FF-CAM:基于通道注意机制前后端融合的人 群计数

张宇倩 李国辉 雷军 何嘉宇

国防科技大学系统工程学院重点实验室,长沙,410073

摘 要 单个图像中的人群计数在计算机视觉领域中备受关注,因为其在公共安全方面具有重要作用。例如,在人群聚集的场景中监控设备可以实时监测人群数量变化,对过度拥挤和异常情况进行预警以预防安全事故的发生。然而,由于受到遮挡,透视扭曲,尺度变化,背景干扰的严重影响,在单个图像中对人群计数的预测要达到较高精确度是极其困难的,其面临着巨大的挑战。在本文中,我们提出了一个创新的名为 FF-CAM 的模型来计算图像中的人群数量。它首先将主网络低层的特征图与高层的特征图合并,实现不同尺度的特征融合,且无需额外的分支或子任务,解决了由于透视导致的尺度多样性问题。随后融合的特征图被送入通道注意力模块以优化不同特征的融合过程,并进行特征通道的重新校准以充分使用全局和空间信息。此外,我们在网络的末端利用扩张卷积来获得高质量的人群密度图,扩张卷积层扩大了感受野,其输出包含更详细的空间信息和全局信息,不会降低空间分辨率。最后,我们加入基于 SSIM 的损失函数用于比较估计人群密度图和真值的局部相关性,和基于回归人数的损失函数用于比较估计人群数量与真实人数之间的差异。我们的 FF-CAM 在 UCF_CC_50 数据集, Shanghai Tech 数据集和 UCF_QRNF 数据集中进行训练并测试,获得了出色的结果。在 UCF_CC_50 数据集上比现有方法的 MAE 提高了 4.5%, MSE 提高了 3.8%。

关键词 人群计数;特征融合;通道注意力;扩张卷积;高质量密度图 中图法分类号 TP391

FF-CAM: Crowd counting based on frontend-backend fusion through channel-attention mechanism

Yuqian Zhang Guohui Li Jun Lei Jiayu He

Science and Technology on Information Systems Engineering Laboratory, National University of Defense Technology, Changsha, P.R.China, 410073 **Abstract** Crowd counting has attracted much attention in computer vision owing to its contribution in public security. For example, in a crowd gathering scenario, the monitoring device can monitor changes in the number of people in real time, and provide early warning of overcrowding and abnormal conditions to prevent the occurrence of safety accidents. But on account of occlusions, perspective distortions, scale variations and background interference it faces a great challenge to achieve high accuracy on prediction of crowd counting in a single image. In this paper we propose a novel model to count crowds named FF-CAM. It merges the front-end feature map with the backend feature map in the baseline, achieving a fusion of various scale features without additional branches or extra subtasks. The fusion is fed into the channel-attention block to optimize the procedure, and to conduct feature recalibration to use global and spatial information. Furthermore, we utilize dilated layers to obtain a high-quality density map. The dilated convolutional layer expands the receptive field, and its output contains more detailed spatial information and global information without reducing the spatial resolution. The SSIM-based loss function is added to compare the local correlation between the estimated density

本课题得到国家自然科学基金(No. 71673293)和国家自然科学基金(No.61806215)资助. **张字倩**,性别女,1996年生,硕士研究生,主要研究领域为计 算机视觉,深度学习以及信息系统工程. E-mail: 446579794@qq.com. **李国辉(通信作者)**,性别男,1963年生,博士,教授,博士生导师,主要研 究领域为计算机视觉,信息系统工程,数据挖掘及虚拟现实技术.E-mail: guohli@nudt.edu.cn. **雷军**,性别男,1989年生,博士,讲师,主要研究领域 为计算机视觉,深度学习,数据挖掘及虚拟现实技术.E-mail: leijun1987@nudt.edu.cn. **何嘉宇**,性别男,1997年生,硕士研究生,主要研究领域为深 度学习,数据挖掘及虚拟现实技术.E-mail: 451567413@qq.com.

map and the ground truth, meanwhile the regression-based loss function is added to compare the difference between the estimated number and the actual number of crowd. Our FF-CAM is verified in the UCF_CC_50 dataset, the ShanghaiTech dataset, and the UCF_QRNF dataset, getting brilliant estimations. Compared to state-of-the-art, MAE is improved by 4.5% and MSE is improved by 3.8% in the UCF_CC_50 dataset.

Key words crowd counting; features fusion; channel-attention; dilated convolutions; high-quality density map

1 引言

近年来,随着生活水平的提高和交通的快速发展,人群计数因其在公共安全方面的贡献而备受关注。例如,在人群聚集的场景中监控设备可以实时 监测人群数量变化,预防过度拥挤和异常情况。然 而,由于受到遮挡,透视扭曲,尺度变化,背景干 扰的严重影响,在单个图像中对人群计数的预测要 达到较高精确性是极其困难的。

在大量的研究和努力之下,人群计数已经取得 了较大的进展。早期的工作主要是检测人群中的每 个行人[1],或使用多个人工提取的特征回归得到人 数[2]。但是在拥挤的场景中由于严重的遮挡难以准 确检测到行人,故会存在较大误差。近年来,主流 的方法由直接计算人数转为生成人群密度图进而 得到总人数以解决严重遮挡问题,基于 GAN 的[3,4] 和基于 CNN[5—13]的方法已经发展并且得到了明 显的改善。此外,人群密度图还包含了空间位置信 息,可更好的应用于安全领域。

然而,由于距监控相机的距离不同和透视问 题,同一幅图像中会存在不同大小的人群,因此人 头尺度多样性是抑制计数准确度的主要难点。一些 工作[5,6,7,8,9]使用具有不同卷积核或是多列的卷 积结构来解决尺度变化的问题,而有些方法[10,11] 则是用相同大小的卷积核堆叠来替换不同的卷积 核。此外,得到的人群密度图由于背景干扰会存在 较大偏差, [12,13]在训练过程中增加了额外的信息 来强调图像中的人群以解决该问题。但这些方法仍 然存在很多不足,不能很好的解决尺度多样性的问 题。Li 等人[10]证明了多列结构中不同分支中的每 列学到的是几乎相同的特征,对尺度变化的贡献很 小。当网络变得复杂时,计算量和计算复杂性急剧 增加,也会导致训练速度的延迟和梯度爆炸。基于 这个问题,为了学习到不同尺度的特征,同时排除 背景噪声的影响,我们考虑采用单列单卷积核的网 络结构,融合低层和高层的特征图。由于网络中不 同级别的层包含不同的比例特征信息,且多个相同 大小卷积核叠加后与大的卷积核具有相同的特征 学习效果。此外,不同级别的层还包含不同级别的 语义信息,低层卷积可以提取细节边缘图案,有效 的回归拥塞区域得到密度图,高层则可以选择性地 获得有用的语义信息,将人头与背景噪声区分开 来。这样做在获得不同尺度信息的同时不增加计算 量和网络结构复杂度。

