

ShrinkNets

Guillaume Leclerc*
Massachusetts Institute of
Technology
Cambridge, Massachusetts, USA
leclerc@mit.edu

Raul Castro Fernandez
Massachusetts Institute of
Technology
Cambridge, Massachusetts, USA
raulcf@csail.mit.edu

Samuel Madden
Massachusetts Institute of
Technology
Cambridge, Massachusetts, USA
madden@csail.mit.edu

1 INTRODUCTION

When designing Neural Networks, finding the appropriate size (width and depth) is key. In particular, these hyper parameters have a strong correlation with over/underfitting **Sam: give reference** **GI: I did some literature review and actually this statement is quite controversial, should we rephrase ?**. We have no reliable way to find them but decades of experimentation led to some heuristics [1] that try to prune that immense space of possible network sizes and suggest range of values to explore. To select the best candidates, researchers have used strategies such as random search **GI: Should we also move the references here if we remove the paragraph from the related work ?**, meta-gradient descent [14], gaussian processes [2], and Parzen Estimators [2] to reduce find good parameters without exhaustively exploring the entire hyper-parameter space. Although these strategies help with finding a good set of hyper-parameters, they still require a compute-intensive search of the space.

In this paper we present a method to automatically find an appropriate network size from a single parameter, drastically reducing the hyper-parameter dimensionality. The key idea is to *learn* the right network size at the same time that the network is learning the main task. For example, for an image classification task, with our approach we can provide the training data to a network—without sizing it a priori—and expect to end up with a network that has learnt the task without overfitting **GI: too strong, we do not really have any guarantee on the generalization, how about changing it to: to end up with a network that learnt a tradeoff between size and accuracy**. This approach has two main benefits. First, we do no longer need to choose a network size before training. Second, the final network size will be appropriate for the task at hand, and not larger **GI: Should we talk about the fact that it automatically gets rid of dead neurons ?**. This is important because oversized networks have a lower inference throughput and higher memory footprint.

Our approach has two main challenges. First, on how to dynamically size the network during training. Second, on how to find a loss function that optimizes for the additional task of sizing, without deteriorating the learning performance of the main task. We describe next ShrinkNets, which cope with both challenges:

2 SHRINKNETS

Our approach consists of starting the training process for the task of interest with an explicitly oversized network. Then, as training progresses, we learn which neurons are not contributing to learning the main task and remove them dynamically, thus shrinking the network size. This method requires two building blocks. First,

a way of identifying neurons that are not contributing to the learning process, and second a way of balancing the network size and the generalization capability for the main task. We introduce a new layer, called Filter, which takes care of *deactivating* neurons. We also modify existing loss functions to incorporate a new term that takes care of balancing network sizing and generalization capability appropriately. We explain them next:

Filter Layers: They have weights in the range $[0, +\infty]$ and are placed after linear and convolutional layers of traditional deep networks **GI: This is how I used it but with the new implementation we can be more creative and put them anywhere, should we say "usually placed" instead ?**. The *Filter Layer* takes an input of size $(B \times C \times D_1 \times \dots \times D_n)$, where B is the batch size, C the number of features (or channels), and D any additional dimension. This structure makes it compatible with fully connected layers with $n = 0$ or convolutional layers with $n = 2$. Their crucial property is a parameter $\theta \in \mathbb{R}^C$, defined as follows:

$$\text{Filter}(I; \theta) = I \circ \max(0, \theta) \quad (1)$$

Where \circ is the Hadamard product (pointwise multiplication **GI: Should we keep only one of them to be more concise ?**), and θ is expanded in all dimensions except the second one to match the input size. It is easy to see that if for any k , if $\theta_k \leq 0$, the k^{th} feature/channel will forever be 0. When this happens, we say the Filter layer disables the neuron. These disabled neurons/channels can be removed from the network without changing its output (we describe how we perform this removal below). We explain next how the weights of the Filter Layer are adjusted during training.

Training Procedure: Once Filter layers are placed in a deep network, we could train it directly and it would be equivalent to a normal neural network. However, our goal is to find the smallest network with reasonable performance. We achieve that by introducing sparsity in the parameters of the *Filter Layers*. Indeed, having a negative component in the θ parameter of the filter layer permanently disable its associated feature **GI: Maybe redundant ? we talked about that in the previous paragraph**. To obtain this sparsity, we simply redefine the loss function:

$$L'(x, y; \theta) = L(x, y) + \lambda |\theta| \quad (2)$$

This choice come from the Lasso loss [17] and we use it for the same reason: it introduces sparsity. The bigger the λ the more entries in θ are set to zero, and since zero entries in θ correspond to dead neurons, lambda effectively control the number of neurons/channels in the entire network. Next, we explain how to implement ShrinkNets efficiently.

