

# Discriminative Feature Learning with Constraints of Category and Temporal for Action Recognition

Zhize Wu<sup>(✉)</sup>, Shouhong Wan, Peiquan Jin, and Lihua Yue

Key Laboratory of Electromagnetic Space Information,  
School of Computer Science and Technology, Chinese Academy of Sciences,  
University of Science and Technology of China, Hefei, Anhui, China  
wuzhize@mail.ustc.edu.cn, {wansh, jpq, llyue}@ustc.edu.cn

**Abstract.** Recently, with the availability of the depth cameras, a lot of studies of human action recognition have been conducted on the depth sequences. Motivated by the observations that each pose has its relative location during a complete action sequence, and similar actions have the fine spatio-temporal differences. We propose a novel method to recognize human actions based on the depth information in this paper. Representations of depth maps are learned and reconstructed using a stacked denoising autoencoder. By adding the category and temporal constraints, the learned features are more discriminative, able to capture the subtle but significant differences between actions, and mitigate the nuisance variability of temporal misalignment. Greedy layer-wise training strategy is used to train the deep neural network. Then we employ temporal pyramid matching on the feature representation to generate temporal representation. Finally a linear SVM is trained to classify each sequence into actions. We compare our proposal on MSR Action3D dataset with the previous methods, and the results shown that the proposed method significantly outperforms traditional model, and comparable to, state-of-art action recognition performance. Experimental results also indicate the great power of our model to restore highly noisy input data.

**Keywords:** Action recognition · Category · Temporal · Feature learning · Stacked denoising autoencoders

## 1 Introduction

Human action recognition has been an active field of research in computer vision. The goal of action recognition is to recognize people's behavior from videos in a given scenario automatically. It has many potential applications including content-based video search, human computer interaction, video surveillance, sports video [22, 28]. Most of these applications require high level understanding of spatial and temporal information from videos that are usually composed of multiple simple actions of persons.

Inferring high-level knowledge from a color video especially in a complex and unconstrained scene is very difficult and costly. However, the recent availability of depth cameras such as Kinect [18] has tremendously improved the abilities to understand human activities. Depth maps have several advantages over traditional intensity sensors. First, depth sensors can obtain the holistic 3D structure of the human body, which is invariant to color and texture. Second, color and texture methods perform worse in the dim lighter and the shadows may bring ambiguity. But depth images are insensitive to changes in lighting conditions. Third, depth sensors greatly simplify the process of foreground extraction, removing plenty of noise and disturbance in the background [12, 14].

Furthermore, the 3D skeleton joint positions can be estimated from the depth map accurately following the work of Shotton *et al.* [18]. The extracted skeleton joints have strong representation power, which is more discerning and compact than depth or color sequences. Although with these benefits, two significant aspects arise when one employ joint features for depth-based action recognition. First, existing skeleton joints are not complete, some of the estimated joints are not reliable when the human body is partly in view, moreover, the overlapping of human limbs in some interactive actions can lead to the missing of some joints as well. Second, the action can be performed at difference paces and thus spanning different time durations, which largely influence the performance of action recognition, but it is very difficult to specify effective temporal alignment on action sequences.

To address these two challenging problems, we focus on learning feature, which is robust to the incomplete skeleton, and is mitigatory to the nuisance variability of temporal misalignment. Inspired by the satisfactory performance of previous work on exploring relative 3D joint features [4, 7, 22], we propose a novel method to learn discriminative features from joint 3D features to recognize human actions. We build a deep neural network and employ denoising autoencoders, which has proved their strong abilities to reconstruct and denoise data, as the basic unit of our architecture. This work is also motivated by the observations that each pose has its relative location during a complete action sequence, similar actions have the subtle spatio-temporal differences, we simply add the category and temporal constraints on denoising autoencoders to fuse time-series, intra- and inter-class information into features. We stack the denoising autoencoders with category constraint and greedy layer-wise training strategy is used to train the model. Then we use temporal pyramid matching on the feature representation to generate temporal representation. Finally a linear SVM is trained to classify each sequence into actions. Experiments show that this algorithm achieves superior results.