另一方面,各种特征通过简单的连接难以很好 地对融合的不同尺度大小的人头区域的特征进行 有选择性的加强。另外,卷积层的通道容易被忽略, 从而导致空间信息的不足。而由于生成的密度值遵 循逐像素预测原则,因此输出的密度图必须包含空 间相干性,以便它们可以呈现最近像素之间的平滑 过渡。所以我们考虑将 SE (Squeeze-and-Excitation) 模块[14]引入为通道注意力模块来优化融合。Hu等 人[14]提出,SE模块可以考虑通道的权重,进行特 征重新校准以捕获空间相关性并有选择地强调信 息性强的特征。如此一来将该模块加在特征融合之 后可以优化连接过程,对学习到的不同尺度的特征 图进行加权,有选择性地强调不同尺度的特征,避 免直接连接造成的损失。同时捕获的空间相关性能 使最终生成的密度图呈现最近像素之间的平滑过 渡,以生成高质量的人群密度图。

此外,经过池化层的特征图降低了空间分辨 率,丢失了空间信息,产生的人群密度图质量不够 高。我们考虑在网络末端运用扩张卷积。Li等人[10] 证明了扩张卷积比使用卷积、池化加反卷积的方案 更好地保持特征映射的分辨率,可以包含更详细的 空间信息和全局信息,在扩大了感受野的同时不增 加参数或计算量。所以,我们运用扩张卷积可以生 成高质量的人群密度图。

最后,在人群场景中,高密度区域的局部模式 和纹理特征与其他区域大不相同,但欧几里德的损 失建立在像素独立性假设上并忽略了它们,密度图 的局部相关性未被考虑。另外,其没有将输入图像 的全局计数错误考虑进去,也与用来衡量准确度的 评估指标没有直接关系。为此,我们考虑在损失函 数中加入结构相似性指数(SSIM),和关于回归人 数的损失函数。结构相似性指数根据局部模式计算 两个图像之间的相似性,可以比较生成人群密度图 与真值之间的相似性。关于回归人数的损失函数直 接衡量估计人群数量与真值之间的差异。通过改进 损失函数,网络将生成适合输入图像整体密度水平 的特征,这有助于产生更准确的密度值。

基于上述讨论,我们提出了一种新型人群计数的结构:FF-CAM(Frontend-backend Fusion Network through Channel-attention mechanism),如图1所示。我们提出的方法在UCF_CC_50数据集中的测试结果优于目前最先进的方法。简而言之,我们的贡献包括以下三个方面:

(1)我们融合了主网络低层和高层的特征图。 网络主干只有一列且只有一个大小的卷积内核,减 去了额外的分支及参数量。不同级别的卷积层不仅 包含不同的语义信息,还包含不同的比例特征信 息。它们的融合可以适应由于透视效应引起的尺度 变化,并且共享更多特征,同时可以排除背景干扰。 它还具有更少的参数和计算量。

(2) 我们引入了 SE 模块[14]作为 FF-CAM 的 通道注意力模块。避免直接连接造成的损失,通道 注意力模块可以对融合的不同尺度大小的人头区 域的特征进行有选择性的加强,由此提高网络的表 达能力。另一方面,它可以考虑通道的权重,进行 特征重新校准以捕获空间相关性,使最终生成的人 群密度图呈现最近像素间的平滑过渡。

(3)我们利用一组扩张卷积[10]作为网络的末端。其在增大了感受野的同时保证较少的参数量, 包含了更详细的空间信息和全局信息,可以生成高 质量的人群密度图。此外,我们将 SSIM (结构相 似性)和回归人数加入到损失函数中[7]。SSIM 可 用于估计人群密度图和真值的局部一致性,基于回 归人数的损失函数则衡量估计人群数量与真值之 间的差异。综合后的损失函数可以更好的衡量训练 的估计值与真实值间的差异,产生更准确的密度 值,提高训练准确度。



图 1 FF-CAM 网络的结构图。网络输入的是原始的拥挤人群图像,原始图像依次输入不同的卷积层组合得到不同的特征图, Conv1_2 等即表示从不同的卷积层组合输出的特征图。低层和高层的特征图融合(Concat)后再输入通道注意力模块。最后经 过扩张卷积模块后得到最终的人群密度图。

2 相关工作

在图像和视频中对人群进行计数已经有了很 多年的发展,因为它在视频监控和公共安全中发挥 着重要作用,故而受到计算机视觉领域中人们的长 期关注。但是由于遮挡,透视失真,尺度变化和背 景干扰,计数精度的提高是一个相当大的挑战。目 前人群场景计数的研究大致有以下一些方法:

2.1 传统的方法

2.1.1 基于检测的方法

早期的工作主要是检测单个个体并计算数量。 2012年,Piotr Dollar等人[15]使用类似移动窗口的 探测器来探测人体并计算图像中人的数量。Haar小 波分类器[16]用于从检测到的人体中提取低级特 征,而[17]中则用 HOG(直方图定向梯度)分类器 来提取特征。Pedro F Felzenszwalb 等人[18]尝试检 测身体的一些特定部分而不是整体,因为人体在拥 挤的场景中总是被遮挡。但是所有这些早期工作在 非常拥挤的场景中都得到了较差的结果。

2.1.2 基于回归的方法

随着场景变得越来越拥挤,基于检测的方法存 在很大限制,因此基于回归的方法被提出。Antoni B Chan 等人[19]使用前景和纹理特征生成低级信息, 并在学习了人群与提取的特征相对应的关系后计 算出数量。随后在 2013 年,Idrees 等人[2]引入傅立 叶分析和 SIFT (尺度不变的特征变换)来提取[19] 中提出的特征。但是一些显著的特征很容易被忽 视,从而导致更大的偏差。在[20]中,局部区域中 的特征与其密度图之间的线性映射用来整合显著 性信息。2015年,由于理想线性映射增益的问题, Pham 等人[21]建议通过随机森林回归来学习非线 性映射而不是线性映射。

