



当代学术文丛



量子聚类算法的研究及其在异常检测中的应用
基于语义属性数据离群聚类的异常检测
基于趋势的时间序列数据析取之异常检测算法
基于语义属性数据核分类方法的异常检测

Anomaly Detection Algorithms in non-Numerical Data

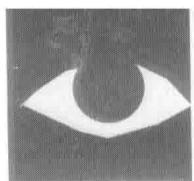
非数值属性数据异常检测算法

李志华 张海涛
孙雅 耿振民 编著



江西人民出版社
Jiangxi People's Publishing House
全国百佳出版社

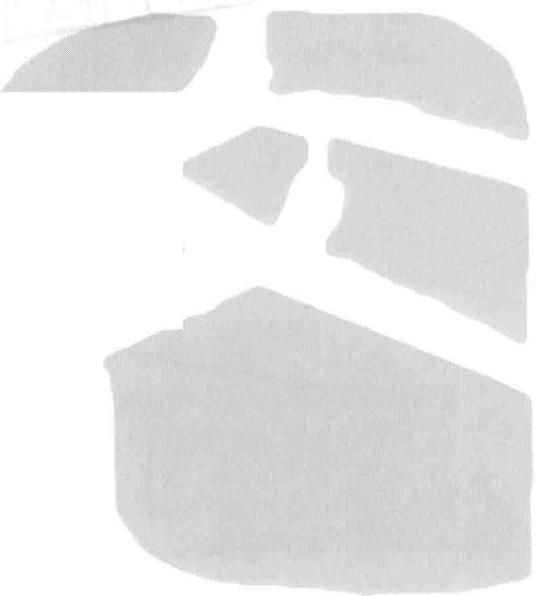
7837



非数值属性数据异常检测算法

Anomaly Detection Algorithms in non-Numerical Data

李志华 张海涛
孙 雅 耿振民 编著



江西人民出版社
Jiangxi People's Publishing House
全国百佳出版社

图书在版编目(CIP)数据

非数值属性数据异常检测算法 / 李志华等编著.
—南昌: 江西人民出版社, 2015.12
ISBN 978-7-210-08140-1

I . ①非… II . ①李… III . ①数据检测 IV .
①TP274

中国版本图书馆 CIP 数据核字(2015)第 305750 号

非数值属性数据异常检测算法

李志华 张海涛 孙雅 耿振民 编著

责任编辑:徐 昊 李月华

出 版:江西人民出版社

发 行:各地新华书店

地 址:江西省南昌市三经路 47 号附 1 号

发行部电话:0791-88629871

编辑部电话:0791-86898815

邮 编:330006

网 址:www.jxpph.com

E-mail:gjzx999@126.com web@jxpph.com

2015 年 12 月第 1 版 2015 年 12 月第 1 次印刷

开 本:787 毫米×1092 毫米 1/16

印 张:14.75

字 数:200 千字

ISBN 978-7-210-08140-1

赣版权登字-01-2015-895

版权所有 侵权必究

定 价:36.00 元

承 印 厂:南昌三联印务有限公司

赣人版图书凡属印刷、装订错误,请随时向承印厂调换

前 言

计算机技术、通信技术、网络技术的迅猛发展将人类带进了大数据时代。

数据挖掘 (Data Mining) 是对所观测的数据对象进行分析,以便找到这些数据对象中便于理解且有用的各种模式、新规律、新知识的过程。近年来,在该领域虽然新的和改进的算法不断涌现、应用领域也正在逐渐扩大,但是同样面临许多新问题,其中更加丰富的数据类型如流、时间序列、图、音频、视频和图像等,这些大多数是非数值属性数据,需要改变传统的数据挖掘算法去适应针对这些数据类型数据的挖掘需要。当然,丰富数据挖掘领域中的算法和技术,以便适应更加广泛的应用需求,依然是数据挖掘发展的主旋律。

