



Performance evaluation of deep feature learning for RGB-D image/video classification



Ling Shao^{a,b,*}, Ziyun Cai^c, Li Liu^b, Ke Lu^{d,e}

^a College of Electronic and Information Engineering, Nanjing University of Information Science and Technology, Nanjing 210044, China

^b School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, UK

^c Department of Electronic and Electrical Engineering, University of Sheffield, Mappin Street, Sheffield S1 3JD, UK

^d University of Chinese Academy of Sciences, Beijing 100049, China

^e Beijing Center for Mathematics and Information Interdisciplinary Sciences, Beijing, China

ARTICLE INFO

Article history:

Received 24 August 2016

Revised 15 December 2016

Accepted 2 January 2017

Available online 3 January 2017

Keywords:

Deep neural networks

RGB-D data

Feature learning

Performance evaluation

ABSTRACT

Deep Neural Networks for image/video classification have obtained much success in various computer vision applications. Existing deep learning algorithms are widely used on RGB images or video data. Meanwhile, with the development of **low-cost** RGB-D sensors (such as Microsoft Kinect and Xtion Pro-Live), high-quality RGB-D data can be easily acquired and used to enhance computer vision algorithms [14]. It would be interesting to investigate how deep learning can be employed for extracting and fusing features from RGB-D data. In this paper, after briefly reviewing the basic concepts of RGB-D information and four prevalent deep learning models (i.e., Deep Belief Networks (DBNs), Stacked Denoising Auto-Encoders (SDAE), Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) Neural Networks), we conduct extensive experiments on five popular RGB-D datasets including three image datasets and two video datasets. We then present a detailed analysis about the comparison between the learned feature representations from the four deep learning models. In addition, a few suggestions on how to adjust hyper-parameters for learning deep neural networks are made in this paper. According to the extensive experimental results, we believe that this evaluation will provide insights and a deeper understanding of different deep learning algorithms for RGB-D feature extraction and fusion.

© 2017 Elsevier Inc. All rights reserved.

1. Introduction

Learning good feature representations from input data for high-level tasks receives much attention in computer vision, robotics and medical imaging [19,39]. Image/video classification is a classic and challenging high-level task, which has many practical applications, such as robotic vision [1], video surveillance [42] and image retrieval [18]. The objective is to predict the labels of new coming images/videos. Though RGB image/video classification has been studied for many years, it still faces a lot of challenges, such as complicated background, illuminance change and occlusion. With the invention of the **low-cost** Microsoft Kinect sensor, it opens a new dimension (i.e., depth data) to overcome the above challenges. Compared to

* Corresponding author at: College of Electronic and Information Engineering, Nanjing University of Information Science and Technology, Nanjing 210044, China and School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, UK.

E-mail address: ling.shao@ieee.org (L. Shao).

RGB images, depth images are robust to the variations in color, illumination, rotation angle and scale. It has been proved that combining RGB and depth information in image/video classification tasks can significantly improve the classification accuracy [14]. Therefore, an increasing number of RGB-D datasets have been created as benchmarks [10]. Moreover, Deep Neural Networks for high-level tasks obtain great success in recent years. Different from hand-crafted feature representations such as SIFT and HOG, deep learned features are automatically learned from the images or videos. These neural network models improve the state-of-the-art performance on many important datasets (e.g., the ImageNet dataset), and some of them even overcome human performance [43]. Combining the advantages of RGB-D images and Deep Neural Networks, many researchers are making great efforts to design more sophisticated algorithms. However, no single existing approach can successfully handle all scenarios. Therefore, it is important to comprehensively evaluate the deep feature learning algorithms for image/video classification on popular RGB-D datasets. We believe that this evaluation will provide insights and a deeper understanding of different deep learning algorithms for RGB-D feature extraction and fusion.

1.1. Related work to RGB-D information

In the past decades, since RGB images usually provide the limited appearance information of the objects in different scenes, it is extremely difficult to solve certain challenges such as the partition of the foreground and background which have the similar colors and textures. Besides that, the object appearance described by RGB images is sensitive to common variations, such as illuminance change. This drawback significantly impedes the usage of RGB based vision algorithms in real-world situations. Complementary to the RGB images, depth information for each pixel can help to better perceive the scene. RGB-D images/videos provide richer information, leading to more accurate and robust performance on vision applications.

The depth images/videos are generated by a depth sensor. Compared to early expensive and inconvenient range sensors (such as Konica Minolta Vivid 910), the **low-cost** 3D Microsoft Kinect sensor makes the acquisition of RGB-D data cheaper and easier. Therefore, the research of computer vision algorithms based on RGB-D data has attracted a lot of attention in the last few years [48]. Bo et al. [7] presented a hierarchical matching pursuit (HMP) based on sparse coding to learn new feature representations from RGB-D images in an unsupervised way. Tang et al. [40] designed a new feature called histogram of oriented normal vectors (HONV) to capture local 3-D geometric characteristics for object recognition on depth images. In [6], Blum et al. presented an algorithm that can automatically learn feature responses from the image, and the new feature descriptor encodes all available color and depth data into a concise representation. Spinello et al. introduced an RGB-D based people detection approach which combines a local depth-change detector employing HOD and RGB data HOG to detect the people from the RGB-D data in [38]. In [12], Endres et al. introduced an approach which describes a volumetric voxel representation through optimizing the 3D pose graph using the g^2o framework which can be directly used for path planning, robot localization and navigation.

1.2. Related work to deep learning methods

According to our evaluation, we select four representative deep learning methods including Deep Belief Networks (DBNs), Stacked Denoising Auto-Encoders (SDAE), Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) Neural Networks for our experiments. These methods have been widely applied in numerous contests in pattern recognition and machine learning. DBN is fine-tuned by backpropagation (BP) without any training pattern deformations which receives much success with 1.2% error rate on the MNIST handwritten digits [16]. Meanwhile, it achieved good results on phoneme recognition, with an error rate of 26.7% on the TIMIT core test set. SDAE was first introduced in [41] as an extension of Stacked auto-encoder (SAE) [28]. BP-trained CNNs achieved a new MNIST record of 0.39%. In 2012, GPU-implemented CNNs achieved the best results on the ImageNet classification benchmark [26]. LSTM won the ICDAR handwriting competition in 2009 and achieved a record 17.7% phoneme error rate on the TIMIT natural speech dataset in 2013. More relevant work and history on these four deep learning methods can be found in [34].

Currently, aiming to obtain more robust features from RGB and depth images/videos, various algorithms based on Deep Neural Networks have been proposed. R. Socher et al. presented convolutional and recursive neural networks (CNN-RNN) [37] to obtain higher order features. In CNN-RNN, CNN layers firstly learn low-level translationally invariant features, and then these features are given as inputs into multiple, fixed-tree RNNs. Bai et al. proposed subset based sparse auto-encoder and recursive neural networks (Sub-SAE-RNNs) [2] which first train the RGB-Subset-Sparse auto-encoder and the Depth-Subset-Sparse auto-encoder to extract features from RGB images and depth images separately for each subset. These learned features are then transmitted to RNNs to reduce the dimensionality and learn robust hierarchical feature representations. In order to combine hand-crafted features and machine learned features, Jin et al. used the Convolution Neural Networks (CNNs) to extract the machine learned representation and Locality-constrained Linear Coding (LLC) based spatial pyramid matching for hand-crafted features [23]. This new feature representation method can obtain both the advantages of hand-crafted features and machine learned features. From these above successful methods, we can observe that they are all the extensions of our selected methods (CNNs, DBNs, SDAE or LSTM). Therefore, it is important to explore the performance of our selected methods on different kinds of RGB-D datasets.

1.3. Deep learning methods for RGB-D data analysis

Since deep learning methods have shown to be useful for standard RGB vision tasks like object detection, image classification and semantic segmentation, more works on RGB-D perception naturally consider neural networks for learning representations from depth information. In general, the RGB-D vision problems that can be addressed or enhanced by means of the deep learning methods are summarized from four aspects: object detection and tracking, object and scene recognition, human activity analysis and indoor 3-D mapping. In this paper, our experiments focus on object and scene recognition, and human activity analysis.

1.3.1. Object detection and tracking

The depth information of an object is immune to object appearance changes, environmental illumination and subtle movements of the background. With the invention of the low-cost Kinect depth camera, researchers immediately realized that features based on depth information can significantly improve detecting and tracking objects in the real world where all kinds of variations occur. Depth-RCNN [13] is the first object detector using deep convolutional nets on RGB-D data, which is an extension of the RCNN framework. The depth map is encoded as three extra channels (with Geocentric Encoding: Disparity, Height, and Angle) appended to the color images. Furthermore, Depth-RCNN was extended to generate 3D bounding boxes through aligning 3D CAD models to the recognition results. Tracking via deep learning methods in RGB-D data is also an important topic. In [47], Xue et al. proposed to train a deep convolutional neural network, which improves tracking performance, to classify people in RGB-D videos.

