蓝超 等

半监督多视图学习在大数据分析中的应用探讨

DOI: 10.3969/j.issn.1009-6868.2015.05.009 网络出版地址; http://www.cnki.net/kcms/detail/34.1228.tn.20150923.0923.002.html

半监督多视图学习在大数据分析中 的应用探讨

Semi-Supervised Multi-View Learning in Big Data

中图分类号: TN929.1 文献标志码: A 文章编号: 1009-6868 (2015) 05-0032-003

摘要: 半监督多视图学习是机器学习领域一种极具潜力的大数据处理和分析方 法,该方法能有效处理异构和半监督数据,并能方便地在线化和并行化,适合处理 海量数据。该方法在大数据时代的应用前景值得研究人员和业界关注。指出未来 需要通过引入其他领域新的研究技术和成果,不断丰富和完善半监督多视图学习的 理论体系和算法设计,并在实验和实践中不断检验和探索。

关键词: 半监督;多视图;大数据;并行化

Abstract: This paper introduces a promising machine-learning paradigm called semi-supervised multi-view learning. With this paradigm, information is extracted from heterogeneous and semi-supervised data sets. Lately, multi-view learning has been scaled up online and through parallelization to deal with emerging big data challenges. Due to its successful application in many research domains and the fact that it has been explored and used by leading companies, multi-view learning may have a future in the big-data era as a major data analytic technique. New research techniques should be introduced into this area to improve the theoretical system and algorithm design of semi-supervised multi-view learning.

Keywords: semi-supervised; multi-view; big data; parallelization

蓝超/LAN Chao1 饶泓/RAO Hong² 浣军/HUAN Jun¹

(1. 堪萨斯大学工程学院, 堪萨斯 66045; 南昌大学信息工程学院,江西南昌 330031)

(1.College of Engineering, University of Kansas, Kansas 66045, USA;

2. College of Information Engineering, Nanchang University, Nanchang 330031,

- 如何有效利用无监督数据,是半 监督学习研究的核心问题
- 提高无监督数据的样本复杂度的 效率是值得关注的问题
- 多视图学习的样本复杂度分析主 要集中于主动学习

器学习已成为产业界大数据分 机析的主流工具。在2015年北京 全球软件开发大会上,基于大数据的 机器学习和数据挖掘专题讨论得到 了业界广泛参与,如百度、搜狗、阿里 巴巴、京东、美团、猿题库等著名企业 参与了该专题的讨论。当前,机器学 习在产业界的推广正处于白热化阶 段。与此同时,如何结合大数据的特 点,选择恰当的学习方法仍是值得探 讨的问题。

从数据角度而言,大数据之"大" 并不仅限于容量,也指数据类型多和 收发速度快,三者统称为大数据的

收稿日期:2015-08-27 网络出版时间:2015-09-23 "3V"特征[[]。因此,若想迅速有效地 处理和分析大数据,不仅需要从数据 库和计算机体系结构等入手(如IBM 的 Hadoop, 微软的大数据生态系统), 更需选择合适的数据分析方法,才能 充分挖掘海量数据中潜在的信息和

除了3V,大数据普遍存在却鲜有 讨论的另一特征是半监督性。在传 统机器学习中,用于训练分类系统的 数据大多是有监督的,即数据的实际 分类已知。然而,许多应用问题如自 然语言处理,大部分数据都不知道实 际分类,这类数据称为无监督数据。 如何有效利用无监督数据,使其帮助 有监督数据一起训练分类系统,是半 监督学习四研究的核心问题。随着大 数据时代的到来,大量无监督数据将 被采集,导致越来越多应用问题演变 成大规模的半监督学习问题。

大规模半监督学习在学术界的 探讨可追述到 2005 年 Delalleau 提出 的高效无参归纳法国。在该论文中, 作者提出将用于构造核矩阵的(无监 督)数据集进行下采样,减小核矩阵 大小,从而提高算法的存储和计算效 率。在此思想基础上衍生出来的工 作有很多,其中以大规模化基于核或 基于图的半监督学习方法为主。除 了数据下采样,也有研究员通过构造 混合分布模型以减小参数估计的计 算复杂度。此外,半监督学习的大规

