

Recognizing Human Activities from Raw Accelerometer Data Using Deep Neural Networks

Licheng Zhang, Xihong Wu and Dingsheng Luo*

Key Lab of Machine Perception (Ministry of Education), Speech and Hearing Research Center
Department of Machine Intelligence, School of Electronics Engineering and Computer Science

Peking University, Beijing, 100871, China

Emails: {zhanglc, wxh, dsluo}@cis.pku.edu.cn

* Corresponding author, IEEE Member

Abstract—Activity recognition from wearable sensor data has been researched for many years. Previous works usually extracted features manually, which were hand-designed by the researchers, and then were fed into the classifiers as the inputs. Due to the blindness of manually extracted features, it was hard to choose suitable features for the specific classification task. Besides, this heuristic method for feature extraction could not generalize across different application domains, because different application domains needed to extract different features for classification. There was also work that used auto-encoders to learn features automatically and then fed the features into the K-nearest neighbor classifier. However, these features were learned in an unsupervised manner without using the information of the labels, thus might not be related to the specific classification task. In this paper, we recommend deep neural networks (DNNs) for activity recognition, which can automatically learn suitable features. DNNs overcome the blindness of hand-designed features and make use of the precious label information to improve activity recognition performance. We did experiments on three publicly available datasets for activity recognition and compared deep neural networks with traditional methods, including those that extracted features manually and auto-encoders followed by a K-nearest neighbor classifier. The results showed that deep neural networks could generalize across different application domains and got higher accuracy than traditional methods.

Keywords- activity recognition; deep neural networks; feature learning; accelerometer data

I. INTRODUCTION

Activity recognition from wearable sensor data has been studied for many years by lots of researchers around the world. Previous works usually needed to extract features manually, which were hand-designed by the researchers, and then were fed into the classifiers as the inputs. For instance, Bao and Intille [1] collected a large amount of accelerometer data, and then manually extracted four features: mean, energy, frequency-domain entropy, and correlation, and then fed these features into four classifiers: decision table, K-nearest neighbor (KNN), decision tree, and Naïve Bayes. Ravi et al. [2] also collected accelerometer data and manually extracted four features: mean, standard deviation, energy, and correlation of acceleration data. Then these features were fed into eighteen different classifiers as the inputs.

Due to the blindness of hand-designed features, it was hard to choose suitable features for the specific classification task. Besides, from the above examples, we can find that in different application domains, the researchers needed to extract different features for classification. That is to say, this heuristic method for feature extraction could not generalize across different application domains.

Plötz et al. [3] recommended auto-encoders for automatic feature learning. The learned features were then fed into the K-nearest neighbor (KNN) classifier. Their method could generalize across different application domains. However, the features were learned in an unsupervised manner without making use of the precious label information, which could improve activity recognition performance.

In this paper, we recommend deep neural networks (DNNs) for activity recognition, which were widely used in speech recognition [4], computer vision [5] and natural language processing [6]. DNNs can automatically learn suitable features without relying on the experience of the researchers and then perform classification. Besides, the DNN model has a process of fine-tuning using the precious information of the labels, which can adjust the features to make them more related to the specific classification task and improve classification performance. We compared the DNN model with traditional methods, including eight methods that manually extracted features according to the experience of the researchers, and auto-encoders followed by a KNN classifier. The results of the experiments on three publicly available datasets show that deep neural networks can generalize across different application domains and get higher accuracy than traditional methods.

The reminder of this paper is organized as follows: Section II describes the related work. Section III details how to train the DNN model. Section IV describes three publicly available datasets and presents the details of the experimental design, and gives the recognition results of our method and traditional methods. Section V concludes this paper and gives our future work.

