



Available online at www.sciencedirect.com

ScienceDirect

Procedia Computer Science

Procedia Computer Science 89 (2016) 759 – 763

Twelfth International Multi-Conference on Information Processing-2016 (IMCIP-2016)

3-D Projected PCA based DMM Feature Fusing with SMO-SVM for Human Action Recognition

M. Naveenkumar* and A. Vadivel

National Institute of Technology Trichy, Tamilnadu 620 015, India

Abstract

Action recognition in video sequence is a very important and challenging problem yet. This paper presents an efficient feature extraction method for human action recognition for depth video sequence. For the video sequence acquired by depth sensor, all 3-D projections (xy, yz and zx) are calculated for each depth frame. For each projection view, the difference between each alternative frames have been considered to form the Depth Motion Map (DMM). Principle Component Analysis technique is applied to decrease the facet of DMM-feature. Sequential minimal optimization (SMO) is pre-owned to train the Support Vector Machine (SVM). The proposed approach is evaluated on MSR Action-3D data set and compared with the existing approaches. The empirical results convey that proposed approach achieves good results than the existing approaches.

© 2016 The Authors. Published by Elsevier B.V. This is an open access article under the CC BY-NC-ND license (http://creativecommons.org/licenses/by-nc-nd/4.0/).

Peer-review under responsibility of organizing committee of the Organizing Committee of IMCIP-2016

Keywords: Depth Motion Map; Human Action Recognition; SMO-SVM.

1. Introduction

Recognition of human actions in a video sequence with a conventional camera poses a challenging task. Analysis of conventional videos for action recognition has many limitations such as more computation power and also lack of 3D action data. When low cost depth cameras such as Asus Xtion PRO LIVE¹ and Kinect² are released in market, research interest on action recognition has been increased. Depth cameras provide the 3D depth data of the object and it will be helpful for segmentation and action recognition. Moreover, depth data is insensitive to illumination changes and also provide the shape information so that it will be helpful for recognition of human actions. The real time applications like smart home and video surveillance applications incorporate the theory of action recognition. There are two main challenges for the human action recognition: 1) description and modeling of the action (Feature Extraction) and 2) Classification of actions.

The contributions in the paper are as follows. Initially in view of computing an efficient technique for feature extraction for human action recognition, fusing sequential minimal optimization learning algorithm is taken for reference. And then the computed results are made suitable for the action inputs taken by depth camera. MSR Action-3D data set is taken for consideration in evaluating the proposed work. The data set comprise 20 action types with 10 subjects. Some of the sample collection from data set is depicted in the Fig. 3. From observations, a positive

Peer-review under responsibility of organizing committee of the Organizing Committee of IMCIP-2016 doi:10.1016/j.procs.2016.06.053

^{*}Corresponding author. Tel.: +91 -9791238410. *E-mail address:* mnaveenmtech@gmail.com

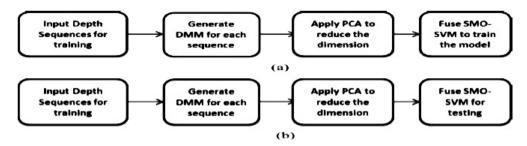


Fig. 1. Work Flows for the Proposed Work: (a) Training Work Flow; (b) Testing Work Flow.

response is seen when compared with the existing ones. The remaining portion in the paper is structured in following manner. The relevant work is reviewed in Section two. The proposed method is delineated in Section three and allied results are reported in Section four. Eventually the paper is culminated in Section five.

2. Related Work

The human action is a spatio-temporal pattern and the representation of feature for action have to grab these patterns. Spatio-Temporal Interest Points (STIP) with Histogram of Gradients (HOG) is used to capture the spatio-temporal pattern³. Since depth map doesn't contain texture, this approach fails to classify the actions. Klaser *et al.*⁴ explained a spatio-temporal descriptor using 3D-gradients. Actionlet ensemble model is built to depict actions without intra class variance⁵. In literature⁶, two features are used: First, a 3D joint feature based on skeleton points captured by depth camera. Second, Local Occupancy Pattern (LOP). After extracting the two features from each frame, temporal dynamics is depicted by fourier temporal pyramid. In literature⁷, informative points from the skeleton joints are extracted and used them to form a descriptor for action recognition. Tang *et al.*⁸ devised a method to describe a shape of the object efficiently by using normal vectors rather than the gradients. Random occupancy patterns are used for robust action recognition by Wang *et al.*⁹. Oreifej and Liu¹⁰ have introduced a descriptor based on Histogram of 4D normal, HON4D, for perceiving the activities using depth sequences.

An hardware implementation for Histogram of Oriented 4D Normals (HON4D) feature extraction is given by Hsu *et al.*¹¹ in real time action recognition context. The hardware features less computation and high speed feature extraction with the adoption of sliding histogram concept. The noticeable part of Sliding histogram is that it permits looped classification without video segmentation in advance. In this regard, a bag of 3D points (like a bag of words) from the silhouette of the depth image is recommended by Li *et al.*¹². Keceli and Can¹³ recommended a multimodal approach for action recognition. In this approach, skeleton model extracts angle and displacement features and depth data infers HOG features. To interpret the semantics of actions, the model is trained using Random Forest algorithm. Latterly, Devanne *et al.*¹⁴ recommended a 3-D joint-based framework to seize both the dynamics and shape of the human body simultaneously.

3. Proposed Method

The training and testing frameworks for the proposed work is depicted in Fig. 1. For the video sequence acquired by depth sensor, all 3-D projections (xy, yz and zx) are calculated for each depth frame^{15, 16}. Unlike literature^{15, 16}, for each projection view, the difference between each alternative frames have been considered to form the Depth Motion Map feature (DMM-f). Assume the depth video sequence contains N frames, the calculation of the DMM-f is as follows. frameⁱ – frameⁱ is the depth map of motion energy.

$$DMM-f_{xy} = \sum_{i=8}^{N-5} (|frame_{xy}^{i} - frame_{xy}^{i-2}|)$$
 (1)

Download English Version:

https://daneshyari.com/en/article/4962225

Download Persian Version:

https://daneshyari.com/article/4962225

<u>Daneshyari.com</u>