
Learning Network Size while Training with ShrinkNets

Anonymous Authors¹

Abstract

Let's write the abstract at the end

1. Tentative outline

- We want to figure out what is the proper network size
- If we had a oversized network and for each neuron we would have an on/off switch, we could optimize the state of each switch to achieve any size/accuracy tradeoff (we can prove that, do we care ?)
- Solving such problem is NP-Hard
- We could approximate the binary switch by relaxing them and doing L1 loss instead of L0 loss
- This is the definition of Shrinknets
- group sparsity tries to achieve a similar goal but with a different formulation
- How are they related ?
- If we add a specific constraint on our formulation we obtain the group sparsity one (proven)
- We conclude that without this constraint our formulation has more degree of freedom
- What happen when we drop this constraint
- The problem becomes non-convex and without a global minimum
- Not having a global minimum is bad, how can we solve that ?
- Adding an extra regularization parameter allow us to have a global minimum (a lot of them actually)
- Now we have a framework to switch on and of neurons

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

- We are still doing useless computation for neurons that will stay off forever, if we pruned them then we would speed up training
- Define the strategy to detect "forever dead" neurons
- (should we talk about the software architecture ?)
- now we evaluate the system
- First we show that our system has indeed more freedom than Group sparsity by showing that for any given amount of sparsity it can fit closer than the previous method
- On the same problem we show that removing neurons on the fly does not dramatically reduce accuracy (we might not be able to see any difference in performance for a linear/logistic regression though)
- We show that works with Bigger networks (multi layers perceptrons and CNNs), describe the training dynamics, discuss the shape we obtain
- Show that it is easier to do hyper-parameter optimization on the regularization parameter that the size of each layer. I think a nice experiment would be to plot the evolution of the final size of a two layer network (so going from a lambda to two sizes). And plot the trajectory in the "size space" as we decrease lambda. It would show how the system trades off the budget between the two layer in a nice and visual way
- Then we evaluate performance

2. Outline2

- Introduction
- Related Work
- Our Approach
- Experiments
- Discussion
- Conclusion

3. Introduction

3.1. ShrinkNets

3.2. Motivation

4. Our Approach

5. Theoretical analysis

5.1. Notations

In order to avoid any potential ambiguity, in this section we will describe in details the mathematical notations used in this article. Non-bold letters represent scalar values, while bold lowercase and upper case respectively denote vectors and matrices. \mathbf{A}^T stands for the transpose of the matrix \mathbf{A} . Subscripts are used to index particular elements of vectors and matrices. \mathbf{x}_i , \mathbf{A}_i , $(\mathbf{A}^T)_j$ and $\mathbf{A}_{i,j}$ respectively correspond to the i^{th} component of \mathbf{x} , the i^{th} row of \mathbf{A} , the j^{th} column of \mathbf{A} and the j^{th} component of the i^{th} row of \mathbf{A} . All the following definitions assume \mathbf{A} to be an $n \times p$ matrix. For any vector \mathbf{b} with n components, we define $\text{diag}(\mathbf{b})$ a $n \times n$ matrix such that: $\forall 1 \leq i \leq n$, $\text{diag}(\mathbf{b})_{i,i} = \mathbf{b}_i$ and 0 otherwise. For any $l \in [0, +\infty]$ we define the norm: $\|\mathbf{A}\|_l = \left(\sum_{i=1}^n \sum_{j=1}^p |\mathbf{A}_{i,j}|^l \right)^{\frac{1}{l}}$. For the rest of this paper and unless stated otherwise, \mathbf{y} will represent the output of a network, \mathbf{x} the input, \mathbf{b} a bias, λ regularization factors, Ω regularization methods, $\boldsymbol{\theta}$ general model parameters and a will stand for any element-wise activation function. The only constraint that we want to enforce is that $a(0) = 0$. We use $\llbracket u, v \rrbracket$ to denote interval of integers, $\mathbf{0}$ is the null vector (size depending on the context). $\#S$ is meant to represent the cardinality of a set S . To simplify the notation of function composition we use the following operator: $g(f(\mathbf{x})) = (f \circ g)(\mathbf{x})$ and for a long sequence of functions from f_1 to f_n we use: $f_n(\dots f_1(\mathbf{x})) = (\bigcirc_{k=1}^n f_k)(\mathbf{x})$.