II. RELATED WORK

Activity recognition is an important research problem, drawing wide attention around the world. Bao and Intille [1]

collected accelerometer data of 20 subjects wearing five biaxial accelerometers positioned on different parts of the body. Then they extracted four features: mean, energy, frequency-domain entropy, and correlation. And then they fed these features into four classifiers: decision table, K-nearest neighbor, decision tree, and Naïve Bayes. Twenty activities were recognized. The results showed that decision tree performed best. Ravi et al. [2] also collected accelerometer data of 2 subjects wearing a tri-axial accelerometer near the pelvic region. Eight activities were performed. Four features were extracted: mean, standard deviation, energy, and correlation of acceleration data. Then these features were fed into eighteen different classifiers for activity recognition in four different settings. The results showed that Plurality Voting performed best in three settings and Boosted support vector machine performed best in the remaining setting. Parkka et al. [7] collected context data using different sensors, including accelerometers, a Global Positioning System recorder, etc. They extracted six features for classification. Among them, three features were peak frequency, median and peak power of the recorded up-down chest acceleration. The other three were the variance of the recorded back-forth chest acceleration, sum of variances of three-dimensional wrist accelerations and power ratio of 1-1.5Hz and 0.2-5Hz frequency bands, which were measured using the left-right magnetometer on chest. Seven activities were recognized. They compared three classifiers: automatic decision tree, custom decision tree and artificial neural network. The results showed that automatic decision tree outperformed custom decision tree and artificial neural network.

Phone-based activity recognition has also been researched for many years. Kwapisz et al. [8] collected the accelerometer data of 29 users using android smartphones worn in the front pant pocket when the subjects were performing six different activities. They extracted six features: the average acceleration, standard deviation, average absolute difference, average resultant acceleration, time between peaks, and the binned distribution, and chose three frequently-used classifiers: decision tree, logistic regression and three-layer neural network. The classification results showed that three-layer neural network performed best. They indicated that their method could not be applied to real-time activity recognition. Kose et al. [9] achieved a real-time activity recognition system on smartphones. Four features were extracted: average, minimum, maximum, and standard deviation. They compared the performance of two classifiers: Naïve Bayes and Clustered KNN. The results showed that Clustered KNN displayed a much better performance than the Naïve Bayes method in terms of the recognition accuracy on mobile phones with limited resources.

The works described above had two common problems: they needed to extract features manually, which were hand-designed by the researchers. Due to the blindness of manually extracted features, it was hard to choose suitable features for the specific classification task. Besides, from the above works, we can find that different application domains needed to extract different features. Therefore, extracting features manually could not generalize across different application domains.

Plötz et al. [4] made use of deep learning for feature extraction. They used auto-encoders to learn features

automatically. They compared the learned features with manually extracted features. Both of the two types of features used KNN as the classifier. The results showed that automatically learned features outperformed manually extracted features. However, the features were learned in an unsupervised manner without making use of the precious label information, which could improve activity recognition performance.

In this paper, we recommend deep neural networks for activity recognition. They can automatically learn suitable features without relying on the experience of the researchers. In addition, a discriminative fine-tuning on the whole network is performed after the unsupervised feature learning, using the information of the labels by back-propagation, which adjust the features slightly to get the category boundaries right and make the features better for classification.

III. METHOD

Deep neural networks are a kind of the neural network model. The graphical representation of the DNN model is shown in Fig. 1. It has an input layer, two or more hidden layers, and an output layer. The nodes represent the variables and the links between the nodes represent the weight parameters. Arrows denote the information flow direction through the network. Due to the large amount of parameters, DNNs possess the ability of automatically learning suitable features from the raw data. The parameters of a DNN are usually randomly initialized and then the whole network is trained using back-propagation.

The DNNs that we chose for activity recognition were the deep belief networks (DBNs) [10], whose parameters were initialized by a generative pre-training and then were fine-tuned using back-propagation. We detailed how to train deep belief networks as follows.

