

基于元伪标签和光照不变特征的人脸反欺诈算法

冯浩宇¹, 胡永健², 王宇飞³, 刘琲贝², 余翔宇², 钟睿²

¹ 华南理工大学计算机科学与工程学院 广州 中国 510641

² 华南理工大学电子与信息学院 广州 中国 510641

³ 广东警官学院刑事技术系 广州 中国 510440

摘要 人脸反欺诈(Face anti-spoofing, FAS)在防止人脸识别系统遭受欺诈攻击方面起着至关重要的作用,得益于深度学习网络强大的特征提取能力,基于深度学习的FAS算法取得比基于传统手工特征算法更好的性能,成为近期的研究热点。尽管大多数基于深度学习的FAS算法能在库内达到很好的检测效果,但是跨库检测性能欠佳,主要原因是库内和库外数据往往在不同条件下采集,例如拍摄设备、环境光照和攻击呈现设备不同,导致库内和库外数据的分布不同,两者之间存在域位移。当训练数据的多样性不足时,容易在库内学习过程中过拟合,跨库泛化性能不好。尽管我们可以判断起因,然而在真实世界的应用过程中解决上述问题并不容易。一方面,人脸反欺诈模型难以收集所有场景下的有标签训练样本;另一方面,不同应用场景使得同一因素产生不同的影响,例如,不同场景的光照导致域位移,影响了分类模型对本质性欺诈纹理的提取。为此,本文将元伪标签引入人脸反欺诈任务,提出一种基于元伪标签的人脸反欺诈方法。主要贡献包括:第一,提出一种基于图像块的“教师生成伪标签,学生反馈”半监督学习框架,挖掘局部图像的高区分度特征,解决有标签样本不足的问题;第二,基于局部重力模式(Pattern of local gravitational force, PLGF),设计一种带有注意力模块的光照不变特征分支,抑制应用场景中最容易影响特征提取的光照因素;第三,将元学习与半监督学习框架相结合,优化教师生成伪标签的过程,提高算法的跨库检测能力。与现有流行算法相比,在三个公开的测试数据集(包括CASIA、Replay-Attack和MSU)上,所提出方法在库内测试和跨库测试下均有突出的表现,尤其是泛化性能得到显著提高。在样本数量中等时,在不同库中的半总错误率保持最低。

关键词 人脸反欺诈;元学习;半监督学习;光照不变特征;元伪标签;深度学习

中图分类号 TP309.1 DOI号 10.19363/J.cnki.cn10-1380/tn.2023.09.05

Face Anti-Spoofing Based on Meta-pseudo-label and Illumination-invariant Feature

FENG Haoyu¹, HU Yongjian², WANG Yufei³, LIU Beibei², YU Xiangyu², ZHONG Rui²

¹ School of Computer Science and Engineering, South China University of Technology, Guangzhou 510641, China

² School of Electronic and Information Engineering, South China University of Technology, Guangzhou 510641, China

³ Department of Criminal Science and Technology, Guangdong Police College, Guangzhou 510440, China

Abstract Face anti-spoofing (FAS) plays a vital role in preventing face recognition systems from presentation attacks. Benefitted from powerful capability of feature extraction of deep learning (DL) network, FAS algorithms based on deep learning are much superior to those ones based on traditional handcrafted features in detection performance and thus become a research hotspot. Although most DL-based FAS algorithms can achieve good performance for intra-database test, the performance decreases greatly for cross-database test. The main reason is that samples are often collected under different capturing settings for intra-database and cross-database, for example, different cameras, environment illuminations, presentation medium, and thus their distributions are different, which can lead to the domain shift problem. When the diversity of training samples is insufficient, a model trained with such samples can be easily overfitting for intra-database while not being able to perform well for cross-database. Although the reasons of poor generalization are clear, the solution to them cannot be easily achieved in real-world applications. On the one hand, it is difficult for the FAS model to collect labeled training samples for all scenarios; on the other hand, different application scenarios make the same factor behave differently, which affects the extraction of intrinsic spoofing textures. In this paper, we introduce meta-pseudo-label into the FAS task, and propose a FAS method based on meta-pseudo-label. There are three major contributions. First, we propose a semi-supervised learning framework of “teacher generates pseudo-label and student feedbacks” based on image patches, which extracts the highly discriminative features of local images to solve the problem of insufficient labeled samples. Second, based on the Pattern of Local Gravitational Force (PLGF), we design an illumination-invariant feature

通讯作者: 胡永健, 博士, 教授, Email: eeyjhu@scut.edu.cn。

本课题得到国家重点研发计划项目(No. 2019QY2202)、广州开发区国际合作项目(No. 2019GH16)和中新国际联合研究院项目(No. 206-A018001)资助。

收稿日期: 2022-01-30; 修改日期: 2022-04-08; 定稿日期: 2023-06-12

branch with an attention module to suppress the illumination influence on feature extraction in application scenarios. Third, the meta-learning is combined with the semi-supervised learning framework to optimize the process of generating pseudo-label by teacher and improve the generalization ability of the algorithm. Compared with state-of-the-art algorithms, the proposed method performs quite well in both intra-database test and cross-database test on three public datasets including CASIA, Replay-Attack, and MSU. Specifically, the performance in cross-database is greatly improved. It can achieve the lowest HTER values for middle-sized sample number.

Key words face anti-spoofing; meta-learning; semi-supervised learning; illumination-invariant feature; meta-pseudo-label; deep learning

1 引言

人脸欺诈攻击方式一般来说可分为四类: 照片攻击: 攻击者利用打印的照片或者显示屏人脸图像来欺诈认证系统; 视频重放攻击: 攻击者利用提前拍摄的被攻击者的视频; 人脸面具攻击: 攻击者戴着精心制作的被攻击者的人脸面具, 或者利用 3D 打印制作 3D 面具; 对抗样本攻击: 攻击者通过生成式对抗网络(Generative adversarial networks, GAN)生成特定的样本噪声去干扰人脸认证系统, 使其产生错误的定向身份验证。相比于人脸面具攻击和对抗样本攻击, 照片攻击和视频重放攻击成本更低廉且更易于实施, 是目前主要的人脸欺诈方式。

人脸反欺诈(Face anti-spoofing, FAS)方法主要分为基于传统手工特征的方法和基于深度学习的方法。大多数传统的 FAS 技术专注于设计手工特征。由于照片打印和视频重放攻击在二次成像过程中会造成图像质量失真, 留下一些欺诈痕迹, 例如模糊、摩尔纹和打印噪声, 有文献使用纹理描述子来表现这种差别^[1-4]。此外, 由于重拍过程造成的扰动, 欺诈样本通常相比于真人图像有更低的视觉质量, 也有文献提出通过评估输入图像质量的方法来检测欺诈攻击样本^[5-6]。除了分析单张图像, 还有一些工作基于视频帧序列, 提出利用时域信息提取动态特征的方法^[7-8]。相比于基于手工特征的传统方法, 深度神经网络具有更强的特征学习能力, 已经广泛应用在 FAS 方法。Yang 等人在文献[9]首先应用 CNN(Convolutional neural network, CNN)到 FAS 中, 其利用 VGG-Net^[10]作为特征提取器, 然后将模型的全连接层的输出用来训练支持向量机(Support Vector Machine, SVM)分类器。

之后, 更多基于深度学习的方法被提出。深度学习的方法根据技术动机可大致分为 7 类^[11]。第一类基于纹理、颜色、高低频图像特征等^[12-14], 典型文献[12]受传统 LBP 特征提取方法的启发, 提出一种中心差分卷积(Central difference convolution, CDC)运算符, 提取具有高区分度的亮度级语义信息和梯

度级的细节信息。第二类基于生物活体和伪造人脸材料特征等^[15-17], 例如体温、脉搏、五官运动等活体人脸的特征和反光、伪影、摩尔纹等伪造人脸的特征。典型文献[15]提出将时域的人脸远程光电容积脉搏波(Remote photoplethymography, rPPG)和空域的深度信息结合起来进行欺诈检测。rPPG 信号是人脸皮肤下血液流动变化产生的时序信号, 活体在一段时间内一定会产生生理特性, 比如每隔一段时间会产生眨眼或者轻微的摇头, 这样打印的照片或者屏幕重放的视频和真实人脸在 rPPG 信号上有较大的区分度。第三类基于软件的辅助信息, 例如, 估计深度图、运动信息等作为检测方法的辅助信息输入^[18-20]。典型文献[20]利用估计得到的深度图和像素位置信息构成 3D 点云(3D point cloud), 作为辅助信息。第四类基于硬件的多模态信息, 通过联合利用可见光图像、红外图像、深度图等多种模态的人脸信息, 挖掘更为丰富有效的分类特征^[21-23]。第五类利用半监督、自监督/无监督学习进行域调整^[24-26]。典型文献[24]提出了一个半监督学习的网络框架, 用有限的标签训练数据来应对人脸欺诈攻击。其在训练模型时渐进地采用带有可靠伪标签(Pseudo-label)的未标注数据来丰富训练数据的多样性, 并利用时域的一致性来保证选定图像的伪标签的可靠性。第六类利用解耦合/解纠缠学习分离欺诈痕迹^[27-29]。典型文献[27]提出分离“欺诈痕迹”(例如, 颜色畸变, 3D 面膜边缘, 摩尔纹)的思路, 即去除与欺诈无关的因素(如光照等)后反映活体人脸和欺诈人脸之间本质的区分特征。第七类源于信号异常检测思路, 将真脸看成一类数据, 而将所有的非真脸看成是异常数据^[30-32]。

总的来说, 基于深度学习网络的欺诈检测逐渐成为主流, 大多数算法能在库内达到很好的检测效果, 但是跨库检测性能欠佳, 主要原因是库内和库外的数据往往在不同条件下采集, 例如拍摄设备、环境光照和呈现设备不同, 导致库内和库外数据之间存在域位移^[11], 当训练数据的多样性不足时, 容易在库内学习过程中过拟合, 泛化性能不好。尽管我们可以判断起因, 但在真实世界的应用过程中解决上述问题并不容易。一方面, 人脸反欺诈模型难以收集

所有场景数量足够多的有标签训练样本; 另一方面, 应用场景的客观条件使得同一因素产生不同的影响, 例如, 不同场景的光照导致域位移, 影响了分类模型对本质性欺诈纹理(如, 不同场景中的同一材质)的提取。

针对上述两个方面的问题, 本文将元伪标签引入人脸反欺诈任务, 提出一种基于元伪标签的人脸反欺诈方法。主要贡献包括: (1)提出一种基于图像块的“教师生成伪标签, 学生反馈”半监督学习框架, 挖掘局部图像的高区分度特征, 解决有标签样本不足的问题; (2)基于局部重力模式(Pattern of local gravitational force, PLGF)^[33], 设计一种带有注意力模块的光照不变特征分支, 抑制应用场景中最容易影响特征提取的光照因素; (3)将元学习与半监督学习框架相结合, 优化教师生成伪标签的过程, 提高算法的跨库检测能力。实验结果表明, 与现有流行算法相比, 在有标签样本不充足情况下, 本文检测模型不仅在库内具有很高的准确率, 在跨库检测中, 半总错误率也有显著下降。

2 预备知识及相关工作介绍

本节介绍与 FAS 相关的基于教师学生方法的半监督学习、光照不变特征、元学习基本知识和相关工作, 以及与文本算法构建注意力机制相关的局部聚集描述子向量的知识。

2.1 基于教师学生方法的半监督学习

基于深度学习的 FAS 方法一般采用监督学习, 也就是用事先对训练集标注的标签来监督分类模型的训练, 让其拟合先验目标分布, 但这种方式限制了模型的自主学习能力; 同时, 获取“标记”也需耗费大量的人力和物力; 而且要采集所有摄像头在不同拍摄环境下的样本并获得所有的欺诈种类是不现实的。上述几个因素限制了 FAS 网络在实际使用中的性能。半监督学习可以让学习器不依赖外界交互, 自动地利用未标记样本^[34]。在半监督学习中一个替代手工标签的常见方案是引入额外的专家知识, 通过与外界的交互来将部分未标记样本转变为有标记样本, 例如训练一个相对重量级的教师模型来学习训练数据, 然后使用这个表现良好的教师模型来监督一个相对轻量的学生模型的学习。又因为这类方法模仿教师的教学过程, 也被称为教师-学生方法。文献[35]最早提出教师-学生方法, 利用教师的输出分数来监督学生的训练。研究者经常让教师网络比学生模型更大来保证教师模型的容量和性能, 但文献[36]表明教师模型不一定要比学生模型大, 即使教

师和学生有一样的网络, 仍能提高学生的学习能力。

文献[37]提出的伪标签是一种简单有效的半监督学习方法, 其利用有标签和无标签数据在分布上的连续性和一致性, 通过在损失函数中添加与无标签样本相关的正则项, 增强学习器对未知数据的泛化能力。针对伪标签的置信度问题, 目前也有一些研究, 如标签平滑^[38]、置信度偏差^[39]。文献[40]在教师-学生方法框架下引入了伪标签, 而文献[41]使用了图像增强及正则化方法, 增强了伪标签数据的丰富性, 使学生比教师学习得更好。

