
Deep Learning using Linear Support Vector Machines

Yichuan Tang

TANG@CS.TORONTO.EDU

Department of Computer Science, University of Toronto. Toronto, Ontario, Canada.

Abstract

Recently, fully-connected and convolutional neural networks have been trained to achieve state-of-the-art performance on a wide variety of tasks such as speech recognition, image classification, natural language processing, and bioinformatics. For classification tasks, most of these “deep learning” models employ the softmax activation function for prediction and minimize cross-entropy loss. In this paper, we demonstrate a small but consistent advantage of replacing the softmax layer with a linear support vector machine. Learning minimizes a margin-based loss instead of the cross-entropy loss. While there have been various combinations of neural nets and SVMs in prior art, our results using L2-SVMs show that by simply replacing softmax with linear SVMs gives significant gains on popular deep learning datasets MNIST, CIFAR-10, and the ICML 2013 Representation Learning Workshop’s face expression recognition challenge.

1. Introduction

Deep learning using neural networks have claimed state-of-the-art performances in a wide range of tasks. These include (but not limited to) speech (Mohamed et al., 2009; Dahl et al., 2010) and vision (Jarrett et al., 2009; Ciresan et al., 2011; Rifai et al., 2011a; Krizhevsky et al., 2012). All of the above mentioned papers use the softmax activation function (also known as multinomial logistic regression) for classification.

Support vector machine is an widely used alternative to softmax for classification (Boser et al., 1992). Using SVMs (especially linear) in combination with convolutional nets have been proposed in the past as part of a

multistage process. In particular, a deep convolutional net is first trained using supervised/unsupervised objectives to learn good invariant hidden latent representations. The corresponding hidden variables of data samples are then treated as input and fed into linear (or kernel) SVMs (Huang & LeCun, 2006; Lee et al., 2009; Quoc et al., 2010; Coates et al., 2011). This technique usually improves performance but the drawback is that lower level features are not been fine-tuned w.r.t. the SVM’s objective.

Other papers have also proposed similar models but with joint training of weights at lower layers using both standard neural nets as well as convolutional neural nets (Zhong & Ghosh, 2000; Collobert & Bengio, 2004; Nagi et al., 2012). In other related works, Weston et al. (2008) proposed a semi-supervised embedding algorithm for deep learning where the hinge loss is combined with the “contrastive loss” from siamese networks (Hadsell et al., 2006). Lower layer weights are learned using stochastic gradient descent. Vinyals et al. (2012) learns a recursive representation using linear SVMs at every layer, but without joint fine-tuning of the hidden representation.

In this paper, we show that for some deep architectures, a linear SVM top layer instead of a softmax is beneficial. We optimize the primal problem of the SVM and the gradients can be backpropagated to learn lower level features. Our models are essentially same as the ones proposed in (Zhong & Ghosh, 2000; Nagi et al., 2012), with the minor novelty of using the loss from the L2-SVM instead of the standard hinge loss.

Compared to nets using a top layer softmax, we demonstrate superior performance on MNIST, CIFAR-10, and on a recent Kaggle competition on recognizing face expressions. Optimization is done using stochastic gradient descent on small minibatches. Comparing the two models in Sec. 3.4, we believe the performance gain is largely due to the superior regularization effects of the SVM loss function, rather than an advantage from better parameter optimization.

2. The model

2.1. Softmax

For classification problems using deep learning techniques, it is standard to use the softmax or 1-of-K encoding at the top. For example, given 10 possible classes, the softmax layer has 10 nodes denoted by p_i , where $i = 1, \dots, 10$. p_i specifies a discrete probability distribution, therefore, $\sum_i^{10} p_i = 1$.

Let \mathbf{h} be the activation of the penultimate layer nodes, \mathbf{W} is the weight connecting the penultimate layer to the softmax layer, the total input into a softmax layer, given by \mathbf{a} , is

$$a_i = \sum_k h_k W_{ki}, \quad (1)$$

then we have

$$p_i = \frac{\exp(a_i)}{\sum_j^{10} \exp(a_j)} \quad (2)$$

The predicted class \hat{i} would be

$$\begin{aligned} \hat{i} &= \arg \max_i p_i \\ &= \arg \max_i a_i \end{aligned} \quad (3)$$