5.2. Definition

$$\Omega_{\lambda}^s = \lambda \|\boldsymbol{\beta}\|_1 \quad (1)$$

Our goal when designing ShrinkNets was to be able to remove inputs and outputs of layers. For classic fully connected layers, which are defined as :

$$f_{\mathbf{A},\mathbf{b}}(\mathbf{x}) = a(\mathbf{A}\mathbf{x} + \mathbf{b}) \quad (2)$$

removing an input neuron j is equivalent to have $(\mathbf{A}^T)_j = \mathbf{0}$ and removing an output neuron i is the same as having $\mathbf{A}_i = \mathbf{0}$ and $\mathbf{b}_i = 0$. Solving optimization problems while trying to set entire groups of parameters to zero has been already studied and the most popular method is without doubt the group sparsity regularization [ref]. For any partition-

ning of the set of parameter defining a model in p groups: $\boldsymbol{\theta} = \bigcup_{i=1}^P \boldsymbol{\theta}_i$ we define it the following way:

$$\Omega_{\lambda}^{gp} = \lambda \sum_{i=1}^p \sqrt{\#\boldsymbol{\theta}_i} \|\boldsymbol{\theta}_i\|_2 \quad (3)$$

In the context of a fully-connected layer, the groups are either: columns of \mathbf{A} if we want to remove inputs, or rows of \mathbf{A} and the corresponding entry in \mathbf{b} if we want to remove outputs. For simplicity, we will focus our analysis in the simple one-layer case. In this case filtering outputs does not make a lot of sense, this is why we will only consider the former case. The group sparsity regularization then becomes:

$$\Omega_{\lambda}^{gp} = \lambda \sum_{j=1}^p \left\| (\mathbf{A}^T)_j \right\|_2 \quad (4)$$

Because $\forall i, \#\boldsymbol{\theta}_i = n$, To make the notation simpler, we now embed \sqrt{n} inside λ .

Since group sparsity and ShrinkNets try to achieve the same goal we will try to understand their similarities and differences. First let's recall the two problems. The ShrinkNet problem is:

$$\min_{\mathbf{A}, \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A} \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \Omega_{\lambda}^s \quad (5)$$

And the Group Sparsity problem is:

$$\min_{\mathbf{A}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \Omega_{\lambda}^{gp} \quad (6)$$

We can prove the under the condition: $\forall j \in \llbracket 1, p \rrbracket, \left\| (\mathbf{A}^T)_j \right\|_2 = 1$ the two problems are equivalent (??). However if we relax this constraint then shrinknet becomes non-convex and has no global minimum (????). Fortunately, by adding an extra term to the ShrinkNet regularization term we can prove that:

$$\min_{\mathbf{A}, \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A} \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \Omega_{\lambda}^s + \lambda_2 \|\mathbf{A}\|_p^p \quad (7)$$

has many global minimum (??) for all $p > 0$. This is the reason we define the *regularized ShrinkNet penalty*:

$$\Omega_{\lambda, \lambda_2, p}^{rs} = \lambda \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\theta}\|_p^p \quad (8)$$

In practice we observed that $p = 2$ or $p = 1$ are good choice, while the latter will also introduce additional sparsity in the parameters.

6. Evaluation

6.1. Multi-Target Linear and Multi-Class Logistic regressions

We will start evaluating ShrinkNets with the simplest problems possible: linear and logistic regression. As we showed, Group sparsity share similarities with our method, this is why we will use it as baseline. We decided to focus on multi-target linear regression because in the single target case, groups in the Group Sparsity problem would have a size of one (\mathbf{A} would be a vector in this case).

The evaluation will be done on two datasets `scml` and `oes97` [ref] for linear regressions and we will use `gina_prior2` [ref] and the *Gas Sensor Array Drift Dataset* [ref] (that we shorten in `gsadd`) for logistic regressions.

