基于异构哈希网络的跨模态人脸检索方法

董 震 裴明涛

(北京理工大学智能信息技术北京市重点实验室 北京 100081)(北京理工大学计算机学院 北京 100081)

摘 要 该文提出一种基于异构哈希网络的跨模态人脸检索方法.异构哈希网络能够将位于不同空间的人脸图像和人脸视频映射到一个公共且有判别力的二值空间上,以获得有效的二值哈希表示.该网络包含图像分支、视频分支和哈希函数三个部分,首先图像和视频分支分别将人脸图像和人脸视频映射到一个公共空间,然后在公共空间中学习非线性哈希函数.网络的训练使用了三种损失函数:Fisher 损失、softmax 损失和三元排序损失(triplet ranking loss),其中的 Fisher 损失关注于公共空间的判别力,softmax 损失强调公共空间上表达的可分性,三元排序损失旨在提升最终的检索性能.在多个人脸视频数据集上的跨模态人检索实验结果表明了所提出方法的有效性.

关键词 异构哈希网络;跨模态;人脸检索;深度学习;损失函数 中图法分类号 TP391 **DOI**号 10.11897/SP.J.1016.2019.00073

Cross-Modality Face Retrieval Based on Heterogeneous Hashing Network

DONG Zhen PEI Ming-Tao

(Beijing Laboratory of Intelligent Information Technology, Beijing Institute of Technology, Beijing 100081) (School of Computer Science, Beijing Institute of Technology, Beijing 100081)

Abstract Cross-modality face retrieval is to retrieve faces of a particular person in one modality given his/her face information in another modality, such as retrieving video shots containing particular person given one image of him/her (query-by-image video retrieval), or retrieving the face images of one person by using his/her video clip as query (query-by-video image retrieval). It is an important problem in computer vision with wide range of applications. Take the criminal investment as an example. The "query-by-image video retrieval" task plays an important role in rapid locating and tracking a suspect from masses of surveillance videos with the ID card, passport, or driver license photo of the suspect as query. The "query-by-video image retrieval" task helps to determine the identity of an unknown suspect by retrieving a huge mug-shot image database given his/her video shot taken by the surveillance cameras in the crime scene. In this paper, we present a cross-modality face retrieval method which uses the heterogeneous hashing network to generate effective and compact hash representations for both face images and face videos. The network contains an image branch and a video branch to project face images and videos into a common discriminative space, respectively. Each channel are equipped with two modules: feature extractor module and non-linear mapping module. The feature extractor modules aim to represent face images or videos via appropriate features, and the non-linear mapping modules are designed to transform the heterogenous image and video feature spaces into the common space. On the common space, the similarity between a face image and a face video can be measured through the distance of their corresponding discriminative features, but these features are still

收稿日期:2017-05-04;在线出版日期:2018-07-11.本课题得到国家自然科学基金(61472038)资助. 董 震,男,1990年生,博士,主要研究方向为计算机视觉、人脸识别、人脸检索. E-mail: dongzhen@bit.edu.cn. 裴明涛(通信作者),男,1977年生,博士,副教授,中国计算机 学会(CCF)会员,主要研究方向为计算机视觉、人脸识别. E-mail: peimt@bit.edu.cn.

high-dimensional vectors of floating point numbers, which cannot satisfy the requirements of low computation and storage complexities in the retrieval task. The non-linear hash functions are thus learned in the common space to obtain the corresponding binary hash representations. To catch the compatability and the effectiveness of the branches and the hash functions, the heterogeneous hashing network is trained with three loss functions. Fisher loss, softmax loss, and triplet ranking loss. Our Fisher loss uses the difference form of the inter-class and the intra-class scatter where the mean vectors are learnable variables, which is feasible for the mini-batch based optimization method. The Fisher loss and the softmax loss are jointly exploited to enhance the discriminative power of the common space. The triplet ranking loss is enforced to the final binary space for the improvement of the retrieval performance. Experiments on a large-scale face video dataset and two challenging TV-series datasets demonstrate the effectiveness of the proposed method. The contributions of the paper are three-folds: (1) We propose an effective cross-modality face retrieval method based on the heterogeneous hashing network. Our network is able to generate isomorphic discriminative compact binary representations of both face images and videos. (2) The proposed heterogeneous hashing network provides a general framework for deep learning based cross-modality hashing methods, and can be easily adopted in many other cross-modality retrieval tasks. (3) The proposed method achieves excellent results of face retrieval across image and video modalities on a large scale face video dataset and two challenging TV-series datasets.

Keywords heterogeneous hashing network; cross-modality; face retrieval; deep learning; loss function

1 引 言

现实世界中采集到的人脸数据可能存在多种模 态,按照载体类型,可以分为人脸图像和人脸视频; 按照具体表现形式,可以分为照片,肖像画,素描画, 漫画等;按照采集设备类型,可以分为可见光人脸数 据和红外人脸数据等.给定一个特定人物某种模态 的人脸信息,跨模态人脸检索任务是指从包含另一 种模态人脸信息的数据库中找寻到该人,例如给定 某人的人脸图像,在人脸视频数据库中检索到包含 该人的所有人脸视频(基于图像的人脸视频检索); 或者根据某人的人脸视频,在人脸图像数据库中检 索含有该人的所有人脸图像(基于视频的人脸图像 检索).伴随着互联网和自媒体技术的不断发展,互 联网上多种模态的人脸数据越来越丰富,这使得跨 模态人脸检索具有非常广泛的应用前景.以罪犯追 踪和身份分析为例,基于图像的人脸视频检索的应 用场景是:根据嫌疑犯的身份证、护照或者驾驶证等 证件上的人脸图像信息,从含有大量监控视频的数 据库中快速定位和追踪该嫌疑犯;基于视频的人脸 图像检索的应用场景是:根据案发现场监控拍摄到 的含有罪犯的视频,对公安机关备案的人脸图像数 据集进行检索以确定罪犯的身份.

以人脸图像和人脸视频为例,跨模态人脸检索 需要度量一张人脸图像和一个人脸视频的相似性. 最直接的度量方法是将人脸视频看成帧数据的无序 集合,利用人脸图像与视频所有帧的相似性的最大 值或平均值作为该人脸图像和视频的相似性的最大 值或平均值作为该人脸图像和视频的相似性.然而, 这种度量方法没有考虑人脸视频各帧之间的相关 性,浪费了宝贵的相关信息,同时在人脸视频包含的 帧数较多时还面临着较大的计算和存储消耗问题. 对一个人脸视频片段中的所有帧联合建模可以充分 利用各帧之间的相关信息,并生成一个与帧数无关 的人脸视频表示,然而,对人脸视频所有帧的联合建 模将导致视频特征和图像特征处于异构的特征空 间,使得跨模态人脸检索任务变得非常困难.

