

文章编号: 1672-3961(2010)05-0060-06

一种基于 SVM 的快速特征选择方法

戴平, 李宁*

(南京大学计算机软件新技术国家重点实验室, 南京 210093)

摘要: 针对现有特征选择方法计算量大、速度慢的缺点, 提出了一种基于 SVM 的快速特征选择算法。该算法使用 SVM 作为分类器, 并利用粒子群优化算法进行搜索。通过利用 SVM 线性核与多项式核函数的特性, 减少了在特征选择中训练分类器的次数, 降低了计算复杂度。实验结果表明在不损失分类精度的情况下, 能显著提高特征选择的速度。

关键词: 特征选择; 支持向量机; 粒子群优化

中图分类号: TP319 **文献标志码:** A

A fast SVM-based feature selection method

DAIPing LINing

(National Laboratory for Novel Software Technology, Nanjing University, Nanjing 210093, China)

Abstract: Aiming at large computation and slow convergence speed of the traditional feature selection methods, a fast SVM-based feature selection method was proposed. A support vector machine was employed as the classifier, and particle swarm optimization method was employed as the searching strategy. The proposed method reduced the iterations of training classifiers by taking advantage of the characteristics of linear and polynomial kernel functions, so that it could reduce the complexity of the calculation. Experimental results showed that the method accelerates feature selection in the case of no loss of the classification performances.

Key words: feature selection; support vector machine; particle swarm optimization

0 引言

支持向量机 (support vector machine, SVM)^[1] 是一种用于分类和回归的监督式学习算法。该方法是建立在统计学习理论基础上的机器学习方法, 相比其它机器学习方法, SVM 方法具有良好的分类性能和泛化能力^[2]。目前, SVM 方法被应用到了物体识别、医学辅助诊断等许多领域^[3]。

但是, 已有的研究显示, SVM 在有冗余或不相关特征的数据集上, 分类性能会下降, 并不是越多的特征就越有利于提高学习算法的分类精度^[4]。而且, 随着特征规模的增大, 对机器学习在大规模数据处理上的研究提出了迫切要求。因此, 使特征选择成为机器学习中的一个重要研究内容。

粒子群优化算法 (particle swarm optimization algorithm, PSO) 是一种基于群智能的进化计算方法。粒子群算法不像遗传算法依靠交叉算作、变异算子等操作个体, 而是依靠个体间的信息交换完成进化演算。目前, PSO 已经被应用在特征选择^[5]、分类^[6]等许多方面。

Wesley Mao Sindhwan 等提出了一系列特征选择方法来提高 SVM 分类器的准确性^[4, 7-8], Shih-Wei Lin

收稿日期: 2010-04-02

基金项目: 国家自然科学基金资助项目 (60875011)

作者简介: 戴平 (1986-), 男, 上海人, 硕士研究生, 主要研究方向为机器学习。E-mail: dpl130@163.com

* 通讯作者: 李宁 (1968-), 女, 江苏沐阳人, 副教授, 硕士, 主要研究方向为机器学习与图像处理。E-mail: ln@nju.edu.cn

等提出了利用 PSO 进行特征选择和 SVM 参数联合优化的方法^[9]。但是, 上述 Wrapper 特征选择方法计算量大、速度慢。因此, 本文利用 SVM 线性核与多项式核函数的特性, 结合二进制 PSO 方法, 提出了一种基于 SVM 的快速特征选择方法 (fast SVM-based feature selection method FSVMSF)。

1 基本概念

1.1 特征选择

特征选择问题可以描述为在 D 维的特征集中, 选择一组 d 维的特征子集使分类器的分类性能最好。有 2 个问题需要解决: 一是选择的标准, 即评价准则, 用来评价不同特征子集的性能; 二是搜索策略, 在 d 确定的情况下, 共有 C_D^d 种可能的组合。如果把所有可能的组合都计算一遍, 计算量是巨大的。因此需要有一种有效的搜索方法, 在允许的时间能找出一组满足要求的特征。

特征选择方法可以分成 2 类: 一类是 Filter 方法; 另一类是 Wrapper 方法。Filter 方法是一种高效的方法, 通常选用信息熵、相关性、类间 (内) 距离等做为评价标准, 可以快速去除大类的噪声, 但是当分类器和特征的关联较大时, 此类方法并不能保证选出一个较好的特征子集; Wrapper 方法虽然在速度上比 Filter 方法慢, 但此类方法与相应的学习算法是紧密相连的, 因此避免了 Filter 方法会出现选出的特征子集不理想的情况, 选出的特征子集效果更好, 规模更小。现在, 越来越多的研究者选用了此类方法进行特征选择^[10]。为了在进行特征选择的同时, 获得尽可能高的分类精度, 本文选择 Wrapper 方法进行特征选择。