中兴通讯技术 **32** 2015年10月 第21卷第5期 Oct. 2015 Vol.21 No.5

模化也可通过提高优化算法效率或 并行程度来实现。

尽管有着近十年的研究历史,目 前大部分的大规模半监督学习方法 仍具有局限性。一方面,它们大都基 于传统的单视图机器学习方法,忽略 了大数据的异构性,容易导致信息丢 失;另一方面,规模化算法的策略较 为传统,如下采样或加速优化,无法 有效满足大数据带来的新的挑战,如 数据的在线化和分布化等。

文章结合大数据的特点介绍机 器学习领域的一种新兴的半监督学 习方法——多视图学习。它不仅在 许多研究领域已得到成功应用,更因 具备有效利用无监督异构数据的能 力和迅速发展的大规模化算法,有望 成为大数据时代最具潜力的数据分 析工具之一。

1多视图学习优点多

1.1 多视图学习及其优点

多视图学习℡是指专门针对多视 图数据而进行建模和学习。其中,多 视图数据是指由多组(往往具有不同 意义的)特征进行描述的数据,而每 一组特征称为一个视图。多视图学 习的主要思想是基于无监督数据的 视图一致性,即分类器在同一无监督 数据不同视图下的分类结果应基本 一致。将此约束加入学习法则,多视 图学习便能巧妙利用无监督数据帮 助分类器的训练。

多视图学习的一个显著优点是 缓解过学习问题,即由于模型过于复 杂而将数据噪声也学入分类法则的 现象。传统的机器学习方法大多将 数据的所有视图堆砌成一个高维的 单视图数据,进行建模和学习。此 时,如果视图间存在冗余信息,为高 维数据所建的模型将比实际需求更 为复杂,容易导致过学习。而多视图 学习则为每个视图分别建模,有效降 低了模型复杂度。

多视图学习的另一优点是提升

数据的总体分类能力。当数据特征 所蕴含的分类信息总体较弱时,可将 特征集拆分成多个视图进行多视图 学习的。通过各视图的弱分类器协同 训练,达到各视图"单独学习弱,集成 学习强"的目的。

另外,多视图学习还能有效处理 异构数据6。大数据时代,数据的异 构性越来越强。比如,客服中心为了 提高工作效率,需要根据来电客户的 个人信息和语音信息对其来电目的 进行快速预测。这里,用户的个人信 息是静态文本数据,而语音信息则是 动态时序数据,两者不仅数据类型不 同,也往往服从不同的后验(预测)概 率分布。强行统一或堆砌两类特征 不仅给建模带来困难,也容易导致信 息丢失,失去大数据分析的优势。多 视图学习则允许各视图分别选择合 适的分类器,再通过协同训练提高视 图整体的分类能力。

1.2 多视图学习的理论研究

多视图学习的理论分析主要基 于其分类器的泛化误差。一个分类 器的泛化误差是其在指定数据分布 下,错分数据的概率。许多理论研究 证明:多视图中的协同学习方法在满 足条件时,分类器的泛化误差将以极 大的概率收敛到极小的范围之内。

多视图学习的早期理论分析基 于两个假设四:视图充分性,即每个视 图能分别提供充分的(但不必很强) 分类信息。该假设在大数据中较容 易被满足;各视图间条件独立,即给 定数据分类,其各视图间统计独立。

多视图学习另一个理论研究是 其样本复杂度。样本复杂度指通过 多少数据的训练就能保证分类器收 敛到预期的泛化误差。目前,多视图 学习的样本复杂度分析主要集中于 主动学习,即选取哪些无监督数据进 行人工分类,才能使分类器的训练最 快收敛。在此问题中,有监督数据的 样本复杂度被证明与泛化误差的倒 数成 log 比^[8], 而无监督数据样本复杂 度则与泛化误差的倒数成正比門。在 大数据时代,无监督数据的容量常常 远大于有监督数据,是计算机的存储 和计算的主要负担。因此,提高无监 督数据的样本复杂度的效率是值得 关注的问题。

