

Multi-Loss Regularized Deep Neural Network

Chunyan Xu, Canyi Lu, Xiaodan Liang, Junbin Gao, Wei Zheng, Tianjiang Wang, Shuicheng Yan

Abstract—A proper strategy to alleviate overfitting is critical for deep neural network (DNN). In this work, we introduce the cross-loss-function regularization for boosting the generalization capability of the DNN, which results in the Multi-Loss regularized Deep Neural Network (ML-DNN) framework. For a particular learning task, e.g., image classification, only a single loss function is used for all previous DNNs, and the intuition behind the multi-loss framework is that the extra loss functions with different theoretical motivations (e.g., pairwise loss and LambdaRank loss) may drag the algorithm away from the overfitting to one particular single loss function (e.g., softmax loss). In the training stage, we pre-train the model with the single core loss function, and then warm-start the whole ML-DNN with the convolutional parameters transferred from the pre-trained model. In the testing stage, the outputs by ML-DNN from different loss functions are fused with average pooling to produce the ultimate prediction. The experiments conducted on several benchmark datasets (CIFAR-10, CIFAR-100, MNIST and SVHN) demonstrate that the proposed ML-DNN framework, instantiated by the recently proposed Network-in-Network, considerably outperforms all other state-of-the-art methods.

Index Terms—Deep neural network, overfitting, multi-loss, visual classification.

I. INTRODUCTION

Due to the increasing computing power and availability of large training data, there has been a resurgence of interest in neural networks. Especially, the deep neural network learning framework [1], [2], [3], [4] has drawn much attention recently, which has achieved very promising performance on the kind of vision tasks, e.g., image classification [5], [6], pedestrian detection [7], and scene labeling [8], etc. The successes of the deep neural network frameworks have been largely attributed to the considerable number of model parameters in many convolution layers. For example, the winner of ILSVRC 2014 classification challenge, i.e., GoogLeNet [9], employed 22 layers deep network. Although the Deep Neural Network (DNN) with a large number of parameters has a powerful machine learning capability, it often suffers from overfitting.

In recent literatures, some regularization techniques have been proposed to prevent the heavy overfitting during training

DNNs, such as data augmentation [10], [11], weight decay [12], dropout [13], [14], drop-connect [15] and stochastic pooling [16]. Specifically, the most simple method to alleviate overfitting [10], [11] is to manually enlarge the training images with the label-preserving transformations, such as horizontal/vertical reflection, and image translation. The weight decay [12] is also a technique to help prevent overfitting by adding a penalty to the maximum likelihood. In [13], [14], the dropout technique stochastically sets half of the activations within a layer to zero for each training sample in the training stage, while the drop-connect method [15] is proposed to regularize large fully-connected layers in neural networks. The stochastic pooling method [16], which replaces the conventional deterministic pooling operation with a stochastic procedure, and randomly picks up the activations within each pooling region according to a multinomial distribution of the activations within the pooling region. However, these approaches may suffer from the certain limitation (e.g., overfitting) for a particular loss function. The optimization target (e.g. loss function) of visual recognition has not been well explored to prevent overfitting, which is truly critical for training the network.

In this work, we propose the Multi-Loss regularized Deep Neural Network (ML-DNN) framework to harness the regularizations among a set of different loss functions. The intuition behind this multi-loss framework is that the loss functions (e.g., pairwise ranking loss and LambdaRank loss) with different theoretical motivations may prevent the algorithmic overfitting to one single loss function (e.g., softmax loss). Specifically, the softmax loss function [1] can be used for minimizing the cross-entropy loss over all training samples. The pairwise ranking loss [18] takes the preference order of label pairs into count in the classification problems, while the LambdaRank loss function [19], [20] is introduced for optimizing the top- k classification accuracy. Our ML-DNN learning framework can effectively combine all the valuable properties of these loss functions by the cross-loss function regularization. Each kind of loss function in our MK-DNN is selected based on different theoretical motivations. Multiple loss functions of our ML-DNN can constrain the parameters of a neural network from different aspects, and help to regularly learn a DNN for boosting its discriminative capability.

As indicated in Figure 1, our proposed framework proposes to use the cross-loss-function regularization in the output layers, which combines the softmax loss, pairwise ranking loss LambdaRank top-1 loss and LambdaRank top-2 loss functions. In this work, we instantiate our ML-DNN framework under the architecture of Network In Network (NIN) [4], which is one of the most popular architectures used in computer vision problems due to its model discriminability for local patches within the receptive field and its reduced number

Email address: (xuchunyan01, canyilu, xdliang328)@gmail.com (C. Xu, C. Lu, X. Liang), tjwang@hust.edu.cn (T. Wang), jbgao@csu.edu.au (J. Gao), w0209.zheng@samsung.com (W. Zheng) and eleyans@nus.edu.sg (S. Yan).

C. Xu is with the School of Computer Science and Technology from Nanjing University of Science and Technology, Nanjing, 210094, China and the Department of Electrical and Computer Engineering at National University of Singapore, 117583, Singapore.

S. Yan, C. Lu and X. Liang are with Department of Electrical and Computer Engineering, National University of Singapore, 117583, Singapore.

J. Gao is with School of Computing and Mathematics, Charles Sturt University, Bathurst, NSW 2795, Australia.

T. Wang are with School of Computer Science and Technology, Huazhong University of Science and Technology, Wuhan, Hubei, 430074, China.

W. Zheng is with Beijing Samsung Telecom R&D Center, China.

This work was done when Chunyan Xu was a visiting scholar at National University of Singapore.

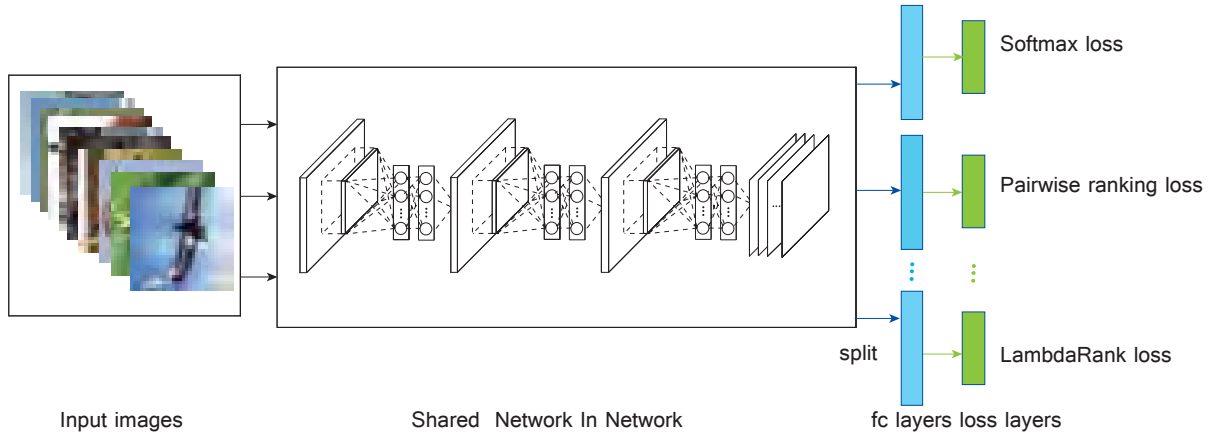


Fig. 1. Illustration of the structure of the proposed multi-loss regularized deep neural network (ML-DNN). The shared Network In Network (NIN) has a similar network structure to [4], except for the last loss layer. We feed some given images into the shared NIN and pass the output feature maps from the shared NIN into multiple loss branches. Each loss branch is composed of a full-connected (fc) layer and a loss layer. Example images are from the CIFAR-10 dataset [17].