1.3.2. Object and scene recognition

The conventional RGB-based deep learned features may suffer from the distortions of an object. RGB information is less capable of handling these environmental variations. Fortunately, the combination of RGB and depth information can potentially enhance the robustness of the deep learned features. Zaki et al. [49] presented an RGB-D object recognition framework which employed a CNN pre-trained on RGB data as feature extractors for both color and depth channels. Then they proposed a rich coarse-to-fine feature representation scheme, called Hypercube Pyramid, which can capture discriminatory information at different levels of detail. Zhu et al. [50] introduced a novel discriminative multi-modal fusion framework for RGB-D scene recognition which simultaneously considered the inter- and intra-modality correlation for all samples and meanwhile regularizing the learned features to be discriminative and compact. Then the results from the multimodal layer can be back-propagated to the lower CNN layers.

1.3.3. Human activity analysis

Apart from outputting both RGB and depth information, another contribution of Kinect is a fast human-skeletal tracking algorithm. This tracking algorithm can provide the exact location of each joint of the human body over time, which makes the representation of complex human activities easier. Wu et al. [45] proposed a novel method called Deep Dynamic Neural Networks (DDNN) for multimodal gesture recognition, which learns high-level spatiotemporal representations using deep neural networks suited to the input modality: a Gaussian-Bernoulli Deep Belief Network (DBN) to handle skeletal dynamics, and a 3D Convolutional Neural Network (3DCNN) to manage and fuse batches of depth and RGB images. Li et al. [30] proposed a feature learning network which is based on sparse auto-encoder (SAE) and principal component analysis for recognizing human actions.

1.3.4. Indoor 3-D mapping

The emergence of Kinect boosts the research for indoor 3-D mapping through deep learning methods due to its capability of providing depth information directly. Zhang et al. [24] proposed an approach to embed 3D context into the topology of a neural network trained for the performance of holistic scene understanding. After a 3D scene is depicted by a depth image, the network can align the observed scene with a predefined 3D scene template and then reason about the existence and location of each object within the scene template. To recover full 3D shapes from view-based depth images, Wu et al. [46] proposed to represent a geometric 3D shape as a probability distribution of binary variables on a 3D voxel grid through a Convolutional Deep Belief Network.

Aiming to make a comprehensive performance evaluation, we collect five representative datasets including two RGB-D object datasets [9,27], an RGB-D scene dataset [35], an RGB-D gesture dataset [31] and an RGB-D activity dataset [44] which can be divided into four categories: object classification, scene classification, gesture classification and action classification. This is the first work to comprehensively focus on the performance of deep learning methods on popular RGB-D datasets. In our experiments, in order to make the comparison of CNNs, DBNs, SDAE and LSTM under a fair environment, the pre-trained CNNs model through abundant RGB data and the RGB-D coding methods are not included. It is because that not all of these four deep learning methods can use other RGB data for pre-training and the particular RGB-D coding methods may not be suitable for all of the four kinds of deep learned features. Therefore, the design of our experiments is in a traditional way for providing insights and a deeper understanding of different deep learning algorithms for RGB-D feature extraction and fusion, which is introduced in detail in Section 4. In addition, besides results of the classification accuracies, our evaluation also provides a detailed analysis including confusion matrices and error analysis. Some tricks about adjusting hyper-parameters that we observed during our experiments are also given in this evaluation.

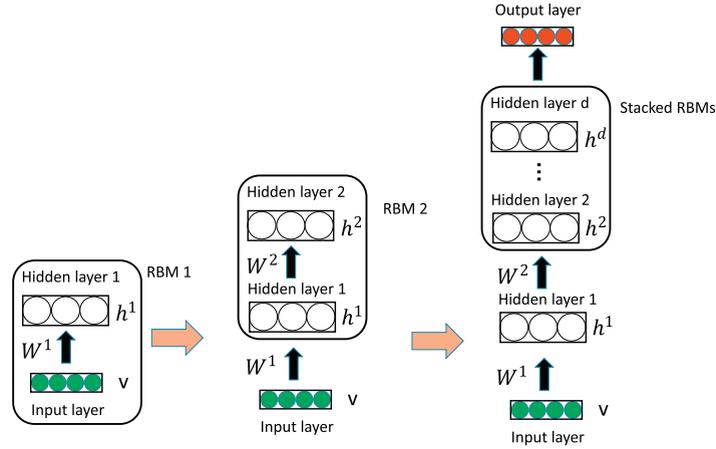


Fig. 1. The schematic representation of DBNs. It is just an example of DBNs structure. In practice, the number of units on each hidden layer is flexible.

The rest of this paper is organized as follows. In [Section 2](#), we briefly review the deep learning models which we use for evaluation in our experiments. In [Section 3](#), we present the data pre-processing techniques on deep learned features. [Section 4](#) describes experimental analysis, results and some tricks on our selected RGB-D datasets. Finally, we draw the conclusion in [Section 5](#).

2. Deep learning models

In recent years, many successful deep learning methods [15,41] as efficient feature learning tools have been applied to numerous areas. The aim of deep nets is to learn high-level features at each layer from the features learned at the previous layers. Some methods (such as DBNs [15] and SDAE [41]) have something in common: they have two steps in the training procedure - one is unsupervised pre-training and the other is fine-tuning. In the first step, through an unsupervised algorithm, the weights of the network are able to be better than random initialization. This phase can avoid local minima when doing supervised gradient descent. Therefore, we can consider that unsupervised pre-training is a regularizer. In the fine-tuning step, the criterion (the prediction error which uses the labels in a supervised task) is minimized. These two approaches for learning deep networks are shown to be essential to train deep networks. Other methods like CNNs [26] contain more connections than weights. The model itself realizes a form of regularization. The aim of this kind of neural networks is to learn filters, in a data-driven fashion, as a tool to extract features describing inputs. This is not only used in 2D convolutions but also can be extended into 3D-CNNs [22].

In this section, we will briefly introduce four deep learning models which are used in our experiments, DBNs, SDAE, CNNs and LSTM.

2.1. Deep Belief Networks

Deep Belief Networks (DBNs) stack many layers of unsupervised Restricted Boltzmann Machines (RBMs) in a greedy manner which was first introduced by Hinton et al. [15]. An RBM consists of visible layers and hidden layers. Each neuron on the layers is fully connected to all the neurons on the next layer. But there are no connections in the same layer. The learned weights are used to initialize a multi-layer neural network and then adjusted to the current task through supervised information for classification. A schematic representation of DBNs can be found in [Fig. 1](#).

In practice, the joint distribution $p(\mathbf{v}, \mathbf{h}; \theta)$ over the visible units \mathbf{v} and hidden units \mathbf{h} can be expressed as:

$$p(\mathbf{v}, \mathbf{h}; \theta) = \frac{\exp(-E(\mathbf{v}, \mathbf{h}; \theta))}{Z}, \quad (1)$$

where the model parameters $\theta = \mathbf{w}, \mathbf{a}, \mathbf{b}$ and $Z = \sum_{\mathbf{v}} \sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))$ is the normalization factor. The energy $E(\mathbf{v}, \mathbf{h}; \theta)$ of the joint configuration (\mathbf{v}, \mathbf{h}) is defined as:

$$E(\mathbf{v}, \mathbf{h}; \theta) = - \sum_{i=1}^V \sum_{j=1}^H w_{ij} v_i h_j - \sum_{i=1}^V b_i v_i - \sum_{j=1}^H a_j h_j, \quad (2)$$

where V and H are the numbers of the visible and hidden units. w_{ij} is the symmetric interaction between visible unit v_i and hidden unit h_j . b_i and a_j are the bias terms.

The marginal probability of the model to a visible vector \mathbf{v} is expressed as:

$$p(\mathbf{v}; \theta) = \frac{\sum_{\mathbf{h}} \exp(-E(\mathbf{v}, \mathbf{h}; \theta))}{Z}. \quad (3)$$

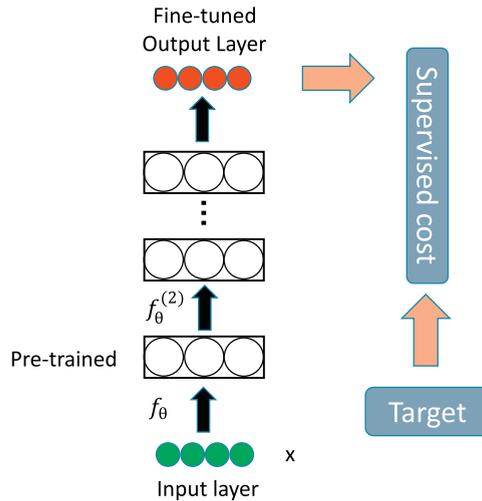


Fig. 2. A diagram of Stacked Denoising Auto-Encoders which includes an unsupervised pre-training step and a supervised fine-tuning step. Through performing gradient descent, the parameters are fine-tuned to minimize the error with the supervised target.