近期文献[42]提出了一种利用生成对抗网络训练高效深度神经网络的新框架 DAFL(Data-free learning), 通过训练生成器来逼近原始数据集, 相比于利用教师网络训练学生网络的蒸馏模型, 可在无训练数据的情况下压缩深度神经网络, 且性能无太大下降。文献[43]提出非对称分支, 以不同的方式提取局部和全局特征, 增强特征的多样性, 其次提出跨分支监督, 允许一个分支从另一个分支获得监督, 对不同知识进行迁移, 增强教师-学生网络之间的权重多样性, 也在一定程度消除了伪标签的噪声。

2.2 与光照不变特征相关的 FAS 方法

应用场景中背景的复杂性, 拍摄设备的多样性以及不同的光照条件是影响 FAS 算法泛化性能的 3 个重要因素。由于实验室仿真前面 2 个因素的影响成本高昂, 因此, 通过算法降低光照条件的影响成为最具价值的研究。

近期文献[44]认为数据库间的光照差异是导致算法泛化性差的重要因素, 并且真脸和假脸在深度信息和材质属性方面有明显差异, 故提出利用深度图、反射系数图和反射比图等人脸本质属性进行人脸欺诈检测。通过线性迭代的方法从单张 RGB 图像去估计图片中人脸的深度图, 借此得到人脸各像素表面的法向量, 再利用朗伯反射模型去估计出人脸的反射系数和反射比图。该算法简便, 但存在一个缺陷: 从单张图片去估计反射系数的方法鲁棒性不高, 且生成的深度图是类似于纹理特征的伪深度图, 与真实深度有差距。文献[45]将多尺度视网膜特征(Multi-scale retinex, MSR)作为光照不变特征引入到 FAS, 通过双流网络将 RGB 图像和 MSR 图中的特征进行融合, 相比于仅使用 RGB 图提取特征的方法, 提升了算法的鲁棒性。本质上, 该方法所用的 MSR 类似于对 RGB 图像进行高斯平滑滤波, 留下的是与光照无关的反射系数成分, 在一定程度上消除了光照成分, 但相比于 RGB 图像, 损失了较多的面部细节和纹理信息。

与文献[44]采用的雅可比线性迭代近似估计深

度图及文献[45]用高斯滤波器得到光照成分这两种方法获得的光照不变特征相比, 文献[33]借鉴万有引力的相关理论, 提出的 PLGF 特征既属于光照不变性特征, 又保留了较丰富的纹理细节, 适用于 FAS 场景。具体来说, PLGF 假设每个图像像素都类似于一个行星体, 它与其局部相邻像素有相互作用, 可用大小和方向的形式进行描述。对于输入图像的每个像素, 可计算局部重力的大小和方向。在描述符中以引力角和大小的形式, 结合了梯度方向和梯度大小的优点。幅度是图像感兴趣区域的量度, 幅度大通常代表梯度较高的局部特征, 例如边缘、线条、纹理。另一方面, 局部重力方向表示局部的几何结构信息。利用二进制描述符对它们编码, 可捕获嵌入在基于局部引力表示的幅度和方向分量中的判别信息。

2.3 与元学习相关的 FAS 方法

元学习是改善 FAS 泛化性能很有前景的一个研究方向。元学习希望模型获取一种“学会学习”(Learning to learn)的能力, 使其可在获取已有“知识”的基础上快速学习新的任务, 它的意图在于通过少量的训练实例设计能够快速学习新技能或适应新环境的模型。与传统的机器学习专注于通过一个任务的样本来学习一个任务不同, 元学习面向任务, 将每一个任务看作是元学习的样本, 然后从多个任务中学习一个学习的策略, 也叫做“元知识”(Meta-knowledge), 它可以是初始参数的估计^[46], 也可以是整个学习模型^[48]等。对于元知识的训练, 主要基于双层(bi-level)优化的思想。在元训练的过程中, 有两个优化过程: 基础学习(Base learning)过程, 也称作内层 (Inner)(或者叫低级/基级); 元学习 (Meta-learning)过程, 也称作外层(Outer)(或者叫高级/元级)。在内层中使用测试集更新模型, 然后在外层基于更新后的模型优化元知识。元学习的方法被广泛应用, 例如 Finn 等人^[46]提出模型无关元学习 (Model-agnostic meta-learning, MAML), 旨在归纳出表征能力最佳的网络模型参数。

在 FAS 领域, 为了更好地泛化到未知攻击, 文献[47]提出通过在元学习过程中寻找通用的学习方向的框架。该算法采用了细粒度的学习策略, 模型学习更灵活、泛化能力更强。不过, 该算法在一定程度上依赖于辅助信息的监督。文献[48]结合元学习和神经架构搜索(Neural architecture search, NAS), 提出了一个由中心差分卷积和池化算子组成的新型搜索空间, 利用有效的静态-动态呈现来充分挖掘时空差异, 并提出了利用跨域/类型知识进行鲁棒搜索的元 NAS。该算法充分利用了中心差分卷积对于提取细

粒度特征的作用, 有效结合了时空域信息, 能搜索出可泛化到不可见的数据域的网络架构。但该算法进行空间搜索时要求较大的算力, 设计网络时也需要较多的技巧来提升性能。文献[49]提出利用像素图来代替二元标签进行监督, 为检测器提供一个“元教师”来指导其学习。通过元教师指导的像素级标签将传统标签抽象成反映真人和攻击区别的本质属性标签, 有助于探索更大范围的欺诈痕迹, 可应对二维、三维的欺诈攻击方式, 但算法并没有对元教师的初始化作太多的探讨, 指导能力还有待提升。文献[30]将 FAS 问题转化为一个一类自适应问题, 提出一类自适应人脸欺诈检测 (One-class-adaptation FAS, OCA-FAS), 通过在 OCA 任务上训练一个元学习器, 以学习适应真脸(活体脸), 并提出一个元损失函数搜索策略, 用于搜索一个更好的损失函数, 帮助元学习器完成 OCA 任务。该算法充分利用新环境中的真脸, 使模型能快速适应新环境。文献[25]认为可将人脸欺诈检测定义为学习一个用于零样本或小样本 FAS 的元模型, 检测器应该首先学习有区分度的特征, 使其能从预先定义好的攻击类型泛化到未见过的欺诈种类, 并通过从预先定义好的攻击和几个新的攻击类型学习, 快速适应到新的欺诈类型。该方法在具体实现是通过训练元学习器来求解深度回归。

2.4 局部聚集描述子向量

局部聚集描述子向量(Vector of locally aggregated descriptors, VLAD)^[50]是一种把传统方法(例如尺度不变特征变换)获得的一张图像的若干局部特征压缩为一个特定大小的全局特征的方法, 它通过聚类实现特征降维。具体来说, N 个 D 维的局部特征 x_i 作为输入, K 个聚类中心 c_k 作为 VLAD 的参数, 输出 $K \times D$ 维的特征 V , V 的计算公式如下:

$$V(j, k) = \sum_{i=1}^N \alpha_k(x_i)(x_i(j) - c_k(j)),$$

$$i \in \{1, 2, \dots, N\}, j \in \{1, 2, \dots, D\}, k \in \{1, 2, \dots, K\} \quad (1)$$

其中 $x_i(j)$ 表示第 i 个局部特征第 j 维的值, $c_k(j)$ 表示第 k 个聚类中心第 j 维的值。 $\alpha_k(x_i)$ 是指示函数, 表示第 i 个局部特征是否属于第 k 类, 如果 x_i 与第 k 个聚类中心最近, 则 $\alpha_k(x_i)$ 的值为 1, 与其他聚类中心的值为 0。也就是说, $V(j, k)$ 计算的是属于某一类的所有局部特征与对应聚类中心的残差和。然后对特征 V 在 D 这一维上执行内归一化(Intra-normalization)操作, 即将每个中心点的特征分别做归一化, 再转换成向量, 最后执行 L2 范数归一化(L2-normalization)得到最终的全局特征向量。

该特征表达了聚类范围内局部特征的某种分布。对于高层次的聚类中心来说,其周围局部特征的分布具有语义层面的含义,而其中设计的残差 $x_i(j) - c_k(j)$ 抹去了不同聚类本身特征分布的差异,从而能够只考虑局部特征与聚类中心的不同所带来的特征分布,更能代表这一簇特征的独特性。

文献[51]尝试将 VLAD 嵌入到可训练的 CNN 架构中,把聚类中心作为网络的参数进行训练,使得聚类中心不再是狭义的聚类中心,而是更能体现特征分布的语义上的中心,设计了神经网络 NetVLAD 代替传统的 VLAD。具体来说,用神经网络卷积得到的特征作为局部特征 x_i ,为了让离散的函数指示函数 $a_k(x_i)$ 可导,把 $a_k(x_i)$ 看作对残差的加权,令函数 $z(x_i) = -\alpha \|x_i - c_k\|^2$ 为评估 x_i 与各个聚类的吻合程度的函数,然后用 $\text{softmax}(z(x_i))$ 作为权重,即表示特征离聚类中心越近, z 越大,权重越高;离聚类中心越远, z 越小,权重越低。用新的权重替代 $a_k(x_i)$,再经过一系列的化简得到:

$$V(j, k) = \sum_{i=1}^N \frac{e^{W_k^T x_i + b_k}}{\sum_{k'} e^{W_{k'}^T x_i + b_{k'}}} (x_i(j) - c_k(j)) \quad (2)$$

这样就实现了全流程可导,进而通过设计网络 NetVLAD 将传统方法神经网络化。

更进一步,文献[52]提出 NeXtVLAD 以解决

NetVLAD 参数量过多而导致的模型不易训练和容易过拟合的问题,从而在视频分类任务中把帧级别的特征降维成视频级别的特征进行分类。NeXtVLAD 的做法是先把高维的特征分解成一组低维的特征,分解思路来源于 ResNeXt,然后加入注意力机制,再进行特征的编解码,最终达到降维的效果。

3 算法介绍

本文算法整体流程包括网络训练和样本测试。

网络训练部分如图 1 所示,包括数据预处理得到 RGB 颜色通道图和 PLGF 图;RGB 颜色通道图划分出有标签样本、无标签样本、增强后的无标签样本后,送入教师学习模块得到教师半监督损失、无标签样本的伪标签、增强无标签损失,然后将无标签样本的伪标签送入学生元学习模块,更新学生模型参数,并得到学生元学习损失,与教师半监督损失、增强无标签损失构成教师损失,更新教师模型参数;将 PLGF 图送入光照不变特征提取网络,得到特征向量和分类向量,利用三元组损失和交叉熵损失监督训练后保存光照不变特征提取网络的模型和参数;利用验证集确定阈值。样本测试部分如图 2 所示,包括模型测试加载数据到学生模型和光照不变特征提取网络,得到相应的 RGB 分类分数和 PLGF 分类分数,加权求和得到分类分数,根据阈值判决分类结果。下面对关键环节进行描述。

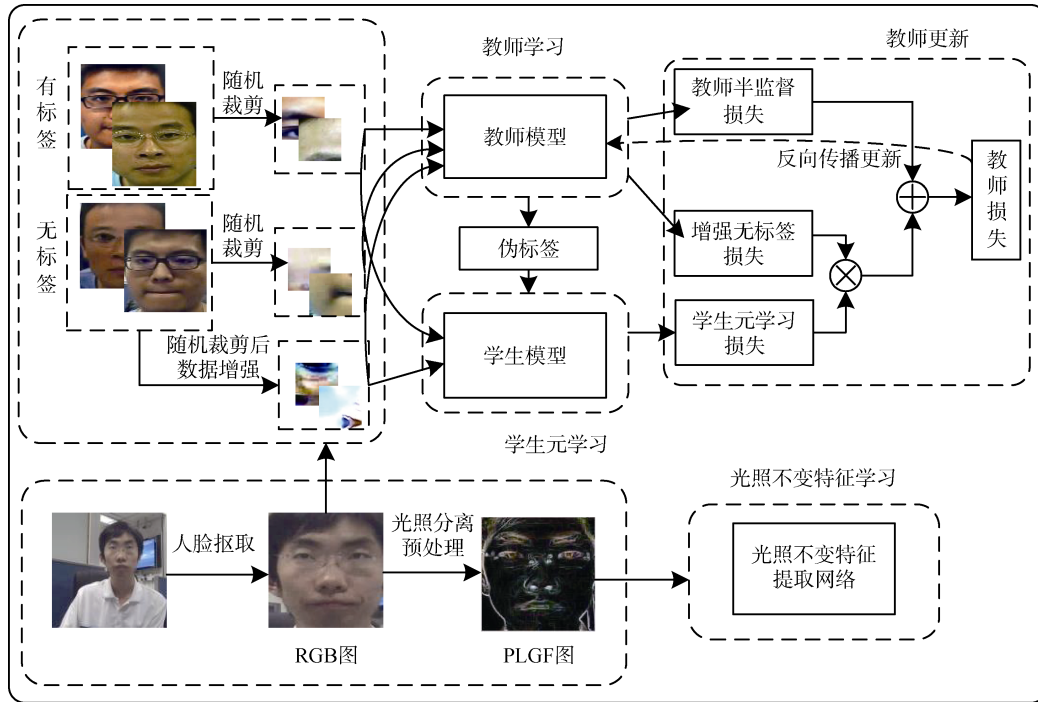


图 1 基于元伪标签和光照不变特征的人脸反欺诈算法训练框架

Figure 1 Face Anti-Spoofing algorithm training framework based on meta-pseudo-label and illumination-invariant feature