of parameters. During the ML-DNN training, the training images are firstly fed into the shared NIN and several parallel fully-connected (fc) layers, each of which corresponds to some different loss layers. We pre-train the model with the single loss function, and then warm-start the whole ML-DNN with the convolutional parameters transferred from the pre-trained model. We feed the convolutional feature maps of the shared NIN into multiple branches of the loss functions. Based on these multiple loss functions, the network is trained by simultaneously optimizing all loss functions with back-propagation. In the testing stage, the outputs by ML-DNN from different loss functions are fused with average pooling to produce the ultimate prediction.

The major contributions of this work can be summarized as follows: (1) We propose a novel ML-DNN framework, which gracefully optimizes the architecture of DNN based on the cross-loss-function regularization. (2) For the image classification task, we present some loss functions (e.g, pairwise loss and LambdaRank top- k loss) for learning a DNN. Multiple loss functions are simultaneously optimized with the stochastic gradient descent (SGD) learning method. (3) Our ML-DNN is a very general framework for alleviating the overfitting during learning DNN. Any CNN architectures and any loss functions for different vision tasks can be conveniently incorporated into our framework. (4) The classification results on several standard datasets well verify the effectiveness of our proposed ML-DNN framework.

II. RELATED WORK

Deep neural network learning has been long studied and applied in the field of computer vision [1], [21], [22], [23], [24]. More than a decade ago, LeCun *et al.* [2] trained multilayer neural networks with the back-propagation algorithm and the gradient learning technique, and then demonstrated its effectiveness on the handwritten digit recognition task. Recently, there has been a resurgence of research interest in neural networks.

Neural Network Structure for Deep Learning: A classic convolutional network is composed of alternatively stacked convolutional layers and spatial pooling layers. The convolutional layer is to extract feature maps by linear convolutional filters followed by nonlinear activation functions (e.g., rectifier, sigmoid, tanh, etc.) Spatial pooling is to group the local features together from spatially adjacent pixels, which is typically done to improve the robustness to slight deformations of objects. The Convolutional Neural Network (CNN) [1] is a special type of neural network that consists of five convolution layers, some of which are followed by max-pooling layers, and three fully-connected layers with a final 1000-way softmax. The deep CNN has exhibited good generalization power in image-related applications. Recently, Krizhevsky *et al.* [1] achieved a breakthrough, outperforming the existing hand-crafted features on ILSVRC 2012 which contains 1000 object classes. Another deep network structure, namely “Network In Network” (NIN) [4] is proposed to build a micro network with more complex structures to abstract the data within the receptive field. It enhances model discriminability for local patches within the receptive field.

Regularization Techniques for Deep Learning: Learning neural network models is prone to overfitting because of the large number of parameters of the models. Some regularization techniques [12], [13], [11] are necessary for learning deep neural networks. Specifically, some data augmentation techniques [10], [11] are usually used to enlarge the training data, such as image translations, horizontal reflections, image rotation and scaling, etc. Dropout, recently proposed by Hinton *et al.* [13], is another regularization approach that stochastically sets half the activations within a layer to zero for each training sample during training. Wan *et al.* [15] proposed the drop-connect regularization, which sets a randomly selected subset of weights within the network to zero. Weight decay [12] was also proposed to prevent overfitting by adding a penalty to the maximum likelihood. The stochastic pooling method [16] was proposed to randomly select the pooled map responses by sampling from a multinomial distribution formed

from the activations of each pooling region.

Loss Functions for Deep Learning: Some loss functions have also been introduced for learning a single DNN, e.g., softmax loss function for the classification task [5], [4], ℓ_2 -norm loss function for the detection task, and sigmoid cross entropy loss function for the multi-label classification task. Recently, Gkioxari *et al.* [25] trained a single Convolutional Neural Network (CNN) [1] jointly for solving person detection, pose estimation, and action classification tasks, where each task is associated with a loss function. Li *et al.* [26] simultaneously learned a pose-joint regressor and a sliding-window body-part detector in a deep neural network, which can be seen as a heterogeneous multi-task learning method. The above works can be classified into two categories: single task learning with a single loss function and multi-task learning with multiple loss functions, both in a single deep neural network. Recently, Szegedy *et al.* [9] also proposed a deep convolutional neural network architecture, in which three softmax loss functions are adopted, one is a total loss of the network and another two can be regarded as the auxiliary classifiers in the training process. At the testing state, the two auxiliary loss functions and their corresponding sub-networks are discarded. This work mainly improves the discriminative capability of a neural network by appending multiple same loss functions in different stages of neural network, which alleviate the problem of vanishing gradient in a deeper CNN architecture to some extent. However, different from their usage of same loss functions, we use multiple different loss functions with different theoretical motivations in the end of a neural network, which allow each training sample to impose much more constraint on the parameters of a neural network and prevent the neural network from being too sure. From the aspect of regularization, it also should be noted that our ML-DNN can be posed as a good regularizer for preventing overfitting by using multiple loss function in one network.

III. MULTI-LOSS REGULARIZED DNN

In this section, we first introduce the network structure of the proposed ML-DNN. Second, several loss functions are investigated for cross-loss-function regularization. Third, we show how to learn a ML-DNN by simultaneously optimizing multiple loss functions. Finally, testing with ML-DNN is conducted with the average pooling algorithm.

A. Network Structure

The structure of the overall ML-DNN framework is mainly composed of two parts, namely the shared NIN and multi-loss branches, as shown in Figure 1.

• **Shared NIN:** The shared Network In Network (NIN) has a similar network architecture as NIN [4]. By building micro neural networks with more complex structures to abstract the data within the receptive field, Lin *et al.* [4] proposed a novel deep network structure called Network In Network (NIN) to enhance model discriminability for local patches within the receptive field. As illustrated in Figure 1, the micro neural network is instantiated with a multilayer perceptron (MLP), which is a potent function approximator. The feature maps

are obtained by sliding the micro networks over the input in a similar way as CNN [1], [2]. The MLP-conv maps the input local patch to the output feature vector with a MLP, which consists of multiple fully connected layers with nonlinear activation functions. The shared NIN is implemented by stacking three above MLP-conv layers. For different loss functions, we use one shared NIN in the structure of ML-DNN. The reason is that the MLP-conv layers form the general hierarchical feature representation for image data, which should be shared by all loss functions.

