

基于互信息熵-近红外光谱的过程模式故障检测

高爽, 栾小丽*, 刘飞

江南大学自动化研究所, 轻工过程先进控制教育部重点实验室, 江苏无锡 214122

摘要 近红外光谱分析在工业过程故障检测方面具有独特的优势, 是一种准确且高效的方法。结合互信息熵和传统的主成分分析, 对近红外光谱特征信息进行提取, 通过构建过程的模式来刻画工业过程的运行状态。利用近红外光谱数据, 从有机分子含氢基团振动信息中获取工业系统的过程模式, 从微观分子层面探索提高工业过程故障检测准确率的有效方法, 结合贝叶斯统计学习技术, 提出了基于近红外光谱数据的工业过程故障检测技术。针对近红外光谱信息量丰富, 谱带较宽, 特征性不强的特点, 首先对工业过程不同运行状态下的近红外光谱吸光度数据进行一阶导数预处理, 采用主成分分析法(principal component analysis, PCA)压缩光谱数据量, 扩大不同运行状态下光谱特征信息的差异性, 提取光谱的内部特征信息。然后采用互信息熵(mutual information entropy, MIE)作为光谱特征信息相关性度量函数, 基于最小冗余最大相关算法进一步减少光谱特征信息间的冗余并最大化光谱特征信息与类别的相关性, 弥补了PCA无监督特征波长选择的不足, 提出一种基于PCA-MIE的过程模式构建方法, 获得的过程模式子集更紧凑更具类别表现力。再利用贝叶斯统计学习算法, 根据后验概率对构建的模式子集进行决策, 判别生产过程的正常状态和故障状态。由于过程模式子集结合了PCA浓聚方差的优势和互信息熵相关性测度的特征信息选择方法, 蕴含了更多的近红外光谱的本质信息与内在规律, 从而更能刻画工业过程的运行状态。接着, 设置测试准确率TA作为评估标准, 用以评价故障检测方法的性能效果。最后利用某化工厂提供的原油脱盐脱水过程近红外光谱数据对所提方法进行验证, 并与传统近红外光谱特征信息提取方法PCA和MIE方法性能进行对比分析, 结果表明基于PCA-MIE的过程模式故障检测方法几乎在所有维数子集上性能都优于其他两种方法, 在特征维数为18维时获得最高的准确率94.6%, 证明了方法的优越性。

关键词 近红外光谱; 互信息熵; 过程模式; 故障检测; 贝叶斯

中图分类号: O657.3 **文献标识码:** A **DOI:** 10.3964/j.issn.1000-0593(2019)06-1736-06

引言

故障检测技术在过程控制系统中发挥越来越重要的作用, 对过程工业的节能减排、安全稳定运行以及提高产品质量具有重要的意义^[1-3]。传统故障检测有三大类方法: 基于机理模型、基于知识以及基于数据驱动的故障检测方法^[4], 其中基于机理模型的方法是最直接有效的, 但由于数学模型的建立及求解十分困难, 使其应用受到限制, 基于知识的方法便应运而生, 它不需要对象的精确数学模型, 其局限表现在对专家知识的依赖、自适应能力和学习能力差等方面。而基于数据驱动的故障检测方法既不需要系统精确的数学模型, 也不依赖工业知识和生产经验, 在石油化工、农业、医

药、食品等诸多领域得到了广泛应用^[5]。但上述研究大多集中在基于过程宏观变量的故障检测。

随着先进测量仪器的使用, 基于分子振动的近红外光谱分析, 具有无损、快速、高通量、低成本等优点, 为工业过程的故障检测提供了技术手段^[6]。由于近红外光谱是从分子层面获取过程信息, 相比较于基于过程宏观变量的故障检测方法, 对早期故障的判断更为灵敏, 因此近红外作为一种新兴的测量仪器, 在故障检测领域展现了其极大的优势与应用前景。由于光谱信息量(待选自变量数目)远远大于样本量(采样次数), 光谱数据之间存在大量的冗余和共线性信息。为了降低模型复杂度, 对光谱信息进行变量选择(剔除无效信息、保留对样品品质指标有显著影响的变量)格外重要。