A. Pre-training

The parameters of DBNs are initialized by a generative pre-training. Different from traditional deep neural networks, DBNs treat every two adjacent layers as a Restricted Boltzmann Machine (RBM). The input layer and the first hidden layer form a Gaussian-binary RBM, whose input units are real-valued data and liner with Gaussian noise and the hidden units are binary. Two consecutive hidden layers form a binary-binary RBM, in which all the units are binary and stochastic. A RBM is a fully connected, bipartite undirected graph, which consists of a visible layer V and a hidden layer H , as is shown in Fig. 2. There are no visible-visible or hidden-hidden connections. For a RBM, the visible units represent the observations and are connected to the hidden units using the undirected weighted connections.

The hidden units act as feature detectors. Both V and H are vectors, representing the visible units and hidden units respectively. Again, the links between the nodes represent the weight parameters, which, however, are undirected. Suppose the variables (V, H) take values (v, h) . Then the RBM is trained as follows: sampling a new state h for the hidden units based on $p(h|v)$, then sampling a new state v for the visible units based on $p(v|h)$, and then repeating the above process [11].

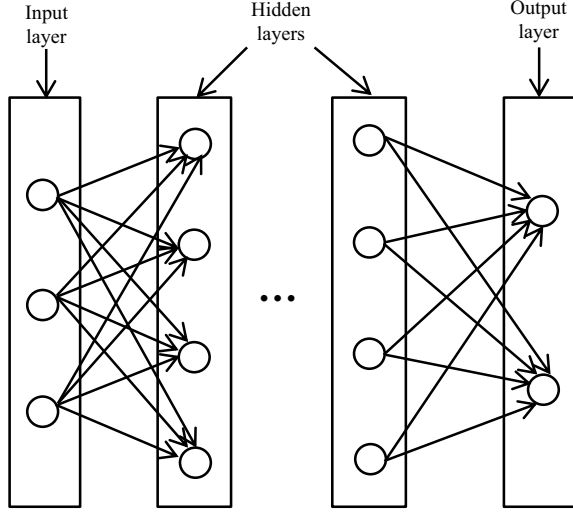


Figure 1. The graphical representation of a deep neural network

Noticing that a RBM is an undirected graph and the weighted connections are undirected, alternate updating v and h is feasible. Once one RBM is trained, we treat the hidden units of this RBM as the visible units of the next RBM and the next layer as the hidden units, and train the next RBM. Once all the RBMs of a deep belief network are trained, the generative pre-training is finished.

B. Fine-tuning

Pre-training learns a generative model, which does not make use of the information of the labels, and the learned features (the hidden units) may be irrelevant for classifying different categories. So after pre-training, the whole network is treated as a feed-forward, deterministic neural network. And then a discriminative fine-tuning is performed on the whole network using the back-propagation algorithm, making use of the precious information in the labels, to slightly adjust the features in each layer, which were discovered by the unsupervised pre-training, and get the category boundaries right for better classification. After the fine-tuning process, the training of a deep neural network is completed. The back-propagation algorithm can be conjugate gradients algorithm or stochastic gradient decent algorithm.

IV. EXPERIMENTS

In the following subsection, we described three datasets, did experiments on these datasets to compare the performance of DNNs with that of traditional methods, and finally gave the results.

A. Datasets

We chose three publicly available datasets to do experiments. The details of the three public datasets are presented as follows.

Opportunity Chavarriaga et al. [12] collected a large amount of sensor data of daily activities in a sensor-rich environment. We chose the accelerometer data recorded by the

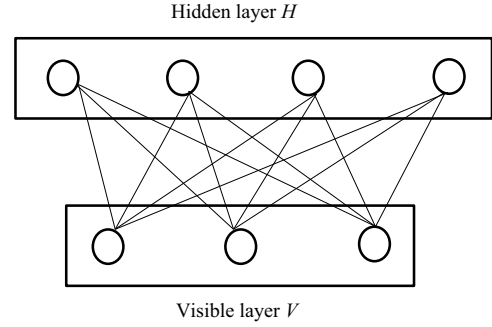


Figure 2. The undirected graph of a RBM

accelerometer attached to the right arm of one subject for experiments. ADL refers to activity of daily living. The data of ADL1, ADL2, ADL3 and Drill run formed the training set and the data of ADL4 and ADL5 formed the testing set. We chose 11 low-level activities, including a null activity. The sampling rate of this accelerometer is 64 Hz. We set the frame length to be 1s. Two consecutive frames overlapped by 50%. Finally, the training set contains 4800 frames and the testing set contains 1800 frames.