异常 (Anomaly) 指的是在数据集中偏离大部分数据,即这些偏离数据与大部分数据存在明显的差异。人们研究发现这些偏离的数据并非全部由随机因素或偶然误差而产生,而是产生于某些完全不同的机制,所以这些异常数据中通常隐藏有价值的重要信息和可以帮助人们深入认知被检测对象的重要线索。异常检测 (Anomaly Detection) 是数据挖掘领域研究的最基本问题之一,主要用来发现数据集中与其他数据存在明显不同的异常。

虽然,数据挖掘研究与应用在近年来得到了前所未有的快速发展。但是不论是数据挖掘还是异常检测都需要进一步与数理统计、神经网络、模式识

别、机器学习、知识库系统、高性能计算和可视化处理技术等学科领域相结合,特别是与这些前沿学科中出现的新理论和新技术紧密结合,形成学科交叉,面对大数据的挑战。

为了进一步推动异常检测技术向非数值属性数据等其他领域拓展,以适应不断涌现的新型应用,并反映该领域的最新研究成果,让相关专业的研究生、高年级本科生和工程技术人员比较深入地理解并快速掌握新兴的异检测技术,特编写此书。本书的另外两名作者是本人指导的硕士,他们协助完成本书部分章节的工作。本书撰写过程中融合了作者在攻读博士学位期间所做的主要研究成果和近几年在时间序列异常检测领域的研究工作,同时也吸纳了国内外具有代表性的研究成果。并且在撰写过程中,围绕异常检测这个主题,从不同角度进行展开,力求以点带面。

本书共分八章。第一章主要介绍异常检测的相关概念、检测方法的研究现状、检测方法的评估指标和经典的异常检测方法等;第二章主要介绍了常见非数值属性数据如语义数据、异构数据、时间序列数据等的概念、特点、表达形式、度量方法等;第三章主要介绍了量子理论、量子理论中的粒子分布机制与数据样本分布的同质性分析,以及基于量子理论的聚类、分类等数据挖掘算法和基于这些数据挖掘算法的异常检测方法;第四章主要通过分析离群数据与异常的内在关系,介绍了一种基于度量的离群数据聚类算法和基于该算法的异常检测方法;第五章介绍了语义属性数据在核空间中的分布情况和内积计算方法、语义属性数据的核分类方法,以及基于该方法的异常检测方法;第六章介绍了信息论相关的知识和用熵表达数据样本集中的结构特征,通过结构的挖掘实现数据样本的聚类,以及基于该算法的异常检测方法;第七章介绍了时间序列的表示方法、析取方法,以及基于这些数据析取方法的异常检测方法;第八章在介绍了时间序列的聚类和分类方法基础上,结合时间序列异常检测的特殊性,进一步介绍了时间序列的可视化异常检测技术。

感谢实验室的研究生张海涛、孙雅、刘庭绪、胡振宇、陈超群等同学,他们

有的参加了部分章节的撰写,有的做了大量的编辑排版工作。

由于作者水平有限,书中难免存在不妥之处,恳请广大读者、同行和专家批评指正。

李志华

2015 年 10 于江南大学

|目 录|

前 言	—	1
第一章 绪 论	—	1
1.1 引 言	—	1
1.2 异常检测相关概念	—	2
1.3 异常检测方法	—	2
1.4 异常检测的评估指标	—	7
1.5 本章小结	—	8
参考文献	—	9
第二章 相关基础知识	—	12
2.1 引 言	—	12
2.2 常见非数值属性数据及度量	—	13
2.2.1 语义属性数据	—	13
2.2.2 异构属性数据	—	13
2.2.3 数据样本的度量	—	14
2.2.3.1 距离度量方法	—	15
2.2.3.2 相似性测度	—	16
2.2.3.3 非度量方法	—	17