2.2 基于深度学习的方法

随着深度学习的快速发展,卷积神经网络在人 群计数领域显示出了很大的优势。

2015 年, Zhang 等人[22]训练卷积神经网络对 人群密度图进行回归。他们使用密度和透视信息重 新得到图像,然后使用它们微调训练好的网络并预 测密度图。然而,其适用性受到透视图的要求和每 个测试场景微调的限制。2016 年, Zhang 等人[9] 使用多尺度卷积神经网络架构来解决人群场景中 的大规模变化,并使用 1×1 卷积操作融合来自每 个特定尺寸的卷积网络训练的特征图以回归得到 密度图。它解决了尺度变化导致的问题。在此之后 后,多列[8]或多尺度[6.11.17]网络架构经常被用于 人群计数问题。具体而言, Sam 等人[7]引入了一个 分类器,根据密集级别选择指定的训练列。Cao 等 人[8]使用尺度融合模块作为编码器来提取不同尺 度的特征,并使用一组转置的卷积作为解码器来生 成高质量的密度图,还提出了局部模式一致性损失 函数。Zhang 等人[11]结合了多层的特征图来适应 行人规模和视角的变化,引入了多任务损失,增加 了相对人头数量损失函数。但是一些工作[10]则建 议用相同大小的卷积核堆替换不同的卷积核。Li 等 人[10]验证了使用多列卷积的有效性可能并不突 出,这种分支结构中的每一列学到的都是几乎相同 的特征。因此它使用 VGG16 作为基线,并在后端 引入了扩张层,得到了很大的改进。此外,[12,13] 在训练过程中增加了额外的信息以排除背景干扰。 Shi 等人[12]将透视信息整合到人群密度图中,提供 有关图像中人物尺度变化的附加信息,这十分有效 地提高了小尺寸的人群区域的密度回归的精度。Liu 等人[13]提出了一项自监督的任务以改进人群计数 网络的训练, 在训练时利用未标记的人群图像以显 著提高效果。它可以生成子图像的排名,其可以用 于训练网络来估计一个图像是否包含比另一个图 像更多的人。但额外的信息或任务可能会导致更多 的资源和计算量的需求。

在 2019 年,更多解决方案被提出。Qi Wang 等 人[23]构建了一个大尺度、多样化的合成人群计数 数据集来预先训练他们设计的空间全卷积网络。 Weizhe Liu 等人[24]引入了端到端架构,该架构结 合了使用多个大小的感受域得到的特征,并学习在 每个图像位置的每个特征的权重。Chenchen Liu 等人[25]将检测到的模糊的图像区域放大到高分辨率以进行重新训练,并添加了本地化任务。几乎所有方法都添加了额外的信息或任务来增强单一人群计数的任务。

3 主要方法论述

许多先前的方法引入了多列融合的网络结构, 以减少由于透视效应导致的头部尺度变化引起的 误差。它们可以融合各种不同尺寸的卷积核或不同 列的各种感受野的特征图。但是不同大小的内核可 能会导致更多的参数量和计算量,而多列架构可能 使网络更复杂。受[11]的启发,我们提出基于单一 大小卷积核的单列网络,通过通道注意机制融合低 层和高层的特征图。该网络对于头部尺度变化和背 景噪声将更具鲁棒性,同时保持结构的简洁。此外, 我们网络最后的部分利用扩张卷积模块,并且将基 于 SSIM 和基于回归人数的两个损失函数添加到综 合损失函数中。

我们提出的网络结构模型如图 1 所示,该模型 被称为 FF-CAM (Frontend-backend Fusion Network through Channel-attention mechanism)。我们将从四 个方面详细阐述该模型。

3.1 低层一高层融合

在人群场景的采集过程中,由于同一场景下人 与摄像机的距离不同,会因为透视效应导致人头大 小不同,也就是存在尺度多样性的问题。为了提取 不同尺度大小的特征,解决尺度多样性带来的问 题,并排除背景干扰,我们提出了低层一高层特征 图融合的方法。

如图 1 所示,我们网络的主干采用 VGG-16 结构,它具有强大的特征表示能力且易于连接。我们运用 VGG-16 的前 13 层来提取多尺度的特征图。 组成 FF-CAM 的所有卷积核大小均为 3×3 (除一个3×3卷积之前的1×1卷积用于降低计算复杂度和最后一层1×1卷积层用于代替全卷积层外),多个3×3 的卷积核堆叠与大尺度的卷积核具有相同的效果,例如2个3×3 的卷积核堆叠的效果相当于1个5×5 的卷积核,3个3×3卷积核则相当于1个7×7 的卷积核,以此类推。因此其可以学习到不同尺度的特征,但计算量要少得多,并且可以构建更深的网络。

网络中不同级别的特征层不仅包含不同级别

的语义信息,还包含不同的比例特征信息。低层可 以提取细节边缘图案,这对于在人群密度图中回归 拥塞区域的值具有重要意义。但它无法捕捉细节, 这可能会导致杂乱的背景干扰,从而导致不正确的 回归。高层则可以选择性地获得有用的语义信息, 因此网络可以将人群与背景噪声区分开来。

鉴于它们的特性,我们通过通道注意模块融合 低层和高层的特征图,以从主干网络中获取并融合 足够多的特征。

如图1所示,我们使用来自VGG-16主干网络 中的 Conv1 2, Conv2 2, Conv3 3, Conv4 3 和 Conv5 3 层的特征图,其中卷积层参数设置与 VGG-16 相同。这些不同层级特征图的输入有助于 提取多尺度的特征。通过最大池化层后,这些输出 特征图对应的大小分别为原始输入图像的 1/2.1/4.1/8 和 1/16。首先,使用最近邻插值对 Conv4 3 输出的特征图进行上采样,并与 Conv3 3 输出的特征图融合,再将融合后的特征图输入通道 注意力模块,调整两层不同特征信息融合时的权 重,提高网络的表征能力。随后,Conv5 3 输出的 特征图和 Conv2 2 输出的特征图的融合操作类似于 Conv4 3 和 Conv3 3, 融合得到的特征图同样输入 通道注意力模块。经通道注意力模块处理后的特征 图输入一组卷积层: Conv1×1×512, Conv3×3× 512 和 Conv3×3×512。3×3 卷积之前的 1×1 卷积 用于降低计算复杂度。我们将该组卷积层输出的特 征图定义为 Conv6 3 层, 其同样被上采样并与 Convl 2 的输出融合,然后以相同的方式输入到通 道注意力模块。最后,输出的特征图通过扩张卷积 模块后生成人群密度图。接下来我们将具体介绍通 道注意力模块和扩张卷积模块,具体结构如图2和 图 3。