The contributions of this paper are manifold. First, a new discriminative feature learning algorithm is proposed to recognize depth-based videos. Second, a novel category and temporal constraints are added into denoising autoencoders to preserve temporal, intra-and inter class information. Third, the experiments show that our model has a strong capacity to reconstruct and denoise corrupted data. The remainder of this paper is organized as follows. Section 2 reviews the related work. Section 3 describes the entire flow of our methodology to recognize actions. Section 4 discusses the experimental results. Section 5 concludes the paper.

## 2 Related Work

Since human motion can be considered as a continuous evolution of the spatial configuration of the segments or body posture. Therefore, effective representation of the body configuration and its dynamics over time has been the central to the research of human activity recognition.

Recently, low-level hand-crafted features have been designed to recognize human actions. Spatio-temporal salient points like STIP [7] or some local features, like Cuboids [4], HOG/HOF [8] and SIFT [24] have been widely used. However, directly employ these original methods for color sequences on depth data is infeasible, mainly because of contamination of undefined depth points. In [14] Omar and Zicheng conducted an experiment using MSR-Daily Activity Dataset [22] and found that 60% of the detected Dollar [4] interest points were fired on locations irrelevant to the action of interest and the corresponding classification accuracy is very low (52%). Therefore, recent methods for action recognition in depth sequences explore alternative features particularly for depth-based videos. Li *et al.* [11] projected the depth map into three orthogonal planes and sampled representative 3D points to obtain a bag of 3D points. An action graph was deployed to model the dynamics of the salient postures. Lu *et al.* [25] extracted spatio-temporal interest points from depth videos and built a cuboid similarity feature. Similarly, in [14], Omar and Zicheng quantized the 4D space and represented the possible directions for the 4D normal in order to build a histogram in the 4D space. Due to the temporal misalignment, Su *et al.* [19] introduced a metric to analyze the rate-invariant of trajectories on Riemannian Manifolds with application in visual speech recognition, Jiang Wang [23] proposed a learning-based temporal alignment method, called maximum margin temporal warping (MMTW), to align two action sequences and measure their matching score.

As mentioned before, skeletal information has strong representation power. Lu *et al.* [26] computed histograms of 3D joint locations, reprojected the extracted features using LDA [17], and clustered them into visual words. The temporal evolutions of these words were modeled by HMMs [15]. Jiang *et al.* [22] combined skeleton and depth information to obtain Local Occupancy Patterns (LOP) at each joint and built a Fourier Temporal Pyramid, an actionlet ensemble was learn to represent the actions. Jiajia [12] proposed a dictionary learning algorithm adding the group sparsity and geometry constrains, obtain an over complete set of the input skeletal features. The Temporal Pyramid Matching was used for keeping the temporal information.

In a view of research of unsupervised feature learning [1, 13, 20], which is a set of algorithms that attempt to learn a hierarchy of features by building high-level features from low-level ones. Some models such as CNN [10], DBN [5] and Autoencoders [6] have been shown to yield excellent results in several tasks, e.g. object recognition, natural language processing, and audio classification. One reason for the success of deep learning methods is that they usually learn to capture the posterior distribution of the underlying explanatory factors for the data [2]. However, their extension to the depth maps case is still an open

issue. Therefore, rather than elaborately designing the hand-crafted features as in [14, 16], we choose to learn high level features from data, during the process of learning, we simply add the category and temporal constraints, in order to capture the small but significant differences between actions, and mitigate the nuisance variability of temporal misalignment. The experimental results further prove the feasibility and validity of this feature learning architecture.

### 3 Proposed Method

In this section, we will first describe the basic Denoising Autoencoders. Next, we will extend the model by adding the category and temporal constraints, to make the learned features more discriminative and obtain better accuracies for recognizing actions. Then we introduce the stacking techniques to build a deep architecture. Finally, we employ TPM (temporal pyramid matching) to generate the representation and do classification.

