
Differentiable Sparsification for Deep Neural Networks

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep neural networks have relieved the feature engineering burden on human
2 experts. However, comparable efforts are instead required to determine an effective
3 architecture. In addition, as the sizes of networks have over-grown, a considerable
4 amount of resources is also invested in reducing the sizes. The sparsification of an
5 over-complete model addresses these problems as it removes redundant parameters
6 or connections. In this study, we propose a fully differentiable sparsification
7 method for deep neural networks, which allows parameters to be zero during
8 training with the stochastic gradient descent. Thus, the proposed method can
9 simultaneously learn the sparsified structure and weights of networks in an end-
10 to-end manner, which can be directly applied to modern deep neural networks
11 and imposes minimum overhead to the training process. To the authors' best
12 knowledge, it is the first fully [sub-]differentiable sparsification method that zeroes
13 out components, and it provides a foundation for future structure learning and
14 model compression methods.

15 1 Introduction

16 The success of deep neural networks has changed the paradigm of machine learning and pattern
17 recognition from feature engineering to architecture engineering [16, 14, 22, 7, 32]. Although deep
18 neural networks have relieved the burden of feature engineering, comparable human efforts are instead
19 required to determine an effective architecture, such as the number of neurons or layers and the
20 connections between nodes. In addition, as deep neural networks have over-grown (even up to 10–68
21 million parameters) [8, 10, 15, 32], considerable effort is also being invested in reducing existing
22 model sizes and in meeting the demands of deploying such networks on constrained platforms at
23 inference time [26, 25].

24 These problems can be addressed by the sparsification of an over-complete model [20]. A network
25 structure can be carved out of an over-complete model by removing redundant blocks [3, 29] or
26 deleting unnecessary connections between nodes or blocks [2, 18, 30], which also reduces the network
27 size. Among several approaches, pruning has long been adapted [17, 6, 26, 19, 5]. It typically requires
28 a pre-trained model and several steps (select unimportant parameters of a pre-trained model, delete
29 the parameters, and retrain the pruned model) and may repeat the process multiple times. Another
30 approach is a sparsity regularizer with the proximal gradient [21, 3, 29, 33] which shrinks redundant
31 parameters to zero during training and requires no pre-trained model. Among the most popular ones
32 is l_1 -regularizer [28]. However, as it acts on an individual parameter, it often produces unstructured
33 irregular models; thus, it diminishes the benefit of parallel hardware computation, such as GPUs [29].
34 In order to obtain regular sparse structures, a sparse regularization with l_2 -norm [3, 29] was adopted
35 on a group of parameters so that all parameters under the same group are either retained or zeroed-
36 out together. By zeroing-out parameters at a group level, the number of neurons or layers can be
37 automatically determined as a part of training. However, the optimization of a regularization term is
38 performed as a separate step separately from the gradient descent-based optimization for prediction

39 loss. The update rules should be implemented manually and the approach is limited to the cases
40 where closed form solutions for the proximal operation are known.

41 In this work, we propose a fully [sub-]differentiable sparsification method for deep neural networks,
42 which directly optimizes a regularized objective function and allows parameters to be exactly zero
43 during training with the stochastic gradient descent. Thus, it can simultaneously learn the sparsified
44 structure and weights of deep neural networks in an end-to-end manner. It leads to simpler implementa-
45 tion and it does not require to manually code a pruning step or an update rule like a soft-thresholding
46 operator. It can adopt various norms as a regularizer regardless of whether their closed form solutions
47 for the proximal operator are known or not. Another advantage of the proposed method is that it
48 can be easily applied on a group of parameters or a building block; thus, it can produce a structured
49 model and maximize the benefits of parallel hardware computation (e.g., GPUs) and it well suits the
50 trend of a modularized design in deep learning [22, 27, 7, 32].

51 2 Related Work

52 2.1 Proximal Gradient

53 Our proposed method is related to sparsity regularization with the proximal gradient [21]. A regular-
54 ized objective function is written as

$$\mathcal{L}(D, W) + \lambda \mathcal{R}(W), \quad (1)$$

55 where \mathcal{L} denotes a prediction loss, \mathcal{R} is a regularization term, D is a set of training data, W is a set
56 of model parameters, and λ controls the trade-off between prediction loss and model complexity. The
57 most popular regularizer is l_1 -norm,

$$\mathcal{R}(W) = \sum_i |w_i|, \quad (2)$$

58 where w_i is an individual element of W . To optimize the regularization term, parameter updating is
59 performed with a proximal operator,

$$w_i \leftarrow \text{sign}(w_i) (|w_i| - \eta\lambda)_+, \quad (3)$$

60 where \leftarrow denotes an assignment operator, η is a learning rate, and $(\cdot)_+$ represents $\max(\cdot, 0)$. As this
61 approach acts on an individual parameter, it often produces unstructured irregular models. In order
62 to obtain regular sparse structures, the sparse regularization with $l_{2,1}$ -norm can be adopted [3, 29].
63 All parameters in the same group are either retained or zeroed-out together. The regularization with
64 $l_{2,1}$ -norm is written as

$$\mathcal{R}(W) = \sum_g \|\mathbf{w}_g\|_2 = \sum_g \sqrt{\sum_i w_{g,i}^2}, \quad (4)$$

