

文章编号:1007-5321(2020)05-0048-09

DOI:10.13190/j.jbupt.2020-068

基于前向学习网络的人脸欺诈检测

宋 昱, 孙文赞, 陈昌盛

(1. 深圳大学 电子与信息工程学院, 深圳 518060; 2. 深圳大学 深圳市媒体信息内容安全重点实验室, 深圳 518060;
3. 深圳大学 广东省智能信息处理重点实验室, 深圳 518060; 4. 深圳大学 广东省人工智能与数字经济实验室, 深圳 518060;
5. 深圳市人工智能与机器人研究院, 深圳 518060)

摘要: 为了克服现有人脸欺诈检测方法在少样本应用场合下的局限性, 将前向学习网络用于欺诈检测. 通过前向学习的方式从图像中无监督地学得卷积滤波器, 在人脸欺诈检测应用场合下, 对前向学习网络进行了改进, 改进后的网络使用了面向人脸欺诈检测任务的卷积滤波器. 使用主成分分析变换所得的最小特征值对应的特征向量作为卷积滤波器提取图像的特征. 将所提方法在 CASIA-FASD、Idiap Replay-Attack 和 OULU-NPU 数据集上进行了验证, 实验结果表明, 在少样本跨攻击类型实验中, 所提方法显著提升了欺诈人脸检测的准确率.

关键词: 人脸欺诈检测; 前向学习网络; 表示学习

中图分类号: TP309; TP391.41 **文献标志码:** A

Few-Shot Face Spoofing Detection Using Feedforward Learning Network

SONG Yu, SUN Wen-yun, CHEN Chang-sheng

(1. College of Electronics and Information Engineering, Shenzhen University, Shenzhen 518060, China;
2. Shenzhen Key Laboratory of Media Security, Shenzhen University, Shenzhen 518060, China;
3. Guangdong Key Laboratory of Intelligent Information Processing, Shenzhen University, Shenzhen 518060, China;
4. Guangdong Laboratory of Artificial Intelligence and Digital Economy, Shenzhen University, Shenzhen 518060, China;
5. Shenzhen Institute of Artificial Intelligence and Robotics for Society, Shenzhen 518060, China)

Abstract: In order to overcome the limitations of the existing face spoofing detection methods under few-shot face anti-spoofing applications, this paper proposes to use feedforward learning network for face anti-spoofing. The convolutional filters are learned unsupervisedly from the images in a feedforward manner. The feedforward learning network is adapted in the spoof face detection applications by using face anti-spoofing task-oriented convolutional filters learned from the training images. The eigenvectors that correspond to the smallest eigenvalues obtained from the principle component analysis transform are used as convolution filters for extracting features from images. The method is evaluated on some benchmark datasets including CASIA-FASD dataset, Idiap Replay-Attack dataset and OULU-NPU dataset. Experiments show that under the cross presentation attack detection experiments, the proposed method significantly improves the classification accuracy of existing methods.

Key words: face spoofing detection; feedforward learning network; representation learning

收稿日期: 2020-06-15

基金项目: 中国博士后科学基金项目(2019M663068); 广东省基础与应用基础研究基金项目(2019A1515110425); 广东省自然科学基金项目(2020A1515010563); 深圳市科技计划项目(JCYJ20180305124550725)

作者简介: 宋 昱(1988—), 男, 博士后, E-mail: songy@szu.edu.cn.

随着数字设备和人脸识别技术的发展,人脸识别技术在信息安全、访问控制和身份验证方面得到了越来越多的应用。然而,人脸识别系统对于人脸欺诈攻击比较脆弱。人脸欺诈攻击主要包括照片、视频和3维面具攻击。攻击者使用欺诈人脸模仿目标对象的真实人脸,从而达到对人脸识别系统非法入侵的目的。

在机器学习的研究中,当训练样本数量较少时,研究者提出使用前向学习网络^[1]来自动学得特征提取算子。通过实验可以证明,在少样本学习场景下,前向学习网络的推广性优于基于深度学习的方法^[1]。以主成分分析网络(PCANet)^[1]为例,通过对训练样本的图像块进行主成分分析(PCA, principal component analysis)得到卷积滤波器。在少样本学习问题中,PCANet^[1]有着较好的性能。受到前向学习网络处理少样本学习问题能力的启发,笔者提出使用前向学习网络用于少样本人脸欺诈检测。所提算法有以下3点主要贡献。

1) 将前向学习网络结构用于少样本人脸欺诈检测任务中。与基于人工特征的算法不同,所提算法自动地从样本中学得特征;与基于深度学习的方法不同,所提方法无监督地从训练样本中学得特征。这一特性使得所提算法适用于少样本人脸欺诈检测任务中。

2) 从训练样本的图像块中通过PCA变换学得面向人脸欺诈检测任务的卷积滤波器。在PCANet^[1]中,使用的卷积滤波器是协方差矩阵最大特征值对应的特征向量。与之相对的,在所提方案中,使用的卷积滤波器是协方差矩阵最小特征值对应的特征向量。所提方案中使用的卷积滤波器与PCANet^[1]中使用的卷积滤波器有着本质不同。

3) 考虑到人脸欺诈检测任务的特殊性,对前向网络的最后一层进行了改进。去除了块直方图的计算步骤,所提方案使用整幅图像计算直方图。

1 相关工作

为了保证访问的合法性,人脸欺诈检测应该作为人脸识别系统的一个必要部分,可以根据算法原理将其大致区分为2类:第1类算法根据人工特征检测欺诈人脸;第2类算法根据从数据中学得的特征检测欺诈人脸。

基于人工特征的人脸欺诈检测算法具有如下特征:区分真实人脸和欺诈人脸的特征是通过人工设

计的特征提取算子提取的,需要根据特定领域的知识设计特征提取算子。基于人工特征的人脸欺诈检测算法根据使用的特征可大致分为以下4类。

1) 基于纹理特征的算法。Boulkenafet等^[2-3]通过提取不同色彩空间中互补的低层特征得到了亮度通道和色彩通道的颜色-纹理联合特征;Boulkenafet等^[4]提取了人脸图像不同色彩空间的加速鲁棒特征。