图 2 通道注意力模块的结构图。其中, CONCAT 表示两层 特征图的融合,得到空间维数为*h×w×c*的特征图。

3.2 通道注意力模块

注意力模型现在已经成为神经网络中的一个 重要概念,在不同的领域中被研究和应用。[14]介 绍了 SE 模块,它模拟了卷积特征图的通道之间的 相互依赖性,从而提高了网络的表征能力。

由于大多数先前的工作直接组合来自不同卷 积层的特征图,没有考虑融合时它们各自的权重。 另一方面,由于空间信息的不足,卷积层的通道总 是被忽略。SE 模块可以进行特征重新校准,选择性 地强调有用信息,并且抑制不太有用的特征,网络 可以学习使用全局信息。此外,它还有助于捕获空 间相关性,而无需额外的监督。最后一点,它在计 算上很轻巧。有如此多的好处,它却只会略微增加 模型复杂性和计算负担。

此外, SE 模块已被证明可以改善网络性能,并可以通过整个网络进行累积[14]。因此,我们将 SE 块转换为我们的通道注意力模块。具体结构如图 2 所示。通道注意力模块包括三个过程:挤压S,激励 E 和重新缩放 R。

首先,两个卷积层融合后输出的特征图 N 进行 挤压操作 S。挤压操作在空间维度上聚合特征图, 并通过全局平均池化层来生成通道统计量。给定特 征 图 的 空 间 维 数 为 $h \times w \times c$,挤压操作后变为 $1 \times 1 \times c$ 。每一个通道的特征图 $N_x(x = 1, 2, ..., c)$ 对应 的通道描述符 D_x 由以下公式计算:

$$D_{x} = S(N_{x}) = \frac{1}{h \times w} \sum_{i=1}^{h} \sum_{j=1}^{w} n_{x}(i, j)$$
(1)

其中, $n_x(i, j)$ 表示特征图 N_x 上第 i 行第 j 列的元素的值。

特征图 N 通过挤压操作后生成了通道描述符 $D = \{D_x, x = 1, 2, ..., c\}$ 。通道描述符嵌入了通道特征 响应的全局分布,因此其较低层能够利用全局感受 野的信息。

然后,我们将 *D* 送入激励操作 *E*,产生提取描述符 *T*。它由基于非线性的两个完全连接层,一个 Relu 函数和一个 sigmoid 函数组成。将其表示为:

$$T = E(D; FC) = \sigma(g(D; FC))$$

= $\sigma(FC, \delta(FC, D))$ (2)

其中 FC_1 是具有缩小率 k 的降维层, FC_2 是维数增加层。 k 是一个超参数,它可以改变模型中块的容量和计算成本。根据[14],我们设置 k = 16,以实现准确性和复杂性之间的良好平衡。 δ 是 Relu函数, σ 是 Sigmoid 函数。两个完全连接层可以通过减小维度来限制模型复杂性,极大地减少了参数量和计算量。并且其能更多的学习通道之间的非线性相互作用,可以更好地拟合通道间复杂的相关性,提高泛化性。此外,与 one-hot 激活函数相反,Sigmoid 激活函数强调多个通道,故整个激励操作能完全捕获通道依赖性并控制每个通道的激励,获得 0—1 之间归一化的权重。

最后,通道注意力模块的输入*N*由提取描述符 *T*重新加权:

$$F = R(N;T) = T \cdot N \tag{3}$$

其中 *R* 表示输入特征图 *N* 和提取描述符 *T* 之间的 通道乘法,即通过乘法将 *T* 逐通道的权重加权到 *N* 中对应的每个通道特征图的每个特征点上,完成在 通道维度上的对原始特征的重标定。模块的最终输 出 *F* 可以直接被送入下一层。



图 3 扩张卷积模块的结构图。其中扩张卷积第一行的参数 分别表示卷积核大小和通道数。

3.3 扩张卷积模块

在我们的网络中,输入的人群图像由最大池化 层下采样再经上采样融合之后,生成的特征图为原 始输入的 1/2。特征图在经过池化层后,虽然控制 了过拟合同时保持了不变性,但降低了空间分辨 率,丢失了部分空间信息,产生的密度图质量不够 高。

Li 等人[10]证明了扩张卷积比使用卷积、池化 加反卷积的方案更好地保持特征映射的分辨率。虽 然反卷积层可以减轻信息的丢失,但会增加额外的 复杂性,且会导致执行延迟。基于此,我们在网络 的末端利用扩张卷积层。扩张卷积层扩大了感受 野,而不增加参数或计算量。同时,经过扩张卷积 的输出可以包含更详细的空间信息和全局信息,不 会降低空间分辨率。所以,我们运用扩张卷积可以 生成高质量的人群密度图,同时提高人群估计准确 率。

我们在网络的末端运用扩张卷积,如图3表示 网络末端的扩张卷积模块。它由具有扩张率为2的 四层扩张卷积层和一层 1×1 的卷积层组成。每个 扩张卷积层的通道数都不同,每一层后都会通过批 量标准化层和 Relu 层。1×1 卷积层用来输出最终 的人群密度图,相较于全连接层其参数量更少,计 算量更小。最后,网络输出高分辨率的人群密度图。

3.4 综合损失函数

主流工作将像素上的欧几里德损失设置为训 练过程中的损失函数。在人群场景中,高密度区域 的局部模式和纹理特征与其他区域(低密度区域或 背景)大不相同,但欧几里德的损失建立在像素独 立性假设上并忽略了它们,密度图的局部相关性未 被考虑。此外,该损失函数与用来衡量准确度的 MAE及 MSE 没有直接关系,也没有将输入图像的 全局计数错误考虑进去。为了解决上述问题,我们 将基于结构相似性指数(SSIM)的损失函数、基于 回归人数的损失函数与欧几里德损失相结合作为 我们的最终损失函数,该函数可用于估计人群密度 图和真值的局部一致性,并估计人群数量与真实人 数之间的差异,从而使综合后的损失函数更好的表 示训练产生的估计值与真实值间的差异,以生成高 质量的人群密度图,提高训练准确度。

3.4.1 欧几里德损失函数

欧几里德损失用于在像素级别上衡量输出密 度图与相应真值之间的差异,其定义如下:

$$L_{2}(\Theta) = \frac{1}{N} \sum_{i=1}^{N} ||F_{d}(I_{i};\Theta) - D_{i}||^{2} \qquad (4)$$

其中 Θ 表示网络训练时的一组参数, *N* 是训练样本的数量。 $F_a(I_i; \Theta)$ 表示具有参数 Θ 的网络输入图像 I_i 后输出的估计密度图, 而 D_i 是对应的真值密度 图。

