文章编号:1673-5005(2008)01-0147-05

基于鲁棒主元分析的故障诊断方法

邓晓刚,田学民

(中国石油大学 信息与控制工程学院,山东 东营 257061)

摘要:针对工业过程的建模数据中含有离群点的情况,提出了一种基于鲁棒主元分析(PCA)的故障诊断方法。该方 法使用广义极大似然估计(M估计)代替最小二乘估计,将传统的主元分析问题转化为一个加权的重构误差优化问 题,然后通过改进的非线性迭代部分最小二乘(NIPALS)算法来求得问题的最优解,在此基础上建立主元模型并构造 监控统计量检测过程故障。在连续搅拌反应器(CSTR)仿真系统上的应用结果表明,鲁棒 PCA 方法能够消除离群点 对主元模型的影响,比 PCA 方法分析过程数据更为准确,能更有效地诊断过程故障。 关键词:故障诊断;鲁棒主元分析;离群点; NIPALS 算法; M估计

中图分类号:TP 277 文献标识码:A

Fault diagnosis method based on robust principal component analysis

DENG Xiao-gang, TIAN Xue-min

(College of Information and Control Engineering in China University of Petroleum, Dongying 257061, Shandong Province, China)

Abstract: A robust principal component analysis (PCA) method was proposed to analyze the model data with outliers in process monitoring. By replacing the least squares estimator with a robust M-estimator, the traditional principal component analysis problem was transformed into a weighted reconstructed error optimization problem. The problem can be solved by improved nonlinear iterative partial least squares (NIPALS) algorithm so that precise principal component model was available and monitoring statistics were used to detect faults. The simulation results on a continuous stirred tank reactor (CSTR) system show that the proposed robust PCA method can remove the influence of outliers, analyze the process data more accurately and diagnose process faults more effectively than traditional PCA method.

Key words: fault diagnosis; robust principal component analysis; outliers; NIPALS algorithm; M estimator

基于数据分析的故障诊断和过程监控方法是近 年来该领域的一个热点问题,目前使用比较成熟的 是主元分析(PCA)方法。该方法要求建模数据的噪 声服从正态分布,离群数据的存在会严重影响模型 的准确性。然而在实际生产过程中,由于传感器故 障和过程偶然波动等原因,系统记录的数据中往往 会有离群点存在。为克服传统 PCA 的不足,研究人 员提出了一系列的鲁棒 PCA 方法^[1],主要可以分为 鲁棒协方差分解算法和投影寻踪算法。鲁棒协方差 分解方法使用鲁棒估计器如最小协方差行列式值 (MCD)来计算鲁棒的协方差阵,使得 PCA 分析具有 鲁棒性^[2],但是该算法要求样本数目大于变量数 目。投影寻踪算法则是通过对优化目标的鲁棒化处 理来消除离群点的影响^[34],从而能够准确地计算出 最优解。Xie 等提出一种基于粒子群寻优算法的投 影寻踪方法^[5],Yang 等分析了基于模糊处理的鲁棒 PCA 方法^[67]。笔者结合广义极大似然估计(M 估 计)和非线性迭代部分最小二乘(NIPALS)算法,提 出一种新的鲁棒 PCA 方法。

1 传统 PCA 算法与鲁棒 PCA 算法

1.1 传统 PCA 方法

假设一段正常工况下的过程历史数据为 X ∈

收稿日期:2007-06-10

基金项目:国家"863"计划项目(2004AA412050)

作者简介:邓晓刚(1981-),男(汉族),山东广饶人,博士研究生,主要从事故障诊断、过程监控等方面的工作。

 $\mathbf{R}^{n\times m}$,描述 m 个变量的 n 次采样值。通过 PCA 分析,可 以用 k(k < m) 个不相关的新变量来描述数据的主要 信息,从而达到在低维空间中分析高维数据的目的。