For each dataset we fit with different regularization parameters and measure the error and sparsity obtained after convergence. In this context we define sparsity as the ratio of columns that have all their weight under 10^{-3} in absolute value. Parameters were chosen in order to obtain the widest sparsity spectrum. Loss is normalized depending on the problem to be in the $[0, 1]$ range. We summarized the results in ???. From our experiments it is clear that ShrinkNets can fit the data closer than Group Sparsity for the same amount of sparsity. The fact that we are able to reach very low loss demonstrate that even if our objective function is non convex, in practice it works as good or better as convex alternatives. Details about these datasets and the parameters used are available in ???.

6.2. Neuron Removal strategies

In our previous experiment, we showed that the ShrinkNet loss was reasonable in practice, we are now interested in the impact on early pruning. The method we suggest for early pruning uses a parameter γ that control the aggressiveness of neuron removal so we will try to evaluate its impact on the final loss achieved by the model and the cost required to train the model. Our cost model is simple and hardware independant, we sum the number of input neurons at each epoch. In theory the cost in time should be asymptotically linear to this metric. To have a baseline we also train the same model but without neuron removal. Keep in mind that this is just in order to have some reference. Indeed, if we were to remove the neurons with small weights it would deteriorate the loss (and picking the threshold would be completely arbitrary). Therefore the baseline is evaluated with all neurons. One could consider it as a "theoretical lower bound" of the best achievable loss.

We picked multiple combinations of dataset and regularization parameters (λ) and for each we fit with different aggressiveness parameters (γ). We measure the loss after

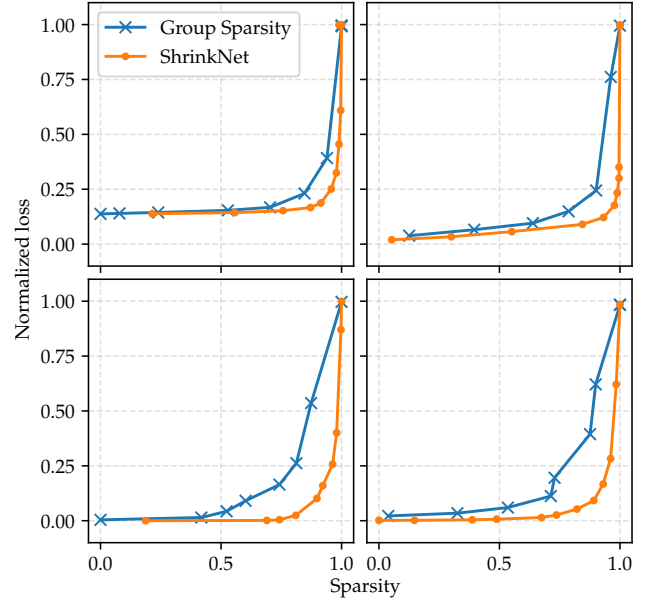


Figure 1. Loss/Sparsity trade off comparison between Group Sparsity and Shrinknet on linear and logistic regression. From top to bottom and left to right we show the results for `scml`, `oes97`, `gina_prior2` and `gsadd`.

convergence and the total cost and report the result in ???. In order to reduce the noise in the result, each experiment was performed 30 times and we display the range around ± 1 standard deviation.

TODO: Interpret the results

6.3. Convergence and training dynamics of Neural networks

6.4. Hyper-optimization of ShrinkNets

6.5. Performance

7. Speeding up training with pruning

8. Related Work

9. Discussion

10. Conclusion

A. Appendix

A.1. Proofs

Unless specified, all the proofs consider the Multi-Target linear regression problem

Proposition A.1. $\forall (n, p) \in \mathbb{N}_+^2, \mathbf{y} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{R}^p, \lambda \in \mathbb{R}$