#### 3.1 Denoising Autoencoders

Autoencoders were proposed by Hinton [6] to recognize handwritten digits, which achieved the state of the art at that time. An autoencoder is a special kind of neural networks whose target values are equal to the input ones. A single-layer Autoencoder comprises two parts: **encoder** and **decoder**.

Encoder: The transformation function maps an input vector  $x$  into a hidden layer feature vector  $h$ . Its typical form is a non-linearity function. For each example  $x^{(i)}$  from a data set  $\{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ , we define:

$$f_{\theta}(x^{(i)}) = s(Wx^{(i)} + b) \quad (1)$$

Decoder : The parameterized function maps the hidden layer feature vector  $h$  back to the input space, producing a reconstruction vector:

$$g_{\theta}(h^{(i)}) = s(W'h^{(i)} + c) \quad (2)$$

The set of parameters of this model is  $\theta = \{W, W', b, c\}$ , where  $W$  and  $W'$  are the encoder and decoder weight matrices and  $b$  and  $c$  are the encoder and decoder bias vectors. It is worth mentioning the input vector  $x^{(i)}$  and the reconstruction vector  $r^{(i)}$  have the same dimension  $d_x$ , the hidden layer  $h^{(i)}$  has the dimension  $d_h$ , thus the size of  $W$  is the same as the size of transpose of  $W'$ , which is  $d_h * d_x$ .

The basic autoencoders aim to minimize the reconstruction error of all samples:

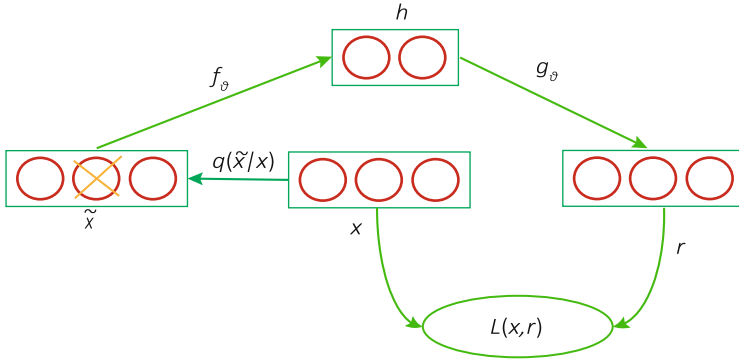
$$L_{AE}(\theta) = \sum_i L(x^{(i)}, g_{\theta}(f_{\theta}(x^{(i)}))) \quad (3)$$

In practice, the choice of function  $s$  is usually a sigmoid function  $s(x) = \frac{1}{1+e^{-x}}$  or a tanh function  $s(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$  and the loss function  $L$  is usually a square loss function  $L(x, r) = \|x - r\|^2$ .

Vincent [20] proposed Stacked Denoising Autoencoders (SDA), exploring a strategy to denoise corrupted version of input data. The input  $x$  is first corrupted into  $\tilde{x}$  using stochastic mapping  $\tilde{x} \sim q(\tilde{x}|x)$ . This is like randomly selecting some nodes of the input and blinding them, that is, every node in the input layer has a possibility  $q$  to be switched to zero. The stochastic corrupted data is regarded as the input of next layer, see Fig. 1. This yields the following objective function:

$$L_{DAE}(\theta) = \sum_i \mathbb{E}_{q(\tilde{x}|x)} \left[ L(x^{(i)}, g_\theta(f_\theta(x^{(i)}))) \right] \quad (4)$$

where  $\mathbb{E}_{q(\tilde{x}|x)}[\cdot]$  is the expectation over corrupted examples  $\tilde{x}$  drawn from the corruption process  $q(\tilde{x}|x)$ .