PCA 分析将 X 分解为主元空间数据 X 与残差空间数据 E 之和,即

$$X = \sum_{i=1}^{n} t_i p_i^{\mathsf{T}} + \sum_{i=k+1}^{m} t_i p_i^{\mathsf{T}} = T_k P_k^{\mathsf{T}} + E = \hat{X} + E. \quad (1)$$

式中,t_i为得分向量;p_i为载荷向量。

假设残差数据服从正态分布,式(1)的求解是 一个基于最小二乘估计的优化问题,即

$$\min \sum_{i=1}^{n} \|\boldsymbol{e}_{i}\|^{2} = \sum_{i=1}^{n} \sum_{j=1}^{m} e_{ij}^{2}, \qquad (2)$$

其中

 $e_i = x_i - \hat{x}_i = x_i - x_i p_i p_j^{\mathsf{T}}$, 优化的过程是为了依次求解载荷向量 p_j 。

1.2 鲁棒 PCA 算法

从式(2)可以看出,传统 PCA 方法是基于信号 重构误差最小准则的优化过程,数据中的离群点会 严重影响优化问题的结果。为增强 PCA 方法的鲁棒 性,使用鲁棒性强的 M 估计来定义上述问题^[89],即

$$\min \sum_{i=1}^{n} \sum_{i=1}^{m} \rho(e_{ij}) .$$
 (3)

如果 ρ 是一个凸函数,并且 $\psi(e) = \rho'(e)$ 是连续的,由式(3)首先来解得 p_1 ,依次得

$$\sum_{i=1}^{n} \sum_{j=1}^{m} \psi(e_{ij}) \frac{\mathrm{d}e_{ij}}{\mathrm{d}p_{1}} = 0, \qquad (4)$$

$$\sum_{i=1}^{n} \sum_{j=1}^{m} \frac{\psi(e_{ij})}{e_{ij}} e_{ij} \frac{\mathrm{d}e_{ij}}{\mathrm{d}p_{1}} = 0, \qquad (5)$$

$$\sum_{i=1}^{n} \sum_{j=1}^{m} \frac{\psi(e_{ij})}{e_{ij}} \frac{\mathrm{d}e_{ij}^{2}}{\mathrm{d}p_{1}} = 0.$$
 (6)

定义 $w(e_{ij}) = \frac{\psi(e_{ij})}{e_{ij}}$, 对于迭代优化求解过程 中固定的 $w(e_{ij})$, 式(6) 又可以描述为如下优化问

$$\min \sum_{i=1}^{n} \sum_{j=1}^{m} w(e_{ij})(e_{ij})^{2}, \qquad (7)$$

即

$$\min \sum_{i=1}^{n} \sum_{j=1}^{m} (\sqrt{w(e_{ij})}e_{ij})^{2}.$$
 (8)

定义 $z_{ij} = \sqrt{w(e_{ij})}x_{ij}$,则式(8)的根本优化目的 在于

$$\min \sum_{i=1}^{n} \| \boldsymbol{z}_{i} - \boldsymbol{z}_{i} \boldsymbol{p}_{1} \boldsymbol{p}_{1}^{\mathsf{T}} \|^{2}.$$
 (9)

根据对数据分布的不同假设,可以使用不同的-M估计器。本文中选用 Huber 提出的一种 M 估计 器^[10],具有下述形式的ρ函数:

$$\rho(e) = \begin{cases} e^{2/2}, \ |e| \le s; \\ s|e| - s^{2/2}, \ |e| > s. \end{cases}$$
(10)

式中,s为调节参数。该估计器重视中间样本的作用, 削弱极端样本的影响。由于 $w(e) = \frac{\rho'(e)}{e}$, 与 ρ 函数 对应的加权系数 w 的计算公式为

 $w(e) = \begin{cases} 1, \ |e| \le s; \\ s/|e|, \ |e| > s. \end{cases}$ (11)

