

中山大学本科生
优秀毕业论文

(中文题目) 面向视觉识别的自主学习框架

(英文题目) Self-Learning Framework for
Visual Recognition

学位申请人: 严肖朋

导 师: 林惊

专 业: 自动化

学 院: 电子与信息工程学院

2017年6月

编号: 2017230

论文题目：面向视觉识别的自主学习框架

专 业：自动化

学生姓名：严肖朋

学 号：13351054

指导教师：林惊 教授

摘 要

近年来，随着深度卷积神经网络（DCNN）的发展，物体检测这个传统的视觉识别问题在大数据集上取得了很大的进步，应用也越来越广泛。物体检测现行通用做法是通过 CNN 提取图片的感兴趣区域（Region of interest, ROI），再将其转化为分类问题，可以取得很高的识别精度，例如 RCNN, Faster RCNN 等流行网络。然而，这些方法都是一种监督式的学习模型，往往要求大量的高质量训练集。这种高质量数据集一般是很难获取的，所以本课题提出了一种在物体检测领域下的利用少量标签样本的自主学习框架。本课题利用 R-FCN 全卷积的网络结构，将大量无标签的训练样本用 CNN 提取出的 ROI，通过当前分类器，分类出按预测置信度排序的样本，将得分高者通过自步学习（Self-paced Learning, SPL）的方式进行伪标签标记，再送入分类器中训练；将得分低者利用一种交互式的主动学习（Active Learning, AL）的方式，用户主动地对该样本进行标签标记，再送入分类器中训练。通过这种迭代训练的方式，可以提高分类器的抗噪性和稳定性，节省大量劳力，检测效果达到了许多最新研究方法的水平。

关键词：物体检测；自主学习；深度卷积神经网络

Title: Self-Learning Framework for Visual Recognition
Major: Automation
Name: Xiaopeng Yan
Student ID: 13351054
Supervisor: Prof. Liang Lin

Abstract

Benefiting from the amazing performance of deep Convolutional Neural Net with large scale training data, object detection, one of the key targets of computer vision, has made incredible progress in the recent years. The general practice of object detection is to extract the region of interest (ROI) by CNN, and then convert it into a classification problem, which can achieve high recognition precision, such as RCNN, Faster RCNN and other popular networks. However, these methods are a kind of supervised learning model, which often requires a large number of high-quality training sets.

In this paper, we propose an self-learning framework which is based on a small number of labeled samples in the field of object detection. This subject uses the network structure of R-FCN fully convolution. A large number of unlabeled training samples are extracted the ROI by CNN using the R-FCN fully convolution network. Through the current classifier, the samples sorted by the predictive confidence are classified and the scores are obtained. The highers are pseudo-labeled by self-paced learning(SPL), and then be sent into the classifier to retrain. For the lowers, the user uses a way of independent active learning (AL). The user autonomously label the samples of low confidence, and then send them into the classifier to retrain. Through this iterative training way , we can improve the classifier's anti-noise and stability, save a lot of labor, and achieve the performances of some state-of-the-arts.

Keywords: Object detection, Self-learning, deep Convolutional Neural Net

目 录

摘 要.....	I
ABSTRACT.....	II
第 1 章 引言.....	1
1.1 选题背景与意义.....	1
1.2 国内外研究现状和相关工作.....	2
1.3 本课题的研究内容与主要工作.....	3
1.4 本课题的论文结构与章节安排.....	3
第 2 章 物体检测领域综述.....	4
2.1 物体检测的定义.....	4
2.2 传统的物体检测算法.....	5
2.3 深度学习的物体检测算法.....	5
2.4 本章小结.....	9
第 3 章 自主学习框架.....	10
3.1 自步学习和主动学习方法研究进展.....	10
3.2 本课题提出的自主学习网络框架.....	11
3.3 数学问题描述与分析.....	15
3.4 本章小结.....	19
第 4 章 实验结果及分析.....	20
4.1 数据集及评测方法.....	20
4.2 实验结果展示.....	22
4.3 实验分析.....	27
4.4 本章小结.....	35
第 5 章 总结与展望.....	36
5.1 工作总结.....	36
5.2 研究展望.....	36

参考文献.....	38
相关的科研成果目录.....	40
致 谢.....	41

第 1 章 引言

1.1 选题背景与意义

近年来，以 DCNN^[1]为里程碑的深度学习革命，正在席卷人工智能界。DCNN^[1]在处理高层视觉问题中有着突出的表现。它被广泛地应用于图像分类，人脸识别，物体检测，图像描述和视觉问答系统等应用中。去年，基于深度神经网络的 AlphaGO 在围棋领域击败了人类围棋大师李世石，攻克了人类在棋类游戏领域的最后一个堡垒，同时也给人工智能领域带来了一丝曙光。然而，深度学习的强大能力依赖于大量有标签的训练样本。目前深度神经网络的训练大多是有监督的，需要大量的带有标签的训练样本进行训练。有监督训练神经网络的训练数据主要来自于手工标注，不过，手工标注需要耗费巨大的人力物力，而移动互联网技术的迅速发展和普及，带来了海量的图像和视频数据。这些数据蕴含着大量的语义信息，诸如物体、场景、人类活动行为等。海量多媒体数据的冲击和其蕴含信息的多样性充分发挥了深度神经网络的潜力，让现代的深度网络可以在多元数据支撑下完成各种智能任务。但是仅仅依赖于手工标注已经无法满足当前的需求，如何利用已有的标注好的数据结合新产生的海量无标签数据共同解决当前的数据需求，是一个亟待解决的问题。

针对上述问题，本课题围绕基于如何在半监督学习环境下进行深度神经网络训练，提出了一种动态标注海量数据的机制。基于该机制，重点研究如何在每阶段选取可靠的数据以赋予相应标签的深度学习网络。深度学习的横空出世对于解决特征表达问题给出了令人满意的答案。随着新数据的加入，微调后的模型提出的特征是随之变化的，本课题旨在挖掘动态变换的特征空间对于研究无标签和有标签数据关系提供重要的指示和线索，从而更好的实现无标签数据的自动标注，和有标签数据相结合，更好的实现深度神经网络的半监督（Semi-Supervised Learning）学习，充分利用与日俱增的海量数据，发挥深度学习的强大优势。

1.2 国内外研究现状和相关工作

近年来,在视觉识别领域下的物体检测已由传统的物体检测算法(SIFT, HOG)逐渐转向基于深度学习的物体检测算法。2014年, Ross B. Girshick 提出了基于 Region Proposal 和 CNN 结合的深度学习算法 R-CNN^[2] 代替使用滑动窗口和手工设计特征的传统物体检测算法,使得物体检测精度大幅度提升,并引领了后来一系列的深度学习物体检测方向,如随后出现的 SPP-Net^[3], Fast-RCNN^[4], Faster-RCNN^[5], R-FCN^[6]等。而后又有一些基于回归方法的深度学习物体检测算法,如 YOLO^[7], SSD^[8]等被提出,其检测速度可以达到实时的效果,检测性能与基于 Region Proposal 的方法精度相近。然而,这些方法都是一种有监督的学习方法,需要大量的有标签(类别标签和位置标签)的数据样本,而标记大量的标签数据是需要耗费大量的人力物力的,所以本课题在此基础上提出一种利用少量标签数据进行物体检测的方法。

使用少量有标签的数据样本,充分利用大量的无标签数据资源进行学习的方法属于一种半监督的学习方法,目前比较主流的半监督学习方法有“阶梯”网络(Ladder network),基于对抗式生成网络(GAN)的半监督学习方法等。但是事实上目前的半监督学习方法还存在以下局限性:(1)所依据的假设,不支持特征学习。(2)对大规模的未标注数据,难以高效地选择样本以进行半监督学习。针对上述问题,本课题围绕基于如何在半监督学习环境下进行深度神经网络训练,提出了一种动态的自主学习标注海量数据的机制。在主动学习下,如何选取“困难”的无标签样本数据往往是问题的关键,以前^{[9],[10],[11]}都是基于 SVM (Support Vector Machine)等方法选择不确定性和确定性的样本,随后有一些算法^{[12],[13]}是根据未标记数据的分布和多样性来选择样本。而在自步学习^[14]中,2009年, Bengio 等人根据来自人类和动物的认知原则,提出了一个课程学习^[15](Curriculum Learning)的概念,一个模型的学习过程是通过逐渐加入从简单到复杂样本的训练过程。随后,2010年,有研究人员基于此提出了自步学习^[16]的方法,其中包括对所有样本的正则化项和在样本权重加权损失项。基于此框架各种变种的自步学习方法被提出,例如有 SPaR^[17], SPLD^[18], SPCL^[19]等。

1.3 本课题的研究内容与主要工作

本课题在全监督的网络模型 R-FCN^[6]的基础上改进为一种半监督的深度学习的方式，提出了一套不依赖于额外先验信息，在少量标签数据下的半监督训练框架。本课题的研究出发点始于深度学习的半监督多分类问题，将其思想运用到物体检测领域，重点在于考虑跨特征空间之间转换的语义局部保留。在原始的训练模型上引入来自外界的海量的无标签数据，模型通过自步学习的方式对外来数据进行预测，对每个无标签样本在每一次特征空间转换后都会被赋予一个置信度和标签值。我们可以根据当前需要来选取置信度高的样本赋予对应标签加以训练，每一阶段训练后，会有部分具有高置信度的无标签数据被赋予标签，然后被选入训练更新模型。更新后的模型产生新的特征空间会在和上一阶段的特征空间相比较。这样循环过程中不断更新特征空间。同时，在此基础上用户通过主动学习的方式选择少量低置信度的样本人工标注标签，再加以训练。通过这种迭代交替的自主学习方式，模型性能逐步增加，在少量标注数据的情况下，识别精度达到相同数据情况下全标签的识别精度。

1.4 本课题的论文结构与章节安排

本课题共分为五章，章节内容安排如下：

第一章：引言——主要介绍了本课题的选题背景及意义，国内外研究状况及相关工作和本课题研究的内容和主要工作。

第二章：物体检测领域综述——主要介绍了物体检测的定义，传统的物体检测算法和基于深度学习的物体检测算法。

第三章：自主学习框架——主要介绍了自步学习和主动学习方法的研究进展，本课题提出的自主学习网络框架及其数学问题描述和分析。

第四章：实验结果及分析——主要介绍了相关数据集和评测标准，实验结果展示和实验分析。

第五章：总结与展望——主要介绍了本课题的工作总结和该研究方向的未来工作展望。

第 2 章 物体检测领域综述

2.1 物体检测的定义

2.1.1 物体检测 (Object detection) 是计算机视觉领域一个基础性的经典课题。物体检测任务分为两个子任务：目标分类和目标定位。目标分类是判断输入的图像中是否有感兴趣的类别的目标出现，并输出带分数的类别标签来表示图像中检测到的感兴趣的目标最可能属于的类别；目标定位是定位出输入的图像中感兴趣的类别的目标的位置，并输出目标的可能的坐标，一般用矩形盒 (bounding boxes) 表示，如图 2-1 所示。