Therefore, according to the gradient of the joint likelihood function of data and labels, we can get the update rule of the **v-h** weights as

$$\Delta w_{ij} = \langle v_i h_j \rangle_{data} - \langle v_i h_j \rangle_{model}. \tag{4}$$

The greatest advantage of DBNs is the capability of “learning features” in a “layer-by-layer” manner. The higher-level features are learned from the previous layers. These features are believed to be more complicated and can better reflect the information which is contained in the structures of input data. Another advantage of DBNs is that it learns the generative model without imposing subjective selection of filters. Factored RBM is able to learn the filters while learning the feature activities in an unsupervised learning manner. It solves the concern of the legality of the selected filters. Meanwhile, it shows the biological implementation of visual cortex, namely, the receptive fields for cells in the primary visual cortex. However, a well-performing DBN requires a lot of empirically decided hyper-parameter settings, e.g., learning rate, momentum, weight cost number of epochs and number of layers. Inadequate selection of hyper-parameters will result in over-fitting and blow up DBNs. The property of DBNs that is sensitive to the empirically selected parameters has also been proved in our experiments. An improper set of hyper-parameters results in a huge difference from the best performance. To some extent, this disadvantage compromises the potential of DBNs.

2.2. Stacked Denoising Auto-Encoders

The Stacked Denoising Auto-Encoders (SDAE) [41] is an extension of the Stacked auto-encoder [28]. This model works in much the same way with DBNs. It also uses the greedy principle but stacks denoising auto-encoders to initialize a deep network. An auto-encoder consists of an encoder $h(\cdot)$ and a decoder $g(\cdot)$. Therefore, the reconstruction of the input x can be expressed as $Re(x) = g(h(x))$. Through minimizing the average reconstruction error $loss(x, Re(x))$, the reconstruction accuracy is able to be improved. This unsupervised pre-training is done on one layer at one time.

Same as DBNs, after all layers have been pre-trained, the parameters which can describe levels of representation about x are used as initialization to the deep neural network optimized with a supervised training criterion. In the fine-tuning stage, an output logistic regression layer is added to the top of the unsupervised pre-trained machine. Then, the classifier is fine-tuned through the design data set $D_x = \{d_{x_1}, \dots, d_{x_n}\}$ and the corresponding set of label codes $L_y = \{l_{y_1}, \dots, l_{y_n}\}$ to minimize the entropy loss function between the correct labels and the classifier’s predictions. A schematic diagram of Stacked Denoising Auto-Encoders is shown in Fig. 2.

For binary \mathbf{x} , the cross-entropy loss of the input vector $\mathbf{x} \in \{0, 1\}^d$ and the reconstructed d -dimensional vector $\hat{\mathbf{x}}$ is expressed as:

$$CEL(\mathbf{x}||\hat{\mathbf{x}}) = \sum_i CEL(x_i||\hat{x}_i) = - \sum_i (x_i \log \hat{x}_i + (1 - x_i) \log (1 - \hat{x}_i)), \tag{5}$$

where $\hat{\mathbf{x}} = \text{sigmoid}(c + w^T h(c(x)))$, c is the bias, and w is the transpose of the feed-forward weights. Additionally, another option is to use a Gaussian model.

SDAE makes use of different kinds of encoders to transform the input data, which can preserve a maximization of the mutual information between the original and the encoded information. Meanwhile, it utilizes a noise criterion for minimizing the transformation error. We mentioned that DBNs and SDAE have something in common: they have two steps in

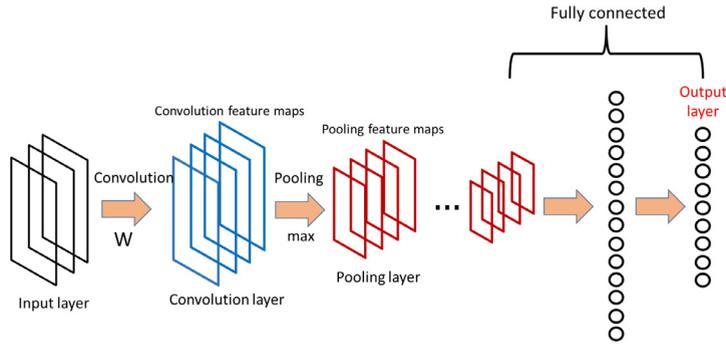


Fig. 3. The classical schematic representation of CNNs which includes an input layer, convolutional layers, max-pooling layers and an output layer. The fully connected part is also presented in the figure.

the training procedure—one is unsupervised pre-training and the other is fine-tuning. The advantage of using auto-encoders instead of RBMs as the unsupervised building block of a deep architecture is that as long as the training criterion is continuous in the parameters, almost any parametrization of the layers is possible [3]. However, in SDAE, training with gradient descent is slow and hard to parallelize. The optimization of SDAE is inherently non-convex and dependent on its initialization. Besides, since SDAE does not correspond to a generative model, unlike DBNs which is with generative models, samples cannot be drawn to check qualitatively what has been learned.

2.3. Convolutional Neural Networks

Convolutional Neural Networks [29] obtain much success in many visual processing tasks in recent years. This deep learning model is motivated by Hubel and Wiesel's work [20] on the cat's visual cortex. This visual cortex includes some cells which are sensitive to small sub-regions of the visual field. It can be called a receptive field. In practice, these cells can be considered as filters on the input space in the CNNs model. It has been proved that it is well-suited to extract the local correlation in natural images/videos.

Convolutional Neural Network consists of one image processing layer, one or more convolutional layers and fully connected layers and one classification layer. A classical schematic representation of CNNs is shown in Fig. 3. The image processing layer is a designed pre-processing layer which can keep being fixed in the training step. We introduce the pre-processing layer in Section 3 in detail. The convolutional layer applies a set of kernels of size $n \times n \times c$ that are able to process small local parts of the input. For most of the 2D-CNNs experiments, the input color images are often processed into gray images to enhance the efficiency and accuracy, therefore, the kernel size is often expressed as $n \times n$. Pooling is another important concept. It is a form of non-linear down-sampling where each map is sub-sampled with mean or max pooling over $m \times m$ contiguous regions (usually, m is from 2 to 5). It can improve translation invariance and tolerance to small differences of positions about object parts, at the same time, lead to faster convergence. The classification layer is fully connected which combines the outputs from the topmost convolutional layer into a feature vector, with one output unit per class label. Additionally, weight sharing is a significant principle since it is able to reduce the number of trainable parameters. More details concerning CNNs can be found in [8]. For a multi-label classification problem with F training examples and M classes, the squared-error is expressed as:

$$E^F = \frac{1}{2} \sum_{f=1}^F \sum_{m=1}^M (t_m^f - y_m^f)^2, \quad (6)$$

where t_m^f is the value of the m th dimension about f th pattern's corresponding label, and y_m^f is the m th output layer unit related to f th input pattern. In our experiments, for better results, we use 2D-CNNs for image datasets and 3D-CNNs for video datasets. Due to the space limitation, we do not give a detailed review of 3D-CNNs. More details can be found in [22].

One major advantage of CNNs is the use of shared weights in convolutional layers. The same filter is used for each pixel in the layer, which leads to the reduction of memory footprint and the improvement of result performance. For image classification applications, CNNs use relatively little pre-processing, which means that the network in CNNs is responsible to learn the filters. Without dependence on prior knowledge and human effort for designing features is another major advantage of CNNs. Besides, compared to traditional neural networks, CNN is more robust toward variation of input features. The neurons in the hidden layers are connected to the neurons that are in the same spatial area instead of being connected to all of the nodes in the previous layer. Furthermore, the resolution of the image data is reduced when calculating to higher layers in the network. However, besides a complex implementation, CNNs have another significant disadvantage that they require very large training data and consume an often impractical amount of time to learn the parameters of the network, which always take several days or weeks. Though the framework for accelerating training and classification of CNNs on

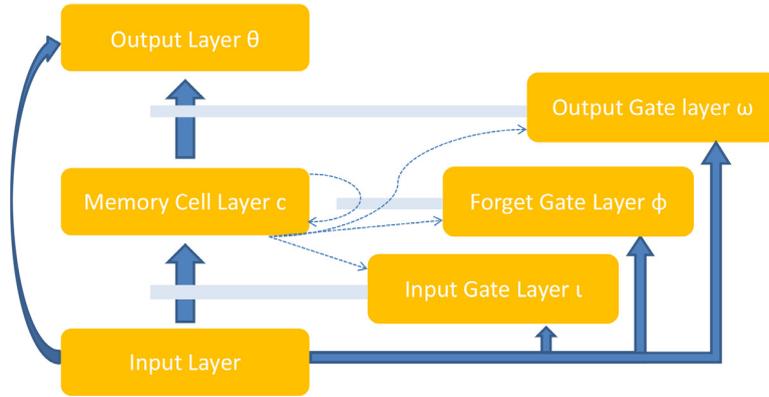


Fig. 4. The standard LSTM architecture. The memory block assemblies contain separate layers of memory cells, input gates, forget gates and output gates, in addition to the input layers and output layers. Blue solid arrows show full all-to-all connectivity between units in a layer. Blue dashed arrows mean connectivity only between the units in the two layers that have the same index. The light gray bars denote gating relationships. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Graphic Processing Units (GPUs) has been implemented and performs nearly hundreds of times faster than on the CPU, it is still not enough for real-world applications.