1.3 多视图学习的常见算法及 大规模化算法

大部分多视图算法为每个视图 分别建立分类器,基于协同训练方式 的不同,有两种主要的多视图算法: 一是基于迭代的协同训练算法四,另 一个是基于协同正则化的算法[10]。 基于迭代的协同训练算法的每个回 合由一个视图的分类器标注一部分 无监督数据的类别,并将它们加入有 监督数据集,一起重新训练其他视图 的分类器,以达到视图一致,提高视 图总体的分类能力。这类算法直观 有效,也较容易处理异构数据,但因 其迭代的训练方式,学习效率相对较 低;基于协同正则化的算法为所有视 图的分类器统一建立一个目标函数, 通过特定约束条件达到视图一致,并 一次性地训练完所有分类器。该类 方法避免了迭代训练,计算效率往往 更高。

为各视图分别建模能最大程度 保留视图的个性特征,但对视图一致 性的要求也更高。当一致性不能被 较好满足时,该类算法的效果会有所 下降門。另一类多视图学习算法则 侧重于将多个视图进行有机融合成 单视图,再进行单视图学习[12]。这类 算法有利于提高特征的鲁棒性和泛 化能力,不同于简单的堆砌视图,但 也具有丢失视图个性特征的风险。 为此,有学者提出同时学习视图的共 性特征和个性特征[13]。

除了改良传统算法,多视图学习 也正沿着在线化和并行化的方向发 展,以迎接大数据的新挑战。在线多 视图学习假设训练数据分批,依次呈 现给分类器进行学习。此时,如何为 分类器设计高效的更新算法,是在线

多视图学习的研究问题。目前研究 的主流算法是基于协同正则化的算 法,而其在线化的主要策略是通过引 入随机梯度下降法,实现目标函数的 在线优化四。该策略的优点是计算 速度快,存储空间小,但对学习参数 的敏感性较高。

目前的并行多视图学习研究主 要基于 MapReduce 的框架。如爱立信 研究院提出四将不同数据的不同视 图被分配到不同计算单元进行计算, 以提高计算效率。但由于视图一致 性的要求,不同计算单元间往往需要 频繁通信,这成为系统效率的决定性 瓶颈,也是当前分布式机器学习的一 个研究热点。并行或分布式多视图 学习中通信量如何增长? 有哪些影 响因素?如何设计低通信量的高效 算法?这些都是大规模多视图学习 需要探索的问题。

2 多视图学习应用广

多视图学习在许多研究领域都 已获得成功应用,包括网页分类、自 然语言处理,计算机视觉、医疗诊断、 药物分析和化学分析。在网页分类 中,每个网页的内容和链表可表示为 两个视图;在文本分类中,同一文本 的不同语言版本可表示为不同视图; 在用户决策识别系统中,用户的语音 和姿势可表示成两个视图;在图像检 索和标识中,图像本身和周边的文本 信息可表示成两种视图。即使在一 些数据并不自然展示出多个视图的 问题中,也可通过从单视图中提取多 个视图进行多视图学习,以充分发挥 多视图学习的优势。

多视图学习在工业界也被积极 探索和推广,涉及领域包括机器翻 译、情感预测、图像检索和推荐系统 等。在微软研究的机器翻译中[16],测 试文章在不同解码器下的译文被视 作不同视图下的准参考译文加入训 练集参与翻译器的特征权重的训练, 从而使翻译器更全面的适应测试集 领域的文章的特点,达到领域自适应 的目的。

3 结束语

半监督多视图学习是大数据时 代极具潜力的分析工具。其在许多 研究领域已获得成功应用,并在工业 界被积极推广。但同时也需指出,大 数据的大容量、分布式和在线化等特 征为半监督多视图学习带来了新的 挑战。比如,如何有效降低无监督数 据的样本复杂度,如何降低分布式视 图间的通信量等。要解决这些问题, 需要通过引入其它领域新的研究技 术和成果,不断丰富和完善半监督多 视图学习的理论体系和算法设计,并 在实验和实践中不断检验和探索。