本文构造了一类称为异构哈希网络(Heterogeneous Hashing Network, HHN)的深度神经网络, 将位于异构空间的特征映射到一个公共的且有足够 判别力的二值空间上,实现不同模态数据之间的度 量,从而实现跨模态的人脸检索任务.异构哈希网络 的结构如图1所示.网络包含图像和视频两个分支, 分别用来处理人脸图像数据和人脸视频数据,将它

75

们从异构特征空间映射到一个公共空间来度量相似 性.每个分支包含两个模块:特征提取模块和同构表 示生成模块.特征提取模块完成人脸图像或人脸视 频的特征提取和表示,例如人脸图像的卷积神经网 络特征^[1-2]、人脸视频的 3D 卷积网络特征^[3]或者神 经聚合网络特征[4]等;同构表示生成模块将位于两 个异构空间的人脸图像特征和人脸视频特征映射到 一个公共且有判别力的空间,在公共空间上,人脸图 像和人脸视频之间的相似性可以通过它们在公共空 间上的相应特征表示之间的距离来度量,然而,公共 空间上的特征表示仍然是高维的浮点数向量,不能满 足检索任务对检索算法低时间和空间复杂度的要求. 为了解决这个问题,在公共空间上进一步学习非线性 的哈希函数来获得人脸图像和视频数据的二值哈希 表示.哈希表示的二值属性,可以以极低的时间复杂 度和极少的存储消耗实现人脸图像和视频之间的相 似性计算,从而减少检索过程的时间和空间花费.



Li 等人^[5]提出了一种跨欧式空间和黎曼流形的哈希方法,用于学习人脸图像和视频的统一哈希表示.该方法首先利用核函数将异构特征空间映射 至高维再生希尔伯特空间,然后在再生希尔伯特空间上学习哈希函数进而获得原始数据的哈希表达, 并利用学得的哈希表达直接度量人脸图像和视频数 据之间的距离.尽管该方法取得了较好的检索结果, 但该方法将特征提取和哈希函数的学习视作两个独 立的过程,这可能会导致特征表示与哈希函数的不 兼容,从而对最终的检索性能造成影响.与之不同, 本文方法将特征学习和哈希函数学习融合到一个统 一的学习框架中联合进行,使得两个部分可以相互 促进,从而更好地为最终的检索任务服务.

本文提出的基于异构哈希网络的跨模态人脸检 索方法是首个利用深度神经网络来实现跨模态人脸 检索的方法,所提出的异构哈希网络能够为位于异 构空间的人脸图像和视频数据生成同一空间上的紧 致二值哈希表示;异构哈希网络为跨模态检索问题 提供了一套深度网络的通用框架,本文设计的网络 不仅可以完成跨图像和视频模态的检索任务,还可 以扩展到其他跨模态的检索任务上;在多个人脸视 频数据集上进行的两类跨模态人脸检索实验,包括 基于图像的人脸视频检索和基于视频的人脸图像检 索,均取得了非常好的检索结果.

2 相关工作

2.1 单模态哈希方法

Gionis 等人^[6]提出的局部敏感哈希(LSH)方法 是数据无关的哈希方法的代表. LSH 方法使用随机 投影作为哈希函数将高维特征表示投影到低维二值 空间,哈希函数的设计原则是:原始特征空间中距离 近的点的哈希表示发生碰撞的概率大于距离远的点 的碰撞概率,所有满足这一设计原则的函数构成 LSH 哈希函数族. Weiss 等人^[7]提出的谱哈希方法 期望在原始特征空间中欧式距离近的数据点在二值 空间中海明距离依然近,针对此要求,谱哈希方法将 哈希函数的求解建模为图的分割过程,其中图的顶 点表示数据点,顶点之间边的权值表示数据点间的 欧式距离,之后可利用该图的拉普拉斯矩阵的特征 向量作为求解的哈希函数. Gong 等人^[8]提出的迭代 量化哈希方法首先对中心化后的数据进行 PCA 降 维,然后迭代执行数据集合的旋转和量化误差的最 小化,最终获得最优解. Wang 等人^[9]提出了半监督 哈希方法,该方法在带标签的数据上最小化经验误 差,同时在所有数据上包括有标签数据和无标签数 据最大化哈希码的方差和均衡性,既可以建模度量 信息,也可以建模语义相似性信息.Liu 等人^[10]探究 了二值空间的海明距离与哈希码内积表示的相似性 度量之间的关系,并提出了有监督的核哈希方法,将 目标函数中海明距离巧妙地转化为二值哈希编码的 内积形式,从而可以方便地进行优化.

2.2 多模态哈希方法

单模态的哈希方法在很多领域取得了成功,但 是只能检索同一个模态的数据,而不能处理多模态 数据检索的问题.多模态检索问题在近年来受到了 很大的关注^[11-22].已有的多模态哈希方法可以分为 两类:无监督方法和有监督方法.代表性的无监督多 模态哈希方法包括对应自编码器哈希方法^[11]、跨视 角哈希方法^[12]和可预测对偶视角哈希方法^[13].这三 类方法的学习目标分别是重构误差最小化、基于图 的相似性保持和可预测性的维持.典型的有监督多 模态哈希方法有跨模态相似敏感哈希方法[14]、多模 态隐二值化嵌入[15]、多模态神经网络哈希方法[16]、 参数化局部多模态哈希方法[17]、语义相关最大化哈 希方法[18] 和语义保持哈希方法[19]. 与无监督的多模 态哈希方法相比,有监督的多模态哈希方法能够充 分地利用训练数据的语义信息来降低多模态数据之 间的差异化程度,并减小不同模态数据之间的语义 鸿沟,有监督的多模态哈希方法也因此比无监督的 方法取得了更好的性能. Jiang 和 Li^[21]提出了深度 跨模态哈希方法,并将其应用在了图像和文本之间 的检索问题上. Cao 等人[22] 提出了深度视觉-语义哈 希方法,该方法利用复杂深度神经网络来挖掘文本 数据和图像数据之间的异构相关信息,并利用学习 到的信息提升跨模态检索的性能.这两个工作利用 深度神经网络在跨图像域和文本域的检索问题中取 得了较好的效果.与之不同,本文的异构哈希网络提 供了一个多模态检索的深度学习框架,并基于此框 架实现了人脸的跨图像域和视频域的检索,取得了 很好的检索性能.

3 异构哈希网络

3.1 网络结构

如图 1 所示,异构哈希网络由五部分组成:图像 分支的特征提取模块 $\chi^{I}(\cdot)$ 、图像分支的同构表示生 成模块 $\eta^{I}(\cdot)$ 、视频分支的特征提取模块 $\chi^{V}(\cdot)$ 、视 频分支的同构表示生成模块 $\eta^{V}(\cdot)$ 以及哈希编码生 成模块 $\phi(\cdot)$.这 5 个部分都可以利用神经网络实 现,因此 HHN 建立了一套用于跨模态人脸检索且可 以进行端到端学习的深度网络结构. HHN 的应用范 围不仅仅局限于跨模态人脸检索,只要能够设计出 合适的特征提取网络,HHN 可以很容易的用于图 像与视频之外的其它模态数据的检索任务.