2) 基于摩尔纹特征的算法。Garcia等^[5]使用了因数字网格重叠产生的摩尔纹,其可通过频域的峰值检测得到。

3) 基于图像质量评价的算法。Galbally等^[6]使用了从一幅图像中提取的25个图像质量特征区分真实人脸和欺诈人脸。

4) 基于组合特征的算法。Wen等^[7]提出了一种基于图像失真分析的人脸欺诈检测算法。

基于人工特征的方法有如下缺点:首先,需要具有特定领域知识的专家设计这些特征提取算子,这些关于真实人脸和欺诈人脸特征差异的知识可能并不全面;其次,一旦设计完成了特征提取算子,它们就是固定的。它们可从特定类型的欺诈攻击中提取特征,但是对于不同类型的欺诈攻击,它们可能会失效。这些方法不能很好地处理未见攻击类型。另外,通过实验可以证明,这些方法在少样本人脸欺诈检测场景中表现不佳,这说明需要大量训练样本才能使得这些方法在实际场合中具有足够的推广性。

随着深度学习的发展,更多的研究者开始使用基于数据驱动的方法进行人脸欺诈检测。这一类方法使用深度神经网络从训练样本中自动学得判别特征。基于数据驱动特征的人脸欺诈检测算法可以大致分为以下2类。

1) 基于学得判别特征的算法。Rehman等^[8]在CASIA-FASD数据集上训练了一个11层的视觉几何研究组(VGG, visual geometry group)网络。Nagpal等^[9]在MSU移动人脸欺诈数据集上训练了一个Inception网络和残差网络,用于评价不同卷积神经网络(CNN, convolutional neural networks)结构在反欺诈人脸应用场景中的性能。考虑到人脸欺诈检测任务的局部性,Atoum等^[10]使用了辅助监督信息,除了全局深度图,该方法还使用了从人脸图像中提取的局部特征;Liu等^[11]将逐像素点人脸深度监督信息和序列rPPG信号监督信息作为欺诈检测深度网络的辅助监

督信息;George等^[12]在CNN框架中使用了深度逐像素点监督信息,其中监督信息包含了二值信息而不是深度信息.

2) 基于学得其他特征的方法. Li等^[13]提出在CNN框架中使用人工特征,从卷积特征图中提取色彩局部二值模式,卷积特征通过在VGG-face模型上精调得到;Li等^[14]提出了一种无监督域自适应欺诈人脸检测方案.

一些基于人工特征的人脸欺诈检测算法的缺点被基于数据驱动的方法克服了. 使用从训练样本中自动学得特征,人脸欺诈检测的能力得到提升. 在未见攻击类型检测场合中,基于数据驱动方法的性能比基于人工特征的方法好. 在测试阶段,基于数据驱动的方法对于未见样本更加鲁棒. 然而,基于数据驱动的算法有一个主要的缺点,即其通常需要大量的训练样本才能学得神经网络的权重. 当训练样本数量不够时,由于其容易过拟合而导致性能不佳. 通过实验可以证明,在少样本欺诈检测场合中,基于数据驱动方法的性能并不好.

2 基于前向学习网络的人脸欺诈检测算法

2.1 从PCA滤波器到面向人脸欺诈检测任务的卷积滤波器

PCA变换的目标函数可以写为

$$\min_{V \in \mathbb{R}^{D \times L}} \|X - VV^T X\|_F^2, \text{ s. t. } V^T V = I_L \quad (1)$$

其中: X 为包含训练样本的 $D \times N$ 矩阵, D 为 X 中训练样本的维度, L 为主成分的个数. 因为 X 的每一列都是一个训练样本, X 可以写成 $X = [x_1, x_2, \dots, x_N]$,将其代入式(1)中,可以得到

$$\min_{V \in \mathbb{R}^{D \times L}} \| [x_1, x_2, \dots, x_N] - VV^T [x_1, x_2, \dots, x_N] \|_F^2, \text{ s. t. } V^T V = I_L \quad (2)$$

式(2)可以进一步写为

$$\min_{V \in \mathbb{R}^{D \times L}} \sum_{i=1}^N \|x_i - VV^T x_i\|_F^2, \text{ s. t. } V^T V = I_L \quad (3)$$

由式(2)可以推出式(3)的原因如下:根据矩阵Frobenius范数的定义,有

$$\|X\|_F^2 = \sum_{i=1}^D \sum_{j=1}^N x_{ij}^2 = \sum_{j=1}^N \left(\sum_{i=1}^D x_{ij}^2 \right) = \sum_{j=1}^N \|x_j\|_F^2 \quad (4)$$

其中: x_j 为矩阵 X 的第 j 列. 式(3)中的 $x_i - VV^T x_i$ 恰好是式(2)中矩阵的第 i 列,所以可以由式(2)推

出式(3). 按照列对矩阵 V 进行分块,即 $V = [v_1, v_2, \dots, v_L]$,然后代入式(3),可以得到

$$\min_{v_i \in \mathbb{R}^D, i=1,2,\dots,L} \sum_{i=1}^N \|x_i - [v_1, v_2, \dots, v_L] \times [v_1, v_2, \dots, v_L]^T x_i\|_F^2, \text{ s. t. } [v_1, v_2, \dots, v_L]^T [v_1, v_2, \dots, v_L] = I_L \quad (5)$$

式(5)可以进一步写为

$$\min_{v_i \in \mathbb{R}^D, i=1,2,\dots,L} \sum_{i=1}^N \left\| x_i - \sum_{j=1}^L (v_j^T x_i) v_j \right\|_F^2, \text{ s. t. } v_j^T v_i = \delta_{ij} \quad (6)$$

其中:当 $i=j$ 时, δ_{ij} 取1;当 $i \neq j$ 时, δ_{ij} 取0. 根据式(6),可以将矩阵 V 的每一列看作表示训练样本 x_i 的一个原子,表示系数是训练样本 x_i 在相应原子上的投影. V 中的原子构成了一组标准正交基. 目标是希望用一组标准正交基中的 L 个原子,使得重构 X 中训练样本的误差最小. 为了最小化重构误差, V 中的原子必须能够表示 X 中主要的变化,换句话说, X 中的“低频”特征通过这组基重构出来. 从而, V 中的原子是矩阵 XX^T 的主成分.

