

基于深度图的人体动作分类自适应算法

蒋韦晔, 刘成明

(郑州大学 软件学院 河南 郑州 450002)

摘要: 由于深度相机成本的降低,越来越多的研究人员使用 RGB-D(red, green, blue and depth) 视频进行人类动作识别(human activity recognition, HAR)。使用深度运动图的局部二值模式进行特征提取,利用自适应差分进化极限学习机(self-adaptive differential evolution extreme learning machine, SaDE-ELM)用于动作分类,其中隐藏节点的学习参数通过自适应差分进化的方法进行修改。为了验证所提出方法的有效性,用 3 个公共数据集(MSR Action3D, MSRDaily Activity3D, MSRGesture3D)进行了实验。仿真结果表明,该方法优于基于内核的极限学习机(kernel extreme learning machine, KELM)的方法。

关键词: 人类动作识别; 深度运动图; 差分进化; 自适应差分进化极限学习机

中图分类号: TP391

文献标志码: A

文章编号: 1671-6841(2021)01-0016-06

DOI: 10.13705/j.issn.1671-6841.2019584

0 引言

HAR 可以覆盖多种应用,例如基于内容的视频分析和检索、视觉监视和人机交互。对 HAR 的初步研究基于 RGB 摄像机。然而,由于 RGB 摄像机的灵敏度受混乱背景和光照条件的影响,其鲁棒性较差。低成本的深度相机对基于深度 HAR 领域的研究产生了深远影响,在该领域中,相机拍摄的图像对照明变化和背景杂乱不敏感,从而增强了动作识别效果。

对于动作识别问题,深度运动图(deep motion map, DMM)会从输入深度数据中捕获形状和运动提示,然后将数据从 3D 转换为 2D。DMM 是 3 个正交笛卡尔平面上的投影深度框架之间的累积差异。文献[1]中,深度图像通过投影和累积来生成 DMM 定向梯度直方图(DMM-histogram of oriented gradients, DMM-HOG)特征描述符^[1]。在 DMM-HOG 中,对均匀间隔的单元格计算梯度方向直方图,并执行局部对比度归一化,以减少类内差异。文献[2]使用了深度运动图的局部二值模式(deep motion map-local binary pattern, DMM-LBP)特征描述符,其性能优于基于 DMM-HOG 的特征描述符^[2]。LBP 运算简单,并且具有旋转不变性。DMM 和 LBP 的组合提供了有效的基于块的重叠功能。特征提取后,可以使用主成分分析(principal component analysis, PCA)减小特征尺寸。对于动作分类,通常使用 KELM 方法,与支持向量机的多类别分类相比,它具有更好的泛化性能,并且训练时间较少, KELM 模型阶数随样本数量线性增长^[3]。文献[4]使用 ELM 方法进行动作分类。在 ELM 的训练阶段未修改隐藏参数,因而可能存在许多节点,这些节点在最小化目标函数方面的贡献较小。通过优化 ELM 的网络参数可以提高它的性能。近年来,差分进化(differential evolution, DE)在增强 ELM 的性能方面已获得普遍应用。基于 DE 和 ELM,文献[5]提出了一种称为进化极限学习机(evolutionary-extreme learning machine, E-ELM)的新算法,在该算法中,DE 用于网络参数优化,并使用 ELM 算法计算网络的输出权重。在 E-ELM 中,必须通过实验手动选择 DE 的试验载体生成策略和控制参数。选择不适当的策略和控制参数值可能会对网络泛化性能产生不利影响。自适应的差分进化极限学习机(self-adaptive differential evolution extreme learning machine, SaDE-ELM)是一种无梯度方法,是 ELM 的改进版本。在自适应差分进化算法中,使用 DE 对 ELM 的隐藏节点参数进行优化,并以自适应方式选择策略和控制参数。本文

收稿日期:2019-12-26

基金项目:国家自然科学基金青年基金项目(6140240)。

作者简介:蒋韦晔(1995—),男,硕士研究生,主要从事人体动作识别研究,E-mail:837384806@qq.com;通信作者:刘成明(1979—),男,副教授,主要从事计算图形学、数字图像处理研究,E-mail:cmliu@zzu.edu.cn。