3.2 损失函数

同构表示生成模块将原本属于不同特征空间的 人脸数据映射到同一公共空间.为了增强映射后的 特征表示的判别力,对公共空间添加判别性约束,则 引入 Softmax 损失函数和 Fisher 损失函数,来提升 属于同一个人不同模态的人脸数据之间的关联程 度.同时,在哈希函数的学习过程中利用监督信息对 哈希函数添加约束,即引入三元排序损失函数,使得 同一个人的不同模态人脸数据的特征表示之间的差 异进一步减小. 3.2.1 Softmax 损失函数

为了保证人脸数据在公共空间上的特征表示具 有较强的可分性,引入 softmax 损失:

$$\mathcal{L}_{s} = -\sum_{i=1}^{c} \sum_{K \in \{I,V\}} \sum_{j=1}^{n_{i}^{K}} \log \frac{e^{W_{i}^{\mathrm{T}} r_{i,j}^{K} + b_{i}}}{\sum_{k=1}^{c} e^{W_{k}^{\mathrm{T}} r_{i,j}^{K} + b_{k}}}$$
(1)

其中 W_i 和 b_i 是 softmax 损失的参数, c 是要处理 的个体的总数(类别总数), n_i^K 为第i 个类别中图像 (K=I)或视频(K=V)的数量, $r_{i,j}^K$ 为第i 个个体的 第j 个图像或视频对应的公共空间中的特征表示. 3.2.2 Fisher 损失函数

为了在增强公共空间判别力的同时减小同一个 人的人脸视频与人脸图像数据特征表示的不一致 性,则需要对共同空间添加判别性约束.Fisher 损失 可以在最小化类内差异的同时最大化类间差异,是 一个合适的选择.将人脸图像和视频的公共空间上 的特征放在一起组成一个集合,该集合上的Fisher 损失通过最小化类内散度矩阵 S_W 和类间散度矩阵 S_B 的 Rayleigh 熵(Rayleigh quotient)来实现:

$$\min \frac{tr(\mathbf{S}_{W})}{tr(\mathbf{S}_{B})}$$
(2)

其中 tr(•)表示矩阵的迹.

为了保证 S_W和 S_B的非奇异性以获得最优解, 式(2)所显示的 Rayleigh 熵需要满足特征表示的维 度小于个体数量.然而这在实际应用中通常是难以 满足的,因为用于求解哈希表示的特征向量的维度 通常是比较高的.另外,由于均值向量是在整个训练 集上计算的,Rayleigh 熵形式的 Fisher 损失不能利 用基于批次(batch)的算法进行优化.为了解决上述 的问题,使用差分形式的 Fisher 损失^[23]来替代 Rayleigh 熵形式:

$$tr(\boldsymbol{S}_{W}) - tr(\boldsymbol{S}_{B}) + \|\boldsymbol{R}\|_{F}$$
(3)

其中 || • || _F表示矩阵的 Frobenius 范数,加入此范数 项的目的是保证式(3)所示的目标函数是凸的.进一 步地,为了使差分形式的 Fisher 损失可以利用基于 批次的算法进行优化,将原始 Fisher 损失中作为定 值的均值向量转变为可学习的变量,保持均值向量 的可学习性.这样,在每一次迭代中,均值向量都在 训练数据的每个批次上进行更新.因此,公共空间上 的判别损失表示为

$$\mathcal{L}_{d} = \|\boldsymbol{R}\|_{F}^{2} + \frac{1}{2n} \sum_{i=1}^{c} \sum_{K \in \langle I, V \rangle} \sum_{j=1}^{n_{i}^{K}} \|\boldsymbol{r}_{i,j}^{K} - \boldsymbol{\mu}_{i}\|_{2}^{2} - \frac{1}{2c} \sum_{i=1}^{c} \|\boldsymbol{\mu}_{i} - \boldsymbol{\mu}\|_{2}^{2}$$
(4)

其中 μ_i表示第 i 个个体的所有样本(无论人脸视频 还是图像)在公共空间上特征表示的均值向量,μ表 示所有 c 个个体的所有样本在公共空间上特征表示 的均值向量.

3.2.3 三元排序损失函数

为了学得对检索更有效的二值哈希表示,鼓励 哈希函数能够在学得的正样本对的哈希表示距离 和负样本对的哈希表示距离之间保持一个较大的 间隔,我们采用三元排序损失^[24](triplet ranking loss).给定一个三元组(h, \hat{h}, \hat{h}),其中(h, \hat{h})构成正 样本对,(h, \hat{h})构成负样本对.三元排序损失函数的 优化目标是使得人脸哈希表示之间的相对相似性满 足"与人脸 \hat{h} 相比,人脸 \tilde{h} 与人脸更相似".在本文方 法中,正样本对与负样本对通过样本的语义信息来 进行划分,即属于同一个个体的人脸数据(无论是 人脸图像数据还是视频数据)可以构成正样本对,属 于不同个体的人脸数据可以构成负样本对.三元组 (h, \hat{h}, \hat{h})的损失表示为

 $t(h, \tilde{h}, \hat{h}) = \max(d(h, \tilde{h}) - d(h, \hat{h}) + \zeta, 0)$ (5) 其中, $d(\theta_1, \theta_2)$ 表示二值空间中的海明距离,定义为

$$d(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \frac{l - \boldsymbol{\theta}_1^{\mathrm{T}} \boldsymbol{\theta}_2}{2}$$
(6)

(7)

l 是哈希编码的长度, ζ 表示学得的正样本对的哈希 表达距离和负样本对的哈希表达距离之间所要保持 的间隔.将所有批次中选取的三元组的损失加和,最 终的三元排序损失表示为

$$\mathcal{L}_{t} = \sum_{\substack{i_{1/2}=1, \\ i_{1}\neq i_{2}}}^{c} \sum_{K_{1/2/3} \in \langle I, V \rangle} \sum_{\substack{j_{1/2/3}=1, \\ (K_{2}, j_{2})\neq (K_{1}, j_{1})}}^{n_{i_{1/1/2}}^{K_{1/2/3}}} t(\boldsymbol{h}_{i_{1}, j_{1}}^{K_{1}}, \boldsymbol{h}_{i_{1}, j_{2}}^{K_{2}}, \boldsymbol{h}_{i_{2}, j_{3}}^{K_{3}})$$

其中的 $h_{i,j}^{K}$ 表示的是公共空间特征 $r_{i,j}^{K}$ 对应的二值哈希表示,即第 i 个个体的第 j 个图像或视频(由 K 的取值来决定)对应的二值哈希表示.

同时考虑公共特征空间的判别性约束、可分性 约束和二值空间的三元排序损失,异构哈希网络损 失函数的最终形式为

$$\min_{\mathbf{L}=\alpha} \mathcal{L}_d + \beta \mathcal{L}_s + \mathcal{L}_t \tag{8}$$

其中 $\Gamma = \{ \chi^{I}(\cdot), \eta^{I}(\cdot), \chi^{V}(\cdot), \eta^{V}(\cdot), \phi(\cdot) \}$ 表示 异构哈希网络的 5 个部分所包含的参数的, $\Theta = \{ W_{i}, b_{i} | i = 1, 2, \dots, c \}$ 表示 softmax 损失层的参数, $\Lambda = \{ \mu_{i} | i = 1, 2, \dots, c \}$ 是 Fisher 损失层的参数, α 和 β 是用于平衡各项重要性的超参数.