利用信息论中的互信息熵^[7](MIE), 在光谱信息量大于

收稿日期: 2018-05-06, 修订日期: 2018-10-11

基金项目: 国家自然科学基金项目(61473137, 61722306)资助

作者简介: 高爽, 女, 1995年生, 江南大学自动化研究所轻工过程先进控制教育部重点实验室硕士研究生

e-mail: gao_shuang1995@126.com * 通讯联系人 e-mail: xlluan@jiangnan.edu.cn

样本量且存在大量冗余和共线性的背景下,提出基于 PCA-MIE 的光谱特征波长选择方法,从所选择的特征波长子集中获取工业系统的过程模式。另外区别于传统利用投影空间中的 T^2 统计量和残差空间中的 Q 统计量为监控指标的故障检测技术^[8],利用贝叶斯统计学习判别法,对所构建的过程模式进行决策和划分,判断系统是否处于故障状态。由于所构建的过程模式不仅度量了特征与目标类别之间的相关程度,而且也度量了待选特征项和已选特征项子集的相关程度,并通过一个前向顺序搜索算法将两个标准结合起来,同步优化两个指标,因此基于过程模式的故障检测方法具有更高的准确率。

1 过程模式的构建

设矩阵 $\mathbf{X} \in \mathbf{R}^{N \times M}$ 为所测样本的光谱数据矩阵, N 为样本数, M 为变量数,对光谱数据矩阵进行 PCA 分解^[9],可得

$$\mathbf{X} = \mathbf{T}\mathbf{P}^T = t_1 p_1^T + t_2 p_2^T + \dots + t_M p_M^T \quad (1)$$

式中: t_i 为得分向量; p_i 为负荷向量, $i=1, \dots, M$ 。得分向量和负荷向量均是两两正交,且负荷向量的模是 1。每一个得分向量 t_i 是矩阵 \mathbf{X} 与在此得分向量相应的负荷向量 p_i 方向上的投影,即主元

$$t_i = \mathbf{X}p_i \quad (2)$$

对光谱数据矩阵进行主成分分析实质上是对矩阵的协方差矩阵进行特征向量分析,矩阵 \mathbf{X} 的协方差矩阵可以表示为

$$\mathbf{E} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad (3)$$

对 \mathbf{E} 做特征向量分析,即求解式(4)

$$\lambda_i p_i = \mathbf{E}p_i \quad (4)$$

取前 k 个负荷向量 $[p_1, p_2, \dots, p_k]$ 构成新的主元空间 \hat{U} ,余下的 $M-k$ 维负荷向量 $[p_{k+1}, \dots, p_M]$ 构成残差子空间,则光谱数据 \mathbf{X} 在主元空间 \hat{U} 的投影可将数据特征由原来的 M 维降到 k 维,记为 $\mathbf{F} = [f_1, f_2, \dots, f_k]$,且保留了原光谱中 $\sum_{i=1}^k \lambda_i / \sum_{i=1}^M \lambda_i$ 的信息。

主成分分析虽然在不丢失主要光谱信息的前提下选择为数较少的新变量来代替原来较多的变量,解决了由于谱带的重叠而无法分析的困难,识别出最重要的多个特征。然而作为无监督特征提取算法,PCA 舍弃已有类别标记信息,将所有数据当成无类别标记数据,识别出的不一定是所需要的特征,会损失有用信息。为弥补此不足,创新性地提出了一种混合算法 PCA-MIE 来构成过程模式,获得更优的特征波长子集。

信息熵(IE)是信息论中一个表示变量取值不确定性程度的指标。设两个离散变量 A, B 均取有限值,如果信源 a 与随机变量 b 不是相互独立的,它们的联合分布为 $p(a, b) = P\{A=a, B=b\}$,边缘分布为 $p(a) = P\{A=a\}$, $p(b) = P\{B=b\}$ 。则信源 a 的初始不确定度可以用熵 $H(A)$ 表示^[10]

$$H(A) = - \sum_a p(a) \log p(a) \quad (5)$$

已知 B 下, a 的条件熵定义为

$$H(A|B) = - \sum_{a,b} p(a, b) \log p(a|b) \quad (6)$$

式中 $p(a|b)$ 为条件概率, $p(a|b) = \frac{p(a, b)}{p(b)}$, $p(b) > 0$ 。在已

知随机变量 B 取值的条件下,随机变量 A 的条件熵 $H(A|B)$ 总是不大于初始熵 $H(A)$,即 B 的存在减小了 A 的不确定度。这种不确定度的减少量被定义为变量 A 和 B 之间的互信息熵(MIE)记为 $I(A; B)$