65 where $W = \{\mathbf{w}_g\}$ and \mathbf{w}_g represents a group of model parameters. The regularization term is
66 optimized with a proximal operator,

$$w_{g,i} \leftarrow \left(\frac{\|\mathbf{w}_g\|_2 - \eta\lambda}{\|\mathbf{w}_g\|_2} \right)_+ w_{g,i}. \quad (5)$$

67 When a group has only one single parameter, it degenerates to the l_1 -norm regularization.

68 Another group regularization is exclusive lasso with $l_{1,2}$ -norm [34, 33]. Rather than either retaining
69 or removing an entire group altogether, it promotes the competition or sparsity within a group. The
70 regularization term is written as

$$\mathcal{R}(W) = \frac{1}{2} \sum_g \|\mathbf{w}_g\|_1^2 = \frac{1}{2} \sum_g \left(\sum_i |w_{g,i}| \right)^2, \quad (6)$$

71 and its updating rule is derived as

$$w_{g,i} \leftarrow \text{sign}(w_{g,i}) \left(|w_{g,i}| - \eta\lambda \|\mathbf{w}_g\|_1 \right)_+. \quad (7)$$

72 These proximal operators consist of weight decaying and thresholding steps, and they are performed
 73 at every mini-batch or epoch in a sperate step after the optimization of a prediction loss. It requires
 74 some extra efforts to implement the update rules manually and applications are limited to the cases
 75 where closed form solutions for the proximal operator are known. In contrast, our approach optimizes
 76 a regularized objective function directly and allows parameters to be zero during training with the
 77 stochastic gradient descent. It leads to simpler implementations and it can employ various norms
 78 regardless of whether closed form solutions for proximal operators are known or not, for example
 79 p -norm with $p < 1$.

80 2.2 Differentiable Approach

81 Similar to our work, previous differentiable approaches [2, 18] learn the structure of a neural
 82 network by optimizing architecture parameters in a relaxed continuous domain, where architecture
 83 parameters represent the importance scores of building blocks or the connection strengths between
 84 them. However, as the architecture parameter magnitudes cannot be zero during training, the
 85 top k connections or components are stochastically or deterministically selected according to the
 86 architecture parameter values to derive a discretized architecture. Therefore, the approach may suffer
 87 from the discrepancy between a learned architecture and a final discretized one. Moreover, the
 88 value of k should be pre-specified manually; thus, the same value is set for all blocks or modules,
 89 which may be sub-optimal. Our approach drives the architecture parameters to zero by optimizing a
 90 regularized objective function. This can minimize the model discrepancy, and a network can choose
 91 the different numbers of components or connections in each module through training.

92 3 Proposed Approach

93 3.1 Base Model

94 We assume that there are n components in a module. A component can be any building block for a
 95 deep neural network or its output. For example, it can be a channel or a layer of convolutional neural
 96 networks such as ResNet [7, 8] and DenseNet layer [10]. It can also represent a node in a neural
 97 graph [30] or a convolutional neural network [11]. A module represents a composite of components,
 98 such as a group of channels or nodes. For illustration purposes, we assume that a module \mathbf{y} can be
 99 written as the linear combination of components \mathbf{f}_i :

$$\mathbf{y}(\mathbf{x}) = \sum_{i=1}^n a_i \mathbf{f}_i(\mathbf{x}; \mathbf{w}_i), \quad (8)$$

100 where \mathbf{x} denotes a module input, \mathbf{w}_i model parameters for component \mathbf{f}_i , and a_i an architecture
 101 parameter. Model parameters \mathbf{w}_i denote ordinary parameters, such as a filter in a convolutional layer
 102 or weight in a fully connected layer. The value of a_i represents the importance of component i , and it
 103 represents the connection strength between nodes in another context. Enforcing a_i to be zero amounts
 104 to removing component \mathbf{f}_i or zeroing-out \mathbf{w}_i . Thus, by creating the competition between elements of
 105 a and driving them to be zero, we can eliminate unnecessary components or connections. The sample
 106 model is simple, but we will show that it can be applied to various cases.

107 3.2 Differentiable Sparse Parameterization

108 First, we show how to parameterize architecture parameters with non-negative constraints, which is
 109 useful for the attention mechanism with the softmax. To set up the competition between the elements
 110 of a and allow them to be zero, we parameterize the architecture parameters as follows:

$$\gamma_i = \exp(\alpha_i) \quad (9)$$

$$\tilde{\gamma}_i = (\gamma_i - \sigma(\beta) \cdot \|\gamma\|_1)_+ \quad (10)$$

$$a_i = \frac{\tilde{\gamma}_i}{\sum_{j=1}^n \tilde{\gamma}_j}, \quad (11)$$

113 where α_i and β are unconstrained free parameters, $\sigma(\cdot)$ denotes a sigmoid function, and $(\cdot)_+$ rep-
 114 resents $relu(\cdot) = \max(\cdot, 0)$. When a parameter is non-negative, the proximal operator of Eq. (7)
 115 is reduced to Eq. (10). Although the forms are similar, they have completely different meanings.
 116 The proximal operator is a learning rule, whereas Eq. (10) is the parameterized form of architecture
 117 parameters, which is part of a neural network.