3.3 优 化

利用经典的反向传播算法来训练异构哈希网络 模型,反向传播算法通过利用随机梯度下降算法优 化式(8)来实现.其中 \mathcal{L}_s 关于参数 Θ 和输入特征 $\mathbf{r}_{i,j}^{K}$ ($K \in \{I,V\}$)的梯度可以较为容易地计算出来,许多 深度学习的开源工具,如 Caffe^[25]等,已经提供了该 层的前向和h反向计算过程,这里不再赘述. \mathcal{L}_d 对 μ_i 和 $\mathbf{r}_{i,j}^{K}$ 的梯度为

$$\frac{\partial \mathcal{L}_{d}}{\partial \boldsymbol{\mu}_{i}} = \frac{1}{n} \sum_{K \in \{I,V\}} \sum_{j=1}^{n_{i}^{K}} (\boldsymbol{\mu}_{i} - \boldsymbol{r}_{i,j}^{K}) - \frac{1}{c} (\boldsymbol{\mu}_{i} - \boldsymbol{\mu}),$$
$$\frac{\partial \mathcal{L}_{d}}{\partial \boldsymbol{r}_{i,j}^{K}} = \frac{1}{n} (\boldsymbol{r}_{i,j}^{K} - \boldsymbol{\mu}_{i}) + 2\boldsymbol{r}_{i,j}^{K}$$
(9)

 \mathcal{L}_{t} 对输入(h, \tilde{h}, \hat{h})的梯度可通过如下公式计算:

$$\frac{\partial t}{\partial h} = \frac{1}{2} (\hat{h} - \tilde{h}) \times I\{A\},$$

$$\frac{\partial t}{\partial \tilde{h}} = -\frac{1}{2} h \times I\{A\},$$

$$\frac{\partial t}{\partial \hat{h}} = \frac{1}{2} h \times I\{A\},$$
(10)

$$A \triangleq d(\boldsymbol{h}, \boldsymbol{\tilde{h}}) - d(\boldsymbol{h}, \boldsymbol{\tilde{h}}) + \boldsymbol{\zeta} > 0,$$

其中,I{•}是一个指示函数,当其内部的条件为真时返回1,反之则返回0.根据式(10)可以计算式(7)所示的损失 \mathcal{L}_i 对 $r_{i,j}^{\kappa}$ 的梯度.

如图 1 所示,由于 *L*_a 和 *L*_s 是添加在公共空间的特征表示上的损失,*L*_i 是添加在最终的二值哈希表示上的损失,所以采用先优化 *L*_a 和 *L*_s 再联合优化三者的训练策略,即整个优化方法分为两个步骤:预训练与微调.

在预训练步骤中,对异构哈希网络中除哈希函 数以外的部分,即图像分支和视频分支进行优化,以 便为下一步网络整体微调提供一个较好的初值.预训 练步骤之后的异构哈希网络可以提供两个有效的特 征表示生成器 $\chi^{I}(\cdot) + \eta^{I}(\cdot)$ 和 $\chi^{v}(\cdot) + \eta^{v}(\cdot)$, 它们不仅能够分别为人脸图像和人脸视频数据提 取合适的特征表示,而且能够将位于异构空间中的 数据映射到一个具有较强判别力的公共空间中,从 而实现异构数据的特征表示同构化.由于在此过程 中,只有 $\chi^{I}(\cdot), \eta^{I}(\cdot), \chi^{v}(\cdot)$ 和 $\eta^{v}(\cdot)$ 是可学习的, 而 $\phi(\cdot)$ 是固定不动的,因此目标函数中只利用 \mathcal{L} a和 \mathcal{L}_{s} 即可.在预训练过程中使用了 softmax 和 Fisher 损失,所以利用预训练之后的异构哈希网络 计算得到的人脸图像与人脸视频的特征表示是可分 且具有足够判别力的,这就显著减少了同一个体不 同模态数据间的差异.

微调步骤利用预训练步骤获得的初始值对整个 异构哈希网络进行优化,这一步骤的目的是将图像 分支、视频分支与哈希函数融合至一个统一的优化 框架中,使网络的各组成部分能够通力协作、相互促 进,且都能很好地为最终的检索任务服务.图像与视 频分支对原始数据处理的情况能够影响后面哈希函 数的性能,反过来,哈希函数的性能也能够指导这两 个分支的学习.因此,微调过程在预训练过程的基础 之上,进一步减少同一个体不同模态数据间的差异, 从而实现跨模态人脸检索任务.

算法1总结了异构哈希网络的训练过程,其中, $\Phi = \{\chi^{I}(\bullet), \eta^{I}(\bullet), \chi^{V}(\bullet), \eta^{V}(\bullet)\}$ 表示图像与视频 分支的所有参数, $\Gamma = \Phi \cup \{\phi(\bullet)\}$ 表示整个网络的参 数, γ 和 λ 是网络更新时的学习率.

算法1. 异构哈希网络训练算法.

输入:人脸图像数据 X¹、人脸图像数据 X^V 以及它们相 对应的标签信息、超参数 α 和 β

输出:HHN 的所有参数 Γ

步骤 1. 预训练.

t **←**0.

随机初始化网络和损失层的参数: $\Psi' = \Phi' \cup \Theta' \cup \Lambda'$. 迭代.

前向过程:计算 Fisher 和 softmax 的联合损失

 $\mathcal{L} = {}_{\alpha}\mathcal{L}_{d} + \beta \mathcal{L}_{s}.$

反向过程:根据式(9)逐层计算联合损失对输入的梯度 以及对网络和损失层参数的梯度 $\partial L / \partial \Psi'$.

更新过程: $\Psi^{t} = \Psi^{t+1} - \gamma^{t} (\partial \mathcal{L} / \partial \Psi^{t}).$

 $t \leftarrow t+1$.

```
至收敛
```

```
步骤 2. 微调.
```

s **←**0.

根据步骤 1 结果初始化 Ψ ,并随机初始化哈希函数 $\varphi(\bullet)$, 即 $\Omega^{s} = \Psi' \bigcup \{ \phi^{s}(\bullet) \}.$

迭代:

- 前向过程:计算 Fisher、softmax 和三元排序的联合 损失 $\mathcal{L} = \alpha \mathcal{L}_d + \beta \mathcal{L}_s + \mathcal{L}_t$.
- 后向过程:根据式(9)和式(10)逐层计算联合损失对 输入的梯度以及对网络和损失层参数的 梯度∂L/∂Ω^s.

更新过程: $\Omega^{s+1} = \Omega^s - \lambda^s (\partial \mathcal{L} / \partial \Omega^s)$.

 $s \leftarrow s + 1$.