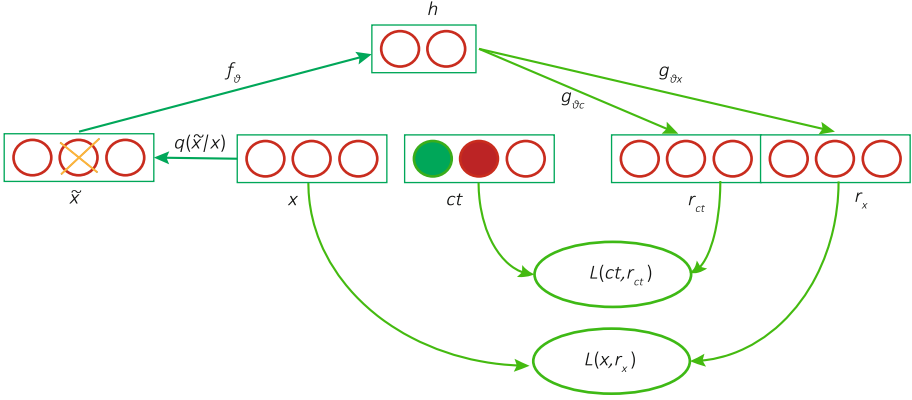
**Fig. 1.** The architecture of the denoising autoencoder. The input data  $x$  is stochastic corrupted into  $\tilde{x}$  by mapping function. The autoencoder then maps  $\tilde{x}$  to  $h$  and maps back  $h$  to  $r$ , the reconstruction result.  $L(x, r)$  is the reconstruction error measurement function.

The reason why DAE can denoise corrupted data is that the training data usually concentrate near a lower-dimensional manifold, yet most of the time the corruption vector is orthogonal to the manifold. The model learns to project the corrupted data back onto the manifold, thus denoising the input.

### 3.2 Adding the Category and Temporal Constraints

Though the features learned by the denoised autoencoders can be highly expressive, as we use the frame-level joint features as the input, all the temporal and category information are discarded. Merely using the model mentioned above, the unsupervised learned features cannot distinguish the significant small differences between similar actions. We modify the denoising autoencoders, adding the category and temporal constraints, to make the model capable of emphasizing the imparities in different actions.

Figure 2 demonstrates our modified autoencoder. Based on the structure of denoising autoencoders, we add an extra target  $ct$  to the network where  $ct$  is a



**Fig. 2.** The architecture of the denoising autoencoder after adding category and temporal constraints. The green solid circular of  $ct$  indicates the temporal information, remaining is a standard unit vector, indicating the category of the video where the frame belongs. The hidden layer  $h$  attempts to reconstruct  $x$  and  $ct$  together, producing the reconstruction vector  $r_x$  and  $r_{ct}$ . The objective error function is  $L(ct, r_{ct}) + L(x, r_x)$ .

vector whose length equals to the action class number and 1,  $(d_c + 1)$ . The first element of vector  $ct$  is the current frame’s relative temporal location in a action sequence, it is simply assigned with the proportion of the frame number and the length of sequence. And the rest of  $ct$  has only one nonzero element whose index indicates the action type of the video where the example frame belongs. In consequence, a category vector  $r_{ct}$  has to be reconstructed by the hidden layer  $h$  using a new mapping function  $g_{\theta ct}$ . Similarly,  $r_x$  is the reconstruction vector of  $x$  by the mapping function  $g_{\theta x}$ . The new training objective of the denoised autoencoder with category and temporal constraints (DAE\_CCT) is:

