

图像分类中的深度主动学习研究综述

江显森

(贵州师范大学 大数据与计算机科学学院, 贵州 贵阳 550025)

摘要:深度学习已在众多领域如图像分类中取得突破性发展,其成功依赖大量标注数据。然而很多领域中数据标注代价昂贵。主动学习主要是通过合适的查询策略选择信息量大的未标注数据交由专家或者工作人员进行标记,试图以尽可能少的高质量标注数据训练高性能的模型。从不同角度详细地对基于预设计策略和基于学习损失的主动学习方法的研究现状进行了分析和比较,最后对现有的主动学习进行了总结和进一步指出了一些值得研究的方向。

关键词:深度学习;主动学习;深度主动学习;预设计策略

中图分类号:TP391

文献标识码:A

文章编号:2096-9759(2021)03-0161-03

1 研究背景及意义

随着 21 世纪的到来,人工智能,尤其是深度学习技术得到了高速的发展。基于深度学习的计算机视觉应用不断涌现。计算机视觉最重要的研究内容之一就是图像分类,其是目标检测、图像分割、图像搜索等众多应用的基础。具体来说,图像分类就是把分类标签集合中的一个标签分配给输入图像。目前,图像分类被广泛应用于医疗图像处理、人脸识别、智能视频分析、交通场景识别等场景。

用于图像分类的传统机器学习方法有 SVM^[1]、KNN^[2] 等等,这些传统机器学习算法比较适合数据量比较少情况,一旦遇上海量的数据,它们往往会显得力不从心。伴随着计算机技术和传感器技术的发展,获取海量的数据已然不是一件难事,由此,深度学习引起了很多研究人员的关注,并且得到了迅速发展。Alex Krizhevsky 在 2012 年 ILSVRC 比赛中提出的 CNN 模型引起了很多人的兴趣,其效果比大多数传统方法都要好,并且获得了 ILSVRC2012 冠军,该经典模型被称作 AlexNet^[3],这也是首次在大规模图像分类中引入深度学习算法。从 AlexNet 之后,陆续出现了许多优秀的 CNN 模型,例如 VggNet^[4]、GoogleLeNet^[5]、ResNet^[6]、DenseNet^[7] 等等,在许多的挑战任务上取得了优异的表现。尽管深度学习在许多的领域上取得了突破性的发展,但是这些都归功于大量标注数据集的公开。对于大多数 CNN 模型,大量标注的数据在训练中是非常有必要的,因此,获取大量标注的数据成为了一项很重要的任务,然而,获取大量的标注数据在现实生活的很多场景中并没有那么容易。实际任务中,有监督学习往往拥有很多未标注数据,而缺少有标注的数据,但注释过程却可能非常昂贵且受限制。另外,大量高质量的未标注数据在一些需要很高专业知识的领域很难找到有经验的知识工作者去标注,尤其是在医疗图像处理、金融风险预测、工业缺陷检测等应用场所。由此,如何使用尽可能少的高质量标注数据去训练高性能的 CNN 模型引起了广大研究人员的关注^[8-12]。因此,主动学习^[13]AL 逐渐受到了应有的重视。

在保留深度学习强大的学习能力之外,主动学习能够降低数据标注的成本。与深度学习不同,主动学习主要是从数据集入手,设计优异的查询规则从大量未标记的数据集选择最具有价值的样本,并且交由专家或者工作人员标记样本,试

图尽可能的降低标注成本。在主动学习系统中:首先使用少量标注好的训练样本来训练初始模型,以达到一定的模型初始分类能力,然后在接下来的每一轮迭代中,主动学习根据某种合适的查询策略从海量的未标注数据集中寻找一批最具有价值的样本以提供给专家或者工作人员,待专家或者工作人员人工标记后,将这些被标注过的数据加入到已标注的数据集中,最后在更新后的已标注数据集中训练模型,不断提高模型的性能,重复上面的过程,直到标注的预算被耗尽或者达到了预先设置好的条件为止。主动学习系统中最重要的内容就是样本选择策略的设计,相关的研究也相当丰富。目前,样本选择策略有单独的基于不确定性的,单独基于多样性的,也有结合基于不确定性和多样性的组合策略。尽管上面的样本选择策略在一定程度上被证实对于模型的性能提升有效,且能显著降低标注的成本。对于现有基于深度学习的 CNN 模型,单一的查询策略对性能的改进非常有限,单一的基于不确定性的查询策略在挑选样本时可能会造成采样偏差,并且易于查询异常值(离群点),即被挑选到的样本集可能代表不了样本的多样性;仅使用单一的基于多样性的采样方法可能会导致标注成本的增加,即可能有一部分信息量比较低的样本会被挑选到。除此之外,还有一些工作^[14-18]研究了结合样本的不确定性和多样性的组合策略,并且试图在其中找到平衡点,这种结合样本的不确定性和多样性的方法克服了单一策略的不足,在一定程度上即提高了模型的性能,又减轻了标注数据的负担,但是这也造成了挑选样本时指标计算复杂,且平衡点难以寻找的问题。近几年来,基于学习损失的主动学习方法也开始涌现^[8-10],这种新的主动学习方法与任务无关,可以很好的推广到其它的领域。综上所述,研究一种更加合理的样本主动学习选择策略具有重要的意义。

