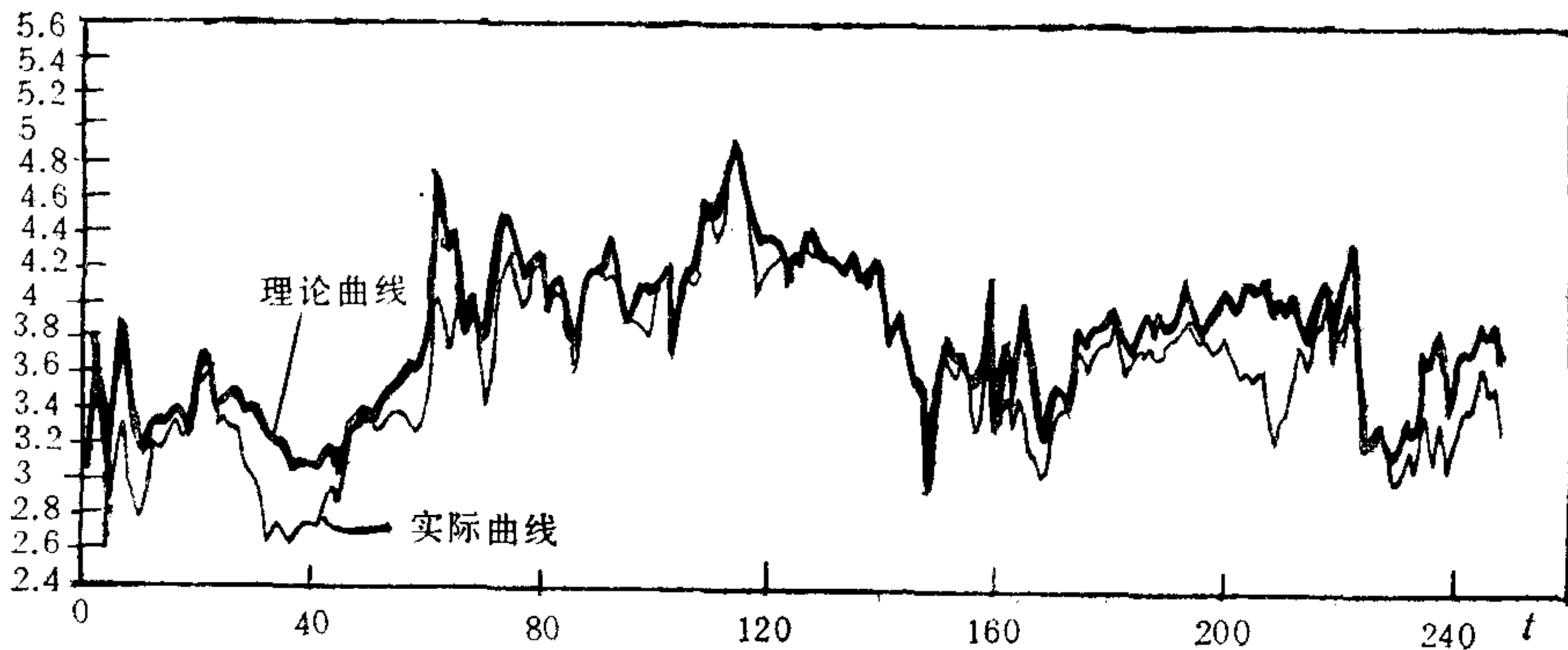


图 3 25 圈实际与理论曲线比较

3. 碳化过程多输入多输出动态模型

针对选择和提取得到的 3 个特征分量, 从机理的热平衡方程出发, 结合系统辨识, 确定了模型的控制通道和模型的 V 型结构, 建立了多输入多输出动态模型

$$X(k+1) = AX(k) + B_1U(k) + B_2U(k-1).$$

其中 X 为 3 维的输出变量(温度分布的特征变量), U 为 3 维输入变量(控制变量). A, B_1, B_2 为系数阵. 采用递推增广最小二乘法和采集的 150 组数据, 对模型进行参数估计, 并用 250 组数据进行了验证. 文中仅给出 25 圈例子(见图 3). 这是 3 个输入通道共同产生的输出, 可以看出在系统剧烈变化时模型的理论曲线与实际曲线仍然吻合较好.

三、小 结

对复杂系统, 本文通过高维特征空间的分析, 提出用中间特征状态变量进行建模, 并给出了提取与选择中间特征状态变量的一种 $K-L$ 正交变换法. 这样不仅能够抓住系统的主要特征, 使建模工作简化, 而且为控制系统的设计创造了条件.

参 考 文 献

- [1] Johson, A. F., Khaligh, B., Dynamic Parameter Estimation by Mean of Pattern Recognition, Proc. IAS-TED Symposium, ACI'83, Copenhege, 1(1983).
- [2] Batchelor, Bruce G., ed., Pattern Recognition Ideas in Practice, New York, 1987.
- [3] 兰特, Z., 索尔维法制碱, 化学工业出版社, 1983.

MODELLING COMPLEX SYSTEM — A METHOD OF HIGH DIMENSION FEATURE SPACE

Han Jianxun Rao Xin

(Tianjin University)

ABSTRACT

This paper proposes a method for the analysis and selection of complicated system variables by applying pattern recognition technique. By means of this method, system modelling and control scheme can be simplified. As the demonstration of the method, pattern recognition is first applied to analysing the inherent laws of the complicated carbonation process. The quantity relation between the crystal quality and the temperature distribution has been found. On this basis, carbonation mechanism and system identification have been combined to establish the dynamic model of the carbonation process successfully.

Key words: Pattern recognition; system identification; mathematical modelling; carbonation process.