• **Multi-loss Regularization:** The multi-loss functions correspond to multiple loss branches of our proposed ML-DNN framework. Multi-loss functions may be able to prevent the algorithm away from overfitting to one single loss function, which may occur with only one loss function used. As shown in Figure 1, multiple loss branches are embedded into the shared NIN. For each loss branch, one fully-connected layer and one output layer are included. The fully-connected layer is followed by the ReLu Nonlinearity and 0.5 dropout. We optimize each output layer with different loss functions. During the network training, we feed the convolutional feature maps of the shared NIN into multiple branches of the loss functions. Based on these multiple loss functions, the network is trained by simultaneously optimizing all loss functions with back-propagation. The last convolution layer in the shared NIN is connected to multiple fully-connected layers. The parameter gradients from different loss functions can thus be conveniently used to optimize the parameters of the shared NIN. Different kind of loss function is based on different theoretical motivations. For example, the softmax loss, pairwise ranking loss and LambdaRank top-1 loss belong to the loss of top-1, while the LambdaRank top-2 loss is a kind of top-2 loss. Thus different loss functions have a certain complementarity and the gradients brought by them help iteratively learn parameters from different aspects. In this way, the whole ML-DNN is able to consider all these multiple loss functions simultaneously and avoid the over-fitting problem during training. It is worth noting that our framework can be instantiated by any CNN architectures and any loss functions which may be designed for different application scenarios in computer vision. In this paper, we used several most widely used loss functions for the visual recognition.

B. Loss Functions

In this section, we introduce several widely used loss functions in our proposed multi-loss regularization with more details. We use the softmax loss, pairwise ranking loss and LambdaRank loss, which enjoy quite different regularization properties, to regularize the training of DNN. It is expected that using these loss functions will effectively alleviate the overfitting issue in the previous DNNs.

For clarity, some notations used in the following are introduced here. We use c_i^+ and c_i^- to denote the positive and negative labels of the image X_i , $i = 1, \dots, N$, where N is the total number of training images, and $|c_i^+|$ is the number of positive labels. We use $q_j(X_i)$ to specify the discrete

probability distribution of the image X_i in the j -th class, $j = 1, \dots, C$, and C is the number of possible classes.

1) *Softmax Loss*: Most of previous DNNs [5], [1], [4] employ the softmax loss function for optimizing DNN parameters. With the softmax function, the normalized probability of the image X_i in j -th class can be computed by

$$p_{i,j} = \frac{\exp(q_j(X_i))}{\sum_{j=1}^C \exp(q_j(X_i))}. \quad (1)$$

To minimize the KL-Divergence between the predictions and the ground-truth probabilities, the softmax cost function is defined as:

$$J_{\text{softmax}} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C \bar{p}_{i,j} \log(p_{i,j}), \quad (2)$$

where $\bar{p}_{i,j}$ denotes the ground truth probability between the image X_i in class j . Each image X_i corresponds to a label vector $y_{i,j} \in \mathbb{R}^C$, with $y_{i,j} = 1$ and $\bar{p}_{i,j} = 1$ indicating the presence of a label for the image X_i , $y_{i,j} = 0$ and $\bar{p}_{i,j} = 0$ indicating the absence of a label for the image X_i .

2) *Pairwise Ranking Loss*: The second loss function we consider is the pairwise ranking loss [18], [27]. The pairwise ranking method can transform the classification problem into a task of classifying the preference order of label pairs: given any two labels, it will decide which label should be ranked first. The pairwise ranking loss method minimizes the classification error on the pairs of labels. Specifically, the goal of the pairwise ranking method is to rank all the labels so that positive labels always have higher prediction scores than negative labels. The pairwise ranking loss J_{pairwise} is defined as

$$J_{\text{pairwise}} = \frac{1}{N} \sum_{i=1}^N \sum_{m=1}^{c_i^+} \sum_{n=1}^{c_i^-} \max(0, 1 - q_m(X_i) + q_n(X_i)), \quad (3)$$

where $m = 1, \dots, c_i^+$ and $n = 1, \dots, c_i^-$ represent indexes of the positive labels and negative labels, respectively. By minimizing J_{pairwise} , we compute the sub-gradient of this loss function during the training of ML-DNN.

3) *LambdaRank Loss*: In order to directly optimize the top- k classification accuracy, the third loss function, LambdaRank loss, is introduced to learn ML-DNN. The LambdaRank method [28] obtains the desired gradients directly, rather than deriving them from a cost loss function. In this way, it can bypass the difficulties brought by the ranking loss function. The key observation of LambdaRank is that in order to train a DNN model, we do not need the losses themselves; we only need the gradients of the losses with respect to the DNN model outputs.

Before discussing the LambdaRank loss function, we first introduce some information about RankNet [19]. For a given image X_i , the set of the label pairs is denoted as S , $(m, n) \in S$ if and only if $m \in c_i^+$ and $n \in c_i^-$. The two outputs $q_m(X_i)$ and $q_n(X_i)$ are mapped to a learned probability that the m -th label is ranked higher than the n -th label by a sigmoid function, thus

$$p_{m,n} = \frac{1}{1 + \exp(-\gamma(q_m(X_i) - q_n(X_i)))}, \quad (4)$$

where the parameter γ is directly related to scaling. We adopt the cross entropy loss function, which penalizes the deviation of the model output probabilities from the desired probabilities. Then the cost is

$$J_{mn} = \log(1 + \exp(-\gamma(q_m(X_i) - q_n(X_i)))). \quad (5)$$

The gradient of the above loss function is

$$\begin{aligned} \frac{\partial J_{mn}}{\partial q_m(X_i)} &= \frac{-\gamma}{1 + \exp(-\gamma(q_m(X_i) - q_n(X_i)))}, \\ \frac{\partial J_{mn}}{\partial q_n(X_i)} &= -\frac{\partial J_{mn}}{\partial q_m(X_i)}. \end{aligned} \quad (6)$$

From the RankNet to LambdaRank loss, the gradient is multiplied by η_{mn} , which calculates the changes in evaluation measures by swapping the rank position of the m -th label and the n -th label.

The model can directly optimize the evaluation measures [20], and the gradient of $q_m(X_i)$ is

$$\frac{\partial J_{mn}}{\partial q_m(X_i)} = \sum_{(m,n) \in S} \eta_{mn} \frac{-\gamma}{1 + \exp(\gamma(q_m(X_i) - q_n(X_i)))}. \quad (7)$$

Similarly,

$$\frac{\partial J_{mn}}{\partial q_n(X_i)} = \sum_{(m,n) \in S} \eta_{mn} \frac{\gamma}{1 + \exp(\gamma(q_m(X_i) - q_n(X_i)))}. \quad (8)$$

In order to optimize the loss function with top- k evaluation measures, $\eta_{mn}(k)$ is defined as the order change in top- k . It is easily observed that $\eta_{mn}(k) = \frac{1}{|c_i^+|}$, if m is in the top- k prediction ($m \leq k$), and n is not in the top- k prediction ($n > k$). Otherwise, $\eta(m, n)$ is zero. In this work, we only consider the top-1 and top-2 evaluation measures, then name them as LambdaRank top-1 loss function and LambdaRank top-2 loss function, respectively.