2 研究现状及存在的问题

本文根据以下要点对主动学习算法的最新研究进展进行了综合分析:①基于预设计策略的主动学习方法;②基于学习损失的主动学习方法。

2.1 基于预设计策略的主动学习方法

样本选择策略的质量会直接影响到整个主动学习方法在图像分类中的效果。目前已经有非常多的预设计策略不断被提出并且被应用到主动学习方法中。在实际应用中,对大量的

收稿日期 2021-01-11

作者简介 江显森(1994-)男,广东茂名人在读研究生,主要研究方向:计算机视觉与图像处理,深度学习,主动学习。

数据进行标注非常枯燥、耗时,对于一些特殊图像,如合成孔径雷达(SAR)图像,对其内容判读非常困难,故能够获得的数据非常有限,针对以上问题,文献[19]提出了一种基于BvSB(最优标号和次优标号)+CST(带约束的自学习)的图像分类算法,通过BvSB主动学习将对当前分类器最具信息量、最有价值的样本挑选出来,同时利用CST半监督学习从剩余的大量未标注数据中挑选出一部分兼具信息性和分类确信度高的样本进行补充,实验结果表明,新算法能够有效地减轻训练过程中的人工标注负担,并获得较好的分类性能。尽管大多数不确定查询选择策略在许多的情况下是非常有效的,但是它们无法考虑大量未标记实例中的信息,并且会查询到异常值(离群值),可能会降低模型的性能,考虑到以上问题,Li等^[20]提出了一种新型的自适应主动学习方法,引入信息密度度量(information density),能够衡量候选样本和未标注样本之间的相互信息,该方法将信息密度度量 and 不确定性度量结合在一起,以选择关键实例来标记,在针对图像分类问题的实验中,与许多的现有主动学习方法相比,该方法能够明显减少学习良好分类器所需的训练数据。文献[18]提出了一种用于视觉概念识别的半监督批处理模式多类主动学习算法 USDM,该算法利用整个活动池来评估数据的不确定性,考虑到不确定数据总是彼此相似,建议所选的数据尽可能多样化,故对目标函数明确施加了多样性约束,实验结果表明其在动作识别、对象分类、场景识别等领域具有优势。目前的AL给出的带标签训练样本对于CNN来说是不够的,因为AL通常会忽视大多数的未标记样本,另外,AL和CNN的处理管道彼此不一致。针对以上两个问题,Wang等^[12]为深度图像分类任务提出了一个经济有效的主动学习框架CEAL,该框架采用互补的样本选择策略,在样本挑选阶段,选择少数信息量最大的样本的同时,选择多数高置信度的样本来自动标注,然后进行模型更新。从整体上看,少数标记样本有利于分类器的决策边界,大多数伪标记样本为鲁棒特征学习提供了足够的训练数据。在两个公开的具有挑战性的基准上的广泛实验结果证明,CEAL具有有效性。但是,CEAL使得模型的训练数据骤然增加,违背了主动学习的初衷,最少样本获取最大性能,并且样本过多会导致迭代的时间加长。生物医学成像中缺少大型已经标记好的医疗数据,需要具有相关专业知识的医生,成本高,周期长,AIFT网络^[17]被提出。AIFT把主动学习和迁移学习集成到一个框架,直接使用一个预训练的AlexNet模型,根据熵和多样性结合的选择策略从未标注数据里找一些比较值得标注的样本,然后模型持续地加入新标注的数据,一直做增量微调。在三种不同的生物医疗数据上测试,证明该方法至少能够减少一半的标注代价。文献[21]提出一种主动模型自适应方法,用于有效地训练深度CNN。其提出了一种动态平衡显著性和不确定性的主动选择准则,以选择最适合标签查询的实例,其中显著性度量了实例对改进网络的贡献。结果表明,该方法能够以较低的训练成本实现有效的深度网络训练。Pan等^[16]提出一种新颖的主动学习框架ALBS,ALBS使用融合样本选择策略EDPC。该策略结合了“最具判别性”和“最具代表性”,可以从未标记的数据集中寻找有价值的样本,并逐步更新模型以不断提高模型性能,解决了模型训练带来的灾难性遗忘问题。实验结果表明了ALBS的有效性和稳定性。为了解决使用单一采样策略的不足,文献[15]设计了一种新的解