假设现在有一幅CASIA数据集中的图像,需将彩色图像转为灰度图像. 从图像中提取互相之间有重叠的 5×5 的图像块,然后计算图像块协方差矩阵的特征向量. 将特征向量按照特征值从大到小的顺序排列,然后将这些特征向量转换为 5×5 的图像,如图1所示. 假设 $L_1=8$,即第1层使用了8个滤波器. 在PCANet^[1]中使用的8个滤波器是前8个特征值对应的特征向量,用实线框标出. 这些是主成分,它们可以提取图像中的“低频”特征. 被这些滤波器提取的特征在图像识别中非常有用,因为图像“低频”特征是分类的重要依据. 然而,人脸欺诈检测和图像分类却是不同的任务. 图2给出了一幅真实人脸图像和一幅欺诈人脸图像,攻击类型是照片攻击. 从图2可以看出,图像的主要特征是一致的,差别在于细节部分. 如果要进行人脸欺诈检测,必须利用图像中的“高频”特征,这些“高频”特征对应的是图像中的纹理和噪声信息. 在传统的基于人工特征的方法中,经常使用纹理信息进行人脸欺诈检测. 受到人脸欺诈检测特殊性的启发,提出使用较小特征值对应的特征向量作为卷积滤波器. 这些特征向量可以称为次成分. 在所提的方法中,使用了图1中虚线框内的特征向量. 可以看出,这些特征向量和局部二值模式特征提取算子类似. 这些特征向量对于提取图像中的“高频”特征非常有效,从而

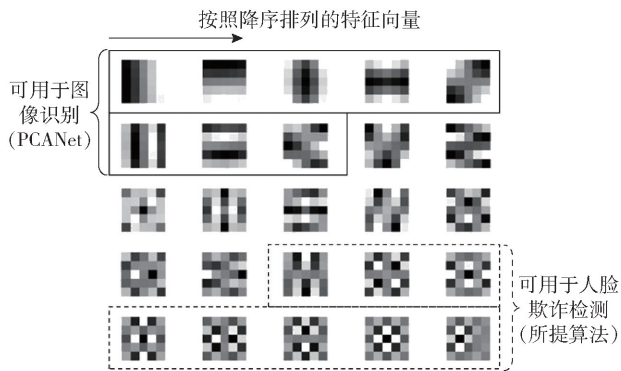
图1 PCANet^[1]和所提算法使用的滤波器结果

图2 真实人脸和欺诈人脸

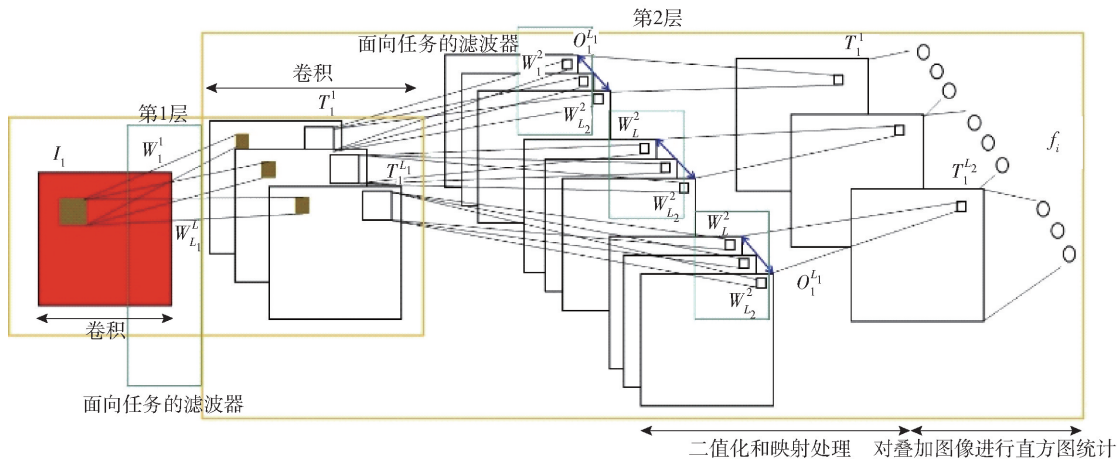


图3 所提算法的结构框图

$\cdots, \mathbf{x}_{i, \tilde{m}'\tilde{n}'} \in \mathbb{R}^{k_1 k_2}$, 其中 $\mathbf{x}_{i,j}$ 表示 I_i 中第 j 个向量化 (向量化表示将图像块拉成一个列向量) 的图像块, \tilde{m}' 和 \tilde{n}' 分别表示垂直方向和水平方向提取图像块的个数. 将这些图像块进行拼接, 得到 $\mathbf{X}_i = [\mathbf{x}_{i,1}, \mathbf{x}_{i,2}, \cdots, \mathbf{x}_{i, \tilde{m}'\tilde{n}'}]$. 将所有训练图像的矩阵进行拼接, 得到

$$\mathbf{X} = [\mathbf{X}_1, \mathbf{X}_2, \cdots, \mathbf{X}_N] \in \mathbb{R}^{k_1 k_2 \times N \tilde{m}' \tilde{n}'} \quad (8)$$

图像块的协方差矩阵可通过 $\mathbf{X}\mathbf{X}^T$ 计算. 对该矩阵进行特征值分解, 然后提取对应于 L_1 个最小特征

有利于进行人脸欺诈检测. 得到次成分的目标函数为

$$\max_{\mathbf{V} \in \mathbb{R}^{D \times L}} \|\mathbf{X} - \mathbf{V}\mathbf{V}^T \mathbf{X}\|_F^2, \text{ s. t. } \mathbf{V}^T \mathbf{V} = \mathbf{I}_L \quad (7)$$

由式(7)可以看出, 在所提方法中, 得到的卷积滤波器使得图像重构误差最大. 这些卷积滤波器提取的特征对应于图像中最不明显的变化, 这些特征将作为欺诈检测的重要依据.

2.2 基于面向人脸欺诈检测任务滤波器的人脸欺诈检测网络

基于1.1节中的分析, 提出使用面向人脸欺诈检测任务滤波器的用于人脸欺诈检测的前向学习网络. 假设现在有 N 幅训练图像 $\{I_i\}_{i=1}^N$, 每幅图像的大小为 $m \times n$, 对应的标签是 $\{l_i\}_{i=1}^N$. 图像块和卷积滤波器的大小为 $k_1 \times k_2$. 所提算法的结构框图与 PCANet^[1] 相似, 如图3所示. 在学习特征的过程中, 只使用了训练图像 $\{I_i\}_{i=1}^N$. 所提算法可以应用于灰度图像和彩色图像, 下面给出所提方法的学习过程.

1) 第1层(滤波器学习). 从第 i 幅图像中提取大小为 $k_1 \times k_2$ 的互相有重叠的图像块, 即 $\mathbf{x}_{i,1}, \mathbf{x}_{i,2},$

值的特征向量. 这些卷积滤波器称为面向人脸欺诈检测任务的卷积滤波器, 可以表示为

$$\mathbf{V}_l^1 = \text{mat}_{k_1, k_2}(\mathbf{p}_l(\mathbf{X}\mathbf{X}^T)) \in \mathbb{R}^{k_1 \times k_2}, l = 1, 2, \cdots, L_1 \quad (9)$$

其中: $\text{mat}_{k_1, k_2}(\mathbf{v})$ 是一个函数, 将向量 $\mathbf{v} \in \mathbb{R}^{k_1 k_2}$ 映射为矩阵 $\mathbf{V} \in \mathbb{R}^{k_1 \times k_2}$, $\mathbf{p}_l(\mathbf{X}\mathbf{X}^T)$ 表示提取矩阵 $\mathbf{X}\mathbf{X}^T$ 第 l 个最小的特征向量.