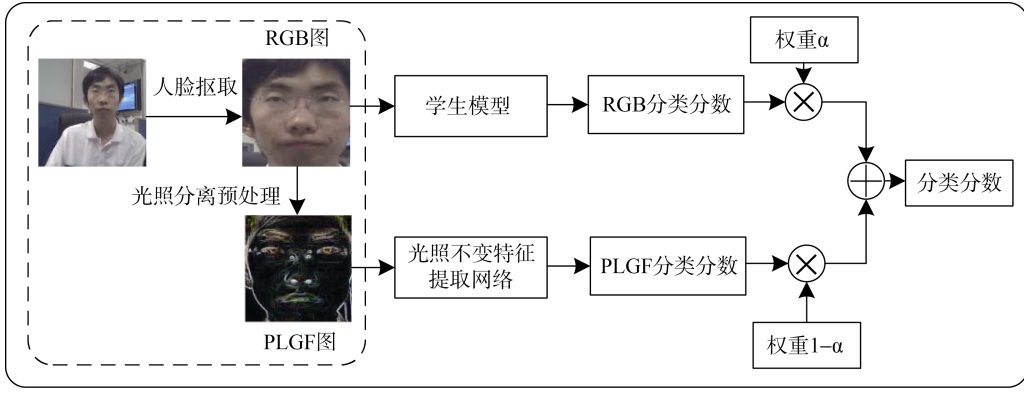


图 2 基于元伪标签和光照不变特征的人脸反欺诈算法测试框架

Figure 2 Face Anti-Spoofing algorithm testing framework based on meta-pseudo-label and illumination-invariant feature

在详述算法之前, 先对本文的函数及符号表示做如下规定:

让 $CE(q, p)$ 表示真实分布 q 和预测分布 p 的交叉熵损失; 如果 p 是一个标签值则先经过 One-Hot 编码变为标签向量; 如果 q 和 p 是一个批次, 则 $CE(q, p)$ 表示批次内各个样本的交叉熵损失的均值, 为便于阐述, 若不作特殊说明, 后文针对的是批次内一个样本展开讨论。设 k 是标签类别数, 以一个样本为例, 交叉熵损失如式(3)所示。

$$CE(q, p) = -\sum_{i=0}^k q_i \log p_i, i \in \{0, 1, \dots, k\} \quad (3)$$

让 $\arg \max(v)$ 表示取出向量 v 中最大值所在的索引, $\max(v)$ 表示取出向量 v 中最大值。例如一个 CNN 作为二元分类器, 输出分类向量 $p = [p_0, p_1]$, 其中 p_0 和 p_1 分别为样本的假体和真人的分类分数,

如果 $p_0 < p_1$, 则 $\arg \max(p) = 1$, $\max(p) = p_1$ 。

规定网络输出结果的变量名的格式统一为“网络结构名(输入变量; 网络参数)”, 例如, 送入有标签数据 $\{x_i, y_i\}$ 到参数为 θ_T 的教师模型 T 的预测结果表示为 $T(x_i; \theta_T)$ 。规定损失变量名的格式统一为“ \mathcal{L} (下标: 损失名字)(自变量: 模型参数)”, 例如, 教师有标签损失为 $\mathcal{L}_T(\theta_T)$ 。

3.1 基于伪标签构建半监督学习的 RGB 分支

3.1.1 基于伪标签构建教师-学生学习框架

作为 RGB 分支的主体, 基于伪标签构建教师-学生学习框架, 如图 3 所示, 由三个阶段组成: 1) 数据预处理: 采用裁剪的图像块作为 RGB 分支的训练输入; 2) 教师学习: 在有标签图像上训练教师模型; 3) 学生学习: 使用教师模型在无标签图像上来生成伪标签, 在有标签和带伪标签的无标签图像上训练学生模型。

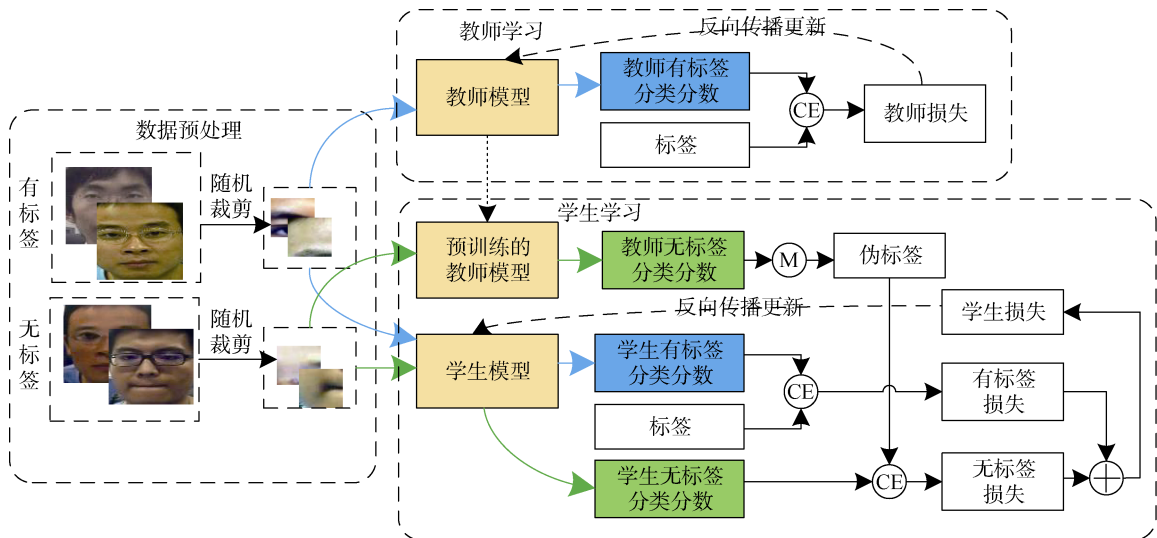


图 3 基于伪标签的教师-学生学习框架

Figure 3 Teacher-student studying method framework based on pseudo-label

第一阶段是数据预处理。在 RGB 分支训练阶段, 数据按比例划分为有标签数据和无标签数据, 并将 RGB 图像随机裁剪出的图像块作为输入。一方面, 因为目前公开数据集的视频数量有限, 同一段视频帧的信息存在较大关联性, 而图像块有助于增加训练数据多样性。另一方面, 图像块便于网络学习细粒度的特征, 从而可更容易适应分布差异大的库外样本。这样在测试阶段时采用整张图像作为输入时, 网络能侧重对一般性的分布无关的欺诈信息的提取。

第二阶段是教师学习。给定若干个批次的有标签数据集 $D\{(x_l, y_l)\}$, 送入教师模型进行预训练, 输出预测结果, 与真实标签 y_l 计算交叉熵得到教师有标签损失 $\mathcal{L}_l(\theta_T)$, 如式(4)所示:

$$\mathcal{L}_l(\theta_T) = CE(y_l, T(x_l; \theta_T)) \quad (4)$$

然后以最小化教师有标签损失为目标, 更新教师模型参数。

第三阶段是学生学习。给定无标签数据集 $D\{x_u\}$ 和有标签数据集 $D\{(x_l, y_l)\}$, 每次迭代分别从中采样一个批次的无标签数据 $x_u \sim D\{x_u\}$ 及有标签数据 $(x_l, y_l) \sim D\{(x_l, y_l)\}$ 。先将 x_u 送入预训练好的参数固定为 θ_T 的教师模型 T 生成伪标签 $T(x_u; \theta_T)$, 然后将 x_u 及 (x_l, y_l) 送进参数为 θ_S 的学生模型 S 分别得到无标签预测结果 $S(x_u; \theta_S)$ 及有标签预测结果 $S(x_l; \theta_S)$ 。接着计算学生模型在无标签数据的预测结果和伪标签的交叉熵, 得到学生无标签损失 $\mathcal{L}_u(\theta_T, \theta_S)$, 如式(5)所示:

$$\mathcal{L}_u(\theta_T, \theta_S) = CE(T(x_u; \theta_T), S(x_u; \theta_S)) \quad (5)$$

再计算有标签预测结果与真实标签 y_l 的交叉熵, 得到学生有标签损失 $\mathcal{L}_l(\theta_S)$, 如式(6)所示:

$$\mathcal{L}_l(\theta_S) = CE(y_l, S(x_l; \theta_S)) \quad (6)$$

学生无标签损失与学生有标签损失求和得到学生损失。这一阶段的优化目标为最小化学生损失, 更新学生模型参数, 最终得到优化后的学生模型。

3.1.2 使用无监督数据增强构建半监督损失作为教师模型的辅助损失

为了提高教师的指导能力, 对教师学习过程采用无监督数据增强(Unsupervised Data Augmentation, UDA)^[41]的方式, 将其得到的半监督损失作为教师模型的辅助损失。UDA 采用了具针对性的数据增强, 以产生多样化和真实的扰动, 相当于模拟不同光照、拍摄设备的影响, 有助于提高数据多样性。另一方面, UDA 通过联合优化标注数据的监督损失和未标注数据的无监督损失来计算最终的损失, 可将标注信息

从已标注数据传播到未标注数据。

具体来说, 给定一个参数为 θ_T 的教师模型 T , 对于一个批次的有标签数据 (x_l, y_l) , 经过教师模型 T , 由 3.1.1 的教师学习阶段得到了有标签损失 $\mathcal{L}_l(\theta_T)$ 。对于无标签数据 x_u , 先做一次数据增强得到 \tilde{x}_u , 然后将增强前的数据 x_u 和增强后的数据 \tilde{x}_u 分别送进教师模型 T , 得到相应的预测结果 $T(x_u; \theta_T)$ 和 $T(\tilde{x}_u; \theta_T)$, 将这两个结果计算得到的交叉熵损失作为无标签损失 $\mathcal{L}_u(\theta_T)$, 如式(7)所示:

$$\mathcal{L}_u(\theta_T) = CE(T(x_u; \theta_T), T(\tilde{x}_u; \theta_T)) \quad (7)$$

最后将有标签损失 $\mathcal{L}_l(\theta_T)$ 和无标签损失 $\mathcal{L}_u(\theta_T)$ 相加得到半监督损失 \mathcal{L}_{ssl} 作为教师损失, 进而反向传播更新教师模型参数。

此外, 本算法的伪标签为 0/1 的硬标签, 即将前述教师模型对无标签数据得到的预测结果 $T(x_u; \theta_T)$ 取出最大值所属类别作为伪标签 \hat{y}_u 。受到文献[38]利用标签平滑保证置信度的启发, 我们对伪标签 \hat{y}_u 进行平滑, 如式(8)表示:

$$S = Y(1 - \alpha) + \alpha / 2 \quad (8)$$

其中, α 表示平滑系数, Y 表示 \hat{y}_u 经过 One-Hot 编码的伪标签向量, S 表示平滑后的伪标签向量。

更进一步, 还设立了置信度机制, 其得到的置信度分数与无标签损失 $\mathcal{L}_u(\theta_T)$ 相乘。置信度机制具体为: 如果教师无标签预测结果 $T(x_u; \theta_T)$ 最大值大于预设的阈值 t 才认为是可信的, 输出置信度分数 C 为 1; 否则为 0。判决公式如下:

$$C = \begin{cases} 0, \max(T(x_u; \theta_T)) < t \\ 1, \max(T(x_u; \theta_T)) \geq t \end{cases} \quad (9)$$

教师模型和学生模型采用相同的网络结构 WideResNet^[53]。一方面, 文献[36]表明, 软标签形式的伪标签的熵比硬标签的熵大, 包含了更丰富信息; 另一方面, UDA 通过提升数据多样性能让教师模型学习到更多信息。因而, 本文采用和学生模型一样大的教师模型, 在保证教师模型指导能力的情况下, 降低训练时对内存的需求。如图 4 所示, WideResNet 由堆叠的残差卷积块组成, 通道逐渐加宽; 输入大小逐层降低到原始尺寸的 1/4, 然后经过全局平均池化层(Global average pooling, GAP)压缩为 1×1 的尺寸大小, 最后经过全连接层得到分类向量。其中, 由于通道数多, 采用 GAP 在减少输出维度的同时也能保留不同通道的特征。

