



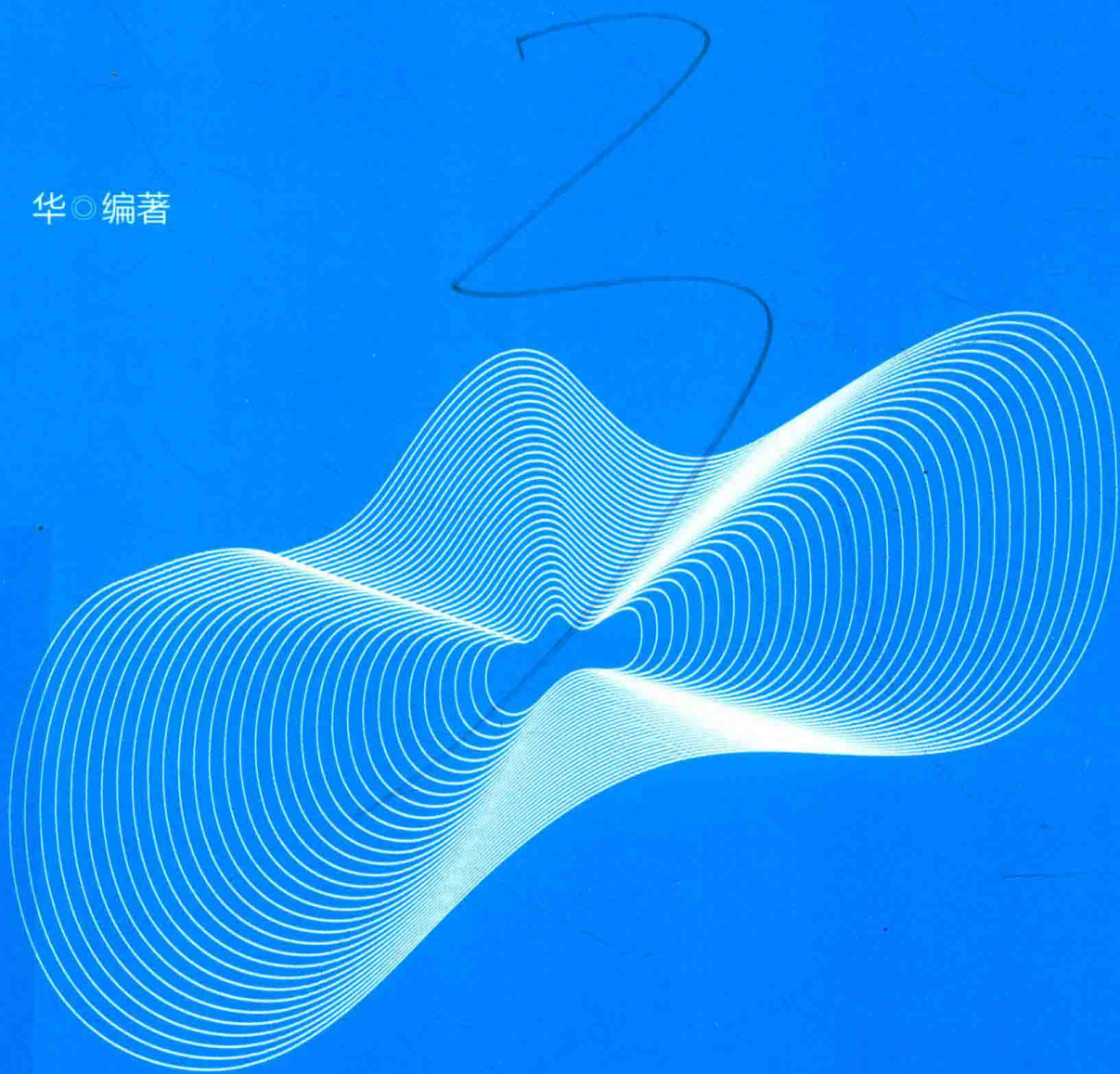
高等院校信息类新专业规划教材

大数据和人工智能技术丛书

EVOLUTIONARY
MACHINE LEARNING

演化机器学习

徐 华◎编著



北京邮电大学出版社
www.buptpress.com



高等院校信息类新专业规划教材
大数据和人工智能技术丛书

演化机器学习

徐 华 编著



北京邮电大学出版社
www.buptpress.com

内 容 简 介

近年来,演化计算作为计算智能中传统的优化技术,已经广泛应用于求解各种数据挖掘问题,形成了一种基于遗传的机器学习新范式学习分类器。一方面,在真实场景中采集的原始数据不可避免地包含着冗余乃至噪声属性的信息,这些不相关的特征将对学习分类器算法的学习性能与计算效率造成负面影响。另一方面,学习分类器以显式规则表示目标概念,在监督学习或强化学习机制的基础上,利用演化算法对规则空间进行搜索,从而完成学习任务。规则空间的有效搜索是影响学习分类器性能的关键。针对上述问题,本书的主要探讨内容:一是学习分类器与特征选择方法,重点是做两者的整合研究,将学习分类器的分类模型构建过程与特征选择的特征子集搜索过程统一集成在基于遗传的机器学习框架下,同时改善分类算法的预测性能与运行效率;二是从提高规则空间的搜索质量出发,着眼于分类问题,介绍了基于分布估计算法的学习分类器。

本书可作为大数据及人工智能等相关专业的教材与参考用书。

图书在版编目(CIP)数据

演化机器学习 / 徐华编著. -- 北京: 北京邮电大学出版社, 2021. 4

ISBN 978-7-5635-6330-2

I. ①演… II. ①徐… III. ①机器学习 IV. ①TP181

中国版本图书馆 CIP 数据核字(2021)第 027485 号

策划编辑: 姚 顺 刘纳新 责任编辑: 刘 颖 封面设计: 七星博纳

出版发行: 北京邮电大学出版社

社 址: 北京市海淀区西土城路 10 号

邮政编码: 100876

发行部: 电话: 010-62282185 传真: 010-62283578

E-mail: publish@bupt.edu.cn