118 We can easily verify that a_i is allowed to be zero and is also differentiable from a modern deep
 119 learning perspective. The free parameters α_i and β are real-valued and they do not restrict a training
 120 process with the stochastic gradient descent. Thus, we can train a_i through α_i and β . The exponential
 121 function in Eq. (9) ensures that the architecture parameters are non-negative. Typically, a_i cannot
 122 be zero due to the exponential function of Eq. (9). However, $\tilde{\gamma}_i$ in Eq. (10) can be zero by the
 123 thresholding operation; hence, a_i can be zero as well. The term $\sigma(\beta) \cdot \|\gamma\|_1$ plays the role of a
 124 threshold, and the thresholding operation is interpreted as follows: if the strength of component i
 125 in a competition group is small compared to the total strength, it is dropped from the competition.
 126 Note that the scalar parameter β in Eq. (10), which determines the magnitude of a threshold, is
 127 not a hyper-parameter, but its value is automatically determined through training. Mathematically,
 128 the thresholding operator is not differentiable, but this should not pose an issue considering the
 129 support of $relu$ as a built-in differentiable function in a modern deep learning tool. Additionally, γ is
 130 non-negative; thus, its l_1 -norm is simply the sum of γ_i (i.e., $\|\gamma\|_1 = \sum \gamma_i$). The softmax of Eq. (11)
 131 is also differentiable. The softmax is optional but useful when we need to promote the competition
 132 between components as in the attention mechanism.

133 By relaxing the differentiability, singed architecture parameter can be similarly formed as

$$a_i = sign(\alpha_i) (\|\alpha_i\| - \sigma(\beta)\|\alpha\|_1)_+, \quad (12)$$

134 where α and β are free parameters. The gradient of the $sign$ function is zero almost everywhere, but
 135 it does not cause a problem for a modern deep learning tool. The equation can be rewritten as

$$a_i = \begin{cases} (\alpha_i + \sigma(\beta)\|\alpha\|_1)_- & \text{if } \alpha_i < 0 \\ (\alpha_i - \sigma(\beta)\|\alpha\|_1)_+ & \text{otherwise,} \end{cases}$$

136 where $(\cdot)_- = \min(\cdot, 0)$. The gradient can be computed separately according to whether its value
 137 is negative or not. To the authors' understanding, $sign$ function is already taken care in the ex-
 138 plained manner by TensorFlow [1], and thus $tf.math.sign()$ can be simply used without the manual
 139 implementation of the conditional statement.

140 3.3 Sparsity Regularizer

141 In the proposed approach, an objective function is written as

$$\mathcal{L}(D, W, a) + \lambda \mathcal{R}(a), \quad (13)$$

142 where a denotes the vector of architecture parameters. Sparsifying a is equivalent to sparsifying a
 143 deep neural network. Therefore, we can use the regularization term on a to encourage the sparsity of
 144 a . The proposed method is not limited to a particular norm and we can drive different sparsity patterns
 145 depending on the types of norms. For example, the most popular choice for parameter selection is
 146 l_1 -norm, but it is unsuitable on a in Eq. (11) as it is normalized using the softmax. Its l_1 -norm is
 147 always 1, i.e., $\|a\|_1 = \sum_{i=1}^n |a_i| = 1$. Therefore, we should employ p -norm with $p < 1$:

$$\mathcal{R}(a) = \left(\sum_{i=1}^n |a_i|^p \right)^{\frac{1}{p}} = \left(\sum_{i=1}^n a_i^p \right)^{\frac{1}{p}}, \quad (14)$$

148 where the second equality holds as a_i is always non-negative. To the authors' best knowledge, a closed
 149 form solution for the proximal operator of p -norm with $p < 1$ is unknown, but the regularization term
 150 is differentiable almost everywhere as $relu$ is. Thus, the proposed approach can directly optimize the
 151 regularized objective function and zero-out components with the stochastic gradient descent.

152 By simply switching one norm to another one for a regularizer, different sparsity patterns can be
 153 derived. For example, an individual component can be removed with l_1 -norm (Eq. 2) and a group of
 154 components or an entire module can be zeroed-out with $l_{2,1}$ -norm (Eq. 4). Note that we do not need
 155 to manually implement different updating rules as in the proximal gradient approach. We just need to
 156 rewrite a regularization term in the objective function. Examples codes and experiment results are
 157 shown in the supplementary material.

158 **3.4 Rectified Gradient Flow**

159 If γ_i Eq. (10) or α_i in Eq. (12) is less than the threshold, the gradient will be zero, and it will
 160 not receive learning signals. However, note that it does not necessarily mean that a component
 161 dies permanently once its importance score is less than the threshold. The component still has a
 162 chance to recover because the threshold is adjustable and the importance scores of other components
 163 may decrease. Nevertheless, to ensure that the architecture parameters of dropped components
 164 continuously receive learning signals, we propose approximating the gradient of the thresholding
 165 function. As in [31] where the gradient of a step function was approximated using that of *leaky relu*
 166 or *soft plus*, we suggest employing *elu* [4] as a variant of the proposed method: *relu* is used in the
 167 forward pass, but *elu* is used in the backward pass.