1.2 粒子群优化算法

PSO 最早是由 Eberhard 和 Kennedy 在 1995 年提出的^[11]。PSO 和其它群算法一样, 通过群中每个个体不停的运动来搜索最优解。每个粒子由自己当前的最优解 P_{best} 和所有粒子的最优解 G_{best} 两方面决定它的运动方向。

每一个粒子代表了 D 维解空间中的一个点, 其下一个位置由自己的当前位置和速度所决定。第 i 个粒子在第 t 次迭代下的位置可以由 $X_i^t = (x_i^1, x_i^2, \dots, x_i^D)$ 表示。而速度 $V_i^t = (v_i^1, v_i^2, \dots, v_i^D)$ 也是 D 维空间中的一个矢量。

设 $P_i^t = (p_i^1, p_i^2, \dots, p_i^D)$ 表示 X_i 的当前最优解, $P_g^t = (p_g^1, p_g^2, \dots, p_g^D)$ 表示全局最优解。则第 t+1 次迭代中, 所有粒子的位置和速度更新公式如下:

$$V_{id}^{t+1} = V_{id}^t + \zeta [r_1 (P_{id}^t - X_{id}^t) + \zeta_2 r_2 (P_{gd}^t - X_{id}^t)]$$
 (1)

$$X_{id}^{t+1} = X_{id}^t + V_{id}^{t+1}$$
 (2)

其中, ζ, ζ_2 是学习因子, 通常在 [0 2] 之间取值, r_1, r_2 是 [0 1] 间的随机数, $d=1, 2, \dots, D$

PSO 可以用来处理连续空间中的搜索问题, 但是并不能用来解决优化组合问题。因此, Eberhard 和 Kennedy 在 1997 年又提出了用于解决离散问题的 BPSO 算法^[12] (binary particle swarm optimization algorithm, BPSO)。在 BPSO 中, X_{id} 取 0 或 1, V_{id} 为 X_{id} 取 1 的概率, 通过转换函数限制在 [0 1] 之间, 通常都采用 Sigmoid 函数。BPSO 的粒子更新公式:

$$X_{id}^t = \begin{cases} 1, & \text{rand}() \leq S(V_{id}^t) \\ 0, & \text{otherwise} \end{cases}$$
 (3)

为了克服 Wrapper 特征选择速度慢的缺点, 通过启发式搜索方法来提高 Wrapper 的速度是一个重要的研究方向。PSO 算法作为一种启发式搜索方法, 有算法简洁, 参数少, 收敛快等特点, 故本文使用 BPSO 算法做为特征选择中的搜索算法。

1.3 支持向量机

Vapnik 等人在多年研究统计学习理论的基础上, 对线性分类器提出了一种新的最佳设计准则。其原理以线性可分为基础, 扩展到线性不可分的情况, 甚至扩展到使用非线性函数, 这种分类器被称为支持向量机。

设样本集为 $(X_i, Y_i), i=1, 2, \dots, N$, N 为训练集样本数量。 $X_i \in R^d$ 为样本的特征, $Y_i \in \{-1, 1\}$ 为样本类别。 SVM 将样本映射到一个更高维的空间里, 在这个空间里建立 1 个分类超平面。在分开样本的超平面的 2 边建有 2 个互相平行的超平面。分类超平面使 2 个平行超平面的距离最大化。 d 维空间中的分类超平面方程为

$$\langle w^{\circ} \cdot x_i \rangle + b = 0 \quad i = 1, 2, \dots, N$$
 (4)

如果训练样本是线形可分的, 那么 w 和 b 又可由下式表示:

$$y_i (\langle w^{\circ} \cdot x_i \rangle + b) - 1 \geq 0 \quad i = 1, 2, \dots, N$$
 (5)

因此, 满足式 (5) 且使 $\|w\|^2$ 最小的超平面就是分类超平面。求分类超平面的问题最终可以表示成对 α_i 求 $Q(\alpha)$ 最大值的约束优化问题:

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\langle x_i^{\circ} \cdot x_j^{\circ} \rangle), \quad \sum_{i=1}^N y_i \alpha_i = 0, \quad \alpha_i \geq 0$$
 (6)

