

HYPER-PARAMETER OPTIMIZATION OF DEEP CONVOLUTIONAL NETWORKS FOR OBJECT RECOGNITION

Sachin S. Talathi

Qualcomm Research Center, 5775 Morehouse Dr, San Diego CA 92121

ABSTRACT

Recently sequential model based optimization (SMBO) has emerged as a promising hyper-parameter optimization strategy in machine learning. In this work, we investigate SMBO to identify architecture hyper-parameters of deep convolution networks (DCNs) object recognition. We propose a simple SMBO strategy that starts from a set of random initial DCN architectures to generate new architectures, which on training perform well on a given dataset. Using the proposed SMBO strategy we are able to identify a number of DCN architectures that produce results that are comparable to state-of-the-art results on object recognition benchmarks.

Index Terms— hyper-parameter optimization, deep convolution networks, sequential model based optimization

1. INTRODUCTION

The primary task for a supervised machine learning algorithm is to use training dataset $\{x_{tr}, y_{tr}\}$ to find a function $f : x \rightarrow y$, that also generalize well across the test (or the hold out) dataset $\{x_h, y_h\}$. Very often f is obtained through the optimization of a training criterion, C , with respect to a sets of parameters, θ . The learning algorithm used to optimize C usually contains its own set of free parameters λ_l , referred to as the learning algorithm hyper-parameters. These hyper-parameters are often estimated using grid search cross validation. In addition to the learning algorithm hyper-parameters, λ_l , neural network models such as deep convolution networks (DCNs) also comprise of hyper-parameters $\lambda_a \notin \lambda_l$, that define the architectural configuration of the network. Grid search techniques are prohibitively expensive to tune λ_a , given the fact that there are a few tens of these architectural hyper-parameters. As a result, many of the state-of-the-art DCNs are manually designed, making the task of tuning these hyper-parameters more of an art than a science [1].

In recent years, there has been a concerted effort in the machine learning community to develop better algorithms to solve the hyper-parameter optimization problem [1, 2, 3, 4, 5]. Many of these works have successfully applied direct search methods for non-linear optimization such as the sequential model based optimization (SMBO) to generate better results on various supervised machine learning tasks than were previously reported. Motivated by these works, in this paper we attempt to address the question: *Can SMBO be used to determine superior architectural configurations for DCNs?* The paper is organized as follows: In the Methods Section 2, we will briefly formulate the problem of hyper-parameter optimization for DCNs. We then present the general strategy of sequential model based optimization (SMBO) [4] and summarize our approach to SMBO for designing DCN architectures. The results of image classification training and evaluation on the benchmark CIFAR-10 dataset are then presented in the Results Section 3, which is followed by the Conclusion.

2. METHODS

2.1. Formulation of the problem

Let $M_\lambda(x, w)$ represent the DCN model that operates on input data $x \in \mathbb{R}^D$ and generates an estimate \hat{y} for the output data $y \in \mathbb{Z}_2^C$. The DCN model, M , is parameterized by two set of parameters, the first being w , which are obtained through the optimization of a training criterion, C , using a gradient descent type learning algorithm such as the back-propagation algorithm and the second being $\lambda = \{\lambda_a, \lambda_l\}$, which represent the set of the so called hyper-parameters. The hyper-parameters λ_a define the DCN architecture and the hyper-parameters λ_l , are associated with the learning algorithm used to optimize C . The objective for DCN hyper-parameter optimization is

to solve the joint optimization problem as stated below:

$$\begin{aligned} \{w, \lambda\} &= \underset{\lambda}{\operatorname{argmin}} [\Psi(\hat{y}_h, y_h)] \quad \text{where,} \\ \hat{y}_h &= M_\lambda(\underset{w}{\operatorname{argmin}} (C(x_{tr}, y_{tr}, \theta)), x_h) \quad (1) \end{aligned}$$

where $\{(x_{tr}, y_{tr}), (x_h, y_h)\} \in (x, y)$, are the input and output training and hold-out (or the test) data set respectively, $\Psi = \sum_{\{y_h\}} \mathbb{I}_{y_h \neq \hat{y}_h}$.