至收敛

4 实 验

4.1 数据集与实验设置

为了评价所提出方法的性能,我们在 ICT-TV^[26]

和 Celebrity-1000^[27] 两个公共数据集上进行了实 验. ICT-TV 中的人脸视频片段均收集于两部美剧的 第一季:《生活大爆炸》(the Big Bang Theory, BBT) 和《越狱》(Prison Break, PB). BBT 数据集共包含了 14个人的共4458段人脸视频片段,每个人的视频 片段数量最小为11,最大为1528.PB数据集共包含 了 19 个人的共 7500 段人脸视频片段,每个人的视 频片段数量最小为 49,最大为 1965. 两部美剧所提 供的人脸视频片段都是经过标注仅包含人脸区域, 目以逐帧方式存储的每帧大小为150×150 像素的 人脸图像. ICT-TV 数据集是根据美剧建立的,其中 包含的人物数量较少.与 ICT-TV 数据集不同, Celebrity-1000(Celeb1K)是一个大规模的人脸视频 数据集,包含1000个人的159726个人脸视频片段, 共含有大约240万帧数据(每个人脸视频片段约包 含15帧),这些人脸视频片段都是在真实场景中拍 摄的.为了消除真实场景中的背景干扰,人脸视频的 所有帧都经过了人脸检测和人脸对齐的操作,并归 一化为大小为 64×48 像素的人脸灰度图像.

由于数据集中每个人的视频数量相差较大,且 有些视频帧数较少,所以在 BBT、PB、Celeb1K 数据 集上都分别选择一个子集作为实验数据.具体来讲, 先根据数据集中每个个体包含人脸视频的数量由大 到小对个体进行排序,之后选取含有视频较多的个 体的数据用来进行实验,其中 BBT 数据集选取视频 数最多的5个个体,PB数据集选取视频数最多的 7个个体, Celeb1K 选取视频数最多的 15 个个体. 对于每个个体,从其视频中随机选取 200 个作为训 练数据.由于多模态检索方法在训练时需要提供各 个模态的数据,所以,对于每个个体选取出的 200 个 视频片段,每个视频片段中随机选取一帧作为训练 图像,即每个个体提供 200 个训练视频和 200 个训 练图像.因此,BBT 数据集的训练数据为 1000 张人 脸图像和 1000 个人脸视频, PB 数据集的训练数据 为1400张人脸图像和1400个人脸视频, Celeb1K 数据集的训练数据为 3000 张人脸图像和 3000 个人 脸视频.训练集选取之后,将数据子集中剩余的人脸 视频作为人脸视频测试集,仍然从每个视频中随机 选取一帧组成人脸图像测试集.针对不同的检索任 务,测试集中对查询集和待查询集的划分略有不同. 对于"基于图像的人脸视频检索"任务,从每个个体 的测试人脸图像中随机选取 5 个作为查询集,整个 人脸视频测试集作为待查询集;对于"基于视频的人 脸图像检索"任务,从每个个体的测试人脸视频中随 机选取 5 个作为查询集,整个人脸图像测试集作为 待查询集.

4.2 实现细节

为了与其它哈希方法进行公平的比较,本文方 法与对比方法针对人脸图像和视频使用相同的特征 表示.具体来讲,人脸图像用一个 p 维的向量表示, 这个向量是由原始图像的像素特征经过主成分分析 (Principal Component Analysis,PCA)获得的,称为 灰度 PCA 特征;对于一个人脸视频片段,将该视频 片段看成各帧数据的无序集合,首先对每一帧提取 和人脸图像一样的 p 维灰度 PCA 特征,然后利用 这些帧的灰度 PCA 特征计算一个大小为 $p \times p$ 的 RBF 核矩阵特征 **K**,称为灰度 RBF 核矩阵特征:

$$\boldsymbol{K}(i,j) = \exp\left(-\frac{\|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2^2}{2\sigma^2}\right),$$

$$\sigma = \frac{1}{p^2} \sum_{i=1}^p \sum_{j=1}^p \|\boldsymbol{x}_i - \boldsymbol{x}_j\|_2$$
(11)

其中 x_i 与 x_j 分别表示矩阵 $X \in \mathbb{R}^{p \times q}$ 的第i和第j行,X是一个含有q帧的人脸视频的各帧p维灰度 PCA 特征矩阵,K(i,j)表示矩阵K的第i 行策j列的元素.

在本文方法的实现中,图像分支由全连接层和 ReLU激活层相互堆叠构成,该分支的大小为"100-512-1024-100",这里的数字代表输入层和全连接 层的所含有的神经元的个数,每个全连接层后面跟 随着一个 ReLU激活层;视频分支首先使用一个向 量化操作层将核矩阵转化为向量,之后是全连接层 和 ReLU激活层的堆叠,因此视频分支的大小为 "100×100-5050-100-512-1024-100".向量化 操作层^[28]执行的操作是

$$\operatorname{vec}(\mathbf{K}) = \begin{bmatrix} \mathbf{V}_{1,1}, \sqrt{2}\mathbf{V}_{1,2}, \cdots, \sqrt{2}\mathbf{V}_{1,p}, \mathbf{V}_{2,2}, \sqrt{2}\mathbf{V}_{2,3}, \cdots, \\ \sqrt{2}\mathbf{V}_{2,p}, \cdots, \sqrt{2}\mathbf{V}_{p-1,p}, \mathbf{V}_{p,p} \end{bmatrix}^{\mathrm{T}}$$
(12)

通过这样的方式,向量化操作层输出 p(p+1)/2维的向量,其中 $V = \log(K)$ 是 RBF 核矩阵 K 的矩 阵对数.哈希函数也是由全连接层和 ReLU 激活层 组成,大小为"100-100-*l*",每个全连接层后面都 接着一个 ReLU 激活层,其中 *l* 表示哈希编码的长 度.最后一个全连接层之后,添加一个双曲正切层来 近似 sgn(•)符号函数来实现哈希码的量化.以上的 网络结构是在 ICT-TV 数据集上(包括 BBT 和 PB) 的.Celebrity 数据集包含的个体更多,数据集更复 杂,因此网络的结构也设计得更大,图像分支大小 为"100-256-512-512-1024-2048-100",视频 分支大小为"100×100-5050-100-256-512512-1024-2048-100",哈希函数的大小为"100-100-*l*".

预处理与调优过程都是在 Caffe 框架^[25]下实现 的.预训练过程通过随机梯度下降方法进行,其中 momentum 和 weight decay 分别设置为 0.9 和 5× 10^{-4} ,学习率初始化为 0.01 并按照多项式规则(polynomial policy)进行递减,递减的参数 power value 设置为 0.8.训练集的批次大小设置为 256,迭代的 总次数设置为 10^5 .微调过程中的 momentum 和 weight decay 分别设置为 0.8 和 5×10⁻⁵,学习率初 始化为 0.001 并按照多项式规则进行递减,递减的 参数 power value 设置为 0.8,全部的迭代次数为 5×10⁴.

4.3 对比方法和评价指标

我们与7个单模态哈希方法和6个多模态哈希 方法进行了对比:

(1) 单模态哈希方法:LSH^[6]、SH^[7]、ITQ^[8]、
 SITQ^[8]、RR^[8]、SSH^[9]和KSH^[10];

(2) 多 模 态 哈 希 方 法: CCA^[8]、PDH^[13]、 CMSSH^[14]、MMNN^[16]、SCM^[18]和 HER^[5].

为实现公平的比较,我们的方法和所有对比方 法都采用相同的数据集划分策略和相同的数据特征 表示,即人脸图像的 100 维灰度 PCA 特征和人脸视 频的大小为 100×100 的灰度 RBF 核矩阵特征.由 于一些对比的哈希方法仅能生成长度比输入特征小 的哈希编码,因此,在实验中,哈希码的长度分别设 置为 8、16、32 和 64 来显示哈希方法性能随哈希编 码位数不同的变化情况.为了公平的比较,所有对比 方法中的重要参数都是根据原始参考文献中的建议 设置的.