3.4.2 基于 SSIM 的损失函数

SSIM 是一种广泛用于图像质量评估领域的指标。它根据局部模式(包括均值,方差和协方差) 计算两个图像之间的相似性。SSIM 值的取值范围 是[-1,1]。两个图像越相似,其值越大。当两个图像 相同时,它等于1。

受 SANet[7]启发,我们将 SSIM 加入损失函数。 首先,使用标准偏差为 1.5 的 11×11 归一化高斯核 来 估 计 局 部 统 计 量 。 然 后 , 权 重 由 $W = \{W(r) | r \in R, R = \{(-5,5), ..., (-5,5)\}\}$ 定义,其中 r为中心, R包含所有位置内核。因此,对于每个 位置t,计算密度图 F_d 和相应的真值D的局部统计 量。

首先计算 F_d 的局部均值 μ_E 和方差 σ_E^2 :

$$\mu_{F_{d}}(t_{F_{d}}) = \sum_{r_{F_{d}} \in R_{F_{d}}} W(r_{F_{d}}) \cdot F(t_{F_{d}} + r_{F_{d}})$$
(5)

$$\sigma_{\mathbf{F}_{d}}(t_{\mathbf{F}_{d}}) = \sum_{\mathbf{r}_{\mathbf{F}_{d}} \in \mathbf{R}_{\mathbf{F}_{d}}} W(\mathbf{r}_{\mathbf{F}_{d}}) \cdot \left[F(t_{\mathbf{F}_{d}} + \mathbf{r}_{\mathbf{F}_{d}}) - \mu_{\mathbf{F}_{d}}(t_{\mathbf{F}_{d}})\right]^{2} \tag{6}$$

其次,是D的局部均值 μ_p 和方差 σ_p^2 :

$$\mu_D(t_D) = \sum_{r_D \in R_D} W(r_D) \cdot F(t_D + r_D)$$
(7)

$$\sigma_{D}^{2}(t_{D}) = \sum_{r_{D} \in R_{D}} W(r_{D}) \cdot [F(t_{D} + r_{D}) - \mu_{D}(t_{D})]^{2}$$
(8)

由此我们可以计算 F_{d} 和D间的局部协方差 $\sigma_{F,D}$:

$$\sigma_{F_{t}D}(t) = \sum_{r \in \mathbb{R}} W(r) \cdot [F(t+r) - \mu_{F_{t}}(t_{F_{t}})] \cdot [Y(t+r) - \mu_{D}(t_{D})]$$
(9)

根据这些指标,SSIM 逐点计算如下:

$$SSIM = \frac{(2\mu_{F_d}\mu_D + Q_1)(2\sigma_{F_dD} + Q_2)}{(\mu_{F_d}^2 + \mu_D^2 + Q_1)(\sigma_{F_d}^2 + \sigma_D^2 + Q_2)}$$
(10)

其中, Q_1 和 Q_2 是随机的非常小的常数, 以避免被零除, 我们依照[7]的设置来给它们赋值。

最后,基于 SSIM 的损失函数定义为:

$$L_{s} = 1 - \frac{1}{M} \sum_{t}^{M} SSIM(t)$$
(11)

其中M 是密度图中的像素总数。

3.4.3 基于回归人数的损失函数

大多数基于密度估计的计数算法通过测量预 测密度图和地面实况密度图之间的每像素误差来 优化其计数模型。然而,这种方法与用来衡量准确 度的评估指标 MAE 和 MSE 没有直接关系,也没有 将输入图像的全局计数错误考虑进去。为此,我们 新增了另一个关于回归人数的损失函数,它直接衡 量估计人群数量与真实人数之间的差异。通过增加 该损失函数,网络将生成适合输入图像的整体密度 水平的特征,这有助于产生更准确的密度值。其定 义如下:

$$L_{c} = \left\| \hat{C} - C \right\|^{2} \tag{12}$$

其中, Ĉ和C分别是训练得到的人群数量和真 实的人群数量。

3.4.4 综合损失函数

将基于 SSIM 的损失函数和基于回归人数的损 失函数加入到训练过程中,最终的综合损失函数表 示如下:

$$L = L_2 + \alpha L_c + \beta L_s \tag{13}$$

其中 α 和 β 分别是基于回归人数的损失函数和基于 SSIM 的损失函数的权重,用作三个函数的平衡。 我们根据[7]的经验设定 β =0.001,在实验验证后设定 α =1,具体实验见第4.7节。

4 实验

我们的实验是在 4 块 TITAN Xp GPU 上进行 的。该网络基于 Pytorch 框架,我们使用 Adam 优 化器来优化参数并将原始学习速率设置为 le-5。参 数通过高斯分布随机初始化,平均值为零,标准差 为 0.01。除了输出层之外,我们还在每个卷积层之 后使用批量标准化层和 Relu 层,以提高训练速度并 有效地避免梯度的消失和爆炸。

4.1 真值的生成

现有的数据集一般都给定了原始图像以及其 对应的人群在图像中的坐标位置及总人数。和[9] 一样,我们同样用高斯自适应核来生成密度图的真 值。高斯自适应核的定义如下:

$$F(x) = \sum_{i=1}^{N} \delta(x - x_i) \times G_{\sigma_i}(x), \sigma_i = \beta \overline{d_i}$$
(14)

其中,在真值 δ 中,对于其中任意位置x和每一个 人头目标 x_i ,i=1,2,...,N,定义 $\delta(x-x_i)$ 是标准差 为 σ_i 的高斯核,而 d_i 是 k 个最近邻的平均距离。根 据[9]的经验,我们设置 $\beta=0.3$,k=3。对于每幅输 入的人群场景图像,高斯核可将其中所有标注的人 头模糊化,生成人群密度图的真值。

4.2 评估指标

大多数现有工作使用两个度量指标来衡量人 群计数的准确性,平均绝对误差(MAE)和均方误 差(MSE)。MAE 表示估计的准确性,而 MSE 反 映估计的鲁棒性。定义如下:

$$MAE = \frac{1}{N} \sum_{i=1}^{N} |F_{di} - D_i|$$
 (15)

$$MSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left| F_{di} - D_{i} \right|^{2}}$$
(16)