C. Parameter Learning for ML-DNN

The DNN with a large number of parameters has a powerful machine learning capability, but it is difficult to prevent it from overfitting. To address this problem, some regularization techniques are proposed for training DNNs, such as data augmentation [10], [11], weight decay [12], dropout [13], [14], drop-connected [15], stochastic pooling [16], etc. These approaches may suffer from the certain limitation for a particular loss function. The optimization of loss functions has not been well explored to prevent the overfitting problem, which is truly critical for learning parameters.

In order to alleviate overfitting, we propose to train an ML-DNN model with the multiple loss functions in this work. For learning ML-DNN, we adopt four kinds of loss functions, including Softmax loss, Pairwise loss, LambdaRank top-1 loss and LambdaRank top-2 loss. The motivation of this ML-DNN framework is that other loss functions, e.g., pairwise ranking loss and LambdaRank loss, may have the potential to prevent the algorithm overfitting to one softmax loss function. Therefore, the key idea of our ML-DNN learning framework is to use the cross-loss-function regularization for boosting the generalization capability of the ML-DNN model. Based on the

above loss functions, we propose the multi-loss regularization method for learning an effective ML-DNN model.

Firstly, we warm-start the shared NIN by initializing it with the pre-trained parameters by [4]. In order to improve the generalization capability of DNNs, the fine-tuning is then performed to adjust the parameters of ML-DNN with the cross-loss-function regularization method. We learn an ML-DNN model by simultaneously optimizing multiple loss functions in the training stage. The underlying reason for warm-start is that our offline experiments show that when the step-size is large in the early stages, the gradients from different objects are quite diverse, which makes the optimization quite slow in convergence. Warm-start can well avoid this issue. We then use the back-propagation technique [2] to update the parameters of the ML-DNN model. Given some training images, the predicted probabilities from the multi-loss function are calculated, and the corresponding gradients are obtained by back-propagation through the network. For the shared NIN, the gradients from the fully-connected layers are summed together for parameter updating. In each iteration, the neurons in dropout layers will be stochastically selected with probability 0.5 to forward their activation to the output units, and only the selected neurons will participate in the back-propagation during this iteration. Similarly, we adopt all the neurons for prediction with their activation value multiplied by 0.5 for normalization.

D. Testing with Average Pooling

For each test image, we feed it into our ML-DNN framework and output several probability outputs corresponding to the different loss functions in the training. We then fuse the probabilities from different loss functions. We adopt the average pooling technique to produce the ultimate prediction of the testing image. The probability of the image X_i with the j -th label is given by

$$P_j(X_i) = \frac{\sum_{l=1}^L P_{l,j}(X_i)}{L}, \quad (9)$$

where L is the total number of loss functions, and $P_{l,j}(X_i)$ denotes the probability of the image X_i from the l -th loss function.

Our testing process with the ML-DNN model is different from simply combining the predictions of several DNNs. We can analyze their differences from three characteristics. 1) Previous combination methods often separately train different DNNs with same/different loss functions, and their predictions are then combined with the pooling technique in the testing. It is difficult to alleviate the overfitting problem for learning DNNs. 2) The training time of training one ML-CNN is only slightly more than training a DNN with a single loss function, that is, only the additional fully-connect layer and output layers are trained instead of training all convolution layers many times for all lost functions. 3) The ML-DNN uses the cross-loss-function regularization, which can boost the generalization capability of the DNN. The diverse properties of different loss function can be combined in our unified framework. Therefore, learning the ML-DNN model is not only

faster than training multiple DNNs, but also can improve the performance of DNN by alleviating the overfitting problem.

IV. EXPERIMENTS

We evaluate the effectiveness of our proposed ML-DNN method on four standard benchmark datasets: CIFAR-10 [17], CIFAR-100 [17], MNIST [2], and SVHN [31].

A. Experimental Settings

To ensure the fairness of comparisons with previous baselines, all experiments are conducted based on the following experimental setups. Test error is used as the evaluation metric. We use four different loss functions for learning an ML-DNN model: softmax loss, pairwise ranking loss, LambdaRank top-1 and LambdaRank top-2 loss functions. For the shared NIN, we follow the same network definitions used in [4], which are publicly available¹. For each loss branch of our ML-DNN, the size of the fully-connected (fc) layer is 256 for the CIFAR-10, MNIST and SVHN datasets, and 2560 for CIFAR-100 dataset, followed by the ReLU Nonlinearity and dropout (50%) layer. For updating the ML-DNN parameters, we use the stochastic gradient descent (SGD) learning method with the mini-batch size of 128 at a fixed constant momentum value of 0.9. The weight decay is set to 0.001. The global learning rate of learning the ML-DNN model is set to 0.01. We decrease the learning rate every 100K iterations by 10. For the parameter learning of ML-DNN, the local learning rate of the shared NIN is set to 0.01. For each loss branch, the learning rate for all fully-connected layers and all loss layers is the same and set as 1.0. The maximum of iterations is set to 120,000.

B. CIFAR-10

The CIFAR-10 dataset [17] consists of 60,000 32×32 color images of 10 classes, with 6,000 images per class. There are 50,000 training images and 10,000 test images. The dataset is divided into five training batches and one test batch, each with 10,000 images. The test batch contains 1,000 randomly-selected images from each class. Following the same experiment settings used in [5], we process the data with the same global contrast normalization and Zero Components Analysis (ZCA) whitening. To fairly compare with the other previous algorithms [4], [5], we have conducted extensive experiments under two settings (i.e., without data augmentation and with data augmentation). In terms of data augmentation, we augment the images by zero padding 4 pixels on each direction, then perform corner cropping and random flipping during training.

Under above two settings, we compare our proposed ML-DNN method with previous DNN regularization techniques, including stochastic pooling [16], Maxout Networks [5], Network In Network (NIN) [4], and probabilistic maxout units [29]. Detailed classification error results are shown in Table I. The comparison results on CIFAR-10 indicate that the proposed ML-DNN method achieves 9.55% test errors without

¹<https://github.com/mavenlin/cuda-convnet/tree/master/NIN>

TABLE I
CIFAR-10 CLASSIFICATION ERRORS OF VARIOUS METHODS.

Method	Test Error (%)
Without data augmentation	
Stochastic Pooling [16]	15.13
Maxout Networks [5]	11.68
Maxout Units [29]	11.35
NIN [4]	10.41
ML-DNN	9.55
With data augmentation	
Maxout Networks [5]	9.38
Drop-Connect [15]	9.32
NIN [4]	8.81
ML-DNN	8.12

TABLE III
MNIST CLASSIFICATION ERRORS OF VARIOUS METHODS.