面向人脸欺诈检测任务的卷积滤波器可以提取人脸图像中的“高频”特征, 并进行欺诈检测. 将滤

滤波器和训练图像进行卷积,第1层的第 l 个滤波器的输出可以表示为

$$\mathbf{I}_i^l = \mathbf{I}_i * \mathbf{V}_i^l, i = 1, 2, \dots, N \quad (10)$$

其中 $*$ 表示卷积运算.

2) 第2层和输出层(滤波器学习和后处理).

与第1层类似,从 \mathbf{I}_i^l 中提取互相有重叠的图像块,并构成 $\mathbf{Y}_i^l = [\mathbf{y}_{i,l,1}, \mathbf{y}_{i,l,2}, \dots, \mathbf{y}_{i,l,\tilde{m}'\tilde{n}'}] \in \mathbb{R}^{k_1 k_2 \times \tilde{m}'\tilde{n}'}$,其中 $\mathbf{y}_{i,l,j}$ 表示第 j 个向量化的图像块.将第 l 个滤波器输出的所有矩阵进行拼接,得到 $\mathbf{Y}^l = [\mathbf{Y}_1^l, \mathbf{Y}_2^l, \dots, \mathbf{Y}_N^l] \in \mathbb{R}^{k_1 k_2 \times N\tilde{m}'\tilde{n}'}$.将所有滤波器的所有矩阵进行拼接,得到

$$\mathbf{Y} = [\mathbf{Y}^1, \mathbf{Y}^2, \dots, \mathbf{Y}^{L_1}] \in \mathbb{R}^{k_1 k_2 \times L_1 N\tilde{m}'\tilde{n}'} \quad (11)$$

第2层的面向人脸欺诈检测任务的卷积滤波器可以表示为

$$\mathbf{V}_k^2 = \text{mat}_{k_1, k_2}(p_k(\mathbf{Y}\mathbf{Y}^T)) \in \mathbb{R}^{k_1 \times k_2}, k = 1, 2, \dots, L_2 \quad (12)$$

对于第2层的每一幅输入图像 \mathbf{I}_i^l ,都有 L_2 幅输出图像,大小为 $m \times n$,可以表示为

$$\mathbf{O}_i^l = \{\mathbf{I}_i^l * \mathbf{V}_k^2\}_{k=1}^{L_2} \quad (13)$$

第2层输出图像的个数为 $L_1 L_2$.每一幅输入图像 \mathbf{I}_i^l 在第2层都有 L_2 幅输出图像 $\{\mathbf{I}_i^l * \mathbf{V}_k^2\}_{k=1}^{L_2}$.将这些图像二值化,即 $\{H(\mathbf{I}_i^l * \mathbf{V}_k^2)\}_{k=1}^{L_2}$.将 \mathbf{O}_i^l 中的 L_2 幅输出图像转换为1幅图像:

$$\mathbf{T}_i^l = \sum_{k=1}^{L_2} 2^{k-1} H(\mathbf{I}_i^l * \mathbf{V}_k^2) \quad (14)$$

上述编码过程保留了特征的完整性.在PCANet^[1]中, \mathbf{T}_i^l 中的每幅图像都被分为 B 个块.在所提算法中,将这一步省去了,图像 \mathbf{T}_i^l 不被分块.计算图像 \mathbf{T}_i^l 的直方图,表示为 $\text{Hist}(\mathbf{T}_i^l)$.输入图像 \mathbf{I}_i^l 的特征表示为直方图的集合,即

$$\mathbf{f}_i = [\text{Hist}(\mathbf{T}_i^1), \dots, \text{Hist}(\mathbf{T}_i^{L_1})]^T \in \mathbb{R}^{(2^{L_2})L_1} \quad (15)$$

对于彩色输入图像,将输入图像的3个通道看作3幅灰度图像,其中每一幅作为前向学习网络的输入图像.构建3个前向学习网络,分别处理这3幅图像.在输出层,将这3个前向学习网络的特征拼接成彩色图像最终的特征.在得到输入图像的特征 $\{\mathbf{f}_i\}_{i=1}^N$ 后,结合 $\{\mathbf{I}_i\}_{i=1}^N$,可以训练一个支持向量机或一个单层全连接的神经网络,从而得到一个分类器.在测试时,通过前向学习网络提取特征,然后对分类器进行分类,即可得到预测的标签.

在统计直方图时,没有对输出图像做分块统计,

这里给出一个解释.在PCANet^[1]中,任务是图像识别.因为属于不同类的图像在不同的位置可能有相似的特征,所以需要考虑特征的位置.例如,在手写数字识别的分类中,数字“6”和“9”在不同位置有类似的特征.通过考虑特征的位置,可以将这些图像分开.然而,在人脸欺诈检测任务中,特征的位置反而是不重要的.整幅图像中都存在重复的欺诈特征,从而去除了分块步骤,将整幅图像进行直方图的统计.

3 实验结果与分析

3.1 数据集

在实验中使用3个人脸欺诈检测数据集,即CASIA-FASD数据集^[15]、Idiap Replay-Attack数据集^[16]和OULU-NPU数据集^[17].这3个数据集经常作为标准数据集用于衡量不同算法的性能.这些数据集提供了不同光照条件下不同质量的视频.在CASIA-FASD数据集中有50个人,每个人录了12段视频,共有600段视频,其攻击类型包括弯曲照片攻击、剪切照片攻击和视频攻击. CASIA-FASD数据集中视频的总帧数为111 166.在Replay-Attack数据集中有50个人,每个人录了26段视频,共有1 300段视频. Replay-Attack数据集中的攻击类型包括数字照片攻击、打印照片攻击和视频攻击,视频的总帧数为347 498.在OULU-NPU数据集中有55个人,每个人录了90段视频,共有4 950段视频,攻击类型包括打印照片攻击和视频攻击,视频的总帧数为661 905.