我们使用准确率-召回率(Precision Recall, PR)曲线、准确率-优先返回样本数量(Precision curve w. r. t. Number of top returned samples, PN)曲线、召回率-优先返回样本数量(Recall curve w. r. t. Number of top returned samples,RN)曲线 和平均正确率均值(mean Average Precision,mAP) 进行评价.同时,也列出了 mAP 对哈希码位数的变 化曲线以探究不同哈希方法中哈希码长对检索性能 的影响.实验中报告的各性能指标均为 300 次测试 的平均值.受篇幅所限,我们仅给出哈希码长为 64 位时的 PR、PN 和 RN 曲线.

4.4 实验结果与分析

严格来讲,单模态哈希方法并不能直接应用于 跨模态的人脸检索任务.为了使单模态方法可以适 用于跨模态检索任务,我们将人脸视频片段看成是 各帧数据的无序集合,然后,将一张人脸图像与一个 人脸视频片段的距离用该图像与视频中每一帧的距 离的平均值来表示,从而实现跨图像与视频模态的 检索任务.对于多模态的哈希方法,由于除 HER 以 外的方法都只能处理表示为欧式空间中的向量的数 据,因此利用式(12)先将人脸视频的灰度 RBF 核矩 阵特征向量化,然后再执行这些方法.

为了充分评价本文方法在跨模态人脸检索任务上的性能,我们进行了两方面的实验:基于图像的人 脸视频检索和基于视频的人脸图像检索.对于基于 图像的人脸视频检索任务,表1和表2分别列出了 本文方法与对比方法在ICT-TV数据集和 Celeb1K 数据集上的 mAP 比较;对于基于视频的图像检索 任务,表3和表4分别列出了本文方法与对比方法在 ICT-TV 数据集和 Celeb1K 数据集上的 mAP 比较. 图 2、图 3和图 4分别显示了这些方法在 ICT-TV 的 BBT 数据集、ICT-TV 的 PB 数据集以及 Celeb1K 数据集上两种跨模态检索任务的 PR、PN、RN 以及 mAP 曲线的对比.在这3个图中,左栏为"基于图像 的人脸视频检索"任务的相应曲线, 右栏为"基于视频的人脸图像检索"任务的相应曲线.

表 1 ICT-TV 数据集上本文方法与其它方法在基于图像的 人脸视频检索任务上的 mAP 比较

数据集	类型	方法	8位	16 位	32 位	64 位
		LSH	0.2448	0.2919	0.3162	0.3726
	単樹	RR	0.2916	0.3551	0.3571	0.3929
	医态	ITQ	0.3546	0.3599	0.3898	0.3998
	哈	SH	0.2881	0.3421	0.3145	0.3177
	希	SSH	0.3432	0.2761	0.2541	0.2447
	力 注	KSH	0.7825	0.7988	0.8638	0.8262
DDT	144	SITQ	0.4703	0.3818	0.3837	0.3379
DDI	4	CCA	0.3664	0.3067	0.2763	0.2440
	多齿	CMSSH	0.5582	0.5985	0.5794	0.5489
	医态	PDH	0.2252	0.2223	0.2817	0.2669
	哈希方法	MMNN	0.6322	0.7296	0.7860	0.7819
		SCM	0.7659	0.8461	0.8690	0.8605
		HER	0.8017	0.8121	0.8677	0.8673
		HHN	0.8218	0.8541	0.9051	0.9112
	24	LSH	0.1764	0.1728	0.1868	0.1857
	早 档	RR	0.1689	0.1786	0.1872	0.1896
	(态哈希方注	ITQ	0.1746	0.1728	0.1794	0.1860
		SH	0.1685	0.1752	0.1903	0.2028
		SSH	0.1838	0.1840	0.1799	0.1823
		KSH	0.3420	0.4006	0.4260	0.4414
PB	14	SITQ	0.1893	0.1818	0.1824	0.1831
ТD	kı	CCA	0.2193	0.2363	0.2198	0.1982
	多樟	CMSSH	0.2567	0.2397	0.2162	0.2078
	态	PDH	0.2316	0.2318	0.2131	0.2142
	哈	MMNN	0.4076	0.5522	0.5328	0.7126
	希	SCM	0.4946	0.6165	0.6247	0.6927
	万法	HER	0.4601	0.5666	0.6123	0.6268
	14	HHN	0.7193	0.7608	0.7956	0.8162

表 2 Celeb1K 数据集上本文方法与其它方法在基于图像的 人脸视频检索任务上的 mAP 比较

数据集	类型	方法	8位	16 位	32 位	64 位
-	单	LSH	0.0912	0.1130	0.1545	0.1915
		RR	0.1156	0.1447	0.1904	0.2094
	太	ITQ	0.1169	0.1686	0.2054	0.2317
	哈	SH	0.1276	0.1648	0.1935	0.2063
	希方法	SSH	0.1342	0.1765	0.1823	0.1568
		KSH	0.2796	0.3468	0.3841	0.4383
Colob 1K		SITQ	0.1541	0.1937	0.2280	0.2341
Celebik	多模态	CCA	0.1646	0.1950	0.1814	0.1470
		CMSSH	0.1738	0.1652	0.1500	0.1402
		PDH	0.1256	0.1264	0.1191	0.1229
	哈	MMNN	0.1735	0.3656	0.3918	0.3672
	希	SCM	0.2654	0.3370	0.4157	0.4809
	力 辻	HER	0.7324	0.7740	0.7659	0.7970
	1Z	HHN	0.7705	0.7831	0.8028	0.8254

表 3 ICT-TV 数据集上本文方法与其它方法在基于视频的 人脸图像检索任务上的 mAP 比较

数据集	类型	方法	8位	16 位	32 位	64 位
	**	LSH	0.2300	0.2969	0.3410	0.3426
	甲樟	RR	0.3013	0.3537	0.3723	0.3865
	态	ITQ	0.3213	0.3524	0.3779	0.3986
	哈	SH	0.2728	0.3336	0.3120	0.3098
	希	SSH	0.3266	0.2809	0.2571	0.2510
	力	KSH	0.7296	0.7852	0.8459	0.8023
DDT	厶	SITQ	0.4451	0.3787	0.3895	0.3600
DDI	M	CCA	0.3688	0.3057	0.2739	0.2423
	シ 樟	CMSSH	0.5905	0.6109	0.6096	0.5834
	太	PDH	0.2298	0.2278	0.2271	0.2517
	哈	MMNN	0.6694	0.7552	0.8193	0.8256
	希	SCM	0.8079	0.8814	0.8966	0.8907
	力	HER	0.6567	0.7214	0.7354	0.7888
	14	HHN	0.9080	0.9274	0.9364	0.9457
1	出	LSH	0.1641	0.1817	0.1899	0.1850
	▶ 岸	RR	0.1776	0.1728	0.1876	0.1934
	态	ITQ	0.1824	0.1745	0.1779	0.1869
	哈	SH	0.1699	0.1762	0.1898	0.2014
	希方	SSH	0.1900	0.1859	0.1832	0.1853
		KSH	0.3337	0.3920	0.4246	0.4309
PR	14	SITQ	0.1951	0.1847	0.1810	0.1847
ТD	1	CCA	0.2184	0.2375	0.2203	0.1977
	夕樟	CMSSH	0.2084	0.1946	0.1817	0.1693
	态	PDH	0.2290	0.2292	0.2083	0.2076
	哈	MMNN	0.3980	0.5522	0.5500	0.7067
	希	SCM	0.5195	0.6400	0.6458	0.7110
	力 注	HER	0.4551	0.5773	0.6580	0.6947
	14	HHN	0.7265	0.7787	0.8040	0.8160