其中 N 是测试图像的数量, D_i 是第 i 个图像中的真实人群数, F_{di} 是第 i 个图像中的估计人群数。

4.3 在UCF_CC_50数据集上的实验

Idrees 等人[2]提出的 UCF_CC_50 数据集包括 50 个具有不同视角和分辨率的图像。这是一个非常 拥挤的数据集,平均人数达到了 1280 人,最多的 一幅图片中有 4543 人。由于包含各种人群场景且 图像总数有限,这是一个非常具有挑战性的数据 集。因此,我们按照[2]中的标准设置执行 5 倍交叉 验证,最大程度地利用样本:将数据集随机均分成 五等份,以其中的四份作为训练集,剩下的一份作 为测试集,共进行五次训练和测试,五次实验的结果如表1所示。最后再取误差指标的平均值作为实验的最终结果。

表 1 UCF_CC_50 数据集 5 倍交叉验证结果

测试集序号	MAE	MSE
1	383.65	579.99
2	144.33	183.25
3	293.15	337.64
4	257.15	317.47
5	155.54	192.51
均值	246.764	322.172

我们将结果与最先进的方法进行比较,表2中 列出了 MAE 和 MSE 比较的结果。我们的 FF-CAM 的估计误差 MAE 和 MSE 在所有模型中是最小的, 这表明我们得到了对 UCF_CC_50 数据集计数的最 佳估计,相比于效果最好的[10],我们的 MAE 提高 了 4.5%, MSE 提高了 3.8%。该结果验证了 FF-CAM 模型的准确性和鲁棒性。

表 2 UCF_CC_50 数据集的估计误差

方法	MAE	MSE
MCNN[9]	377.6	509.1
CMTL[8]	322.8	397.9
Switch-CNN[7]	318.1	439.2
SaCNN[11]	314.9	424.8
CSRNet[10]	266.1	397.5
FF-CAM	246.8	322.2

训练好的模型在 UCF_CC_50 数据集上得到的 部分密度估计图如图 4 所示。由图 4 可以看出,我 们的模型对极度拥挤的场景能进行很好的预测并 生成分布较为准确的密度图,且预测人数更接近真 实人数,好于 CSRNet[10]模型。由这些图可以看出, 第二张由于透视存在人头尺度大小不一的问题,而 得到的密度图很好的解决了该问题,在不同人头大 小的位置生成的密疏程度不一。第三张具有干扰的 楼房背景,而得到的密度图很好的排除了干扰,未 将其统计入人数。





图 4 FF-CAM 模型在 UCF_CC_50 数据集上的实验对比 密度图。第一行为原始图像,总人数分别为 1997, 2960, 1045。第二行为密度图的真值,第三行为 CSRNet[10]结构 得到的密度估计图,预测总人数分别为 2100, 3430, 1185。 第四行为 FF-CAM 得到的密度估计图,预测总人数分别为 2006, 2600, 1022。

4.4 在ShanghaiTech数据集上的实验

ShanghaiTech 数据集是一个多样且拥挤的数据 集,由Y. Zhang 等人[9]提出。该数据集包括 part A 和 part B 两部分,part A 是从网上收集而来,共有 482 张图片;part B 则是从上海的拥挤繁忙的街道 上收集而来,共有 716 张图片。两个部分都是十分 拥挤的数据集,part A 平均人数达到了 501 人,最 多的一幅图片中有 3139 人。而 part B 相对不那么 拥挤,平均人数为 124 人,最多的一幅图片中有 578 人。在 part A 数据集中,300 张图片用来训练,剩 下的 182 张则用来测试。part B 数据集里的 400 张 图片用来训练,316 张用于测试。

表 3 中列出了我们将估计结果的误差 MAE 和 MSE 与最先进的方法进行比较的结果。从表中可以 看出,我们的方法在 part B 数据集中测试的结果优 于其他的方法, MAE 和 MSE 分别提高了 2.8%和 1.3%。这说明我们的方法在 part B 数据集上表现得 很好,证明了 FF-CAM 的优越性。同时其在 part A 数据集上的 MSE 提高了 4.5%,说明模型的鲁棒性 较强。但 MAE 则略差于 CSRNet[10],这反映出我 们的方法可能需要更多的训练和实验来提高其预 测的准确性。

表 3	Shangha i Tech	数据集的估计误差
-----	----------------	----------

方法	Par	t A	Pa	rt B
	MAE	MSE	MAE	MSE
MCNN[9]	110.2	173.2	26.4	41.3
Switch-CNN[7]	90.4	135.0	21.6	33.4
SaCNN[11]	86.8	139.2	16.2	25.8

CSRNet[10]	68.2	115.0	10.6	16.0
FF-CAM	71.0	109.8	10.3	15.8

图 5 和图 6 展示了训练好的模型在 Shanghai Tech 数据集上估计得到的部分密度估计图。可以看出,我们的模型在这两个部分的数据集上都有较好的表现,生成了分布较为准确的密度图,预测的结果更接近于真值,且分辨率也较高。比较图 5 和图 6, ShanghaiTech part A 数据集财相对稀疏,这说明在极度拥挤的数据集上我们的网络还需要更多的图片进行训练以提高模型的准确度。



图 5 FF-CAM 模型在 Shanghai Tech A 数据集上的实验对比密 度图。第一行为原始图像,总人数分别为 239,1005,1174。 第二行为密度图的真值,第三行为 CSRNet[10]得到的密度 估计图,预测总人数分别为 379,741,1448。第四行为 FF-CAM 得到的密度估计图,预测总人数分别为 346,870, 1402。





图 6 FF-CAM 模型在 ShanghaiTech B 数据集上的实验对比密 度图。第一行为原始图像,总人数分别为 28,130,467。 第二行为密度图的真值,第三行为 CSRNet[10]得到的密度 估计图,预测总人数分别为 24,117,418。第四行为 FF-CAM 得到的密度估计图,预测总人数分别为 27,121,429。

4.5 在UCF ONRF数据集上的实验

UCF_QNRF 数据集由 H. Idrees[26]等人提出,同样是一个多样且拥挤的数据集,但图片总数量有1535 张,人的总数多达1251642,远多于其他两个数据集。其是从三个不同的数据集来源收集而来,包含了全球各个场景,且同时拥有拥挤和稀疏的人群场景。我们取1201 张图片用来训练,剩下的334 张则用于测试。

表 4 中列出了我们将估计结果的误差 MAE 和 MSE 与最先进的方法进行比较的结果。从表中可以 看出,我们方法的 MAE 提高了 13.3%,这说明预 测效果有了明显提升,估计误差较小。但 MSE 则 略逊于现有方法,可能是预测结果还不够稳定,存 在少量误差较大的图片。

方法	MAE	MSE
MCNN[9]	277	426
CMTL[8]	252	514
Switch-CNN[7]	228	445
H. Idrees et al[26]	132	191
FF-CAM	114.5	200.5