如果有 $\alpha^0 = (\alpha_1^0, \dots, \alpha_N^0)$ 为使式 (6) 最大的解, 那么 $w^0 = \sum_{i=1}^N \alpha_i^0 y_i x_i^{\circ}$, $\alpha_i^0 > 0$ 的样本就是支持向量。 b^0 可由任意一个支持向量通过式 (5) 求得。而最后得到的分类超平面的分类函数为

$$f(x) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i (\langle x_i^{\circ} \cdot x \rangle + b^0) \right\}.$$
 (7)

在大多数情况下, 样本并不是线形可分的。因此, SVM 需要将样本映射到高维特征空间中来分类。只要一个函数 $K(x, y)$ 满足 Mercer 条件, 就可以被定义成核函数进行映射。此时, 式 (6) (7) 分别变为

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j K(x_i, x_j),$$
 (8)

$$f(x) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i K(x_i, x) + b^0 \right\}.$$
 (9)

2 基于 SVM 的快速特征选择方法 (FSVMFS)

现有的 Wrapper 特征选择方法, 在每次对特征子集进行评价时, 都要重新训练分类器, 计算量大, 如果能减少这部分的运算, 将大大减少 Wrapper 方法的时间开销。

在 SVM 的分类过程中, 由式 (9) 可知, 由每个样本的权重 α_i 和偏移量 b 就可以对样本进行分类, 与特征维数不相关。利用该特性, 可以使用训练好的 SVM 分类器对特征选择后的样本进行分类。本文在此基础上, 结合 PSO 算法提出了一种基于 SVM 的快速特征选择方法。

2.1 FSVMFS 的核函数及其特性

对于线性不可分的情况, 可以设法通过非线性变换转换为另一个空间中的线性问题, 在这个变换空间中求分类超平面。由式 (8) (9) 可知, 无需知道采用的非线性变换的形式, 只需要知道它的内积运算即可。变换空间的内积运算 $K(x, y)$ 就是 SVM 分类器的核函数。

设 SVM 是在样本集 $X = (x_1, \dots, x_N)$ 上训练后得到的解为 $\alpha^0 = (\alpha_1^0, \dots, \alpha_N^0)$, b^0 , $\sigma = [\sigma_1, \dots, \sigma_D]$ 为对角矩阵, $\sigma_i = \{0, 1\}$ 表示在 D 维特征空间上的一个特征子集。那么, 在这个特征子集上的分类函数为

$$f(x) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i K(\sigma x_i, \sigma x) + b^0 \right\}.$$
 (10)

现在常用的 SVM 核函数包括:

$$K_{\text{lin}}(x, y) = \langle x^{\circ} \cdot y \rangle,$$
 (11)

$$K_{\text{poly}}(x, y) = (\gamma \langle x^{\circ} \cdot y \rangle + 1)^d,$$
 (12)

$$K_{\text{RBF}}(x, y) = \exp(-\gamma \|x - y\|^2).$$
 (13)

由矩阵运算的性质可知, 对于对角矩阵 σ , 有 $\sigma x^{\circ} \cdot \sigma y = x^{\circ} \cdot y$ 。利用这个性质, 将线性核与多项式核分别代入式 (12) 后, 可以得到核函数的以下性质:

$$f(x) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i (\sigma x_i^{\circ} \cdot \sigma x) + b^0 \right\} = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i (\langle x_i^{\circ} \cdot \sigma x \rangle + b^0) \right\} = f(\sigma x),$$
 (14)

$$f(x) = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i ((\gamma (\sigma x_i^{\circ} \cdot \sigma x) + 1)^d) + b^0 \right\} = \text{sign} \left\{ \sum_{i=1}^N \alpha_i^0 y_i ((\gamma (\langle x_i^{\circ} \cdot \sigma x \rangle + 1)^d) + b^0) \right\} = f(\sigma x).$$
 (15)

若样本 x 在特征选择后为 x' , 由以上两式可知, 对于一个 SVM 分类器, 在 x' 和在样本 σx 上的分类结果是一样的。因为使用原有的 D 维特征空间上的支持向量可以对 σx 分类, 计算 $f(\sigma x)$ 比计算 $f(x)$ 更简