结合特征选择方法、DMM-LBP 的优点和 SaDE-ELM 分类器的自适应特性,提出了一种改进的基于深度的人体行为分类方法。

1 本文方法

HAR 任务包括根据给定的输入视频自动识别动作标签。图 1 显示了基于深度数据集的 HAR,首先,使用 DMM 和 DMM-LBP 从深度输入中提取特征,然后通过 PCA 优化特征,使用 SaDE-ELM 方法对动作分类。我们将所提出的方法与 KELM 分类的性能在数据集 MSR Action3D^[6],MSRDaily Activity3D^[7]和 MSRGesture3D^[8]上进行了比较。

1.1 特征提取

特征提取是将时间图像序列转换为分类器可用的特征集的过程。一般采用诸如深度图、3D 关节直方图、DMM-HOG 和基于骨骼的特征之类的特征提取方法,其中 DMM-LBP 是更加紧凑而有用的特征。对于 N 帧给定的深度视频,二维投影视图 m_{front} 、 m_{side} 、 m_{top} 分别对应于正面视图、侧面视图和顶视图。这些图是通过将视频帧投影到 3 个正交的笛卡尔平面上而生成的。计算方法为

$$DMM_{\{front,side,top\}} = \sum_{k=1}^{N-1} |m_{\{front,side,top\}}^{k+1} - m_{\{front,side,top\}}^k|, \text{ 其中 } k \text{ 是帧索引。}$$

LBP 是基于旋转不变纹理的算子。为了产生 LBP,对具有中心像素的邻居执行阈值化。如图 2 所示,用于两个手波动作的 DMM 序列正投影 (DMM_{front}) 和 LBP 编码图像。PCA 由于其简单性和低噪声敏感性而被广泛用于处理高维数据^[9],它是用原始数据的大部分方差来代替低维特征子空间。

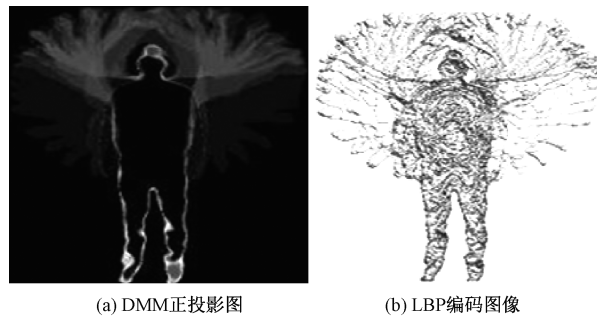


图 2 两个手波动作的 DMM 的正投影图像和 LBP 编码图像

Figure 2 DMM_{front} image and LBP coded image for the two hands wave

1.2 特征分类

分类器的目标是根据训练数据集将特征映射到动作标签。

1.2.1 ELM 分类器 ELM 对隐藏节点参数进行随机选择,使用等式 $\beta = L^+ T$ 计算输出权重 β ,其中 T 表示目标输出; L 是隐藏层输出矩阵, $L^+ = L(L^T L)^{-1}$ 是摩尔彭罗斯 L 的广义逆矩阵。

用于训练样本的单层前馈神经网络的通用 ELM 架构如图 3 所示。其输入特征标记为 $\{x_i, y_i\}_{i=1}^a, I_q (q = 1, \dots, m)$ 表示输入层节点, $N_r (r = 1, \dots, n)$ 表示隐藏层节点,而 $O_v (v = 1, \dots, c)$ 表示输出层节点。隐藏层输出矩阵由 $L = [h(x_1) \ h(x_2) \ \dots \ h(x_i)]^T$ 表示,其中 $h(x_q) = g(x_q, w_r, b_r)$, $g(\cdot)$ 是非线性激活函数, w_r 表示将第 r 个隐藏节点连接到第 q 个输入节点的权重向量, b_r 表示第 r 个隐藏节点的偏置。

当特征映射 $h(x_q)$ 未知时,使用内核 ELM,ELM 的内核矩阵为

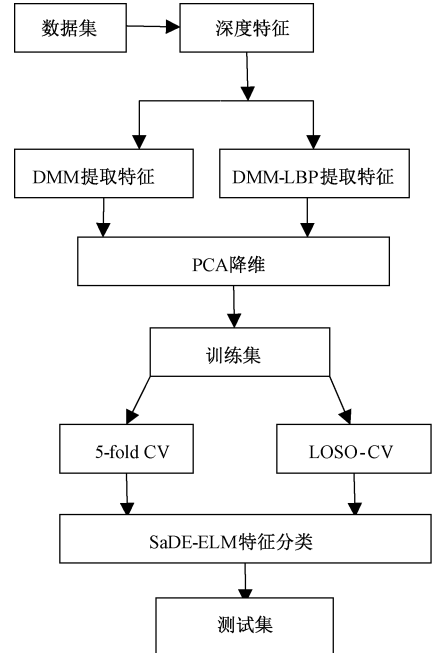


图 1 SaDE-ELM 算法流程图

Figure 1 Flow chart of SaDE-ELM algorithm

$$\lambda_{ELM} = LL^T : \lambda_{EMT_{i,j}} = h(\mathbf{x}_i) \dot{h}(\mathbf{x}_j) = R(\mathbf{x}_i, \mathbf{x}_j), \quad (1)$$