综合 3.1.1 小节和 3.1.2 小节, 基于伪标签和 UDA 的 RGB 分支的算法伪代码如表 1 所示。

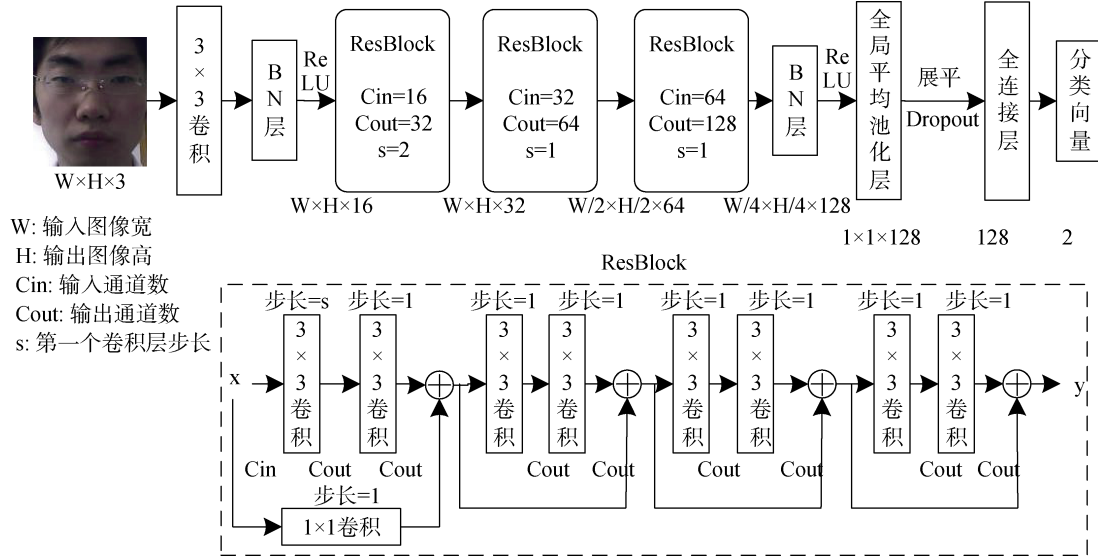


图 4 WideResNet 结构图
Figure 4 Structure of WideResNet

表 1 基于伪标签和 UDA 的 RGB 分支算法伪代码

Table 1 Algorithm's pseudo codes of RGB branch based on pseudo-label and UDA

算法 1 基于伪标签和 UDA 的 RGB 分支的算法过程

输入: 无标签数据集 $D\{x_u\}$ 和有标签数据集 $D\{(x_l, y_l)\}$, 随机数据增强方法 $\zeta(\cdot)$

输入: 教师模型 T 和学生模型 S , 教师模型和学生模型的学习率 η_T 和 η_S , 教师学习和学生学习的代数 N_T 和 N_S , 无监督系数 λ , 平滑度 α , 置信度阈值 t , 批次大小 s

- (1) 无标签数据集 $D\{x_u\}$ 的每个样本 x_u 进行一次随机数据增强得到增强版本 $\hat{x}_u = \zeta(x_u)$, 组成配对的无标签数据集 $D\{(x_u, \hat{x}_u)\}$
- (2) 随机初始化学生模型 S 的参数 θ_S 和教师模型 T 的参数 θ_T
- (3) For $n=1$ to N_T do:
- (4) 采样一个批次的 $(x_u, \hat{x}_u) \sim D\{(x_u, \hat{x}_u)\}$ 和 $(x_l, y_l) \sim D\{(x_l, y_l)\}$
- (5) 将 (x_l, y_l) 送入参数为 θ_T 的教师模型 T 得到教师有标签预测结果 $T(x_l; \theta_T)$
- (6) 计算教师有标签损失: $\mathcal{L}_l(\theta_T) = \text{CE}(y_l, T(x_l; \theta_T))$
- (7) 将 x_u 送入参数为 θ_T 的教师模型 T 得到教师无标签预测结果 $T(x_u; \theta_T)$
- (8) 将 \hat{x}_u 送入参数为 θ_T 的教师模型 T 得到教师增强无标签预测结果 $T(\hat{x}_u; \theta_T)$
- (9) $T(x_u; \theta_T)$ 取出最大值索引并经过平滑得到平滑无标签预测结果 $\bar{T}(x_u; \theta_T)$:
 $\bar{T}(x_u; \theta_T) \leftarrow \text{OneHot}(y_h) \cdot (1 - \alpha) + \alpha / 2, s.t. y_h = \arg \max(T(x_u; \theta_T))$
- (10) 对 $T(x_u; \theta_T)$ 的每个样本 $T(x_u; \theta_T)^{(i)}$ 计算置信度分数得到置信度向量 \mathbf{C} :
 $\mathbf{C} = [C_1, C_2, \dots, C_s]$, 其中 $C_i = 0$ 如果 $\max(T(x_u; \theta_T)^{(i)}) < t$ 否则 $1, i = \{1, 2, \dots, s\}$
- (11) 计算教师无标签损失: $\mathcal{L}_u(\theta_T) = \text{CE}(\bar{T}(x_u; \theta_T), T(\hat{x}_u; \theta_T)) \cdot \mathbf{C}$
- (12) 计算教师半监督损失: $\mathcal{L}_{ssl} = \mathcal{L}_l(\theta_T) + \lambda \min\{1, n / N_T\} \cdot \mathcal{L}_u(\theta_T)$
- (13) 用梯度下降法更新教师模型: $\theta_T \leftarrow \theta_T - \eta_T \nabla_{\theta_T} \mathcal{L}_{ssl}$
- (14) end For
- (15) For $n=1$ to N_S do:
- (16) 采样一个批次的 $(x_u, \hat{x}_u) \sim D\{(x_u, \hat{x}_u)\}$ 和 $(x_l, y_l) \sim D\{(x_l, y_l)\}$
- (17) 将 (x_l, y_l) 送入参数为 θ_S 的学生模型 S 得到学生有标签预测结果 $S(x_l; \theta_S)$

续表

算法 1 基于伪标签和 UDA 的 RGB 分支的算法过程

(18)	计算学生有标签损失: $\mathcal{L}_t(\theta_s) = \text{CE}(y_i, S(x_i; \theta_s))$
(19)	将 x_u 送入参数为 θ_T 的教师模型 T 得到教师无标签预测结果 $T(x_u; \theta_T)$, 取最大值索引并经过平滑得到伪标签 \hat{y}_u : $\hat{y}_u \leftarrow \text{OneHot}(y_{uh}) \cdot (1 - \alpha) + \alpha / 2, s.t. y_{uh} = \arg \max(T(x_u; \theta_T))$
(20)	将 x_u 送入参数为 θ_s 的学生模型 S 得到预测结果 $S(x_u; \theta_s)$
(21)	计算学生无标签损失: $\mathcal{L}_u(\theta_T, \theta_s) = \text{CE}(\hat{y}_u, S(x_u; \theta_s))$
(22)	计算学生损失: $\mathcal{L}_s = \mathcal{L}_t(\theta_s) + \mathcal{L}_u(\theta_T, \theta_s)$
(23)	用梯度下降法更新学生模型参数: $\theta_s \leftarrow \theta_s - \eta_s \nabla_{\theta_s} \mathcal{L}_s$
(24)	end For

输出: 参数为 θ_s 的学生模型 S 和参数为 θ_T 的教师模型 T

3.2 基于 PLGF 设计光照不变特征提取网络, 构建 PLGF 分支

3.2.1 进行光照分离处理得到 PLGF 图作为 PLGF 分支的输入

受到文献[33]提出的具有光照不变性特征 PLGF 的启发, 本文设计光照分离处理, 从 RGB 颜色通道数据获取 PLGF 图。PLGF 图包含了仅与反射系数相关的材质特征以及光照成分中颜色损失引入的欺诈噪声等区分性线索, 减少了不同光照环境对模型性能的影响, 再经过带有注意力模块的光照不变特征提取网络对本质特征的提取, 能够提高活体检测模型的鲁棒性。

具体来说, 首先对三个颜色通道的人脸特征在水平方向和垂直方向分别与 PLGF 算子进行 PLGF 卷积得到水平梯度 G_{hor} 和垂直梯度 G_{ver} 。PLGF 卷积具体表达如下列式子所示:

$$G_d[x, y] = \sum_{u=-1}^1 \sum_{v=-1}^1 f_d[u, v] I[x-u, y-v], d \in \{hor, ver\},$$

$$f_{hor} = \begin{bmatrix} -1 & 0 & 1 \\ -\sqrt{2} & 0 & \sqrt{2} \\ -1 & 0 & 1 \end{bmatrix}, f_{ver} = \begin{bmatrix} -1 & -\sqrt{2} & -1 \\ 0 & 0 & 0 \\ 1 & \sqrt{2} & 1 \end{bmatrix} \quad (10)$$

其中 f_{hor} 和 f_{ver} 分别为文献[33]中 PLGF 的水平方向和垂直方向的 3×3 卷积核。 $I[x, y]$ 为坐标 (x, y) 的像素值, $G_d[x, y]$ 为坐标 (x, y) 的水平/垂直方向梯度。

然后根据朗伯模型, 对水平方向和垂直方向的梯度进行光照分离得到水平光照分离梯度和垂直光照分离梯度 ISG_{hor} 和 ISG_{ver} 。为防止分母为 0, 将表达式分母加上极小值 ε 。由于在很小的区域内光照强度变化缓慢, 故可认为是恒值 L 。消除光照分量 L 的影响则可得仅与反射系数相关的人脸材质性特征。光照分离具体表达如下列式子所示:

$$ISG_d[x, y] = \frac{G_d[x, y]}{I[x, y] + \varepsilon} = \frac{\sum_{u=-1}^1 \sum_{v=-1}^1 f_d[u, v] I[x-u, y-v]}{I[x, y] + \varepsilon} =$$

$$\frac{\sum_{u=-1}^1 \sum_{v=-1}^1 f_d[u, v] R[x-u, y-v] L[x-u, y-v]}{R[x, y] L[x, y] + \varepsilon},$$

$$d \in \{hor, ver\} \quad (11)$$

其中, $I[x, y]$ 为坐标 (x, y) 的像素值, $R[x, y]$ 为该坐标像素的反射系数, $L[x, y]$ 为该坐标像素成像的光照强度。

接着对水平方向和垂直方向的光照分离梯度进行线性激活操作得到合成梯度 ISG , 组成 PLGF 图, 如下式所示:

$$ISG = \arctan(\sqrt{(ISG_{hor})^2 + (ISG_{ver})^2}) \quad (12)$$

得到的 PLGF 图如图 5 所示。

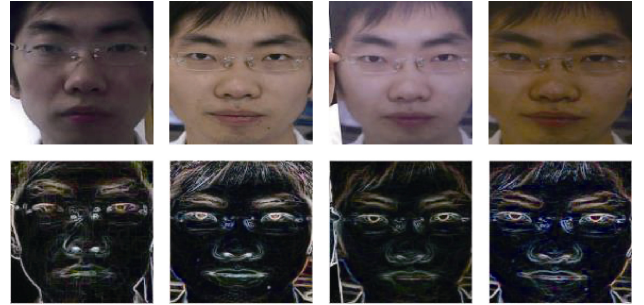


图 5 PLGF 示意图。第一行为原图, 第二行为 PLGF 图, 从左至右分别是场景 1 的真样本、场景 2 的真样本、场景 1 的假样本、场景 2 的假样本。

Figure 5 Illustration of PLGF. The first row and second row indicate original images and PLGF images, respectively, and from left to right, real samples in scene 1, real samples in scene 2, fake samples in scene 1, fake samples in scene 2

3.2.2 设计带注意力机制的光照不变特征提取网络

为提取 PLGF 图的光照不变特征, 如图 6 所示, 我们设计了带有注意力机制的光照不变特征网络 AttVLAD。首先, 用含跳接、丢弃(Dropout)、池化的 4 层卷积层搭建轻量的特征提取主干网, 得到多通

3.3 设计元标签优化教师模型生成的伪标签

3.1 节设计的教师-学生方法是使用预训练好的教师模型来指导学生, 专注于先让教师自身在训练集上取得良好的分类性能, 而没有对教师的指导能力作评估, 缺少根据学生学习的情况通过反馈来优化。Pham 等人^[55]认为尽管伪标签方法在多个任务中表现较好, 但有一个主要的缺陷: 如果伪标签是不准确的, 学生会学到不准确的数据, 因而, 学生可能并不比教师好, 这个缺陷也称为伪标签的置信偏差问题。其在传统的伪标签方法基础上引入了元(Meta)建模的过程, 提出了元伪标签(Meta-pseudo-label, MPL)。传统的基于伪标签的蒸馏方法是基于预训练好的教师模型, 利用教师模型提供的伪标签作为学生模型的目标进行训练, 而该方法可通过学生模型在有标签数据上的表现来帮助教师模型优化。具体来说, 利用来自学生的反馈来优化教师生成的伪标签, 借鉴元学习的双层更新方式, 学生从一个由教师标注的伪标签批次数据中学习并更新的过程作为内层, 以及教师从学生在有标签数据集上表现的反馈信号中学习作为外层。受此启发, 本文采用双层更新的元学习方式训练学生和教师, 可让教师根据学生学习的情况调整自身的教学。

