

学校编码: 10384
学号: 23220060153381

分类号__密级__
UDC__

厦 门 大 学

博 士 学 位 论 文

带置信度分类器的研究与应用

Research and Application of Classifier with Confidence

王 华 珍

指导教师姓名: 林 成 德 教 授
专 业 名 称: 控制理论与控制工程
论文提交日期: 2009 年 4 月
论文答辩时间: 2009 年 5 月
学位授予日期: 2009 年 月

答辩委员会主席: _____
评 阅 人: _____

2009 年 5 月

厦门大学学位论文原创性声明

本人呈交的学位论文是本人在导师指导下,独立完成的研究成果。本人在论文写作中参考其他个人或集体已经发表的研究成果,均在文中以适当方式明确标明,并符合法律规范和《厦门大学研究生学术活动规范(试行)》。

另外,该学位论文为()课题(组)的研究成果,获得()课题(组)经费或实验室的资助,在()实验室完成。(请在以上括号内填写课题或课题组负责人或实验室名称,未有此项声明内容的,可以不作特别声明。)

声明人(签名):

年 月 日

厦门大学学位论文著作权使用声明

本人同意厦门大学根据《中华人民共和国学位条例暂行实施办法》等规定保留和使用此学位论文，并向主管部门或其指定机构送交学位论文（包括纸质版和电子版），允许学位论文进入厦门大学图书馆及其数据库被查阅、借阅。本人同意厦门大学将学位论文加入全国博士、硕士学位论文共建单位数据库进行检索，将学位论文的标题和摘要汇编出版，采用影印、缩印或者其它方式合理复制学位论文。

本学位论文属于：

1. 经厦门大学保密委员会审查核定的保密学位论文，于 年 月 日解密，解密后适用上述授权。

2. 不保密，适用上述授权。

（请在以上相应括号内打“√”或填上相应内容。保密学位论文应是已经厦门大学保密委员会审定过的学位论文，未经厦门大学保密委员会审定的学位论文均为公开学位论文。此声明栏不填写的，默认为公开学位论文，均适用上述授权。）

声明人（签名）：

年 月 日

摘要

高风险领域的分类问题对模式分类算法提出以下三个挑战：

- 1) 能否设计一种分类器模型，使得它的输出结果能够附带置信度；
- 2) 预测输出的置信度是有效的，应该使得算法的准确率能够被置信度所控制。
- 3) 算法应能够独立地对每个测试数据提供相应的置信度评估，也就是说，能够根据指定的置信度产生相应的预测结果。

针对上述三个挑战，我们引入了基于转导推理和随机性检验的置信预测方法来解决这些问题。近年新发展起来的一致性预测器（CP）是这种方法的典型代表。但是，CP 在实践中的实用性较差，主要是其固有的运算效率低下、以及对样本奇异函数的设计缺乏指导性准则。我们的工作主要是改进了 CP 的理论模型，提出了混合压缩一致性预测器（HCCP）的算法框架及其实现技术，使其更适用于实际应用。

HCCP在预测性能与计算效率间取得了一个较好的折衷，它在处理大数据集学习问题时，在保持算法的预测效率的同时大大提高了CP的运算效率。HCCP的运作特点是将用于学习的样本序列划分成两个部分，并采用两阶段混合压缩：第一阶段先将前一部分序列样本压缩成一个模型，并以知识的形式保存；第二阶段再将上述知识传递给后续的检验样本序列用于置信预测。在算法实现技术方面，HCCP采用有监督的度量学习方法来实现有效信息在两个子序列（训练样本序列和检验样本序列）中的传递。并分别通过有监督核学习方法（HCCP-KerNN）和随机森林技术（HCCP-RF）实现了度量学习和样本奇异函数的设计。我们还从实验角度展现了HCCP-RF算法在田纳西—伊斯曼化工过程（TEP）这样的流程工业大系统的在线故障置信检测中的适用性和有效性。

针对小数据集的学习问题，我们也提出了一种无划分的 HCCP-RF 算法，它取消了对学习样本集的划分，更加适用于处理小样本数据。通过对慢性胃炎中医诊断数据集的实验，验证了该算法的有效性。

最后，对本文的工作进行了总结，并对今后的研究工作提出了展望。

关键词：分类问题;置信预测;一致性预测器

Abstract

There are three challenges to the researchers on the classification in the high-risk areas:

- 1) Can we develop a classification algorithm that outputs predictions coupled with confidence level?
- 2) Are these confidences for the predictions really valid, i.e., could the accuracy rate be guaranteed by the confidence level?
- 3) Could the algorithm give a prediction with a confidence level tailored for each individual instance, in other words, could it provide a prediction corresponding to the confidence level predefined?