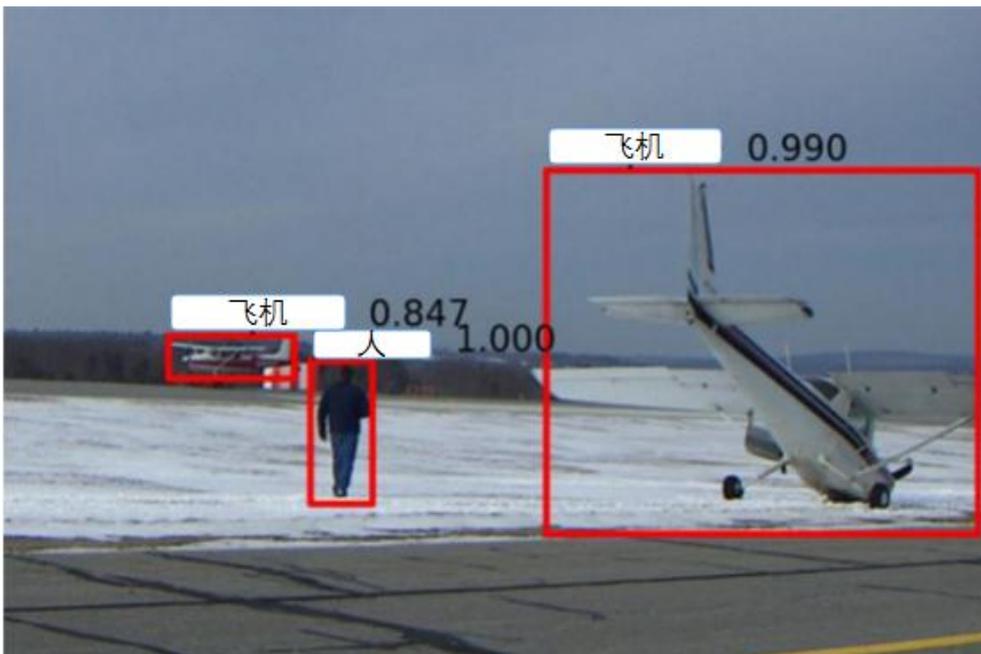


图 2-1 物体检测

2.1.2 物体检测对计算机视觉领域和实际应用具有重要意义，它是很多高级视觉任务的基础，如目标跟踪，事件识别，视频理解等，而且也被应用到很多的实际任务中，如智能检测，无人车，机器人导航等。尽管每年都有新的研究成果发表，准确率一步步的被刷新，然而，当前物体检测算法在实际应用中精度仍

然较低，且无法达到实时性的要求，不能应用于现实中手机的检测任务中。因此，当前物体检测任务仍旧是重要的挑战性的研究课题。

2.2 传统的物体检测算法

传统的物体检测算法主要分为三步：第一步，在输入的图像中提取出一些候选区域；第二步，对这些候选区域提取特征；第三步，使用分类器根据提取的特征对候选区域分类。如下图 2-2 所示：

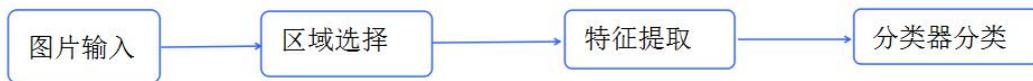


图 2-2 传统检测算法流程

2.2.1 候选区域的选择

传统的对候选区域的选择一般是采用滑动窗口的策略，对整张输入图像进行遍历，并设置不同的长宽比和尺度大小。这种穷举的方法虽然可以包含目标可能出现的位置。但是复杂度高，冗余性大，速度慢。

2.2.2 特征提取

这一步的特征提取一般是按照人工经验手动设计特征，常用的特征有 SIFT, HOG 等，然而这种特征一般不具有很强的鲁棒性，不能很好地适应图像对光照，形态和背景等多样性变化。

2.2.3 分类器分类

提取出的特征好坏直接影响分类器对目标类别的分类效果，目前常用的分类器是 SVM, Adaboost 等。

2.3 深度学习的物体检测算法

基于传统的物体检测算法使用滑动窗口选择候选区域复杂度高，冗余性大和手工设计的特征没有很好的鲁棒性的缺点，且随着大规模数据集的发展，如

ImageNet^[20], Caltech^[21], Pascal (Pattern Analysis, Statical Modeling and Computational Learning) VOC^{[22][23]} 和 COCO^[24] 等数据集和计算机计算能力的提高, GPU (Graphics Processing Unit) 的使用使得计算精度和计算速度都有了很大提升, 使用 CNN 的深度学习物体检测算法逐渐发展起来。深度学习物体检测算法主要分为基于 Region Proposal 的方法和基于回归的方法。

2.3.1 基于 Region Proposal 的物体检测算法

从 2014 年开始, 物体检测取得巨大进步。Ross B. Girshick 使用 Region Proposal 和 CNN 结合的方法, 提出 R-CNN^[2] 的网络框架, Region Proposal 方法的深度学习算法快速发展, 先后出现了 SPP-Net^[3], Fast-RCNN^[4], Faster-RCNN^[5] 和最新的 R-FCN^[6] 网络。Region Proposal 方法的主要步骤是 (i) 提取 ROI (ii) 类别分类和位置回归。各种研究方法提取 ROI 和分类定位方法不太一样。例如, R-CNN^[2] 网络是使用选择性搜索 (Selective Search) 方法提取 ROI, 再使用 CNN 提取特征, 最后使用 SVM 分类器做类别分类, 使用线性回归的方法微调位置坐标。R-CNN^[2] 具有以下一些缺点: (i) 训练过程复杂: 微调网络+训练 SVM 分类器+训练边框回归器。(ii) 训练速度慢: 使用 GPU, VGG16^[1] 模型处理一张图片需要 47s, 每一个候选区域都要经过 CNN 网络。(iii) 训练耗时, 占用存储空间大: 5000 张图片产生几百 G 的特征文件。训练过程如下图 2-3^[2]。

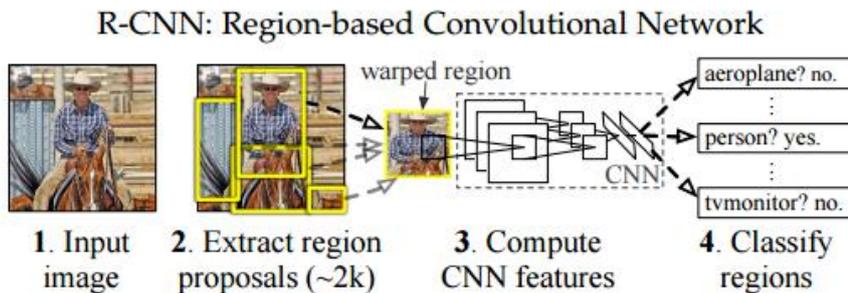
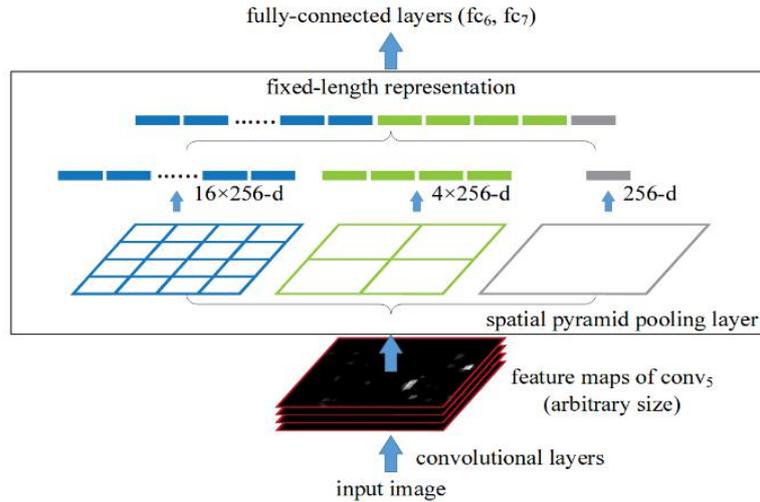
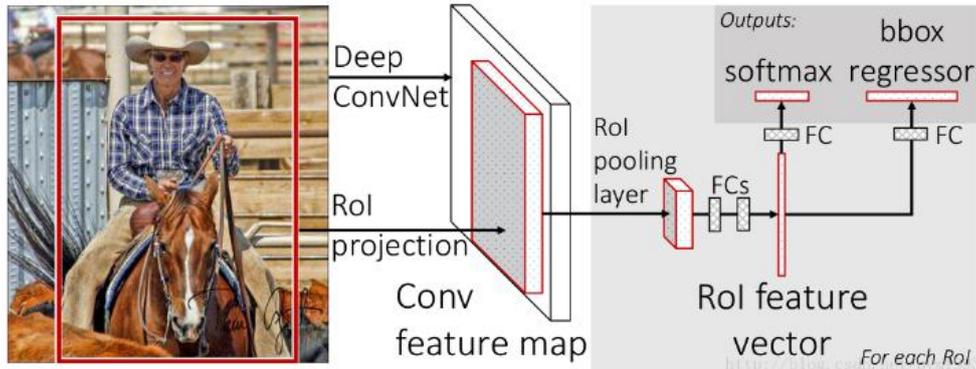


图 2-3 R-CNN^[2] 架构

针对 R-CNN^[2] 检测速度慢的问题, SSP-Net^[3] 被提出, 它使用一整张图片进行卷积, 将特征图对应的区域映射到候选区域, 再通过 spp-net 层将所候选区域统一到一个维度, 输入到全连接层, 检测速度和精度都得到了提高, 如图 2-4^[3]。

图 2-4 SPP-Net^[3]架构

Fast-RCNN^[4]集合了 R-CNN^[2]和 SPP-Net^[3]的精髓, SPP-Net^[3]和 R-CNN^[2]都需要多阶段的训练, 设计了 roi-pooling 层, 它是 spp-net 层一种精简版, 设计了多任务损失函数(multi-task loss), 使用 softmax 分类器和 box regression 回归方法, 将分类任务和边框回归统一到了一个框架之内, 图 2-5^[4]。

图 2-5 Fast-RCNN^[4]架构

Fast-RCNN^[4]缺陷在于仍然没有解决 selective search 进行候选框选择的时候计算速度慢的问题。Faster-RCNN^[5]网络是以整张图像为输入, 直接利用 RPN (Region Proposal Networks)网络来计算候选框, 代替了 selective search 选择候选框, 并使用 bounding boxes regression 微调候选框位置与大小, 最后使用 softmax 方法进行目标分类, Faster-RCNN^[5]如图 2-6。