*Visiting Student

Software architecture¹: *Filter Vectors* can easily be implemented in a few lines of code in many existing deep learning frameworks. However, in ShrinkNets, we assume that we start with obviously oversized networks. If disabled neurons are not quickly removed, the overhead might cause the training process to be significantly slower than classic neural networks. This is why we want to support what we call garbage collection: removing the weights in the layers that are responsible (or using) features that are disabled.

To be able to capture this relationship between layers, we augmented *Pytorch* [13] (the framework we used as the foundation of our implementation) with a graph structure very similar to the one available in *Keras* [3]. In this graph, edges are effectively event hubs responsible for propagating *feature removal* events to its endpoints. ShrinkNet layers are designed to emit and react to these events but also propagate the event further if it has an impact at the other endpoint of the layer (input or output) **GI: What I mean here is: if the event comes from the output, it might be propagated to the inputs, and the other way around.**

This event-based implementation coordinated by edges makes it very easy to integrate new layers in the library. We even provide automatic wrapping for layers from *PyTorch* if they have no internal state. For more complex ones, we provide utilities that makes the implementation very concise.

3 EVALUATION

3.1 Convergence

To demonstrate that the approach is viable, we will first show that Shrinking Networks converge properly. For this experiment we trained a one hidden layer neural network with one filter layer to control the number of hidden units. We initialized the models with 10000 neurons and trained them on MNIST [10] using different regularization factors (λ). We summarized the results in Figure 1. We can see on this plot is that for a given λ the number of hidden units seems to always converge to some value. It has two implications: it seems that λ is, as we would think, a proxy of the network size (bigger λ imply smaller networks). Secondly, it indicates that the regular spikes we see in the loss function that occur when neurons are disabled will eventually disappear because the number of neurons reach a plateau. **GI: I feel we should give a conclusion for this paragraph but I am not sure what to say**

3.2 Relevance in the context of hyper-parameter optimization

One of our main goal is to help finding neural network architectures that perform reasonably well faster than existing techniques. For simplicity, we will try to see how ShrinkNets perform when doing random search over a sub set of the hyper-parameter space (we will fix everything except the size-related parameters). However, we believe that the results generalize to more complex methods, most of them might actually benefit a lot from the reduction of the dimensionality of the hyper-parameter space.

To isolate the effect of ShrinkNets we will consider two very simple architectures: multi-layer perceptrons (with three hidden

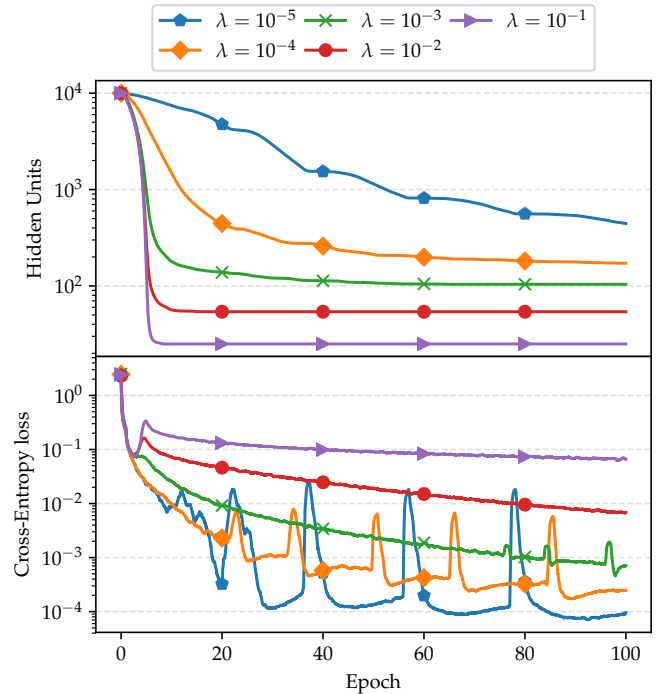


Figure 1: Evolution of the number of hidden units and loss over time on the MNIST dataset for different λ values

layers) and the *LeNet-5* model [10]. For this experiment we assume we have no prior information about the size except an upper bound of 50 channels per convolutional layer and 5000 neuron for fully connected layers. We try to find the appropriate size of the network using classical neural networks and ShrinkNets. To make sure that we are as fair as possible and since one might argue that ShrinkNets is also regularizing the network in addition to finding the size, we will also considered regularized classic networks using the most commonly used technique: the L_2 penalty [12].