168 This heuristic approach leads to a similar learning mechanism proposed in [30], where the gradient
 169 flows to dropped (or zeroed out) edges but does not flow through these dropped edges. The architecture
 170 parameters in our proposed method correspond to the edges in [30]. Note that our approach can
 171 be easily implemented, and it does not require additional codes to control the gradient flow. The
 172 implementation codes are shown in the appendix.

173 **4 Application and Experiment**

174 In this section, we show that the proposed approach can be applied to reduce the size of a network as
 175 well as to learn the structure. Our aim is not to achieve state-of-the-art performance but to validate the
 176 idea and the broad applicability of the proposed approach. In order to show the broad applicability
 177 with limited computing resources, we perform experiments with relative small datasets such as
 178 CIFAR-10/100 [13]. Our implementations closely follow those of baseline models, including model
 179 structures and hyper-parameter settings.

180 **4.1 Channel Pruning in Convolutional Network**

Table 1: Performance on CIFAR-10, DenseNet-100-BC-K12

Model	Sparsity(%)		Top-1 Error(%)		Parmas		FLOPs	
	Avg.	Std.	Avg.	Std.	Avg.	Std.	Avg.	Std.
Base	00.0	0.0	5.44	0.11	7.6×10^5	0.0	5.8×10^8	0.0
NS	70.0	0.0	6.53	0.19	2.9×10^5	9.3×10^2	1.9×10^8	5.0×10^6
NS	80.0	0.0	8.39	0.28	2.0×10^5	2.0×10^3	1.4×10^8	3.4×10^6
DS	70.3	0.1	5.77	0.09	2.7×10^5	2.3×10^3	1.8×10^8	1.6×10^6
DS	80.4	0.1	6.64	0.11	1.7×10^5	0.6×10^3	1.3×10^8	2.0×10^6

181 Network-slimming(NS) [20] prunes unimportant channels in convolutional layers by leveraging the
 182 scaling factors in batch normalization. Let x_i and y_i be the input and output of batch normalization
 183 for channel i and then the operation can be written as

$$\tilde{x}_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}; y_i = a_i \tilde{x}_i + b_i,$$

184 where μ_i and σ_i denotes the mean and standard deviation of input activations, a and b are scale
 185 and shift parameters, ϵ is a small constant for numerical stability. The scaling parameter a can be
 186 considered as an importance score or an architecture parameter, and the affine transformation can be
 187 re-written as

$$y_i = a(\tilde{x}_i + b_i).$$

188 By pushing a_i to be zero, a corresponding channel can be removed. Network-slimming trains an
 189 initial network with l_1 -regularization on a to identify insignificant channels. After the training,
 190 channels with small values of a are pruned. To compensate the damage caused by pruning, a pruned
 191 network is fine-tuned. In our approach, we parameterize the scaling parameter using Eq. (12) and train
 192 a network with l_1 -norm on a using the stochastic gradient descent without pruning and fine-tuning.

193 We perform comparison experiments on CIFAR-10/100 [13]. The training and test sets contain
 194 50,000 and 10,000 samples respectively, and the final test error is reported at the end of training or

195 fine-tuning on all training images. We adopt a standard data augmentation scheme (random shifting
 196 and flipping) as in [7, 8]. In network-slimming, λ in Eq. 13 is fixed to 10^{-5} , but in our approach
 197 we vary its value to induce different level of sparsity. DenseNet-BC-K12 with 100 layers [10] and
 198 ResNet with 164 layers [7, 8] are employed as base networks. In network-slimming, pre-trained
 199 models are obtained by training networks for 160 epochs with the initial learning of 0.1. The learning
 200 rate is divided by 10 at 50% and 75% of the total number of training epochs. After the training,
 201 channels with small values of a are pruned and a slimmed network is fine-tuned for another 160
 202 epochs with the same setting as in the initial training, but learned weights are not re-initialized. In our
 203 approach, we train networks for 320 epochs without fine-tuning or re-training. Network-slimming
 204 initializes the scaling factor to be 0.5 and we set $\alpha = 0.5(n+1)/n$ and $\beta = \log(n^2 + n - 1)$ in
 205 Eq. 12 so that a starts with 0.5.

206 Table 1 shows experiment results on CIFAR-10 with DenseNet. More experiments, including ResNet
 207 and CIFAR-100, are given in the supplementary materials due to page limitation. The authors strongly
 208 urge readers to see the supplementary materials. We ran each experiments 5 times and showed the
 209 average and the standard deviation. Our proposed method is denoted by *Differentiable Sparsification*
 210 (DS). We controlled the value of λ such that similar pruning rate with that of the network-slimming
 211 approach. In the tables, sparsity denotes the number of removed channels in hidden layers. The
 212 experiments show that the proposed differentiable approach more effectively learns slimmed models.