2.4. Long short-term memory neural networks

Long short-term memory (LSTM) is an extension of recurrent neural network (RNN) architecture which was first proposed in [17] for addressing the vanishing and exploding gradient problems of conventional RNNs. Different from traditional RNNs, when there exist long time lags of unknown size among important events, an LSTM network can classify, predict and process time series from experience. LSTM provides remedies for the RNN's weakness of exponential error decay through adding constant error carousel (CEC) which allows for constant error signal propagation along with the time. Besides, taking advantages of multiplicative gates can control the access to the CEC.

An LSTM architecture consists of an input layer, an output layer and a layer of memory block cell assemblies. A classical schematic representation of standard LSTM architecture is shown in Fig. 4. Fig. 4 shows that the memory block assemblies are composed of multiple separate layers: the input gate layer (ι), the forget gate layer (ϕ), the memory cell layer (c), and the output gate layer (ω). The input layer projects all of the connections to each of these layers. The memory cell layer projects all of the connections to the output layer (θ). Moreover, each memory cell c_j projects a single ungated peephole connection to each of its associated gates. The memory cell and the gates receive a connection from every neuron in the input layer. Through gated control, the network can effectively maintain and make use of past observations. An LSTM network computes a mapping from an input sequence $x = (x_1, \dots, x_T)$ to an output sequence $y = (y_1, \dots, y_T)$ through computing the network unit activations through the following equations iteratively from $t = 1$ to T [32]:

$$i_t = \sigma(W_{ix}x_t + W_{im}m_{t-1} + W_{ic}c_{t-1} + b_i), \quad (7)$$

$$f_t = \sigma(W_{fx}x_t + W_{fm}m_{t-1} + W_{fc}c_{t-1} + b_f), \quad (8)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot g(W_{cx}x_t + W_{cm}m_{t-1} + b_c), \quad (9)$$

$$o_t = \sigma(W_{ox}x_t + W_{om}m_{t-1} + W_{oc}c_t + b_o), \quad (10)$$

$$m_t = o_t \odot h(c_t), \quad (11)$$

$$y_t = W_{ym}m_t + b_y, \quad (12)$$

where the W terms denote weight matrices, the b terms denote bias vectors, σ is the logistic sigmoid function, and i, f, c and o represent the input gate, forget gate, cell activation vectors and output gate respectively, all of which are the same size as the cell output activation vector m . \odot is the element-wise product of the vectors. g and h are the cell input and cell output activation functions, generally \tanh .

LSTM can solve the vanishing gradient point problem in RNN. Meanwhile, LSTM has the capability of bridging long time lags between inputs, which can remember inputs up to 1000 time steps in the past. This advantage makes LSTM learn long sequences with long time lags. Besides, it appears that there is no need for parameter fine tuning in LSTM [17]. LSTM can work well over a broad range of parameters such as learning rate, input gate bias and output gate bias. However, in LSTM, the explicit memory adds more weights to each node, and all of these weights have to be trained. This increases the dimensionality of the task and potentially makes it harder to find an optimal solution.

Note that different from DBNs, SDAE and CNNs, LSTM is a sequence learning method which is hardly applied to image classification and object detection. Therefore, in our experiments, we only show the performance about LSTM on a gesture recognition dataset (SKIG dataset) and an action recognition dataset (MSRDailyActivity3D dataset).

3. Data preprocessing on deep learned features

Data preprocessing is an important part of the procedure of learning deep features. In practice, through a reasonable choice of preprocessing steps, it will result in a better performance according to the related task. Common preprocessing methods include normalization and PCA/ZCA whitening. Generally, one without much working experience about the deep learning algorithms will find it hard to adjust the parameters for raw data. When the data is processed in a small regular range, tuning parameters will become easier [11]. However, in the whole process of our experiments, we find that not every dataset is suitable to be either normalized or whitened. Therefore, we will have a test on the dataset and then choose the preprocessing steps according to the situations. Additionally, before we test the algorithms on the datasets, we will first observe properties of the data itself to gain more information which will help us to save more time.

3.1. Normalization

General normalization approaches include simple rescaling, per-example mean subtraction and feature standardization. The choice of these methods mainly depends on the data. In our experiments, since feature standardization is able to set every dimension of raw data to have zero-mean and unit-variance, at the same time, deep features will work with the linear SVM classifier, we choose feature standardization to normalize our data. Therefore, our data is normalized through first subtracting the mean of each dimension from each dimension and then dividing it by its standard deviation.

3.2. PCA/ZCA whitening

Following the step of feature standardization, we apply PCA/ZCA whitening to the entire dataset [21]. This is commonly used in deep learning tasks (e.g., [25]). Whitening cannot only make the deep learning algorithm work better but also speed up the convergence of the algorithm. However, in our experiments, for SDAE and DBNs, the results after whitening did not show an obvious improvement. To make the experiments under a fair environment, as long as whitening does not lead to a worse result, we choose to do ZCA whitening to the normalized data. Since we transfer RGB images to gray-scale images to make the data have the stationary property in our experiments and the data has been scaled into a reasonable range, the value of epsilon in ZCA whitening is set large (0.1) for low-pass filtering. More details about PCA/ZCA whitening can be found in [21].

4. Experiments on deep learning models

In this section, we evaluate four deep feature learning algorithms (DBNs, CNNs, SDAE and LSTM) on three popular image recognition datasets and two video recognition datasets including 2D&3D object dataset [9], RGB-D object dataset [27], NYU Depth v1 indoor scene segmentation dataset [35], Sheffield Kinect Gesture dataset (SKIG) [31] and MSRDailyActivity3D dataset [44]. Example images from these datasets are given in Fig. 5. Note that in our experiments, we only show the performance about LSTM on SKIG dataset and MSRDailyActivity3D dataset. In all of these five datasets, we follow the standard setting procedures according to the authors of their respective datasets. Over all of the datasets, we process raw RGB images into gray-scale images and choose the first channel of the depth images as training and test data. According to DBNs, CNNs, SDAE and LSTM, after weights are learned in the deep neural networks, we are able to extract the image or video features from the preprocessed images/videos. Then a linear SVM classifier is trained and tested on the related test sets. To make the results comprehensive, we compare the final results computed on deep features from RGB data only, deep features from depth data only, RGB-D features concatenation and deep features from RGB-D fusion. In RGB-D features concatenation experiments, we concatenate the feature vectors which are extracted from RGB data and depth data respectively into new vectors. Different from concatenation experiments, according to RGB-D fusion experiments, we firstly concatenate RGB images/frames and relative depth images/frames together, and then extract features from deep learning models. Illustration about these two experimental procedures is shown in Fig. 6. Detailed experimental settings, some important parameters, tricks and experiences about adjusting hyper-parameters are shown in the following subsections. All experiments are performed using Matlab 2013b and C++ on a server configured with a 16-core processor and 500G of RAM running the Linux OS.