$$I(A; B) = H(A) - H(A|B) \quad (7)$$

即

$$I(A; B) = \sum_b \sum_a p(a, b) \log \frac{p(a, b)}{p(a)p(b)} \quad (8)$$

最小冗余最大相关(mRMR)算法^[11]中最大相关性是 PCA 降维后的特征 $f_i (i=1, \dots, k)$ 与类别 C 之间的相关性,采用互信息熵 $I(f_i; C)$ 来衡量,需要最大的优化函数如式(9)

$$\max D(S, C), D = \frac{1}{|S|} \sum_{f_i \in S} I(f_i; C) \quad (9)$$

其中 S 是候选特征子集, $|S| = r$ 表示选择特征的个数。在这样选择的子集中特征之间可能具有很多冗余的特征,也就是特征之间的依赖性很强。对于两个具有强依赖性的特征,如果去掉一个,不会对分类造成很大影响。因此,引进最小冗余度

$$\min R(S), R = \frac{1}{|S|^2} \sum_{f_i, f_j \in S} I(f_i; f_j) \quad (10)$$

将式(9)和式(10)进行组合,构成定义目标函数 $\Phi(\cdot)$,使其最大化

$$\max \Phi(D, R), \Phi = D/R \quad (11)$$

采用前向搜索法搜索最优特征子集,假设在当前步骤下已选出 $r-1$ 个特征,记为特征子集 S_{r-1} ,任务是从剩下的 $\{F - S_{r-1}\}$ 中,找到第 r 个特征 x_j ,使 $\Phi(\cdot)$ 最大,相应的前向搜索算法的实质是优化下面条件

$$\max_{f_j \in F - S_{r-1}} \left[I(f_j; C) / \frac{1}{r-1} \sum_{f_i \in S_{r-1}} I(f_j; f_i) \right] \quad (12)$$

即可得到过程模式 $\mathbf{Z} = [z_1, z_2, \dots, z_r]$, $r \leq k < M$ 。

2 基于过程模式的贝叶斯分类

贝叶斯分类模型是一种典型的基于统计方法的分类模型^[12]。贝叶斯定理是贝叶斯理论中最重要的一个公式,是贝叶斯学习方法的理论基础,它将事件的先验概率与后验概率巧妙地联系起来,利用先验信息和样本数据信息确定事件的后验概率^[13]。

在近红外故障检测问题中,原始光谱数据 \mathbf{X} 经过主成分分析得到数据矩阵 $\mathbf{F} = [f_1, f_2, \dots, f_k]$,再经过互信息熵变量选择后获得过程模式 $\mathbf{Z} = [z_1, z_2, \dots, z_r]$, $r \leq k < M$ 作为贝叶斯统计学习算法的输入。假定有 m 个类 C_1, C_2, \dots, C_m ,针对过程模式 \mathbf{Z} ,分类法将预测过程模式 \mathbf{Z} 属于最后后验概率的类。贝叶斯决策将未知的过程模式分配给类 C_i ,当且仅当

$$P(C_i | \mathbf{Z}) > P(C_j | \mathbf{Z}), 1 \leq j \leq m, j \neq i \quad (13)$$

其中 $P(C_i | \mathbf{Z})$ 为后验概率,可由如下贝叶斯定理计算获得

$$P(C_i | \mathbf{Z}) = \frac{P(\mathbf{Z} | C_i) P(C_i)}{\sum_{j=1}^m P(\mathbf{Z} | C_j) P(C_j)} \quad (14)$$

对过程模式 Z 进行判别, 对于每个类 C_i , 可转化为计算 $P(Z|C_i)P(C_i)$ 。过程模式 Z 被指派到类 C_i , 当且仅当

$$P(Z|C_i)P(C_i) > P(Z|C_j)P(C_j), 1 \leq j \leq m, j \neq i \quad (15)$$

令

$$l_{12}(Z) = \frac{P(Z|C_1)}{P(Z|C_2)} \quad (16)$$

$$\theta_{12} = \frac{P(C_2)}{P(C_1)} \quad (17)$$

其中 l_{12} 为似然比, θ_{12} 为阈值。对于原油脱盐脱水工业过程, 结合贝叶斯判别准则, 其故障检测原理可以表述为

$$\text{故障检测: } \begin{cases} l_{12}(Z) > \theta_{12} & Z \in C_1 \\ l_{12}(Z) < \theta_{12} & Z \in C_2 \\ l_{12}(Z) = \theta_{12} & Z \in C_1 \text{ or } C_2 \end{cases} \quad (18)$$