传统的 PCA 分析可以通过 NIPALS 算法来求 解^[11],对于本文中的鲁棒 PCA 优化问题,需要对 NIPALS 算法进行改进。将鲁棒 M 估计和 NIPALS 算 法相结合,可以得到如下鲁棒 PCA 算法步骤:

①采集过程数据X,对数据进行标准化;

② 对第一个载荷向量 t₁ 进行随机初始化;

③ 计算 $p_1, p_1 = X^T t_1 / t_1^T t_1$,并进行归一化: $p_1 = p_1 / \|p_1\|$;

④计算 $t_1, t_1 = Xp_1;$

⑤ 返回步骤 ③ 进行迭代计算,比较两次计算 出的 p₁,如果误差在允许的范围内,则算法已经收 敛,继续向下计算:

⑥ 计算残差 $E = X - t_1 p_1^T$;

⑦ 按照式(11) 计算加权系数 $w(e_{ij})$, 对 x 进行 加权更新 $x_{ii} = \sqrt{w(e_{ii})} x_{ii}$;

⑧返回步骤③进行更新计算,如果最近两次计算出p,之差在允许的范围内,则计算完毕,得到p,;

⑨在 p_1 的正交补空间中寻找满足条件的 p_2 ,依次类推,可以找出所有的载荷向量;

④为保证载荷向量更为准确,在所有载荷向量 求出后,可返回步骤②重新寻找所有载荷向量,直 到两次计算出的P在误差允许的范围内结束。

1.3 主元模型的比较

为验证鲁棒 PCA 方法的有效性,需要使用一些 指标来比较两个主元模型的相似性。

如果两个 PCA 模型的主元方向完全对应相同或 相近,则可以认为两模型是同一模型。主元方向的相 似性比较通常使用夹角余弦来描述^[12],如

$$|\cos\theta| = \frac{|a'b|}{\|a\|\cdot\|b\|}.$$
 (12)

如果夹角余弦绝对值趋于1,则两方向趋于同一主 元方向,否则趋于0则两方向趋于正交。

还可以通过计算两个主元子空间的距离来评价 两个主元模型的相似性^[13]。两个主元子空间 *S*₁ 和 **S**₂ 的主元方向分别为 W 和 V,则它们之间的距离 D(S₁,S₂)为

 $D(S_1, S_2) = \sqrt{1 - \lambda_{\min}(W^T V V^T W)}$. (13) 如果该距离趋于0,表示 $S_1 \approx S_2$ 描述了同一个子空 间,趋于1则表示不是同一个子空间。

2 基于鲁棒 PCA 的故障检测

与传统 PCA 方法相似,使用鲁棒 PCA 算法来进 行故障检测需要构造 T² 和 Q 两个统计量,定义如 下:

 $T^{2} = xP_{k}\Lambda^{-1}P_{k}^{\mathsf{T}}x^{\mathsf{T}}, Q = ee^{\mathsf{T}}.$ 统计阈值分别为 $T^{2}_{a} = \frac{k(n^{2}-1)}{(n-k)}F_{a}(k,n-k),$

$$Q_{\alpha} = \theta_1 \Big[\frac{h_0 c_{\alpha} \sqrt{2\theta_2}}{\theta_1} + 1 + \frac{\theta_2 h_0 (h_0 - 1)}{\theta_1^2} \Big]^{1/h_0}.$$

其中 k 为主元空间主元数目,其他参数的定义可以 参考文献[14]。

值得注意的是,使用过程新数据进行故障检测 之前,需要对数据进行标准化处理,本文中使用中值 估计和 Q,散度估计^[1]代替传统的均值和标准差完 成标准化处理。

3 仿真结果分析

以一个连续搅拌反应釜(CSTR)系统作为仿真 对象,验证鲁棒 PCA 算法用于过程监控的有效性。 CSTR 控制系统如图 1 所示(图中, C_{AF} 为人口物流中 A 的浓度, Q_F 为人口物流的流量, T_F 为人口物流的温 度, C_A 为出口物流中 A 的浓度,Q 为出口流量, Q_C 为 冷却水流量, T_{CF} 为冷却水人口的温度, T_C 为冷却水出 口温度,T 为反应釜温度,h 为反应釜液位,LT 为液位 变送器,LC 为液位控制器,FT 为流量变送器,FC 为流 量控制器,TT 为温度变送器,TC 为温度控制器)。