决方案“多准则主动学习”MCADL。该样本选择策略同时考虑多个标准(即密度、相似性、不确定性和基于标签的度量)来选择信息丰富的样本。实验结果表明MCADL是合理且有效的。然而,MCADL选择信息样本的速度比基于熵的策略、最小置信度方法和边缘抽样方法要慢,需要从多个标准计算信息值,增加了计算时间;另外,需要手动设置一些超参数以获得足够的性能。在基于池的主动学习中,可利用的未标记数据通常不用于模型训练,Gao等^[14]建议使用半监督的AL框架统一模型训练和样本选择,提出了一个基于样本一致性的简单而有效的选择指标,在挑选时,该指标隐性地平衡了样本不确定性和多样性。针对基于不确定性的采样方法的信息冗余问题,有研究人员建议通过考虑代表性来缓解这个问题,文献[22]提出新的基于池的主动学习方法DAAL,该方法能够选择同时对分类器具有较高不确定性和对未标记样本具有良好代表性的样本。

2.2 基于学习损失的主动学习方法

基于预设计策略的主动学习方法难以适应不同情况,通用性较差。近年来,有不少研究人员开始关注如何创造一种通用的主动学习方法,该方法与任务无关,可以很好地推广到其它领域。基于学习损失的主动学习方法^[8-10]的提出对主动学习的发展具有重要意义。文献[8]提出一种新颖的主动学习框架,通过引入针对主动学习任务的损失函数,以选择信息量最大的未标记样本来训练深度信念网络模型。迄今为止,这是将主动学习的准则集成到用于训练深度信念网络的损失函数中的第一项研究工作。对各种单模态和多模态视觉数据集的广泛实证研究证实了该方法在现实世界中图像识别应用中的潜力。目前,大多数的主动学习方法存在一个共同的问题:针对其目标任务而人工设计样本选择策略,这种方式非常低效且无法满足大多数实际需求。针对以上问题,Yoo等^[9]提出一种新的主动学习方法,该方法简单,但与任务无关,并能有效地与深层网络协同工作。其在目标网络上添加一个名为“损失预测模块”的小参数模块,并学习它来预测该未标记数据可能造成的目标损失,然后,该模块建议为目标模型标记哪些可能产生较大目标损失的数据。在图像分类、目标检测和人体姿态估计任务中的实验结果表明该方法性能始终优于以前的方法。虽然该方法可以推广到多种任务,但是没有考虑数据的多样性或密度等其他有用信息。Li等^[10]提出了一种基于学习损失的主动学习方法,该方法同样是任务无关的。其和文献^[9]一样附加一个学习模块来预测未标记数据的目标损失,并选择损失最大的数据进行标记。其证明了基于学习损失的主动学习算法实际上是一个学习排序问题,因此使用一个简单而有效的列表方法,通过优化Spearman的秩相关度量来训练损失预测模块。

3 总结与展望

最近的大多数工作显示了主动学习在深度图像分类任务上获得了成功,主动学习凭借其减轻数据标注负担的能力引起广大研究人员的兴趣。主动学习方法主要可以分为基于预设计策略的和基于学习损失的主动学习方法:基于预设计策略的AL方法由单一的采样方法逐步发展到多标准结合的采样方法,期间涌现出很多优秀的经典方法。但是单一的采样方法有可能会给训练集样例带来冗余问题,多标准结合的采

样方法也有可能造成指标计算复杂,平衡点难以寻找的问题,并且这些以往的AL方法主要是针对具体目标任务而设计,难以作为一种通用框架推广到其它领域。基于学习损失的AL方法简单且与任务无关,能很好地和深层网络协同工作,这使得这种方法可以应用于任何使用深度网络的机器学习任务中。但是这种方法并没有考虑数据的多样性或密度等其他有用信息,有一定的局限。主动学习的理论研究已经日益成熟,但是依然还是有一些值得研究的方向:①增量训练是一个比较值得关注的方向。在每次AL迭代中重新利用已标注样本集训练深度模型是难以接受的,因为这样会浪费宝贵的时间。增量训练会带来极大的计算资源和时间上的节省,但是简单的增量训练会产生遗忘灾难问题,且会带来采样偏差;②目前的大多数AL方法主要是针对特定任务而预设计采样策略,并不适合推广到其它的场景,基于学习损失的AL方法与任务无关,更有利于被广泛应用到各种任务中,是未来的一个重要发展方向。另外,显式或隐式地将数据的多样性或密度引入基于学习损失的AL方法中,也是一个比较值得关注的问题,当然也会面临新的挑战。