如 3.1 小节的学生学习阶段所述, 给定无标签数据集 $D\{x_u\}$ 和有标签数据集 $D\{(x_l, y_l)\}$, 每次迭代分别从中采样一个批次的无标签数据 $x_u \sim D\{x_u\}$ 及有标签数据 $(x_l, y_l) \sim D\{(x_l, y_l)\}$ 。先将 x_u 送入预训练好的参数固定为 θ_T 的教师模型 T 生成的伪标签 $T(x_u; \theta_T)$, 然后将 x_u 送进参数为 θ_S 的学生模型 S 的得到无标签预测结果 $S(x_u; \theta_S)$ 。如式(13)所示, 通过最小化学生模型在无标签数据的预测结果和教师模型提供的伪标签的交叉熵损失, 可得到优化后的学生模型参数 θ_S^{PL} 。

$$\theta_S^{PL} = \arg \min_{\theta_S} \underbrace{CE(T(x_u; \theta_T), S(x_u; \theta_S))}_{:= \mathcal{L}_u(\theta_T, \theta_S)} \quad (13)$$

假定有一个固定参数为 θ_T 的教师模型能提供准确伪标签, 那么通过伪标签的监督, 优化后的参数为 θ_S^{PL} 的学生模型最终会在有标签数据集 $\{x_l, y_l\}$ 表现良好, 取得到较低的学生有标签损失 $\mathcal{L}_l(\theta_S^{PL})$, 如式(14)所示:

$$\mathcal{L}_l(\theta_S^{PL}) = CE(y_l, S(x_l; \theta_S^{PL})) \quad (14)$$

学生有标签损失 $\mathcal{L}_l(\theta_S^{PL})$ 越低, 间接地表明教师生成的伪标签越准确。由式(13)可知经过伪标签优化

后的学生参数 θ_S^{PL} 依赖于教师参数 θ_T , 则 θ_S^{PL} 可记为 $\theta_S^{PL}(\theta_T)$, 学生有标签损失亦即学生在有标签数据上的表现, 可重写为关于 θ_T 的函数 $\mathcal{L}_l(\theta_S^{PL}(\theta_T))$, 因而可通过学生有标签损失反向传播优化教师参数 θ_T 。这样伪标签的优化目标变成最小化学生有标签损失 $\mathcal{L}_l(\theta_S^{PL}(\theta_T))$, 结合式(13), 如式(15)所示:

$$\min_{\theta_T} \mathcal{L}_l(\theta_S^{PL}(\theta_T)) = \min_{\theta_T} \mathcal{L}_l(\arg \min_{\theta_S} \mathcal{L}_u(\theta_T, \theta_S)) \quad (15)$$

至此伪标签引入了元学习的框架, 建立了元伪标签的概念。本文算法进行的是单任务的元学习, 也就是每一个批次只包含一个任务, 一个任务由采样的一组无标签数据 $x_u \sim D\{x_u\}$ 及有标签数据 $(x_l, y_l) \sim D\{(x_l, y_l)\}$ 组成, 前者作为支持集, 后者作为查询集; 内层就是学生模型通过伪标签的监督在无标签数据上优化的过程, 外层就是将优化后的学生模型在有标签数据上的损失反馈给教师模型, 进而优化教师模型的过程。

下面讨论具体的优化算法。由于式(15)包含 $\arg \min$ 函数, 需要等到 θ_S 达到最优才能进行下一步, 没法用梯度方式来直接优化, 为此, 我们借鉴文献[46]的思想, 采用一步(One-step)梯度更新 θ_S 近似代替多步的 $\theta_S^{PL}(\theta_T)$:

$$\theta_S^{PL}(\theta_T) \approx \theta_S - \eta \nabla_{\theta_S} \mathcal{L}_u(\theta_T, \theta_S) \quad (16)$$

其中 η 是学习率。代入到式(15), 可得到元伪标签优化目标:

$$\min_{\theta_T} \mathcal{L}_l(\theta_S - \eta \nabla_{\theta_S} \mathcal{L}_u(\theta_T, \theta_S)) \quad (17)$$

借鉴元学习的双层优化方式, 设定学生模型学习率 η_S , 教师模型学习率 η_T , 二层优化过程如下:

第一层(优化学生): 首先将一个批次的无标签数据 x_u 送入到参数为 θ_T 的教师模型得到教师无标签预测结果 $T(x_u; \theta_T)$, 对其取出最大值所属类别及平滑, 伪标签 \hat{y}_u , 并计算教师无标签损失 $\mathcal{L}_u(\theta_T)$, 如式(18)所示:

$$\mathcal{L}_u(\theta_T) = CE(\hat{y}_u, T(x_u; \theta_T)) \quad (18)$$

然后再将 x_u 送入参数为 θ_S 的学生模型得到学生无标签预测结果, 与伪标签 \hat{y}_u 计算交叉熵, 得到学生无标签损失 $\mathcal{L}_u(\theta_T, \theta_S)$, 如式(19)所示:

$$\mathcal{L}_u(\theta_T, \theta_S) = CE(\hat{y}_u, S(x_u; \theta_S)) \quad (19)$$

接着用梯度下降法优化学生模型:

$$\theta_S' = \theta_S - \eta_S \nabla_{\theta_S} \mathcal{L}_u(\theta_T, \theta_S) \quad (20)$$

第二层(优化教师): 给定一个批次的有标签数据

(x_l, y_l) , 在文献[55]中, 将 (x_l, y_l) 送入优化后的参数为 θ_S' 的学生模型得到有标签损失 $\mathcal{L}_l(\theta_S')$ 作为教师元学习损失来更新教师模型, 为区分于后面的改进工作, 将其记为 $\mathcal{L}_{mpl}^{baseline}$, 如式(21)所示:

$$\mathcal{L}_{mpl}^{baseline} = \mathcal{L}_l(\theta_S') = CE(S(x_l; \theta_S'), y_l) \quad (21)$$

我们做的改进是将 (x_l, y_l) 分别送入优化前的参数为 θ_S 的学生模型和优化后的参数为 θ_S' 的学生模型, 得到旧的有标签损失 $\mathcal{L}_l(\theta_S)$ 和新的有标签损失 $\mathcal{L}_l(\theta_S')$, 然后二者作差得到学生元学习损失 \mathcal{L}_{up} , 为了利用教师的伪标签和预测分数之间的误差反馈, 学生元学习损失 \mathcal{L}_{up} 再与式(18)得到的教师无标签损失 $\mathcal{L}_u(\theta_T)$ 相乘, 得到改进的教师元学习损失 \mathcal{L}_{mpl} , 如

下式所示:

$$\mathcal{L}_l(\theta_S) = CE(y_l, S(x_l; \theta_S)) \quad (22)$$

$$\mathcal{L}_l(\theta_S') = CE(y_l, S(x_l; \theta_S')) \quad (23)$$

$$\mathcal{L}_{up} = \mathcal{L}_l(\theta_S) - \mathcal{L}_l(\theta_S') \quad (24)$$

$$\mathcal{L}_{mpl} = \mathcal{L}_{up} \cdot \mathcal{L}_l(\theta_T) \quad (25)$$

最后与 3.1.2 小节的教师半监督损失 \mathcal{L}_{ssl} 相加, 构成教师损失 \mathcal{L}_T , 用梯度下降法优化教师模型:

$$\theta_T' = \theta_T - \eta_T \nabla_{\theta_T} \mathcal{L}_T \quad (26)$$

这样通过学生模型更新前后的表现差异的反馈, 能让教师模型更准确了解学生模型对所学知识的掌握程度。

结合 3.1.2 小节的 UDA, 得到基于元伪标签和 UDA 的 RGB 分支的算法过程, 如表 3 所示。

表 3 基于元伪标签和 UDA 的 RGB 分支算法伪代码

Table 3 Algorithm's pseudo codes of RGB branch based on meta-pseudo-label and UDA

算法 3 基于元伪标签和 UDA 的 RGB 分支算法过程

输入: 无标签数据集 $D\{x_u\}$ 和有标签数据集 $D\{(x_l, y_l)\}$, 随机数据增强方法 $\zeta(\cdot)$	
输入: 教师模型 T 和学生模型 S , 教师模型和学生模型的学习率 η_T 和 η_S , 训练代数 N , 无监督系数 λ , 平滑度 α , 置信度阈值 t , 批次大小 s	
(1)	无标签数据集 $D\{x_u\}$ 的每个样本 x_u 进行一次随机数据增强得到增强版本 $\hat{x}_u = \zeta(x_u)$, 组成配对的无标签数据集 $D\{(x_u, \hat{x}_u)\}$
(2)	随机初始化学生模型 S 的参数 θ_S 和教师模型 T 的参数 θ_T
(3)	For $n=1$ to N do:
(4)	采样一个批次的 $(x_u, \hat{x}_u) \sim D\{(x_u, \hat{x}_u)\}$ 和 $(x_l, y_l) \sim D\{(x_l, y_l)\}$
(5)	将 (x_l, y_l) 送入参数为 θ_T 的教师模型 T 得到教师有标签预测结果 $T(x_l; \theta_T)$
(6)	计算教师有标签损失: $\mathcal{L}_l(\theta_T) = CE(y_l, T(x_l; \theta_T))$
(7)	将 x_u 送入参数为 θ_T 的教师模型 T 得到教师无标签预测结果 $T(x_u; \theta_T)$
(8)	将 \hat{x}_u 送入参数为 θ_T 的教师模型 T 得到教师增强无标签预测结果 $T(\hat{x}_u; \theta_T)$
(9)	$T(x_u; \theta_T)$ 取出最大值索引并经过平滑得到平滑无标签预测结果 $\bar{T}(x_u; \theta_T)$:
	$\bar{T}(x_u; \theta_T) \leftarrow \text{OneHot}(y_h) \cdot (1 - \alpha) + \alpha / 2, s.t. y_h = \arg \max(T(x_u; \theta_T))$
(10)	对 $T(x_u; \theta_T)$ 的每个样本 $T(x_u; \theta_T)^{(i)}$ 计算置信度分数得到置信度向量 C :
	$C = [C_1, C_2, \dots, C_s]$, 其中 $C_i = 0$ 如果 $\max(T(x_u; \theta_T)^{(i)}) < t$ 否则 $1, i = \{1, 2, \dots, s\}$
(11)	计算教师无标签损失: $\mathcal{L}_u(\theta_T) = CE(\bar{T}(x_u; \theta_T), T(\hat{x}_u; \theta_T)) \cdot C$
(12)	计算教师半监督损失: $\mathcal{L}_{ssl} = \mathcal{L}_l(\theta_T) + \lambda \min 1, \{n / N_T\} \cdot \mathcal{L}_u(\theta_T)$
(13)	计算教师增强无标签损失: $\mathcal{L}_{ua}(\theta_T) = CE(T(\hat{x}_u; \theta_T), \arg \max(T(x_u; \theta_T)))$
(14)	将 (x_l, y_l) 送入参数为 θ_S 的学生模型 S 得到学生有标签预测结果 $S(x_l; \theta_S)$
(15)	计算学生有标签损失: $\mathcal{L}_l(\theta_S) = CE(y_l, S(x_l; \theta_S))$
(16)	将 x_u 送入参数为 θ_T 的教师模型 T 得到教师无标签预测结果 $T(x_u; \theta_T)$, 取最大值索引并经过平滑得到伪标签 \hat{y}_u :
	$\hat{y}_u \leftarrow \text{OneHot}(y_{uh}) \cdot (1 - \alpha) + \alpha / 2, s.t. y_{uh} = \arg \max(T(x_u; \theta_T))$
(17)	将 x_u 送入参数为 θ_S 的学生模型 S 得到预测结果 $S(x_u; \theta_S)$
(18)	计算学生无标签损失: $\mathcal{L}_u(\theta_T, \theta_S) = CE(\hat{y}_u, S(x_u; \theta_S))$
(19)	用梯度下降法更新学生模型参数: $\theta_S' = \theta_S - \eta_S \nabla_{\theta_S} \mathcal{L}_S$
(20)	将 (x_l, y_l) 送入参数为 θ_S' 的学生模型 S 得到学生新的有标签预测结果 $S(x_l; \theta_S')$
(21)	计算学生新的有标签损失: $\mathcal{L}_l(\theta_S') = CE(y_l, S(x_l; \theta_S'))$

续表

算法 3 基于元伪标签和 UDA 的 RGB 分支算法过程

- (22) 计算学生元学习损失: $\mathcal{L}_{up} = \mathcal{L}_i(\theta_s) - \mathcal{L}_i(\theta'_s)$
- (23) 计算教师损失: $\mathcal{L}_T = \mathcal{L}_{ssl} + \mathcal{L}_{up} \cdot \mathcal{L}_{ua}(\theta_T)$
- (24) 用梯度下降法更新教师模型: $\theta_T \leftarrow \theta_T - \eta_T \nabla_{\theta_T} \mathcal{L}_T$
- (25) end For
- 输出: 参数为 θ_s 的学生模型 S 和参数为 θ_T 的教师模型 T

4 实验设置

4.1 人脸欺诈数据库

为验证算法的有效性, 本文用三个流行的公开数据集进行实验, 具体如表 4 所示。

表 4 实验中用到的人脸欺诈数据库

Table 4 Database for FAS used in the experiment

数据库	CASIA ^[54]	Replay Attack ^[3]	MSU ^[11]
年份来源	2012-中国 CBSR	2012-瑞士 Idiap	2015-美国 MSU
拍摄人数	50	50	35
视频数量	600	1200	280
拍摄设备	3	2	2
攻击类型	打印攻击/ 重放攻击	打印攻击/ 重放攻击	打印攻击/ 重放攻击
拍摄场景	3	2	1