经 销: 各地新华书店

印 刷: 北京玺诚印务有限公司

开 本: 787 mm×1 092 mm 1/16

印 张: 15

字 数: 314 千字

版 次: 2021 年 4 月第 1 版

印 次: 2021 年 4 月第 1 次印刷

ISBN 978-7-5635-6330-2

定价: 45.00 元

· 如有印装质量问题, 请与北京邮电大学出版社发行部联系 ·

大数据和人工智能技术丛书

顾问委员会

吴奇石 黄永峰 吴 斌 欧中洪

编委会

名誉总主编：马少平

总 主 编：许云峰 徐 华

编 委：康艳梅 朱卫平 沈 炜 冷 飏
孙 艺 高 慧 高 崇 刘 刚

总 策 划：姚 顺

秘 书 长：刘纳新

伴随着信息技术和移动互联网技术的发展以及社会经济数字化进程的加快，以计算机、智能硬件为代表的数字化移动设备已经渗透到社会领域的方方面面，正在深刻地改变着传统的工作和生活方式。一方面，网络移动互联技术的普及为信息化革命注入了强大活力，覆盖全球的移动互联网为人们提供了一个方便快捷又廉价的信息发布平台，从而带来了数据规模的爆炸性增长。另一方面，云存储技术与数据库技术的发展也为收集并存储这些海量数据提供了便利。然而，虽有海量的数据，若缺乏强有力的数据分析工具则只能带来信息的缺位。时至今日，面对数据的汪洋大海，人们已经将自己置身于“数据泛滥但信息贫乏”的窘境当中。20世纪80年代末，一门致力于在实现数据分析自动化的交叉学科应运而生，这就是至今方兴未艾的数据挖掘（Data Mining, DM）。简言之，数据挖掘指的是从大规模数据中提取有用的信息（知识），因此也被称为从数据中发现知识（Knowledge Discovery from Data, KDD）。作为一门交叉学科，数据挖掘涉及多个领域，既有负责海量数据存储的数据库技术，也有提供强大的运算能力的高性能云计算技术，不过其中最为核心的是进行知识抽取的数据分析技术。从早期的统计学到模式识别（Pattern Recognition, PR），直至当下更为盛行的机器学习（Machine Learning, ML）、深度学习（Deep Learning），数据挖掘从多个领域汲取养分，增强了本学科海量数据分析的实用性和有效性。

