

---

# Learning Network Size while Training with ShrinkNets

---

Anonymous Authors<sup>1</sup>

## Abstract

Let's write the abstract at the end

## 1. Tentative outline

- Introduction
  - Finding an appropriately sized network is challenging
  - ML practitioners spend a large amount of time tuning the size of layers in a network in order to get the best possible accuracy.
  - Many techniques have been proposed for hyperparam opt but these are computationally expensive and take long to train
  - We propose ShrinkNets, an approach to tune the size of the network as it is being trained without incurring the significant costs of hyperparameter optimization.
  - Thus, our contribution are: ...
- Our Approach
  - Assign an on/off switch to each neuron. But this is np-hard, so we consider a relaxation
  - Formulation
  - Relationship to sparsity
  - Strategies to kill neurons (relation to theory above?)
  - Implementation details
- Related Work
  - Hyperparameter optimization: random, bayesian opt, bandit methods
  - distillation techniques
  - post-training compression techniques
  - group sparsity, non-parametric neural networks

- training dynamics paper: first overfitting and then randomization?

- Experiments
  - Accuracy obtained by Shrinknets
  - Time taken to reach that accuracy compared with other hyperopt methods
  - Characterize method wrt params
  - Other experiments
- Discussion
  - Where does the method shine?
  - Side-effects of method: smaller networks, time to train
  - Potential extensions/limitations
- Conclusion

## 2. Introduction

One of the key determinants of network accuracy is the shape of the network, i.e., the number of layers, number of neurons per layer, and connections between layers. An under-sized network is likely to have low accuracy because of insufficient capacity while an over-sized network is difficult to train due to additional parameters and is computationally inefficient during training as well as inference. Consequently, many techniques have been proposed to determine the optimal size of a neural network. This problem is also referred to as *hyperparameter optimization*. The most popular techniques include hyperparameter optimization strategies ranging from random search (?), what-is-this-paper (?), meta-gradient descent (?), Gaussian processes (?) etc. All of these techniques without exception require a compute-intensive search of parameter space and often the training of many tens or hundreds of models. As a result, tuning models via hyperparameter optimization takes many times longer than the training of a single model and also requires more compute power for the training of every model with different model shape.

In this paper we present a novel method to automatically find an appropriate network size, drastically reducing optimization time. The key idea is to learn the right network size at the same time that the network is learning the main

---

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

task. Our strategy, called **ShrinkNets**, is blah blah **MV: Short blurb about the technique** Our approach has two main benefits. First and foremost, we no longer need to choose a network size before training. We can merely set an initial size for the network and then the algorithm will determine the best network that is smaller or equal in size to the initial size. Second, our technique trains a single model as opposed to the hundreds trained during hyperparameter optimization, as a result, we can find the best model faster as well as with less computational overhead.

Thus, our contributions are as follows: (a) we propose a novel technique based on dynamically switching neurons on/off to learn network size as the network is being trained. (b) we show that our technique is a relaxation **MV: ?** of group sparsity and prove **MV: fill in**. (c) we demonstrate the efficacy of our technique on CNNs and FCNs where ShrinkNets finds networks within  $\pm X\%$  of best hand-crafted accuracy in  $XX\%$  of training time compared to competing hyperparameter optimization methods. (d) we also demonstrate that ShrinkNets can achieve this accuracy with only  $YY\%$  of neurons.

### 3. Our Approach

- Assign an on/off switch to each neuron. But this is np-hard, so we consider a relaxation
- Formulation
- Relationship to group sparsity
- Strategies to kill neurons (relation to theory above?)
- Implementation details

Our approach has two main challenges. First, we need a way to dynamically size the network during training. Second, we need a loss function that optimizes for the additional task of sizing, without deteriorating the learning performance of the main task. Our approach, called ShrinkNets, copes with both challenges.