213 4.2 Discovering Neural Wirings

Table 2: Performance on CIFAR-10, Discovering Neural Wirings

Model	Top-1 Error(%)		Params		Mult-Adds	
	Avg.	Std.	Avg.	Std.	Avg.	Std.
MobileNetV1($\times 0.25$)	13.44	0.24	2.2×10^5	0.0	3.3×10^6	0.0
No Update($\times 0.225$)	13.86	0.27	2.2×10^5	3.7×10^1	4.5×10^6	3.7×10^4
DNW($\times 0.225$)	10.30	0.20	1.8×10^5	6.7×10^1	3.1×10^6	4.6×10^4
PG- l_1 -norm	12.17	0.44	2.1×10^4	9.4×10^2	3.3×10^6	1.7×10^5
PG- $l_{1,2}$ -norm	13.62	0.56	9.6×10^4	1.6×10^4	3.4×10^6	8.6×10^4
DS-No Rectified Grad.	10.55	0.23	6.1×10^4	5.7×10^2	3.4×10^6	4.5×10^4
DS-Rectified Grad.	9.36	0.27	4.7×10^4	8.4×10^2	3.3×10^6	6.7×10^4

214 Discovering Neural Wirings(DNW) [30] relaxes the notion of layers and treats channels as nodes
 215 in a neural graph. By allowing channels to learn connections between them, it jointly discovers the
 216 structure and learns the parameters of a neural network. An input to node v , \mathbf{x}_v , is expressed

$$\mathbf{y}_v = \sum_{(u,v) \in \mathcal{E}} w_{u,v} \mathbf{x}_u,$$

217 where \mathbf{x}_u denotes the state of a proceeding node, \mathcal{E} represents a edge set and $w_{u,v}$ is a connection
 218 weight of an edge. The structure of a neural graph can be determined by choosing a subset of edges.

219 At each iteration of training, DNW chooses the top k edges with the highest magnitude, which is
 220 called a real edge set, and refers to the remaining edges as a hallucinated edge set. On the forwards
 221 pass or at inference time, real edges are only used. As DNW allows the magnitude of the weights
 222 in both sets to change throughout training, a hallucinated edge may replace a real edge when it
 223 strengthens enough. The weights of real edges are updated in an ordinary manner with the stochastic
 224 gradient descent, but those of hallucinated edges are updated by a specialized leaning rule: the
 225 gradient flows to hallucinated edges but does not flow through them.

226 The architecture parameters in our proposed method correspond to the edges in DNW. We parame-
 227 terize the edges using Eq. (12) and train a network with l_1 -norm on edges to induce sparsity. The
 228 rectified gradient leads to an update rule which is similar to that of DNW: the rectified gradient
 229 ensures that dropped architecture parameters continuously receive learning signals by approximating
 230 the gradient of the thresholding function. However, we do not need to keep track of the real and
 231 hallucinated edge sets. We simply optimize the objective function with approximated gradients. The
 232 rectified gradient can be implemented in a couple of lines using modern deep learning tools and the
 233 code is shown in the supplementary material.

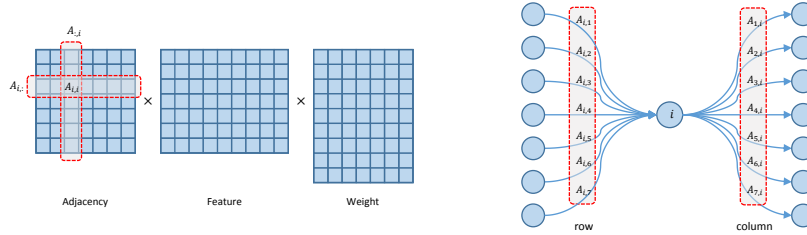


Figure 1: *Left*: Structure of a GCN block. Each block consists of a shared adjacency, an input feature, and a weight matrix. Each row and column of an adjacent matrix are treated as groups to enable the learning of relationship between a node (indexed by i) and its neighbors. *Right*: Row grouping creates the competition between in-coming nodes, and column grouping creates the competition between out-going nodes.

234 We perform experiments on CIFAR-10/100 [13]. The final test error is reported at the end of training
 235 without using separate validation data set. MobileNetV1 ($\times 0.25$) [9] is employed as a base model
 236 and our implementation closely follows that of DNW. We train for 160 epochs with initial learning
 237 rate 0.1. The learning rate is scheduled using Cosine Annealing. DNW chose the value of k such that
 238 a final learned model has similar Mult-Adds with the base model, and we also set the value of λ in
 239 the same manner.

240 Table 2 shows experiment results on CIFAR-10. We ran each experiments 5 times and showed the
 241 average and the standard deviation. Our proposed method is denoted by DS. DNW without the update
 242 rule corresponds to DS without the rectified gradient method. Even without the rectified gradient, the
 243 performance of the proposed method is close to that of DNW with the update rule. It validates the
 244 effectiveness of our approach. We also performed experiments with the proximal gradients of Eq.(3)
 245 and Eq.(7), which is denoted by PG in Table 2. PG- l_1 -norm also uses the l_1 -norm as a regularizer,
 246 but the learning is not as effective as ours. Similarly, PG- $l_{1,2}$ -norm uses the update rule of Eq.(7)
 247 whose shape is similar to our sparse parameterization Eq. (12), but the performance is worse than
 248 ours. More experiments, including CIFAR-100, are given in the supplementary materials due to page
 249 limitation.

250 DNW determines the size and the structure of a network by choosing k edges. However, there is no
 251 clear notion how to choose k for different stages (or blocks) and thus it uses the same pruning rate for
 252 all stages, which may be restrict because each stage may play a different role and need a different
 253 amount of resources. In contrast, our approach controls the model complexity by adjusting the value
 254 of λ in the objective function, and a different amount of resources is allocated for each stage through
 255 training. As shown in Table 2, the proposed method uses model parameters more efficiently.