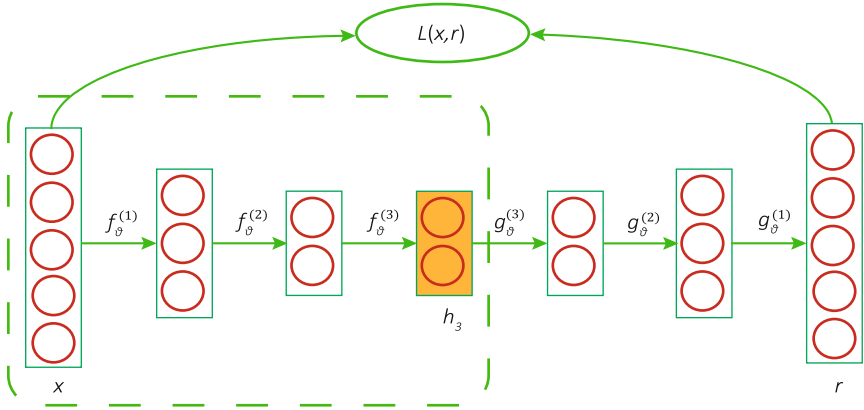
$$L_{DAE\_CCT}(\theta) = \sum_i \mathbb{E}_{q(\tilde{x}|x)} \left[ L(x^{(i)}, g_{\theta x}(f_\theta(x^{(i)}))) + \lambda L(ct^{(i)}, g_{\theta ct}(f_\theta(x^{(i)}))) \right] \quad (5)$$

where  $\lambda$  is a hyper-parameter controlling the strength of the category and temporal regularization. It can be optimized by stochastic gradient descent, analogous to the process of optimizing traditional autoencoders.

The reason why we use a regularization term rather than directly learn the class labels as targets is that the input is the joint vector for one frame, yet the class labels are for the whole video. Apparently there are some similar postures among actions. For example, the *stand and put the hands down* posture appears at the beginning of almost all actions. Training the same posture for different labels will lead to trivial results. The regularization term establishes a trade-off between preserving category and temporal information and reconstructing the input data.

### 3.3 Stacked Architecture

By stacking several layers of denoising autoencoders with the category constraint, we build a deep neural network with great expressive power. Greedy layer-wise training is employed: we first train the first layer to get the encoding function  $f_1$ , then apply it on the clean input to compute the output value, which is used to train the second layer autoencoder to learn  $f_2$ . The process is iteratively conducted. At last we fine-tune the deep neural network as in Fig. 3. We use the output of the last autoencoder as the output for the stacked architecture.



**Fig. 3.** Fine tuning of the stacking architecture. Each layer autoencoder is trained successively to obtain the encoding and decoding functions, which are used to initialize the parameters of the stacking architecture. All parameters are fine tuning to minimize the reconstruction error  $L(x, r)$ , by performing gradient descent. The structure inside the dotted box is the model to extract features and the deepest hidden layer  $h_3$  is the final representation we seek.

### 3.4 Feature Representation and Classification

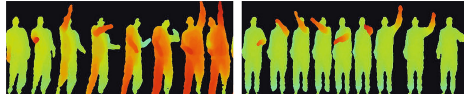
To represent dynamic feature of a whole action sequence, a temporal pyramid matching (TPM) [12] is employed. Motivated by Spatial Pyramid Matching (SPM) [9], a max pooling function is used to generate the multiscale structure. We recursively partition the video sequence into increasingly finer segments along the temporal direction and use max pooling to generate histograms from each sub-region. Typically, 4 levels with each containing 1, 2, 4 and 8 segments are used. The final feature is the concatenation of histograms from all segments. After the final representation for each video is obtained, a multi-class linear SVM [3] is used to speed up the training and testing, results will be discussed in the next section.

## 4 Experimental Results

We evaluate our proposal on a depth-based action recognition dataset, MSR Action3D dataset [11]. We compare our algorithm with the previous methods, and the experimental results shown that the proposed method significantly outperforms traditional model, and attain better than, or comparable to, state-of-art action recognition performance. We also reveal the strong denoising capability of our method to reconstruct noisy 3D joint sequences. In all experiments, we train a deep architecture stacked by two autoencoders, where the first one contains 200 nodes in the hidden layer and the second one contains 400 nodes in the hidden layer. We penalize the average output  $\bar{h}_j$  of the second autoencoder and pushing it to 0.1, in order to add some sparsity to the model and learn an over-completed representation of joint features. The parameter  $\lambda$  is experimentally assigned to 1.3.