物料 A 进入反应釜发生一级不可逆反应,生成物质 B,同时放出大量的热,冷却剂通过夹套把热量带走。

在对 CSTR 系统进行仿真的过程中,采集 10 个 测量变量(见图 1)并添加正态分布的噪声,得到正 常工况和 10 种故障情形(见表 1)的仿真数据。在正 常工况数据的基础上,随机抽取了 10% 的数据在原 测量值的基础上增加或减少 5%,作为离群数据。

表1 故障类型

故障	故障描述
F1	进料量突然发生变化
F2	进料温度逐渐发生变化
F3	进料浓度逐渐发生变化
F4	热交换故障,冷水散热能力降低
F5	催化剂失活
F6	冷却水进入,温度变化
F7	反应器内温度设定值变化
F8	进料温度传感器出现故障
F9	反应温度传感器出现故障
F10	冷却水调节阀出现故障

使用 PCA 方法分析无离群点的正常数据,建立 主元模型,该模型是标准模型,记为模型I。分别使用 鲁棒 PCA 和传统 PCA 方法分析含有离群点的正常 工况数据,建立主元模型 II 和模型 II,三种模型的 故障检测阈值均为95% 置信限。鲁棒 PCA 模型中误 差加权系数根据式(11) 计算,其中调节参数 s 根据 重构误差数据的 Q。散度估计进行适当选取。表 2 中 列出模型 I 的主元描述方差情况,它是正常工况数 据进行 PCA 模型分析的真实描述。模型 II 的方差描 述列于表3中,各个主元描述的方差与表2 的情况非 常相近,前5 个主元可以描述大约70% 的变化。

表2 模型 I 方差描述结果

	•••				
主元 序号	描述 方差	累积 方差(%)	主元 序号	描述 方差	累积 方差(%)
1	1.9780	19.7802	6	0.9628	82.9748
2	1.9289	39.0694	7	0.9089	92.0639
3	1.3394	52.4636	8	0.6863	98.9272
4	1.0649	63.1125	9	0.0884	99 . 811 7
5	1.0234	73.3465	10	0.0188	100.0000

表3 模型Ⅱ方差描述结果

主元 序号	描述 方差	累积 方差(%)	主元 序号	描述 方差	累积 方差(%)
1	1.8167	18.1669	6	0.9724	79.7173
2	1.7632	35.7990	7	0.9006	88.7234
3	1.3044	48.8428	8	0.7200	95.9233
4	1.0957	59.7997	9	0.2268	98. 191 8
5	1.0194	69.9938	10	0.1808	100.0000

表4为模型 Ⅲ 的主元方差分析,前5个主元仅

描述 60% 的数据方差,这是因为直接对含有离群点 的数据进行 PCA 分析,模型 Ⅲ 的方差分布受到离群 点的影响,各个主元所解释方差偏离了真实模型。图 2 中给出3 个模型描述方差的比较。可以看出,模型 Ⅱ 和模型 Ⅱ 的方差曲线接近,模型 Ⅲ 的方差曲线则 表现出明显的不同。

表4 模型Ⅲ方差描述结果

主元 序号	 描述 方差	累积 古差(%)	主元 序号	描述	 累积 古美(%)
71.2	<u></u>	7足(ル)	11.2	7/2	7足(~)
1	1.4208	14.2082	6	0.9717	70.0038
2	1.3850	28.0583	7	0.9526	79.5294
3	1.1472	39.5301	8	0.8223	87.7523
4	1.0840	50.3702	9	0.6311	94.0631
5	0.9917	60.2870	10	0.5937	100.0000
	2.0 1.5 概 规 行 授 1.0 - 概 一 0.5 - 0 0	2 4	模型 I 製型 Ⅲ →	6 8	10