Method	Test Error (%)
2layer CNN+2layer NN [16]	0.53
Stochastic Pooling [16]	0.47
NIN [4]	0.47
Maxout Networks [5]	0.45
ML-DNN	0.42

data augmentation and 8.12% test error with data augmentation, both of which perform better than all the baselines. The ML-DNN improves the performance of the NIN [4] of 10.41% by more than 0.8% without data augmentation, and 8.81% by higher than 0.69% with data augmentation. The baseline methods achieve the classification error of 15.13% for the stochastic pooling method [16], 11.68% for maxout networks [5], and 11.35% for maxout Units [29], all of which are much higher than the test error of our ML-DNN when without data augmentation. Similarly, we also improve a test error by 0.69%, compared to the baselines NIN [4]. This well verifies the superiority of our ML-DNN method, which is very effective for achieving better generalization capability during the ML-DNN training.

C. CIFAR-100

This CIFAR-100 dataset [17] has the same number of images as the CIFAR-10 database, but contains 100 classes, with only one tenth as labeled examples per class. There are 500 training images and 100 testing images per class. For the CIFAR-100 dataset, we adopt the same network settings and data preprocessing (including the data augmentation technique) as the CIFAR-10, described in [4].

Table II shows the performances of our proposed ML-DNN and other state-of-the-art methods. It can be seen that the ML-DNN method gets a misclassification test error of 34.18%, which significantly improve the performance over the NIN [4] by less 1.5% errors. In particular, our method also beat all other regularization methods with a large margin, e.g., 42.51%, 38.57% and 38.14% for stochastic pooling [16], maxout networks [5] and Probabilistic Maxout Units [29], respectively. Moreover, we also have conducted an experiment with data augmentation and achieved a test error 31.47%, which is lower than NIN [4] by 2.06%. Our proposed ML-

TABLE II
CIFAR-100 CLASSIFICATION ERRORS OF VARIOUS METHODS.

Method	Test Error (%)
Without data augmentation	
Stochastic Pooling [16]	42.51
Maxout Networks [5]	38.57
Maxout Units [29]	38.14
NIN [4]	35.68
ML-DNN	34.18
With data augmentation	
NIN [4]	33.53
ML-DNN	31.47

TABLE IV
SVHN CLASSIFICATION ERRORS OF VARIOUS METHODS.

Method	Test Error (%)
Stochastic Pooling [16]	2.80
Maxout Networks [5]	2.47
NIN [4]	2.35
Multi-digit num. recog. [30]	2.16
ML-DNN	1.92

DNN method can obtain better performances over other state-of-the-arts on classifying 100 classes.

D. MNIST

The MNIST handwritten digit classification dataset [2] consists of 28×28 pixel greyscale images of hand-written digits (from 0 to 9). There are 60,000 training images and 10,000 testing images in total.

The comparison results are presented in Table III. It can be observed that our proposed method outperforms the competing methods on the MNIST dataset. It is noted that the test error of our ML-DNN is lower than all the competing methods (such as 2 layer CNN + 2 layer NN [16], stochastic pooling [16], maxout networks [5], and NIN [4]). Our ML-DNN method gets 0.42% classification error, and gets lower classification error than the baseline method NIN by 0.05%. The performance on MNIST also demonstrates the advantage of the ML-DNN method.

E. SVHN

The Street View House Numbers (SVHN) dataset [31] is composed of 73,257 images for training, 26,032 images for testing, and 531,131 extra training color images of 32×32 . The task of the dataset is to classify the digit located at the center of each image. Following the previous work [5], we select out 400 samples per class from the training set, and 200 samples per class from the extra set. The rest of the training set and the extra set are used for training. We preprocess the dataset by local contrast normalization.

As reported in Table IV, the classification error of ML-DNN significantly outperforms the other state-of-the-arts including the stochastic pooling [16], maxout networks [5], multi-digit number recognition method [30] and NIN [4]. The proposed ML-DNN method improves the performance by 0.24% and

TABLE V
CLASSIFICATION ERROR COMPARISONS BETWEEN SINGLE-LOSS DNNs AND ML-DNN.

Method	Test Error (%)			
	Without data augmentation		With data augmentation	
	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-100
NIN + Softmax loss	10.41	35.68	8.81	33.53
NIN + Pairwise ranking loss	10.57	36.76	9.27	34.47
NIN + LambdaRank top-1 loss	10.37	36.83	9.34	34.62
NIN + LambdaRank top-2 loss	9.94	36.24	9.37	33.87
Average results (four model combination)	9.58	34.78	8.34	32.35
ML-DNN	9.55	34.18	8.12	31.47

TABLE VI
CLASSIFICATION ERRORS OF EACH BRANCH OF ML-DNN AND ML-DNN.

Method	Test Error (%)			
	Without data augmentation		With data augmentation	
	CIFAR-10	CIFAR-100	CIFAR-10	CIFAR-100
ML-DNN's Softmax loss branch	9.66	35.54	8.35	32.34
ML-DNN's Pairwise ranking loss branch	9.75	36.66	8.82	33.27
ML-DNN's LambdaRank top-1 loss branch	9.84	35.62	8.37	31.93
ML-DNN's LambdaRank top-2 loss branch	9.93	35.68	8.64	31.82
ML-DNN	9.55	34.18	8.12	31.47

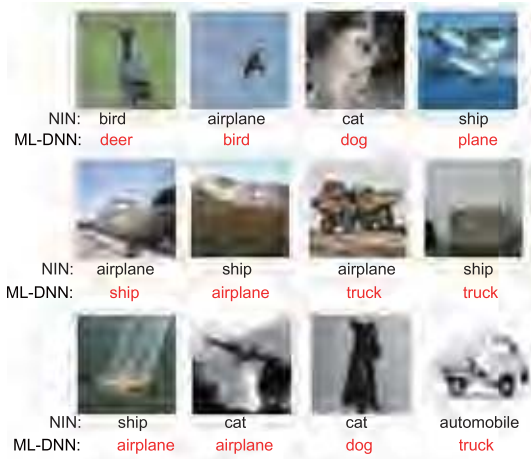


Fig. 2. Exemplar test images from CIFAR-10 [17]. These images are wrongly classified by NIN, but correctly classified by ML-DNN. There are two label names marked below each image, one of which is its wrong class with black color by NIN, the other is its right class with red color by ML-DNN.

0.43% over the state-of-the-arts [30] and [4], respectively. Compared with the existing stochastic pooling and maxout networks methods, the ML-DNN significantly improves the performance by 0.88% and 0.55%, respectively. It demonstrates that the method effectively works for boosting the generalization capability of the ML-DNN model.

F. Algorithm Analysis

As observed from Table I to Table IV, the classification errors achieved by our ML-DNN method are reported on the above four datasets, 9.55% for CIFAR-10, 34.18% for CIFAR-100, 0.42% for MNIST, and 1.92% for SVHN. It can be concluded from the tables that: 1) The ML-DNN method does improve the performance of the current DNN classification framework, and even on the difficult CIFAR-

100 dataset; 2) Our regularized technique with the cross-loss-function plays an important role in alleviating overfitting problem for learning a capable DNN model and the resulting ML-DNN can obviously improve the classification performance; 3) Compared with some existing regularization techniques (e.g. dropout, stochastic pooling and maxout network), the ML-DNN does improve the capability of the DNN model by multi-loss regularization.