### 4.1 MSR Action3D Dataset

MSR Action3D dataset [11] is an action dataset of depth sequences captured by a depth camera. The dataset contains 20 actions: *high arm wave*, *horizontal arm wave*, *hammer*, *hand catch*, *forward punch*, *high throw*, *draw x*, *draw tick*, *draw circle*, *hand clap*, *two hand wave*, *sideboxing*, *bend*, *forward kick*, *side kick*, *jogging*, *tennis swing*, *tennis serve*, *golf swing*, *pick up & throw*. Each action is performed by 10 subjects, each subject performs each action 2 or 3 times. There are 567 depth map sequences in total. Some examples of the depth map sequences are shown in Fig. 4. The provided skeleton data is used to train and test our model.



**Fig. 4.** Sample frames of the MSR-Action3D dataset

We use the same experimental setting as in [22], half of the subjects are used for training and the rest half for testing. We compare with several recent methods and summarize the results in Table. 1. The recognition accuracy is computed by running the experiments 10 times and taking the average of each experiments accuracy. We obtain a recognition accuracy of 88.6%, which is slightly lower than the state-of-the-art result (88.9% [14]) by 0.3%.

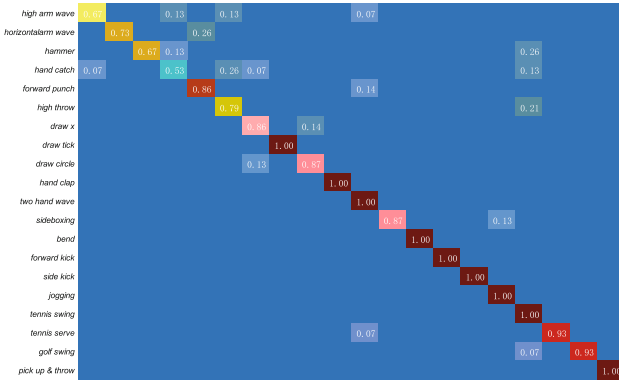
The confusion matrix is shown in Fig. 5. The proposed method performs very well on most of the actions. Some actions, such as catch and wave, are too similar to each other for the proposed method to capture the difference.

We also compare the recognition accuracy for each action of our stacked denoising autoencoders with and without the category and temporal constraints, correspondingly named DAE\_CCT, DAE, in Fig. 6, and superiority is apparent for the most of the actions. The recognition accuracy rate is significantly improved 7% from 81.6% after adding the category and temporal constraints.

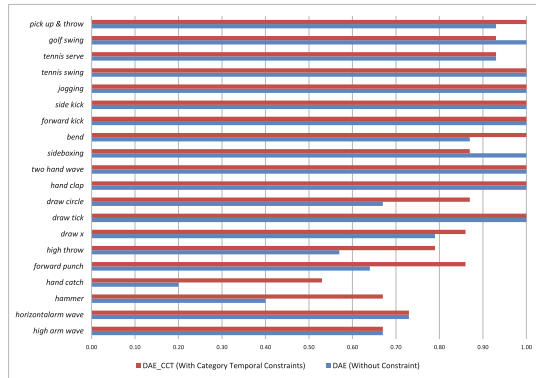


**Table 1.** Comparison of recognition rate on MSR Action3D Dataset

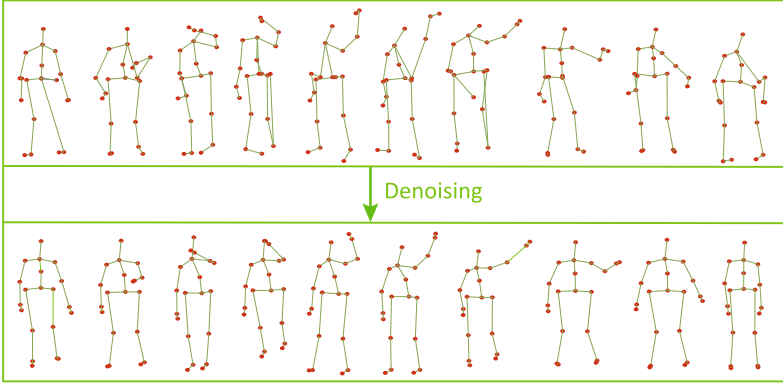
Method	Accuracy
Recurrent Neural Network [13]	0.425
STIP [7] + BOW	0.696
Action Graph on Bag of 3D Points [11]	0.747
Eigenjoints [27]	0.823
Random Occupy Pattern [21]	0.865
Actionlet Ensemble [22]	0.882
HON4D [14] + $D_{desc}$	0.889
Hidden Markov Model	0.63
Dynamic Temporal Warping	0.54
Proposed Method (without constraints)	<b>0.820</b>
Proposed Method (DAE_CCT)	<b>0.886</b>