单,所以可以在 D 维特征空间上,通过 $\alpha^0=(\alpha_1^0, \cdots, \alpha_l^N)$ 和 b^0 对 $\sigma(X)$ 进行分类而不重新训练分类器。即对于使用线形核或多项式核的 SVM分类器,可以将一个在 D 维特征空间上训练得到的 SVM对经过特征选择后的样本集进行分类。

2.2 FSVMFS算法

基于上述思想,本文提出了一种 SVM快速特征选择方法:

设样本集为 $(x_i, y_i), i=1, 2, \cdots, N$ N 为训练集样本数量, $x_i \in R^d, y_i \in \{-1, 1\}$ 。粒子 $a=(a_1, a_2, \cdots, a_d), a_i$ 表示特征 i 是否被选上。进行特征选择时,不仅希望得到的特征子集规模小,而且分类器性能要尽可能高。所以采用如式 (16)所示的粒子适应度函数定义:

$$fitness=P-\sum_{i=1}^d a_i/d,$$

(16)

其中, P 是 SVM分类器的分类精度。

基于 SVM的快速特征选择算法的具体流程如下:

- (1)初始化粒子群,在训练样本集 X 上训练得到 SVM分类器,令 $i=0$
- (2)对于每一个粒子 $x_i=(x_1^i, x_2^i, \cdots, x_d^i)$ 计算其表示的特征子集上的新样本 $x_i^j Y$
- (3)对样本集 $x_i^j Y$ 进行分类,计算适应值 $fitness$
- (4)如果新适应值大于局部最优解,令粒子 X_i 的局部最优值 $P_i=x_i$
- (5)如果新适应值大于全局最优解,令全局最优值 $P_g=x_i$
- (6)计算粒子移动速度 $v_i^{j+1}=v_i^j+\zeta_1(P_i-x_i^j)+\zeta_2(P_g-x_i^j)$
- (7)计算粒子新的位置 $x_i^{j+1}=\begin{cases} 1 & \text{rand}() \leq S(v_i^{j+1}) \\ 0 & \text{otherwise} \end{cases}$
- (8)自增,如果没有到最大迭代次数,则转 (2),否则算法结束。

3 实验结果与分析

3.1 实验数据

为了对算法的有效性进行评估,在多个 UC的数据集^[13]上进行了实验研究。表 1给出了实验中用到的数据集的信息。

表 1 实验中用到的 UC 数据集
Table 1 UCI datasets used in the experiment

| 编号 | 数据集 | 类别数 | 样本数 | 特征数 |
|----|---------------|-----|-------|-----|
| 1 | Ionosphere | 2 | 351 | 34 |
| 2 | Kr vs kp | 2 | 3 196 | 36 |
| 3 | Satapg(heart) | 2 | 270 | 13 |
| 4 | Soybean | 19 | 683 | 35 |
| 5 | Wine | 3 | 178 | 13 |
| 6 | Sonar | 2 | 208 | 60 |

本实验使用了 K 倍交叉验证方法,实验中 K 取 10。数据集被随机等分成 10份,每份都保持类别比尽量不变。每次取其中 1份作为测试集,另 9份合并为训练集。

本实验对一些参数的设置如下: $\zeta_1=\zeta_2=2$ 粒子数 $PN=20$ 最大迭代次数 $N=250$ 。

3.2 与无特征选择的 SVM方法的比较

为了考察 FSVMFS方法对于提高分类精度的有效性,实验一对未做特征选择的 SVM和 FSVMFS方法进行了比较,实验结果如表 2。从表 2中可以看出,FSVMFS方法与未做特征选择的 SVM方法相比,在所有数据集上,经过特征选择之后,都能够大幅提高分类器的分类精度。最少能提高 1.8%的精度,最多则达 16%。说明 FSVMFS方法对于提高分类精度是有效的。

表 2 特征选择对分类精度的影响
Tabel2 Inference of feature selection on classification accuracy

| 数据集 | 分类精度 /% | |
|----------------|----------|------------------|
| | FSVMFS | SVM(无特征选择) |
| Ionosphere | 88.034 2 | 95.735 3 |
| Kr vs kp | 94.180 2 | 96.048 3 |
| Saaplog(heart) | 84.074 1 | 90.946 5 |
| Soybean | 93.704 2 | 97.205 2 |
| Wine | 95.505 6 | 100.000 0 |
| Sonar | 80.288 5 | 96.190 5 |