USC-HAD Zhang et al. [13] collected a human activity dataset of 12 basic human activities using a high performance inertial sensor device, which was located at the subjects' front right hip. We chose the accelerometer data of one subject for experiments. The sampling rate is 100 Hz. We set the frame length to be 1s. Two consecutive frames overlapped by 50%. We randomly chose samples of each activity from the dataset and used them to compose a testing set. For each activity, its percentage in the testing set is equal to that in the training set. Eventually, the training set contains 4400 frames and the testing set contains 780 frames.

Daily and Sports Activities Altun et al. [14] collected sensor data of 19 activities using inertial sensors and magnetometers positioned on different parts of the body. The sampling rate is 25 Hz. We chose the accelerometer data recorded by the sensor attached to the chest of one subject for experiments. We set the frame length to be 1.2s. Two consecutive frames overlapped by 50%. We randomly chose samples of each activity from the dataset and used them to compose a testing set. In the training set, all activities have equal samples. So does the activities in the testing data. Eventually, the training set contains 8000 frames and the testing set contains 1500 frames.

B. Experiment Design

We compared our method with eight frequently-used methods: Bayesian nets, Naïve Bayes, Support vector machine, Logistic regression, K-nearest neighbor, Decision tree, Random forest, and three-layer feed-forward neural network (3-layer NN). These classifiers were more often used and could be found in [1, 2, 3, 8, 9, 15-17]. 3-layer NN has an input layer, a hidden layer and an output layer. Due to its small amount of parameters, it still needed to extract features manually and took these features as inputs. The collected accelerometer data

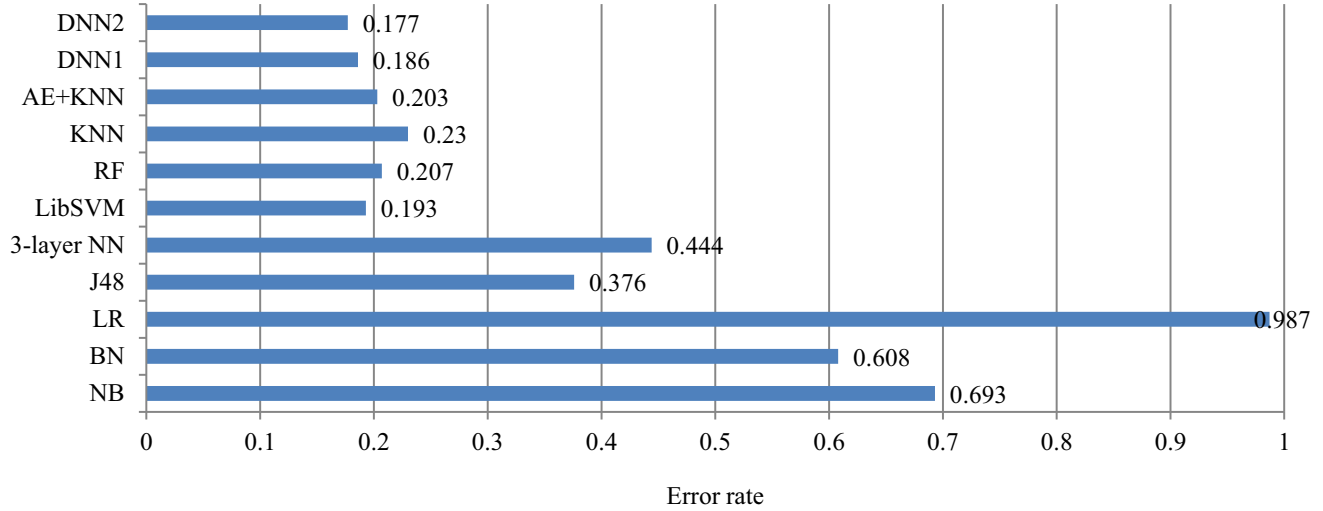


Figure 3. Classification results on the Opportunity dataset

included three dimensions: x , y , and z . We calculated the magnitude of acceleration denoted as mag :