**Fig. 5.** Confusion matrix of the **DAE\_CCT** on MSR Action3D dataset.



**Fig. 6.** Comparison of the recognition accuracy for each action before and after adding the category and temporal constraints, named **DAE**, **DAE\_CCT** correspondingly.



**Fig. 7.** Examples showing the capability of our model to denoise corrupted data. Top: the corrupted input 3D joint sequence *high arm wave* from MSR Action3D dataset. Bottom: the reconstructed 3D joint sequence

## 4.2 Capability to Denoise Corrupted Data

The proposed model also has the strong capability to reconstruct realistic data from corrupted input. The top row of Fig. 7 is an action sequence *high arm wave* selected from MSR Action3D dataset. In order to better demonstrate our algorithm efficiency, we add some Gaussian noise to the joint positions and leave out joints randomly. The bottom row is the reconstruction action sequence, where we can observe that the missing joints are all restored via our model and the motions are more natural and fluent than before.

## 5 Conclusion

This paper presented a novel feature learning methodology for human action recognition with depth cameras. To better represent the 3D joint features, a deep stacked denoising autoencoder that incorporated with the category and temporal constraints were proposed. The proposed model is capable of capturing subtle spatio-temporal details between actions, robust to the noises and errors in the joints positions. The experiments demonstrated the effectiveness and robustness of the proposed approach. In the future, we aim to integrate the more precise temporal information into our feature learning architecture.

**Acknowledgments.** This work was supported by the National Natural Science Foundation of China (Grant No. 61272317) and the General Program of Natural Science Foundation of Anhui of China (Grant No. 1208085MF90).