Faced to these challenges, we have introduced a method which uses the transductive inference and the randomness test of i.i.d. sequences to develop our solution. The recently emerged Conformal Predictor (CP) is an alternative solution which can output prediction with valid confidence. However, there are still certain disadvantages in the framework of CP, such as the inherent computational costliness and the lack of guidance for the design of the example nonconformity measure. We have focused on the improvement and the enhancement of CP, and have then proposed a new Hybrid-Compression Conformal Predictor (HCCP) which performs better in practice.

HCCP aims to obtain a good balance between the predictive performance and the computational efficiency. It can maintain a relatively high predictive performance while improving greatly the computational efficiency in dealing with large data sets. HCCP divides the whole training examples into two subsets (called as the training set and the validation set, respectively) and executes the predicting process in two stages. Firstly, it abstracts a compression model M based on the training set; secondly, it designates, for each example in the validation set, the new features which are generated by M and would then be applied by the classical CP algorithm to output the prediction with confidence level. We have proposed a method based on the supervised

metric learning to transfer the useful information from the first stage to the second stage. In detail, we have incorporated the adaptive kernel-based distance metric learning method (as in HCCP-KerNN) and the random forest algorithm (as in HCCP-RF), respectively, to realize the supervised metric learning and the example nonconformity measure. The application is simulated on the standard large data set as Tennessee Eastman Process (TEP). The applicability and effectiveness of the proposed HCCP-RF algorithm are illustrated on this online fault detection of large-scale industrial process.

To deal with the problem of small-sample classification, we have also put forward the non-partition HCCP-RF algorithm, which disclaims the partition of the whole learning set of examples. The application is simulated on the traditional Chinese chronic gastritis data set, which is a typical small-sample problem. The informative as well as effective predictions of the non-partition HCCP-RF algorithm have been shown in the experiment.

Finally, the summary of our work and the future research are presented.

Keywords: classification problem, prediction with confidence, conformal predictor

目录

第 1 章 引言	1
1.1 基于机器学习的分类问题	1
1.1.1 分类问题描述	1
1.1.1 分类问题研究进展	2
1.2 高风险领域模式分类的挑战	4
1.2.1 算法能否输出预测置信度	4
1.2.2 置信度估计是否有效	5
1.2.3 算法能否对单个预测输出特定置信度下的预测结果	6
1.3 论文内容和结构安排	7
第 2 章 机器学习算法的预测置信度分析	9
2.1 置信度	9
2.1.1 社会科学领域的置信度	9
2.1.2 统计学领域的置信度	10
2.1.3 机器学习领域的置信度	10
2.2 贝叶斯方法的置信度分析及其局限性	12
2.2.1 贝叶斯分类器	12
2.2.2 贝叶斯分类器的可校准性	12
2.3 统计学习理论的置信度分析及其局限性	13
2.3.1 统计学习理论的泛化误差	14
2.3.2 PAC 误差界的局限性	14
2.4 基于特定数据集的误差率的可校准性及其局限性	16
2.4.1 泛化误差估计技术	16
2.4.2 基于特定数据集的误差率的可校准性	18
2.5 本章小结	21

第 3 章 一致性预测器	22
3.1 带置信度的统计预测方法	23
3.2 转导推理	24
3.2.1 转导推理的意义.....	24
3.2.2 转导推理分类方法.....	25
3.3 算法随机性理论	27
3.3.1 对象的 Kolmogorov 算法复杂性.....	28
3.3.2 序列的算法随机性 Martin – Lof 检测.....	29
3.3.3 序列的算法随机性 p -值检验.....	31
3.4 一致性预测器算法原理	32
3.4.1 一致性预测器的算法思路.....	32
3.4.2 一致性预测器的序列随机性检验方法.....	34
3.4.3 一致性预测器的算法流程.....	36
3.4.4 一致性预测器的可校准性.....	38
3.5 样本奇异函数的设计	40
3.5.1 样本奇异函数的内含算法.....	40
3.5.2 基于 SVM 的样本奇异函数设计.....	41
3.5.3 基于 KNN 的样本奇异函数设计.....	43
3.6 本章小结	44
第 4 章 混合压缩一致性预测器模型研究	45
4.1 归纳式一致性预测器 (ICP)	46
4.1.1 ICP 算法原理.....	46
4.1.2 ICP 算法流程.....	47
4.1.3 ICP 算法的特点.....	49
4.2 混合压缩一致性预测器 (HCCP) 的算法原理	50
4.2.1 HCCP 算法思想.....	50
4.2.2 HCCP 的算法流程.....	51
4.2.3 HCCP 的可校准性.....	53