现实世界中各个行业、各个部门对于数据挖掘提出的应用需求五花八门，总的说

来, 数据挖掘任务可以分为两大类: 描述性 (Descriptive) 挖掘与预测性 (Predictive) 挖掘。前者主要挖掘数据集的一般性质, 对应于机器学习领域的无监督学习 (Unsupervised Learning); 后者则根据当前数据进行推断并做出预测, 其中主要用到的是有监督学习 (Supervised Learning) 技术。作为一种典型的预测性任务, 分类 (Classification) 旨在找出描述和区分数据类别 (或称概念) 的分类模型, 以便使用该模型预测类别标签值未知的数据实例, 模型构建过程中主要的可利用信息是已知类别标签值的训练数据集 (Training Dataset)。现实中人们感兴趣的预测值有很大一部分都可以抽象为离散型的类别信息 (如文字识别、故障诊断、入侵检测等), 故而分类成为一种主要的分析手段。为此, 在过去的几十年中, 研究人员也提出了许多应用广泛的分类算法, 包括决策树 (Decision Tree)、朴素贝叶斯分类器 (Naive Bayesian Classifier)、 k -近邻算法 (k -Nearest Neighbor, k -NN)、人工神经网络 (Artificial Neural Network, ANN)、支持向量机 (Support Vector Machine, SVM) 和深度神经网络 (Deep Neural Network, DNN)。

近年来, 一种新的名为学习分类器 (Learning Classifier System, LCS) 的机器学习范式正吸引越来越多研究者的注意力。总的来说, 学习分类器基于规则归纳的思想, 主要致力于解决分类问题。在学习分类器中, 处于中心地位的规则学习单元通常是演化计算 (Evolutionary Computation, EC) 中的遗传算法 (Genetic Algorithm, GA), 故而在一部分文献中又被称为基于遗传的机器学习 (Genetic-Based Machine Learning, GBML)。通过将学习性能指标定义为优化目标函数, 学习分类器实质上将作为学习问题的分类任务转化为传统的优化问题进行求解, 继而基于遗传算法的全局优化搜索能力确保规则知识表示的假设空间中一定强度的随机化搜索, 以期在合理的运算时间内收敛到较优性能的问题解, 从而更好地平衡了算法性能与计算效率这对矛盾。

在实际的分类问题中, 一方面, 在真实场景中采集的原始数据不可避免地包含着冗余乃至噪声的属性信息, 这些不相关的特征将对学习分类器算法的学习性能与计算效率造成负面影响。另一方面, 学习分类器以显式规则表示目标概念, 在监督学习或强化学习机制的基础上, 利用演化算法对规则空间进行搜索, 从而完成学习任务。规

则空间的有效搜索是影响学习分类器性能的关键。针对上述问题，本书从内容体例上分为上、下两篇，各自自成体系，分别以相对独立的方式对学习分类器中的重要内容加以介绍，以便于读者根据需要分别研读。上篇重点介绍了学习分类器与特征选择方法，重点是介绍两者的整合研究内容，将学习分类器的分类模型构建过程与特征选择的特征子集搜索过程统一集成在基于遗传的机器学习框架下，同时改善分类算法的预测性能与运行效率。下篇从提高规则空间的搜索质量出发，着眼于分类问题，介绍了基于分布估计算法的学习分类器。本书的相关内容主要来自作者所在清华大学团队近年来在演化机器学习研究方面的工作积累，力图以比较全面的视角对学习分类器及相关重要问题的研究做一个深入的介绍和探讨。在本书完成之际感谢清华大学计算机系人工智能研究团队成员温贇、杨甲东、袁源、王勃、黄嘉宇、余文梦、孟繁阳等同学为此付出的辛勤努力。