During training, our approach starts with an explicitly oversized network. As training progresses, we learn which neurons are not contributing to learning and remove them dynamically, effectively shrinking the network. This method requires two key components: first, we need a way to identify neurons that are not contributing to the learning process, and second we need a way to balance the network size and the generalization capability for the main task. We introduce a new Filter **MV: Choose a different name if possible bc this can be confused w/convolutional filters. Maybe Shrinkage layer?** layer that takes care of *deactivating* neurons. We also modify existing loss functions to incorporate a new term that takes care of balancing network size and generalization capability appropriately.

**Filter Layers:** Filter layers have weights in the range  $[0, +\infty]$  and are usually placed after linear and convolutional layers. 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, in the case of convolutional layers), 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  $\theta \in \mathbb{R}^C$ . The output is defined as follows:

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

Where  $\circ$  is the pointwise multiplication, and  $\theta$  is expanded in all dimensions to match the input size (except the second one since they are equal by definition). It is easy to see that if for any  $k$ , if  $\theta_k \leq 0$ , the  $k^{\text{th}}$  input feature/channel is multiplied by zero and have no influence on the output. If this happens, we say the Filter layer deactivates the neuron. These disabled neurons/channels can be removed from the network without changing its output. Before explaining how that is achieved, we explain next how the weights of the Filter Layer are initialized and adjusted during training.

**Training Procedure:** Once Filter layers are placed in a network and initialized (sampled from the Uniform $[0, 1]$  distribution), we could train the network directly using our standard loss function, and we could achieve performance 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*, thus forcing the deactivation of neurons. To obtain this sparsity, we simply redefine the loss function:

$$L'(x, y; \theta) = L(x, y) + \lambda |\max(0, \theta)| \quad (2)$$

The additional term  $\lambda |\max(0, \theta)|$  introduces sparsity (see Lasso loss (?)). The second component of the loss increases the gradient with respect to  $\theta$ , thus pushing its value towards zero. Neurons with little impact on the original loss (gradient lower than  $\lambda$ ), will not be able to compete against this attraction towards zero. Because the entries in  $\theta$  with a value of 0 or less correspond to dead neurons,  $\lambda$  effectively controls the number of neurons/channels in the entire network. We introduced the  $\max(\dots)$  into the loss to make sure that neurons are permanently disabled when performing gradient descent based optimization. Next, we explain how to implement ShrinkNets efficiently.

#### 3.1. Strategies to Remove Neurons

**Implementation** It is possible to reduce the overhead of the training process by removing neurons as soon as they become deactivated by  $\theta$  going to 0. To do this, we implemented a *neural garbage collection* mechanism which prunes deactivated neurons on-the-fly, reducing the processing time and memory overhead. To support this feature, it is

crucial to understand the information flow between neurons and layers in the neural network. We achieve this by representing such information flow as a graph. Vertices represent layers, and edges are event-hubs responsible for propagating information about disabled neurons to the relevant layers.

## 4. Theoretical analysis

- Comparison to group sparsity
- 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)

### 4.1. 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,  $\lambda$  regularization factors,  $\Omega$  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:  $f_n(\dots f_1(\mathbf{x})) = (\bigcirc_{k=1}^n f_k)(\mathbf{x})$ .

### 4.2. Definition

$$\Omega_\lambda^s = \lambda \|\boldsymbol{\beta}\|_1 \quad (3)$$

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 :

$$f_{\mathbf{A},\mathbf{b}}(\mathbf{x}) = a(\mathbf{A}\mathbf{x} + \mathbf{b}) \quad (4)$$

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 partitioning 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:

$$\Omega_\lambda^{gp} = \lambda \sum_{i=1}^p \sqrt{\#\boldsymbol{\theta}_i} \|\boldsymbol{\theta}_i\|_2 \quad (5)$$

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:

$$\Omega_\lambda^{gp} = \lambda \sum_{j=1}^p \left\| (\mathbf{A}^T)_j \right\|_2 \quad (6)$$

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

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:

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

And the Group Sparsity problem is:

$$\min_{\mathbf{A}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \Omega_\lambda^{gp} \quad (8)$$