3.3 与相关工作的比较

实验二对 FSVMFS方法和 ShihWeiLi的方法^[9]进行了比较。ShihWeiLi的方法使用 RBF核作为 SVM核函数,同时在算法中对 SVM的两个参数 γ 和 C 进行了调整。FSVMFS方法选用线性核作为 SVM的核函数。两种方法都使用了 BPSO进行特征选择。为了方便比较,对部分实验数据进行归一化。实验结果如表 3。从表 3可以看出,FSVMFS方法在 3 个数据集上分类精度优于 ShihWeiLi方法,在一个数据集上相同,有两个数据集的分类精度不如 ShihWeiLi的方法。考虑到 ShihWeiLi的方法对 SVM的参数也进行了调整,从总体上说,FSVMFS的分类精度与 ShihWeiLi的方法持平,但大大降低了算法运行的时间。

表 3 两种算法的分类性能比较
Table 3 Classification performance of two algorithms

| 数据集 | FSVMFS | | | ShihWeiLi的方法 | | |
|----------------|-----------------|------|------|-----------------|------|----------|
| | 分类精度 | 特征数 | 耗时 | 分类精度 | 特征数 | 耗时 |
| Ionosphere | 95.735 3 | 9.9 | 1.00 | 96.111 1 | 3.9 | 11.008 2 |
| Kr vs kp | 96.048 3 | 9.5 | 1.00 | 97.777 2 | 8.8 | 33.906 2 |
| Saaplog(heart) | 90.946 5 | 4.8 | 1.00 | 90.123 5 | 2.7 | 8.013 0 |
| Soybean | 97.205 2 | 18.2 | 1.00 | 97.174 7 | 10.5 | 8.699 6 |
| Wine | 100.000 0 | 3.9 | 1.00 | 100.000 0 | 2.1 | 15.766 6 |
| Sonar | 96.190 5 | 5.0 | 1.00 | 95.673 1 | 5.4 | 10.087 9 |

为了分析两种算法运算时间差异的具体原因,给出了算法运行时间的比较,如表 4。从表 4中可以看出,ShihWeiLi方法将大量的时间耗在了训练分类器上,而 FSVMFS方法只消耗了很少的时间用来训练分类器。在实验数据集上,ShihWeiLi方法的总耗时平均是 FSVMFS的 14.5 倍。此外,由于 FSVMFS在每次分类时,只需要改变测试集,而 ShihWeiLi方法不仅要改变测试集,还要改变训练集。所以,消耗在改变特征集上的时间也比 FSVMFS要多。综合上述优点,使得 FSVMFS比 ShihWeiLi方法速度更快。

表 4 两种算法运行时间比较
Tabel4 Time consumption of two algorithms

| 数据集 | FSVMFS | | | | ShihWeiLi的方法 | | | |
|----------------|---------|---------|---------|---------|--------------|---------|----------|---------|
| | 总耗时 | 改变特征集 | 训练分类器 | 计算适度值 | 总耗时 | 改变特征集 | 训练分类器 | 计算适度值 |
| Ionosphere | 1.000 0 | 0.254 9 | 0.011 6 | 0.733 5 | 11.008 2 | 1.6541 | 8.582 6 | 0.771 5 |
| Kr vs kp | 1.000 0 | 0.105 8 | 0.008 8 | 0.885 4 | 33.906 2 | 1.154 5 | 29.653 8 | 3.097 9 |
| Saaplog(heart) | 1.000 0 | 0.095 7 | 0.292 5 | 0.611 8 | 8.013 0 | 0.388 5 | 6.853 7 | 0.707 8 |
| Soybean | 1.000 0 | 0.151 3 | 0.004 6 | 0.844 1 | 8.699 6 | 0.624 1 | 7.308 7 | 0.766 8 |
| Wine | 1.000 0 | 0.165 9 | 0.165 2 | 0.668 9 | 15.766 6 | 0.846 8 | 13.673 0 | 1.246 8 |
| Sonar | 1.000 0 | 0.674 7 | 0.014 8 | 0.310 5 | 10.087 9 | 3.006 9 | 6.429 9 | 0.651 1 |

4 结束语

本文提出了一种基于 SVM的快速特征选择方法。该方法在不损失分类精度的情况下,使得进行特征选择的时间大大减少。如何进一步提高其分类精度,是需要进一步研究的问题。

参考文献:

[1] VAPNIK V N. An overview of statistical learning theory[J]. IEEE Transactions on Neural Networks, 1999, 10(5): 988-999