对于随机输入的样本 \mathbf{x}_i 和 $\mathbf{x}_j, R(\mathbf{x}_i, \mathbf{x}_j)$ 表示给定输入的径向基函数(RBF)。

KELM 算法步骤如下。

输入：动作识别训练集 T ；测试集 T' ，内核函数 Ω ，常数 C 。

输出：KELM 网络的输出用于测试集 T' 的输出。

1) 根据文献[10]中的方法初始化内核参数。

2) 使用式(1)计算隐藏层输出。

3) 使用公式 $f_r(\mathbf{x}_i) = [R(\mathbf{x}_i, \mathbf{x}_1) \cdots R(\mathbf{x}_i, \mathbf{x}_a)] (\frac{1}{C} +$

$\lambda_{ELM})^{-1} T$ 计算输出。

4) 将标签分配给测试样品。

1.2.2 SaDE-ELM 分类器 对于全局优化,DE 是进化算法类别中的强大技术。由于 DE 使用简单、控制参数数量少、空间复杂度低、性能高和收敛速度快等优点,该方法已被广泛用于 SLFN 的参数优化。通过使用自适应差分进化算法^[11],可以避免在 DE 中手动选择试验向量生成策略及其相关的控制参数。

SaDE-ELM 算法步骤如下。

输入：数据集,总体大小,隐藏节点数。

输出：权重和最小误差。

初始化：设置索引代数 $G=0$,随机初始化 NP 个个体。

1) 使用公式 $pr_{l,G} = \begin{cases} 1/4 & \text{if } G \leftarrow LP \\ S_{l,G} / \sum_{l=1}^4 S_{l,G} & \text{otherwise} \end{cases}$, 为每个目标向量选择一种策略。

2) 生成一个新的组,其中每个试验向量都是基于目标向量生成的策略。

3) 如果有任何变量超出定义的边界,则重新初始化试验向量。

4) 用公式 $\theta_{k,G+1} = \begin{cases} u_{p,G+1} & \text{if } RMSE_{\theta_{p,G}} - RMSE_{u_{p,G+1}} > \epsilon \cdot RMSE_{\theta_{p,G}}, \\ & \text{if } |RMSE_{\theta_{p,G}} - RMSE_{u_{p,G+1}}| < \epsilon \cdot RMSE_{\theta_{p,G}}, \\ & \text{and } \|\beta_{u_{p,G+1}}\| < \beta_{\theta_{p,G}} \text{,} \\ \theta_{p,G} & \text{otherwise,} \end{cases}$ 最小化 ELM 的目标函数。

公式 $\theta_{p,G} = [\omega_{1,(p,G)}^T, \dots, \omega_{n,(p,G)}^T, b_{1,(p,G)}, \dots, b_{n,(p,G)}]$ 给出了第一代总体 θ 的初始化,其中:权重 ω 和偏差 b 是 NP 向量的随机分配参数; n 表示隐藏节点的数量; $p=1,2,\dots,NP$; G 表示代数。

输出权重 $\beta_{p,G}$ 用式 $\beta_{p,G} = L_{p,G}^+ T$ 计算, $L_{p,G}^+$ 是 $L_{p,G}$ 的广义逆矩阵。RMSE 用式

$$RMSE_{p,G} = \sqrt{\frac{\sum_{i=1}^a \left\| \sum_{j=1}^n \beta_{jG}(\omega_{(j,(p,G))}, b_{(j,(p,G))}, I_i) - t_i \right\|^2}{c * a}}$$

来计算,其中: a 表示训练样本的数量; c 表示总类数和; t_i 表示目标输出。

根据概率 pr 从目标向量生成实验向量,对于选择策略 $l(l = \{1,2,3,4\})$, $S_{l,G}$ 由下式给出,

$$S_{l,G} = \left(\sum_{g=G-LP}^{G-1} S_{l,g} \right) / \left(\sum_{g=G-LP}^{G-1} SS_{l,g} + \sum_{g=G-LP}^{G-1} f_{l,g} \right) + \epsilon,$$