参考文献

- [1] LABRINIDIS A, and JAGADISH H V. Challenges and Opportunities with Big Data: A community white paper developed by leading researchers across the United States [R]. 2012
- [2] ZHU X J. Semi-supervised learning literature survey [EB/OL]. http://pages.cs.wisc.edu/ ~jerryzhu/research/ssl/semireview.html
- [3] DELALLEAU O, YOSHUA B and ROUX N. Efficient non-parametric function induction in semi- supervised learning [C]// Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics, Bridgetown, Barbados, 2005: 96-103
- [4] XU C, TAO D C, and CHAO C. A Survey on Multi-view Learning [Z]. arXiv preprint: 1304.5634, 2013
- [5] CHEN M M, CHEN Y X and WEINBERGER K Q. Automatic feature decomposition for single view co-training [C]//Proceedings of the 28th International Conference on Machine Learning, Washington, USA, 2011: 953-960
- [6] LIAN W Z, RAI P, SALAZAR E and CARIN L. Integrating Features and Similarities: Flexible Models for Heterogeneous Multiview Data [C]//Twenty-Ninth AAAI Conference on Artificial Intelligence, Austin Texas, USA,
- [7] BLUN A and MITCHELL T. Combining labeled and unlabeled data with co-training [C]//Proceedings of the eleventh annual conference on Computational learning theory, Madison, USA, 1998: 92-100
- [8] WANG W and ZHOU Z H. On multi-view active learning and the combination with semi-supervised learning[C]//Proceedings of the 25th international conference on Machine learning, Helsinki, Finland, 2008: 1152-1159
- [9] BALCAN M F and BLUM A. A Discriminative Model for Semi-Supervised Learning [J]. Journal of the ACM (JACM), 2010,57(3):19. doi: 10.1145/1706591.1706599

- [10] SINDHWANI V, NIYOGI P and BELKIN M. A co-regularization approach to semisupervised learning with multiple views[C]// Proceedings of ICML workshop on learning with multiple views, Bonn, Germany, 2005: 74-79
- [11] CHRISTOUDIAS C. URTASUN R and DARRELL T, Multiview leanringin the presence of view disagreement [C]// Conference on Uncertainty in Artificial Intelligence, California, USA, 2012
- [12] CHEN N, ZHU J and XING E P, Predictive subspace learning for multi-view data: a large margin approach[C]//Advances in neural information processing systems, Hyatt Regency, Vancouver Canada, 2010
- [13] JING X Y, HU R M, ZHU Y P, Wu S S. LIANG C. and YANG J Y. Intra-View and Inter-View Supervised Correlation Analysis for Multi-View Feature Learning[C]// Twenty-Eighth AAAI Conference on Artificial Intelligence, Québec City, Qu é bec, Canada, 2014
- [14] Ruijter T D, Tsivtsivadze E and Heskes T. Online co-regularized algorithms [M]. Germany: Springer Berlin Heidelberg, 2012
- [15] HARIHARAN C. and SUBRAMANIAN S. Large scale multiview learning on mapreduce[C]// 19th International Conference on Advanced Computing and Communications, Chennai, India, 2013
- [16] MAEIREZO B, LITMAN D and HWA R. Cotraining for predicting emotions with spoken dialogue data[C]// Proceedings of the ACL 2004 on Interactive poster and demonstration sessions, Barcelona, Spain,

作者简介



蓝超,美国堪萨斯大学计算 机系在读博士;研究方向为 机器学习和模式识别;在国 际会议及期刊上发表论文 10余篇。



饶泓,南昌大学信息工程学 院教授;主要研究方向为机 器学习、数据挖掘、智能信 息处理;主持多项国家及省 市科技项目;在国内外学术 刊物及国际会议发表论文 30余篇。



浣军,美国堪萨斯大学计算 机系教授,美国国家卫生 局、航空航天局以及自然科 学基金评审委员会、委员 Elisver 及 Springer 大数据 期刊编委等;研究方向为机 器学习和大数据挖掘,以及 在生物信息学中的应用;于 2009年获美国国家科学基 金会教师早期职业发展奖;

在《自然-生物技术》等国际学术期刊及会议上 发表论文100余篇。

中兴通讯技术 34 2015年10月 第21卷第5期 Oct. 2015 Vol.21 No.5