To validate the regularization capability of ML-DNN, we show the comparison results of different DNN using different single loss functions, their combined average results (namely as “Average results from four model combination”), and our ML-DNN results, presented in Table V. The architecture of single loss DNNs also utilizes the NIN architecture, followed by a single loss function. The difference between single loss DNNs (e.g. “NIN + Softmax loss”, “NIN + Pairwise ranking loss”, “NIN + lambdaRank top-1 loss” and “NIN + LambdaRank top-2 loss”) and ML-DNN model is the strategy of using loss function, (e.g., single loss or multiple different loss functions). The ML-DNN method outperforms four kinds of the single loss DNNs, and even surpasses the most popular softmax loss for learning a DNN model. Moreover, our ML-DNN outperforms the average results by combining four DNNs with single loss by 0.03% and 0.6% over CIFAR10 and CIFAR-100 without data augmentation, 0.22% and 0.98% over CIFAR10 and CIFAR-100 with data augmentation. This superiority can be observed on different datasets, which can further validate the effectiveness of ML-DNN.

We furthermore analyze the results of each loss branch in the framework of our ML-DNN. Each loss branch of ML-DNN corresponds to a fully-connected layer and a loss layer. The detailed classification test errors are shown in Table VI. The result of our ML-DNN, which combines the classification results of four branches, is better than the results from each branch of ML-DNN. Compared the results of single loss DNN in Table V, each branch of ML-DNN can get a higher classifi-

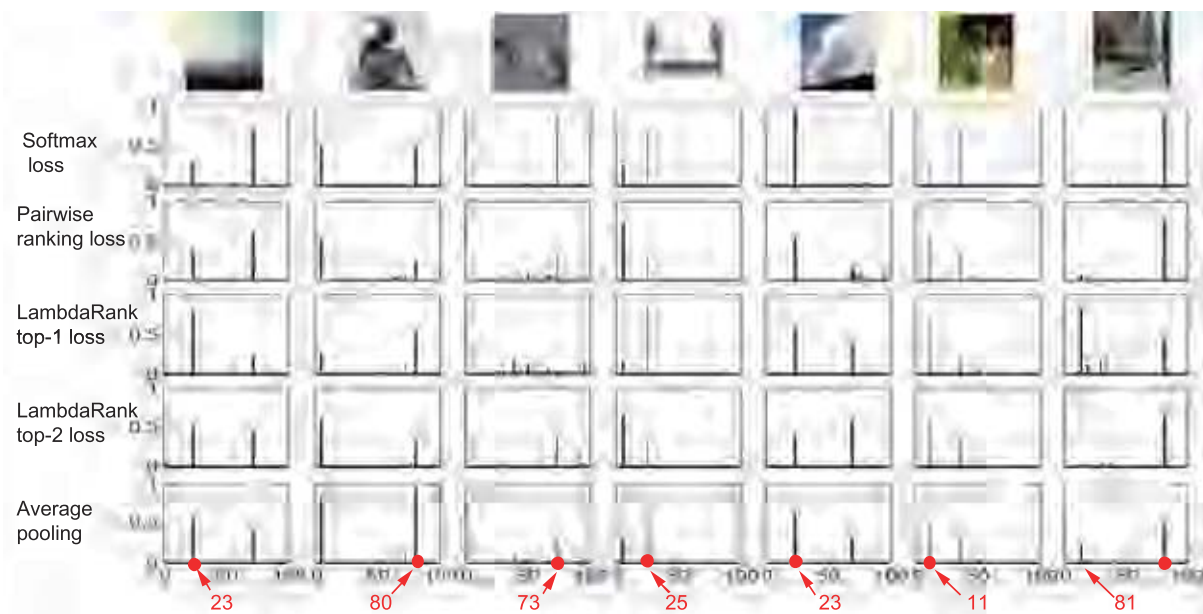


Fig. 3. Exemplar images with the predicted probabilities on CIFAR-100 test dataset [17] using the ML-DNN. Some testing images and 100 predicted probabilities for all labels are shown. The predictions from the second to fifth rows are from four different loss functions. The ultimate results on the last row are produced by fusing the predictions from four different loss functions with the average pooling technique. And the red points in the last row are marked for the ground truth label of the corresponding test images.

cation result. For example, the softmax loss branch of our ML-DNN outperforms the single loss DNN with softmax loss, e.g., 10.41% vs 9.66% on CIFAR-10 without data augmentation, 33.53% vs 32.34% on CIFAR-100 with data augmentation. It demonstrates that compared with the single loss DNN, the ML-DNN method can boost the performance for the image classification task on four benchmark datasets.

Some test images from the CIFAR-10 dataset [17] are shown in Figure 2. These images are wrongly classified by NIN, but correctly classified by ML-DNN. There are two label names marked below each image, one of which is its wrong class with black font by NIN, the other is its right label (i.e. the ground-truth label) by ML-DNN. For example, the bird and airplane have the blue sky, two wings and similar color. The automobile and truck are also similar in appearance (such as wheels, seat, color, etc.). Our proposed ML-DNN learning employs cross-loss-function regularization, while the NIN learns the parameters with only a softmax loss function. The effectiveness of our proposed ML-DNN again speaks well that our method can successfully recognize the confusing and difficult objects.

As shown in Figure 3, we further report some predicted probabilities from four different loss functions and the fused results using the ML-DNN model. The output layers predict different label probabilities, due to the intrinsic properties of different loss functions. For example, for the image of the first column, the predicted label by the softmax loss and pairwise ranking loss are different from that of the LambdaRank top-1 loss and the LambdaRank top-1 loss functions, while the ultimate predicted label is the same with that of the LambdaRank top-1 loss and the LambdaRank top-1 loss functions. By employing the average pooling technique to fuse all the probabilities from all loss layers, we can achieve

the best results by balancing their probabilities. The average pooling results can overall be considered the outputs by ML-DNN from four different loss functions and improve the generalization capability of the ML-DNN.

To further analyze the complementary among multiple loss functions, we show the gradient values brought by four different loss functions in our ML-DNN framework. As can be seen in Fig. 4, four different curves represent the gradient values brought by four different loss branches, such as softmax loss (blue curve), pairwise ranking loss (green curve), LambdaRank top-1 loss (magenta curve) and LambdaRank top-2 loss (light blue curve). The light blue curve, which presents the gradient values brought from the LambdaRank top-2 loss branch, is different from other three curves of top-1 loss branches (e.g., softmax loss, pairwise ranking loss and LambdaRank top-1 loss). We can thus say that the top-2 loss function is complementary to these top-1 loss functions to some extent. Moreover, it can be observed that gradient values (blue, green and magenta curves) brought by three top-1 loss functions also have a certain degree of difference. Therefore, multiple loss functions of our ML-DNN can help constrain the parameters of a neural network from different aspects. From the analyses above and the practical numerical experiments conducted in the paper, we can conclude that the multi-loss can actually guide learning a better network and improve the classification performance of it.