表 4 UCF_QNRF 数据集的估计误差

图 7 展示了训练好的模型估计得到的部分密度 估计图。可以看出,我们的模型对图 7 后两张的估 计值较 Switch-CNN[7]更为准确,且生成的密度图 的分布也更加精准,分辨率更高,这反映出我们模 型对拥挤和相对稀疏的场景都能进行很好的预测 并生成分布较为准确的密度图,较接近于真值。同 时我们可以看到三幅图都具有房屋和树木的背景 干扰,预测生成的密度图则避免了此干扰,进一步 验证模型的抗干扰性。但是,第一张图的估计值相 对真值有一定的偏差,是我们模型的测试中少量的 误差较大的图片,这也可以解释模型的 MSE 略逊 于现有方法。下一步需要更多的训练来提高模型的 鲁棒性,排除大的误差。



图 7 FF-CAM 模型在 UCF-QNRF 数据集上的实验对比密度 图。第一行为原始图像,总人数分别为 349,435,1017。 第二行为密度图的真值,第三行为 Switch-CNN[7]得到的密 度估计图,预测总人数分别为 365,477,1069。第四行为 FF-CAM 得到的密度估计图,预测总人数分别为 393,440,

 $1017\,{}_{\circ}$

4.6 消融实验

我们在 ShanghaiTech part A 数据集上进行了消融实验来验证 FF-CAM 结构的有效性,图 8 给出了消融实验的结果对比。

我们首先在 VGG-16 基线上进行了训练和测试的实验。从图 8 可以看出,FF-CAM 的估计误差明显优于 VGG-16 基线的结果。与 VGG-16 网络相比,FF-CAM 模型的 MAE 提高了 17.3%,MSE 提高了 12.7%,证明 FF-CAM 的网络结构很好的提高了预测精度。

随后我们在保持 FF-CAM 的其他结构不变时, 分别去掉其中的通道注意力模块,扩张卷积模块, 基于 SSIM 的损失函数和基于回归人数的损失函 数,进行训练并测试。每一个消融实验得到的 MAE 和 MSE 的对比如图 8。

在去掉所有的通道注意力模块后,模型的 MAE 下降了 11.4%, MSE 下降了 7.7%, 验证了通道注 意力模块对整个模型的增益。

在去掉扩张卷积模块后,模型的 MAE 下降了 7.6%, MSE 提高了 4.0%,证明了扩张卷积的有效 性。 相对于其他模块,基于 SSIM 的损失函数和基于回归人数的损失函数对整个模型的影响较小,但 去掉后模型的 MAE 和 MSE 也有所下降,说明其在 一定程度上提高了预测精度。具体来说,基于回归 人数的损失函数提高效果略高于基于 SSIM 的损失 函数。

消融实验结果表明,分别去掉各个模块后预测 精度都有一定的下降,这说明每个模块都对网络性 能有一定的提升作用,验证了我们提出的方法的有 效性和合理性。



4.7 参数实验

我们在 ShanghaiTech part B 数据集上对综合损 失函数中参数α的取值进行了消融实验,来得到最 优取值的参数有效性,图9给出了参数实验的结果 对比。

由图 9 可看出,误差评估指标 MAE 和 MSE 关 于不同参数 α 取值的曲线先递减后递增,当 α=1 时 误差最小,故取 α=1。



(b) MSE 结果对比图。

图 9 参数 α 消融实验的结果对比图。其中,横轴表示 α 的 取值变化,纵轴表示评估指标值的变化。

5 结论

在本文中,我们提出了一个用于人群计数的 FF-CAM 框架。它基于单列网络,仅利用单一大小 的卷积内核,但性能卓越。我们提出了主干网络的 低层和高层的特征图融合,然后输入通到道注意力 模块,最后将得到的特征图馈送到扩张卷积模块中 以产生高分辨率的密度图。我们的 FF-CAM 准确, 稳定,简洁,并具有良好的泛化能力。其在 UCF_CC_50 数据集和 ShanghaiTech part B 数据集 上的测试结果优于现有的方法。在接下来的工作 中,我们将在人群计数的其他公开数据集上进行训 练和测试,并与最先进的方法进行比较,以检验我 们提出的网络在不同的环境和疏密场景下的性能。

致 谢 这项工作得到了国家自然科学基金(No. 71673293)和国家自然科学基金(No. 61806215)的支持。

参考文献

- H. Idrees, K. Soomro, and M. Shah. Detecting humans in dense crowds using locally-consistent scale prior and global occlusion reasoning. IEEE transactions on pattern analysis and machine intelligence, 2015, 37(10):1986-1998.
- [2] H. Idrees, I. Saleemi, C. Seibert, and M. Shah. Multi-source multi-scale counting in extremely dense crowd images. Proceedings of the IEEE conference on computer vision and pattern recognition. Oregan, Portland, 2013: 2547-2554.
- [3] Jianxing Yang, Yuan Zhou, Sun-Yuan Kung. Multi-scale generative adversarial networks for crowd counting. Proceedings of the IEEE international conference on pattern recognition. Beijing, China, 2018: 1051-4651.
- [4] Greg Olmschenk, Hao Tang, Zhigang Zhu. Crowd counting with minimal data using generative adversarial networks for multiple target regression. Proceedings of the IEEE winter conference on applications of computer vision. lake tahoe, CA, USA, 2018: 1151-1159.
- [5] V. A. Sindagi and V. M. Patel. Generating high-quality crowd density maps using contextual pyramid cnns. Proceedings of the IEEE international conference on computer vision. Venice, Italy, 2017: 1879-1888.
- [6] X. Cao, Z. Wang, Y. Zhao, and F. Su. Scale aggregation network for accurate and efficient crowd counting. Proceedings of the European

conference on computer vision. Munich, Germany, 2018: 734-750.