The experiment consist in sampling parameters randomly for each class of network, train them for the same amount of epochs (100 for MLPs and 200 for CNNs). Pick the epoch that performed the best on the validation set and evaluate it on the testing set. We repeat the process for 50 models. This evaluation was done on MNIST [10], FashionMNIST [18] and CIFAR10 [9]. The distribution of the testing accuracy we obtained is summarized on Figure 2.

We can see that the distributions for ShrinkNets are consistently better than the others for all datasets and both architectures. It shows that by training fewer networks, we are likely to obtain on-par or even better models which effectively reduces the cost of hyper-parameter optimization. In some cases the best ShrinkNets models even outperforms the best model of the other techniques.

4 RELATED WORK

There are many methods in the litterature that try to simplify the neural network structure. Most of them focuses on removing connections (eg: [4], [6]). However, most hardware (GPU's, TPUS [8],

¹The code is available as Python/PyTorch library on <http://github.com/mitdbg/fastdeepnets> **GI: Should we rename the repository ?**

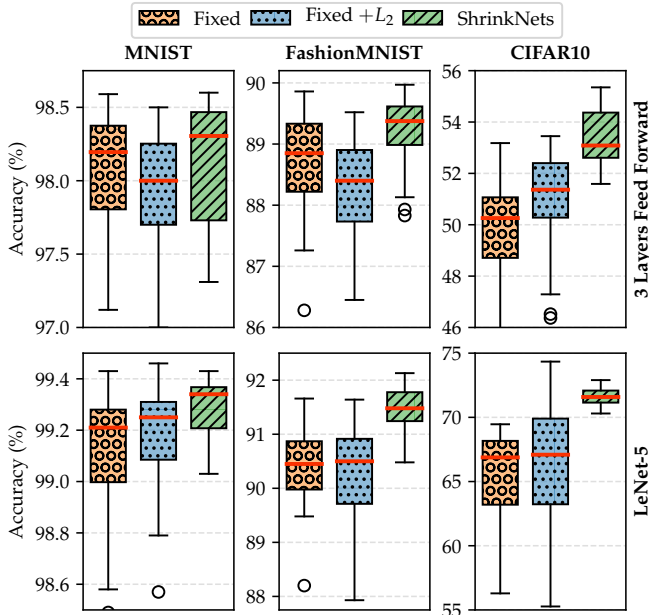


Figure 2: Distributions of the testing accuracy for different training methods, datasets and architectures using random search

TensorCores) do not really leverage sparsity in the connections efficiently (as explained in [5]). On the other side, ShrinkNets and some others [16] and [15] try to remove entire neuron instead of connections. This is very useful because it reduces the size of the matrices, producing speed-up on every device. ShrinkNets improves on [16] because it removes neurons during training, speeding up the rest of the process and beats [15] on convergence speed and because it does not need to change the optimizer. To the best of our knowledge this is also the first contribution that tries to learn the number of channels of a convolutional neural network.

5 CONCLUSION

In this paper we presented a novel technique to guess a reasonable network size based on single parameter λ that control the tradeoff between loss and the size. We demonstrated that it works both on fully connected networks and convolutional neural networks. Even though the firsts results seem promising, there are many ways we could improve. In the current implementation we only "learn" the number of features (neurons or channels). We could try to augment it with dynamic number of layers as seen in [11] to be able to determine the entire architecture.

We saw on Figure 1 that the loss temporarily suffers from the removal of neurons. It is likely that the loss would be more stable if the number of neurons converged faster or neurons disappeared slower. For this reason we plan to explore proximal gradient methods to optimize the filter vectors and/or randomize neuron removals.

During our evaluation we picked small datasets mainly to be able to train many models and have statistically significant distributions. With more computation resources and time, we could see if it generalizes to bigger datasets and other architectures like ResNet [7] (small modifications to the existing code base are required to support them)