[2] CRISTIANINI N, SHAW E-TAYLOR J. Introduction to support vector machine and other kernel based learning machine[M]. Cambridge: Cambridge University Press, 2000

[3] VAPNIK V N, VAPNIK V. Statistical learning theory[M]. New York: Wiley New York, 1998

[4] WESTON J, MUKHERJEE S, CHAPELLE O, et al. Feature selection for SVMs[J]. Advances in Neural Information Processing Systems, 2001, 668-674

[5] WANG L, YU J. Fault feature selection based on modified binary PSO with mutation and its application in chemical process fault diagnosis[J]. Lecture Notes in Computer Science, 2005, 3612: 832

[6] DE FALCO J, DELLA CIPPA A, TARANTINO E. Facing classification problem with Particle swarm optimization[J]. Applied Soft Computing Journal, 2007, 7(3): 652-658

[7] MAO K Z. Feature subset selection for support vector machines through discriminative function pruning analysis[J]. IEEE Transactions on Systems, Man, and Cybernetics: Part B, 2004, 34(1): 60-67

[8] SINDHWANI V, RAKSHIT S, DEODHARE D, et al. Feature selection in MLPs and SVMs based on maximum output information[J]. IEEE Transactions on Neural Networks, 2004, 15(4): 937-948

[9] LINSW, YING K C, CHEN S C, et al. Particle swarm optimization for parameter determination and feature selection of support vector machines[J]. Expert Systems with Applications, 2008, 35(4): 1817-1824

[10] SHMAK, TODORIKIM, SUZUKI A. SVM-based feature selection of latent semantic features[J]. Pattern Recognition Letters, 2004, 25(9): 1051-1057

[11] KENNEDY J, EBERHART R C. Particle swarm optimization[C] // Proceedings of IEEE International Conference on Neural Networks, Piscataway, New Jersey, USA, IEEE Computer Society Press, 1995, 4: 1942-1948

[12] KENNEDY J, EBERHART R C. A discrete binary version of the particle swarm algorithm[C] // Proceedings of the IEEE International Conference on Systems, Man and Cybernetics, Orlando, Florida, USA, IEEE Computer Society Press, 1997, 5: 4104-4108

[13] BLAKE C L, MERZ C J. UCI Repository of machine learning databases[DB/OL]. Irvine, California, USA: University of California, Department of Information and Computer Science, 1998[2010-03-20]. <http://www.ics.uci.edu/ml/Repository.html>

(编辑: 孙培芹)

(上接第 59 页)

[6] OHSHIMA M, ZHONG N, YAO Y Y, et al. Relational peculiarity oriented mining[J]. Data Mining and Knowledge Discovery, 2007(15): 249-273

[7] HODGE V J, AUSTIN J. A survey of outlier detection methodologies[J]. Artificial Intelligence Review, 2004, 22(2): 85-126

[8] CHANDOLA V, BANERJEE A, KUMAR V. Anomaly detection: a survey[J]. ACM computing Survey, 2009, 41(3): 1-58

[9] 陈斌, 陈松灿, 潘志松, 等. 异常检测综述[J]. 山东大学学报: 工学版, 2009, 39(6): 13-23
CHEN Bin, CHEN Songcan, PAN Zhisong, et al. Survey of outlier detection technologies[J]. Journal of Shandong University Engineering Science, 2009, 39(6): 13-23

[10] YANG J, ZHONG N, YAO Y Y, et al. Local peculiarity factor and its application in outlier detection[C] // Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Nevada, USA, the ACM, 2008, 776-784

[11] 薛安荣. 空间离群点挖掘技术的研究[D]. 镇江: 江苏大学, 2008
XUE Anrong. Study on technology for spatial outlier detection[D]. Zhenjiang, China: Jiangsu University, 2008

[12] ZHONG N, YAO Y Y, OHSHIMA M, et al. Interest-ness, peculiarity and multi-database mining[C] // Proceedings of the 2001 IEEE International Conference on Data Mining, Washington, DC, USA, IEEE Computer Society, 2001, 566-573

[13] YANG J, ZHONG N, YAO Y Y, et al. Peculiarity analysis for classifications[C] // Proceedings of the 2009 IEEE International Conference on Data Mining, IEEE Computer Society, Washington, DC, USA, IEEE Computer Society, 2009, 607-616

[14] WU M, JERMANE C. Outlier detection by sampling with accuracy guarantees[C] // Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, New York, NY, USA, the ACM, 2006, 767-772

(编辑: 陈斌)