其中: l 表示策略; g 表示生成; $f_{l,g}$ 表示在下一代中丢弃的 g^{th} 和 l^{th} 策略处的试验向量的数量。

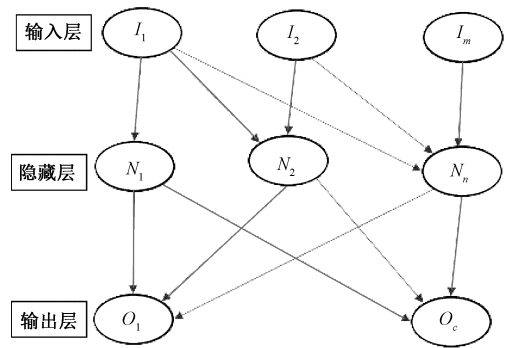


图 3 SLFN 的极限学习机架构

Figure 3 Extreme learning machine architecture of SLFN

2 实验设置及实验分析

实验是在 MSR Action3D,MSRDaily Activity3D 和 MSRGesture3D 公开获得的数据集上进行的。

2.1 实验设置

所有实验均在 Matlab 2016b 上进行。DMM-LBP 参数设置类似于文献[10]。在提取 DMM 并计算 LBP 之后,使用 PCA 降低数据集的维数。表 1 给出了平均特征大小和降维后的平均特征大小。分类任务由 SaDE-ELM 算法执行。SaDE-ELM 方法的参数选择为 20,下界和上限设置为 -1 和 1,最大世代数为 60,隐藏节点数为 80。专门针对特定的训练数据进行了 5 折交叉验证和“遗漏一名受试者”(LOSO)。在 5-fold CV(5 折交叉验证)中,从 10 名随机选择的受试者中选择 5 名进行训练,休息时进行测试。

表 1 特征描述

Table 1 Feature description

数据集	特征名称	交叉验证	平均特征大小	降维后平均特征大小
MSR Action3D	DMM	LOSO	600	508
MSR Action3D	DMM-LBP	LOSO	600	501
MSR Action3D	DMM	5-fold	600	272
MSR Action3D	DMM-LBP	5-fold	600	271
MSRDaily Activity3D	DMM	LOSO	322	165
	DMM-LBP	LOSO	312	164
MSRDaily Activity3D	DMM	5-fold	200	94
	DMM-LBP	5-fold	315	153
MSRDaily Activity3D	DMM	LOSO	350	300
	DMM-LBP	LOSO	350	300
MSRDaily Activity3D	DMM	5-fold	350	160
	DMM-LBP	5-fold	350	316

2.2 实验结果及分析

混淆矩阵用于根据精确率、召回率、 f 值和准确率评估所提出方法的性能,计算公式为 $Precision = TP / (TP + FP)$, $Recall = TP / (TP + FN)$, $f\text{-score} = (2 * precision * recall) / (precision + recall)$, $Accuracy = (TP + TN) / (TP + TN + FP + FN)$, 准确率和召回率是使用“真阳性”(TP)、“假阳性”(FP)和“假阴性”(FN)来计算的,这些值指示正确分类的操作数。被错误标记为肯定类别的操作数,动作应分别标记为肯定类,但错误地标记为否定类的操作数。

表 2 和表 3 分别提供了 MSRGesture3D、MSR Action3D 和 MSRDaily Activity3D 数据集的 5-fold CV(5 折交叉验证)和 LOSO 结果。在表 2 和表 3 中,EAP1、P-AP1、E-AP2 和 P-AP2 对应于由 DMM-KELM、DMM-SaDE-ELM,现有方法 DMM-LBP-KELM 和提出的方法 DMM-LBP-SaDE-ELM。我们对方法 E-AP1 和 P-AP1, E-AP2 和 P-AP2 进行了成对比较。从结果中可以看出,与现有方法相比,该方法的中值精确率提高了 1%~15%。我们使用 Wilcoxon 符号秩检验对获得的结果进行了统计分析。对于成对比较(E-AP1、P-AP1、E-AP2 和 P-AP2),测试数据集 DS1(MSRGesture3D)、DS2(MSR Action3D)和 DS3(MSRDaily) Activity3D)得出的 p 值小于 0.05,表明我们提出的方法比现有方法具有更好的性能。所提方法的性能与文献[10]中提出的方法在精确率、准确率、召回率和 f 值方面进行了比较。表 2 和表 3 列出了 3 个基于深度模态的数据集(包括手势、动作和活动)的这些性能指标的最小值、最大值、平均值和中值。