4.2 实验环境及参数设置

4.2.1 实验环境

算法在 Linux 系统上进行, 主要基于深度学习框架 Pytorch1.6.1 来实现, 所用显卡为 GTX1080Ti, CUDA 版本为 10.1.105, cudnn 版本 7.6.4。

4.2.2 实验的训练参数

训练样本总数为 16000, 其中有标签样本数量为 4000, 真实样本和攻击样本各 2000 张, 批次大小为 32, 无监督系数 λ 为 8。输入大小为 $256 \times 256 \times 3$ 的原图, 分为两个支路: 原图经过随机裁剪得到 $64 \times 64 \times 3$ 的图像块, 并采用 UDA 的方式, 送入教师模型和学生模型构成的 RGB 分支进行训练; 原图经过光照分离卷积得到 PLGF 图, 送入 PLGF 分支训练光照不变特征提取网络。教师模型和学生模型均采用带 nesterov 动量的 SGD 优化器, 其中动量 μ 为 0.9, 学习率 ε 初始值为 0.05, 参数更新公式为:

$$\begin{cases} v_{t+1} = \mu v_t - \varepsilon \nabla f(\theta_t) \\ \theta_{t+1} = \theta_t + \mu v_{t+1} - \varepsilon \nabla f(\theta_t) \end{cases} \quad (27)$$

其中 v 为动量速度, t 为当前训练的迭代次数, $\nabla f(\theta_t)$ 为参数 θ_t 的梯度, ε 为当前的学习率。

同时学习率随着训练迭代次数衰减, 学习率更新公式为:

$$lr(t) = \begin{cases} 0, t < s_{wt} \\ \frac{t}{\max\{1, s_{wt} + s_{wu}\}}, s_{wt} \leq t < s_{wt} + s_{wu} \\ \max\{0, 0.5 + \frac{t - s_{wt} - s_{wu}}{\max\{1, s_{tl} - s_{wu} - s_{wt}\}} \cos(\frac{\pi}{2})\}, t \geq s_{wt} + s_{wu} \end{cases} \quad (28)$$

其中 t 为当前训练的迭代次数, $lr(t)$ 为第 t 次迭代时对应的学习率。 s_{wu} 为预热步数, 设为 3000。 s_{wt} 为等待步数, 其中教师模型优化器和学生模型优化器分别设为 0 和 3000。 s_{tl} 为总步数, 设为 64000。

光照不变特征提取网络采用带动量的 SGD 优化器, 其中动量 μ 为 0.9, 学习率 ε 初始值为 0.01, 参数更新公式为:

$$\begin{cases} v_{t+1} = \mu v_t - \varepsilon \nabla f(\theta_t) \\ \theta_{t+1} = \theta_t + v_{t+1} \end{cases} \quad (29)$$

其中 v 为动量速度, t 为当前训练的迭代次数, $\nabla f(\theta_t)$ 为参数 θ_t 的梯度, ε 为当前的学习率。

同时学习率随着训练迭代次数衰减, 设置第一门限步数 s_1 为 200, 第二门限步数 s_2 为 2000, 学习率更新公式为:

$$lr(t) = \begin{cases} \varepsilon, t \leq s_1 \\ 0.1\varepsilon, s_1 < t \leq s_2 \\ 0.01\varepsilon, t > s_2 \end{cases} \quad (30)$$

其中 t 为当前训练的迭代次数, $lr(t)$ 为迭代 t 时对应的学习率。

4.2.3 实验的测试参数

测试时输入大小为 $256 \times 256 \times 3$ 的原图, 分为两个支路: 原图送入训练好的学生模型输出 RGB 分类分数; 原图经过光照特征分离得到 PLGF 图, 送入训练好的光照不变特征提取网络输出 PLGF 分类分数。RGB 分类分数 p_{RGB} 与 PLGF 分类分数 p_{PLGF} 加权求和得到总的分类分数 p_{cls} :

$$p_{cls} = 0.8p_{RGB} + 0.2p_{PLGF} \quad (31)$$

4.3 算法评价指标

主要有错误接受率(False accept rate, FAR), 指攻击人脸判断成真实人脸数占标签为攻击人脸数的比率; 错误拒绝率(False refuse rate, FRR), 指真实人脸判断成攻击人脸数占标签为真实人脸数的比率; 等错误率(Equal error rate, EER), 指 FRR 和 FAR 相等时阈值决定的 FRR 或 FAR 数值。半总错误率(Half Total Error Rate, HTER), 指 FRR 与 FAR 的均值; EER 和 HTER 主要用来衡量算法在库内协议的准确率和跨库协议的泛化性能。

4.4 数据增强方式设置

为增强无标签数据多样性, 我们选取的数据增强方法包括有 15 种, 具体为: 最大化对比度、调节亮度、调节色彩平衡、调节对比度、调节锐度、裁剪、直方图均衡化、反色、像素值最低 0~4 个比特位置 0、随机旋转、水平方向错切、垂直方向错切、水平方向平移、垂直方向平移、将高于一定阈值的像素值反转。每个样本随机选取两种数据增强方法进行增强。

5 实验结果

5.1 与其他算法的对比

在库内实验中, 对于不含验证集的数据库 CASIA 和 MSU, 在其对应的训练集上训练, 然后在其对应的测试集上测试得到 EER; 对于含有验证集的数据库 Replay-Attack(以下简称 Replay), 在其对应的训练集上训练, 然后在其对应的验证集上测试得到 EER, 再根据 EER 所在的阈值在其对应的测试集上测试得到 HTER, 实验结果如表 5 所示。在跨库实验中, 各个数据库在对应的训练集上训练, 然后在对应的测试集上测试得到 EER 所在的阈值, 根据这个阈值在其他数据库的测试集上测试得到 HTER, 实验结果如表 6 所示。可以看到, 与传统手工特征文献[1,4]、基于深度学习文献[9,45,26]相比, 本方法库

内实验在 CASIA 和 MSU 中的 EER 仅次于文献[26], 在 Replay 的 EER 及 HTER 均仅次于文献[26]和[45]。而本方法跨库的 HTER 则是所对比算法中最低的, 证明本文方法的有效性和很好的稳健性能。

表 5 在库内与其他算法的 EER 及 HTER 对比(%)
Table 5 EER and HTER(%) of intra-test of different algorithms

数据库	CASIA	Replay		MSU
算法	EER	EER	HTER	EER
CNN ^[9]	7.4	6.1	2.1	—
Color-LBP ^[1]	6.2	0.4	2.9	—
Color Texture ^[4]	2.1	0.4	2.8	4.9
MSR-MobileNet ^[45]	4.2	0.1	0.3	3.0
DRL-FAS ^[26]	0.2	0.0	0.0	0.2
MPL+UDA+AttVLAD (本算法)	1.9	0.3	0.4	1.1

表 6 在跨库与其他算法的 HTER 对比(%)
Table 6 HTER(%) of cross-test of different algorithms

训练数据库	CASIA	Replay		MSU
测试数据库	Replay	CASIA	MSU	Replay
CNN ^[9]	48.5	45.5	48.6	37.1
Color-LBP ^[1]	37.9	35.4	33.0	44.8
Color Texture ^[4]	30.3	37.7	34.1	33.9
Auxiliary ^[15]	27.6	28.4	—	—
DeSpoofing ^[28]	28.5	41.1	27.8	33.2
MSR-MobileNet ^[45]	30.0	33.4	26.5	38.6
DRL-FAS ^[26]	28.4	33.2	15.6	29.7
MPL+UDA+AttVLA D(本算法)	8.5	23.2	10.1	12.1

5.2 消融实验

5.2.1 UDA、PLGF 分支对伪标签的作用

本小节主要讨论 3.1.2 小节的 UDA 及 3.2 的 PLGF 分支对 3.1.1 小节 RGB 分支伪标签的作用。如表 7 所示, 对于 RGB 分支的单流网络, PL+UDA 与 PL 相比, 在库内任务中 HTER 持平或下降, 而跨库,

表 7 UDA、PLGF 分支对伪标签的作用的消融实验的 HTER(%)

Table 7 HTER(%) of ablation experiments of the efforts on PL by UDA and PLGF branch

训练数据库	CASIA		Replay		MSU	
测试数据库	CASIA (库内)	Replay (跨库)	Replay (库内)	CASIA (跨库)	MSU (跨库)	MSU (库内)
PL	2.7	31.0	0.1	45.6	37.7	0.1
AttVLAD	1.9	25.6	0.3	20.5	18.1	1.1
PL+UDA	1.8	7.1	0.1	32.7	22.7	0.1
PL+UDA+ AttVLAD	1.8	6.9	0.5	23.4	19.3	0.6

(注: PL 表示输入 RGB 图像块到采用伪标签的 WideResNet; UDA 表示 RGB 分支采用伪标签; AttVLAD 表示输入 PLGF 图到 AttVLAD。)

时 HTER 大幅度下降, 说明 UDA 增加训练数据多样性对提高跨库性能的有效性; 与 PL+UDA 相比 PL+UDA+AttVLAD 加入了 AttVLAD 的 PLGF 分支, 在库内 HTER 处于较低水平, 同时进一步降低跨库任务的 HTER, 说明提取出 PLGF 图的光照不变特征有助于辅助 RGB 图像的判决。

5.2.2 元伪标签的作用

本小节主要讨论 3.3 小节的元伪标签对于 3.1 小

节半监督学习的改进作用。如表 8 所示, 在 RGB 分支的单流网络中, 元伪标签相比于伪标签, 库内 HTER 稍稍上升但整体较低, 跨库 HTER 在 Replay 跨 CASIA 任务中大幅下降, 在其他任务中大致持平; 在双流网络中, 元伪标签相比于伪标签在保证库内 HTER 较低的同时在部分跨库任务中(Replay 跨 CASIA、Replay 跨 MSU)的 HTER 上有一定程度的下降。

表 8 元伪标签的作用的消融实验的 HTER(%)

Table 8 HTER(%) of ablation experiments of the efforts by meta-pseudo-label

训练数据库	CASIA		Replay			MSU	
测试数据库	CASIA(库内)	Replay(跨库)	Replay(库内)	CASIA(跨库)	MSU(跨库)	MSU(库内)	Replay(跨库)
PL+UDA	1.8	7.1	0.1	32.7	22.7	0.1	12.1
MPL+UDA	2.8	7.4	0.8	27.5	22.8	0.1	11.8
PL+UDA+AttVLAD	1.8	6.9	0.5	23.4	19.3	0.6	7.0
MPL+UDA+AttVLAD	1.9	8.5	0.3	23.2	10.1	1.1	12.1

(注: PL 表示输入 RGB 图像块到采用伪标签的 WideResNet; MPL 表示输入 RGB 图像块到采用元伪标签的 WideResNet; UDA 表示 RGB 分支采用伪标签; AttVLAD 表示输入 PLGF 图到 AttVLAD。)

5.2.3 平滑度及置信度阈值的作用及参数选取

本小节主要讨论 3.1.2 小节采用的标签平滑和置信度阈值的作用以及参数选取。如表 9 所示, 分别在不同平滑度 α 及置信度阈值 t 下进行对比实验。可以看到, 采用标签平滑比不采用标签平滑($\alpha=0$)在跨库时 HTER 有小幅下降, 采用置信度阈值比不采用置信度阈值($t=0$)在跨库时 HTER 有小幅下降。在平滑度 α 取 0.1 和置信度阈值 t 取 0.5 时虽然没有在所有任务取得最低的 HTER, 但综合效果最好。

5.2.4 VLAD 对 PLGF 分支的作用

本小节主要讨论 3.2.2 小节的光照特征提取网络 AttVLAD 中基于 VLAD 的注意力模块的作用。如表 10 所示, 加入了基于 VLAD 的注意力模块, 在保持

库内低 HTER 的同时, 除了 CASIA 跨 Replay 任务, 在大多跨库任务中的 HTER 有较大幅度下降。

5.2.5 不同数量有标签样本的对比实验

本小节主要讨论在 RGB 分支半监督学习的训练过程中, 在样本总量同为 16000 时有标签样本数量的影响。如表 11 所示, 综合效果在 4000(即中等有标签样本数量)时效果最好。

5.2.6 元伪标签改进工作与原始工作的对比实验

本小节主要讨论在 3.3 节的元伪标签的第二层优化教师中, 改进的教师元学习损失 \mathcal{L}_{mpl} 相比文献 [55] 的教师元学习损失 $\mathcal{L}_{mpl}^{baseline}$ 的作用。如表 12 所示, 改进工作在保持库内低 HTER 的同时提升跨库性能。

表 9 不同平滑度及置信度阈值的消融实验的 HTER(%)

Table 9 HTER(%) of ablation experiments of different smoothness and confidence threshold