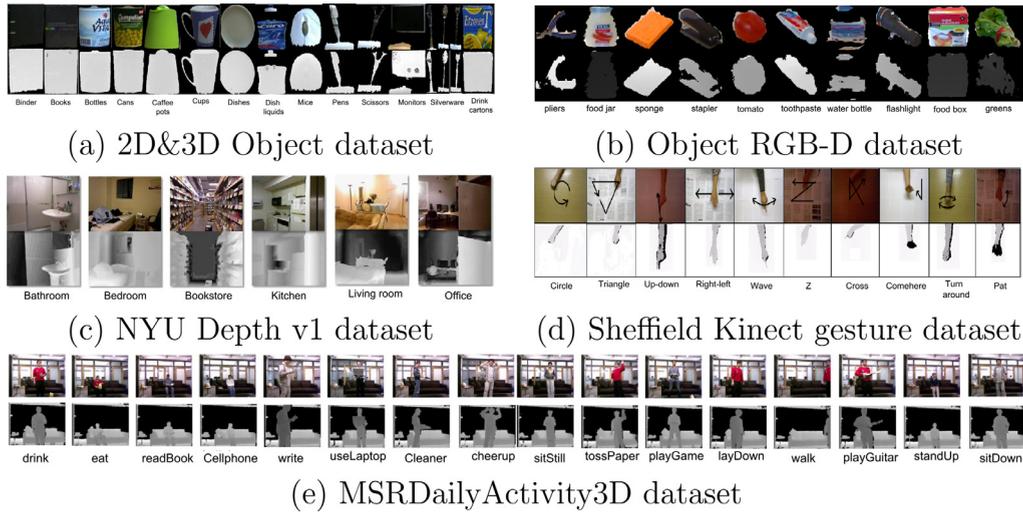
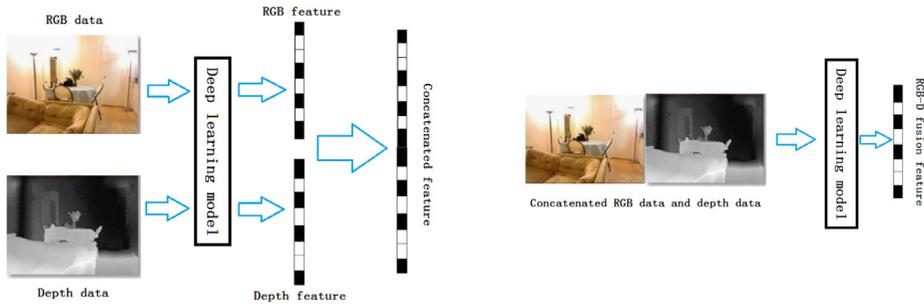


Fig. 5. Some example images in our selected datasets. From left to right in order, (a) 2D&3D Object dataset, (b) object RGB-D dataset, (c) NYU Depth v1 dataset, (d) Sheffield Kinect Gesture dataset and (e) MSRDailyActivity3D dataset.



(a) RGB-D features concatenation (b) Deep features from RGB-D fusion

Fig. 6. Illustration about two experimental procedures used in our evaluation work.

Table 1

The final comparison results between neural-network classifier and SVM on the 2D&3D object dataset. The second, fourth and seventh columns are the results of RGB test images, depth test images and RGB-D fusion test images on the neural-network classifier separately. The third, fifth, sixth and eighth columns are the results of RGB test images, depth test images, concatenated RGB-D image features and RGB-D fusion test images on SVM separately.

Method	RGB	RGB (SVM)	Depth	Depth (SVM)	RGB-D Concatenation (SVM)	RGB-D fusion	RGB-D fusion (SVM)
DBNs	72.1	74.5	75.7	78.6	82.3	78.3	79.1
CNNs	77.3	79.1	81.0	83.5	83.6	82.7	84.6
SDAE	73.0	74.5	74.2	75.6	79.3	77.6	78.4

4.1. 2D&3D object dataset

We evaluate deep feature learning for object category recognition on the 2D&3D object dataset [9]. This dataset includes 18 different categories (i.e., binders, books and scissors) with each of them containing 3–14 objects resulting in 162 objects. The views of each object are recorded every 10° along the vertical axis. Therefore, there are totally $162 \times 36 = 5832$ RGB images and $162 \times 36 = 5832$ depth images respectively. For the consistency with the setup in [9], since the low number of examples of classes perforator and phone, our experiments do not include them. Meanwhile, knives, forks and spoons are combined into one category ‘silverware’. We choose 6 objects per category for training, and the left are used for testing. If the number of objects in a category is less than 6 (e.g., scissors), 2 objects are added into the test. Since images are cropped in different sizes, we resize each image into 56×56 pixels. We give the final comparison results between neural-network classifier and SVM in Table 1.

The hyper-parameters of the DBNs, SDAE and CNNs models are described in Tables 2–4. Fig. 7 shows confusion matrixes about our three deep learning models across 14 classes on the 2D&3D dataset.

Table 2
Hyper-parameters about DBNs experiments on the 2D&3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	2
Units for each layer	100/100/100	100/100/100	100/100
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.009	0.009	0.008
Number of unsupervised epochs	13	13	13
Number of supervised epochs	17	30	24

Table 3
Hyper-parameters about SDAE experiments on the 2D&3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	2	2	2
Units for each layer	100/100	100/100	100/200
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.1	0.1	0.1
Number of unsupervised epochs	10	10	15
Number of supervised epochs	10	10	30

Table 4
Hyper-parameters about CNNs experiments on the 2D&3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of convolution layers	2	2	2
Number of sub-sampling layers	2	2	2
Kernel size	5	5	5
Learning rate	0.1	0.06	0.1
Number of epochs	30	60	30

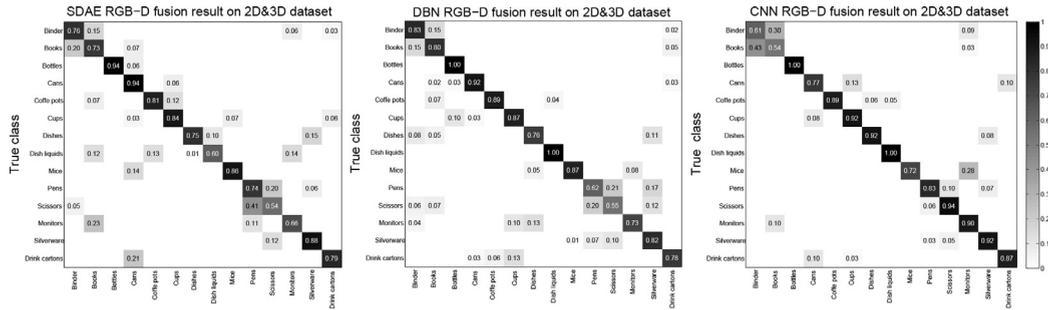


Fig. 7. Confusion matrixes about three deep learning models on the 2D&3D dataset. The labels on the vertical axis express the true classes and the labels on the horizontal axis denote the predicted classes.

From the comparison results of our experiments about three selected deep learning models on 2D&3D dataset in Table 1, it can be seen that the accuracy of RGB, depth and RGB-D fusion results through SVM outperforms that through the neural-network classifier. In each deep learning method, accuracies of RGB-D concatenation through SVM and RGB-D fusion features through SVM are higher than deep features from RGB data only and deep features from depth data only. In these three methods (DBNs, CNNs and SDAE), CNNs obtain the highest performance (84.6%). From the comparison of three confusion matrixes in Fig. 7, we can see that our three deep learning models all have the lowest error rates in bottles, cans, coffee pots and cups. Binders, books, pens and scissors have higher error rates. The main reason is that binders and books are similar in shape and color. Pens, scissors and silverware are similar in shape. It is worth to note that the error rates of binders and books in SDAE and DBNs are much lower than that of binders and books in CNNs, and the error rates of pens and scissors in SDAE and DBNs are much higher than that of pens and scissors in CNNs. The error rates of other categories are approximately similar. This interesting phenomenon may be due to the principle of the three different deep learning methods. In addition, it proves that in general SDAE and DBNs are more in common than CNNs.

4.2. Object RGB-D dataset

We test these deep learning algorithms on the second dataset called RGB-D object dataset. This dataset contains 41,877 images which are organized into 51 categories about 300 everyday objects such as apples, mushrooms and notebooks. All

Table 5

The final comparison results between neural-network classifier and SVM on object RGB-D dataset. The second, fourth and seventh columns are the results of RGB test images, depth test images and RGB-D fusion test images on the neural-network classifier separately. The third, fifth, sixth and eighth columns are the results of RGB test images, depth test images, concatenated RGB-D image features and RGB-D fusion test images on SVM separately.

Method	RGB	RGB (SVM)	Depth	Depth (SVM)	RGB-D concatenation (SVM)	RGB-D fusion	RGB-D fusion (SVM)
DBNs	80.9	81.6	75.1	78.6	84.3	82.4	83.7
CNNs	82.4	82.5	75.5	78.9	83.4	83.2	84.8
SDAE	81.4	82.0	71.9	73.7	82.3	82.6	84.2

Table 6

Hyper-parameters about DBNs experiments on object RGB-D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	3
Units for each layer	110/100/20	110/100/20	110/100/20
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.009	0.009	0.009
Number of unsupervised epochs	13	13	13
Number of supervised epochs	8	10	22

Table 7

Hyper-parameters about SDAE experiments on object RGB-D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	2	2	2
Units for each layer	100/100	130/100	110/200
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.1	0.08	0.05
Number of unsupervised epochs	10	15	15
Number of supervised epochs	15	30	30

Table 8

Hyper-parameters about CNNs experiments on object RGB-D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of convolution layers	2	2	2
Number of sub-sampling layers	2	2	2
Kernel size	5	5	5
Learning rate	0.1	0.06	0.03
Number of epochs	30	60	80

of the objects are segmented from the background through combining color and depth cues. Following the setup in [27], we choose to run category recognition experiments by randomly selecting one object from the categories for testing. Each image in object RGB-D dataset is resized into 56×56 pixels for consistency with the 2D&3D dataset. Table 5 summarizes the comparison between neural-network classifier and SVM.