考虑到表 2 和表 3 的中位数精度测度,提出的仅使用带有 SaDE-ELM 分类器的 DMM 功能的方法的性能要优于 DMM-ELM 方法。结合 DMM 和 LBP 功能比仅 DMM 功能具有显著的性能改进,因为与 DMM 功能相比,它可以生成紧凑的功能,并可以增强图像的边缘。与使用 ELM 分类器相比,使用 SaDE-ELM 方法优化前馈网络的隐层权重和偏差,可以使准确性性能的中位数提高 1%~15%(使用成对比较)。这是因为初始化参数是在 ELM 算法中随机生成的,这对输出权重矩阵有重大影响。为了改进 ELM 算法,采用优化的 DE 算法(SaDE-ELM)自适应地收集 DE 算法的未知参数,从而提高了 HAR 的准确性。

表2 5-fold CV 下的实验结果

Table 2 Experimental results under 5-fold CV

单位:%

方法		准确率				精确率				召回率				f值			
5-fold		最小值	最大值	平均值	中值	最小值	最大值	平均值	中值	最小值	最大值	平均值	中值	最小值	最大值	平均值	中值
E-AP1	DS1	53	77	68	69	56	80	71	71	53	77	68	68	59	81	73	73
E-AP1	DS1	57	81	71	72	59	83	74	74	58	81	71	72	62	84	75	75
E-AP2	DS1	72	94	82	81	72	95	83	88	70	94	81	88	75	95	83	89
E-AP2	DS1	73	95	86	87	75	95	87	88	71	94	85	89	73	94	86	89
E-AP1	DS2	68	87	79	81	70	89	82	82	68	86	77	78	79	89	84	88
E-AP1	DS2	72	90	80	83	75	90	82	82	72	89	80	88	79	91	85	85
E-AP2	DS2	76	95	85	88	79	95	87	89	76	94	85	85	84	96	89	90
E-AP2	DS2	82	95	87	88	81	95	88	89	79	94	87	89	84	97	89	88
E-AP1	DS3	21	48	35	36	13	54	31	31	17	38	28	29	39	71	57	58
E-AP1	DS3	33	54	44	47	25	64	40	39	28	48	37	38	43	77	60	60
E-AP2	DS3	49	56	49	54	37	59	48	47	41	55	49	48	58	80	68	67
E-AP2	DS3	50	58	52	53	44	60	51	53	45	58	52	51	58	81	70	71

表3 LOSO CV 下的实验结果

Table 3 Experimental results under LOSO CV

单位:%

方法		准确率				精确率				召回率				f值			
LOSO		最小值	最大值	平均值	中值	最小值	最大值	平均值	中值	最小值	最大值	平均值	中值	最小值	最大值	平均值	中值
E-AP1	DS1	50	94	77	78	54	95	77	78	50	94	77	78	63	94	86	90
E-AP1	DS1	67	97	85	84	72	98	86	88	67	97	85	85	74	98	91	92
E-AP2	DS1	68	100	87	88	68	100	88	88	68	100	87	86	81	100	94	94
E-AP2	DS1	77	100	94	96	75	100	94	97	77	100	95	95	84	100	97	98
E-AP1	DS2	64	95	84	86	66	91	83	88	61	92	83	86	82	98	92	93
E-AP1	DS2	79	100	93	94	74	100	92	96	77	100	92	94	88	100	96	97
E-AP2	DS2	81	98	92	92	81	99	91	91	82	98	91	91	92	99	96	96
E-AP2	DS2	88	100	97	96	88	100	95	98	89	100	95	97	94	100	98	99
E-AP1	DS3	31	58	46	47	26	51	40	39	31	57	46	46	63	87	77	77
E-AP1	DS3	56	66	62	62	52	69	62	63	56	66	62	62	74	91	82	82
E-AP2	DS3	47	68	60	62	44	74	58	58	53	68	61	62	75	94	85	85
E-AP2	DS3	67	80	70	73	65	80	71	68	67	81	73	75	81	93	87	87