训练数据库	CASIA		Replay			MSU	
测试数据库	CASIA(库内)	Replay(跨库)	Replay(库内)	CASIA(跨库)	MSU(跨库)	MSU(库内)	Replay(跨库)
$t=0.6$	$\alpha=0$	3.2	8.1	2.5	27.9	20.1	0.5
	$\alpha=0.1$	2.8	8.9	1.0	26.5	19.7	0.7
	$\alpha=0.15$	2.8	7.4	0.8	27.5	22.7	0.7
$t=0.5$	$\alpha=0$	2.5	7.5	2.0	27.7	19.4	0.3
	$\alpha=0.1$	2.8	7.8	0.9	26.6	16.6	0.5
	$\alpha=0.15$	3.5	8.6	2.3	26.1	21.0	0.9
$t=0$	$\alpha=0$	2.3	7.5	2.0	28.1	20.6	0.4
	$\alpha=0.1$	3.1	8.7	1.4	26.8	15.7	0.5
	$\alpha=0.15$	3.0	10.0	1.9	23.8	18.9	0.7

(注: α 表示平滑度, t 表示置信度阈值, $\alpha=0$ 表示没有进行标签平滑, $t=0$ 表示没有使用置信度机制。)

表 10 VLAD 对 PLGF 分支的作用的消融实验的 HTER(%)

Table 10 HTER(%) of ablation experiments of the efforts on PLGF branch by VLAD

训练数据库	CASIA		Replay			MSU	
测试数据库	CASIA (库内)	Replay (跨库)	Replay (库内)	CASIA (跨库)	MSU (跨库)	MSU (库内)	Replay (跨库)
Baseline	1.1	20.4	0.2	30.7	31.4	1.1	20.3
AttVLAD	1.9	25.6	0.3	20.5	18.1	1.1	9.4

(注: baseline 表示没有基于 VLAD 的注意力模块的光照不变特征提取网络。)

表 11 RGB 分支在不同数量有标签样本的对比实验的 HTER(%)

Table 11 HTER(%) of experiments of different amounts of labeled samples

训练数据库	CASIA		Replay	
测试数据库	CASIA	Replay	Replay	CASIA
有标签样本				
3000	4.2	9.8	1.0	29.2
4000	2.8	7.4	0.8	27.5
5000	3.4	8.8	1.6	29.6
6000	3.4	9.2	2.0	26.5

表 12 不同教师元学习损失的对比实验的 HTER(%)

Table 12 HTER(%) of experiments of different loss function of meta-learning for teacher

训练数据库	CASIA		Replay	
测试数据库	CASIA	Replay	Replay	CASIA
原始工作	3.2	10.9	0.0	29.6
改进工作	2.8	7.8	0.9	26.6

5.2.7 图像块对 RGB 分支训练的作用及可视化

本小节主要讨论输入图像块对 RGB 分支训练的作用。

1) 实验结果: 分别将完整图像和图像块送入 RGB 分支的主干网 WideResNet 进行训练, 如表 13 所示, 裁剪的 64×64 大小的图像块与整张图像相比, 在库内两种情况效果相当, 而在跨库时, 图像块的效果更好。

2) Grad-CAM 可视化: 训练时输入 256×256 大小的 RGB 原图随机裁剪出 64×64 大小的图像块, 在测试时分别输入 256×256 大小的 RGB 原图和随机裁剪

出的 64×64 大小的图像块。以 CASIA 跨 Replay 任务为例, 在 CASIA 的库内测试及在 Replay 的跨库测试时使用 Grad-CAM 得到的热力图分别如图 7 和图 8 所示。从可视化结果来看, 更深层次卷积层提取的高层次特征更能表征分类结果。模型更侧重对五官外平坦人脸区域进行判别, 说明其认为这些区域更能体现真假区别。将图像块作为训练输入, 对感兴趣区域提取的特征粒度更小、分散, 有助于模型做综合全面的判别。

3) 测试样本分类分数分布图: 测试输入尺寸为 256×256 时, 在 CASIA 跨 Replay 任务的库内测试及跨库测试的分类分数分布图如图 9 所示, 模型能够较好区分样本, 且在跨库任务性能没有太大下降。

5.2.8 网络的时间复杂度比较

本小节主要讨论 RGB 分支 WideResNet 和 PLGF 分支 AttVLAD 的时间复杂度, 并与其他主流用于人脸反欺诈的主干网络结构进行对比。实验环境均相同, 选取参数量、浮点运算次数(Floating-Point Operations, FLOPs)和平均帧检测时间三个指标来评估不同主干网络结构的复杂度, 其中 FLOPs 以百万(M)或十亿(G)为单位, 平均帧检测时间以 ms 为单位。具体选取 Replay 数据库上 480 段视频共 15040 帧上利用训练完成的模型进行逐帧测试, 并计算模型测试总时长, 最后根据测试总时长和测试总帧数计算平均帧检测时间。对比结果如表 14 所示, 与其他主流网络结构相比, RGB 分支 WideResNet 的参数减少到 1.47M, 但由于通道数增加, FLOPs 增加到 10.57G, 平均帧检测时间增加到 4.55ms, 所以计算复杂度有所增加。PLGF 分支 AttVLAD 的参数量为 96K,

表 13 WideResNet 在不同尺寸作训练输入的对比实验的 HTER(%)

Table 13 HTER(%) of experiments of feeding different size of training input to WideResNet

训练数据库	CASIA		Replay			MSU	
测试数据库	CASIA (库内)	Replay (跨库)	Replay (库内)	CASIA (跨库)	MSU (跨库)	MSU (库内)	Replay (跨库)
输入尺寸							
64×64	0.2	28.2	3.0	41.0	28.0	0.5	16.4
256×256	0.1	31.2	3.2	43.1	37.1	2.2	16.4

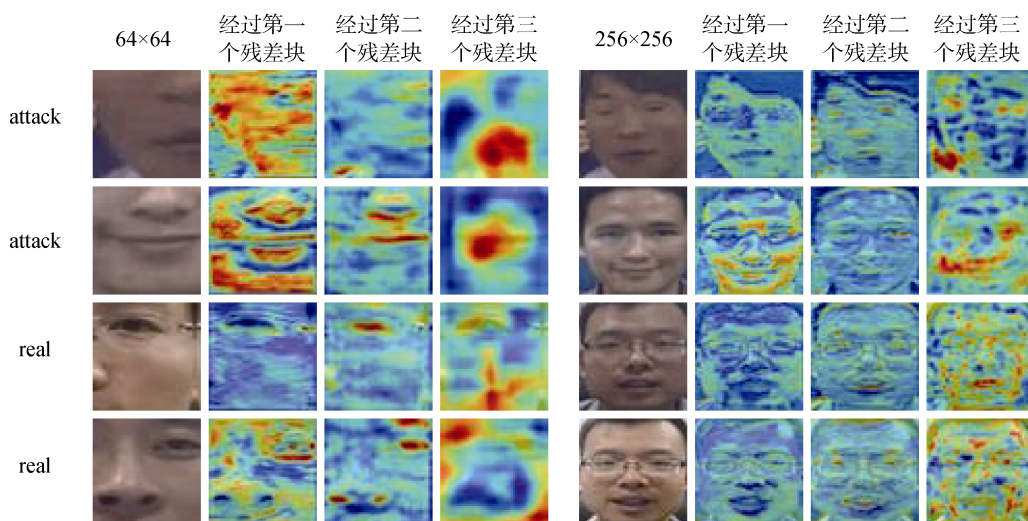


图 7 测试时分别输入 CASIA 原图和图像块的热力图

Figure 7 Heatmap of testing input of CASIA original images and patches

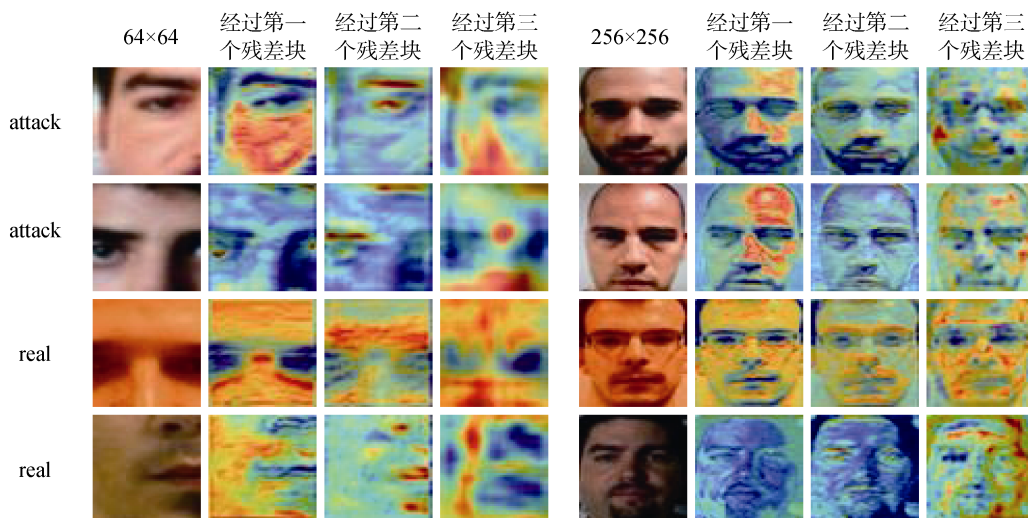


图 8 测试时分别输入 Replay 原图和图像块的热力图

Figure 8 Heatmap of testing input of Replay original images and patches

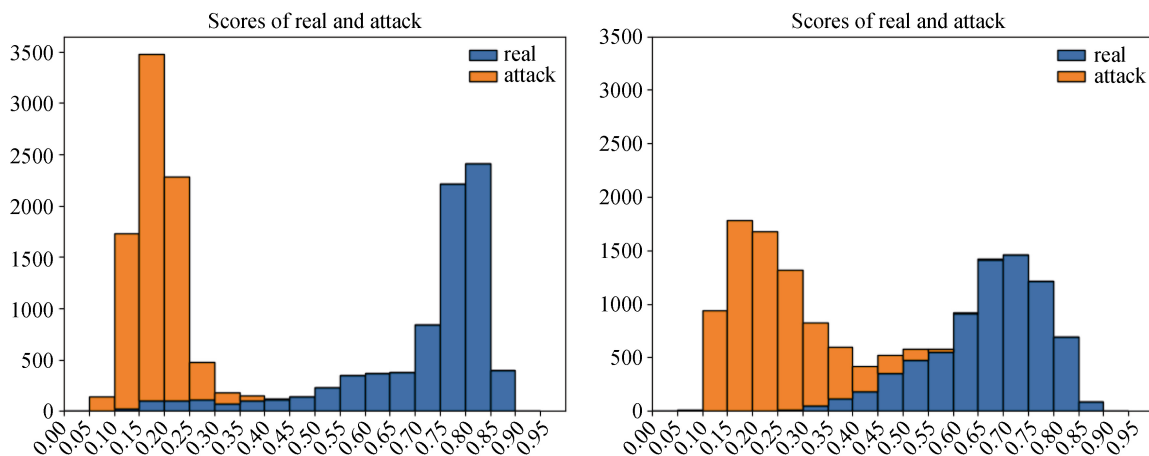


图 9 CASIA 跨 Replay 任务的库内测试(左)及跨库测试(右)的分类分数分布图

Figure 9 Distribution of classification scores of intra test(left) and cross test(right) for task CASIA to Replay

表 14 不同网络结构的时间复杂度对比

Table 14 Time complexity of different network architecture

网络结构	参数量 (Params)	浮点运算次数 (FLOPs)	平均帧检测 时间/ms
ResNet18	11.69M	1.82G	1.95
Inception3	23.83M	2.85G	2.15
MobileNetV2	3.50M	320.24M	1.94
DenseNet	7.79M	2.88G	2.30
WideResNet (本算法)	1.47M	10.57G	4.55
AttVLAD (本算法)	96K	44.12M	1.98

FLOPs 为 44.12M, 平均帧检测时间 1.98ms, 可以在增加少量时间复杂度的同时辅助 RGB 分支提升性能。综上, 本文网络能够适用于在线人脸识别系统, 而对于嵌入式终端、移动设备等应用场景, 则需进一步做模型轻量化。

6 结论

针对有标签样本数据成本不菲、获取不易以及光照条件直接影响检测模型的泛化性能这两个问题, 本文提出了一种“教师生成伪标签, 学生反馈”的半监督学习框架, 与元学习相结合, 通过带有伪标签的未标注数据来丰富训练数据的多样性, 并利用标签平滑和置信度机制来提高置信度; 另一方面, 设计了带有注意力机制的 PLGF 分支, 提取光照不变特征, 在提高判别准确度的同时改善泛化性能。在多个公开活体检测数据集上的实验结果表明, 与同类最新算法相比, 本文方法在库内检测保持高准确率的同时大幅降低了跨库检测半总错误率, 算法泛化性能得到显著提高。未来的改进方向包括用像素图替代二元标签监督、模型轻量化及损失函数的优化等。