3.2 数据集预处理

在实验中,考虑了少样本人脸欺诈检测场景,也就是训练样本的数量较少,并且没有其他额外的数据.少样本欺诈检测是实际中经常遇到的,此时收集大量的真实样本和欺诈样本的成本较高.为了衡量所提算法相比其他算法在少样本欺诈检测场合下的性能,从每段视频的前几秒中提取视频帧.对于CASIA-FASD数据集,提取每段视频的前1/10,共提取了10 824帧视频.对于Replay-Attack数据集,提取每段视频的前1/10,共提取了34 599帧视频.对于OULU-NPU数据集,提取每段视频的前1/20,共提取了36 038帧视频.对于每个人,假设录像时间远远短于原来的录像时间.对于CASIA-FASD和Replay-Attack数据集,将录像时间缩短至原来的1/10.对于OULU-NPU数据集,

将录像时间缩短至原来的 1/20, 每个人的录像时间在 5 s 左右, 现在只使用大约前 1/4 s 中的帧. 对于 CASIA-FASD 和 Replay-Attack 数据集, 每个人的录像时间在 10 s 左右, 现在只使用大约前 1 s 中的帧. 在实际中, 录制大量不同人、不同类型的长时间视频的成本较高. 表 1 中给出了每个数据集的详细情况. 经过人脸检测^[18]后, 将检测的人脸区域调整为一幅 200 × 200 的图像.

表 1 实验中使用的欺诈检测数据集的详细情况

数据集	人数	视频 总数	总帧数	使用的 帧数	比例
CASIA-FASD ^[15]	50	600	111 166	10 824	1/10
Replay-Attack ^[16]	50	1 300	347 498	34 599	1/10
OULU-NPU ^[17]	55	4 950	661 905	36 038	1/20

3.3 训练和测试的设置

在实验中考虑了跨攻击类型的场合. 跨攻击类型实验是指训练时只见过某一种攻击类型, 没有见过其他的攻击类型. CASIA-FASD 中的攻击图片分为弯曲照片攻击集 C_{wp} 、剪切照片攻击集 C_{cp} 和视频

攻击集 C_v . Replay-Attack 中的攻击图片分为视频攻击集 R_v 、数字照片攻击集 R_{dp} 和打印照片攻击集 R_{pp} . OULU-NPU 中的攻击图片分为打印照片攻击集 O_p 和视频攻击集 O_v . 对于 CASIA-FASD, 考虑了 3 类跨攻击类型的实验, 分别是 $C_v \rightarrow (C_{cp}, C_{wp})$ 、 $C_{cp} \rightarrow (C_v, C_{wp})$ 和 $C_{wp} \rightarrow (C_v, C_{cp})$. 对于 Replay-Attack, 考虑了 3 类跨攻击类型的实验, 分别是 $R_v \rightarrow (R_{dp}, R_{pp})$ 、 $R_{dp} \rightarrow (R_v, R_{pp})$ 和 $R_{pp} \rightarrow (R_v, R_{dp})$. 对于 OULU-NPU, 考虑了 2 类跨攻击类型的实验, 分别是 $O_p \rightarrow O_v$ 和 $O_v \rightarrow O_p$. 在训练时, 选择 3 个人作为训练, 其他人作为测试, 用于训练和测试的人之间没有重叠. 对该实验重复了 10 次, 记录了评价指标的平均值和标准差. 在实验中使用了 3 个评价指标, 分别是攻击人脸分类错误率 (APCER, attack presentation classification error rate)、真实人脸分类错误率 (BPCER, bonafide presentation classification error rate) 和半总错误率 (HTER, half total error rate). 在人脸欺诈检测实验中, 经常使用上述指标. 将帧级别的测试准确率按照视频进行平均, 得到视频级别的准确率. 表 2 ~ 表 4 中给出的是视频级别的测试准确率.

表 2 CASIA-FASD 上的跨攻击类型实验各算法的分类错误率 %

方法	APCER	BPCER	HTER	APCER	BPCER	HTER	APCER	BPCER	HTER
所提算法 (HSV, 6)	24.27 ± 15.08	8.99 ± 7.88	16.63 ± 4.91	31.49 ± 17.48	6.98 ± 5.43	19.23 ± 7.28	42.59 ± 17.45	5.54 ± 6.28	24.06 ± 6.06
所提算法 (HSV, 7)	21.31 ± 16.32	13.17 ± 14.30	17.24 ± 5.97	28.87 ± 15.47	6.26 ± 4.56	17.57 ± 7.21	48.40 ± 21.38	6.76 ± 7.85	27.58 ± 7.57
所提算法 (HSV, 整幅, 8)	21.93 ± 13.79	9.71 ± 10.32	15.82 ± 4.88	27.11 ± 14.04	6.19 ± 4.14	16.65 ± 6.19	50.12 ± 20.43	7.19 ± 8.12	28.66 ± 7.13
所提算法 (HSV, 非重叠块)	23.80 ± 13.47	17.41 ± 10.08	20.60 ± 4.93	40.84 ± 19.97	12.80 ± 11.34	26.82 ± 5.77	59.24 ± 14.36	11.15 ± 7.65	35.20 ± 6.22
所提算法 (YCbCr, 6)	26.00 ± 18.76	16.98 ± 18.87	21.49 ± 6.04	30.61 ± 18.23	12.66 ± 11.74	21.64 ± 8.44	49.10 ± 20.25	8.20 ± 8.62	28.65 ± 6.55
所提算法 (YCbCr, 7)	21.86 ± 18.13	18.42 ± 16.19	20.14 ± 6.03	28.73 ± 19.89	10.65 ± 7.98	19.69 ± 9.17	44.60 ± 19.69	7.99 ± 6.45	26.29 ± 6.89
所提算法 (YCbCr, 8)	23.08 ± 16.16	16.55 ± 14.94	19.81 ± 5.82	28.59 ± 17.03	10.43 ± 7.25	19.51 ± 7.39	47.13 ± 19.29	7.48 ± 6.48	27.30 ± 6.93
CTA (HSV) ^[3]	25.46 ± 18.72	16.62 ± 16.30	21.04 ± 7.35	34.38 ± 15.18	12.66 ± 10.65	23.52 ± 5.57	52.51 ± 19.67	7.63 ± 5.76	30.07 ± 8.07
CTA (YCbCr) ^[3]	37.65 ± 20.52	14.68 ± 13.55	26.17 ± 6.31	53.01 ± 12.19	7.99 ± 6.65	30.50 ± 5.50	60.19 ± 23.22	8.71 ± 10.23	34.45 ± 8.41
IQM ^[6]	78.67 ± 39.66	18.78 ± 38.65	48.72 ± 1.68	8.42 ± 26.62	90.79 ± 29.12	49.60 ± 1.25	66.12 ± 42.33	30.72 ± 40.42	48.42 ± 1.90
IDA ^[7]	49.98 ± 17.45	28.06 ± 11.57	39.02 ± 5.81	58.77 ± 14.47	26.62 ± 12.03	42.69 ± 5.27	80.67 ± 10.53	9.86 ± 5.22	45.27 ± 3.96
CNN (FC) ^[8]	44.62 ± 31.53	38.23 ± 27.49	41.43 ± 5.64	45.49 ± 32.05	39.34 ± 28.49	42.42 ± 4.23	69.51 ± 27.31	27.04 ± 24.25	48.27 ± 2.68
CNN (GAP) ^[8]	58.56 ± 14.39	33.35 ± 10.32	45.96 ± 2.32	48.03 ± 8.73	40.70 ± 6.54	44.36 ± 2.42	50.00 ± 16.46	41.71 ± 11.94	45.85 ± 2.93
CNN (BS) ^[12]	57.05 ± 11.83	38.88 ± 7.12	47.96 ± 3.90	55.46 ± 9.75	38.80 ± 7.18	47.13 ± 2.26	56.80 ± 15.49	35.53 ± 12.38	46.16 ± 2.56