2.2. Sequential model based optimization

SMBO is a direct search method for non-linear optimization, in which one begins by selecting a meta-model of the function for which an extrema is sought. One then applies an active learning strategy to select a query point that provides the most potential to approach the extrema. More specifically, let us assume that we have a database $D_{1:t} = \{\lambda_{1:t}, e_{1:t}\}$ of t DCN models, where $\lambda_i|_{i=1}^t$ represents the set of hyper-parameters and $e_i|_{i=1}^t$ is the validation error on the hold-out dataset generated by each of the t DCN models. The basic idea underlying SMBO is to replace the original optimization problem of finding extrema of a given function, such as $\Psi(\lambda)$, which is time consuming and computationally expensive, with an equivalent problem of optimization of expected value of an utility function, $u(e)$ [4]. As we will see below, optimizing over the expected value of the utility function is computationally much cheaper and faster than solving the original problem. Under SMBO, one usually begins by assigning a prior distribution $p(e)$ on e . One then uses the database $D_{1:t}$ to obtain an estimate for the likelihood function $p(\lambda_{1:t}|e)$. The prior and the likelihood function are used to obtain an estimate for the posterior distribution $p(e|\lambda_{1:t}) \propto p(\lambda_{1:t}|e)p(e)$. The objective then is to choose λ_{t+1} , which maximizes $u(e)$ under $p(e|\lambda_{1:t})$, i.e., $\lambda_{t+1} = \underset{\lambda}{\operatorname{argmax}} [E(u(e))]$, where $E(u(e))$ is given as:

$$\begin{aligned} E(u(e)) &= \int u(e)p(e|\lambda_{1:t})de \\ &= \int u(e) \frac{p(\lambda_{1:t}|e)p(e)}{p(\lambda_{1:t})} de \quad (2) \end{aligned}$$

A common choice for the utility function $u(e)$, is the Expected Improvement function [4], $u(e) = \max((e^* - e), 0)$, in which case, Eq. 4 becomes

$$E(u(e)) = \int_0^{e^*} (e^* - e)p(e|\lambda_{1:t})de \quad (3)$$

Then under SMBO we have,

$$\lambda_{t+1} = \underset{\lambda}{\operatorname{argmax}} \left[\frac{e^*}{p(\lambda_{1:t})} \int_0^{e^*} p(\lambda_{1:t}|e)p(e)de \right] \quad (4)$$

HyperOpt ($D_{1:t}, p, N$):

While $t \leq N$ **do**:

1. Estimate e^* s.t. $p(e_{1:t} < e^*) = 0.5$
2. Use e^* and $\lambda_{1:t}$ to estimate $l(\lambda)$ and $g(\lambda)$
3. Evaluate λ_{t+1} according to Eq. 5 and Eq. 6
4. Train the new DCN model with λ_{t+1} to estimate e_{t+1}
5. $t \leftarrow t + 1$

Table 1. SMBO algorithm for estimating architecture hyper-parameters for DCN.

If $p(e < e^*) = \gamma$ (a constant), i.e., choose e^* to be some quantile of observed e values, and we define two density functions; $l(\lambda) = p(\lambda_{1:t}|e)$ when $e \leq e^*$ and $g(\lambda) = p(\lambda_{1:t}|e)$ when $e > e^*$ as proposed in [1], then,

$$\lambda_{t+1} = \underset{\lambda}{\operatorname{argmax}} \left[\frac{e^* \gamma l(\lambda)}{\gamma l(\lambda) + (1 - \gamma)g(\lambda)} \right] \quad (5)$$

Bergstra et.al., [1], proposed an adaptive Parson estimator algorithm to evaluate Eq. 5 so as to maximize the ratio $l(\lambda)/g(\lambda)$. In this work, we propose a simple strategy to evaluate Eq. 5 as follows: let $l(\lambda) + g(\lambda) = U(\lambda) = k$ (a constant). In other words, λ is drawn from a uniform distribution. Then, if $\gamma = 0.5$, we have from Eq. 5

$$\lambda_{t+1} = \frac{2e^*}{k} \underset{\lambda}{\operatorname{argmax}} [l(\lambda)] \quad (6)$$