图 2 模型 Ⅰ、Ⅱ和Ⅲ的方差比较

对3个模型的载荷向量(主元方向)进行比较。 以模型 I 作为标准模型,分别计算其与模型 II 和 模型 II 的各个主元方向的夹角余弦的绝对值 |cosθ_{1.I}|,|cosθ_{1.I}|(表5)。从表5中可以看出,模 型 I 和模型 II 的10个主元之间夹角余弦绝对值有 8个在0.9以上,而模型 I 和模型 II 的10个主元之 间夹角余弦绝对值只有前2个大于0.9,其他的均 小于0.85,说明模型 I 和模型 II 的10个主元方向 基本相同,而模型 II 由于没有消除离群点对主元模 型的影响,其主元方向偏离真正的主元模型。

表 5	模型	Π	. Ш	的主元方向比较	
तर २	医尘	ш.	чш	的エルカ凹ん我	

主元 序号	$ \cos \theta_{1,\mathbf{I}} $	cos θ _{I,∎}	主元 序号	$ \cos \theta_{I,I} $	$ \cos \theta_{\mathrm{I},\mathbf{I}} $
1.	0.9669	0.9476	6	0.9810	0.4684
2	0.9671	0.9810	7	0.9085	0.3307
3	0.9745	0.7360	8	0.9615	0.7785
4	0.5782	0.1045	9	0.9784	0.8337
5	0.5449	0.2774	10	0.9789	0.8302

如果3个模型均取前7个主元作为主元子空间, 而后3个主元作为残差子空间,根据子空间距离计 算方法,分别计算模型 Ⅰ 和模型 Ⅱ、Ⅲ 的主元子空 间之间的距离,鲁棒 PCA 模型 Ⅱ 与真实的主元模型 Ⅰ 之间的距离仅为 0. 2715,相比之下,该值远小于 没有考虑离群点影响的主元模型 Ⅲ 与真实的主元 模型 Ⅰ 之间的距离 0.627 1。

从上述对主元方向间夹角的比较和主元子空间 距离的计算结果分析中可以看出,鲁棒 PCA 方法在 建模数据存在离群点的情况下能够建立更为准确的 主元模型。

分别使用模型 Ⅱ、Ⅲ 对仿真故障进行检测,并 对故障F3 和F7 的结果进行分析(见图3~6)。故障 F3 的描述为从第200个采样时刻开始进料浓度逐渐 降低,从而引起系统浓度的变化和冷却剂流量及出 口温度的变化。从图3和图4中可以看出,鲁棒 PCA 方法的 T^e统计量从第310个采样时刻就开始超出阈 值,传统的 PCA 方法直到第605个采样时刻以后才 开始超出阈值,两种方法的 Q 统计量显著超出阈值 的时间都在650个采样时刻附近,但是使用鲁棒 PCA 方法时,从第500个采样时刻起 Q 统计量就开 始在阈值附近波动。



故障 F7 为在第 200 个采样时刻反应温度设定 值变化4 °C,鲁棒 PCA 和传统 PCA 方法的监控结果 分别见图5 和图6。鲁棒 PCA 方法可以显著地检测到 故障,而 PCA 方法虽然两个统计量的曲线在第 200 个采样时刻有明显的变化,但是没有明显地超出阈 值。从 Q 统计量上分析,由于离群点的影响,造成 PCA 模型的统计阈值略大于鲁棒 PCA 模型的统计 阈值,因此虽然两种方法的 Q 统计量变化相近,但是 PCA 方法的 Q 统计量没有超出阈值。从 T^2 统计量上 分析,两模型 T^2 统计量的阈值完全相同,但是由于 离群点使得主元方向发生了偏离,因此 PCA 方法的 T^2 统计量未能显著检测到故障。