256 4.3 Learning relationship between Nodes in Graph

257 In this section, we applied the proposed approach to learn the structure of an adjacency matrix in a
 258 graph convolutional network (GCN). The purpose of this case study is to test whether our approach
 259 can learn semantic structure from data rather than reducing the size of a neural network.

260 We adopted the model of [11], one of the most successful GCN models. A GCN block or a layer is
 261 defined (see Fig. 1) as

$$H^{l+1} = F(AH^lW^l),$$

262 where A is an adjacency matrix; H^l and W^l are an input feature and a weight matrix for layer l ,
 263 respectively; and F is a nonlinear activation function. In general, A is non-negative and shared
 264 across GCN blocks. It is obtained by normalization. For example, $A = \tilde{D}^{-1}\tilde{A}$ or $A = \tilde{D}^{-\frac{1}{2}}\tilde{A}\tilde{D}^{-\frac{1}{2}}$,
 265 where \tilde{A} is an unnormalized adjacency matrix; and \tilde{D} is a diagonal matrix, where $\tilde{D}_i = \sum_j \tilde{A}_{i,j}$.
 266 The adjacency matrix represents the connections or relationships between nodes on a graph and is
 267 usually given by prior knowledge. Learning the value of $A_{i,j}$ amounts to determining the relationship
 268 between nodes i and j . If the value of $A_{i,j}$ is zero, it can be considered that the two nodes are
 269 unrelated.

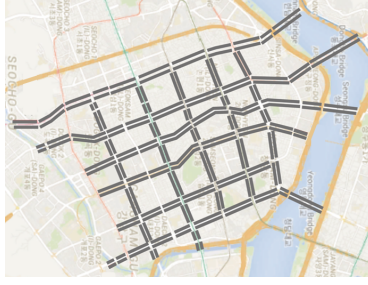


Figure 2: The gray lines represent 170 road links where the experimental data were collected.

270 As shown in Fig. 1, we defined each row and column as a group. Row grouping creates the competition
 271 between in-coming nodes, whereas column grouping creates the competition between out-going
 272 nodes. Each row and column of unnormalized adjacency matrix \tilde{A} can be parameterized similarly as
 273 in $\tilde{\gamma}$ of in Eq. (10):

$$274 \quad \gamma_{i,j} = \exp(\alpha_{i,j})$$

$$\tilde{A}_{i,j} = \left(\gamma_{i,j} - \sigma(\beta_i^r) \cdot \|\gamma_{i,:}\|_1 - \sigma(\beta_j^c) \cdot \|\gamma_{:,j}\|_1 \right)_+.$$

275 The softmax normalization of Eq. (11) is replaced with Sinkhorn normalization [23, 24, 12] to make
 276 A doubly-stochastic: each row and column sum up to 1. Initializing A with \tilde{A} , we can convert \tilde{A}
 277 into a doubly stochastic matrix by iteratively applying the following equations:

$$A = D_r^{-1} \tilde{A} \quad \text{and} \quad A = A D_c^{-1},$$

278 where D_r and D_c are diagonal matrices; $[D_r]_i = \sum_j A_{i,j}$; $[D_c]_j = \sum_i A_{i,j}$. Note that although the
 279 normalization is iterative, it is differentiable. Balanced normalization is also possible by iteratively
 280 applying

$$A = D_r^{-\frac{1}{2}} A D_c^{-\frac{1}{2}}.$$

281 We verified through numerical experiments that iteratively applying the above equation also makes \tilde{A}
 282 to doubly stochastic, but we could not find a theoretical justification. We leave the mathematical proof
 283 as an open question for a future work. As competition groups are created in row- and column-wise
 284 approaches, a regularized objective function can be written as

$$\mathcal{L}(D, W, A) + \frac{\lambda}{2} \sum_{i=1}^N \left\{ \mathcal{R}(A_{i,:}) + \mathcal{R}(A_{:,i}) \right\},$$

285 where $W = \{W^l\}$, N is the size of square matrix A , and $A_{i,:}$ and $A_{:,i}$ denote i th row and column
 286 vector of A , respectively. We employ l_p -norm of Eq. (14) with $p = 0.5$ for a regularizer.

287 To validate our purposed method, we applied a GCN to estimate future traffic speeds in a road network.
 288 The traffic speed data were collected from 170 road segments. Thus, the sizes of an adjacent matrix is
 289 170×170 . The map of the data area collected is shown in 2. One-step ahead observation is estimated
 290 from eight past observations and an output layer generates 170 estimates, one for each road segment.
 291 A prediction loss is measured using the mean relative error (MRE). More detailed specifications of
 292 the experimental data and our GCN model can be found in the supplementary material. A prediction
 293 loss is measured using the mean relative error (MRE).