2.2. Support Vector Machines

Linear support vector machines (SVM) is originally formulated for binary classification. Given training data and its corresponding labels (\mathbf{x}_n, y_n) , $n = 1, \dots, N$, $\mathbf{x}_n \in \mathbb{R}^D$, $t_n \in \{-1, +1\}$, SVMs learning consists of the following constrained optimization:

$$\begin{aligned} \min_{\mathbf{w}, \xi_n} \quad & \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \xi_n \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{x}_n t_n \geq 1 - \xi_n \quad \forall n \\ & \xi_n \geq 0 \quad \forall n \end{aligned} \quad (4)$$

ξ_n are slack variables which penalizes data points which violate the margin requirements. Note that we can include the bias by augment all data vectors \mathbf{x}_n with a scalar value of 1. The corresponding unconstrained optimization problem is the following:

$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \max(1 - \mathbf{w}^\top \mathbf{x}_n t_n, 0) \quad (5)$$

The objective of Eq. 5 is known as the primal form problem of L1-SVM, with the standard hinge loss. Since L1-SVM is not differentiable, a popular variation

is known as the L2-SVM which minimizes the squared hinge loss:

$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \max(1 - \mathbf{w}^\top \mathbf{x}_n t_n, 0)^2 \quad (6)$$

L2-SVM is differentiable and imposes a bigger (quadratic vs. linear) loss for points which violate the margin. To predict the class label of a test data \mathbf{x} :

$$\arg \max_t (\mathbf{w}^\top \mathbf{x}) t \quad (7)$$

For Kernel SVMs, optimization must be performed in the dual. However, scalability is a problem with Kernel SVMs, and in this paper we will be only using linear SVMs with standard deep learning models.

2.3. Multiclass SVMs

The simplest way to extend SVMs for multiclass problems is using the so-called *one-vs-rest* approach (Vapnik, 1995). For K class problems, K linear SVMs will be trained independently, where the data from the other classes form the negative cases. Hsu & Lin (2002) discusses other alternative multiclass SVM approaches, but we leave those to future work.

Denoting the output of the k -th SVM as

$$a_k(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \quad (8)$$

The predicted class is

$$\arg \max_k a_k(\mathbf{x}) \quad (9)$$

Note that prediction using SVMs is exactly the same as using a softmax Eq. 3. The only difference between softmax and multiclass SVMs is in their objectives parametrized by all of the weight matrices \mathbf{W} . Softmax layer minimizes cross-entropy or maximizes the log-likelihood, while SVMs simply try to find the maximum margin between data points of different classes.

2.4. Deep Learning with Support Vector Machines

Most deep learning methods for classification using fully connected layers and convolutional layers have used softmax layer objective to learn the lower level parameters. There are exceptions, notably in papers by (Zhong & Ghosh, 2000; Collobert & Bengio, 2004; Nagi et al., 2012), supervised embedding with nonlinear NCA (Salakhutdinov & Hinton, 2007), and semi-supervised deep embedding (Weston et al., 2008). In this paper, we use L2-SVM's objective to train deep

neural nets for classification. Lower layer weights are learned by backpropagating the gradients from the top layer linear SVM. To do this, we need to differentiate the SVM objective with respect to the activation of the penultimate layer. Let the objective in Eq. 5 be $l(\mathbf{w})$, and the input \mathbf{x} is replaced with the penultimate activation \mathbf{h} ,

$$\frac{\partial l(\mathbf{w})}{\partial \mathbf{h}_n} = -C t_n \mathbf{w} (\mathbb{I}\{1 > \mathbf{w}^\top \mathbf{h}_n t_n\}) \quad (10)$$

Where $\mathbb{I}\{\cdot\}$ is the indicator function. Likewise, for the L2-SVM, we have

$$\frac{\partial l(\mathbf{w})}{\partial \mathbf{h}_n} = -2C t_n \mathbf{w} (\max(1 - \mathbf{w}^\top \mathbf{h}_n t_n, 0)) \quad (11)$$

From this point on, backpropagation algorithm is exactly the same as the standard softmax-based deep learning networks. We found L2-SVM to be slightly better than L1-SVM most of the time and will use the L2-SVM in the experiments section.

3. Experiments

3.1. Facial Expression Recognition

This competition/challenge was hosted by the ICML 2013 workshop on representation learning, organized by the LISA at University of Montreal. The contest itself was hosted on Kaggle with over 120 competing teams during the initial developmental period.