参考文献:

- [1] Cortes Corinna, Vapnik Vladimir. Support-vector networks [J]. Machine Learning, 1995, 20(3): 273-297.
- [2] Altman Naomi S. An introduction to kernel and nearest-neighbor nonparametric regression [J]. The American Statistician, 1992, 46(3): 175-185.
- [3] Krizhevsky Alex, Sutskever Ilya, Hinton Geoffrey E. ImageNet classification with deep convolutional neural networks [M]. Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1. Lake Tahoe, Nevada; Curran Associates Inc. 2012: 1097 – 1105.
- [4] Simonyan Karen, Zisserman Andrew. Very deep convolutional networks for large-scale image recognition [J]. arXiv preprint arXiv:14091556, 2014.
- [5] Szegedy C., Wei Liu, Yangqing Jia, et al. Going deeper with convolutions; proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), F 7-12 June 2015, 2015 [C].
- [6] He K., Zhang X., Ren S., et al. Deep Residual Learning for Image Recognition; proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), F 27-30 June 2016, 2016 [C].
- [7] Huang Gao, Liu Zhuang, Van Der Maaten Laurens, et al. Densely connected convolutional networks; proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, F, 2017 [C].
- [8] Ranganathan H., Venkateswara H., Chakraborty S., et al. Deep active learning for image classification; proceedings of the 2017 IEEE International Conference on Image Processing (ICIP), F 17-20 Sept. 2017, 2017 [C].
- [9] Yoo D., Kweon I. S. Learning Loss for Active Learning; proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), F 15-20 June 2019, 2019 [C].
- [10] Li Minghan, Liu Xialei, van de Weijer Joost, et al. Learning to Rank for Active Learning: A Listwise Approach [J]. arXiv preprint arXiv:200800078, 2020.
- [11] Sener Ozan, Savarese Silvio. Active Learning for Convolutional Neural Networks: A Core-Set Approach; proceedings of the International Conference on Learning Representations, F, 2018 [C].
- [12] Wang K., Zhang D., Li Y., et al. Cost-Effective Active Learning for Deep Image Classification [J]. IEEE Transactions on Circuits and Systems for Video Technology, 2017, 27(12): 2591-2600.
- [13] Settles Burr. Active Learning Literature Survey [J]. University of Wisconsin, Madison, 2010, 52(55-66): 11.
- [14] Gao Mingfei, Zhang Zizhao, Yu Guo, et al. Consistency-based semi-supervised active learning: Towards minimizing labeling cost; proceedings of the European Conference on Computer Vision, F, 2020 [C]. Springer.
- [15] Yuan Jin, Hou Xingxing, Xiao Yaoqiang, et al. Multi-criteria active deep learning for image classification [J]. Knowledge-Based Systems, 2019, 172: 86-94.
- [16] Pan Longfei, Wang Xiaojun. ALBS: An Active Learning Framework Based on Syncretic Sample Selection Strategy [M]. Proceedings of the 2019 11th International Conference on Machine Learning and Computing. Zhuhai, China; Association for Computing Machinery. 2019: 533 – 538.
- [17] Zhou Zongwei, Shin Jae, Zhang Lei, et al. Fine-tuning convolutional neural networks for biomedical image analysis: actively and incrementally; proceedings of the Proceedings of the IEEE conference on computer vision and pattern recognition, F, 2017 [C].
- [18] Yang Yi, Ma Zhigang, Nie Feiping, et al. Multi-Class Active Learning by Uncertainty Sampling with Diversity Maximization [J]. International Journal of Computer Vision, 2015, 113(2): 113-127.
- [19] 陈荣, 曹永锋, 孙洪. 基于主动学习和半监督学习的多类图像分类 [J]. 自动化学报, 2011, 37(08): 954-962.
- [20] Li X., Guo Y. Adaptive Active Learning for Image Classification; proceedings of the 2013 IEEE Conference on Computer Vision and Pattern Recognition, F 23-28 June 2013, 2013 [C].
- [21] Huang Sheng-Jun, Zhao Jia-Wei, Liu Zhao-Yang. Cost-Effective Training of Deep CNNs with Active Model Adaptation [M]. Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. London, United Kingdom; Association for Computing Machinery. 2018: 1580 – 1588.
- [22] Wang Shuo, Li Yuexiang, Ma Kai, et al. Dual Adversarial Network for Deep Active Learning; proceedings of the Computer Vision-ECCV 2020, Cham, F 2020//, 2020 [C]. Springer International Publishing.