4.2.4	HCCP 的计算效率	54
4.3	HCCP 样本奇异函数设计	57
4.3.1	HCCP 的样本奇异函数的特点	57
4.3.2	HCCP 的样本奇异函数的设计思路	57
4.3.3	基于有监督核学习的 HCCP 样本奇异性函数设计	61
4.3.4	基于随机森林的 HCCP 样本奇异函数设计	64
4.4	HCCP 在标准数据集上的可校准性实验	68
4.4.1	数据集创建	69
4.4.2	HCCP 算法的参数设置	70
4.4.3	实验结果与讨论	70
4.5	本章小结	72
第 5 章	混合压缩一致性预测器在大数据集上的应用	74
5.1	大数据集的学习问题	74
5.1.1	大数据集学习问题的特点	74
5.1.2	流程工业大系统的故障检测及其挑战	75
5.1.3	HCCP 在流程工业大系统故障检测中的适用性	77
5.2	HCCP 在 TEP 中的应用	78
5.2.1	田纳西-伊斯曼化工过程 (TEP)	79
5.2.2	TEP 数据集构建	80
5.2.3	HCCP 实验设置	81
5.2.4	HCCP-ProxNN 的框围预测 (Hedged prediction)	82
5.2.5	HCCP-ProxNN 的预测效率	83
5.2.6	HCCP-ProxNN 的计算效率	85
5.3	HCCP-RF 算法的拓展应用	85
5.3.1	RF 相似性度量 <i>Prox</i> 的鲁棒性	86
5.3.2	RF 相似性度量 <i>Prox</i> 与训练样本集规模之间的关系	87
5.3.3	HCCP-RF 的拓展应用	89
5.4	本章小结	89

第 6 章	混合压缩一致性预测器在小数据集上的应用	91
6.1	中医诊断的智能化方法概述	91
6.2	HCCP-RF 算法在小样本集上的修正方案	92
6.2.1	修正思路	92
6.2.2	无划分 HCCP-RF 的算法流程	93
6.3	无划分 HCCP-RF 在慢性胃炎中医诊断中的应用	95
6.3.1	病例收集和数据库构建	95
6.3.2	实验设置	96
6.3.3	实验结果和讨论	97
6.4	本章小结	101
第 7 章	总结与展望	103
	参考文献	106
	致谢	112
	发表论文	113

Contents

Chapter 1 Introduction	1
1.1 Classification Problem in Machine Learning	1
1.1.1 Classification.....	1
1.1.2 Current Research of Classification.....	2
1.2 Challenges of Pattern Recognition with High-Risk	4
1.2.1 Prediction with Confidence.....	4
1.2.2 Confidence being Valid.....	5
1.2.3 Confidence Tailored for Each Individual Instance.....	6
1.3 Main Work and Structure of The Thesis	7
Chapter 2 Confidence Analysis in Machine Learning	9
2.1 Confidence	9
2.1.1 Confidence in Social Sciences.....	9
2.1.2 Confidence in Statistics.....	10
2.1.3 Confidence in Machine Learning.....	10
2.2 Limitations of Bayesian Confidence	12
2.2.1 Bayesian Classifier.....	12
2.2.2 The Calibration of Bayesian Classifier.....	12
2.3 Limitations of Statistical Learning Theory	13
2.3.1 Generalization Error of Statistical Learning Theory.....	14
2.3.2 Limitations of PAC Bound.....	16
2.4 Limitations of Calibration of Error Estimation based on Certain Dataset	
16	
2.4.1 Error Estimation Methods.....	16
2.4.2 The Calibration of Error Estimation based on Certain Dataset....	18
2.5 Summary	21