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

(proposition A.1). However if we relax this constraint then shrinknet becomes non-convex and has no global minimum (propositions A.2 and A.3). Fortunately, by adding an extra term to the ShrinkNet regularization term we can prove that:

$$\min_{\mathbf{A}, \beta} \|\mathbf{y} - \mathbf{A} \text{diag}(\beta) \mathbf{x}\|_2^2 + \Omega_{\lambda}^s + \lambda_2 \|\mathbf{A}\|_p^p \quad (9)$$

has many global minimum (proposition A.4) for all  $p > 0$ . This is the reason we define the *regularized ShrinkNet penalty*:

$$\Omega_{\lambda, \lambda_2, p}^{rs} = \lambda \|\beta\|_1 + \lambda_2 \|\theta\|_p^p \quad (10)$$

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

## 5. Evaluation

- Accuracy obtained by Shrinknets
- Time taken to reach that accuracy compared with other hyperopt methods
- Characterize method wrt params (lambda, starting network size), training dynamics
- Other experiments: compression achieved by network, shape of network obtained
- Extra: Shrinknets have more freedom than group sparsity,

### 5.1. 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

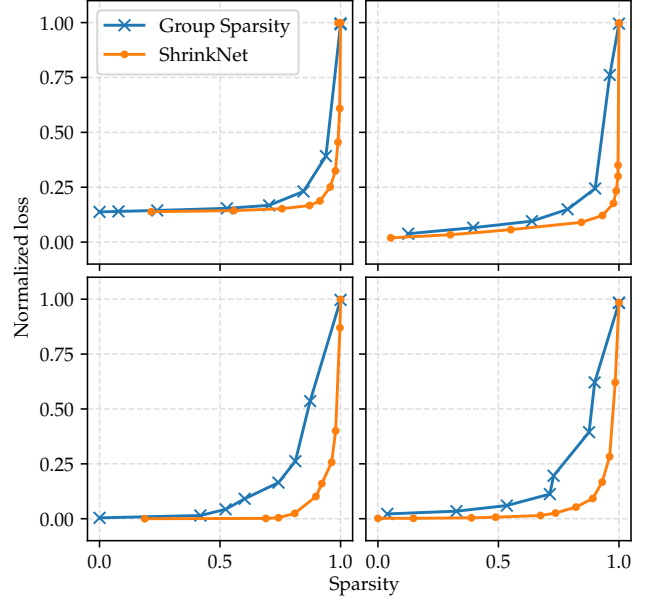


Figure 1. Loss/Sparsity trade off comparison between Group Sparsity and Shrinknet on linear and logistic regression. From top to bottom and left to right we show the results for `scml`, `oes97`, `gina_prior2` and `gsadd`.

widest sparsity spectrum. Loss is normalized depending on the problem to be in the  $[0, 1]$  range. We summarized the results in fig. 1. 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 appendix A.2.1.

### 5.2. 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  $\gamma$  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 ( $\lambda$ ) and for each we fit with different aggressiveness parameters ( $\gamma$ ). We measure the loss after convergence and the total cost and report the result in fig. 2. In order to reduce the noise in the result, each experiment was performed 30 times and we display the range around  $\pm 1$  standard deviation.

## TODO: Interpret the results

### 5.3. Convergence and training dynamics of Neural networks

### 5.4. Hyper-optimization of ShrinkNets

### 5.5. Performance

## 6. Speeding up training with pruning

## 7. Related Work

- Hyperparameter optimization: random, bayesian opt, bandit methods
- distillation techniques
- post-training compression techniques
- group sparsity, non-parametric neural networks
- training dynamics paper: first overfitting and then randomization?