The data consist of 28,709 48x48 images of faces under 7 different types of expression. See Fig 1 for examples and their corresponding expression category. The validation and test sets consist of 3,589 images and this is a classification task.

WINNING SOLUTION

We submitted the winning solution with a public validation score of 69.4% and corresponding private test score of 71.2%. Our private test score is almost 2% higher than the 2nd place team. Due to label noise and other factors such as corrupted data, human performance is roughly estimated to be between 65% and 68%¹.

Our submission consists of using a simple Convolutional Neural Network with linear one-vs-all SVM at the top. Stochastic gradient descent with momentum is used for training and several models are averaged to slightly improve the generalization capabilities. Data preprocessing consisted of first subtracting the mean

¹Personal communication from the competition organizers: <http://bit.ly/13Zr6Gs>



Figure 1. Training data. Each column consists of faces of the same expression: starting from the leftmost column: Angry, Disgust, Fear, Happy, Sad, Surprise, Neutral.

value of each image and then setting the image norm to be 100. Each pixels is then standardized by removing its mean and dividing its value by the standard deviation of that pixel, across all training images.

Our implementation is in C++ and CUDA, with ports to Matlab using MEX files. Our convolution routines used fast CUDA kernels written by Alex Krizhevsky². The exact model parameters and code is provided on by the author at <https://code.google.com/p/deep-learning-faces>.

3.1.1. SOFTMAX VS. DLSVM

We compared performances of softmax with the deep learning using L2-SVMs (DLSVM). Both models are tested using an 8 split/fold cross validation, with a image mirroring layer, similarity transformation layer, two convolutional filtering+pooling stages, followed by a fully connected layer with 3072 hidden penultimate hidden units. The hidden layers are all of the rectified linear type. other hyperparameters such as weight decay are selected using cross validation.

	Softmax	DLSVM L2
Training cross validation	67.6%	68.9%
Public leaderboard	69.3%	69.4%
Private leaderboard	70.1%	71.2%

Table 1. Comparisons of the models in terms of % accuracy. Training c.v. is the average cross validation accuracy over 8 splits. Public leaderboard is the held-out validation set scored via Kaggle’s public leaderboard. Private leaderboard is the final private leaderboard score used to determine the competition’s winners.

²<http://code.google.com/p/cuda-convnet>

We can also look at the validation curve of the Softmax vs L2-SVMs as a function of weight updates in Fig. 2. As learning rate is lowered during the latter

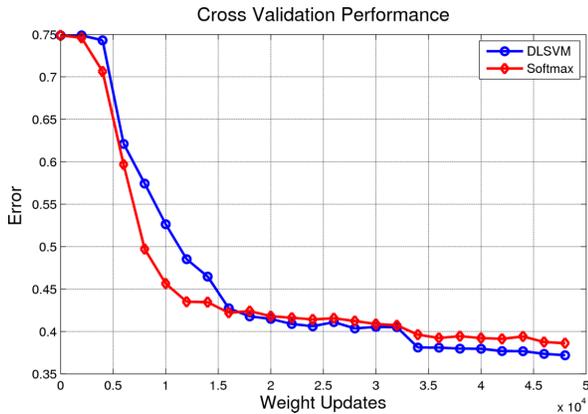


Figure 2. Cross validation performance of the two models. Result is averaged over 8 folds.

half of training, DLSVM maintains a small yet clear performance gain.

We also plotted the 1st layer convolutional filters of the two models:

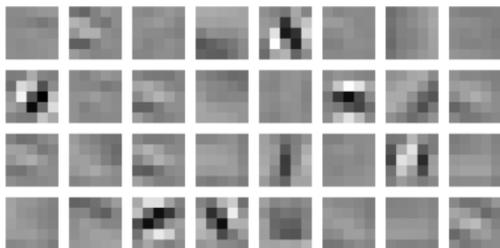


Figure 3. Filters from convolutional net with softmax.

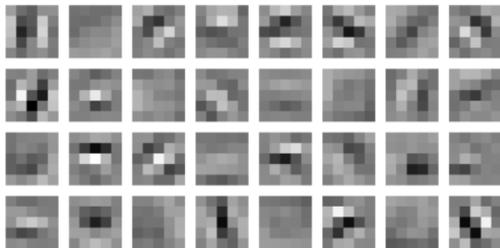


Figure 4. Filters from convolutional net with L2-SVM.