3.4 跨攻击类型的实验

在实验中, 将所提算法与 7 种人脸欺诈检测算法进行了比较. 比较的算法包括基于 HSV 色彩空间颜色纹理分析的检测算法^[3] (CTA (HSV), color texture analysis in the HSV color space)、基于 YCbCr

色彩空间颜色纹理分析的检测算法^[3] (CTA (YCbCr), color texture analysis in the YCbCr color space)、基于图像质量评价的检测算法^[6] (IQM, image quality analysis)、基于图像失真分析的检测算法^[7] (IDA, image distortion analysis)、基于全连接层和全

表 3 Replay-Attack 上的跨攻击类型实验各算法的分类错误率

%

方法	APCER	BPCER	HTER	APCER	BPCER	HTER	APCER	BPCER	HTER
所提算法(HSV)	40.37 ± 11.21	12.73 ± 12.84	26.55 ± 7.14	16.72 ± 10.86	16.31 ± 16.69	16.52 ± 7.23	85.07 ± 8.12	3.19 ± 2.75	44.13 ± 3.01
所提算法(YCbCr)	26.11 ± 11.71	16.95 ± 10.91	21.53 ± 6.82	20.05 ± 14.73	18.33 ± 16.08	19.19 ± 7.15	86.15 ± 6.07	1.92 ± 1.55	44.03 ± 2.80
CTA(HSV) ^[3]	45.45 ± 8.49	15.32 ± 7.42	30.38 ± 5.38	33.35 ± 14.75	18.44 ± 13.21	25.89 ± 5.91	91.19 ± 7.33	6.28 ± 5.23	48.73 ± 1.28
CTA(YCbCr) ^[3]	52.93 ± 6.67	15.99 ± 9.17	34.47 ± 3.79	20.51 ± 11.55	25.39 ± 14.16	22.95 ± 7.27	97.08 ± 3.64	1.67 ± 2.06	49.37 ± 1.01
IQM ^[6]	52.39 ± 46.97	33.62 ± 43.73	43.00 ± 8.14	31.21 ± 35.80	42.66 ± 31.15	36.93 ± 6.13	55.00 ± 41.77	29.47 ± 40.97	42.23 ± 6.74
IDA ^[7]	31.16 ± 15.25	8.33 ± 6.43	19.75 ± 8.90	16.70 ± 9.03	13.05 ± 7.46	14.88 ± 4.27	48.79 ± 18.38	5.99 ± 4.88	27.39 ± 7.49
CNN(FC) ^[8]	2.02 ± 4.36	96.78 ± 7.36	49.40 ± 2.29	7.61 ± 9.22	68.79 ± 27.89	38.20 ± 9.92	78.82 ± 13.72	13.97 ± 12.43	46.40 ± 6.02
CNN(GAP) ^[8]	18.31 ± 10.97	87.42 ± 12.91	52.87 ± 2.06	14.82 ± 7.81	94.87 ± 7.17	54.84 ± 1.50	44.98 ± 11.65	53.06 ± 22.62	49.02 ± 5.93
CNN(BS) ^[12]	7.59 ± 5.56	89.29 ± 10.33	48.44 ± 2.62	4.83 ± 1.63	95.95 ± 1.46	50.39 ± 0.66	60.83 ± 17.92	28.18 ± 17.72	44.51 ± 3.93

表 4 OULU-NPU 上的跨攻击类型实验各算法的分类错误率

%

方法	APCER	BPCER	HTER	APCER	BPCER	HTER
所提算法(HSV)	29.77 ± 8.60	10.94 ± 5.33	20.36 ± 2.53	29.40 ± 10.15	11.10 ± 4.71	20.25 ± 3.59
所提算法(YCbCr)	28.97 ± 9.01	11.02 ± 5.08	19.99 ± 2.64	29.32 ± 8.84	12.14 ± 6.75	20.73 ± 2.93
CTA(HSV) ^[3]	22.68 ± 12.55	12.56 ± 9.22	17.62 ± 4.18	23.13 ± 8.58	10.18 ± 4.29	16.65 ± 3.44
CTA(YCbCr) ^[3]	23.85 ± 8.62	11.51 ± 5.45	17.68 ± 2.43	26.52 ± 7.06	12.53 ± 5.47	19.53 ± 3.20
IQM ^[6]	34.83 ± 45.40	58.77 ± 46.79	46.80 ± 5.00	26.18 ± 34.83	55.52 ± 37.89	40.85 ± 5.70
IDA ^[7]	43.69 ± 11.01	27.58 ± 12.12	35.63 ± 2.32	36.88 ± 10.16	33.82 ± 10.45	35.35 ± 3.48
CNN(FC) ^[8]	39.67 ± 12.31	37.75 ± 12.25	38.71 ± 5.25	25.45 ± 13.05	35.66 ± 9.54	30.55 ± 4.06
CNN(GAP) ^[8]	44.55 ± 4.21	55.13 ± 6.14	49.84 ± 1.20	46.88 ± 4.71	47.41 ± 4.56	47.15 ± 0.62
CNN(BS) ^[12]	47.14 ± 3.27	51.89 ± 5.38	49.52 ± 1.42	46.59 ± 4.31	48.12 ± 3.65	47.36 ± 1.27