The hyper-parameters of three deep learning models DBNs, SDAE and CNNs are shown in Tables 6–8.

As we can see from Table 5, CNNs outperform DBNs and SDAE by 0.5% and 0.3%. Due to the limitation of space, we only give the confusion matrix of the best performance (CNNs RGB-D fusion) in our experiments. Fig. 8 shows the confusion matrix about CNNs across 51 classes over object RGB-D dataset.

4.3. NYU Depth v1

Besides image object classification, we also evaluate these three deep feature learning models on indoor scene classification. NYU Depth v1 dataset consists of 7 different kinds of scene classes totally containing 2347 labeled frames. Since the standard classification protocol removes scene ‘cafe’ from the dataset, we use the remaining 6 different scenes. It is worth noting that since there are so many objects in one scene and the correlation between images in one scene is low, it makes NYU Depth v1 a very challenging dataset. The baseline when only using RGB images is 55% [35]. Table 9 shows the performance comparison between neural-network classifier and SVM on this dataset.

The hyper-parameters of DBNs, SDAE and CNNs can be found in Tables 10–12. Fig. 9 shows confusion matrixes about our three deep learning models across 6 classes over NYU Depth v1 dataset.

As we have mentioned above, NYU depth v1 dataset is very challenging. Therefore, in our three deep learning methods, CNNs achieve the best performance which is only 71.8%. Different from 2D&3D object dataset and object RGB-D dataset, RGB-D fusion through SVM always obtains the higher recognition accuracy (70.5% DBNs, 71.8% CNNs and 71.1% SDAE)

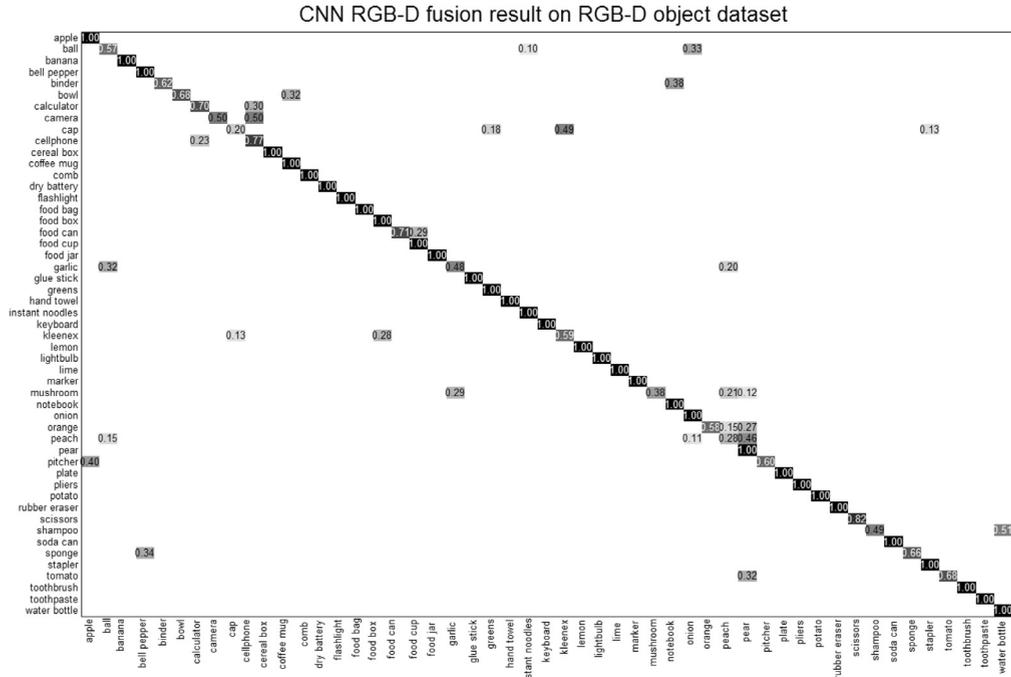


Fig. 8. Confusion matrix about CNNs on object RGB-D dataset. The labels on the vertical axis express the true classes and the labels on the horizontal axis denote the predicted classes.

Table 9

The performance comparison results between neural-network classifier and SVM on NYU Depth v1 dataset. The second, fourth and seventh columns are the results of RGB test images, depth test images and RGB-D fusion test images on the neural-network classifier separately. The third, fifth, sixth and eighth columns are the results of RGB test images, depth test images, concatenated RGB-D image features and RGB-D fusion test images on SVM separately.

Method	RGB	RGB (SVM)	Depth	Depth (SVM)	RGB-D Concatenation (SVM)	RGB-D fusion	RGB-D fusion (SVM)
DBNs	62.4	66.7	57.3	60.8	68.3	65.5	70.5
CNNs	68.4	69.5	56.5	56.9	70.4	70.1	71.8
SDAE	65.2	68.4	51.5	55.0	70.3	69.6	71.1

Table 10

Hyper-parameters about DBNs experiments on NYU Depth v1 dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	3
Units for each layer	120/100/80	120/100/80	110/100/100
Unsupervised learning rate	0.06	0.04	0.1
Supervised learning rate	0.006	0.008	0.008
Number of unsupervised epochs	3	3	3
Number of supervised epochs	35	45	22

Table 11

Hyper-parameters about SDAE experiments on NYU Depth v1 dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	3
Units for each layer	120/100/80	120/100/60	130/200/120
Unsupervised learning rate	0.01	0.01	0.01
Supervised learning rate	0.1	0.1	0.1
Number of unsupervised epochs	15	15	15
Number of supervised epochs	30	35	50

Table 12
Hyper-parameters about CNNs experiments on NYU Depth v1 dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of convolution layers	2	2	2
Number of sub-sampling layers	2	2	2
Kernel size	8	8	8
Learning rate	0.008	0.008	0.004
Number of epochs	50	45	80

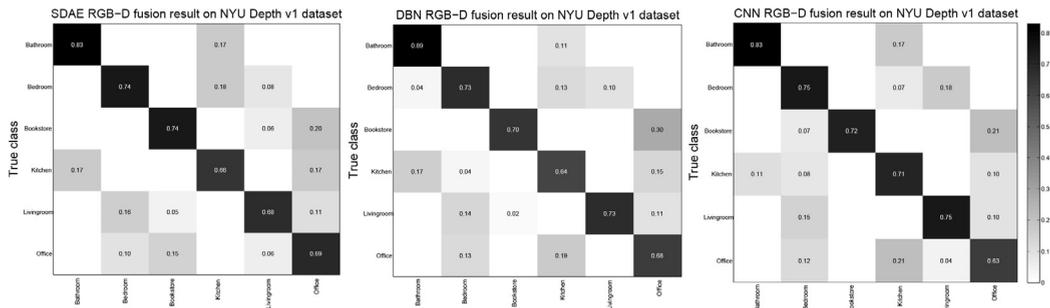


Fig. 9. Confusion matrixes about three deep learning models on NYU Depth v1 dataset. The labels on the vertical axis express the true classes and the labels on the horizontal axis denote the predicted classes.

Table 13

The performance comparison results between neural-network classifier and SVM on SKIG dataset. The second, fourth and seventh columns are the results of RGB test videos, depth test videos and RGB-D fusion test videos on the neural-network classifier separately. The third, fifth, sixth and eighth columns are the results of RGB test videos, depth test videos, concatenated RGB-D video features and RGB-D fusion test videos on SVM separately.

Method	RGB	RGB (SVM)	Depth	Depth (SVM)	RGB-D Concatenation (SVM)	RGB-D fusion	RGB-D fusion (SVM)
DBNs	78.3	83.1	68.9	73.8	84.7	81.5	85.9
3D-CNNs	87.2	91.3	77.5	82.2	92.6	88.1	93.3
SDAE	78.9	79.1	74.4	78.9	81.1	78.3	83.3
LSTM	82.6	83.1	75.7	77.5	87.2	86.7	91.3

Table 14
Hyper-parameters about DBNs experiments on SKIG dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	3
Units for each layer	120/100/100	120/100/100	110/100/100
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.01	0.009	0.006
Number of unsupervised epochs	3	3	3
Number of supervised epochs	30	40	55

compared to RGB-D concatenation (SVM) and RGB-D fusion. It may be because the scene images from NYU depth v1 dataset contain many irregular objects which seem much more complicated than the object images from the previous two datasets. From the confusion matrixes about these three deep learning methods, to a great extent, it can be seen that the distribution of error rates is similar.