While not much can be gain from looking at these filters, SVM trained conv net appears to have more textured filters.

3.2. MNIST

MNIST is a standard handwritten digit classification dataset and has been widely used as a benchmark dataset in deep learning. It is a 10 class classification problem with 60,000 training examples and 10,000 test cases.

We used a simple fully connected model by first performing PCA from 784 dimensions down to 70 dimensions. Two hidden layers of 512 units each is followed by a softmax or a L2-SVM. The data is then divided up into 300 minibatches of 200 samples each. We trained using stochastic gradient descent with momentum on these 300 minibatches for over 400 epochs, totaling 120K weight updates. Learning rate is linearly decayed from 0.1 to 0.0. The L2 weight cost on the softmax layer is set to 0.001. To prevent overfitting and critical to achieving good results, a lot of Gaussian noise is added to the input. Noise of standard deviation of 1.0 (linearly decayed to 0) is added. The idea of adding Gaussian noise is taken from these papers (Raiko et al., 2012; Rifai et al., 2011b).

Our learning algorithm is permutation invariant without any unsupervised pretraining and obtains these results: **Softmax: 0.99%** **DLSVM: 0.87%**

An error of 0.87% on MNIST is probably (at this time) state-of-the-art for the above learning setting. The only difference between softmax and DLSVM is the last layer. This experiment is mainly to demonstrate the effectiveness of the last linear SVM layer vs. the softmax, we have not exhaustively explored other commonly used tricks such as Dropout, weight constraints, hidden unit sparsity, adding more hidden layers and increasing the layer size.

3.3. CIFAR-10

Canadian Institute For Advanced Research 10 dataset is a 10 class object dataset with 50,000 images for training and 10,000 for testing. The colored images are 32×32 in resolution. We trained a Convolutional Neural Net with two alternating pooling and filtering layers. Horizontal reflection and jitter is applied to the data randomly before the weight is updated using a minibatch of 128 data cases.

The Convolutional Net part of both the model is fairly standard, the first C layer had $32 \ 5 \times 5$ filters with Relu hidden units, the second C layer has $64 \ 5 \times 5$ filters. Both pooling layers used max pooling and downsampled by a factor of 2.

The penultimate layer has 3072 hidden nodes and uses Relu activation with a dropout rate of 0.2. The dif-

ference between the Convnet+Softmax and ConvNet with L2-SVM is the mainly in the SVM's C constant, the Softmax's weight decay constant, and the learning rate. We selected the values of these hyperparameters for each model separately using validation.

	ConvNet+Softmax	ConvNet+SVM
Test error	14.0%	11.9%

Table 2. Comparisons of the models in terms of % error on the test set.

In literature, the state-of-the-art (at the time of writing) result is around 9.5% by (Snoeck et al. 2012). However, that model is different as it includes contrast normalization layers as well as used Bayesian optimization to tune its hyperparameters.

3.4. Regularization or Optimization

To see whether the gain in DLSVM is due to the superiority of the objective function or to the ability to better optimize, We looked at the two final models' loss under its own objective functions as well as the other objective. The results are in Table 3.

	ConvNet+Softmax	ConvNet+SVM
Test error	14.0%	11.9%
Avg. cross entropy	0.072	0.353
Hinge loss squared	213.2	0.313

Table 3. Training objective including the weight costs.

It is interesting to note here that lower cross entropy actually led a higher error in the middle row. In addition, we also initialized a ConvNet+Softmax model with the weights of the DLSVM that had 11.9% error. As further training is performed, the network's error rate gradually increased towards 14%.

This gives limited evidence that the gain of DLSVM is largely due to a better objective function.

4. Conclusions

In conclusion, we have shown that DLSVM works better than softmax on 2 standard datasets and a recent dataset. Switching from softmax to SVMs is incredibly simple and appears to be useful for classification tasks. Further research is needed to explore other multiclass SVM formulations and better understand where and how much the gain is obtained.

Acknowledgment

Thanks to Alex Krizhevsky for making his very fast CUDA Conv kernels available! Many thanks to Relu Patrascu for making running experiments possible! Thanks to Ian Goodfellow, Dumitru Erhan, and Yoshua Bengio for organizing the contests.