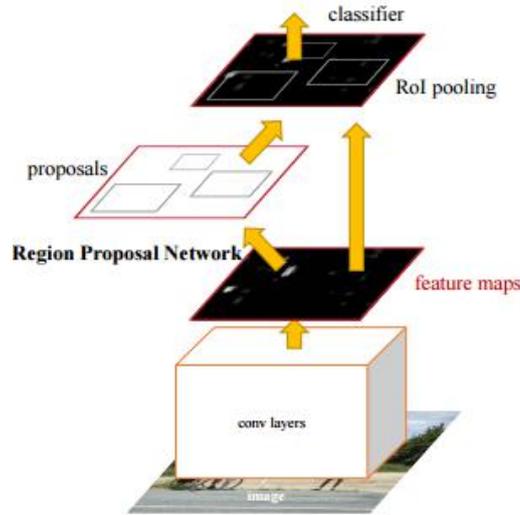


图 2-6 Faster-RCNN^[5]架构

2.3.2 基于回归方法的物体检测算法

尽管基于 Region Proposal 方法的物体检测算法精度逐渐提高，但其速度上并不能满足实时性的要求。因此出现一类基于回归方法的物体检测算法，对于给定的输入图像，直接在图像上回归出每个目标的位置的坐标以及目标类别。基于回归思想的深度学习物体检测算法主要有 YOLO^[7]和 SSD^[8]网络，其检测精度虽稍不及基于 Region Proposal 的方法，但其检测速度可以达到实时性的要求。YOLO^[7]是将一张输入图片划成 7 x 7 的网格，对每个网格，预测每个边框是目标的置信度和每个边框在多个类别上的概率。根据预测出的 7*7*2 的目标窗口，根据阈值除去可能性较低的目标窗口和使用非极大值抑制的方法除去冗余窗口，直接回归出目标类别和位置坐标，网络结构如图 2-7^[7]。

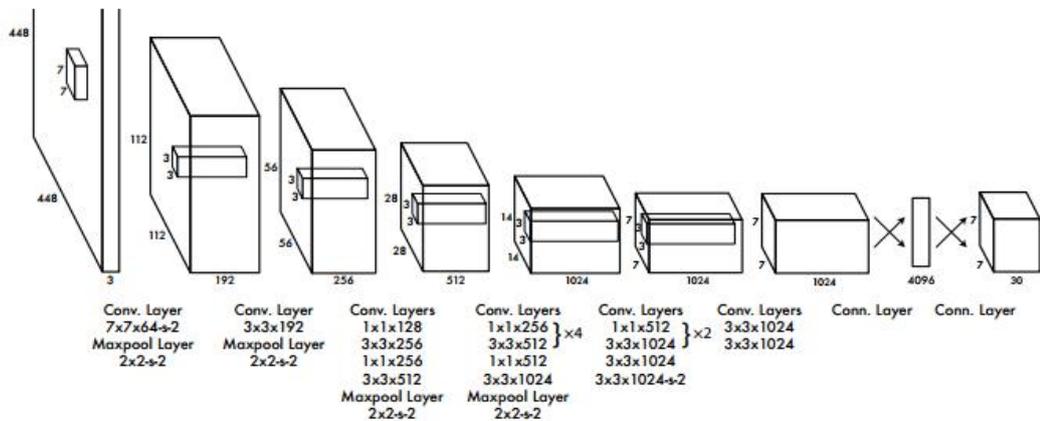


图 2-7 YOLO^[7]架构

鉴于 YOLO^[7] 定位精度低的缺陷, SSD^[8] 结合 YOLO^[7] 的回归思想和 Faster-RCNN^[5] 的 anchor 机制, 使用多尺度区域进行回归, 既保证了检测速度, 提高了检测精度。SSD^[8] 首先第一步都是利用卷积神经网络提取特征, 在最后的卷积层对前面各层每个尺度的特征图运用 anchor 的机制提取候选框, 回归出物体的类别和位置。SSD^[8] 网络在 VOC2007^[22] 上检测精度可达 79.8%, 速度在 GPU 上达 58 帧每秒, SSD^[8] 如图 2-8。

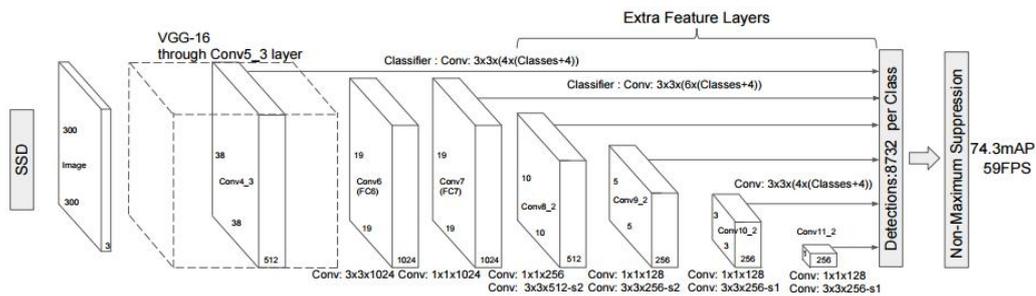


图 2-8 SSD^[8] 架构

2.4 本章小结

本章主要介绍了物体检测的相关定义, 传统的物体检测的方法, 主要有三个步骤——区域选择, 特征提取和分类器分类, 目前的流行的深度学习的检测方法, 基于 Region Proposal 的方法, 如 Faster-RCNN^[5] 和基于回归的方法, 如 SSD^[8]。

第 3 章 自主学习框架

3.1 自步学习和主动学习方法研究进展

在大数据移动互联网时代,针对真实大数据做信息处理和信息挖掘已经成为席卷整个科技界的前沿热点。然而现实中的大数据有三个显著特点:一是“大”,数据量巨大;二是“脏”,数据具有很强的多样性;三是“无监督”,大量数据未经人工刷选和标注。因此,从海量“大”规模的“无监督”“脏”数据中,构建有效的机器学习方法对数据进行挖掘已成为研究热点。

针对这一问题,由深度学习的鼻祖之一, Montreal 大学的 Yoshua Bengio 教授团队于 2009 年提出的课程学习^[15]理论提供了一套行之有效,具有理论意义的方法。该理论的思想是:根据人类或者动物的认知机理,学习是从简单的,普适的课程,然后逐渐增加难度,过渡到学习更复杂,更专业的知识,以此完成对复杂对象的认知。人们受教育学习各种“课程”的过程正是按照这种方式来获得知识和能力的。模拟人类和动物的这种学习过程,我们可以将训练数据按照其对学习目标的难易程度,从易到难的学习,以这种方式完成数据学习和认知过程,如图 3-1 所示。

在 2010 年的 NIPS 会议上 Stanford 大学的 Daphne Koller 教授团队进一步将该理论模型化,初步以数学公式的形式模型化了自步学习^[16]。随后几年,科研人员系统地分析了自步学习的内在机理,提出了自步学习的数学公理化的构造条件,并根据不同的任务目标,延伸出了各种实用的自步学习方法,如:对视频动作识别构造了 SPLD^[18]算法;对多媒体多模态事件检测提出了 SPaR^[17]算法等,这些算法针对各种应用问题,均有良好的表现。

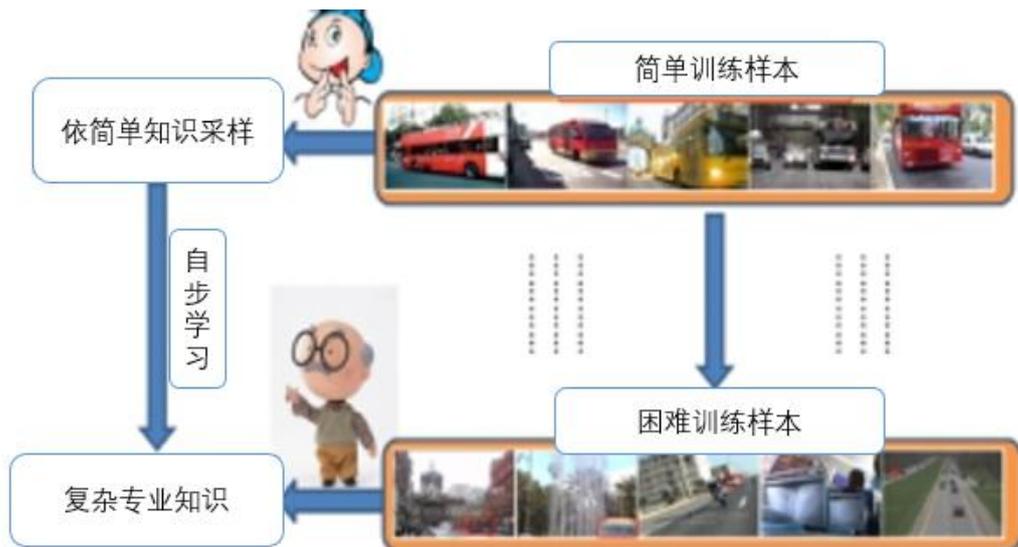


图 3-1 自步学习方法

主动学习这一分支主要关注的是如何选择样本的策略，即如何挑选出最具有意义的无标签的数据用于人工标注。一个最常见的策略是 Lewis 和 Gale 1994^[12] 年和 Tong and Koller 2002^[13] 年提出的确定性选择，确定性的衡量标准是从初始的分类器对新的无标签的数据的预测。几种^{[9],[10],[11]}基于 SVM 确定不确定的样本和测量样本确定性的方法都比较接近决策边界。一些后续的方法在选择样本时利用无标签数据的信息密度测量，例如卡普尔等人采用实例最大化增加候选实例的方法和基于高斯过程模型的利用剩余实例的相互信息的方法，同时考虑在无标签数据选择实例的多样性。2013 年，Elhamifar^[25] 提出了一套基于凸规划的通用框架，同时考虑样本选择的不确定性和多样性。

3.2 本课题提出的自主学习网络框架

如图 3-2 所示为本课题提出的自主学习物体检测的网络框架，本框架包含 CNN 提取特征和模型初始化阶段，分类器更新阶段，通过自步学习标注高置信度无标签样本阶段和通过主动学习标注低置信度的无标签样本阶段，箭头所示为 workflow。仿照在人类教学过程中“学生预习，做作业；老师批改；学生通过批改的作业再学习”的过程，主动学习模型对于“学生”（分类器）的少量“作业”（低置信度样本）进行人工标注，然后使用标注的数据重新训练分类器。下面分别介绍这些阶段。