3 结束语

本文提出了一种通用的 HAR 方法,该方法已应用于动作、活动和手势等模态。针对基于深度的视频提出了一种改进的动作识别方法,包括特征提取、特征组合、特征选择和动作分类。由于 DMM-LBP 特征的紧凑性和旋转不变性,我们采用了 DMM-LBP 特征的组合。对于特征选择,我们使用了主成分分析(PCA)。对于动作分类,我们使用了 SaDE-ELM 方法。在 SaDE-ELM 方法中,采用基于 DE 的自适应控制参数,优化前馈网络的隐节点参数,并利用 ELM 算法确定网络输出权值。该方法不仅利用了 ELM 较好的泛化性能,而且利用自适应 DE 为 SLFN 获得了合适的权重和偏差。我们在 3 个公开可用的数据集上将该方法与 KELM 方法的性能进行了比较。

由于多模态数据将基于人和对象之间的交互来细化动作分类任务,因此可以将来自输入视频的多模态数据(骨骼+深度)用于动作识别,作为将来的工作,由于深层网络会自动从视频中提取特征,因此使用深层 CNN 进行人体行为分类,可以被认为是另一项未来的工作。该工作还可以扩展到数据集,包含从多个摄像机捕获的多个视图。

参考文献:

- [1] YANG X, ZHANG C, TIAN Y. Recognizing actions using depth motion maps-based histograms of oriented gradients[C]//Proceedings of the 20th ACM International Conference on Multimedia. New York, 2012: 1057-1060.

- [2] 程万里. 基于深度数据特征融合的人体动作识别[D]. 郑州: 郑州大学, 2018.
CHENG W L. Human action recognition based on depth data feature fusion[D]. Zhengzhou: Zhengzhou University, 2018.
- [3] WONG C M, VONG C M, WONG P K, et al. Kernel-based multilayer extreme learning machines for representation learning[J]. IEEE transactions on neural networks and learning systems, 2018, 29(3): 757-762.
- [4] 王杰, 刘向晴. 彩色图像分割的FCM预分类核极限学习机方法[J]. 郑州大学学报(理学版), 2018, 50(2): 75-80.
WANG J, LIU X Q. FCM pre-classification kernel extreme learning machine algorithm of color image segmentation[J]. Journal of Zhengzhou university(natural science edition), 2018, 50(2): 75-80.
- [5] ZHU Q Y, QIN A K, SUGANTHAN P N, et al. Evolutionary extreme learning machine[J]. Pattern recognition, 2005, 38(10): 1759-1763.
- [6] MSRAction3D dataset. (2018-01-05)[2019-02-10]. <http://research.microsoft.com/enus/um/people/zliu/actionrecorsrc/>.
- [7] MSRDailyActivity3D dataset. (2018-09-07)[2019-02-10]. http://users.eecs.northwestern.edu/~jwa368/my_data.html.
- [8] DING W W, LIU K, FU X J, et al. Profile HMMs for skeleton-based human action recognition[J]. Signal processing: image communication, 2016, 42: 109-119.
- [9] KARAMIZADEH S, ABDULLAH S M, MANAF A A, et al. An overview of principal component analysis[J]. Journal of signal and information processing, 2013, 4(3): 173-175.
- [10] CHEN C, JAFARI R. Action recognition from depth sequences using depth motion maps-based local binary patterns[C]//2015 15th IEEE Winter Conference on Applications of Computer Vision. Waikoloa, 2015:1092-1099.
- [11] NUNES U M, FARIA D R, PEIXOTO P. A human activity recognition framework using max-min features and key poses with differential evolution random forests classifier[J]. Pattern recognition letters, 2017, 99: 21-31.

Adaptive Algorithm for Human Motion Classification Based on Depth Map

JIANG Weiye, LIU Chengming

(School of Software, Zhengzhou University, Zhengzhou 450002, China)

Abstract: Due to the cost reduction of depth camera, more and more researchers began to use RGB-D (red, green, blue and depth) video for human activity recognition (HAR). The local binary pattern of deep motion map (DMM-LBP) was used for the extraction of features; a self-adaptive extreme learning machine algorithm based on differential evolution algorithm (SaDE-ELM) was used for the classification of actions; and the learning parameters of hidden nodes were also modified by SaDE-ELM. In order to verify the effectiveness of the proposed method, experiments were conducted in three public datasets (MSR Action3D, MSR DailyActivity3D and MSR Gesture3D), and the results of simulation showed that this method was better than the kernel extreme learning machine (KEML).

Key words: human action recognition; deep motion map; differential evolution; self-adaptive differential evolution extreme learning machine

(责任编辑:方惠敏 孔 薇)