Our proposed simplification in Eq. 6, does not require an estimate for both $l(\lambda)$ and $g(\lambda)$. We can simply pick $\lambda \in U(\lambda)$ and evaluate according to empirical distribution of $l(\lambda)$ generated from $D_{1:t}$ to evaluate the next potential λ_{t+1} . Since the proposed algorithm chooses $\lambda \in U(\lambda)$ at every step, a much larger space of potential λ 's are explored, which may in turn slow down the convergence rate to an optimal λ^* . In order to counter this, we adopt a hybrid strategy by combining Eq. 5 and Eq. 6. With probability p , at every iteration, we choose the potential λ_{t+1} according to Eq. 5 and with probability $1 - p$, we choose the potential λ_{t+1} according to Eq. 6. In Table 1, we summarize the steps involved in our proposed algorithm.

3. RESULTS

We evaluate our proposed SMBO algorithm on the CIFAR-10 benchmark, which consists of 60,000 32x32 color images. The collection is split into 50,000 training and 10,000 testing images. All the DCNs generated

by our proposed algorithm were trained using cuda-convnet2¹. We used a 3 step cooling procedure; starting with learning rate $l = 0.01$, the momentum $m = 0.9$, the weight decay parameter $w_c = 0.0005$ for the first 120 epochs followed by another 20 epochs by reducing learning rate by a factor of 10 (keeping other parameters the same) and then training for 10 more epochs by further reducing learning rate by a factor of 10.

Since the primary focus for us in this work is to determine whether SMBO can be used to identify suitable DCN architectures, we fixed the DCN hyper-parameters associated with the back-propagation learning algorithm as described above. The set of DCN architecture hyper-parameters that we consider for optimization are listed in Table 2.

Conv. Layer	1. No. of Conv. layers 2. No. of filters per layer 3. Filter size; 4. Filter stride
Norm. Layer	1. Size
Pool Layer	1. Size; 2. Stride 3. type of pooling: max/avg.
Hidden Layer	1. No. of hidden layers 2. No. of nodes per hidden layer 3. dropout value

Table 2. DCN architecture hyper-parameters. In addition to the parameters listed above; we also consider two additional boolean hyper-parameters to represent the presence or absence of the Norm layer or the pool layer and a third boolean hyper-parameter indicating the presence or absence of dropout. For normalization layer; we only consider local response normalization across filter maps [6], with a scaling factor of 0.75.

It has been reported in the literature [7] that very deep networks are difficult to train primarily suffering from vanishing gradient problem at larger depths. In order to alleviate this problem, for our implementation of SMBO for DCN architectural hyper-parameters, all the DCNs are generated to comprise a local logistic regression (LR) cost function layer at the output of one or more of the convolution block.

For the results presented here, we consider $t = 32$ as the size of our initial database based on our analysis of random search hyper-parameter optimization [5] and we set $p = 0.9$.