- [7] D. B. Sam, S. Surya, and R. V. Babu. Switching convolutional neural network for crowd counting. Proceedings of the IEEE conference on computer vision and pattern recognition. Honolulu, HI, USA, 2017: 6.
- [8] V. A. Sindagi and V. M. Patel. Cnn-based cascaded multitask learning of high-level prior and density estimation for crowd counting. Proceedings of the IEEEinternational conference on advanced video and signal based surveillance. Lecce, Italy, 2017: 1-6.
- [9] Y. Zhang, D. Zhou, S. Chen, S. Gao, and Y. Ma. Single image crowd counting via multi-column convolutional neural network. Proceedings of the IEEE conference on computer vision and pattern recognition. Las Vegas, NV, USA, 2016: 589-597.
- [10] Y. Li, X. Zhang, and D. Chen. Csrnet: dilated convolutional neural networks for understanding the highly congested scenes. Proceedings of the IEEE conference on computer vision and pattern recognition. Salt Lake City, UT, USA, 2018: 1091-1100.
- [11] L. Zhang, M. Shi, and Q. Chen. Crowd counting via scale-adaptive convolutional neural network. Proceedings of the IEEE winter conference on applications of computer vision. Lake Tahoe, CA, USA, 2018: 1113-1121.
- [12] Miaojing Shi, Zhaohui Yang, Chao Xu, Qijun Chen. Revisiting perspective information for efficient crowd counting. Proceedings of the IEEE conference on computer vision and pattern recognition. Long Beach, CA, USA, 2019 :7271-7280.
- [13] X. Liu, J. van de Weijer, and A. D. Bagdanov. Leveraging unlabeled data for crowd counting by learning to rank. Proceedings of the IEEE conference on computer vision and pattern recognition. Salt Lake City, UT, USA, 2018: 7661-7669.
- [14] Jie Hu, Li Shen, Gang Sun. Squeeze-and-Excitation Networks. Proceedings of the IEEE conference on computer vision and pattern recognition. Salt Lake City, UT, USA, 2018: 7132-7141.
- [15] Piotr Dollar, Christian Wojek, Bernt Schiele, and Pietro Perona. Pedestrian detection: an evaluation of the state of the art. Proceedings of the IEEE transactions on pattern analysis and machine intelligence, 2012, 34(4):743-761.
- [16] Paul Viola and Michael J Jones. Robust real-time face detection. International journal of computer vision, 2004, 57(2):137-154.
- [17] Navneet Dalal and Bill Triggs. Histograms of oriented gradients for human detection. Proceedings of the IEEE conference on computer vision and pattern recognition. San Diego, CA, USA, 2005: 886-893.
- [18] Pedro F Felzenszwalb, Ross B Girshick, David McAllester, and Deva Ramanan. Object detection with discriminatively trained part-based models. IEEE transactions on pattern analysis and machine intelligence, 2010, 32(9):1627-1645.
- [19] Antoni B Chan and Nuno Vasconcelos. Bayesian poisson regression for crowd counting. Proceedings of the IEEE 12th international conference on computer vision. Kyoto, Japan, 2009: 545-551.
- [20] Victor Lempitsky and Andrew Zisserman. Learning to count objects in images. Proceedings of the advances in neural Information processing

systems. Cambridge, MA, USA, 2010: 1324-1332.

- [21] Viet-Quoc Pham, Tatsuo Kozakaya, Osamu Yamaguchi, and Ryuzo Okada. Count forest: co-voting uncertain number of targets using random forest for crowd density estimation. Proceedings of the computer vision IEEE international conference on IEEE computer society. Washington, DC, USA, 2015: 3253-3261.
- [22] C. Zhang, H. Li, X. Wang, and X. Yang. Cross-scene crowd counting via deep convolutional neural networks. Proceedings of the IEEE conference on computer vision and pattern recognition. Boston, MA, USA, 2015: 833-841.
- [23] Qi Wang, Junyu Gao, Wei Lin, Yuan Yuan. Learning from synthetic data for crowd counting in the wild. Proceedings of the IEEE conference on computer vision and pattern recognition. Long Beach, CA, USA, 2019: 8190- 8199.
- [24] Weizhe Liu, Mathieu Salzmann, Pascal Fua. Context-aware crowd counting. Proceedings of the IEEE conference on computer vision and pattern recognition. Long Beach, CA, USA, 2019: 5099-5108.
- [25] Chenchen Liu, Xinyu Weng, Yadong Mu. Recurrent attentive zooming for joint crowd counting and precise localization. Proceedings of the IEEE conference on computer vision and pattern recognition. Long Beach, CA, USA, 2019: 1217-1226.
- [26] H. Idrees, M. Tayyab, K. Athrey, D. Zhang, S. Al-Maadeed, N. Rajpoot, and M. Shah. Composition loss for counting, density map estimation and localization in dense crowds. Proceedings of the European conference on computer vision. Berlin/Heidelberg, Germany, 2018: 544-559.



Yuqian Zhang, born in 1996. She is a postgraduate student. She majors in computer vision, deep learning and information systems engineering.

Guohui Li, born in 1963. He is adoctor, a professor and a doctoral supervisor. He majors in computer vision, information system engineering, data mining, and virtual reality technology.

Jun Lei, born in 1989. He is a doctor and a lecturer. He majors in computer vision, deep learning, data mining, and virtual reality technology.

Jiayu He, born in 1997. He is a postgraduate student. He majors in deep learning, data mining, and virtual reality technology.

Background

The paper focuses on the crowd counting in single images

in computer vision. Nowadays in the top computer vision conferences, novel frameworks are proposed to solve challenges and improve the estimation accuracy on the common datasets. State-of-the-art works are introduced in the paper. Our paper improves the estimation errors on one dataset, superior to state-of-the-art, and the estimation errors on other two datasets also get great results.

The study is supported by the National Natural Science Foundation of China (No. 71673293) and (No.61806215).

The two National Natural Science Foundation of China focus on the Analysis of crowd group behavior in public complex place. It is the most concern issue in public security management. However, the crowd behavior in public open area is complex. Crowd behavior is various in different scenarios, thus it is difficult to model them directly. Our previous research found that the crowd behavior can be described by crowd collectiveness when the public safety is considered. The crowd behavior differences in different scenarios and the difficult problems in unified modeling can be resolved by extracting the general crowd collectiveness. We regard the crowd system in public place as a complex system. From the perspective of visual data observation, we investigate the methods to mining and discover the universal crowd collectiveness in public complex place. In addition, we analyze the evolution of the crowd collectiveness, which present the generation mechanism of crowd event and abnormal behaviors. According to this idea, we defined four collectiveness that can be quantitatively descripted and measured: dynamic crowd collectiveness, static crowd collectiveness, conflictive-ness in crowd collectiveness, and stability in crowd collectiveness. The collectiveness is measured by the new graph-based learning method. Robust collectiveness map is generated by multi-view learning method. By analyzing the changes of crowd collectiveness in the timeline, the occurrence and evolution of crowd collectiveness event can be represented in spatio-temporal dimension. By analyzing the evolution of crowd collectiveness, we also research the problem of exploring abnormal crowd collectiveness motion and the recognition problem of different crowd collectiveness motions.

This paper collects different crowd scenarios, and can provides the number of people and the density of the scenario in single images. It can help us to analyze the evolution of the crowd collectiveness and the crowd behavior in public open area.