$$mag = \sqrt{x^2 + y^2 + z^2} \quad (1)$$

Then we extracted features for x , y , z , and mag for the eight methods. We manually extracted five frequently-used features: mean, standard deviation, energy, frequency-domain entropy, and correlation of acceleration data, which were used in [1, 2]. So the inputs of the selected eight classifiers had 20 features. The eight classifiers were performed using the WEKA machine learning toolkit [18, 19]. For the DNN model, the inputs were the raw accelerometer data without any preprocessing. We used the deep belief network code from the homepage of Geoffrey Hinton [20]. The training set is divided into mini-batches with each mini-batch having about 100 frames (there is small difference among three datasets). The DNN model was pre-trained using stochastic gradient descent with a mini-batch at a time. An epoch refers to passing over all the mini-batches, i.e., the entire training set. We ran 100 epochs for the Gaussian-binary RBM at learning rate 0.001 and ran 50 epochs for the binary-binary RBMs at learning rate 0.1.

TABLE I. THE NOTATIONS OF DIFFERENT CLASSIFICATION METHODS

Classification Methods	Notations
Bayesian nets	BN
Naïve Bayes	NB
Support vector machine	LibSVM
Logistic regression	LR
K-nearest neighbor	KNN
Decision tree	J48
Random forest	RF
three-layer feed-forward neural network	3-layer NN

During pre-training, we used a weight-cost of 0.0002 and a momentum of 0.9 [21]. For fine-tuning, we used the conjugate gradients algorithm on larger mini-batches at a time, which contained about 1000 frames (there is small difference among three datasets). For the conjugate gradient fine-tuning, we used Carl Rasmussen’s “minimize” code [22], and performed three line searches for each mini-batch in each epoch. We performed 400 epochs of fine-tuning on the entire training set. We also adjusted the number of units of one layer with the number of units of other layers fixed. Finally, the DNN model with a better recognition result had a S-500-500-2000-T structure. S and T refer to the number of input units and output units respectively. That is to say, it had S units in the input layer, 500 units in the first hidden layer, 500 units in the second hidden layer, 2000 units in the third hidden layer, and T units in the output layer, representing T categories. Specifically, for the Opportunity dataset, the number of input units is 256 and the number of output units is 11. And for the USC-HAD dataset, the number of input units is 400 and the number of output units is 12. And For the Daily and Sports Activities dataset, the number of the input units is 120 and the number of the output units is 19.

We also compared our method with auto-encoders followed by a KNN classifier. We used the same structure as [4] did. It had S units in the input layer, 1024 units in the first and second hidden layer, and 30 units in the top layer. We ran 100 epochs for the Gaussian-binary RBM at a learning rate of 0.001 and ran 50 epochs for binary-binary RBMs at a learning rate of 0.1. During pre-training, we used a weight-cost of 0.0002 and a momentum of 0.9. For their method, the 30 features were then fed into the KNN classifier. We set K from 1 to 25 and gave the best accuracy. Higher values of K had no significant impact on the accuracy. For our method, we then fine-tuned the whole network using the conjugate gradients algorithm in the cross-entropy error between the target labels and the output labels. Again, we performed three line searches for each mini-batch in each epoch and we performed 400 epochs of fine-tuning on the entire training set.