表 4 Celeb1K 数据集上本文方法与其它方法在基于视频的 人脸图像检索任务上的 mAP 比较

数据集	类型	方法	8位	16 位	32 位	64 位
	单档	LSH	0.0768	0.1096	0.1486	0.1946
		RR	0.0989	0.1559	0.1813	0.2160
	太	ITQ	0.1130	0.1632	0.2028	0.2302
	哈	SH	0.1193	0.1645	0.1965	0.2093
	希方法	SSH	0.1196	0.1761	0.1897	0.1638
		KSH	0.2496	0.3366	0.3773	0.4484
Colob1K		SITQ	0.1270	0.1827	0.2225	0.2318
Celebin	多模态	CCA	0.1634	0.1917	0.1801	0.1458
		CMSSH	0.1725	0.1601	0.1395	0.1269
		PDH	0.1193	0.1213	0.1150	0.1159
	哈	MMNN	0.1851	0.3818	0.4018	0.3878
	希	SCM	0.2666	0.3351	0.4094	0.4733
	力	HER	0.5909	0.7260	0.7658	0.8177
	広	HHN	0.7612	0.7786	0.7766	0.8323



图 2 BBT 数据集上本方法与其它方法在两类跨模态人脸 检索任务上的曲线对比图(左栏为基于图像的人脸视 频检索的相应曲线,右栏为基于视频的人脸图像检索 的相应曲线)

可以看出,本文提出的跨模态人脸检索方法在 多数情况下性能都显著优于其它的多模态哈希方 法,原因有以下几方面:异构哈希网络优化过程中同 时考虑了模态内与模态间的判别性;由于同时使用 了 Fisher 损失、softmax 损失和三元排序三种损失, 无论是公共空间的特征表示还是最终学得的二值 哈希表示都具备较强的判别力;采用了两阶段的训

图 3 PB数据集上本方法与其它方法在两类跨模态人脸检索 任务上的曲线对比图(左栏为基于图像的人脸视频检索 的相应曲线,右栏为基于视频的人脸图像检索的相应 曲线)

练方式,这使得网络可以首先将数据从异构的空间 中投影到具有较强判别力的公共空间,然后在公共 空间上学习哈希函数,使得各阶段目的明确,便于提 升特征表示的判别力;本文提出的异构哈希网络对 原始数据及从原空间到二值空间的映射没有做任何 先验性的假设,全部是基于数据学习而来,而 HER 方法的核函数可能在某些数据上并不适用.



图 4 Celeb1K 数据集上本方法与其它方法在跨模态人脸检 索任务上的曲线对比图(左栏为基于图像的人脸视频 检索的相应曲线,右栏为基于视频的人脸图像检索的 相应曲线)

从实验结果中还可以发现,多模态方法要好于 单模态方法,这主要是因为多模态方法利用了灰度 RBF核矩阵特征描述符来表示人脸视频,该描述符 可以充分利用视频各帧之间的相关关系来建模人脸 视频中包含的诸如光照、姿态等的复杂变化.此外, 在单模态哈希方法中,LSH、RR、ITQ和SH方法属 于无监督的哈希方法,SSH、KSH、SITQ属于半监 督或有监督的哈希方法;在多模态哈希方法中, PDH 是无监督方法,CCA、CMSSH、MMNN、SCM、 HER 和本文的 HHN 方法是有监督方法.可以发现,有监督和半监督哈希方法在大多数情况下会好 于非监督的哈希方法,主要原因是有监督和半监督 哈希方法在训练过程中充分利用了标签信息,而标签 信息是对检索系统进行评价时唯一可依赖的信息.

4.5 三元排序损失函数的作用

为了验证三元排序损失函数的作用,我们对比 了使用三元排序损失函数所得到的哈希编码 (HHN)与直接量化得到的哈希编码(HHN-2)在跨 模态人脸检索任务上的性能.表5和表6分别列出 了三元损失函数在基于图像的人脸视频检索任务上 和基于视频的人脸图像检索任务上对mAP的影 响.可以看出,通过三元排序损失函数所优化得到的 哈希编码比直接量化得到的哈希编码取得了更好的 检索结果,说明了三元排序损失函数对于检索的重 要作用.

表 5 三元损失函数在基于图像的人脸视频检索任务上 对 mAP 的影响

	数据集	方法	8位	16 位	32 位	64 位
DDT	HHN-2	0.6438	0.6412	0.6394	0.6435	
	DD I	HHN	0.8218	0.8541	0.9051	0.9112
	np	HHN-2	0.7022	0.6904	0.7001	0.6939
	FD	HHN	0.7193	0.7608	0.7956	0.8162
	C-L-L IF	HHN-2	0.6909	0.7005	0.6971	0.6961
	Celebin	HHN	0.7705	0.7831	0.8028	0.8254

表 6 三元损失函数在基于视频的人脸图像检索任务上 对 mAP 的影响

数据集	方法	8位	16 位	32 位	64 位
BBT	HHN-2	0.6674	0.7210	0.7486	0.7776
	HHN	0.9080	0.9274	0.9364	0.9457
PB	HHN-2	0.6910	0.7210	0.7486	0.7776
	HHN	0.7265	0.7787	0.8040	0.8160
Celeb1K	HHN-2	0.7021	0.7398	0.7545	0.7754
	HHN	0.7612	0.7786	0.7766	0.8323

4.6 检索效率

在本实验中,我们在 CeleblK 数据集上对比了 我们的方法与其它对比方法的检索效率.由于所有 的方法都使用相同的特征,因此没有计算特征提取 的时间,而只计算了检索的时间.表 7 列出了各个方

	表 7 档	佥索效率	(单位:s)
方法	检索时间	方法	检索时间
LSH	2.2863e-06	CCA	1.0901e-06
RR	4.0581e-06	CMSSH	2.1930e-06
ITQ	3.7152e-06	PDH	2.1930e-06
SH	1.0528e-04	MMNN	7.6754e-06
SSH	3.5437e-06	SCM	1.0965e-06
KSH	0.0154	HER	0.0433
SITQ	9.8881e-06	HHN	5.9125e-06

83

法的平均检索时间.可以看出,KSH,HER 以及 SH 所需的检索时间较长,而我们的方法与其它的方法 具有相似的检索效率.