本书的相关工作受国家自然科学基金项目 (No. 61673235, 61175110, 60875073, 60575057) 的持续资助支持，本书的编写工作得到清华大学教育教学改革项目的支持。本书可作为人工智能领域、机器学习领域的学生、研究人员和工程技术人员的教材和参考用书。由于本书相关的内容也处于研究与探讨阶段，不当之处请读者给予指正 (请用电子邮件联系: xuhua@tsinghua.edu.cn)。同时我们持续性的相关研究工作与成果会公开在共享网站: <https://github.com/thuiar>。

徐 华
于清华大学

上篇 演化机器学习——内嵌特征选择的学习分类器

第 1 章 上篇引言	3
1.1 研究背景	3
1.2 上篇主要内容	6
1.3 上篇的结构安排	7
第 2 章 相关工作综述	8
2.1 概述	8
2.2 学习分类器研究	9
2.2.1 进化计算	9
2.2.2 基于遗传的机器学习思想	10
2.2.3 Michigan 式学习分类器研究进展	12
2.2.4 Pittsburgh 式学习分类器研究进展	13
2.3 特征选择方法综述	15
2.3.1 机器学习中的特征选择问题描述	16
2.3.2 特征选择的搜索模型	17
2.3.3 Filter 方法	19
2.3.4 Wrapper 方法	20
2.3.5 Embedded 方法	21
2.4 本章小结	21

第 3 章 基于 Memetic 算法的 Wrapper-Filter 特征选择方法	23
3.1 概述	23
3.2 Memetic 算法	24
3.2.1 Memetic 算法思想起源	24
3.2.2 Memetic 算法框架	25
3.3 基于 MA 的混合式 Wrapper-Filter 特征选择方法	28
3.3.1 算法设计思想	28
3.3.2 算法整体框架	29
3.3.3 全局搜索的 GA-Wrapper	30
3.3.4 局部搜索的 Relief-Filter	32
3.3.5 计算复杂度分析	34
3.4 本章小结	35
第 4 章 基于合作式协同进化内嵌特征选择的学习分类器	36
4.1 概述	36
4.2 协同进化算法	37
4.2.1 协同进化思想起源	37
4.2.2 竞争式协同进化算法模型	38
4.2.3 合作式协同进化算法模型	39
4.3 基于合作式协同进化的学习分类器	41
4.3.1 算法设计思想	41
4.3.2 算法框架	42
4.3.3 分类器演化的 Pittsburgh 式学习分类器算法	44
4.3.4 计算复杂度分析	45
4.4 本章小结	46
第 5 章 算法评估结果与分析	47
5.1 概述	47
5.2 算法比较实验框架	48
5.2.1 Benchmark 数据集	48
5.2.2 性能评估指标	48

5.2.3 实验方法	50
5.3 MFS 算法实验结果及分析	50
5.3.1 算法参数设置	50
5.3.2 实验结果与讨论	51
5.4 CoCoLCS_MFS 算法实验结果及分析	53
5.4.1 算法参数设置	53
5.4.2 实验结果与讨论	53
5.4.3 显著性检验	55
5.5 特征选择对学习分类器的影响分析	55
5.5.1 预测准确率	55
5.5.2 运行时间	56
5.5.3 特征约简率	57
5.5.4 显著性检验	58
5.6 本章小结	59
第 6 章 关于演化搜索的一项扩展性工作——基于混合式 GVNS 算法的多处理器 任务调度研究	60
6.1 概述	60
6.2 多处理器任务调度	61
6.2.1 研究背景	61
6.2.2 多处理器任务调度问题模型	62
6.2.3 调度算法研究进展	64
6.3 基于启发式策略的混合式 GVNS 调度算法的总体设计	67
6.3.1 混合式设计的算法思想	67
6.3.2 算法整体框架	67
6.4 任务优先级定序的启发式策略	69
6.5 全局搜索的遗传算法设计	70
6.5.1 种群个体的基因编码	70
6.5.2 种群初始化与个体适应度评估	70
6.5.3 遗传操作设计	71
6.6 局部搜索的变邻域搜索算法设计	72
6.6.1 算法流程	73