4 结束语

在建模数据中存在离群点的情况下,基于 M 估 计和改进的 NIPALS 算法的鲁棒 PCA 方法可以更为 准确地描述过程,及时地检测到过程中出现的故障。 真实生产过程中的数据往往并不满足理想的正态分 布,因此使用鲁棒 PCA 方法建模有利于准确挖掘数 据信息,从而更好地保证生产的安全。在未来的研究 中,如何将鲁棒 PCA 方法用于在线故障检测是一个 值得研究的问题。

参考文献:

- MOLLER F S, FRESE J, BRO R. Robust methods for multivariate data analysis [J]. Journal of Chemometrics, 2005, 19(10):549-563.
- [2] MORAD K, SVRCEK W Y, MCKAYI. A robust direct approach for calculating measurement error covariance matrix [J]. Computer and Chemical Engineering, 1999, 23(7):64-79.
- [3] CHEN J, BANDONI A, ROMAGNOLI J A. Robust statistical process monitoring [J]. Computer and Chemical Engineering, 1996,20(s1):497-502.
- [4] HUBERT M, ROUSSEEUW P J, VERBOVEN S. A fast

(上接第142页)

- [6] GE S S, HONG F, LEE T H. Adaptive neural network control of nonlinear systems with unknown time delays
 [J]. IEEE Transactions on Automatic Control, 2003,48 (11):2004-2010.
- [7] HO D W C, ZHANG P, XU J. Fuzzy wavelet networks for function learning [J]. IEEE Trans Fuzzy System, 2001,9(1):200-211.
- [8] SJBERG J, ZHANG Q, LJUNG L, et al. Nonlinear black-

method for robust principal components with applications to chemometrics [J]. Chemometrics and Chemical Engineering, 2002,60(1/2):101-111.

- [5] XIE L, ZHANG J M, WANG S Q. A robust statistical batch process monitoring framework and its application
 [J]. Chinese Journal of Chemical Engineering, 2004,12
 (5):682-687.
 - YANG T N, WANG S Q. Robust algorithms for principal component analysis [J]. Pattern Recognition, 1999, 20 (9): 927-933.
- [7] SARBU C, POP H F. Principal component analysis versus fuzzy principal component analysis a case study: the quality of danube water (1985-1996) [J]. Talanta, 2005,65(5):1215-1220.
- [8] WANG D, ROMAGNOLI J A. Robust multi-scale principal component analysis with applications to process monitoring[J]. Journal of Process Control, 2005, 15(8):869-882.
- [9] .LI Y. On incremental and robust subspace learning [J]. Pattern Recognition, 2004,37(7):1509-1518.
- [10] HUBER P J. Robust estimation of a location parameter
 [J]. The Annals of Mathematical Statistics. 1964, 35

 (1): 73-101.
- [11] 张杰,阳宪惠. 多变量统计过程控制[M]. 北京:化 学工业出版社, 2000.
- [12] RAICH A, CINAR A. Diagnosis of process disturbance by statistical distance and angle measures [J]. Computers and Chemical Engineering, 1997,21(6):661-673.
- [13] KANO M, HASEBE S, HASHIMOTO I, et al. A new multivariate statistical process monitoring method using principal component analysis[J]. Computers and Chemical Engineering, 2001,25(7/8):1103-1113.
- [14] CHIANG H L, RUSSELL E L, BRAATZ R D. Fault detection and diagnosis in industrial systems [M]. London: Springer, 2001.

(编辑 修荣荣)

box modeling in system identification: a unified overview [J]. Automatica, 1995, 31(12):1691-1724.

[9] XU J, HO D W C, ZHOU D. Adaptive wavelet networks for nonlinear system identification [C]//American Automatic Control Council. Proceedings of American Control Conference. Piscataway, N. J. USA: American Automatic Control Council, c1999;3472-3473.

(编辑 修荣荣)