2.2.3.4 异构数据的度量	— 19
2.3 时间序列数据及度量	— 20
2.3.1 时间序列的表示	— 20
2.3.1.1 离散傅里叶变换	— 21
2.3.1.2 离散小波变换	— 22
2.3.1.3 分段聚集近似	— 23
2.3.1.4 符号化聚集近似	— 24
2.3.1.5 其他方法	— 25
2.3.2 时间序列的相似性度量	— 27
2.3.2.1 闵可夫斯基距离	— 27
2.3.2.2 动态时间弯曲距离	— 28
2.3.2.3 编辑距离	— 31
2.3.2.4 最长公共子序列	— 32
2.3.2.5 MINDIST 距离	— 33
2.3.2.6 其他方法	— 33
2.4 本章小结	— 34
参考文献	— 35
 第三章 量子聚类算法的研究及其在异常检测中的应用	— 40
3.1 引言	— 40
3.2 量子力学的基本理论	— 42
3.2.1 概率波函数	— 42
3.2.2 薛定谔方程	— 42
3.2.3 量子势能	— 43
3.3 QC 算法与 FCM 算法的比较研究	— 44
3.3.1 量子聚类算法	— 44

3.3.1.1 量子聚类的理论根据	44
3.3.1.2 量子聚类算法	46
3.3.1.3 量子聚类算法的特点	46
3.3.2 FCM 算法及其特点	47
3.3.2.1 FCM 算法简介	47
3.3.2.2 FCM 算法的特点	48
3.3.3 FCM 算法的一种量子理论解释	48
3.3.4 QC 算法与 FCM 算法的仿真实验比较	49
3.3.4.1 仿真实验及分析	49
3.3.4.2 QC 算法与 FCM 算法的实验比较结果	56
3.3.5 结论	56
3.4 基于核宽度调节参数估计的量子聚类算法	57
3.4.1 算法概述	57
3.4.2 PeQC 算法	58
3.4.2.1 算法中参数的分析	58
3.4.2.2 δ 的估计	58
3.4.2.3 PeQC 算法	59
3.4.2.4 PeQC 聚类算法分析	59
3.4.3 仿真实验及分析	60
3.4.4 算法小结	62
3.5 语义属性数据模糊量子聚类算法	63
3.5.1 算法概述	63
3.5.2 距离量子势能	64
3.5.3 语义属性数据量子聚类算法	65
3.5.3.1 语义属性数据的相异性度量	65
3.5.3.2 NQC 算法	66

3.5.3.3 NQC 算法的时空复杂度分析	— 67
3.5.3.4 聚类有效性分析	— 67
3.5.4 仿真实验及分析	— 69
3.5.4.1 仿真实验	— 69
3.5.4.2 soybean disease 样本数据集的聚类结果分析	— 72
3.5.5 算法小结	— 72
3.6 基于量子聚类的异常检测方法	— 73
3.6.1 方法概述	— 73
3.6.2 量子势能中的相异性度量分析	— 74
3.6.3 异构数据的距离量子聚类算法	— 74
3.6.3.1 语义属性的相异性度量	— 74
3.6.3.2 异构样本间的 Mahalanobis 距离	— 74
3.6.3.3 距离量子聚类算法	— 75
3.6.4 基于 MDQC 算法的异常检测方法	— 76
3.6.5 仿真实验及分析	— 76
3.6.5.1 样本的选择及预处理	— 76
3.6.5.2 实验结果及分析	— 77
3.6.6 异常检测方法小结	— 80
3.7 本章小结	— 81
参考文献	— 81
第四章 基于语义属性数据离群聚类的异常检测	— 86
4.1 引言	— 86
4.2 离群聚类算法介绍	— 88
4.2.1 主观发现方法	— 88