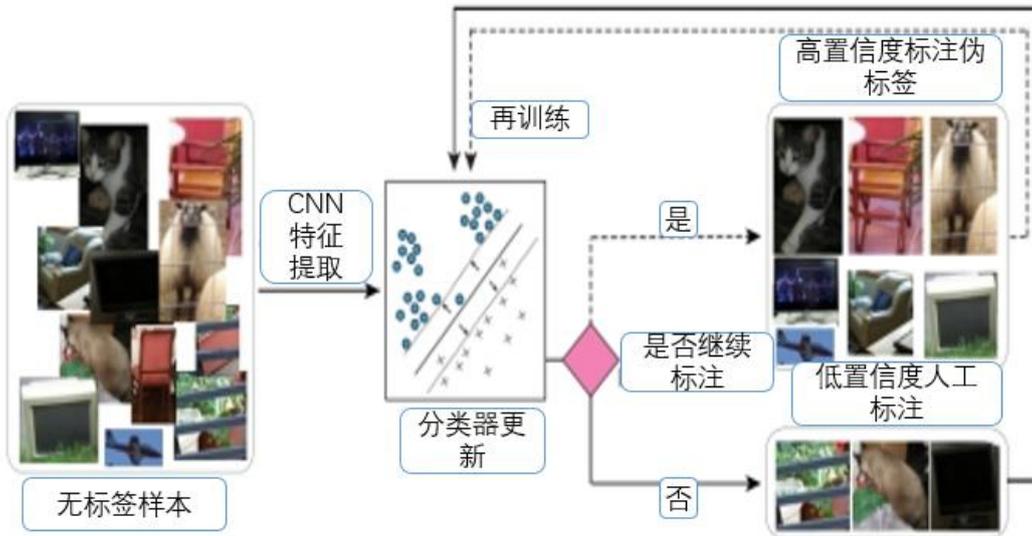


图 3-2 自主学习网络框架

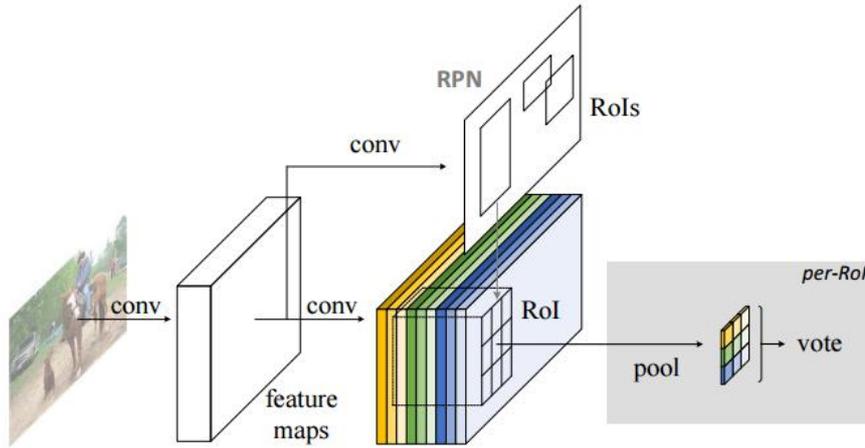
3.2.1 CNN 特征提取和模型初始化

本课题的 CNN 采用的是 R-FCN^[6] 的网络，其是目前在 Pascal VOC2012^[23] 物体检测竞赛上识别精度最高的网络框架。

(1) 它是基于 Region Proposal 方法的一种新的物体检测方法，移除了最后的全连接层的全卷积网络，速度得到了提升，并且用到了著名学者何凯明的 Residual Network-101^[26]。该网络的思想是移除重复计算的子网络，让所有的计算都可以共享。

(2) 在计算机视觉领域，图像分类要求平移不变性，而物体检测要求有平移变化。因此，图像分类受平移不变性的全卷积结构的青睐。目标检测任务要求平移变化的定位表示。例如，图像中目标的位置的移动应该对网络产生响应，这些响应对选择覆盖目标的候选框的好坏是很有意义的。因此其采用了全卷积网络结构，并且引入了平移变化的结构，即用卷积层构建位置敏感分数地图层 (position-sensitive score maps)。每个位置敏感分数地图编码 ROI 的相对空间信息，且加入 RIO pooling 池化层来监控这些位置敏感分数地图。

(3) 网络最后阶段采用投票 (vote) 的方式将每个类的分数送入分类器判断属于哪一类，此处的分类器采用的 softmax。位置定位过程也是如此。网络框架如下图 3-4^[6] 所示。

图 3-3 R-FCN^[6] 架构

对于模型的初始化,我们是先采用该 R-FCN^[6]网络,数据集采用 VOC2007^[22],得到一个初始的模型和相应的参数为我们的初始化模型,初始化的识别精度为 73.9%,识别精度较高,对于自步学习预测无标签样本准确率较好。

3.2.2 分类器的更新阶段

分类器的更新是通过迭代训练的方式来更新模型的参数,从初始化模型开始,通过自步学习阶段,预测出置信度高的无标签的样本数据,并标注图像中预测的目标的类别和定位,置信度低的样本输入到主动学习阶段,通过人工主动标注图像的目标的类别和定位。将这些样本数据再加入原先的样本中,通过前一个模型继续训练这些数据,以此迭代的方式更新模型参数,提高模型的稳定性和抗噪性。在分类器更新阶段,为了使我们的网络能够适应各种图片,提高模型性能,我们对图片进行了数据增强处理。(i) 水平翻转, (ii) 多尺度训练, (iii) 图片旋转, (iv) 批量处理, (v) 色彩变化等。在训练时,我们控制了正负样本的比例大致是 3:1,这有利于提高模型精度。在模型迭代阶段,因为需要从自步学习和主动学习过程中加入样本,所以我们在迭代微调分类器的阶段,需要调整模型的超参数,如设置最大迭代次数,设置学习率,设置 stepsize 等等,这是一个需要靠人工经验的过程,具有经验性,需要多次尝试调整,在测试阶段,我们采用了多尺度和水平翻转的测试方法,这些策略可以有效地提高识别精度。

3.2.3 自步学习阶段

自步学习过程是当前模型预测外来无标签的样本数据并给其标注类别标签和位置标签的过程。在通过当前分类器选择置信度高的无标签样本时,我们采取

了一些策略。(i)若图片中一个物体被分类器预测出两个或以上的类别,虽然其置信度大于我们设置的阈值(0.9),我们仍认为该样本属于“困难”样本,因为其使我们的分类器感到“困惑”,不对其进行伪标签标注。(ii)在设置置信度时,除了人工根据经验设置置信度阈值(0.9),我们还采用了分别给每个类别设置不一样的置信度,具体做法是先给所有类别设置一个较低的置信度,让当前的分类器预测我们的一定数量无标签的样本,并计算预测出的每个类别的数量和置信度之和,再计算每个类别的置信度平均值,然后我们按这个平均值给每个类别分别设置置信度阈值。且为每个类别的阈值分别取平均值乘以1, 1/2, 1/3等不同设置。(iii)在分类器预测样本时,在某些置信度较高的情况下,分类器会误判,比方说“猫”判为“狗”,然后这种情况在总体样本里占比很少,所以我们提出了一个基于大数定理的采样策略来最大可能的排除这些误判样本。在预测的样本中我们采用多次随机采样的方式对样本采样,然后对多次采样求交集作为我们预测的样本,以这种采样方式尽量地去排除误判样本。(iv)为了减少误判样本的数量,我们还采用了多尺度的预测方式进行预测,在多种不同尺度下,对同一张图片进行预测,在一些情况下,不同尺度图片中的物体被预测出的类别不一样,对于被预测出不一样的图片,我们给予排除,减少被误判的概率。

3.2.4 主动学习阶段

主动学习过程是当前模型预测外来无标签的数据样本,将置信度很低,对图像中目标预测错误或者模棱两可,且定位信息不准确的“困难”样本加入人工标注阶段的过程。用户通过标注工具对“困难”样本中的目标进行类别和位置的标注,再将这些标注了的“困难”样本加入模型中继续参与训练。人工标注标签是一个非常耗时费力的过程,按照VOC2007^[22]的数据集格式XML,对一张图片中的属于那二十类的物体进行标注,标注时需要选择物体的类别和框出物体的坐标位置。对于遮挡严重的物体,我们不给予标注;若两种不同物体重叠度较高,难以区分开时,也不给予标注。对于图片中有多个物体时,需要单独对每个物体标注,尽量不要重合。对于预测的图片和人工标注的图片,我们都做了可视化的处理,以确保我们标注的准确性。

3.3 数学问题描述与分析

在物体检测中，我们假设在 m 个样本类别中有 n 个目标候选区域。定义训练样本空间： $D = \{x_i\}_{i=1}^n \subset R^d$ ，其中 x_i 代表第 i 个候选区域的 d 维的特征表达。我们有 m 个分类器对每个候选区域进行类别判断。分类器可以为 SVM, Softmax 等分类器，此处以 SVM 举例说明。于此相关，我们定义样本 x_i 的标签集为： $y_i = \{y_i^{(j)} \in \{-1,1\}\}_{j=1}^m$ ，其中 $y_i^{(j)}$ 代表 x_i 样本的第 j 个类别。例如， $y_i^{(j)}=1$ ，代表 x_i 属于第 j 个类别。

本课题提出的自主学习框架的公式化描述如下：

$$\min_{\{W, b, V, Y\}} \sum_{j=1}^m \frac{1}{2} \|w^{(j)}\|_2^2 + C * L(w^{(j)}, b^{(j)}, D, y^{(j)}, v^{(j)}) + f(v^{(j)}; \lambda_j) \quad (3-1)$$

$$\text{s. t. } V \in \Psi^\lambda, \sum_{j=1}^m |y_i^{(j)} + 1| \leq 2, y_i^{(j)} \in \{-1,1\}, i \notin \Omega^\lambda$$

其中 $W = \{w^{(j)}\}_{j=1}^m \subset R^d$ 和 $b = \{b^{(j)}\}_{j=1}^m \subset R$ 分别代表的是 m 个分类器的权重和偏置参数， C ($C>0$) 是标准正则化参数， $V = \{v^{(j)}\}_{j=1}^m = \{[v_1^{(j)}, v_2^{(j)}, \dots, v_n^{(j)}]\}_{j=1}^m$ 定义了训练样本的重要性权重， λ_j 是一个步长参数，控制第 j 个分类器的学习步长， $f(v^{(j)}; \lambda_j)$ 是自步学习的正则化控制。我们定义所有当前主动学习样本的索引集合为： $\Omega^\lambda = \bigcup_{j=1}^m \{\Omega^{\lambda_j}\}$ ，其中 Ω^{λ_j} 表示第 j 个类别集合的步长 λ_j ， Ω^λ 是为了约束 y_i 。 $\Psi^\lambda = \bigcap_{i=1}^m \{\Psi_i^\lambda\}$ 组成在 m 个分类器的步长为 $\lambda = \{\lambda_j\}_{j=1}^m$ 时的课程约束。特别的，我们定义两种可选的对每个样本 x_i 的课程约束。(i) $\Psi_i^\lambda = [0,1]$ ，对于无标签的样本， $i \notin \Omega^\lambda$ ，其对于所有分类器的重要性权重 $\{v_i^{(j)}\}_{j=1}^m$ 在自步学习过程中优化；(ii) $\Psi_i^\lambda = \{1\}$ 是对于在主动学习过程中有标签的样本，它们的权重在训练过程中被设置。注意，这个 Ψ_i^λ 相对于 m 个分类器的步长 λ 可以动态的改变。

在训练集 X 上定义的损失函数为：

$$\begin{aligned}
L(w^{(j)}, b^{(j)}, D, y^{(j)}, v^{(j)}) &= \sum_{i=1}^n v_i^{(j)} l(w^{(j)}, b^{(j)}; x_i, y_i^{(j)}) \\
&= \sum_{i=1}^n v_i^{(j)} (1 - y_i^{(j)} (w^{(j)T} x_i + b^{(j)}))_+ \\
\text{s. t. } \sum_{j=1}^m |y_i^{(j)} + 1| &\leq 2, y_i^{(j)} \in \{-1, 1\}, i \notin \Omega^\lambda
\end{aligned}$$

其中 $l(w^{(j)}, b^{(j)}; x_i, y_i^{(j)})$ 代表样本 x_i 在第 j 个分类器上的汉明损失，该损失函数 L 是所有分类器的损失之和。约束条件用于惩罚除下列两种外的情况：(i) $y_i^{(j)}$ 被预测为正样本，其余为负样本；(ii) $y_i^{(j)}$ 被所有分类器预测为负样本，例如：背景。