To evaluate the performance of our ML-DNN on the compelling and significantly harder vision task, we perform an experiment on ImageNet dataset [32]. For the ImageNet dataset, we adopt the same setting with the work [33], which extracted random 224×224 patches from 256×256 training images, and then trained a DNN on these patches. In the training process, we also crop 224×224 patches from the middle

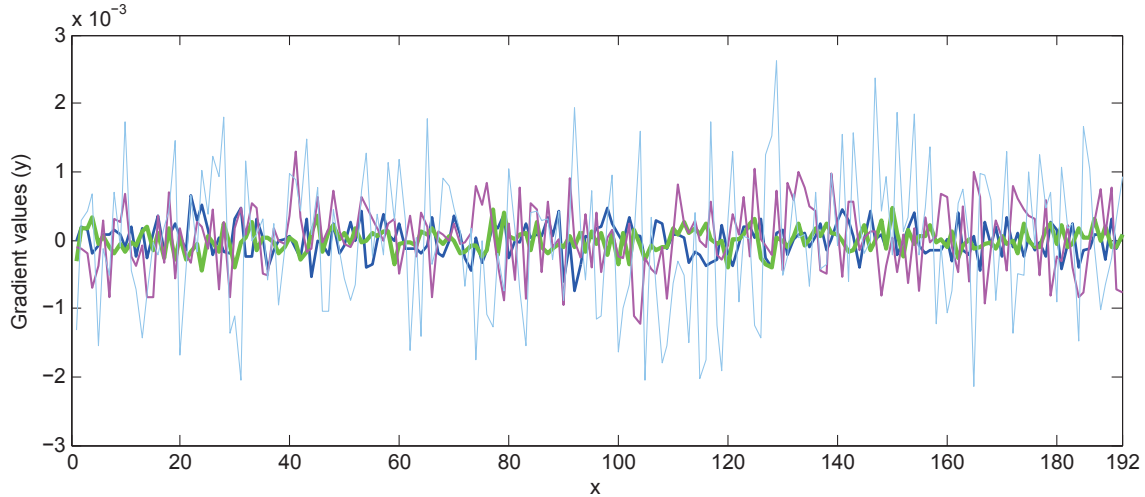


Fig. 4. Gradient values brought by four loss branches on CIFAR-100 dataset. In order to conveniently present the gradient values (y) brought by four loss branches, we just show some gradient values (i.e., the first gradient values of the upper left corner of each feature maps over all channels) in the 100th iteration of the optimization process. The horizontal axis (x=192) is the number of channels. Four different curves (blue, green, magenta and light blue curves) denote the gradient values brought by four different loss branches, e.g., softmax loss, pairwise ranking loss, LambdaRank top-1 loss and LambdaRank top-2 loss, respectively.

TABLE VII
CLASSIFICATION ERRORS OF VARIOUS METHODS ON IMAGENET DATASET.

Method	Test Error (%)	
	Top-1	Top-5
Sparse coding [34]	47.1	28.2
SIFT+FVs [35]	45.7	25.7
Alex-net [1]	37.5	17.0
VGG(deep16) [33]	28.97	10.15
VGG(deep16)+ML-DNN	28.46	9.8

of test images. As described in section III, the structure of the ML-DNN framework is composed of the shared NIN and multi-loss branches, while our ML-DNN on this compelling problem also employs the same multi-loss branches, but adopts a deeper DNN structure, i.e., VGG(deep16) network [33]. Followed the work in [33], we also evaluate the performance from two aspects, top-1 test error and the top-5 test error. Detailed test error of various methods on ImageNet dataset can be seen in Table VII. Our ML-DNN slightly outperforms the VGG(deep16) network with single loss: 28.46% vs 28.97% for the top-1 test error, 9.8% vs 10.15% for the top-5 test error. The main reason for these improvements on this challenging task may be that the ML-DNN can well boost its discriminative capability by considering cross-loss-function regularization, even with the deeper network architecture.

V. CONCLUSION AND FUTURE WORK

In this paper, we proposed a general multi-loss regularized DNN framework for alleviating the overfitting issue of DNN. To boost the generalization capability of DNN, this general scheme allows us to learn an ML-DNN model by simultaneously optimizing multiple loss functions. For the image classification task, we studied the loss functions, pairwise loss and LambdaRank top- k loss, for learning the ML-DNN model. With the average pooling technique, the final prediction

can be simply computed from the outputs of the ML-DNN model from different loss functions. Extensive experimental results on CIFAR-10, CIFAR-100, MNIST and SVHN datasets clearly demonstrate that the proposed ML-DNN framework achieved the state-of-the-art performances. In the future, we plan to further explore the performance of ML-DNN with other vision task, e.g. object detection, image retrieval, and image annotation, etc.

ACKNOWLEDGMENT

This work was supported by the National Natural Science Foundation of China (Grant 61572214 and U1233119). This research was partly supported under Australian Research Council Discovery Projects funding scheme (project DP140102270).

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105. 1, 2, 3, 4, 9
- [2] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," vol. 86, no. 11, 1998, pp. 2278–2324. 1, 2, 3, 5, 6
- [3] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *International Conference on Machine Learning*, 2009, pp. 609–616. 1
- [4] M. Lin, Q. Cheng, and S. Yan, "Network in network," *International Conference on Learning Representations*, 2014. 1, 2, 3, 4, 5, 6, 7
- [5] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," *arXiv preprint arXiv:1302.4389*, 2013. 1, 3, 4, 5, 6
- [6] Y. Wei, W. Xia, J. Huang, B. Ni, J. Dong, Y. Zhao, and S. Yan, "CNN: Single-label to multi-label," *arXiv preprint arXiv:1406.5726*, 2014. 1
- [7] X. Zeng, W. Ouyang, M. Wang, and X. Wang, "Deep learning of scene-specific classifier for pedestrian detection," in *European Conference on Computer Vision*, 2014, pp. 472–487. 1
- [8] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1915–1929, 2013. 1