其中 C_1 为正常模式, C_2 为故障模式。采用贝叶斯统计学习方法进行检测, 其判别规则能对过程模式的变化做出迅速灵敏的判断, 且具有无需确定故障监测的门限的优点。

3 结果与讨论

3.1 样品与仪器

所用试验样品为某化工厂提供的原油脱盐脱水工业过程的光谱数据。采用德国 Bruker 公司生产的 MATRIX-F 型傅里叶红外光谱仪(含 OPUS 定量分析软件包)进行数据采集,

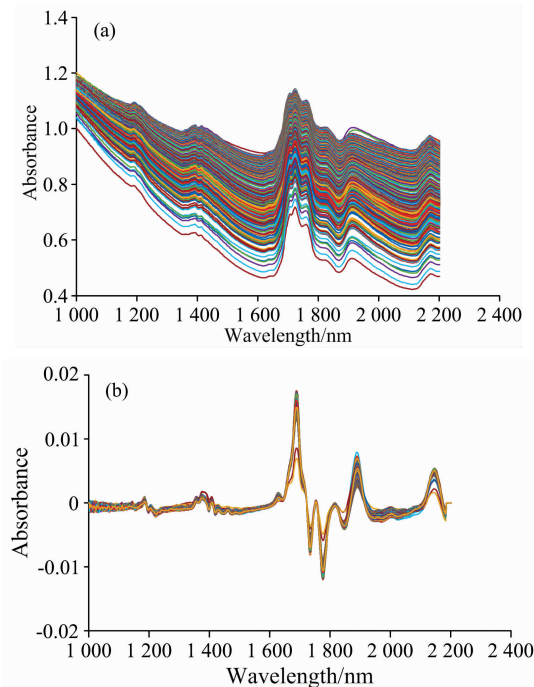


图 1 预处理前后光谱图

(a): 原始光谱; (b): 一阶导数光谱

Fig. 1 Spectra before and after pretreatment

(a): Raw spectra; (b): First derivative spectra

该仪器的光谱测量范围是 $12\ 800 \sim 4\ 000\ \text{cm}^{-1}$, 最小光谱扫描分辨率为 $2\ \text{cm}^{-1}$ 。实验设置光谱测量波长范围为 $10\ 000 \sim 4\ 500\ \text{cm}^{-1}$, 分辨率为 $8\ \text{cm}^{-1}$ 。

3.2 光谱采集与数据处理

利用 1 000 组光谱样本数据进行实验。为消除基线漂移和背景干扰, 分辨重叠峰, 提高分辨率。对原始光谱进行一阶导数预处理, 处理后的光谱如图 1(b)所示, 横坐标为波长, 范围为 $1\ 000 \sim 2\ 222\ \text{nm}$, 纵坐标为吸光度值。然后将测得的 1 000 组样本, 按照 3 : 2 的比例分成训练集和测试集, 分配情况如表 1。

表 1 样品的训练集和测试集

| Table 1 Samples of training and test sets | | | |
|---|-------|-------|-------|
| 样品类型 | 样品数 | 训练样本数 | 测试样本数 |
| 正常 | 500 | 300 | 200 |
| 故障 | 500 | 300 | 200 |
| 总计 | 1 000 | 600 | 400 |

3.3 结果与对比分析

3.3.1 PCA 特征波长选择

将光谱仪提取的光谱特征曲线一阶导数处理后, 得到光谱数据矩阵 $\mathbf{X} \in \mathbf{R}^{N \times M}$, 对其进行 PCA 分析。图 2 为主成分得分贡献率。分别以前 3 个主成分为坐标, 建立样本的三维得分图, 如图 3 所示。红色为故障样品, 蓝色为正常样品, 可以看出大部分故障与正常的样本各自聚为一类。由于三维视图信息维度较小, 且正常样品与早期故障样品内部基团组成差异不明显, 所以其光谱特征信息相似导致部分样品出现重叠现象。因主成分分析仅能呈现样品的聚类趋势, 未考虑类别信息, 识别出的不一定是所需要的特征, 会损失有用信息, 故降维后的主成分不能充分代表光谱的主要信息。