交替最小化策略被采用来解决这个优化。具体而言，通过梯度下降交替更新分类器参数 w, b ；通过自步学习过程更新样本重要性权重 V ；通过重复排序过程更新伪标签数据集 Y 。此外，课程的约束 Ψ^λ 通过主动学习过程逐渐增加的步长参数 λ 。下面，详细介绍优化过程和相关物理解释。

模型初始化：首先我们用预训练好的特征 ImageNet model 代表初始化特征，定义样本集 $\{x_i\}_{i=1}^n$ ，设置 m 个分类器的初始化步长参数 $\lambda = \{\lambda_j\}_{j=1}^m$ ，用当前用户标注的样本 Ω^λ 初始化课程约束 Ψ^λ ， Y 和 V 。

分类器更新：我们通过梯度下降更新分类器的权重和偏置参数 w, b ；固定 $\{\{x_i\}_{i=1}^n, V, Y, \Psi^\lambda\}$ ，原始的自主学习框架的公式化描述方程 (3-1) 可被等价公式化为对每个分类器 $j=1, 2, \dots, m$ 的次优化问题：

$$\min_{\{w^{(j)}, b^{(j)}\}} \sum_{j=1}^m \frac{1}{2} \|w^{(j)}\|_2^2 + C * \sum_{i=1}^n v_i^{(j)} l(w^{(j)}, b^{(j)}; x_i, y_i^{(j)}) \quad (3-2)$$

这是一个标准的分类器的模型，将一个类别当作正样本，其余类别做负样本。当重要性权重 $\{v_i^{(j)}\}_{j=1}^m$ 只取值为 $\{0, 1\}$ 时，对于简单的 SVM 模型，在 $\{v_i^{(j)}\}_{j=1}^m = 1$ 的情况下，该问题可被很多高效求解器求解。这步可以解释为主动学习和自步学习下的分类器的不确定性建模。

自步学习过程：我们通过无标签样本的置信度高低排序标注伪标签，在固定

$\{W, b, \{x_i\}_{i=1}^n, Y, \Psi^\lambda\}$ 的情况下，我们根据 V 给样本排序，通过方程 (3-1) 得到如下方程：

$$\min_V \sum_{j=1}^m C * \sum_{i=1}^n v_i^{(j)} l(w^{(j)}, b^{(j)}; x_i, y^{(j)}) + f(v^{(j)}; \lambda_j)$$

$$\text{s. t. } V \in \Psi^\lambda$$

这个问题成为了标准的自步学习问题，自步正则项 $f(v^{(j)}; \lambda_j)$ 和课程约束 Ψ^λ 都是凸集，可用像基于梯度或者内点的方法解决。这里我们有多种自步正则项的选择，比如正则惩罚样本线性损失的权重。这里，我们有

$$f(v^{(j)}; \lambda_j) = \lambda_j \left(\frac{1}{2} \|v^{(j)}\|_2^2 - \sum_{i=1}^n v_i^{(j)} \right) \quad (3-3)$$

其中 $\lambda_j > 0$ ，方程 (3-3) 对于 V 是凸优化，通过梯度可以找到全局最优解。

考虑到 $\{v_i^{(j)}\} \subset \{0,1\}$ ，我们定义其为：

$$v_i^{(j)} = \begin{cases} -\frac{C l_{ij}}{\lambda_j} + 1, & C l_{ij} < \lambda_j \\ 0, & \text{otherwise} \end{cases}, \quad (3-4)$$

当 $C l_{ij} < \lambda_j$ 时， $v_i^{(j)}$ 取上等式，其余情况取 0， $l_{ij} = l(w^{(j)}, b^{(j)}; x_i, y_i^{(j)})$ 是 x_i 对第 j 个分类器的损失。我们定义重要性样本为 $S = [S_1, S_2, \dots, S_m] (|S_j| \leq k)$ ，固定 $\{W, b, V, \{x_i\}_{i=1}^n, \Psi^\lambda\}$ ，我们优化方程 (3-1) 的 y_i 相当于解下面的方程：

$$\min_{y_i \in \{-1,1\}^m, i \in S} \sum_{j=1}^m v_j^{(j)} l_{ij}, \text{ s.t. } \sum_{j=1}^m |y_i^{(j)} + 1| \leq 2 \quad (3-5)$$

其中 v_i 固定视为常数，方程 (3-5) 可由下面的定理求出最优解。

定理：当不是所有 $w^{(j)T} x_i + b^{(j)} \leq 0 (j = 0, 1, \dots, m)$ ，方程 (3-5) 对于有最优解如：

$$y_i^{(j)} = \begin{cases} +1, & j=j^*, v_i^{(j)} \neq 0 \\ -1, & \text{otherwise} \end{cases}, \text{ 其中 } j^* = \underset{1 \leq j \leq m}{\operatorname{argmax}} (w^{(j)T} x_i + b^{(j)})$$

在方程 (3-5) 中 $\{v_i^{(j)}\} \subset \{0,1\}$ 非负，不影响优化过程。这步的解释为我们

反复地在当前分类器的基础上找到高置信度的样本，然后给其标注伪标签。

主动学习过程：我们通过无标签样本的置信度高低排序，通过人工主动学习的方式更新课程约束 Ψ^λ 。主动学习过程是选择置信度最低的样本将其人工标注为正或负样本，我们选择的标准是基于不确定策略 (Lewis and Gale 1994^[12]; Tong and Koller 2002^[13])。当我们利用当前分类器标注无标签的样本时，若对同一个目标预测出了多个标签，这代表该样本使分类器“困惑”，我们将其称为低置信度，需要人工标注。在人工标注完后，我们通过新标注的样本集 Θ 更新课程约束 Ψ^λ 。对每个被标注的样本，我们的主动学习过程执行如下两步操作：(i) 设置课程约束： $\{\Psi_i^\lambda\}_{i \in \Theta} = \{\emptyset\}$ (ii) 更新样本的标签 $\{y_i\}_{i \in \Theta}$ 并且将其索引加入到当前标注的样本数据集中。

步长参数更新：我们对每个分类器设置初始步长 λ^0 ，利用启发式策略更新参数，对每次迭代我们更新步长参数：

$$\lambda_j^k = \begin{cases} \lambda_j^{(k-1)} + \alpha * \eta_j^k, & 1 \leq k \leq \tau, \\ \lambda_j^{(k-1)}, & k > \tau, \end{cases} \quad (3-6)$$

其中 η_j^k 代表第 j 个分类器在验证集上的平均精度， α 是一个控制步长增长的参数，当所有样本的权重都为 1 时，步长参数停止更新，引入阈值 τ ，当 $k < \tau$ 时， λ 可更新。

整个算法优化过程可用下面算法表示：

输入: $\{x_i\}_{i=1}^n$

输出: 模型的参数 $\{W, b\}$

优化过程: 用预训练的 CNN 初始化 $\{x_i\}_{i=1}^n$, 引入多标签 $\{y_i\}_{i=1}^n$, 课程约束 Ψ^λ , 重要性权重 V 和步长参数 $\lambda = \{\lambda^0\}^m$

While True do:

For all iter=1, ..., T do

通过梯度下降更新分类器 W

基于方程 (3-4) 通过自步学习更新 V

通过再排序更新 S

基于方程 (3-5) 标注伪标签更新 $\{y_i\}_{i \in S}$

End for

更新置信度低的样本集 Θ

If Θ 不为空 do

通过主动学习更新 $\{y_i\}_{i \in \Theta}$ 和 $\{\Psi_i^\lambda\}_{i \in \Theta}$

Else

Break

End if

每隔 t 次迭代通过方程 (3-6) 更新 λ

End while

Return $\{W, b\}$;

3.4 本章小结

本章阐述了自步学习和主动学习方法的研究进展, 然后介绍了本课题提出的自主学习网络框架, 包括了 CNN 特征提取和模型初始化, 分类器更新, 自步学习和主动学习阶段, 最后详细给出了相应的数学公式分析和相关物理解释。

第 4 章 实验结果及分析

4.1 数据集及评测方法

4.1.1 数据集

本课题用到的数据集是 PASCAL VOC 2007^[22]，VOC2012^[23]和 MS COCO^[24]数据集。PASCAL VOC 挑战赛是视觉对象分类和物体检测的一个基准测试，提供了检测算法和学习性能的标准注释数据集和标准的评估方法。VOC2007^[22]中一共有 20 类生活常见目标类别，分别是人类；动物（鸟、猫、牛、狗、马、羊）；交通工具（飞机、自行车、船、公共汽车、小轿车、摩托车、火车）；室内（瓶子、椅子、餐桌、盆栽植物、沙发、电视）。图片都取自于日常生活场景，训练集，验证集和测试集一共有 9963 张图片，包含了 24640 个被标注的物体。VOC2012^[23]中是这 20 类目标类别，增加了目标分类和定位难度，如遮挡，目标变小，光照变化等，训练集和验证集一共有 11540 张图片，包含 27450 个被标注的物体。VOC2012^[23]不提供测试集的标注，测试结果需要提交 PASCAL VOC2012^[23]的服务器对比才知道。所以，在我们的实验中，测试集采用的是 VOC2007^[22]的测试集，共 4952 张图片和 VOC2012^[23]的测试集，共 11992 张图片。图片的标注采用的 xml 的格式，包含的字段主要有物体类别（name），位置坐标（Xmin, Ymin, Xmax, Ymax）和一些其他字段。

4.1.2 评测方法

物体检测需要回答的是“什么物体在图片的哪个位置？”，因此，需要预测图片中的物体的一个带置信度水平的位置坐标（Xmin, Ymin, Xmax, Ymax）和所属

的类别信息。在类别判断正确的情况下，若预测的位置坐标和测试图片的标准的位置坐标（ground truth）重叠度（IOU）大于 0.5 时，计算方式如图 4-1，认为是正确的正样本，否则认为是错误的正样本；在类别判断错误的情况下，若预测的位置坐标和测试图片的标准的位置坐标重叠度（IOU）大于 0.5 时，认为是正确负样本，否则认为是错误的负样本。其实质是计算一个混淆矩阵，相关计算方法如下：

$$\text{IOU 计算: } \text{IOU} = (A \cap B) / (A \cup B)$$

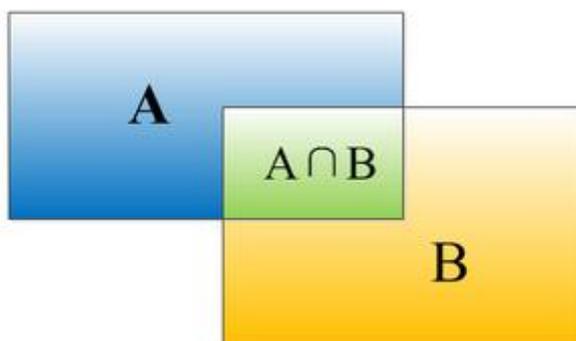


图 4-1 IOU 的示意图

表 4-1 混淆矩阵

	实际正样本	实际负样本
预测正样本	TP	FP
预测负样本	TN	FN

因此，其参数指标的计算方法如下：

$$\text{召回率} = \frac{TP}{TP + FN}$$

$$\text{精确度} = \frac{TP}{TP + FP}$$

$$\text{正样本比率} = \frac{TP}{TP + FN}$$

$$\text{负样本比率} = \frac{FN}{FN + TN}$$

在物体检测的评测中需要计算召回率和精确度，绘制 PR 曲线，如图 4-2。通常，召回率高时，精确度偏低；精确度高时，召回率偏低。学习器把最可能是正例的样本排在前面。按此排序，把样本作为正例进行预测，根据 PR 绘图，如下。如果一个学习器绘制的 PR 曲线的面积（曲线与坐标轴围成）大于另一个学习器，则认为前一个学习器优于第二个。