4.4. Sheffield Kinect Gesture (SKIG) dataset

We also evaluate these four deep learning algorithms on video classification datasets. SKIG is a hand gesture dataset which contains 10 categories of hand gestures with 2160 hand gesture video sequences from six people, including 1080 RGB sequences and 1080 depth sequences respectively. In our experiments, since it has been proved that 5–7 frames (0.3–0.5 s of video) are enough to have the similar performance with the one obtainable with the entire video sequence [33]. Therefore, each video sequence is resized into $64 \times 48 \times 13$. Following the experimental setting in [31], we choose four objects as the training set and test on the remaining data. Table 13 shows the performance comparison between neural-network classifier and SVM on SKIG dataset. Additionally, since 3D-CNNs gain much success in video data classification, we use 3D-CNNs instead of 2D-CNNs in our experiments. We also compare LSTM Neural Networks experimentally in this subsection.

The hyper-parameters of DBNs, SDAE, 3D-CNNs and LSTM can be found in Tables 14–17.

Table 15
Hyper-parameters about SDAE experiments on SKIG dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	2	2	2
Units for each layer	100/80	100/85	100/100
Unsupervised learning rate	0.01	0.01	0.01
Supervised learning rate	0.01	0.015	0.01
Number of unsupervised epochs	12	15	30
Number of supervised epochs	1200	500	500

Table 16
Hyper-parameters about 3D-CNNs experiments on SKIG dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of convolution layers	2	2	2
Number of sub-sampling layers	2	2	2
First Kernel size	$7 \times 7 \times 7$	$7 \times 7 \times 7$	$7 \times 7 \times 7$
Second Kernel size	$7 \times 7 \times 5$	$7 \times 7 \times 5$	$7 \times 7 \times 5$
Initial Learning rate	0.0005	0.0005	0.0004
Number of epochs	40	45	60

Table 17
Hyper-parameters about LSTM experiments on SKIG dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Memory blocks	50	50	60
Output neurons	10	10	10
Learning rate	0.0001	0.0001	0.0001
Number of epochs	2000	2000	2500

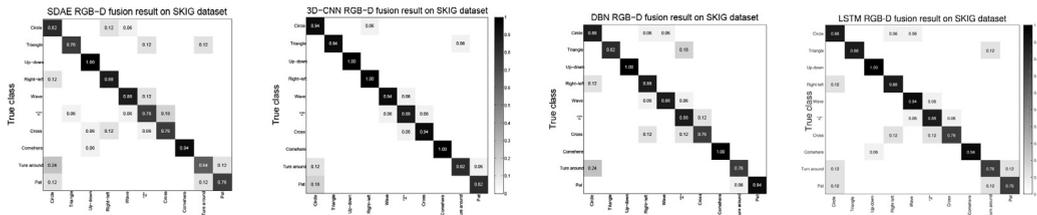


Fig. 10. Confusion matrixes about four deep learning models on SKIG dataset. The labels on the vertical axis express the true classes and the labels on the horizontal axis denote the predicted classes. From left to right in order, SDAE, 3D-CNN, DBN and LSTM.

To get better results in the 3D-CNNs model, we decay the learning rate a half in each epoch.

Fig. 10 shows confusion matrixes about our four deep learning models across 10 classes on the SKIG dataset.

From the comparison of these four deep learning models in Table 13, we can see that 3D-CNNs achieve the best performance among four—93.3%. It may be because that 3D-CNNs consider the more temporal correlation between video frames [22]. Sequence learning method LSTM with raw pixel features achieves 91.3% on the SKIG dataset, which is better than the performances of DBN and SDAE. It is reasonable because LSTM can learn from experience to classify, process and predict time series. Overall, we obtain high accuracies in this dataset. The main reason is that the ten categories in SKIG dataset can be classified easily. Each category is much different from other categories, and every test video in one category is similar to other test videos in the same category. Therefore, in terms of SKIG dataset, inter-class distance is big and intra-class distance is small. The analysis above suggests that deep learning will produce a good performance with less training samples if the experimental dataset is not challenging.

4.5. MSRDailyActivity3D dataset

The last dataset which we test on is MSRDailyActivity3D dataset [44]. It is a daily activity dataset which contains 16 activity types (e.g., drink, eat, play game). There are 10 subjects with each of them performs each activity twice, once in standing position, and once in sitting position. We do the same preprocessing procedure like SKIG and resize each sequence to $64 \times 48 \times 13$. Then subject 1 to subject 5 of “sitting on sofa” and subject 1 to subject 5 of “standing” in this dataset are used as training set and the rest are used for evaluation. Table 18 shows the accuracies of four deep learning methods.

The hyper-parameters of DBNs, SDAE, 3D-CNNs and LSTM are shown in Tables 19–22.

Table 18

The performance comparison results between neural-network classifier and SVM on MSRDailyActivity3D dataset. The second, fourth and seventh columns are the results of RGB test videos, depth test videos and RGB-D fusion test videos on the neural-network classifier separately. The third, fifth, sixth and eighth columns are the results of RGB test videos, depth test videos, concatenated RGB-D video features and RGB-D fusion test videos on SVM separately.

Method	RGB	RGB (SVM)	Depth	Depth (SVM)	RGB-D Concatenation (SVM)	RGB-D fusion	RGB-D fusion (SVM)
DBNs	51.9	62.5	50.6	53.1	66.3	65.0	68.1
3D-CNNs	50.5	65.6	47.3	58.2	61.3	61.3	68.9
SDAE	57.5	59.4	46.3	48.1	64.4	62.5	66.3
LSTM	49.4	64.4	46.3	57.5	63.1	60.0	68.1

Table 19

Hyper-parameters about DBNs experiments on MSRDailyActivity3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	3	3	3
Units for each layer	120/100/100	120/100/100	110/100/100
Unsupervised learning rate	0.1	0.1	0.1
Supervised learning rate	0.004	0.008	0.005
Number of unsupervised epochs	4	4	4
Number of supervised epochs	55	46	60

Table 20

Hyper-parameters about SDAE experiments on MSRDailyActivity3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of hidden layers	2	2	2
Units for each layer	110/80	110/85	100/100
Unsupervised learning rate	0.01	0.01	0.01
Supervised learning rate	0.01	0.015	0.01
Number of unsupervised epochs	15	20	33
Number of supervised epochs	1000	800	800

Table 21

Hyper-parameters about 3D-CNNs experiments on MSRDailyActivity3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Number of convolution layers	2	2	2
Number of sub-sampling layers	2	2	2
First Kernel size	$7 \times 7 \times 7$	$7 \times 7 \times 7$	$7 \times 7 \times 7$
Second Kernel size	$7 \times 7 \times 5$	$7 \times 7 \times 5$	$7 \times 7 \times 5$
Initial Learning rate	0.0003	0.0005	0.0004
Number of epochs	50	45	60

Table 22

Hyper-parameters about LSTM experiments on MSRDailyActivity3D dataset.

Selected hyper-parameters	RGB	Depth	RGB-D fusion
Memory blocks	60	60	70
Output neurons	16	16	16
Learning rate	0.0001	0.0001	0.0001
Number of epochs	2000	2000	2500

To get better results in the 3D-CNNs model, we use the same trick as in the experiments of SKIG Dataset by decaying the learning rate a half in every epoch.

In our deep learning experiments on MSRDailyActivity3D dataset, 3D-CNNs achieve a higher accuracy (68.9%) than DBNs (68.1%), SDAE (66.3%) and LSTM (68.1%). However, compared to the performances of SKIG dataset, we only obtain lower accuracies. There are two main reasons. First, it is a very challenging video dataset. According to this dataset, inter-class distance is small and intra-class distance is big. Second, there are not enough training samples for deep learning models. Therefore, it can be seen that it will show a bad performance with less training samples if the experimental dataset is very challenging. Fig. 11 shows confusion matrixes about our four deep learning models across 16 classes over MSRDailyActivity3D dataset.

Acknowledgments

This work was supported in part by National Natural Science Foundation of China under grant 6152810 and grant U1301251, and in part by the Beijing Natural Science Foundation under grant 4141003.