¹<https://code.google.com/p/cuda-convnet2/>

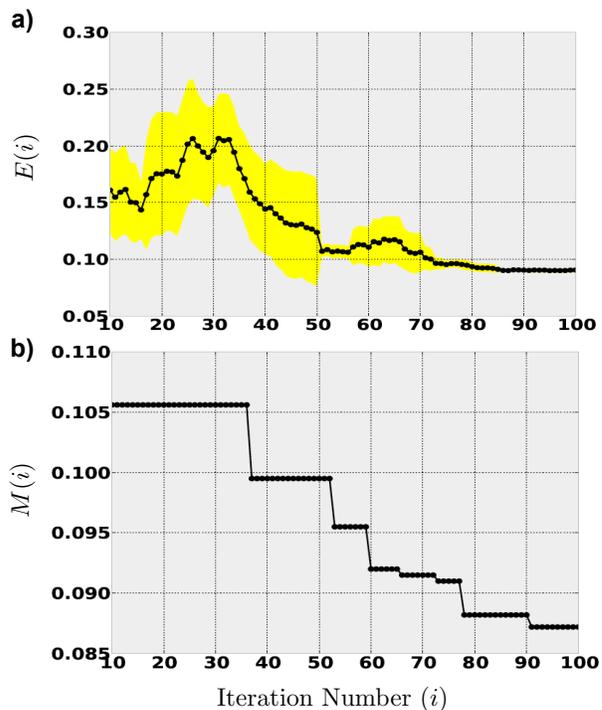


Fig. 1. (a) The mean test error and standard deviation (in yellow) as function of the SMBO iteration number for multi-view mode (b) The minimum multi-view mode test error as function of SMBO iteration number

Architecture	# Trainable Parameters	Parameters	Convolutions							Fully-connected			Test error	
			Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7	Layer 8	Layer 1	Layer 2	Softmax	550 epochs
DCN 1.	≈30.9 M	Filter / size	64x3x3	256x3x3	256x3x3	256x3x3	256x3x3	256x3x3	256x3x3	256x11x11	3,314	4,951	10	7.81%
		Stride	1	1	2	2	1	2	2	10				
		Padding	0	0	1	1	0	1	1	9				
		Pooling (Size, Stride)		(2,2)										
		Normalization					✓			✓	✓			
		Dropout									✓	✓		
DCN 2.	≈4.0 M	Filter / size	128x3x3	128x3x3	128x3x3	256x3x3	256x3x3	256x7x7	NA	NA	NA	NA	10	8.17%
		Stride	1	2	1	1	1	2						
		Padding	0	1	0	0	0	1						
		Pooling (Size, Stride)												
		Normalization		✓			✓	✓						
		Dropout												
DCN 3.	≈3.4 M	Filter / size	256x3x3	128x3x3	256x3x3	256x3x3	256x3x3	128x7x7	NA	NA	NA	NA	10	8.63%
		Stride	1	2	1	1	2	5						
		Padding	0	1	0	0	1	2						
		Pooling (Size,Stride)												
		Normalization	✓											
		Dropout												

Table 3. Architecture for the top 3 DCNs generated by our proposed SMBO algorithm.

In Figure 1a, we plot the mean (std. error, shown in yellow) test error (evaluated in multi-view test mode, [6]) $E(i) = \langle e_{i-10:i} \rangle$ and in Figure 1b, we plot the minimum test error $M(i) = \min(e_{0:i})$ as function of the iteration number i , respectively. We see that the average test error gradually decrease towards an optimal solution, the best minimum found also decreases with increasing iterations. Furthermore, our proposed SMBO procedure generated a large number of “good” DCN architectures that produce test-error of $< 11\%$ even with only 150 training epochs (not-reported). In comparison, the best hand-tuned DCN architecture, produced by [6] exhibits 11% test error in multi-view mode and requires a longer training time on the order of 500 epochs.

Since the state-of-the-art performance numbers for the CIFAR-10 benchmark dataset are usually reported in multi-view mode (with data-augmentation [6]), we report multi-view test error of 7.81% for the best DCN generated by our proposed hyper-parameter optimization strategy, which compares favorably to the current state-of-the-art result on CIFAR-10 of 7.97% [8]. In Table 3, we summarize the top 3 DCN architectures found by our proposed SMBO procedure that produced multi-view test error $< 9\%$ on the CIFAR-10 benchmark. In ??, we summarize the number of parameters in each of these 3 DCN models.