REFERENCES

- [1] Yoshua Bengio. 2012. Practical recommendations for gradient-based training of deep architectures. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* 7700 LECTU (2012), 437–478. https://doi.org/10.1007/978-3-642-35289-8_26 arXiv:1206.5533
- [2] James Bergstra, Rémi Bardenet, Yoshua Bengio, and Balázs Kégl. 2011. Algorithms for Hyper-Parameter Optimization. In *Advances in Neural Information Processing Systems (NIPS)*. 2546–2554. <https://doi.org/2012arXiv1206.2944S> arXiv:1206.2944
- [3] François Chollet et al. 2015. Keras. <https://github.com/fchollet/keras>. (2015).
- [4] Yann Le Cun, John S Denker, and Sara a Solla. 1990. Optimal Brain Damage. *Advances in Neural Information Processing Systems* 2, 1 (1990), 598–605. <https://doi.org/10.1.1.32.7223> arXiv:arXiv:1011.1669v3
- [5] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. 2016. EIE: Efficient Inference Engine on Compressed Deep Neural Network. In *Proceedings - 2016 43rd International Symposium on Computer Architecture, ISCA 2016*. 243–254. <https://doi.org/10.1109/ISCA.2016.30> arXiv:1602.01528
- [6] Song Han, Huizi Mao, and William J Dally. 2015. Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding. (2015). <https://doi.org/abs/1510.00149> arXiv:1510.00149
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 770–778. <https://doi.org/10.1109/CVPR.2016.90> arXiv:1512.03385
- [8] Norman P. Jouppi, Al Borchers, Rick Boyle, Pierre-luc Cantin, Clifford Chao, Chris Clark, Jeremy Coriell, Mike Daley, Matt Dau, Jeffrey Dean, Ben Gelb, Cliff Young, Tara Vazir Ghaemmaghami, Rajendra Gottipati, William Gulland, Robert Hagmann, C. Richard Ho, Doug Hogberg, John Hu, Robert Hundt, Dan Hurt, Julian Ibarz, Nishant Patil, Aaron Jaffey, Alek Jaworski, Alexander Kaplan, Harshit Khaitan, Daniel Killebrew, Andy Koch, Naveen Kumar, Steve Lacy, James Laudon, James Law, David Patterson, Diemthu Le, Chris Leary, Zhuoyuan Liu, Kyle Lucke, Alan Lundin, Gordon MacKean, Adriana Maggiore, Maire Mahony, Kieran Miller, Rahul Nagarajan, Gaurav Agrawal, Ravi Narayanaswami, Ray Ni, Kathy Nix, Thomas Norrie, Mark Omernick, Narayana Penukonda, Andy Phelps, Jonathan Ross, Matt Ross, Amir Salek, Raminder Bajwa, Emad Samadiani, Chris Severn, Gregory Sizikov, Matthew Snelham, Jed Souter, Dan Steinberg, Andy Swing, Mercedes Tan, Gregory Thorson, Bo Tian, Sarah Bates, Daria Toma, Erick Tuttle, Vijay Vasudevan, Richard Walter, Walter Wang, Eric Wilcox, Doe Hyun Yoon, Suresh Bhatia, and Nan Boden. 2017. In-Datcenter Performance Analysis of a Tensor Processing Unit. In *Proceedings of the 44th Annual International Symposium on Computer Architecture - ISCA '17*. 1–12. <https://doi.org/10.1145/3079856.3080246> arXiv:1704.04760
- [9] Alex Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. ... *Science Department, University of Toronto, Tech. ...* (2009), 1–60. <https://doi.org/10.1.1.222.9220> arXiv:arXiv:1011.1669v3
- [10] Y LeCun, L Bottou, Yoshua Bengio, and P Haffner. 2001. Gradient-Based Learning Applied to Document Recognition. In *Intelligent Signal Processing*. 306–351. <https://doi.org/10.1109/5.726791> arXiv:1102.0183
- [11] Benjamin Meier. [n. d.]. Going Deeper: Infinite Deep Neural Networks. ([n. d.]). https://github.com/kutoga/going_deeper/raw/master/doc/going_deeper.pdf
- [12] Andrew Y Ng. 2004. Feature selection, ℓ_1 vs. ℓ_2 regularization, and rotational invariance. In *ICML*. 78–85. <http://www.machinelearning.org/proceedings/icml2004/papers/354.pdf>
- [13] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic differentiation in PyTorch. (2017).
- [14] Fabian Pedregosa. 2016. Hyperparameter optimization with approximate gradient. (2016). arXiv:1602.02355 <https://arxiv.org/pdf/1602.02355.pdf> <http://arxiv.org/abs/1602.02355>
- [15] George Philipp and Jaime G Carbonell. 2017. Nonparametric Neural Network. In *Proc. International Conference on Learning Representations*. 1–27. arXiv:1712.05440 <https://www.cs.cmu.edu/~jggc/publication/NonparametricNeuralNe>
- [16] Simone Scardapane, Danilo Comminiello, Amir Hussain, and Aurelio Uncini. 2017. Group sparse regularization for deep neural networks. *Neurocomputing* 241 (2017), 81–89. <https://doi.org/10.1016/j.neucom.2017.02.029> arXiv:1607.00485

- [17] Robert Tibshirani. 1996. Regression Selection and Shrinkage via the Lasso. (1996), 267–288 pages. <https://doi.org/10.2307/2346178> arXiv:13697412/11/73273
- [18] Han Xiao, Kashif Rasul, and Roland Vollgraf. 2017. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. (2017). arXiv:1708.07747 <https://arxiv.org/pdf/1708.07747.pdf><http://arxiv.org/abs/1708.07747>