294 Three baseline models were used: two were given the road connectivity and the other was not. For
 295 the first baseline, we set the value of $A_{i,j}$ as a constant such that $A_{i,j} = \frac{1}{n_i}$ if node (or road) i and j
 296 are adjacent to each other (n_i is the number of neighbors of node i), and $A_{i,j} = 0$ otherwise. The
 297 second baseline was taken in a similar approach, but we set $\tilde{A}_{i,j} = \exp(\alpha_{i,j})$ if node i and j were
 298 adjacent to each other to ensure that the strengths of the connections were learned. For the third
 299 baseline, the connectivity was not given. However, we set $\tilde{A}_{i,j} = \exp(\alpha_{i,j})$ for all i, j regardless of
 300 the actual connections. For the proposed method, we parameterized the adjacency matrix as in the

301 third baseline but applied the sparsification technique. The balanced normalization was applied to all
 302 cases except the first baseline, for which the row sum is 1.

303 To measure the learned relationship between nodes, we propose the following scoring function:

$$\frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^N \left[(A^r + A^c) \odot M^k \right]_{i,j},$$

304 where $A^r = D_r^{-1}A$, $A^c = AD_c^{-1}$, \odot denotes the element-wise product, and $[M^k]_{i,j} = 1$ if the
 305 geodesic distance between node i and j is less than or equal to k , whereas $[M^k]_{i,j} = 0$ otherwise.
 306 The maximum value is 1, and the minimum is 0. For example, the first and second baselines always
 307 have the maximum value because their adjacency matrices have exactly the same structure of M^1 .
 308 We calculated the scores for $k = 1$ and 2. Note that we used A^r and A^c instead of A because a
 309 sparsified matrix is not guaranteed to be doubly stochastic even if the original Sinkhorn normalization
 310 is adopted.

Table 3: Traffic speed prediction with GCN

Model	#N.Z.	MAPE(%)	L.R.(×100)		λ	#N.Z.	MAPE(%)	L.R.(×100)	
			$k = 1$	$k = 2$				$k = 1$	$k = 2$
I	878	5.6623	100.00	100.00	0.050	1,220	5.4744	87.06	89.94
II	878	5.5160	100.00	100.00	0.075	1,009	5.4957	88.62	91.39
III	28,900	5.6343	13.76	20.87	0.100	835	5.5336	89.79	92.13

(a) Baseline models

(b) Proposed method

311 The performance of the baseline models is shown in Table 4a. We ran each experiment five times and
 312 selected the median among the five lowest validation errors. For the first and second baselines, the
 313 road connectivity is given, and the number of non-zero elements of the adjacent matrices is 878. Note
 314 that the value of the learned relationship for Baselines I and II is constant, but we show it for reference.
 315 The road connectivity is not given for the third baseline, and the number of non-zero elements of
 316 the adjacent matrix is 28,900(= 170 × 170). The performance of the proposed model is shown in
 317 Table 4b. The experiment of Baseline III shows that a GCN finds a nonlinear mapping between input
 318 and target values simply in the way of reducing the prediction loss without learning the semantic
 319 relationships between nodes, but the proposed approach finds actual relationships between nodes. We
 320 further compared the proposed method with the proximal gradient method. The experimental results
 321 are reported in the supplementary material due to page limitation.

322 5 Scope and Limitation

323 Our aim is not to achieve state-of-the-art performance but to validate the idea and the broad applicabil-
 324 ity of the proposed approach. To the authors’ best knowledge, it is the first fully [sub-]differentiable
 325 sparsification method that zeroes out components, and we wish our work would provide a foundation
 326 for future structure learning and model compression methods. The limitation of our approach is that
 327 a sparsity rate cannot be explicitly specified before training as in conventional pruning approaches. If
 328 a specific sparsity rate is required, it should be obtained by a try-and-error. The rectified gradient is
 329 effective as shown in the experiments of discovering neural wirings, but it is not clear in which cases
 330 it is effective or not. We need more theoretical analysis and leave it for a future work.

331 6 Conclusion

332 In this study, we proposed a fully differentiable sparsification method that can simultaneously learn
 333 the sparsified structure and weights of deep neural networks. Our proposed method is versatile in that
 334 it can be seamlessly integrated into different types of neural networks and various problems.

335 Checklist

- 336 1. For all authors...
- 337 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
338 contributions and scope? [Yes]
- 339 (b) Did you describe the limitations of your work? [Yes] See the scope and limitation
340 section.
- 341 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 342 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
343 them? [N/A]
- 344 2. If you are including theoretical results...
- 345 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 346 (b) Did you include complete proofs of all theoretical results? [N/A]
- 347 3. If you ran experiments...
- 348 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
349 mental results (either in the supplemental material or as a URL)? [Yes] Main Codes are
350 included in the supplemental material and a snippet is also given in the appendix(pdf
351 supplemental material)
- 352 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
353 were chosen)? [Yes]
- 354 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
355 ments multiple times)? [Yes] We ran each experiment five times and gave the mean
356 and the standard deviation.
- 357 (d) Did you include the total amount of compute and the type of resources used (e.g.,
358 type of GPUs, internal cluster, or cloud provider)? [No] However, readers may easily
359 estimate required resources because we used well-known data and models in most
360 experiments.
- 361 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 362 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 363 (b) Did you mention the license of the assets? [N/A]
- 364 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 365
- 366 (d) Did you discuss whether and how consent was obtained from people whose data you're
367 using/curating? [N/A]
- 368 (e) Did you discuss whether the data you are using/curating contains personally identifiable
369 information or offensive content? [N/A]
- 370 5. If you used crowdsourcing or conducted research with human subjects...
- 371 (a) Did you include the full text of instructions given to participants and screenshots, if
372 applicable? [N/A]
- 373 (b) Did you describe any potential participant risks, with links to Institutional Review
374 Board (IRB) approvals, if applicable? [N/A]
- 375 (c) Did you include the estimated hourly wage paid to participants and the total amount
376 spent on participant compensation? [N/A]

377 References

- 378 [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro,
379 Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow,
380 Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser,
381 Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek
382 Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal
383 Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete
384 Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-
385 scale machine learning on heterogeneous systems, 2015. URL [https://www.tensorflow.](https://www.tensorflow.org/)
386 [org/](https://www.tensorflow.org/). Software available from tensorflow.org.