References

- [1] P. Allen, Robotic Object Recognition Using Vision and Touch, 34, 2012.
- [2] J. Bai, Y. Wu, J. Zhang, F. Chen, Subset based deep learning for RGB-D object recognition, *Neurocomputing* 165 (2015) 280–292.
- [3] Y. Bengio, Learning deep architectures for ai, *Found. Trends® Mach. Learn.* 2 (1) (2009) 1–127.
- [4] J. Bergstra, Y. Bengio, Random search for hyper-parameter optimization, *J. Mach. Learn. Res.* 13 (1) (2012) 281–305.
- [5] J.S. Bergstra, R. Bardenet, Y. Bengio, B. Kégl, Algorithms for hyper-parameter optimization, in: *Advances in Neural Information Processing Systems*, 2011, pp. 2546–2554.
- [6] M. Blum, J.T. Springenberg, J. Wülfing, M. Riedmiller, A learned feature descriptor for object recognition in RGB-D data, in: *IEEE International Conference on Robotics and Automation*, 2012, pp. 1298–1303.
- [7] L. Bo, X. Ren, D. Fox, Unsupervised feature learning for RGB-D based object recognition, in: *Experimental Robotics*, 2013, pp. 387–402.
- [8] J. Bouvrie, Notes on Convolutional Neural Networks, 2006.
- [9] B. Browatzki, J. Fischer, B. Graf, H. Bulthoff, C. Wallraven, Going into depth: Evaluating 2D and 3D cues for object classification on a new, large-scale object dataset, in: *International Conference on Computer Vision Workshops*, 2011, pp. 1189–1195.
- [10] Z. Cai, J. Han, L. Liu, L. Shao, RGB-D datasets using Microsoft Kinect or similar sensors: a survey, *Multimed. Tools Appl* (2016) 1–43, doi:10.1007/s11042-016-3374-6.
- [11] A. Coates, A.Y. Ng, H. Lee, An analysis of single-layer networks in unsupervised feature learning, in: *International Conference on Artificial Intelligence and Statistics*, 2011, pp. 215–223.
- [12] F. Endres, J. Hess, N. Engelhard, J. Sturm, D. Cremers, W. Burgard, An evaluation of the RGB-D slam system, in: *IEEE International Conference on Robotics and Automation*, 2012, pp. 1691–1696.
- [13] S. Gupta, P. Arbeláez, R. Girshick, J. Malik, Aligning 3D models to RGB-D images of cluttered scenes, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 4731–4740.
- [14] J. Han, L. Shao, D. Xu, J. Shotton, Enhanced computer vision with Microsoft Kinect sensor: a review, *IEEE Trans. Cybern.* 43 (5) (2013) 1318–1334.
- [15] G.E. Hinton, S. Osindero, Y.-W. Teh, A fast learning algorithm for deep belief nets, *Neural Comput.* 18 (7) (2006) 1527–1554.
- [16] G.E. Hinton, R.R. Salakhutdinov, Reducing the dimensionality of data with neural networks, *Science* 313 (5786) (2006) 504–507.
- [17] S. Hochreiter, J. Schmidhuber, Long short-term memory, *Neural Comput.* 9 (8) (1997) 1735–1780.
- [18] R. Hong, Z. Hu, R. Wang, M. Wang, D. Tao, Multi-view object retrieval via multi-scale topic models, *IEEE Trans. Image Process.* 25 (12) (2016).
- [19] R. Hong, L. Zhang, C. Zhang, R. Zimmermann, Flickr circles: aesthetic tendency discovery by multi-view regularized topic modeling, *IEEE Trans. Multimed.* 18 (8) (2016) 1555–1567.
- [20] D.H. Hubel, T.N. Wiesel, Receptive fields and functional architecture of monkey striate cortex, *J. Physiol. (Lond.)* 195 (1) (1968) 215–243.
- [21] A. Hyvärinen, E. Oja, Independent component analysis: algorithms and applications, *Neural Networks* 13 (4) (2000) 411–430.
- [22] S. Ji, W. Xu, M. Yang, K. Yu, 3D convolutional neural networks for human action recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (1) (2013) 221–231.
- [23] L. Jin, S. Gao, Z. Li, J. Tang, Hand-crafted features or machine learnt features? Together they improve RGB-D object recognition, in: *International Symposium on Multimedia*, 2014, pp. 311–319.
- [24] Y.Z.M.B.P. Kohli, S. Izadi, J. Xiao, Deepcontext: context-encoding neural pathways for 3D holistic scene understanding, *arXiv preprint arXiv:1603.04922* (2016).
- [25] A. Krizhevsky, G.E. Hinton, et al., Factored 3-way restricted Boltzmann machines for modeling natural images, in: *International Conference on Artificial Intelligence and Statistics*, 2010, pp. 621–628.
- [26] A. Krizhevsky, I. Sutskever, G.E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [27] K. Lai, L. Bo, X. Ren, D. Fox, A large-scale hierarchical multi-view RGB-D object dataset, in: *International Conference on Robotics and Automation*, 2011, pp. 1817–1824.
- [28] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, Y. Bengio, An empirical evaluation of deep architectures on problems with many factors of variation, in: *International Conference on Machine Learning*, 2007, pp. 473–480.
- [29] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, Gradient-based learning applied to document recognition, *Proc. IEEE* 86 (11) (1998) 2278–2324.
- [30] S.-Z. Li, B. Yu, W. Wu, S.-Z. Su, R.-R. Ji, Feature learning based on SAE-PCA network for human gesture recognition in RGBD images, *Neurocomputing* 151 (2015) 565–573.
- [31] L. Liu, L. Shao, Learning discriminative representations from RGB-D video data, in: *International Joint Conference on Artificial Intelligence*, 2013, pp. 1493–1500.
- [32] H. Sak, A. Senior, F. Beaufays, Long short-term memory based recurrent neural network architectures for large vocabulary speech recognition, *arXiv preprint arXiv:1402.1128* (2014).
- [33] K. Schindler, L. Van Gool, Action snippets: how many frames does human action recognition require? in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [34] J. Schmidhuber, Deep learning in neural networks: an overview, *Neural Networks* 61 (2015) 85–117.
- [35] N. Silberman, R. Fergus, Indoor scene segmentation using a structured light sensor, in: *International Conference on Computer Vision Workshops*, 2011, pp. 601–608.
- [36] J. Snoek, H. Larochelle, R.P. Adams, Practical Bayesian optimization of machine learning algorithms, in: *Advances in Neural Information Processing Systems*, 2012, pp. 2951–2959.
- [37] R. Socher, B. Huval, B. Bath, C.D. Manning, A.Y. Ng, Convolutional-recursive deep learning for 3D object classification, in: *Advances in Neural Information Processing Systems*, 2012, pp. 665–673.
- [38] L. Spinello, C. Stachniss, W. Burgard, Scene in the loop: towards adaptation-by-tracking in RGB-D data, in: *Proc. Workshop RGB-D, Adv. Reason. Depth Cameras*, 2012.
- [39] R. Szeliski, *Computer vision: algorithms and applications*, 2010.
- [40] S. Tang, X. Wang, X. Lv, T.X. Han, J. Keller, Z. He, M. Skubic, S. Lao, Histogram of oriented normal vectors for object recognition with a depth sensor, in: *Asian Conference on Computer Vision*, 2013, pp. 525–538.
- [41] P. Vincent, H. Larochelle, Y. Bengio, P.-A. Manzagol, Extracting and composing robust features with denoising autoencoders, in: *International Conference on Machine Learning*, 2008, pp. 1096–1103.
- [42] S. Vishwakarma, A. Agrawal, A survey on activity recognition and behavior understanding in video surveillance, *Vis. Comput.* 29 (10) (2013) 983–1009.
- [43] L. Wan, M. Zeiler, S. Zhang, Y.L. Cun, R. Fergus, Regularization of neural networks using dropconnect, in: *International Conference on Machine Learning*, 2013, pp. 1058–1066.
- [44] J. Wang, Z. Liu, Y. Wu, J. Yuan, Mining actionlet ensemble for action recognition with depth cameras, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2012, pp. 1290–1297.

- [45] D. Wu, L. Pigou, P.-J. Kindermans, L. Nam, L. Shao, J. Dambre, J.-M. Odobez, Deep dynamic neural networks for multimodal gesture segmentation and recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (8) (2016) 1583–1597.
- [46] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, J. Xiao, 3D shapenets: a deep representation for volumetric shapes, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1912–1920.
- [47] H. Xue, Y. Liu, D. Cai, X. He, Tracking people in RGBD videos using deep learning and motion clues, *Neurocomputing* 204 (2016) 70–76.
- [48] M. Yu, L. Liu, L. Shao, Structure-preserving binary representations for RGB-D action recognition, *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (8) (2016) 1651–1664.
- [49] H.F. Zaki, F. Shafait, A. Mian, Convolutional hypercube pyramid for accurate RGB-D object category and instance recognition, in: *IEEE International Conference on Robotics and Automation*, 2016, pp. 1685–1692.
- [50] H. Zhu, J.-B. Weibel, S. Lu, Discriminative multi-modal feature fusion for RGBD indoor scene recognition, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2969–2976.