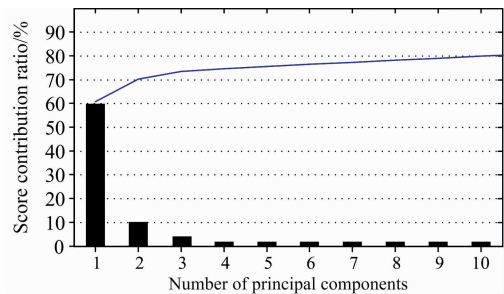


图 2 主成分得分贡献率

Fig. 2 Score contribution ratio of principal components

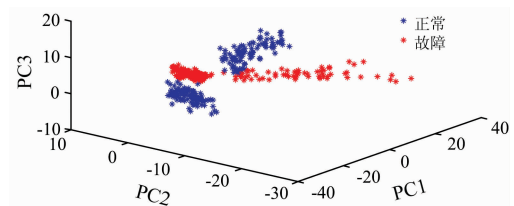


图 3 样品光谱前 3 主成分得分图

Fig. 3 Score of the first three PCs for sample spectral data

3.3.2 MIE 特征波长选择

MIE 方法将预处理后的光谱 699 个特征波长作为原始特征集 $X_{699} = \{x_1, x_2, \dots, x_{699}\}$, 初始最优特征子集为空集

S_0 , 设定选取的特征数量为 m , 根据光谱的特征波长和类别的数据, 计算 X 中各个特征波长 x_i 与类别 C 的互信息熵 $I(x_i; C)$, 结果如图 4 所示。

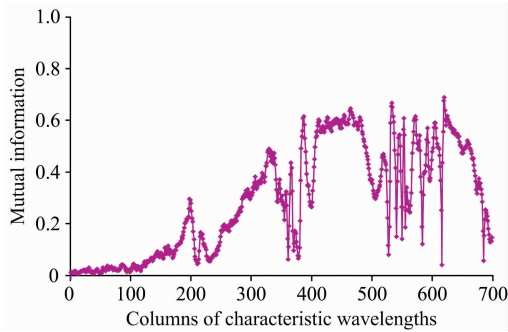


图 4 特征波长所在的列与类别互信息分布

Fig. 4 Distribution of mutual information between the columns for characteristic wavelengths and class

表 2 是按最小冗余最大相关算法降维之后的前 30 个特征排名(即 $m=30$)。该算法步骤如下: 首先, 选取互信息熵值最大的 620 列所在的特征波长 1 927 nm 作为加入最优特征子集 S 的第一个特征波长, 即 $S_1 = \{x_{620}\} X_{698} = X_{699} - \{x_{620}\} = \{x_1, \dots, x_{619}, x_{621}, \dots, x_{699}\}$ 。然后遍历 X_{698} 中每个特征波长, 找到满足式(12)的特征 x , 使得该特征波长既保证与类别之间的相关性充分大, 又能保证与 S 中已有特征波长的冗余度总和充分小, 将得到的特征波长加入 S_1 并从 X_{698} 中去除。以此类推, 直到 S 中特征波长的个数达到设定的 m 值。本文按照特征的排名, 选取特征数量分别为 1, 2, 4, 6, 8, 10, 12, 14, 16, 18 和 20。

3.3.3 PCA-MIE 特征波长选择

PCA-MIE 算法是特征提取的两阶段方法, 第一步先将 PCA 应用于原始特征波长获得简化特征子集 $F = [f_1, f_2, \dots, f_k]$, 当主成分累计贡献率 $> 95\%$ 时, $k=48$, 故先提取 48 个特征波长主成分。第二步将 MIE 应用于简化特征子集 F , 最大限度的减小冗余并最大化特征与类别之间的相关性, 进一步降低特征维数, 得到最优的过程模式子集 Z 。针对构建好的过程模式, 进一步利用贝叶斯统计学习算法进行故障检测, 用测试准确率(test accuracy, TA)作为评估标准, 其定义如下