4.2.2 客观发现方法	— 89
4.3 离群聚类算法及分析	— 90
4.3.1 样本的相异性度量	— 90
4.3.2 离群聚类算法	— 91
4.3.3 算法抗离群点干扰能力分析	— 92
4.3.4 算法的时空复杂度分析	— 95
4.4 基于离群聚类的异常检测研究	— 95
4.4.1 检测方法概述	— 95
4.4.2 异常检测实验及分析	— 96
4.5 本章小结	— 98
参考文献	— 99
 第五章 基于语义属性数据核分类方法的异常检测	— 101
5.1 引言	— 101
5.2 支撑向量机简介	— 102
5.3 核方法分析及支撑向量机中的核函数	— 106
5.3.1 核方法分析	— 106
5.3.2 支撑向量机中的核函数	— 108
5.4 语义属性数据的核分类方法及分析	— 110
5.4.1 样本的相异性度量	— 110
5.4.2 异构属性样本的核分类方法	— 111
5.4.3 标准样本集的仿真实验及分析	— 112
5.5 基于语义属性数据分类方法的异常检测研究	— 116
5.5.1 检测方法概述	— 116
5.5.2 异常检测实验及分析	— 117
5.6 本章小结	— 120

参考文献	—	120
第六章 基于结构熵聚类的异常检测	—	124
6.1 引言	—	124
6.2 信息论基础	—	125
6.2.1 自信息	—	125
6.2.2 信息熵	—	126
6.2.3 互信息	—	128
6.3 连续属性的离散化算法	—	129
6.3.1 离散化问题简述	—	130
6.3.2 离散化算法	—	130
6.4 结构熵聚类算法	—	132
6.4.1 算法概述	—	132
6.4.2 相异性度量	—	133
6.4.2.1 语义属性相异性度量的计算	—	133
6.4.2.2 异构距离的计算	—	134
6.4.3 结构熵聚类算法	—	135
6.4.3.1 异构数据结构熵的计算	—	135
6.4.3.2 结构熵聚类算法	—	137
6.4.3.3 结构熵聚类方法的时空复杂度分析及特点	—	138
6.4.3.4 标准样本集的仿真实验及分析	—	139
6.4.4 算法小结	—	141
6.5 基于结构熵聚类的异常检测研究	—	142
6.5.1 检测方法概述	—	142
6.5.2 异常检测实验及分析	—	142

6.6 本章小结	145
参考文献	145
第七章 时间序列的异常检测	148
7.1 引言	148
7.2 基于层次析取的时间序列数据异常检测算法	149
7.2.1 算法概述	149
7.2.2 时间序列的层次表示	149
7.2.2.1 时间序列的层次标记	149
7.2.2.2 层次序列的获取	150
7.2.3 基于层次的时间序列数据析取	151
7.2.4 仿真实验及分析	152
7.2.4.1 实验数据	152
7.2.4.2 时间序列的分类	153
7.2.4.3 实验结果及分析	154
7.2.5 基于 LETSD 的时间序列异常检测	157
7.2.5.1 基于 LETSD 的时间序列异常检测模型	157
7.2.5.2 仿真实验及分析	160
7.3 基于趋势的时间序列数据析取之异常检测算法	165
7.3.1 算法概述	165
7.3.2 时间序列的趋势符号化	166
7.3.3 时间序列的趋势距离	168
7.3.4 仿真实验及分析	169
7.3.4.1 实验数据	170
7.3.4.2 时间序列的聚类	174
7.3.4.3 实验结果及分析	175

7.3.5 基于 TETSD 的时间序列异常检测	—— 178
7.3.5.1 基于 TETSD 的时间序列异常检测模型	—— 178
7.3.5.2 仿真实验及分析	—— 179
7.4 本章小结	—— 181
参考文献	—— 182