局监督的 CNN^[8] (CNN(FC), CNN with fully connected layer)、基于全局均值池化和全局监督的 CNN^[8] (CNN(GAP), CNN with global average pooling layer) 和基于深度逐像素二值监督的 CNN^[12] (CNN(BS), CNN with binary pixel supervision). CNN 的结构取自文献[11]中的结构. 在训练时, 训练样本只包含视频中 3 个人的某一种攻击类型的图像, 所有的算法都是在这一设定下进行的比较. 将输入图像转换至 HSV 和 YCbCr 色彩空间, 然后使用所提算法进行分类. 所提算法的参数设置如下: 卷积滤波器的大小为 5 × 5, 图像块之间的步长间隔是 3 × 3, 在 2 层中都使用了 8 个卷积滤波器. 在 CASIA-FASD 上的实验结果 $C_{wp} \rightarrow (C_v, C_{cp})$ 、 $C_{cp} \rightarrow (C_{wp}, C_v)$ 和 $C_v \rightarrow (C_{wp}, C_{cp})$ 见表 2 第 2 ~ 4 列、第 5 ~ 7 列和第 8 ~ 10 列. 在 Replay-Attack 上的实验结果 $R_v \rightarrow (R_{dp}, R_{pp})$ 、 $R_{dp} \rightarrow (R_v, R_{pp})$ 和 $R_{pp} \rightarrow (R_v, R_{dp})$ 见表 3 的第 2 ~ 4 列、第 5 ~ 7 列和第 8 ~ 10 列. 在 OULU-NPU 上的实验结果 $O_p \rightarrow O_v$ 和 $O_v \rightarrow O_p$ 见表 4 的第 2 ~ 4 列和第 5 ~ 7 列. 从 CASIA-FASD 上的实

验结果可以看出, 所提算法得到了比其他算法更好的结果. 在 3 个实验中, 所提算法(HSV)的 HTER 是最低的. 比较所提算法(HSV)和 CTA(HSV)的结果可以看出, 所提算法显著降低了 HTER. 在弯曲照片攻击实验中, HTER 降低了 5.22%; 在剪切照片攻击实验中, HTER 降低了 6.87%; 在视频攻击实验中, HTER 降低了 1.41%. 比较所提算法(YCbCr)和 CTA(YCbCr)可以看出, 所提算法得到了更好的结果. 在弯曲照片攻击实验中, HTER 降低了 6.36%; 在剪切照片攻击实验中, HTER 降低了 10.99%; 在视频攻击实验中, HTER 降低了 7.15%. 在 CASIA-FASD 上, 其他方法取得了较差的实验结果, IQM 方法和 IDA 方法的 HTER 较高, 基于深度学习方法的结果也不理想. 上述实验结果说明, 在欺诈检测中, 采用前向学习的方式比采用固定的特征提取算子和基于深度学习的方法好. 从 Replay-Attack 的实验结果可以看出, IDA 得到了最低的 HTER, 所提算法与 IDA 方法的结果接近. 比较所提算法(HSV)和 CTA(HSV)可以看出, 所提算法的结

果更好. 在视频攻击、数字照片攻击、打印照片攻击的实验中, HTER 分别降低了 3.83%、9.37% 和 4.6%. 另外, 所提算法(YCbCr)比 CTA(YCbCr)的结果更好. 在视频攻击实验中, HTER 降低了 12.94%; 在数字照片攻击实验中, HTER 降低了 3.76%; 在打印照片攻击实验中, HTER 降低了 5.34%. 在 Replay-Attack 上, 其他方法的结果较差, IQM 以及基于深度学习方法的 HTER 较高. 上述实验结果再次说明了采用前向学习的方式比采用固定特征提取算子和基于深度学习的方法更好. 从 OULU-NPU 的实验结果来看, CTA(HSV)取得了最低的 HTER, 所提算法与 CTA 的结果类似, 并且远好于其他的方法. 通过仔细分析实验结果可以看出, 所提算法在不同颜色空间中的结果是类似的, 这说明所提算法适应于不同颜色空间的输入图像. 与之相对的, CTA 在不同颜色空间的图像上有着很大的性能差异. 总的来说, 在跨攻击类型欺诈检测任务中, 所提算法在不同数据集上都有较好的表现, 性能优于其他算法.

基于深度学习的人脸欺诈检测方法需要有大量样本才能够训练成功, 当训练样本数量不够时, 很容易出现过拟合现象. 在文中的实验部分, 基于深度学习的方法在训练集上都有不错的分类性能, 但是在测试集上表现不佳. 例如, 在表 3 第 2~4 列中的 CNN(FC)方法的分类有明显的偏向性, 对于欺诈人脸的分类正确率很高, 对于真实人脸分类的正确率很低, 也即网络倾向于将大部分测试样本分类为欺诈人脸, 同样的现象也出现在表 3 第 2~4 列中的 CNN(GAP)和 CNN(BS)方法中以及表 3 第 5~10 列中的 CNN 方法中. 基于人工特征的方法有时也会出现这种现象. 例如, 在表 2 第 2~7 列中的 IQM 方法和表 2 第 8~10 列中的 IDA 方法中, 分类也有明显的偏向性. 造成上述问题的原因是, IQM 提取的特征不够有区分性, 且特征维度只有 18 维, 维度较低. 而所提方法在所有实验场景下均没有出现分类偏向性问题, 说明所提方法有良好的人脸欺诈检测性能.

下面给出所提算法中一些参数和步骤的变化对性能的影响分析. 首先分析当采用非重叠块直方图统计步骤时, 算法性能的变化. 在所提算法的直方图统计步骤中, 将 200×200 的图像分为 50×50 的非重叠块用于统计直方图, 相应的在 CASIA-FASD 上的结果见表 2(表中的所提算法(HSV, 非重叠块)

记录了结果). 从实验结果可以看出, 如果增加了非重叠块直方图的统计步骤, 不但不能提高所提算法的性能, 反而会在一定程度上降低所提算法的性能. 造成分类性能下降的原因正如前文中分析的, 如果对图像进行分块, 将会把空间信息考虑在最后的特征中, 而欺诈检测需要判断图像整体的性质, 不利用空间信息反而是有利的. 如果是一般的图像分类任务, 如人脸识别、手写数字识别等, 则需要利用空间信息.