表 2 最小冗余最大相关特征波长选择

Table 2 TOP30 characteristic wavelengths selected by min-redundancy and max-relevance

| 特征波长排名 | 特征波长所在的列 | 特征波长排名 | 特征波长所在的列 |
|--------|----------|--------|----------|
| 1 | 620 | 16 | 525 |
| 2 | 585 | 17 | 387 |
| 3 | 358 | 18 | 468 |
| 4 | 553 | 19 | 572 |
| 5 | 570 | 20 | 554 |
| 6 | 366 | 21 | 482 |
| 7 | 465 | 22 | 621 |
| 8 | 552 | 23 | 367 |
| 9 | 619 | 24 | 464 |
| 10 | 532 | 25 | 571 |
| 11 | 359 | 26 | 474 |
| 12 | 480 | 27 | 592 |
| 13 | 591 | 28 | 466 |
| 14 | 579 | 29 | 520 |
| 15 | 463 | 30 | 454 |

$$TA = \frac{CD}{Num} \times 100\%$$

其中, Num 表示测试样本总数, CD 表示正确决策的样本数量。

利用光谱 699 个特征波长, 通过 PCA-MIE 算法构造过程的模式子集, 再利用贝叶斯分类器进行故障检测, 表 3 是利用 PCA, MIE 和 PCA-MIE 三种方法获得的不同维数的特征子集和原始特征子集在贝叶斯统计学习算法下的准确率, 其中 NB 表示贝叶斯判别方法。从表 3 可以看出, 全光谱判别的正确率最低, 仅有 50.15%。因为全光谱数据中含有大量的噪声、冗余信息、干扰信息, 大大降低了正确率。利用 PCA 提取特征波长, 可以看出当主成分维数在 14 维时准确度达到最高值 93.11% 之后呈现稳定和下降趋势。这说明了, 选取的特征子集对原油脱盐脱水工业过程故障检测的影响之大以及进行特征波长选择的重要性和必要性。当维数大于 14 维时由于冗余的作用, 不仅不能获得更多的运行信息, 反而会因为特征波长子集中的冗余特征项使判别准确率下降。利用 MIE 选择特征波长, 在特征波长数为 12 维时的准确度最高为 93.56%, 即最优特征子集数为 12。

表 3 贝叶斯算法判别准确率%

Table 3 Discriminant accuracy of Bayesian algorithm%

| 算法 | 主成分数/特征波长数 | | | | | | | | | | | |
|---------|------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| | 1 | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 | 20 | 699 |
| PCA-MIE | 91.33 | 91.9 | 92.15 | 93.47 | 93.53 | 93.91 | 94.1 | 94.21 | 94.41 | 94.6 | 94.2 | — |
| MIE | 87.01 | 89.32 | 92.21 | 92.91 | 93.15 | 93.21 | 93.56 | 93.5 | 92.86 | 93.14 | 93.25 | — |
| PCA | 91.16 | 90.38 | 92.32 | 92.86 | 93.23 | 93.03 | 93.05 | 93.11 | 92.99 | 92.92 | 92.95 | — |
| NB | — | — | — | — | — | — | — | — | — | — | — | 50.15 |

图 5 是 PCA, MIE 和 PCA-MIE 三种方法获得的不同维数子集的贝叶斯判别准确率。从图 5 可以看出在特征波长子集维数较低(低于 5 维)时, PCA 方法分类准确率高于 MIE

算法, 而大于 5 维之后, 随着特征项维数的增加, MIE 方法准确率慢慢赶上并超过 PCA。本算法 PCA-MIE 几乎在所有维数子集上性能都优于其他两种方法, 说明该算法提取的特

征子集更有类别代表能力且冗余最小。在特征维数为 18 维时获得最高准确度 94.6%。

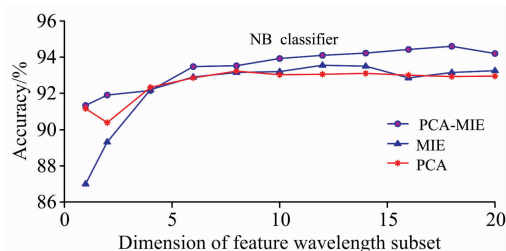


图 5 贝叶斯判别准确率

Fig. 5 Bayesian discriminant accuracy

4 结 论

针对近红外光谱,提出了一种基于互信息熵的 PCA-MIE 算法,该方法结合了传统主成分分析的优势,并从信息论的角度对特征波长的相关度和冗余度进行了综合考量,弥补了无监督特征波长选择的不足。通过 PCA-MIE 获得的模式子集由于包含了更多的过程信息,使得过程模式框架下的贝叶斯分类具有更高的准确率,因此该方法在工业过程的故障检测领域具有良好的应用前景。未来研究将从结构扩展和属性加权两个方面改进贝叶斯统计学习算法,构造组合分类器,以进一步提高故障检测的正确率。