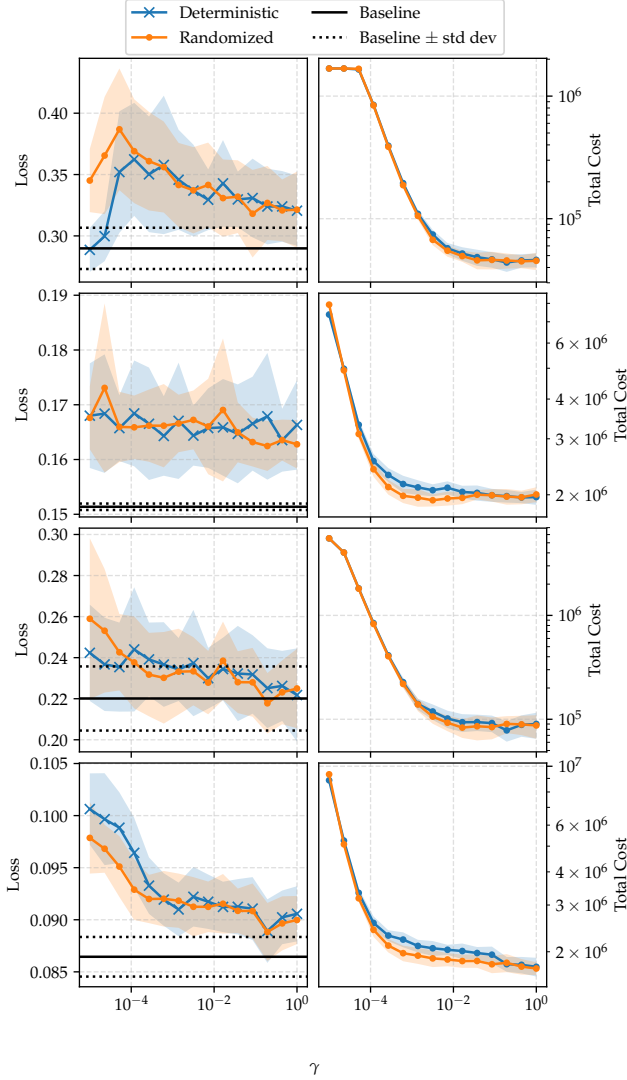


Figure 2. Effect of dynamic neuron removal for different γ . First column is the difference in the final loss in function of the removal factor. We plot theoretical baseline as a reference. Right column is a proxy of the total cost for training the model (i.e. the sum of input neurons at each epoch). Each row is a dataset/ λ combination. From top to bottom we have: scm1d/0.1, oes97/0.1

$$\begin{aligned} & \min_{\mathbf{A}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \sum_{j=1}^p \|(A^T)_j\|_2 \\ &= \min_{\mathbf{A}', \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A}' \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \\ & \text{s.t. } \forall j \in \llbracket 1, p \rrbracket, \|(A^T)_j\|_2^2 = 1 \end{aligned}$$

Proof. In order to prove this statement we will show that for any solution \mathbf{A} in the first problem, there exists a solution in the second with the exact same value, and vice-versa. We now assume we have a potential solution \mathbf{A} for the first problem and we define $\boldsymbol{\beta}$ such that $\beta_j = \|(A^T)_j\|_2^2$, and $\mathbf{A}' = \mathbf{A} (\text{diag}(\boldsymbol{\beta}))^{-1}$. It is easy to see that the constraint on \mathbf{A}' is satisfied by construction. Now:

$$\begin{aligned} & \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \sum_{j=1}^p \|(A^T)_j\|_2^2 \\ &= \|\mathbf{y} - \mathbf{A}' \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \lambda \sum_{j=1}^p \|(A^T)_j \beta_j\|_2^2 \\ &= \|\mathbf{y} - \mathbf{A}' \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \lambda \sum_{j=1}^p |\beta_j| \cdot 1 \\ &= \|\mathbf{y} - \mathbf{A}' \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \end{aligned}$$

Assuming we take an \mathbf{A}' that satisfy the constraint and a $\boldsymbol{\beta}$, we can define $\mathbf{A} = \mathbf{A}' \text{diag}(\boldsymbol{\beta})$. We can apply the same operations in reverse order and obtain an instance of the first problem with the same value. We can now see that the two problems must have the same minimum otherwise we would be able to construct a solution to the other with exact same value. \square

Proposition A.2.

$$\|\mathbf{y} - \mathbf{A} \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2$$

is not convex in \mathbf{A} and $\boldsymbol{\beta}$.