- 387 [2] Karim Ahmed and Lorenzo Torresani. Connectivity learning in multi-branch networks. *CoRR*,
388 abs/1709.09582, 2017. URL <http://arxiv.org/abs/1709.09582>.
- 389 [3] Jose M. Alvarez and Mathieu Salzmann. Learning the number of neurons in deep networks. In
390 *NIPS*, 2016.
- 391 [4] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network
392 learning by exponential linear units (elus). In *ICLR*, 2016.
- 393 [5] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable
394 neural networks. In *ICLR*, 2019.
- 395 [6] Babak Hassibi, David G. Stork, and Gregory J. Wolff. Optimal brain surgeon and general
396 network pruning. In *ICNN*, 1993.
- 397 [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
398 recognition. In *CVPR*, 2016.
- 399 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
400 networks. In *ECCV*, 2016.
- 401 [9] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias
402 Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural
403 networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017. URL <https://arxiv.org/abs/1704.04861>.
404
- 405 [10] Gao Huang, Zhuang Liu, and Laurens van der Maaten. Densely connected convolutional
406 networks. In *CVPR*, 2017.
- 407 [11] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
408 networks. In *ICLR*, 2017.
- 409 [12] P. A. Knight. The sinkhorn-knopp algorithm: convergence and applications. *SIAM Journal on*
410 *Matrix Analysis and Applications*, 30(1), 2008.
- 411 [13] Alex Krizhevsky. Learning multiple layers of features from tiny images. *Technical Report*,
412 2009.
- 413 [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep
414 convolutional neural networks. In *NIPS*, 2012.
- 415 [15] Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. Fractalnet: Ultra-deep neural
416 networks without residuals. In *ICLR*, 2017.
- 417 [16] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel.
418 Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):
419 541–551, Dec 1989. ISSN 0899-7667. doi: 10.1162/neco.1989.1.4.541.
- 420 [17] Yann LeCun, John S. Denker, and Sara A. Solla. Optimal brain damage. In *NIPS*, 1989.
- 421 [18] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In
422 *ICLR*, 2019.
- 423 [19] Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang.
424 Learning efficient convolutional networks through network slimming. In *ICCV*, 2017.
- 425 [20] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. Rethinking the value
426 of network pruning. In *ICLR*, 2019.
- 427 [21] Neal Parikh and Stephen Boyd. Proximal algorithms. *Found. Trends Optim.*, 1(3):127–239,
428 January 2014. ISSN 2167-3888. doi: 10.1561/24000000003. URL [http://dx.doi.org/10.](http://dx.doi.org/10.1561/24000000003)
429 [1561/24000000003](http://dx.doi.org/10.1561/24000000003).
- 430 [22] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
431 image recognition. In *ICLR*, 2015.

- 432 [23] R. Sinkhorn. A relationship between arbitrary positive matrices and doubly stochastic matrices.
433 *The annals of mathematical statistics*, 35(2):876–879, 1964.
- 434 [24] R. Sinkhorn and P. Knopp. Concerning nonnegative matrices and doubly stochastic matrices.
435 *Pacific Journal of Mathematics*, 21(2), 1967.
- 436 [25] Huizi Mao Song Han and William J. Dally. Deep compression: Compressing deep neural
437 network with pruning, trained quantization and Huffman coding. In *ICLR*, 2016.
- 438 [26] John Tran Song Han, Jeff Pool and William Dally. Learning both weights and connections for
439 efficient neural network. In *NIPS*, 2015.
- 440 [27] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov,
441 Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions.
442 In *CVPR*, 2016.
- 443 [28] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal*
444 *Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996.
- 445 [29] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity
446 in deep neural networks. In *NIPS*, 2016.
- 447 [30] Mitchell Wortsman, Ali Farhadi, and Mohammad Rastegari. Discovering neural wirings. In
448 *NeurIPS*, 2019.
- 449 [31] Xia Xiao, Zigeng Wang, and Sanguthevar Rajasekaran. Autoprune: Automatic network pruning
450 by regularizing auxiliary parameters. In *NeurIPS*, 2019.
- 451 [32] Saining Xie, Ross Girshick, Piotr Dollar, Zhuowen Tu, and Kaiming He. Aggregated residual
452 transformations for deep neural networks. In *CVPR*, 2017.
- 453 [33] Jaehong Yoon and Sung Ju Hwang. Combined group and exclusive sparsity for deep neural
454 networks. In *ICML*, 2017.
- 455 [34] Yang Zhou, Rong Jin, and Steven Chu–Hong Hoi. Exclusive lasso for multi-task feature
456 selection. In *AISTATS*, 2010.