- [9] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," *arXiv preprint arXiv:1409.4842*, 2014. 1, 3
- [10] J. Schmidhuber, "Multi-column deep neural networks for image classification," in *Computer Vision and Pattern Recognition*, 2012, pp. 3642–3649. 1, 2, 4
- [11] P. Y. Simard, D. Steinkraus, and J. C. Platt, "Best practices for convolutional neural networks applied to visual document analysis," in *International Conference on Document Analysis and Recognition*, vol. 2, 2003, pp. 958–958. 1, 2, 4
- [12] C. M. Bishop, *Neural networks for pattern recognition*. Clarendon press Oxford, 1995. 1, 2, 4
- [13] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. R. Salakhutdinov, "Improving neural networks by preventing co-adaptation of feature detectors," *arXiv preprint arXiv:1207.0580*, 2012. 1, 2, 4
- [14] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014. 1, 4
- [15] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, and R. Fergus, "Regularization of neural networks using dropconnect," in *International Conference on Machine Learning*, 2013, pp. 1058–1066. 1, 2, 4, 6
- [16] M. D. Zeiler and R. Fergus, "Stochastic pooling for regularization of deep convolutional neural networks," *arXiv preprint arXiv:1301.3557*, 2013. 1, 2, 4, 5, 6
- [17] A. Krizhevsky and G. Hinton, "Learning multiple layers of features from tiny images," *Computer Science Department, University of Toronto, Tech. Rep*, 2009. 2, 5, 6, 7, 8
- [18] T. Joachims, "Optimizing search engines using clickthrough data," in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2002, pp. 133–142. 1, 4
- [19] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, "Learning to rank using gradient descent," in *international conference on Machine learning*, 2005, pp. 89–96. 1, 4
- [20] P. Donmez, K. M. Svore, and C. J. Burges, "On the local optimality of lambdarank," in *ACM SIGIR conference on Research and development in information retrieval*, 2009, pp. 460–467. 1, 4
- [21] S. Mohamed, S. Member, and G. Rubino, "A study of real-time packet video quality using random neural networks," *IEEE Transactions On Circuits and Systems for Video Technology*, vol. 12, pp. 1071–1083, 2002. 2
- [22] H. Choi and C. Lee, "Motion adaptive deinterlacing with modular neural networks," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 21, no. 6, pp. 844–849, June 2011. 2
- [23] K. Kim, S. Lee, J.-Y. Kim, M. Kim, and H.-J. Yoo, "A configurable heterogeneous multicore architecture with cellular neural network for real-time object recognition," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 19, no. 11, pp. 1612–1622, Nov 2009. 2
- [24] N. Sudha, A. Mohan, and P. Meher, "A self-configurable systolic architecture for face recognition system based on principal component neural network," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 21, no. 8, pp. 1071–1084, Aug 2011. 2
- [25] G. Gkioxari, B. Hariharan, R. Girshick, and J. Malik, "R-CNNs for pose estimation and action detection," *arXiv preprint arXiv:1406.5212*, 2014. 3
- [26] S. Li, Z.-Q. Liu, and A. B. Chan, "Heterogeneous multi-task learning for human pose estimation with deep convolutional neural network," pp. 488–495, 2014. 3
- [27] J. Li, Y. Tian, T. Huang, and W. Gao, "Multi-task rank learning for visual saliency estimation," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 21, no. 5, pp. 623–636, 2011. 4
- [28] C. J. Burges, "From RankNet to LambdaRank to LambdaMart: An overview," *Learning*, vol. 11, pp. 23–581, 2010. 4
- [29] J. T. Springenberg and M. Riedmiller, "Improving deep neural networks with probabilistic maxout units," *arXiv preprint arXiv:1312.6116*, 2013. 5, 6
- [30] I. J. Goodfellow, Y. Bulatov, J. Ibarz, S. Arnoud, and V. Shet, "Multi-digit number recognition from street view imagery using deep convolutional neural networks," *arXiv preprint arXiv:1312.6082*, 2013. 6, 7
- [31] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," in *NIPS workshop on deep learning and unsupervised feature learning*, vol. 2011, 2011, p. 4. 5, 6
- [32] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *CVPR*, 2009, pp. 248–255. 8

- [33] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *CoRR*, vol. abs/1409.1556, 2014. 8, 9
- [34] A. Berg, J. Deng, and L. Fei-Fei, "large-scale visual recognition challenge," in *www.image-net.org/challenges*, 2010. 9
- [35] J. Sanchez and F. Perronnin, "High-dimensional signature compression for large-scale image classification," in *Computer Vision and Pattern Recognition (CVPR)*, 2011 *IEEE Conference on*, June 2011, pp. 1665–1672. 9



Chunyan Xu received the B.Sc. degree from Shandong Normal University in 2007 and the M.Sc. degree from Huazhong Normal University in 2010 and the PhD degree in the School of Computer Science and Technology, Huazhong University of Science and Technology in 2015. She was a visiting scholar at National University of Singapore from 2013 to 2015. Currently, she is a lecture in the School of Computer Science and Technology from Nanjing University of Science and Technology, Nanjing, China. Her research interests include deep

neural network, computer vision, manifold learning and kernel methods.



Canyi Lu received the bachelor of mathematics from Fuzhou University in 2009, and the master degree in the pattern recognition and intelligent system in 2012. From August 2013, he was a phd student with the Department of Electrical and Computer Engineering at National University of Singapore. His research interests include computer vision and machine learning. His homepage is <http://sites.google.com/site/canyilu>.



Xiaodan Liang is a Ph.D. student from School of Information Science and Technology, Sun Yat-sen University, China. She is currently working at National University of Singapore as a Research Intern. Her research interests mainly include semantic segmentation, object/action recognition and medical image analysis.



Junbin Gao graduated from Huazhong University of Science and Technology (HUST), China in 1982 with BSc. degree in Computational Mathematics and obtained PhD from Dalian University of Technology, China in 1991. He is a Professor in Computing Science in the School of Computing and Mathematics at Charles Sturt University, Australia. He was a senior lecturer, a lecturer in Computer Science from 2001 to 2005 at University of New England, Australia. From 1982 to 2001 he was an associate lecturer, lecturer, associate professor and professor

in Department of Mathematics at HUST. His main research interests include machine learning, data mining, Bayesian learning and inference, and image analysis.



Wei Zheng received the Bachelor's degree from Tsinghua University, Beijing, China, in 2006. He got the Ph.D. degree at Institute of Computing Technology, Chinese Academy of Sciences, Beijing, in 2013. Currently, he is a researcher at Beijing Samsung Telecom R&D Center, Beijing. His research interests include image categorization, object detection, and scene analysis.



Tianjiang Wang received the B. Sc. degree in computational mathematics in 1982 and the PhD degree in computer science in 1999 from Huazhong University of Science and Technology (HUST), Wuhan, China. He is currently a Professor with the School of Computer Science, Huazhong University of Science and Technology, Wuhan, China. He has finished some related projects and is the author of more than 20 related papers. His research interests include machine learning, computer vision, and data mining.



Shuicheng Yan is currently an Associate Professor at the Department of Electrical and Computer Engineering at National University of Singapore, and the founding lead of the Learning and Vision Research Group (<http://www.lv-nus.org>). Dr. Yan's research areas include machine learning, computer vision and multimedia, and he has authored/co-authored hundreds of technical papers over a wide range of research topics, with Google Scholar citation > 15,000 times and H-index 51. He is ISI Highly-cited Researcher, 2014 and IAPR Fellow 2014. He has been serving as an associate editor of IEEE TKDE, TCSVT and ACM Transactions on Intelligent Systems and Technology (ACM TIST). He received the Best Paper Awards from ACM MM'13 (Best Paper and Best Student Paper), ACM MM12 (Best Demo), PCM'11, ACM MM10, ICME10 and ICIMCS'09, the runner-up prize of ILSVRC'13, the winner prize of ILSVRC14 detection task, the winner prizes of the classification task in PASCAL VOC 2010-2012, the winner prize of the segmentation task in PASCAL VOC 2012, the honourable mention prize of the detection task in PASCAL VOC'10, 2010 TCSVT Best Associate Editor (BAE) Award, 2010 Young Faculty Research Award, 2011 Singapore Young Scientist Award, and 2012 NUS Young Researcher Award.