References

Boser, Bernhard E., Guyon, Isabelle M., and Vapnik, Vladimir N. A training algorithm for optimal margin classifiers. In *Proceedings of the 5th Annual ACM Workshop on Computational Learning Theory*, pp. 144–152. ACM Press, 1992.

Ciresan, D., Meier, U., Masci, J., Gambardella, L. M., and Schmidhuber, J. High-performance neural networks for visual object classification. *CoRR*, abs/1102.0183, 2011.

Coates, Adam, Ng, Andrew Y., and Lee, Honglak. An analysis of single-layer networks in unsupervised feature learning. *Journal of Machine Learning Research - Proceedings Track*, 15:215–223, 2011.

Collobert, R. and Bengio, S. A gentle hessian for efficient gradient descent. In *IEEE International Conference on Acoustic, Speech, and Signal Processing, ICASSP*, 2004.

Dahl, G. E., Ranzato, M., Mohamed, A., and Hinton, G. E. Phone recognition with the mean-covariance restricted Boltzmann machine. In *NIPS 23*. 2010.

Hadsell, Raia, Chopra, Sumit, and Lecun, Yann. Dimensionality reduction by learning an invariant mapping. In *In Proc. Computer Vision and Pattern Recognition Conference (CVPR06)*. IEEE Press, 2006.

Hsu, Chih-Wei and Lin, Chih-Jen. A comparison of methods for multiclass support vector machines. *IEEE Transactions on Neural Networks*, 13(2):415–425, 2002.

Huang, F. J. and LeCun, Y. Large-scale learning with SVM and convolutional for generic object categorization. In *CVPR*, pp. I: 284–291, 2006. URL <http://dx.doi.org/10.1109/CVPR.2006.164>.

Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. What is the best multi-stage architecture for object recognition? In *Proc. Intl. Conf. on Computer Vision (ICCV'09)*. IEEE, 2009.

Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. In *NIPS*, pp. 1106–1114, 2012.

Lee, H., Grosse, R., Ranganath, R., and Ng, A. Y. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *Intl. Conf. on Machine Learning*, pp. 609–616, 2009.

Mohamed, A., Dahl, G. E., and Hinton, G. E. Deep belief networks for phone recognition. In *NIPS Workshop on Deep Learning for Speech Recognition and Related Applications*, 2009.

- Nagi, J., Di Caro, G. A., Giusti, A., , Nagi, F., and Gambardella, L. Convolutional Neural Support Vector Machines: Hybrid visual pattern classifiers for multi-robot systems. In *Proceedings of the 11th International Conference on Machine Learning and Applications (ICMLA)*, Boca Raton, Florida, USA, December 12–15, 2012.
- Quoc, L., Ngiam, J., Chen, Z., Chia, D., Koh, P. W., and Ng, A. Tiled convolutional neural networks. In *NIPS 23*. 2010.
- Raiko, Tapani, Valpola, Harri, and LeCun, Yann. Deep learning made easier by linear transformations in perceptrons. *Journal of Machine Learning Research - Proceedings Track*, 22:924–932, 2012.
- Rifai, Salah, Dauphin, Yann, Vincent, Pascal, Bengio, Yoshua, and Muller, Xavier. The manifold tangent classifier. In *NIPS*, pp. 2294–2302, 2011a.
- Rifai, Salah, Glorot, Xavier, Bengio, Yoshua, and Vincent, Pascal. Adding noise to the input of a model trained with a regularized objective. Technical Report 1359, Université de Montréal, Montréal (QC), H3C 3J7, Canada, April 2011b.
- Salakhutdinov, Ruslan and Hinton, Geoffrey. Learning a nonlinear embedding by preserving class neighbourhood structure. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, volume 11, 2007.
- Vapnik, V. N. *The nature of statistical learning theory*. Springer, New York, 1995.
- Vinyals, O., Jia, Y., Deng, L., and Darrell, T. Learning with Recursive Perceptual Representations. In *NIPS*, 2012.
- Weston, Jason, Ratle, Frdric, and Collobert, Ronan. Deep learning via semi-supervised embedding. In *International Conference on Machine Learning*, 2008.
- Zhong, Shi and Ghosh, Joydeep. Decision boundary focused neural network classifier. In *Intelligent Engineering Systems Through Artificial Neural Networks*, 2000.