下面分析采用不同滤波器数量时, 所提算法性能的变化. 所提算法取 6 和 7 个滤波器时, 在 CASIA-FASD 数据集上的实验结果见表 2. 从实验结果可知, 当采用 8 个滤波器并且色彩空间取为 HSV 时, 所提算法在 $C_{wp} \rightarrow (C_v, C_{cp})$ 和 $C_{cp} \rightarrow (C_{wp}, C_v)$ 实验上取得了最低的 HTER. 具体而言, 在 HSV 色彩空间中, 当采用更少的滤波器时, 所提算法的欺诈检测性能有所降低. 在 YCbCr 色彩空间中出现了同样的现象. 在 $C_v \rightarrow (C_{wp}, C_{cp})$ 实验上的结果与前面 2 个实验不同. 当采用 6 个滤波器, 且色彩空间取 HSV 时, 所提算法取得了最低的 HTER. 具体而言, 在 HSV 色彩空间中, 当增加滤波器数量时, 所提算法欺诈检测性能下降. 在 YCbCr 空间中, 取 7 个滤波器时性能最优. 从这组实验结果中可以得出如下结论: 当采用不同的滤波器数量时, 会对算法性能有一定影响, 但是影响有限. 取不同滤波器数量时, 所提算法的分类性能比较接近. 而对于某些跨攻击类型的实验, 如当训练时见过的是视频攻击, 测试时是照片攻击时, 则采用 6~7 个滤波器时效果更好, 说明此时欺诈检测信息更集中于高频.

这里比较一下所提算法和 CTA 算法的特征维数. CTA 算法的特征维数是 20 001. 当使用 8 个滤波器, 并用于彩色输入图像时, 所提算法的特征维数是 6 144, 远低于 CTA 算法. 然而, 在上述欺诈检测实验中, 所提算法在 CASIA-FASD 和 Replay-Attack 数据集上比 CTA 算法更好, 在 OULU-NPU 数据集上与 CTA 算法接近.

4 结束语

笔者提出了一种新的人脸欺诈检测算法, 对前向学习网络的学习过程做了改进. 使用面向人脸欺诈检测任务的卷积滤波器(这些滤波器是对应于最小特征值的特征向量). 可以有效地提取图像中的高频信息并进行欺诈检测. 与现有算法相比, 所提

算法有很多优势. 所提算法与基于人工特征的算法相比, 不需要利用特定领域知识设计特征提取算子(特征提取算子是从数据中学得的). 所提算法与基于深度学习的方法相比, 可以应用于少样本欺诈检测场合. 在实际应用中, 录制大量人员的真实人脸和欺诈人脸的视频较难, 所以所提算法可以应用于那些深度学习方法效果不好的场合中.

在实验结果部分, 采用了少样本人脸欺诈检测场景, 使用了 CASIA-FASD^[15]、Replay-Attack^[16] 和 OULU-NPU^[17] 3 个数据集. 使用了每个视频中前几秒的数据, 并进行了跨攻击类型的实验, 结果表明, 所提算法比现有算法的性能更好. 在 CASIA-FASD 和 Replay-Attack 数据集上, 与 CTA 算法相比, 所提算法的 HTER 更低.

参考文献:

- [1] Chan T H, Jia K, Cao S, et al. PCANet: a simple deep learning baseline for image classification? [J]. IEEE Transactions on Image Processing, 2015, 24(12): 5017-5032.
- [2] Boulkenafet Z, Komulainen J, Hadid A. Face anti-spoofing based on color texture analysis [C] // International Conference on Biometrics. Phuket: IEEE, 2015: 2636-2640.
- [3] Boulkenafet Z, Komulainen J, Hadid A. Face spoofing detection using color texture analysis[J]. IEEE Transactions on Information Forensics and Security, 2016, 11(8): 1818-1830.
- [4] Boulkenafet Z, Komulainen J, Hadid A. Face antispoofing using speeded-up robust features and fish vector encoding[J]. IEEE Signal Processing Letters, 2017, 24(2): 141-145.
- [5] Garcia D C, Queiroz R L. Face-spoofing 2D-detection based on moire-pattern analysis[J]. IEEE Transactions on Information Forensics and Security, 2015, 10(4): 778-786.
- [6] Galbally J, Marcel S, Fierrez J. Image quality assessment for fake biometric detection: application to iris, fingerprint, and face recognition[J]. IEEE Transactions on Image Processing, 2014, 23(2): 710-724.
- [7] Wen D, Han H, Jain A K. Face spoof detection with image distortion analysis[J]. IEEE Transactions on Information Forensics and Security, 2015, 10(4): 746-761.
- [8] Rehman Y A U, Po L M, Liu M. Deep learning for face anti-spoofing: an end-to-end approach [C] // IEEE Conference on Signal Processing. Poznan: IEEE, 2017: 195-200.
- [9] Nagpal C, Dubey S R. A performance evaluation of convolutional neural networks for face anti spoofing [J/OL]. [2019-03-27]. <https://arxiv.org/abs/1805.04176>.
- [10] Atoum Y, Liu Y, Jourabloo A, et al. Face anti-spoofing using patch and depth-based CNNs [C] // IEEE International Joint Conference on Biometrics. Denver: IEEE, 2017: 319-328.
- [11] Liu Y, Jourabloo A, Liu X. Learning deep models for face anti-spoofing: binary or auxiliary supervision [C] // IEEE Conference on Computer Vision and Pattern Recognition. Salt Lake City: IEEE, 2018: 389-398.
- [12] George A, Marcel S. Deep pixel-wise binary supervision for face presentation attack detection [J/OL]. [2019-07-09]. <https://arxiv.org/abs/1907.04047>.
- [13] Li L, Feng X, Jiang X, et al. Face anti-spoofing via deep local binary patterns [C] // International Conference on Image Processing. Beijing: IEEE, 2017: 101-105.
- [14] Li H, Li W, Cao H, et al. Unsupervised domain adaptation for face anti-spoofing [J]. IEEE Transactions on Information Forensics and Security, 2018, 13(7): 1794-1809.
- [15] Zhang Z, Yan J, Liu S, et al. A face antispoofing database with diverse attacks [C] // International Conference on Biometrics. New Delhi: IEEE, 2012: 26-31.
- [16] Chingovska I, Anjos A, Marcel S. On the effectiveness of local binary patterns in face anti-spoofing [C] // IEEE International Conference of Biometrics Special Interest Group. Darmstadt: IEEE, 2012: 183-194.
- [17] Boulkenafet Z, Komulainen J, Li L, et al. OULU-NPU: a mobile face presentation attack database with real-world variations [C] // IEEE International Conference on Automatic Face and Gesture Recognition. Washington, DC: IEEE, 2017: 612-618.
- [18] Viola P, Jones M. Rapid object detection using a boosted cascade of simple features [C] // IEEE Conference on Computer Vision and Pattern Recognition. Kauai: IEEE, 2001: 511-518.