6.6.2 邻域结构设计	74
6.6.3 局部搜索	75
6.7 算法实验结果评估与分析	75
6.7.1 性能评估指标与参数设置	76
6.7.2 确定图上的结果	76
6.7.3 随机图上的结果	80
6.7.4 显著性检验	80
6.7.5 局部搜索有效性分析	82
6.8 本章小结	83
第7章 上篇总结	84
7.1 上篇的主要工作总结	84
7.2 未来研究工作展望	85
下篇 演化机器学习——分布估计的学习分类器	
第8章 下篇引言	89
8.1 研究背景	89
8.2 研究动机	91
8.3 下篇的主要工作	92
8.4 下篇的结构安排	94
第9章 分布估计算法和学习分类器	96
9.1 概述	96
9.2 进化计算简介	96
9.3 遗传算法	97
9.4 分布估计算法	99
9.4.1 分布估计算法的基本流程	100
9.4.2 分布估计算法分类	100
9.5 学习分类器	101
9.5.1 Michigan 式学习分类器	102
9.5.2 Pittsburgh 式学习分类器	103

第 10 章 基于 L1 正则化贝叶斯网络的分布估计算法	106
10.1 概述	106
10.2 L1 正则化贝叶斯网络	107
10.2.1 贝叶斯网络与 L1 正则化	108
10.2.2 候选链接关系建立	110
10.2.3 剪枝搜索	111
10.3 L1 正则化贝叶斯网络与分布估计算法整合	113
10.4 实验设置	114
10.4.1 测试函数	114
10.4.2 参数配置	116
10.5 优化性能比较	117
10.5.1 第一组实验	117
10.5.2 第二组实验	118
10.6 模型比较	125
10.6.1 测试函数的理想模型	125
10.6.2 模型复杂度比较	126
10.6.3 模型准确度比较	127
10.7 本章小结	129
第 11 章 基于分布估计算法的分类器进化算法	131
11.1 概述	131
11.2 算法框架介绍	132
11.3 分类器的知识表示	133
11.4 规则进化	135
11.4.1 规则评价	135
11.4.2 规则重组	135
11.5 规则集进化	137
11.5.1 规则集评价	137
11.5.2 规则集重组	138
11.6 规则与规则集进化整合	139
11.7 实验设置	139