At the time of writing of this manuscript for camera ready version, we came across a recent paper [9], that reported multi-view test error of 7.25% on CIFAR-10 benchmark, using a hand designed DCN network, that is comprised of only convolution layers and has 1.3 M parameters. Yet another paper [10], reported the

utility of using parametric-relu neuron as opposed to the relu neurons. While none of the optimized DCN networks that we report in Table 3 generate better performance numbers than the latest state-of-the-art numbers reported in [9], we wanted to determine whether the use of parametric-relu neuron can boost the performance of the optimized DCN networks that we have identified through the hyper-parameter optimization approach. Accordingly, we retrained the smallest of the three DCN networks from Table 3 using a version of parametric-rectified non-linear neurons of type $y = ax(x \leq 0) + \sqrt{x}(x > 0)$, where a is a learnable parameter, fixed per neuron layer in the DCN network. We were able to obtain multi-view test error score of 6.9%, which to the best of our knowledge, represents the state-of-the-art score for CIFAR-10 benchmark. In Table 4, we summarize all the known best in class numbers for CIFAR-10 benchmark.

Method	Activation Function Type	Error %	Trainable Params
Maxout [11]	maxout	9.38	>6 M
Dropconnect [12]	relu	9.32	-
dasNet [13]	maxout	9.22	>6 M
Network in Network [14]	relu	8.81	≈1 M
Deeply Supervised [15]	relu	7.97	≈1 M
All-CNN [9]	relu	7.25	≈1.3 M
DCN-3 (Ours)	p-relu	6.9	≈3.4 M

Table 4. Comparison of state-of-the-art results for CIFAR-10 benchmark; relu: rectified linear unit; p-relu: parametric-rectified non-linear unit.

4. CONCLUSION

In this paper, we have proposed a simple SMBO algorithm and a recipe for hyper-parameter optimization of DCN architectures. We have demonstrated that SMBO can be used to generate a large number of “good” DCN architectures, which may then form a backbone for further investigations. Our results suggest that indeed SMBO can be used to identify superior DCNs. In summary, our work in this paper in addition to those from earlier works [1, 3] broaden the scope of the models that can be realistically investigated, without the need for the researchers to be restricted to manual evaluation of a few architectural parameters at any given time.

5. REFERENCES

- [1] J. Bergstra, R. Bardenet, Y. Bengio, and B. Kegl, “Algorithms for hyper-parameter optimization,” in *NIPS*, 2011, vol. 24.
- [2] *Distributed asynchronous hyper parameter optimization*. <https://github.com/hyperopt/hyperopt>.
- [3] J. Snoek, H. Larochelle, and R.P. Adam, “Practical bayesian optimization for machine learning,” in *NIPS*, 2012, vol. 25.
- [4] E. Brochu, V.M. Cora, and N. Freitas, “A tutorial on bayesian optimization of expensive cost functions, with application to active user modeling and hierarchical reinforcement learning,” *Arxiv*, vol. arXiv:1012.2599, 2010.
- [5] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *Journal of Machine Learning Research*, vol. 13, pp. 281–305, 2012.
- [6] A. Krizhevsky, “Learning multiple layers of features from tiny images,” Tech. Rep., Dept. Computer Science, University of Toronto, 2009.
- [7] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large scale image recognition,” *Arxiv*, vol. arXiv:1409.1556, 2014.
- [8] C. Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, “Deeply supervised nets,” in *NIPS*, 2014.
- [9] J.T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, “Striving for simplicity, all convolution net,” in *ICLR*, 2015.
- [10] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers, surpassing human-level performance on imagenet classification,” in *Arxiv:1502.01852*, 2015.
- [11] I.J. Goodfellow, F. Warde, D. Mirza, A. Courville, and Y. Bengio, “Maxout networks,” in *ICML*, 2013.
- [12] L. Wan, M. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, “Regularization of neural networks using dropconnect,” in *ICML*, 2013.
- [13] M. F. Stollenga, J. Masci, F. Gomez, , and J Schmidhuber, “Deep networks with internal selective attention through feedback connections.,” in *NIPS*, 2014.
- [14] M. Lin, Q. Chen, , and S. Yan, “Network in network,” in *ICLR*, 2014.
- [15] C.Y. Lee, S. Xi, P. Gallagher, Z. Zhang, and Z. Tu, “Deeply supervised nets,” in *NIPS*, 2014.