Proof. To prove this we will take the simplest instance of the problem: with only scalars. We have $f(a, \beta) = (y - a\beta x)^2$. For simplicity let's take $y = 1$ and $x > 0$. If we take two candidates $s_1 = (0, 2)$ and $s_2 = (2, 0)$, we have $f(s_1) = f(s_2) = 0$. However $f(\frac{2}{2}, \frac{2}{2}) = x > \frac{1}{2}f(0, 2) + \frac{1}{2}f(2, 0)$, which break the convexity property. Since we showed that a particular case of the problem is non-convex then necessarily the general cannot be convex. \square

Proposition A.3.

$$\min_{\mathbf{A}, \boldsymbol{\beta}} \|\mathbf{y} - \mathbf{A} \text{diag}(\boldsymbol{\beta}) \mathbf{x}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1$$

has no solution if $\lambda > 0$.

Proof. Let's assume this problem has a minimum A^*, β^* . Let's consider $2A^*, \frac{1}{2}\beta^*$. Trivially the first component of the sum is identical for the two solutions, however $\lambda \|\frac{1}{2}\beta\| < \lambda \|\beta\|$. Therefore A^*, β^* cannot be the minimum. We conclude that this problem has no solution. \square

Proposition A.4. *For this proposition we will not restrict ourselves to single layer but the composition of an arbitrary large (n) layers as defined individually as $f_{A_i, \beta_i, b_i}(x) = a(A_i \text{diag}(\beta_i)x + b_i)$. The entire network follows as: $N(x) = (\bigcirc_{i=1}^n f_{A_i, \beta_i, b_i})(x)$. For $\lambda > 0$, $\lambda_2 > 0$ and $p > 0$ we have:*

$$\min \|y - N(x)\|_2^2 + \Omega_{\lambda, \lambda_2, p}^{rs}$$

has at least 2^k global minimum where $k = \sum_{i=1}^n \#\beta_i$

Proof. First let's prove that there is at least one minimum to this problem. The two components of the expression are always positive so we know that this problem is bounded by below by 0. Let's assume this function does not have a minimum. Then there is a sequence of parameters $(S_n)_{n>0}$ such that the function evaluated at that point converges to the infimum of the problem. Since the function is defined everywhere does not have a minimum then this sequence must diverge. Since the entire sequence diverge there is at least one individual parameter that diverges. First case, the parameter is a component k of some β_i for some i . Necessarily $\|\beta_i\|_1$ diverge towards $+\infty$, which is incompatible with the fact that (S_n) converges to the infimum. We can have the exact same argument if the diverging parameter is in A_i or b_i because $p > 0$. Since there is always a contradiction then our assumption that the function has no global minimum must be false. Therefore, this problem has at least one global minimum.

Let's consider one optimal solution of the problem. For each component k of β_i for some i . Negating it and negating the k^{th} column of A_i does not change the the first part of the objective because the two factors cancel each other. The two norms do not change either because by definition the norm is independant of the sign. As a result these two sets of parameter have the same value and are both global minimum. It is easy to see that going from this global minimum we can decide to negate or not each element in each β_i . We have a binary choice for each parameter, there are $k = \sum_{i=1}^n \#\beta_i$ parameters, so we have at least 2^k global minima.

A.2. Experiment details

A.2.1. MULTI-TARGET LINEAR AND LOGISTIC REGRESSIONS

A.2.2. NEURON REMOVAL STRATEGIES

A.2.3. CONVERGENCE AND TRAINING DYNAMICS OF NEURAL NETWORKS

A.2.4. HYPER-OPTIMIZATION OF SHRINKNETS

A.2.5. PERFORMANCE

References

Vanschoren, Joaquin, van Rijn, Jan N., Bischl, Bernd, and Torgo, Luis. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013. doi: 10.1145/2641190.2641198. URL <http://doi.acm.org/10.1145/2641190.2641198>.

\square