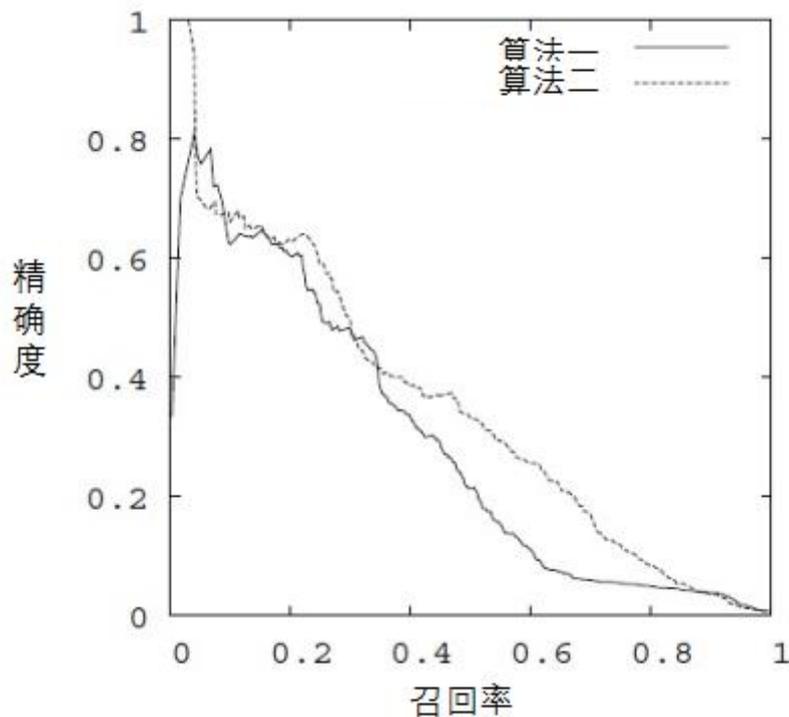


图 4-2 PR 曲线

在我们物体检测任务中有 20 个类别需要计算 20 个 PR 曲线图，计算其围成的面积 AP，然后 20 类求平均得到平均准确度 MAP。该指标作为我们物体检测任务最终评测标准，数值越大表示检测效果越好，模型性能越好。

4.2 实验结果展示

数据集和参数设置：为了验证本课题提出的自主学习框架，我们在公开的标

准数据集 Pascal VOC2007^[22]/2012^[23], COCO^[24]进行了相关实验, 引用 VOC 比赛的相关标准, 预测的坐标和 groundtruth 的 IOU 大于 0.5, 认为的正确的检测, 表现性能用平均精度 (mAP) 来表示。在所有的试验中, 我们设置参数 $\{C, t, \tau, \lambda^0, k, \alpha, T\} = \{0.001, 1000, 5, 0.0003, 100, 0.08, 5\}$

结果展示:

为了验证本课题提出的自主学习框架的有效性, 我们按照控制变量的方法进行了如下实验: 我们运用 R-FCN^[6] 的 ResNet-101^[26] 的网络, 在标准的 VOC2007^[22]/2012^[23], COCO^[24] 上进行, 通过主动学习和自步学习以尽可能少的标注样本进行实验。定义了如下几种实验方案: R-FCN, R-FCN+AL, R-FCN+SPL, R-FCN+AL+SPL, 通过控制主动学习和自步学习标注数据的数量, 进行对比试验。实验直接从预训练的 ImageNet^[20] 的 ResNet-101^[26] 的模型开始实验。本实验以 VOC2007^[22] 为基础数据集, VOC2012^[23] 和 COCO^[24] 数据集设为外来无标签的数据, 分别在测试集 VOC2007^[22] 和测试集 VOC2012^[23] 上进行相关实验。各实验结果如下。

表 4-2 自主学习框架实验结果汇总

测试集为 VOC2007 test，黑体标出为每个类别最高识别精度。”append”：AL 过程增加的 VOC2012 trainval 的数量占初始 VOC2007 trainval 的数量。”pseudo “：SPL 过程增加的 VOC2012 trainval 和 COCO trainval 伪标签的目标区域的数量占初始 VOC2007 trainval 的数量

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
Faster-RCNN	2007		66.9	74.5	78.3	69.2	53.2	36.6	77.3	78.2	82.0	40.7	72.7	67.9	79.6	79.2	73.0	69.0	30.1	65.4	70.2	75.8	65.8
RFCN	0%	0%	73.9	76.0	81.8	76.0	61.3	57.8	80.1	82.9	84.3	56.3	79.7	66.3	85.6	81.3	78.4	78.5	48.9	76.4	73.7	79.8	73.2
RFCN	200%	0%	79.8	82.0	83.8	82.0	71.4	67.9	85.6	87.3	88.3	66.6	85.0	73.0	87.5	87.9	82.9	83.2	57.1	81.6	79.3	85.3	77.8
RFCN+SPL	0%	270%	74.8	78.3	78.2	75.5	66.7	60.0	83.7	83.8	85.2	59.7	80.5	66.5	85.2	82.8	76.8	78.8	50.9	73.8	74.0	82.3	74.3
RFCN+AL	20%	0%	75.8	77.9	83.2	75.8	68.1	60.7	81.9	85.3	86.1	62.2	80.9	68.6	85.4	85.4	77.3	79.2	49.0	73.7	76.9	83.5	75.0
RFCN+AL	60%	0%	77.4	79.9	83.0	77.0	67.6	64.9	85.6	86.1	86.6	62.7	82.0	69.5	87.2	84.7	81.1	79.4	48.9	79.6	76.2	86.6	78.3
RFCN+AL	100%	0%	77.8	80.2	82.1	77.7	69.8	64.8	86.0	87.4	86.0	63.9	84.4	69.4	85.7	85.7	82.2	79.5	53.8	76.9	77.7	86.6	76.4
RFCN+AL+SPL	20%	340%	76.0	76.3	84.0	78.1	62.7	62.4	83.1	86.6	87.2	60.0	79.1	69.4	858.0	86.6	79.4	79.4	51.9	75.7	77.6	79.3	74.8
RFCN+AL+SPL	60%	400%	78.1	79.4	82.7	78.3	69.6	66.3	80.9	86.3	85.5	66.0	85.5	70.5	86.7	87.1	79.9	79.6	55.3	81.3	77.8	83.5	79.0
RFCN+AL+SPL	100%	410%	79.0	80.2	81.5	78.0	73.4	68.4	85.9	87.7	87.2	64.7	85.7	74.0	87.8	87.2	83.4	79.7	54.7	81.9	78.2	82.6	77.0
RFCN+AL+SPL	200%	420%	80.2	81.5	86.2	80.2	72.1	69.1	88.3	87.7	87.4	65.4	86.9	73.6	87.9	87.6	84.9	81.7	56.6	83.0	78.7	86.4	79.1
RFCN+AL+SPL	200%	1000%	81.8	87.0	87.1	88.5	74.7	76.7	87.9	89.2	88.9	69.8	88.4	72.4	88.9	88.1	85.0	86.5	51.6	85.8	77.7	88.0	75.7

表 4-3 自主学习框架实验结果汇总

测试集为 VOC2012 test，黑体标出为每个类别最高识别精度。” append”：AL 过程增加的 VOC2012 trainval 的数量占初始 VOC2007 trainval 和 VOC2007 test 的数量。” pseudo “：SPL 过程增加的 VOC2012 trainval 和 COCO trainval 伪标签的目标区域的数量占初始 VOC2007 trainval 和 VOC2007 test 的数量

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
RFCN	0%	0%	69.1	82.4	77.2	72.8	56.2	52.6	76.1	77.2	85.6	44.7	69.2	45.5	81.6	77.2	80.0	83.7	52.7	69.4	56.8	77.9	64.0
RFCN	100%	0%	77.4	87.4	84.6	78.4	65.5	60.8	81.2	83.3	89.9	61.7	80.4	62.2	87.1	85.1	85.7	85.8	61.2	80.1	67.9	85.2	74.7
RFCN+SPL	0%	105%	72.2	84.4	82.4	76.0	57.9	56.4	77.6	77.9	89.6	49.9	75.4	50.5	87.5	78.2	82.6	83.6	55.5	73.7	58.8	81.1	65.9
RFCN+AL	10%	0%	71.2	83.5	79.2	73.8	56.7	54.9	76.7	75.8	88.4	47.1	73.7	51.6	86.0	80.8	81.1	80.0	52.5	75.0	61.1	79.3	65.8
RFCN+AL	30%	0%	73.8	87.3	77.0	78.7	59.2	54.6	74.9	82.1	92.9	54.5	79.2	61.0	89.9	86.7	79.4	83.8	52.3	69.1	54.8	85.7	71.9
RFCN+AL	50%	0%	75.9	85.4	82.7	77.4	62.3	60.5	81.3	81.1	91.1	58.3	79.2	60.4	88.3	84.5	84.7	82.9	57.2	79.3	66.7	83.7	70.2
RFCN+AL+SPL	10%	130%	72.9	82.8	80.7	76.7	58.6	56.1	77.2	78.7	89.9	51.1	75.7	54.1	87.5	81.3	80.8	83.7	55.6	75.0	60.8	83.2	67.4
RFCN+AL+SPL	30%	175%	75.4	85.1	82.3	79.5	61.5	58.3	77.8	79.4	91.8	54.7	78.9	55.8	89.7	84.6	84.0	84.9	59.8	80.2	62.8	85.5	70.2
RFCN+AL+SPL	50%	190%	76.7	85.6	83.0	78.8	66.3	64.2	76.6	82.1	92.9	53.5	79.2	61.0	90.1	87.7	85.4	87.0	61.8	81.1	61.0	85.4	71.9
RFCN+AL+SPL	100%	200%	77.8	87.3	83.0	80.9	66.3	60.6	82.3	82.1	92.9	59.4	81.4	61.0	91.1	86.7	85.4	86.2	60.1	81.1	68.9	86.5	71.9
RFCN+AL+SPL	100%	500%	78.3	87.4	85.6	80.4	66.2	65.7	83.2	83.1	92.3	60.9	80.0	62.3	90.2	86.9	86.8	86.2	61.4	81.9	67.5	86.1	72.1

下面是我们通过自步学习过程给无标签的 COCO^[24] 数据标注伪标签的可视化图片，第一行红色标出的为高置信度的目标区域，第二行黄色标出的为低置信度的目标区域。主动学习过程人工给置信度低的样本标注标签的标注工具。