## References

1. Bengio, Y.: Learning deep architectures for AI. *Found. Trends<sup>®</sup> Mach. Learn.* **2**(1), 1–127 (2009)
2. Bengio, Y., Courville, A., Vincent, P.: Representation learning: a review and new perspectives. *IEEE Trans. Pattern Anal. Mach. Intell.* **35**(8), 1798–1828 (2013)
3. Chang, C.-C., Lin, C.-J.: Libsvm: a library for support vector machines. *ACM Trans. Intell. Syst. Technol. (TIST)* **2**(3), 27 (2011)
4. Dollár, P., Rabaud, V., Cottrell, G., Belongie, S.: Behavior recognition via sparse spatio-temporal features. In: *2nd Joint IEEE International Workshop on Visual Surveillance and Performance Evaluation of Tracking and Surveillance, 2005*, pp. 65–72. IEEE (2005)
5. Hinton, G., Osindero, S., Teh, Y.-W.: A fast learning algorithm for deep belief nets. *Neural Comput.* **18**(7), 1527–1554 (2006)
6. Hinton, G.E., Salakhutdinov, R.R.: Reducing the dimensionality of data with neural networks. *Science* **313**(5786), 504–507 (2006)
7. Laptev, I.: On space-time interest points. *Int. J. Comput. Vis.* **64**(2–3), 107–123 (2005)
8. Laptev, I., Marszalek, M., Schmid, C., Rozenfeld, B.: Learning realistic human actions from movies. In: *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2008*, pp. 1–8. IEEE (2008)
9. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: spatial pyramid matching for recognizing natural scene categories. In: *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2, pp. 2169–2178. IEEE (2006)
10. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proc. IEEE* **86**(11), 2278–2324 (1998)
11. Li, W., Zhang, Z., Liu, Z.: Action recognition based on a bag of 3d points. In: *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 9–14. IEEE (2010)
12. Luo, J., Wang, W., Qi, H.: Group sparsity and geometry constrained dictionary learning for action recognition from depth maps. In: *2013 IEEE International Conference on Computer Vision (ICCV)*, pp. 1809–1816. IEEE (2013)
13. Martens, J., Sutskever, I.: Learning recurrent neural networks with hessian-free optimization. In: *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 1033–1040 (2011)
14. Oreifej, O., Liu, Z.: Hon4d: histogram of oriented 4d normals for activity recognition from depth sequences. In: *2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 716–723. IEEE (2013)
15. Rabiner, L.: A tutorial on hidden markov models and selected applications in speech recognition. *Proc. IEEE* **77**(2), 257–286 (1989)
16. Sang, R., Jin, P., Wan, S.: Discriminative feature learning for action recognition using a stacked denoising autoencoder. In: Pan, J.-S., Snasel, V., Corchado, E.S., Abraham, A., Wang, S.-L. (eds.) *Intelligent Data Analysis and Its Applications, Volume I. AISC*, vol. 297, pp. 521–531. Springer, Heidelberg (2014)
17. Scholkopf, B., Mullert, K.-R.: Fisher discriminant analysis with kernels. In: *Proceedings of the 1999 IEEE Signal Processing Society Workshop Neural Networks for Signal Processing IX*, Madison, WI, USA, pp. 23–25 (1999)
18. Shotton, J., Sharp, T., Kipman, A., Fitzgibbon, A., Finocchio, M., Blake, A., Cook, M., Moore, R.: Real-time human pose recognition in parts from single depth images. *Commun. ACM* **56**(1), 116–124 (2013)

19. Su, J., Srivastava, A., de Souza, F.D.M., Sarkar, S.: Rate-invariant analysis of trajectories on riemannian manifolds with application in visual speech recognition. In: 2014 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 620–627. IEEE (2014)
20. Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., Manzagol, P.-A.: Stacked denoising autoencoders: learning useful representations in a deep network with a local denoising criterion. *J. Mach. Learn. Res.* **11**, 3371–3408 (2010)
21. Wang, J., Liu, Z., Chorowski, J., Chen, Z., Wu, Y.: Robust 3D action recognition with random occupancy patterns. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) ECCV 2012, Part II. LNCS, vol. 7573, pp. 872–885. Springer, Heidelberg (2012)
22. Wang, J., Liu, Z., Wu, Y., Yuan, J.: Mining actionlet ensemble for action recognition with depth cameras. In: 2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1290–1297. IEEE (2012)
23. Wang, J., Wu, Y.: Learning maximum margin temporal warping for action recognition. In: 2013 IEEE International Conference on Computer Vision (ICCV), pp. 2688–2695. IEEE (2013)
24. Wu, Z., Wan, S., Yue, L., Sang, R.: Discriminative image representation for classification. In: Pan, J.-S., Snasel, V., Corchado, E.S., Abraham, A., Wang, S.-L. (eds.) *Intelligent Data Analysis and Its Applications, Volume II. AISC*, vol. 298, pp. 331–341. Springer, Heidelberg (2014)
25. Xia, L., Aggarwal, J.K.: Spatio-temporal depth cuboid similarity feature for activity recognition using depth camera. In: 2013 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pp. 2834–2841. IEEE (2013)
26. Xia, L., Chen, C.-C., Aggarwal, J.K.: View invariant human action recognition using histograms of 3d joints. In: 2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 20–27. IEEE (2012)
27. Yang, X., Tian, Y.: Eigenjoints-based action recognition using naive-bayes-nearest-neighbor. In: 2012 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 14–19. IEEE (2012)
28. Yang, X., Zhang, C., Tian, Y.: Recognizing actions using depth motion maps-based histograms of oriented gradients. In: Proceedings of the 20th ACM international conference on Multimedia, pp. 1057–1060. ACM (2012)