Chapter 3 Conformal Predictor 22

3.1 Statistical Prediction with Confidence 23

3.2 Transductive Inference 24

 3.2.1 The Significance of Transductive Inference..... 24

 3.2.2 Classification Methods by Transductive Inference..... 25

3.3 Algorithmical Randomness Theory..... 27

 3.3.1 Kolmogorove Complexity of Object 28

 3.3.2 Martin-Lof Randomness Test of Sequence 29

 3.3.3 P-Value Randomness Test of Sequence..... 31

3.4 Conform Predictor Algorithm 32

 3.4.1 The Idea of Conformal Predictor..... 32

 3.4.2 Randomness Test Sequence by Conformal Predictor 34

 3.4.3 Algorithmic Process of Conformal Predictor..... 36

 3.4.4 The Calibration of Conformal Predictor 38

3.5 Nonconformity Measure..... 40

 3.5.1 Underlying Algorithm for Nonconformity Measure..... 40

 3.5.2 Nonconformity Measure based on SVM 41

 3.5.3 Nonconformity Measure based on KNN..... 43

3.5 Summary..... 44

Chapter 4 Model Research of Hybrid-Compression Conformal Predictor 45

4.1 Inductive Conformal Predictor(ICP)..... 46

 4.1.1 ICP Algorithm..... 46

 4.1.2 Algorithmic Process Of ICP..... 47

 4.1.3 The Features of ICP..... 49

4.2 The Algorithm of Hybrid-Compression Conformal Predictor(HCCP) . 50

 4.2.1 The Idea of HCCP 50

4.2.2	Algorithmic Process of HCCP	51
4.2.3	The Calibration of HCCP	53
4.2.4	Computational Complexity of HCCP	54
4.3	HCCP Nonconformity Measure	57
4.3.1	Chacteristics of HCCP Nonconformity Measure	57
4.3.2	Idea of HCCP Nonconformity Measure	57
4.3.3	HCCP Nonconformity Measure based on Supervised Kernel-Metric Learning	61
4.3.4	HCCP Nonconformity Measure based on Random Forest	64
4.4	HCCP Calibration Test on Benchmark Datasets	68
4.4.1	Datasets Creation	69
4.4.2	Experimental Setup	70
4.4.3	Experimental Results and Discussion	70
4.5	Summary	72
 Chapter 5 Application of HCCP on Large-Scale Datasets		74
5.1	Machine Learning on Large-Scale Dataset	74
5.1.1	Problem of Machine Learning on Large-Scale Dataset	74
5.1.2	Challenges in Fault Detection in Large-Scale Industrial Process	75
5.1.3	Applicability of HCCP in Fault Detection in Large-Scale Industrial Process	77
5.2	The Application of HCCP on TEP	78
5.2.1	Tennessee Eastman Process	78
5.2.2	TEP Dataset Creation	80
5.2.3	Experimental Setup	81
5.2.4	Hedged Prediction of HCCP	82
5.2.5	Predictive Efficiency of HCCP	83
5.2.6	Computational Efficiency of HCCP	85
5.3	Expansion of HCCP-RF	85

5.3.1	Robustness of PR Proximity Measure.....	86
5.3.2	Relations between RF Proximity Measure and Size of Training Dataset.....	87
5.3.3	The Expansion of HCCP-RF	89
5.4	Summary	89
Chapter 6 Application of HCCP on Small-Sample Datasets ...		91
6.1	Intelligent Methods in Traditional Chinese Medicine	91
6.2	Modification of HCCP-RF on Small-Sample Dataset.....	92
6.2.1	The Idea of Modification.....	92
6.2.2	Algorithmical Process of Non-Partition HCCP-RF.....	93
6.3	Application of Non-Partition HCCP-RF on Traditional Chinese Chronic Gastritis Dataset	95
6.3.1	Cases Collection and Dataset Creation	95
6.3.2	Experimental Setup	96
6.3.3	Experimental Results and Discussion	97
6.4	Summary	101
Chapter 7 Conclusion and Perspective.....		103
Conferences		106
Acknowledge		112
Publication Paper.....		113

第1章 引言

1.1 基于机器学习的分类问题

1.1.1 分类问题描述

随着信息技术和存储技术的快速发展,政府、商业、企业等各行各业出现了越来越多的复杂非线性高维数据。如何根据用户的特定需求从海量数据中发现有用的知识或者构造从经验中学习的机器,用于对未来数据进行预测成为一个十分迫切的富有挑战性的研究课题。机器学习(Machine Learning, ML)致力于让计算机模仿人类从实例中学习的能力进行数据分析和建模(估计某系统隐藏的、复杂的输入输出关系),使它(这种关系)能够对未知输出做出尽可能准确的预测^[1]。