References

- [1] Minh V T, Afzulpurkar N, Muhamad W M W. *Mathematical Problems in Engineering*, 2015, 2007(1): 32.
- [2] Yu J. *Engineering Applications of Artificial Intelligence*, 2013, 26(1): 456.
- [3] Haghani A, Jeansch T, Ding S X. *IEEE Transactions on Industrial Electronics*, 2014, 61(11): 6446.
- [4] Chen Z, Zhang K, Ding S X, et al. *Journal of Process Control*, 2016, 41: 26.
- [5] Chen Z, Fang H, Chang Y. *IEEE Transactions on Industrial Electronics*, 2016, 63(5): 3290.
- [6] CHU Xiao-li, LU Wan-zhen(褚小立, 陆婉珍). *Spectroscopy and Spectral Analysis(光谱学与光谱分析)*, 2014, 34(10): 2595.
- [7] Paninski L. *Neural Computation*, 2014, 15(6): 1191.
- [8] Lee K, Lee B W, Choi D H, et al. *Journal of the Electrochemical Society*, 2014, 18(3): 53.
- [9] Zhao Q H, Chen M. *Analytical Letters*, 2015, 48(2): 301.
- [10] Zeng X, Wu J, Wang D, et al. *Journal of Hydrology*, 2016, 538: 689.
- [11] Peng H, Long F, Ding C. *IEEE Trans Pattern Anal. Mach. Intell.*, 2005, 27(8): 1226.
- [12] Cheng J, Greiner R. *IEEE Transactions on Vehicular Technology*, 2013, 63(5): 2002.
- [13] Liu Z J, Guo S L, Tian-Yuan L I, et al. *Journal of Hydraulic Engineering*, 2014, 45(9): 1019.

Near Infrared Spectroscopy Process Pattern Fault Detection Based on Mutual Information Entropy

GAO Shuang, LUAN Xiao-li*, LIU Fei

Key Laboratory for Advanced Process Control of Light Industry of Ministry of Education, Institute of Automation, Jiangnan University, Wuxi 214122, China

Abstract The technology of near infrared spectroscopy that has unique advantage in fault detection in industrial processes is an accurate and effective method. Combining the mutual information entropy and the traditional principal component analysis, a new method for extracting the near infrared spectral feature information was first developed. The operating states of the industrial process was described by the constructed process pattern. Near infrared spectroscopy data were used to obtain the process pattern of industrial systems from the vibration information of hydrogen groups in organic molecules in this paper. An effective method to improve accuracy of fault detection in industrial processes was explored from the microscopic molecular level. Combined with Bayesian statistical learning method, an industrial processes fault detection technique based on near infrared spectroscopy data was proposed. Firstly, for the characteristics of rich information, wide spectrum band and weak characteristic, first-order derivative preprocessing of near infrared spectroscopic absorbance data under different operating states of industrial process was applied. Principal component analysis(PCA) was used to compress the amount of spectral data, expand the differences in spectral feature information under different operating states, and extract the internal feature information of the spectrum. Then, mutual information entropy(MIE) was used as correlation measure function of spectral feature information, and the minimum redundancy maximum relevance algorithm was used to further reduce the redundancy between the spectral feature information and maxi-

mize the relevance between the spectral and class. It made up for the deficiency of unsupervised feature wavelength selection of PCA. Therefore, a process pattern construction method based on PCA-MIE was proposed. The obtained process pattern subset was more compact and more expressive. Furthermore, Bayesian statistical learning method was applied to make decisions based on posterior probability of the constructed process pattern subset to identify the normal and accident state of the production process. Because the process pattern subset combines the advantages of PCA in density variance reduction and the feature information selection method of mutual information entropy correlation measure, it contains more essential information and inherent laws of near infrared spectroscopy, which can better describe the operating states of the industrial process. Next, The test accuracy (TA) was set as the evaluation criteria to evaluate the performance of the fault detection method. Finally, the data of crude oil desalination and dehydration process provided by the chemical plant was used to verify the effectiveness of the proposed method. Compared with the performance of traditional near infrared spectral feature information selection methods PCA and MIE, the results showed that the process pattern fault detection based on PCA-MIE outperforms the other two methods on almost all dimensions subsets. The highest accuracy rate is 94.6% when the feature dimensions is 18, which proves the superiority of the proposed method.

Keywords Near infrared spectroscopy; Mutual information entropy; Process pattern; Fault detection; Bayesian

(Received May 6, 2018; accepted Oct. 11, 2018)

* Corresponding author