5 总 结

本文提出了一种使用异构哈希网络的跨模态人 脸检索方法.异构哈希网络能够将位于异构空间的 人脸图像和视频映射到一个公共且有判别力的二值 空间上,获得相应的哈希表示,实现异构数据的特征 表示同构化.网络的训练同时使用了 Fisher 损失, softmax 损失和三元排序损失,其中的 Fisher 损失 使用了类内散度和类间散度的差分形式,并保持均 值向量的可学习性,以方便使用基于批次的优化方 法.Fisher 损失关注于公共空间的判别力,softmax 损失强调公共空间上特征表示的可分性,三元排序 损失旨在提升最终的检索性能,三者相互补充,共同 作用于网络训练.异构哈希网络为基于深度学习的 跨模态哈希方法提供了一套通用的解决框架,不仅 可以完成跨图像和视频模态的检索任务,还可以容 易地扩展到其他跨模态的检索任务上.

参考文献

- [1] Sun Y, Wang X, Tang X. Sparsifying neural network connections for face recognition//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Las Vegas, USA, 2016: 4856-4864
- [2] Taigman Y, Yang M, Ranzato M A, Wolf L. Web-scale training for face identification//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Boston, USA, 2015: 2746-2754
- [3] Tran D, Bourdev L, Fergus R, et al. Learning spatiotemporal features with 3D convolutional networks//Proceedings of the IEEE International Conference on Computer Vision. Santiago, Chile, 2015: 4489-4497
- [4] Yang J, Ren P, Chen D, et al. Neural aggregation network for video face recognition//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Honolulu, USA, 2017: 5216-5225
- [5] Li Y, Wang R, Huang Z, et al. Face video retrieval with imagequery via hashing across Euclidean space and Riemannian manifold//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Boston, USA, 2015: 4758-4767
- [6] Gionis A, Indyk P, Motwani R. Similarity search in high dimensions via hashing//Proceedings of the IEEE Conference on Very Large Data Bases. Edinburgh, UK, 1999: 518-529

- [7] Weiss Y, Torralba A, Fergus R. Spectral hashing//Proceedings of the Advances in neural information processing systems. Vancouver, Canada, 2009: 1753-1760
- [8] Gong Y, Lazebnik S. Iterative quantization: A procrustean approach to learning binary codes//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Colorado Springs, USA, 2011: 817-824
- [9] Wang J, Kumar S, Chang S-F. Semi-supervised hashing for scalable image retrieval//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. San Francisco, USA, 2010; 3424-3431
- [10] Liu W, Wang J, Ji R, et al. Supervised hashing with kernels// Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Providence, USA, 2012: 2074-2081
- [11] Feng F, Wang X, Li R. Cross-modal retrieval with correspondence autoencoder//Proceedings of the ACM International Conference on Multimedia. Orlando, USA, 2014: 7-16
- [12] Kumar S, Udupa R. Learning hash functions for cross-view similarity search//Proceedings of the International Joint Conference on Artificial Intelligence. Barcelona, Spain, 2011: 1360-1365
- [13] Rastegari M, Choi J, Fakhraei S, et al. Predictable dual-view hashing//Proceedings of the International Conference on Machine Learning. Atlanta, USA, 2013: 1328-1336
- [14] Bronstein M M, Bronstein A M, Michel F, Paragios N. Data fusion through cross-modality metric learning using similaritysensitive hashing//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. San Francisco, USA, 2010: 3594-3601
- [15] Zhen Y. Yeung D-Y. A probabilistic model for multimodal hash function learning//Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining. Beijing, China, 2012; 940-948
- [16] Masci J, Bronstein M M, Bronstein A M, Schmidhuber J. Multimodal similarity-preserving hashing. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2014, 36: 824-830
- [17] Zhai D, Chang H, Zhen Y, et al. Parametric local multimodal hashing for cross-view similarity search//Proceedings of the International Joint Conference on Artificial Intelligence. Beijing, China, 2013; 2754-2760
- [18] Zhang D, Li W-J. Large-scale supervised multimodal hashing with semantic correlation maximization//Proceedings of the AAAI Conference on Artificial Intelligence. Québec City, Canada, 2014: 2177-2183
- [19] Lin Z, Ding G, Hu M, Wang J. Semantics-preserving hashing for cross-view retrieval//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Boston, USA, 2015: 3864-3872
- [20] Liu X, He J, Deng C, Lang B. Collaborative hashing// Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Columbus, USA, 2014: 2139-2146

- [21] Jiang Q-Y, Li W-J. Deep cross-modal hashing//Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. Honolulu, USA, 2017: 3270-3278
- [22] Cao Y, Long M, Wang J, et al. Deep visual-semantic hashing for cross-modal retrieval//Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining. San Francisco, USA, 2016: 1445-1454
- [23] Yang M, Zhang L, Feng X, Zhang D. Fisher discrimination dictionary learning for sparse representation//Proceedings of the IEEE International Conference on Computer Vision. Barcelona, Spain, 2011, 543-550
- [24] Norouzi M, Fleet D J, Salakhutdinov R R. Hamming distance metric learning//Proceedings of the Advances in neural information processing systems. Lake Tahoe, USA, 2012: 1061-1069
- [25] Jia Y, Shelhamer E, Donahue J, et al. Caffe: Convolutional

DONG Zhen, born in 1990, Ph. D. His research interests include computer vision, face recognition and retrieval. architecture for fast feature embedding//Proceedings of the ACM International Conference on Multimedia. Orlando, USA, 2014: 675-678

- [26] Li Y, Wang R, Shan S, Chen X. Hierarchical hybrid statistic based video binary code and its application to face retrieval in TV-series//Proceedings of the IEEE International Conference and Workshops on Automatic Face and Gesture Recognition. Ljubljana, Slovenia, 2015; 1-8
- [27] Liu L, Zhang L, Liu H, Yan S. Toward large-population face identification in unconstrained videos. IEEE Transactions on Circuits and Systems for Video Technology, 2014, 24: 1874-1884
- [28] Tuzel O, Porikli F, Meer P. Pedestrian detection via classification on Riemannian manifolds. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2008, 30: 1713-1727

PEI Ming-Tao, born in 1977, Ph. D., associate professor, Ph. D. supervisor. His research interests include computer vision and face recognition.

Background

Cross-modality face retrieval is to retrieve faces of a particular person in one modality given his/her face information in another modality, such as retrieving video shots containing particular person given one image of him/her (query-byimage video retrieval), or retrieving the face images of one person by using his/her video clip as query (query-by-video image retrieval). It is an important problem in computer vision.

The current methods all concentrate on face video retrieval, and only one previous work is about face retrieval across image and video domains as far as we know. Li et al. proposed a hashing method across the Euclidean space and the Riemannian manifold to measure the similarity of face images and videos for cross-domain face retrieval. Although they achieve superior performances than many traditional single-modality and multi-modality methods, the extracted features (DCT and covariance matrix) and the learned hash functions are incompatible due to the independent feature extraction and hash learning procedures, which might influence the effects of the generated hash codes. In contrast, our heterogeneous hashing network integrates the feature extraction and hash learning into a unified optimization framework where two procedures interact for the final retrieval task. Experiments on a large-scale face video dataset and two challenging TV-series datasets demonstrate the effectiveness of the proposed method.

This work was supported by the National Natural Science Foundation of China (No. 61472038).