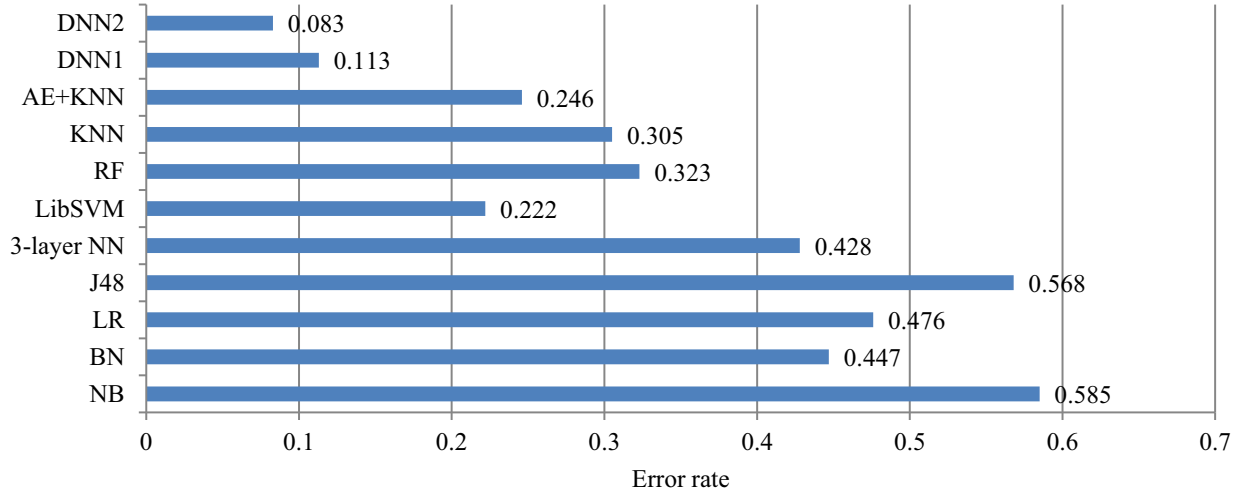


Figure 4. Classification results on the USC-HAD dataset

The notations of different methods are listed in TABLE I. Besides, we use “AE+KNN” to denote the method of feature learning using auto-encoders followed by a KNN classifier. We use “DNN1” to denote deep neural networks having a structure of S-1024-1024-30-T and use “DNN2” to denote deep neural networks that have a structure of S-500-500-2000-T. Again, S and T represent the number of input units and the number of output units respectively.

C. Results

The classification results of different methods on the Opportunity dataset are presented in Fig. 3. And results on the USC-HAD dataset are presented in Fig. 4. And results on the Daily and Sports Activities dataset are presented in Fig. 5. We take accuracy as the evaluation criterion, which is advisable for comparing the performance of different methods. For the KNN

classifier, we set K from 1 to 25 and gave the best accuracy. Higher values of K had no significant impact on the accuracy.

The results showed that compared with the eight methods that extracted features manually, the DNN model (DNN2) performed best on all the three datasets. Besides, the DNN model (DNN1) performed better than auto-encoders followed by a KNN classifier (AE+KNN) on all the three datasets. That is to say, deep neural networks performed better than traditional methods. In addition, for all datasets, deep neural networks work well. Therefore, they can generalize across different application domains.

V. CONCLUSION AND FUTURE WORK

Previous works usually needed to extract features manually, which were hand-designed by the researchers, and then were