## 8. Discussion

- Where does the method shine?
- Side-effects of method: smaller networks, time to train
- Potential extensions/limitations

## 9. Conclusion

## A. Appendix

### A.1. Proofs

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

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

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

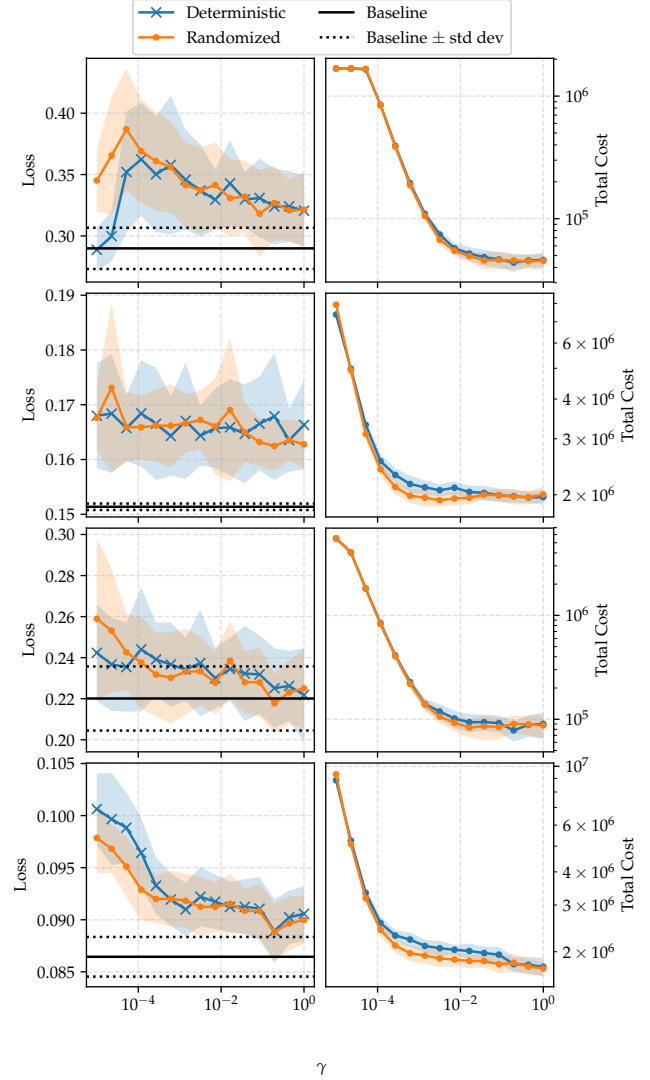


Figure 2. Effect of dynamic neuron removal for different  $\gamma$ . First column is the difference in the final loss in function of the removal factor. We plot theoretical baseline as a reference. Right column is a proxy of the total cost for training the model (i.e. the sum of input neurons at each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

*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  $\beta$  such that  $\beta_j = \|(A^T)_j\|_2^2$ , and  $\mathbf{A}' = \mathbf{A} (\text{diag}(\beta))^{-1}$ . It is easy to see that the constraint on  $\mathbf{A}'$  is satisfied by construction. Now:

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

Assuming we take an  $\mathbf{A}'$  that satisfy the constraint and a  $\beta$ , we can define  $\mathbf{A} = \mathbf{A}' \text{diag}(\beta)$ . 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.

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

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

*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.

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

has no solution if  $\lambda > 0$ .

*Proof.* Let's assume this problem has a minimum  $\mathbf{A}^*, \beta^*$ . Let's consider  $2\mathbf{A}^*, \frac{1}{2}\beta^*$ . Trivially the first component of the sum is identical for the two solutions, however  $\lambda \|\frac{1}{2}\beta^*\| < \lambda \|\beta^*\|$ . Therefore  $\mathbf{A}^*, \beta^*$  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  $f_{\mathbf{A}_i, \beta_i, \mathbf{b}_i}(x) = a(\mathbf{A}_i \text{diag}(\beta_i) \mathbf{x} + \mathbf{b}_i)$ . The entire network follows as:  $N(\mathbf{x}) = (\bigcirc_{i=1}^n f_{\mathbf{A}_i, \beta_i, \mathbf{b}_i})(\mathbf{x})$ . For  $\lambda > 0$ ,  $\lambda_2 > 0$  and  $p > 0$  we have:

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

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  $\beta_i$  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  $\mathbf{A}_i$  or  $\mathbf{b}_i$  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  $\beta_i$  for some  $i$ . Negating it and negating the  $k^{th}$  column of  $\mathbf{A}_i$  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  $\beta_i$ . 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.  $\square$

## 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>.