模式分类是实际应用中普遍存在的问题,也是机器学习领域的基础研究之一。现实生活中存在大量的分类问题,如:机械故障诊断、医学诊断、语音识别、人脸识别、信用评估、文本分类、网络入侵检测、图像识别等。分类的作用和根本目的在于面对某一未知类别的具体事物时,能按照已知的信息将其正确的归于某一类。将某一研究对象正确归入某一类的方法即分类方法。在机器学习领域,模式分类致力于从有限观察发现观测数据中暗含的各种关系,具体说就是从实际问题的一个有限的子集(样本集)出发,探求问题的内在规律(建立模型),从而对未知数据做出正确判断(分类)。这就是机器学习领域的分类问题(classification)。

分类问题用数学语言可以简单描述如下:实际问题的具体对象一般有很多属性,可用高维向量 $x = (x^1, x^2, \dots, x^d)^T \in X$ 表示,其中数据 x 的上标 $1, 2, \dots, d$ 是向量的维数序号,也就是对象的属性序号; $X = \mathbb{R}^d$ 表示 d 维实数空间。假设实际问题对应有限 c 个可选类别,用标识变量 $y \in Y = \{1, 2, \dots, c\}$ 表示对象的类别,其中 Y 叫做类别空间。对象及其类别构成数据对 $z = (x, y) \in Z$, Z 叫做样本空间, $z_i, i = 1, 2, \dots, n, \dots$ 是样本空间中的样本(点),其下标“ i ”表示样本编号。模式分类定义为:

定义 1.1 (分类问题 F) 根据给定的训练样本集 $Z^{(n)} = \{z_1, z_2, \dots, z_n\}$, 其中 $Z^{(n)}$ 表示含有 n 个样本的样本集, $z_i \in Z, i = 1, \dots, n$ 表示样本, 产生一个分类器 $\varphi: \mathbb{R}^d \rightarrow \{1, 2, \dots, c\}$, 使得它对新的待测试数据 $x_{n+1}, x_{n+2}, \dots, x_{n+k}$ 的实际类别值产生相应的预测值 $\hat{y}_{n+1}, \hat{y}_{n+2}, \dots, \hat{y}_{n+k}$ 。我们要求该分类器要能对整个样本空间可能的分布有一个尽可能小的期望判别误差。

从定义 1.1 看, 模式分类技术模仿了人类的逻辑推理过程。逻辑推理遵循的一般性途径大概分为两种方式: 即归纳推理和演绎推理, 它们两者的相互关系见图 1.1

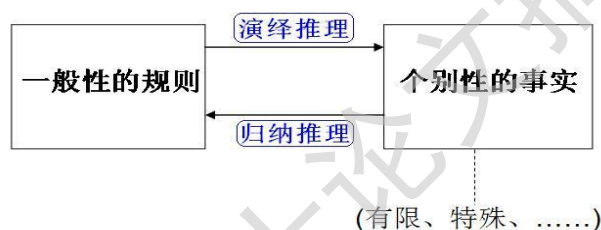


图 1.1 逻辑推理示意图

由上图可知, 模式分类的学习过程是归纳推理, 它从个别性的事实 (一般是有限个的, 具有特殊代表性) 出发归纳出一般性的规则, 而分类器的预测过程是演绎推理, 它从一般性的规则推理出个别性的事实。

1.1.2 分类问题研究进展

随着 1968 年 K.Popper 用不可证伪的概念提出了关于归纳问题的理论^[2], 统计学被看作是归纳推理的一个数学模型, 因此统计学也成为模式分类的重要理论基础之一。在此基础上已经出现了多种以经典统计理论为工具刻画的模式分类方法, 如: 贝叶斯决策、K-近邻、线性判别分析、决策树等等^[3-5]。这些经典的统计机器学习方法在机器学习问题中起着基础性的作用。然而统计分类器偏离了统计推理的本质模式 (它应该是根据观测数据寻求感兴趣的分类器 (函数)), 而变成具有 R.Fisher 参数估计特点的‘模型辨识’方法。Fisher 统计是基于参数估计的统计推理, 他把从给定数据估计函数这个问题 (分类问题、回归问题和密度估计问题) 表达为特定 (参数化) 模型的参数估计问题, 并提出了估计所有模型未知

Degree papers are in the "[Xiamen University Electronic Theses and Dissertations Database](#)". Full texts are available in the following ways:

1. If your library is a CALIS member libraries, please log on <http://etd.calis.edu.cn/> and submit requests online, or consult the interlibrary loan department in your library.
2. For users of non-CALIS member libraries, please mail to etd@xmu.edu.cn for delivery details.

厦门大学博硕士学位论文摘要库