图 4-3 自步学习标注伪标签

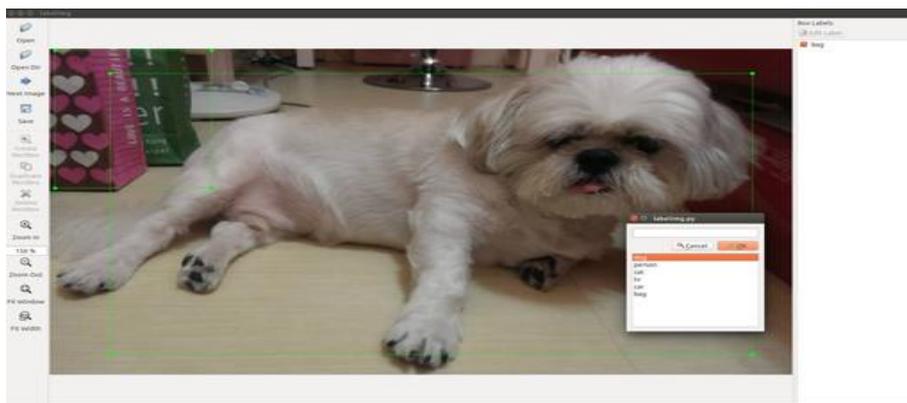


图 (a)

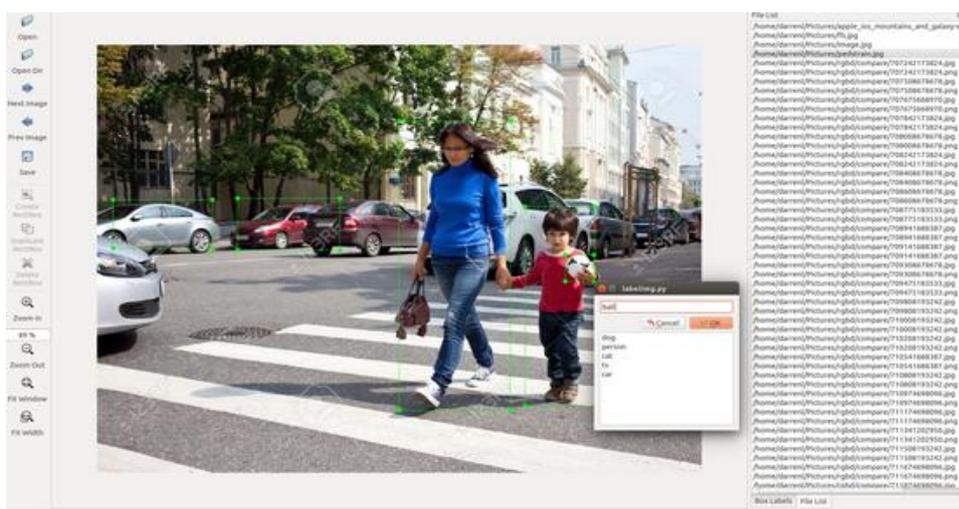


图 (b)

图 4-4 主动学习人工标注标签

4.3 实验分析

上以部分展示了我们通过控制变量方法的实验研究成果,表示在初始预训练的 ImageNet^[20]模型上,通过 R-FCN^[6]结合其他方法在 VOC2007 test 表 4-2 和 VOC2012 test 表 4-3 上得到的 mAP,并且给出了各种类别物体的识别精度。下面我们选择表 4-2 从各阶段对比分析我们的实验结果,表 4-3 的论证分析可类比。

4.3.1 网络模型的选择

表 4-4 模型选择阶段对比实验结果

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
Faster-RCNN	2007		66.9	74.5	78.3	69.2	53.2	36.6	77.3	78.2	82.0	40.7	72.7	67.9	79.6	79.2	73.0	69.0	30.1	65.4	70.2	75.8	65.8
RFCN	0%	0%	73.9	76.0	81.8	76.0	61.3	57.8	80.1	82.9	84.3	56.3	79.7	66.3	85.6	81.3	78.4	78.5	48.9	76.4	73.7	79.8	73.2

在网络的选择过程上,刚开始的时候选择了 Faster-RCNN^[5] 网络架构,因为其是基于 Region Proposal 方法的一个很好的通用流行架构,且在 VOC2007^[22] 数据集和 VGG19^[1] 网络结构上的 mAP 为 66.9%; 但因其训练时间长,测试速度慢,只能单显卡训练,训练模型只有基于 VGG19^[1] 的网络,没有基于最新的效果最好的 ResNet-101^[26] 的网络,模型精度较低,对于自步学习预测伪标签样本易出现误判,导致自步学习伪标注准确率不高。后来经过调研查阅,在 PASCAL VOC2012^[23] 的挑战赛上,第一名的队伍采用的网络结构是 R-FCN^[6],它也是基于 Faster-RCNN^[5] 的网络结构改编而来,采用了全卷积 ResNet-101^[26] 的网络结构,在 VOC2007^[22] 数据集上 mAP 为 73.9%,初始识别精度较高,且支持多显卡训练,大大提高了训练速度,节省了训练时间。

4.3.2 自步学习阶段

表 4-5 自步学习阶段对比实验结果

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
RFCN	0%	0%	73.9	76.0	81.8	76.0	61.3	57.8	80.1	82.9	84.3	56.3	79.7	66.3	85.6	81.3	78.4	78.5	48.9	76.4	73.7	79.8	73.2
RFCN+SPL	0%	270%	74.8	78.3	78.2	75.5	66.7	60.0	83.7	83.8	85.2	59.7	80.5	66.5	85.2	82.8	76.8	78.8	50.9	73.8	74.0	82.3	74.3
RFCN+AL	20%	0%	75.8	77.9	83.2	75.8	68.1	60.7	81.9	85.3	86.1	62.2	80.9	68.6	85.4	85.4	77.3	79.2	49.0	73.7	76.9	83.5	75.0
RFCN+AL+SPL	20%	340%	76.0	76.3	84.0	78.1	62.7	62.4	83.1	86.6	87.2	60.0	79.1	69.4	858.0	86.6	79.4	79.4	51.9	75.7	77.6	79.3	74.8
RFCN+AL	60%	0%	77.4	79.9	83.0	77.0	67.6	64.9	85.6	86.1	86.6	62.7	82.0	69.5	87.2	84.7	81.1	79.4	48.9	79.6	76.2	86.6	78.3
RFCN+AL+SPL	60%	400%	78.1	79.4	82.7	78.3	69.6	66.3	80.9	86.3	85.5	66.0	85.5	70.5	86.7	87.1	79.9	79.6	55.3	81.3	77.8	83.5	79.0
RFCN+AL	100%	0%	77.8	80.2	82.1	77.7	69.8	64.8	86.0	87.4	86.0	63.9	84.4	69.4	85.7	85.7	82.2	79.5	53.8	76.9	77.7	86.6	76.4
RFCN+AL+SPL	100%	410%	79.0	80.2	81.5	78.0	73.4	68.4	85.9	87.7	87.2	64.7	85.7	74.0	87.8	87.2	83.4	79.7	54.7	81.9	78.2	82.6	77.0

在自步学习阶段,我们通过控制自步学习标注伪标签样本的占比验证自步学习方法的有效性。在 RFCN 的初始 mAP 为 73.9%, RFCN+270%自步学习标注的样本的 mAP 为 74.8%,模型精度提高了 0.9%。在 RFCN+20%的主动学习的人工标注数据的 mAP 为 75.8%, RFCN+20%的主动学习的人工标注数据+340%自步学习标注的样本的 mAP 为 76.0%,模型精度提高了 0.2%。在 RFCN+60%的主动学习的人工标注数据的 mAP 为 77.4%, RFCN+60%的主动学习的人工标注数据+400%自步学习标注的样本的 mAP 为 78.1%,模型精度提高了 0.7%。在 RFCN+100%的主动学习的人工标注数据的 mAP 为 77.8%, RFCN+100%的主动学习的人工标注数据+410%自步学习标注的样本的 mAP 为 79.0%,模型精度提高了 1.2%。通过这些控制自步学习阶段的对比实验,我们可以发现在主动学习的人工标注数据相同情况下,增加自步学习标注的伪标签数据可以增加模型的识别精度,且随着伪标签数据数据量的增加,模型精度提高的百分比略有增加。当前模型预测无标签样本的类别和位置信息具有很高的准确性,可能有少量的类别错误信息或者物体位置偏移,随着预测的伪标签数据量的增加,错误预测的数量占比逐渐减少,所以在微调模型时,性能会有所提高,从而验证了自步学习方法的有效性。

4.3.3 主动学习阶段

表 4-6 主动学习阶段对比实验结果

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
RFCN	0%	0%	73.9	76.0	81.8	76.0	61.3	57.8	80.1	82.9	84.3	56.3	79.7	66.3	85.6	81.3	78.4	78.5	48.9	76.4	73.7	79.8	73.2
RFCN+AL	20%	0%	75.8	77.9	83.2	75.8	68.1	60.7	81.9	85.3	86.1	62.2	80.9	68.6	85.4	85.4	77.3	79.2	49.0	73.7	76.9	83.5	75.0
RFCN+AL	60%	0%	77.4	79.9	83.0	77.0	67.6	64.9	85.6	86.1	86.6	62.7	82.0	69.5	87.2	84.7	81.1	79.4	48.9	79.6	76.2	86.6	78.3
RFCN+AL	100%	0%	77.8	80.2	82.1	77.7	69.8	64.8	86.0	87.4	86.0	63.9	84.4	69.4	85.7	85.7	82.2	79.5	53.8	76.9	77.7	86.6	76.4

在主动学习阶段，我们通过控制主动学习人工标注样本的占比验证主动学习方法的有效性。在 RFCN 的初始 mAP 为 73.9%基础上加入各种占比的主动学习标注的数据样本，加入 20%的数据 mAP 为 75.8%，提高了 1.9%。加入 60%的数据 mAP 为 77.4%，提高了 1.6%。加入 100%的数据样本 mAP 为 77.8%，提高了 0.4%。通过这些控制主动学习阶段的对比实验，我们可以发现随着主动学习样本数量的增加，模型的识别精度会逐步得到提高。且随着主动学习数据量的增加，模型识别精度提高的百分比逐渐下降，说明模型性能趋于饱和，逐步达到稳定，这些对比实验验证了主动学习方法的有效性