参考文献

- [1] Boulkenafet Z, Komulainen J, Hadid A. Face Spoofing Detection Using Colour Texture Analysis[J]. *IEEE Transactions on Information Forensics and Security*, 2016, 11(8): 1818-1830.
- [2] Boulkenafet Z, Komulainen J, Hadid A. Face Antispoofing Using Speeded-up Robust Features and Fisher Vector Encoding[J]. *IEEE Signal Processing Letters*, 2017, 24(2): 141-145.
- [3] Chingovska I, Anjos A, Marcel S. On the Effectiveness of Local Binary Patterns in Face Anti-Spoofing[C]. *The International Conference of Biometrics Special Interest Group*, 2012: 1-7.
- [4] Komulainen J, Hadid A, Pietikäinen M. Context Based Face Anti-Spoofing[C]. *2013 IEEE Sixth International Conference on Biometrics: Theory, Applications and Systems*, 2014: 1-8.
- [5] Wen D, Han H, Jain A K. Face Spoof Detection with Image Distor-

- tion Analysis[J]. *IEEE Transactions on Information Forensics and Security*, 2015, 10(4): 746-761.
- [6] Galbally J, Marcel S. Face Anti-Spoofing Based on General Image Quality Assessment[C]. *2014 22nd International Conference on Pattern Recognition*, 2014: 1173-1178.
- [7] de Freitas Pereira T, Komulainen J, Anjos A, et al. Face Liveness Detection Using Dynamic Texture[J]. *EURASIP Journal on Image and Video Processing*, 2014, 2014(1): 1-15.
- [8] Komulainen J, Hadid A, Pietikäinen M, et al. Complementary Countermeasures for Detecting Scenic Face Spoofing Attacks[C]. *2013 International Conference on Biometrics*, 2013: 1-7.
- [9] Yang J W, Lei Z, Li S Z. Learn Convolutional Neural Network for Face Anti-Spoofing[EB/OL]. 2014: arXiv: 1408.5601. <https://arxiv.org/abs/1408.5601>
- [10] Yan Z C, Zhang H, Piramuthu R, et al. HD-CNN: Hierarchical Deep Convolutional Neural Networks for Large Scale Visual Recognition[C]. *2015 IEEE International Conference on Computer Vision*, 2016: 2740-2748.
- [11] Hu Y J, Wang Y F, Liu B B, et al. A Survey on the Latest Development and Typical Methods of Face Anti-Spoofing[J]. *Journal of Signal Processing*, 2021, 37(12): 2261-2277.
(胡永健, 王宇飞, 刘琲贝, 等. 人脸欺诈检测最新进展及典型方法[J]. *信号处理*, 2021, 37(12): 2261-2277.)
- [12] Yu Z T, Zhao C X, Wang Z Z, et al. Searching Central Difference Convolutional Networks for Face Anti-Spoofing[C]. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020: 5294-5304.
- [13] Zhou J W, Shu K, Liu P, et al. Face Anti-Spoofing Based on Dynamic Color Texture Analysis Using Local Directional Number Pattern[C]. *2020 25th International Conference on Pattern Recognition*, 2021: 4221-4228.
- [14] Chen B L, Yang W H, Wang S Q. Generalized Face Antispoofing by Learning to Fuse Features from High- and Low-Frequency Domains[J]. *IEEE MultiMedia*, 2021, 28(1): 56-64.
- [15] Liu Y J, Jourabloo A, Liu X M. Learning Deep Models for Face Anti-Spoofing: Binary or Auxiliary Supervision[C]. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018: 389-398.
- [16] Wu X J, Zhou J H, Liu J, et al. Single-Shot Face Anti-Spoofing for Dual Pixel Camera[J]. *IEEE Transactions on Information Forensics and Security*, 2021, 16: 1440-1451.
- [17] Yu Z, Li X, Niu X, et al. Face anti-spoofing with human material perception[C]. *2020 European Conference on Computer Vision*, 2020: 557-575.
- [18] Wang Z Z, Yu Z T, Zhao C X, et al. Deep Spatial Gradient and Temporal Depth Learning for Face Anti-Spoofing[C]. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020: 5041-5050.
- [19] Atoum Y, Liu Y J, Jourabloo A, et al. Face Anti-Spoofing Using Patch and Depth-Based CNNs[C]. *2017 IEEE International Joint Conference on Biometrics*, 2018: 319-328.
- [20] Li X, Wan J, Jin Y, et al. 3DPC-Net: 3D Point Cloud Network for Face Anti-Spoofing[C]. *2020 IEEE International Joint Conference on Biometrics*, 2021: 1-8.

- [21] Liu W H, Wei X K, Lei T, et al. Data-Fusion-Based Two-Stage Cascade Framework for Multimodality Face Anti-Spoofing[J]. *IEEE Transactions on Cognitive and Developmental Systems*, 2021, 14(2): 672-683.
- [22] Liu A J, Tan Z C, Wan J, et al. Face Anti-Spoofing via Adversarial Cross-Modality Translation[J]. *IEEE Transactions on Information Forensics and Security*, 2021, 16: 2759-2772.
- [23] Chiesa V, Dugelay J L. Advanced Face Presentation Attack Detection on Light Field Database[C]. *2018 International Conference of the Biometrics Special Interest Group*, 2018: 1-4.
- [24] Quan R J, Wu Y, Yu X, et al. Progressive Transfer Learning for Face Anti-Spoofing[J]. *IEEE Transactions on Image Processing*, 2021, 30: 3946-3955.
- [25] Qin Y X, Zhao C X, Zhu X Y, et al. Learning Meta Model for Zero- and Few-Shot Face Anti-Spoofing[J]. *The AAAI Conference on Artificial Intelligence*, 2020, 34(7): 11916-11923.
- [26] Cai R Z, Li H L, Wang S Q, et al. DRL-FAS: A Novel Framework Based on Deep Reinforcement Learning for Face Anti-Spoofing[J]. *IEEE Transactions on Information Forensics and Security*, 2021, 16: 937-951.
- [27] Zhang K Y, Yao T P, Zhang J A, et al. Face Anti-Spoofing via Disentangled Representation Learning[M]. *Computer Vision – ECCV 2020*. Cham: Springer International Publishing, 2020: 641-657.
- [28] Jourabloo A, Liu Y J, Liu X M. Face De-Spoofing: Anti-Spoofing via Noise Modeling[M]. *Computer Vision – ECCV 2018*. Cham: Springer International Publishing, 2018: 297-315.
- [29] Wang G Q, Han H, Shan S G, et al. Cross-Domain Face Presentation Attack Detection via Multi-Domain Disentangled Representation Learning[C]. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020: 6677-6686.
- [30] Qin Y X, Zhang W G, Shi J P, et al. One-Class Adaptation Face Anti-Spoofing with Loss Function Search[J]. *Neurocomputing*, 2020, 417: 384-395.
- [31] Fatemifar S, Arashloo S R, Awais M, et al. Client-Specific Anomaly Detection for Face Presentation Attack Detection[J]. *Pattern Recognition*, 2021, 112: 107696.
- [32] Li Z, Li H L, Lam K Y, et al. Unseen Face Presentation Attack Detection with Hypersphere Loss[C]. *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing*, 2020: 2852-2856.
- [33] Bhattacharjee D, Roy H. Pattern of Local Gravitational Force (PLGF): A Novel Local Image Descriptor[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021, 43(2): 595-607.
- [34] Olivier C, Bernhard S, Alexander Z. Semi-supervised learning [J]. *IEEE Transactions on Neural Networks*, 2009, 20(3): 542-542.
- [35] Hinton G, Vinyals O, Dean J. Distilling the Knowledge in a Neural Network[EB/OL]. 2015: arXiv: 1503.02531. <https://arxiv.org/abs/1503.02531>
- [36] Furlanello T, Lipton Z C, Tschannen M, et al. Born again Neural Networks[EB/OL]. 2018: arXiv: 1805.04770. <https://arxiv.org/abs/1805.04770>
- [37] Lee D. Pseudo-Label: The simple and efficient semi-supervised learning method for deep neural networks[C]. *International Conference on Machine Learning Workshop*, 2013: 7.
- [38] Müller R, Kornblith S, Hinton G. When does Label Smoothing Help? [EB/OL]. 2019: arXiv: 1906.02629. <https://arxiv.org/abs/1906.02629>
- [39] Arazo E, Ortego D, Albert P, et al. Pseudo-Labeling and Confirmation Bias in Deep Semi-Supervised Learning[C]. *2020 International Joint Conference on Neural Networks*, 2020: 1-8.
- [40] Xie Q Z, Luong M T, Hovy E, et al. Self-Training with Noisy Student Improves ImageNet Classification[C]. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020: 10684-10695.
- [41] Xie Q, Dai Z, Hovy E, et al. Unsupervised data augmentation for consistency training[C]. *Advances in Neural Information Processing Systems*, 2020, 3: 6256-6268.
- [42] Chen H T, Wang Y H, Xu C, et al. Data-Free Learning of Student Networks[C]. *2019 IEEE/CVF International Conference on Computer Vision*, 2020: 3513-3521.
- [43] Chen H, Lagadec B, Brémond F. Enhancing Diversity in Teacher-Student Networks via Asymmetric Branches for Unsupervised Person re-Identification[C]. *2021 IEEE Winter Conference on Applications of Computer Vision*, 2021: 1-10.
- [44] Pinto A, Goldenstein S, Ferreira A, et al. Leveraging Shape, Reflectance and Albedo from Shading for Face Presentation Attack Detection[J]. *IEEE Transactions on Information Forensics and Security*, 2020, 15: 3347-3358.
- [45] Chen H N, Hu G S, Lei Z, et al. Attention-Based Two-Stream Convolutional Networks for Face Spoofing Detection[J]. *IEEE Transactions on Information Forensics and Security*, 2020, 15: 578-593.
- [46] Finn C, Abbeel P, Levine S. Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks[EB/OL]. 2017: arXiv: 1703.03400. <https://arxiv.org/abs/1703.03400>
- [47] Shao R, Lan X Y, Yuen P C. Regularized Fine-Grained Meta Face Anti-Spoofing[J]. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, 34(7): 11974-11981.
- [48] Yu Z T, Wan J, Qin Y X, et al. NAS-FAS: Static-Dynamic Central Difference Network Search for Face Anti-Spoofing[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021, 43(9): 3005-3023.
- [49] Qin Y X, Yu Z T, Yan L B, et al. Meta-Teacher for Face Anti-Spoofing[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022, 44(10): 6311-6326.
- [50] Jégou H, Douze M, Schmid C, et al. Aggregating Local Descriptors into a Compact Image Representation[C]. *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2010: 3304-3311.
- [51] Arandjelović R, Gronat P, Torii A, et al. NetVLAD: CNN Architecture for Weakly Supervised Place Recognition[J]. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018, 40(6): 1437-1451.
- [52] Lin R C, Xiao J, Fan J P. NeXtVLAD: An Efficient Neural Network to Aggregate Frame-Level Features for Large-Scale Video Classification[M]. *Lecture Notes in Computer Science*. Cham: Springer International Publishing, 2019: 206-218.

- [53] Zagoruyko S, Komodakis N. Wide residual networks[C]. *British Machine Vision Conference*, 2016, 87: 1-12.
- [54] Zhang Z W, Yan J J, Liu S F, et al. A Face Antispoofing Database with Diverse Attacks[C]. *2012 5th IAPR International Conference*

on Biometrics, 2012: 26-31.

- [55] Pham H, Dai Z H, Xie Q Z, et al. Meta Pseudo Labels[C]. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021: 11552-11563.



冯浩宇 于 2019 年在华中师范大学电子信息工程专业获得学士学位。现在华南理工大学网络空间安全专业攻读硕士学位。研究领域为多媒体信息安全、人工智能及其应用。研究兴趣包括: 人脸欺诈检测、深度学习网络。Email: 201920141513@mail.scut.edu.cn



胡永健 于 2002 年在华南理工大学信息与通信工程专业获得博士学位。现为华南理工大学电子与信息学院教授。研究领域为多媒体信息安全、图像处理、人工智能及其应用。研究兴趣包括: 图像和视频信息隐藏、人脸欺诈检测、图像和视频取证。Email: eeyjhu@scut.edu.cn



王宇飞 于 2018 年在华南理工大学信息与通信工程专业获得博士学位。现任职于广东警官学院刑事技术系讲师。研究领域为多媒体信息安全。研究兴趣包括: 多媒体取证、人工智能应用。Email: 20220710@gdppla.edu.cn



刘琲贝 于 2009 年在中山大学获得博士学位。现为华南理工大学电子与信息学院讲师。研究领域为多媒体信息安全、图像处理、人工智能及其应用。研究兴趣包括: 图像和视频信息隐藏、人脸欺诈检测、图像和视频取证。Email: eebbliu@scut.edu.cn



余翔宇 于 2006 年在武汉大学通信与信息系统专业获得博士学位。现为华南理工大学电子与信息学院副教授。研究领域为多媒体信息安全、图像处理、人工智能及其应用。研究兴趣包括: 信息隐藏。Email: yuxy@scut.edu.cn



钟睿 于 2019 年在暨南大学通信工程专业获得学士学位。现在华南理工大学电子信息专业攻读硕士学位。研究领域为多媒体信息安全、人工智能及其应用。研究兴趣包括: 人脸欺诈检测。Email: 202021012794@mail.scut.edu.cn