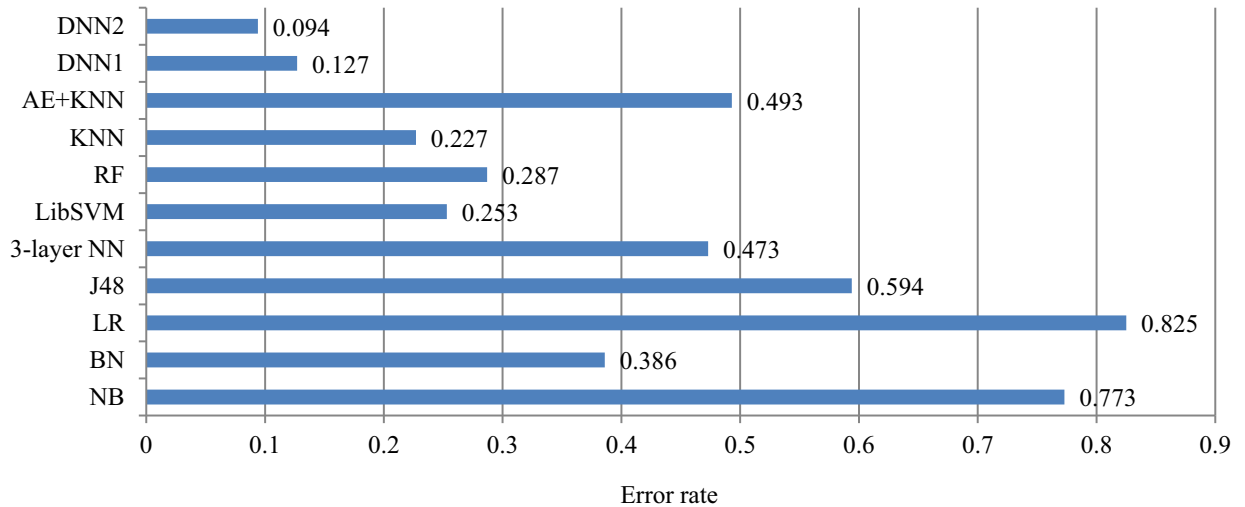


Figure 5. Classification results on the Daily and Sports Activities dataset

fed into the classifiers as the inputs. Due to the blindness of manually extracted features, it was hard to choose suitable features for the specific classification task. Besides, this method of manual feature extraction could not generalize across different application domains. Auto-encoders learned features automatically without relying on the experience of the researchers. However, the features were learned in an unsupervised manner without using the information of the labels, which could improve activity recognition performance. In this paper, we recommend deep neural networks for activity recognition, which can automatically learn suitable features without relying on the experience of the researchers and performed a discriminative fine-tuning after the unsupervised feature learning using the precious information of the labels, which would improve the classification performance. We did experiments on three publicly available datasets and compared our method with traditional methods. The results showed that DNNs could generalize across different application domains and got higher accuracy than traditional methods.

For future work, we will implement the DNN model on smartphones, and test the recognition time of the DNN model on smartphones to determine whether the DNN model could achieve the real-time effect. Besides, we will use other methods, which also automatically learn features, such as recurrent neural network and convolutional neural network, to try to improve activity recognition performance.

ACKNOWLEDGMENT

The work is supported in part by National Basic Research Program (973 Program) of China (No. 2013CB329304), the National Natural Science Foundation of China (No. 90920302, No. 91120001, No.61121002), the "Twelfth Five-Year" National Science & Technology Support Program of China (No. 2012BAI12B01), the Seeding Grant for Med and Info Sciences of Peking University (No.2014-MI-10) and the Key Program of National Social Science Foundation of China (No. 12&ZD119).