第八章 时间序列的可视化异常检测	—— 184
8.1 引言	—— 184
8.2 基于时间序列分类的可视化异常检测算法	—— 185
8.2.1 算法概述	—— 185
8.2.2 时间序列的预处理	—— 185
8.2.3 LMPC 分类方法	—— 188
8.2.3.1 实验数据介绍	—— 188
8.2.3.2 分类评价指标	—— 190
8.2.3.3 实验及结果分析	—— 190
8.2.4 基于 LMPC 的时间序列异常检测方法	—— 195
8.2.4.1 实验数据	—— 195
8.2.4.2 评价标准	—— 197
8.2.4.3 实验结果及分析	—— 197
8.3 基于时间序列聚类的可视化异常检测方法	—— 203
8.3.1 方法概述	—— 203
8.3.2 符号化聚合近似方法的改进	—— 204
8.3.3 基于 LMP_SAX 的时间序列聚类算法	—— 205
8.3.3.1 实验数据	—— 206
8.3.3.2 评价标准	—— 208
8.3.3.3 实验及结果分析	—— 209

8.3.4 基于 LMP_SAX 的时间序列异常检测方法	— 213
8.3.4.1 实验数据	— 214
8.3.4.2 实验结果及分析	— 214
8.4 本章小结	— 220
参考文献	— 220
附录:本书算法的实验环境	— 222

第一章

绪 论

1.1 引 言

随着社会信息化的不断发展,信息技术应用领域的不断拓展,各个应用领域包括经济、医疗、建筑、环境等均积累了越来越多的数据。自 20 世纪 80 年代开始,世界各地的数据总量飞速增长,每年甚至几个月便会增长一倍,然而,面对这些海量级别的复杂数据,如何有效地对这些历史数据进行组织和整合,然后通过数据挖掘^[1](Data Mining)来发现其中隐藏的存在潜在价值的知识和信息成为当前最为迫切的问题,是数据挖掘学科面临的一个巨大的挑战。

所谓数据挖掘是对所观测的数据对象进行分析,拟从大规模的复杂类型数据中挖掘出潜在的、未知的、有价值的知识和信息的过程,其主要目的是发现那些出现频率较高的模式。当今,数据挖掘技术在多种行业、各个领域中均得到了广泛、深入的应用,其中异常检测就是其成功应用的范例之一。基于数据挖掘中各种算法的异常检测方法或检测技术在信息安全领域、网络安全领域、金融领域、医疗领域、气象灾害领域、地质灾害监测领域等都得到了

成功应用,取得了良好的效果,创造了可观的社会效益。如今,异常检测技术已在工业设备诊断、网络安全、自然灾害检测、疾病诊断等国民生产的各行各业中显示出强大的生命力和进一步发展的潜力。

本章首先对异常检测的相关定义、所面临问题、发展现状、评价指标等进行一个比较系统的介绍。

1.2 异常检测相关概念

通常把那些符合期望(正常类或目标类)行为的数据或数据模式称为正常,而那些不符合期望(反类或异常类)行为的数据模式称为异常^{[2][3][4]}。数据挖掘本质上是期望寻找出事物的某种发展规律,所以常常会丢掉那些出现频率很小的模式。而在很多领域中,那些极少出现的模式或者说异常通常比频繁出现的模式有着更重要的科研及实用价值。异常检测能够发现一些真实存在的而又出乎意料的信息,“异常”中蕴含显著的、至关重要的行为信息^[2]。但是在传统数据挖掘领域,由于针对异常的采样代价高昂或者由于采样非常困难,使得对异常数据或异常行为知道的不多,甚至一无所知。然而,异常检测技术的价值潜力不可忽视。截至目前,异常检测在环境监测、网络入侵、信用卡(或保险)欺诈检测、疾病诊断等方面都得到了广泛的应用^{[3][4]}。

1.3 异常检测方法

作为数据挖掘的主要任务之一,异常检测在近年来得到的关注和研究越来越多,研究人员也从不同角度提出了很多方法,从异常检测原理的角度概括如下:

(1) 基于分布的方法

基于分布的方法^[5]是最早被提出来的异常检测方法,主要发展于统计学领域。该类方法假设数据集的分布为已知分布(如正态分布、泊松分布等),