目 录

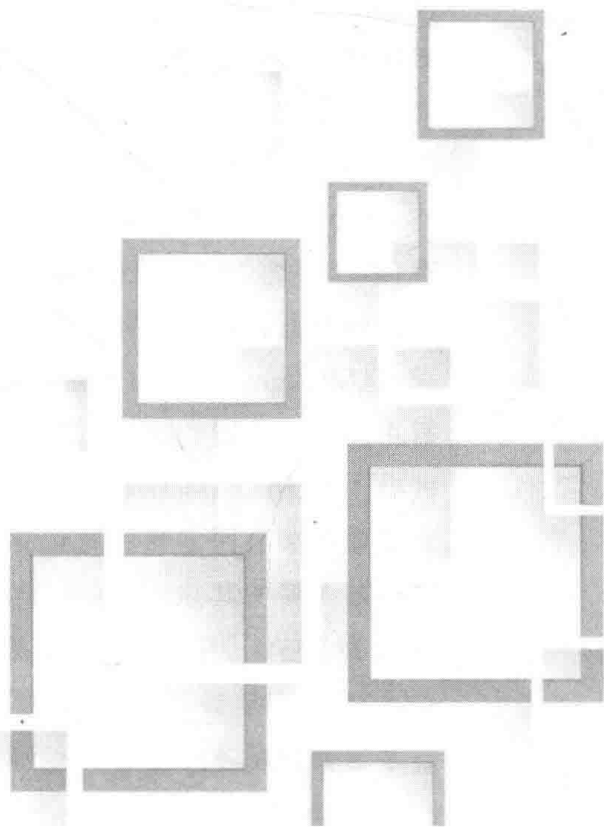
11.7.1	测试数据	141
11.7.2	比较对象	143
11.7.3	评价指标	144
11.7.4	参数设置	145
11.8	实验结果	146
11.8.1	构造问题	146
11.8.2	实际问题	152
11.8.3	参数敏感度分析	157
11.9	本章小结	158
第 12 章	面向学习分类器的嵌入式特征选择算法	159
12.1	概述	159
12.2	本章相关工作	160
12.3	算法框架介绍	161
12.4	嵌入式特征选择算法	162
12.4.1	特征冗余度计算	162
12.4.2	特征关联度计算	163
12.5	特征选择与学习分类器整合	165
12.6	实验设置	167
12.6.1	测试数据	167
12.6.2	比较对象	168
12.6.3	评价指标	169
12.6.4	参数设置	170
12.7	实验结果	170
12.7.1	构造问题	170
12.7.2	实际问题	172
12.7.3	参数敏感度分析	180
12.8	本章小结	181
第 13 章	基于进化纠错输出编码的多类分类算法	182
13.1	概述	182
13.2	纠错输出编码方法简介	183

13.3 编码矩阵设计	184
13.3.1 编码矩阵评价	184
13.3.2 优化方法	187
13.4 解码策略	188
13.5 实验设置	189
13.5.1 测试数据	189
13.5.2 比较对象	189
13.5.3 评价指标与参数设置	190
13.6 实验结果	190
13.6.1 分类准确率	190
13.6.2 训练开销	193
13.6.3 结果复杂度与特征约简率	195
13.6.4 纠错输出编码性能	197
13.7 本章小结	198
第 14 章 下篇总结与展望	200
14.1 下篇总结	200
14.2 后续工作展望	201
参考文献	203

上篇

演化机器学习

——内嵌特征选择的学习分类器



第1章

上篇引言

1.1 研究背景

伴随着信息技术和移动互联网技术的发展以及社会经济数字化进程的加快,以计算机、智能移动硬件为代表的数字化设备已经渗透到日常生活的方方面面,正在深刻地改变着传统的工作和生活方式。一方面,网络移动互联技术的普及为信息化革命注入强大活力,覆盖全球的 Internet 为人们提供了一个方便快捷又廉价的信息发布平台,带来了数据量的爆炸性增长。另一方面,云存储技术与大规模分布式数据库技术的发展也为收集并存储这些海量数据提供了便利。然而,虽有海量的数据,若缺乏强有力的数据分析工具却只能带来信息的缺位或者数据孤岛,时至今日,面对数据的汪洋大海,人们已经将自己置身于“数据泛滥但信息贫乏”的窘境当中^[1]。20 世纪 80 年代末,一门致力于在实现数据分析自动化的交叉学科应运而生,这就是至今方兴未艾的数据挖掘(Data Mining, DM)。简言之,数据挖掘指的是从大规模数据中提取有用的信息(知识),因此也被称为从数据中发现知识(Knowledge Discovery from Data, KDD)^[1]。作为一门交叉学科,数据挖掘涉及多个领域,既有负责海量数据存储的数据库技术,也有提供强大的运算能力的高性能计算技术,不过其中最为核心的是进行知识抽取的数据分析技术。从早期的统计学^[2]到模式识别(Pattern Recognition, PR)^[3],直至当下更为盛行的机器学习(Machine Learning, ML)^[4],数据挖掘从多个领域汲取养分,增强了本学科海量数据分析的实用性和有效性。