4.3.4 分类器更新阶段

表 4-7 分类器更新阶段对比实验结果

method	append	pseudo	mAP	aero	bike	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbike	person	plant	sheep	sofa	train	tv
RFCN	0%	0%	73.9	76.0	81.8	76.0	61.3	57.8	80.1	82.9	84.3	56.3	79.7	66.3	85.6	81.3	78.4	78.5	48.9	76.4	73.7	79.8	73.2
RFCN+AL+SPL	20%	340%	76.0	76.3	84.0	78.1	62.7	62.4	83.1	86.6	87.2	60.0	79.1	69.4	858.0	86.6	79.4	79.4	51.9	75.7	77.6	79.3	74.8
RFCN+AL+SPL	60%	400%	78.1	79.4	82.7	78.3	69.6	66.3	80.9	86.3	85.5	66.0	85.5	70.5	86.7	87.1	79.9	79.6	55.3	81.3	77.8	83.5	79.0
RFCN+AL+SPL	100%	410%	79.0	80.2	81.5	78.0	73.4	68.4	85.9	87.7	87.2	64.7	85.7	74.0	87.8	87.2	83.4	79.7	54.7	81.9	78.2	82.6	77.0
RFCN+AL+SPL	200%	420%	80.2	81.5	86.2	80.2	72.1	69.1	88.3	87.7	87.4	65.4	86.9	73.6	87.9	87.6	84.9	81.7	56.6	83.0	78.7	86.4	79.1
RFCN+AL+SPL	200%	1000%	81.8	87.0	87.1	88.5	74.7	76.7	87.9	89.2	88.9	69.8	88.4	72.4	88.9	88.1	85.0	86.5	51.6	85.8	77.7	88.0	75.7
RFCN	200%	0%	79.8	82.0	83.8	82.0	71.4	67.9	85.6	87.3	88.3	66.6	85.0	73.0	87.5	87.9	82.9	83.2	57.1	81.6	79.3	85.3	77.8

在分类器更新阶段通过控制控制自步学习阶段和主动学习阶段的占比与在全标签的对比来进一步验证自步学习和主动学习方法的有效性。初始 RFCN 的 mAP 为 73.9%，在 RFCN 加入全标签的 200% 的 VOC2012^[23] 的数据的情况下，mAP 为 79.8%；RFCN+20% 的主动学习数据和 340% 自步学习数据的 mAP 为 76.0%；RFCN+60% 的主动学习数据和 400% 自步学习数据的 mAP 为 78.1%；RFCN+100% 的主动学习数据和 410% 自步学习数据的 mAP 为 79.0%；RFCN+100% 的主动学习数据和 410% 自步学习数据的 mAP 为 80.2%；RFCN+200% 的主动学习数据和 1000% 自步学习数据的 mAP 为 81.8%。这些对比实验说明，通过结合主动学习和自步学习的方法可以大幅度提高模型的性能，通过增加 200% 的主动学习数据和 410% 的自步学习数据，模型精度以超过全标签精度 0.4%，再增加伪标签的自步学习数据至 1000%，模型精度达到了 81.8%，超出全标签 2.0%。通过少量的主动学习数据在迭代微调模型时，弥补自步学习过程中少量样本信息错误的情况，将模型往正确的决策边界牵引，使模型精度逐步得到提高。这些对比试验，验证了我们提出的自主学习框架的有效性，在半监督环境下通过少量的标注数据和大量的无标签数据，迭代训练微调，逐步提高模型性能，超过全标签下的识别精度。

通过综上所述各阶段对比分析，可以从多方面验证我们提出的自主学习框架的有效性和实用性。

4.4 本章小结

本章主要阐述了实验过程中使用的数据集 VOC2007^[22]，VOC2012^[23] 的介绍，相关的评测方法，如 IOU 的计算，PR 曲线；展示了相关的实验结果和分析了各阶段的控制变量方法对比试验，验证了我们提出的自主学习框架的有效性和实用性。

第 5 章 总结与展望

5.1 工作总结

在本课题中,我们提出了一个有效的用少量标签数据进行物体检测的自主学习框架,通过结合自步学习和主动学习过程,迭代训练逐步提高分类器的性能。该自主学习框架能够有效地用自步学习的方式选出置信度高的样本并自动的给其标注伪标签,用主动学习的方式选择置信度低的样本供用户人工标注。同时,我们根据自主学习框架的工作流程从模型优化的角度做了合理的理论分析和解释,并且在 VOC 数据集上通过多次对比试验验证了该框架的有效性,另外我们在 SSD^[8]网络模型也验证了我们提出的自主学习框架能够通过该机制提高模型的性能,证明其具有通用性。该自主学习框架基本达到我们提出的通过少量有标签的数据和大量无标签的数据来提高模型的稳定性和抗噪性的初衷。在大数据移动互联网时代,在处理海量的无标签的图像和视频数据时,我们提出了具有很强的通用性的解决方案,节省了大量人力物力。

5.2 研究展望

相较于先前的半监督深度学习训练方法,本课题提出的训练方式不依赖于任何先验知识。我们的训练方法既可以融合基于先验信息的半监督训练技巧,进一步提升模型性能;也可以结合现有的强标签半监督深度学习方法,深化到各种高级语义任务中。最后,我们的模型训练过程是一个动态的子集选取过程,并不排

斥新加入的数据。而更新数据能够保证我们系统训练出来的模型具有很好的适应性和时效性。

本课题的研究都是基于静态图片分析的，在后续的研究中，我们考虑海量无标签图像视频数据的自主学习，视频处理和分析是目前的计算机视觉的主流方向，因为视频信息更加丰富，实用，但其标注更加的困难，所以运用一种半监督的自主学习方式来挖掘视频信息是一个很前沿的研究方向，也具有很强的实用意义。

基于本课题提出的上述关键技术，拟开发一套基于终生学习为基础的大规模半监督深度学习训练系统。在该系统中，我们分阶段提取不同特征空间的数据特征并且构建不同数据特征间的局部联系用于训练网络。为了处理海量样本，我们提供分布式局部图的计算，将样本筛选分布到多台机器上进行处理；同时利用数据并行技术进行深度网络的模型更新。

参考文献

- [1] Simonyan K, Zisserman A. Very deep convolutional networks for large-scale image recognition[J]. arXiv preprint arXiv:1409.1556, 2014.
- [2] Girshick R, Donahue J, Darrell T, et al. Rich feature hierarchies for accurate object detection and semantic segmentation[C]//Proceedings of the IEEE conference on computer vision and pattern recognition. 2014: 580-587.
- [3] He K, Zhang X, Ren S, et al. Spatial pyramid pooling in deep convolutional networks for visual recognition[C]//European Conference on Computer Vision. Springer International Publishing, 2014: 346-361.
- [4] Girshick R. Fast r-cnn[C]//Proceedings of the IEEE International Conference on Computer Vision. 2015: 1440-1448.
- [5] Ren S, He K, Girshick R, et al. Faster r-cnn: Towards real-time object detection with region proposal networks[C]//Advances in neural information processing systems. 2015: 91-99.
- [6] Jifeng Dai, Yi Li, Kaiming He, Jian Sun. {R-FCN}: Object Detection via Region-based Fully Convolutional Networks[J]. arXiv preprint arXiv:1605.06409, 2016.
- [7] Redmon J, Divvala S, Girshick R, et al. You only look once: Unified, real-time object detection[C]//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016: 779-788.
- [8] Liu W, Anguelov D, Erhan D, et al. SSD: Single shot multibox detector[C]//European Conference on Computer Vision. Springer International Publishing, 2016: 21-37.
- [9] Brinker K. Incorporating diversity in active learning with support vector machines[C]//ICML. 2003, 3: 59-66.
- [10] 白龙飞. 基于支持向量机的主动学习方法研究[D]. 山西大学, 2012.
- [11] 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(8): 954-962.
- [12] Lewis D D, Gale W A. A sequential algorithm for training text classifiers[C]//Proceedings of the 17th annual international ACM SIGIR conference on Research and development in information retrieval. Springer-Verlag New York, Inc., 1994: 3-12.

- [13] Tong S, Koller D. Support vector machine active learning with applications to text classification[J]. Journal of machine learning research, 2001, 2(Nov): 45-66.
- [14] Jiang L, Meng D, Yu S I, et al. Self-paced learning with diversity[C]//Advances in Neural Information Processing Systems. 2014: 2078-2086.
- [15] Bengio Y, Louradour J, Collobert R, et al. Curriculum learning[C]//Proceedings of the 26th annual international conference on machine learning. ACM, 2009: 41-48.
- [16] Kumar M P, Packer B, Koller D. Self-paced learning for latent variable models[C]//Advances in Neural Information Processing Systems. 2010: 1189-1197.
- [17] Jiang L, Meng D, Zhao Q, et al. Self-Paced Curriculum Learning[C]//AAAI. 2015, 2(5.4): 6.
- [18] Zhao Q, Meng D, Jiang L, et al. Self-Paced Learning for Matrix Factorization[C]//AAAI. 2015: 3196-3202.
- [19] Jiang L, Meng D, Mitamura T, et al. Easy samples first: Self-paced reranking for zero-example multimedia search[C]//Proceedings of the 22nd ACM international conference on Multimedia. ACM, 2014: 547-556.
- [20] Deng J, Dong W, Socher R, et al. Imagenet: A large-scale hierarchical image database[C]//Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. IEEE, 2009: 248-255.
- [21] Dollár P, Wojek C, Schiele B, et al. Pedestrian detection: A benchmark[C]//Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. IEEE, 2009: 304-311.
- [22] Everingham M, Zisserman A, Williams C K I, et al. The PASCAL visual object classes challenge 2007 (VOC2007) results[J]. 2007.
- [23] Everingham M, Van Gool L, Williams C K I, et al. The pascal visual object classes challenge 2012 (voc2012) results (2012)[C]//URL <http://www.pascal-network.org/challenges/VOC/voc2011/workshop/index.html>. 2010.
- [24] Lin T Y, Maire M, Belongie S, et al. Microsoft coco: Common objects in context[C]//European Conference on Computer Vision. Springer International Publishing, 2014: 740-755.
- [25] Elhamifar E, Sapiro G, Yang A, et al. A convex optimization framework for active learning[C]//Proceedings of the IEEE International Conference on Computer Vision. 2013: 209-216.
- [26] He K, Zhang X, Ren S, et al. Deep residual learning for image recognition[C]//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. 2016: 770-778.

相关的科研成果目录

包括本科期间发表的与毕业论文相关的已发表论文或被鉴定的技术成果、发明专利等成果，应在成果目录中列出。

论文的成果已经投稿到 CCF-A 类顶级会议 ACM Multimedia

致 谢

本课题得以能够顺利完成，得益于实验室导师和师兄的指导和帮助。对于一个从自动化专业转到做计算机视觉方向的我来说，刚开始是什么都不懂的，所以在前期阶段一直是通过阅读大量文献和复现相关论文的实验的工作，慢慢入门这个方向。在此期间，实验室导师和师兄给予了耐心的指导和帮助，让我逐渐成长。

首先我要感谢我的导师，林惊教授，感谢导师的培养和信任，感谢导师悉心指导和严格要求，把我安排到实验室工作学习，提供了实验室良好的学习环境。在刚入门不知所措时，给我指出课题研究方向，并在我遇到问题时，耐心的给予建议和指导，导师对严谨细致、一丝不苟的作风深深的感染了我，让我能够以更负责的态度和更高效的方式进行工作。

其次我要感谢博士生王可泽师兄，感谢师兄一直关注着我的课题进展，经常和我一起讨论交流，在遇到问题时，及时提出建议和给予帮助。对于一个陌生的课题，经常会遇到问题而不知道怎么做，师兄总会在我有问题时给我提供解决的办法。同时也要感谢实验室的其他师兄师姐，经常在实验室组织各种形式的研讨会等学习活动，营造实验室良好的学习氛围，激励着我更加努力学习。

最后我要感谢我的家人和朋友支持和关怀，在课题进展不顺心情低落时，给予的鼓励和支持。

严肖朋

2017年5月