REFERENCES

- [1] L. Bao and S. S. Intille, "Activity recognition from user-annotated acceleration data," *Pervasive Computing, Lecture Notes in Computer Science*, vol. 3001, pp. 1-17, 2004.
- [2] N. Ravi, N. Dandekar, P. Mysore, and M. L. Littman, "Activity Recognition from Accelerometer Data," *Proceedings of the Seventeenth Conference on Innovative Applications of Artificial Intelligence*, vol. 5, pp. 1541-1546, 2005.
- [3] T. Plötz, N. Y. Hammerla, and P. Olivier, "Feature learning for activity recognition in ubiquitous computing," *Proceedings of the Twenty-Second international joint conference on Artificial Intelligence*, vol. 22, pp. 1729-1734, 2011.
- [4] G. E. Hinton, L. Deng, D. Yu, G. Dahl, A. Mohamed, N. Jaitly, A. Senior, V. Vanhoucke, P. Nguyen, T. Sainath, and B. Kingsbury, "Deep neural networks for acoustic modeling in speech recognition," *IEEE Signal Processing Magazine*, vol. 29, pp. 82-97, 2012.
- [5] G. E. Hinton and R. R. Salakhutdinov, "Reducing the dimensionality of data with neural networks," *Science*, vol. 313, pp. 504-507, 2006.
- [6] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th international conference on Machine learning*, pp. 160-167, 2008.
- [7] J. Parkka, M. Ermes, P. Korpiainen, J. Mantyjarvi, J. Peltola, and I. Korhonen, "Activity Classification Using Realistic Data From Wearable Sensors," *IEEE Transactions on Information Technology in Biomedicine*, vol. 10, pp. 119-128, 2006.
- [8] J. R. Kwapisz, G. M. Weiss, and S. A. Moore, "Activity recognition using cell phone accelerometers," *ACM SIGKDD Explorations Newsletter*, vol. 12, pp. 74-82, 2011.
- [9] M. Kose, O. D. Incel, and C. Ersoy, "Online human activity recognition on smart phones," in *Workshop on Mobile Sensing: From Smartphones and Wearables to Big Data*, pp. 11-15, 2012.
- [10] G. E. Hinton, S. Osindero, and Y. W. Teh, "A fast learning algorithm for deep belief nets," *Neural computation*, vol. 18, pp. 1527-1554, 2006.
- [11] A. Fischer and C. Igel, "An introduction to restricted Boltzmann machines," in *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, pp. 14-36, 2012.
- [12] R. Chavarriaga, H. Sagha, A. Calatroni, S. T. Digumarti, G. Troster, J. D. R. Millan, and D. Roggen, "The Opportunity challenge: A benchmark database for on-body sensor-based activity recognition," *Pattern Recognition Letters*, vol. 34, pp. 2033-2042, 2013.
- [13] M. Zhang and A. A. Sawchuk, "Usc-had: a daily activity dataset for ubiquitous activity recognition using wearable sensors," in *Proceedings of the 2012 ACM Conference on Ubiquitous Computing*, pp. 1036-1043, 2012.
- [14] K. Altun, B. Barshan, and O. Tunçel, "Comparative study on classifying human activities with miniature inertial and magnetic sensors," *Pattern Recognition*, vol. 43, pp. 3605-3620, 2010.
- [15] W. Wu, S. Dasgupta, E. E. Ramirez, C. Peterson, and G. J. Norman, "Classification accuracies of physical activities using smartphone motion sensors," *Journal of medical Internet research*, vol. 14, no. 5, e130, 2012.
- [16] O. D. Incel, M. Kose, and C. Ersoy, "A review and taxonomy of activity recognition on mobile phones," *BioNanoScience*, vol. 3, pp. 145-171, 2013.
- [17] M. Shoaib, S. Bosch, O. D. Incel, H. Scholten, and P. J. Havinga, "Fusion of smartphone motion sensors for physical activity recognition," *Sensors*, vol. 14, pp. 10146-10176, 2014.
- [18] S. R. Garner, "Weka: The waikato environment for knowledge analysis," in *Proceedings of the New Zealand computer science research students conference*, pp. 57-64, 1995.
- [19] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I. H. Witten, "The WEKA data mining software: an update," *ACM SIGKDD explorations newsletter*, vol. 11, pp. 10-18, 2009.
- [20] The MATLAB code is available at <http://www.cs.Toronto.edu/~hinton/MatlabForSciencePaper.html>.
- [21] G. E. Hinton, "A practical guide to training restricted Boltzmann machines," *Neural Networks: Tricks of the Trade, Lecture Notes in Computer Science*, vol. 7700, pp. 599-619, 2012.
- [22] Carl Rasmussen's "minimize" code is available at <http://www.kyb.tuebingen.mpg.de/bs/people/carl/code/minimize/>.