

# Learning Local Receptive Fields and their Weight Sharing Scheme on Graphs

Jean-Charles Vialatte<sup>1,2</sup>, Vincent Gripon<sup>2</sup>, Gilles Coppin<sup>2</sup>

<sup>1</sup>City Zen Data

55 rue Charles Nungesser  
29490 Guipavas, France

jean-charles.vialatte@cityzendata.com

<sup>2</sup>IMT Atlantique / CNRS Lab-STICC

Technopole Brest Iroise  
29238 Brest, France

name.surname@imt-atlantique.fr

**Abstract**—We propose a simple and generic layer formulation that extends the properties of convolutional layers to any domain that can be described by a graph. Namely, we use the support of its adjacency matrix to design learnable weight sharing filters able to exploit the underlying structure of signals in the same fashion as for images. The proposed formulation makes it possible to learn the weights of the filter as well as a scheme that controls how they are shared across the graph. We perform validation experiments with image datasets and show that these filters offer performances comparable with convolutional ones.

**Index Terms**—deep learning, convolutional neural networks, local receptive fields, graph signal processing

## I. INTRODUCTION

Convolutional Neural Networks (CNNs) have achieved state-of-the-art accuracy in many supervised learning challenges [1], [2], [3], [4], [5], [6]. For their ability to absorb huge amounts of data with lesser overfitting, deep learning [7] models are the golden standard when a lot of data is available. CNNs benefit from the ability to create stationary and multi-resolution low-level features from raw data, independently from their location in the training images. Some authors draw a parallel between these features and scattering transforms [8].

Obviously CNNs rely on the ability to define a convolution operator (or a translation) on signals. On images, this amounts to learn local receptive fields [9] that are convolved with training images. Considering images to be defined on a grid graph, we point out that the receptive fields of vertices are included in their neighbors – or, more generally, a neighborhood.

Reciprocally, convolution requires more than the neighborhoods of vertices in the underlying graph, as the operator is able to match specific neighbors of distinct vertices together. For instance, performing convolution on images requires the knowledge of coordinates of pixels, that is not directly accessible when considering a grid graph (c.f. [10], [11]). In this paper we are interested in demonstrating that the underlying graph is nevertheless enough to achieve comparable results.

The convolution of a signal can be formalized as its multiplication with a convolution matrix. In the case of images and for small convolution kernels, it is interesting to note that this convolution matrix has the same support as a lattice graph. Using this idea, we propose to introduce a type of layer based on a graph that connects neurons to their neighbors. Moreover,

convolution matrices are entirely determined by a single row, since the same weights appear on each one. To imitate this process, we introduce a weight sharing learning procedure, that consists in using a limited pool of weights that each row of the obtained operator can make use of.

Section II presents related work. Section III describes our methodology and the links with existing architectures. Section IV contains experimental results. Section V is a conclusion.

## II. RELATED WORK

Due to the effectiveness of CNNs on image datasets, models have been proposed to adapt them to other kind of data, e.g. for shapes and manifolds [12], [13], molecular datasets [14], or graphs [15], [16], [13]. A review is done in [17]. In particular, CNNs have also been adapted to graph signals, such as in [18], [19] where the convolution is formalized in the spectral domain of the graph defined by its Laplacian [20]. This approach have been improved in [21], with a localized and fast approximated formulation, and has been used back in vision to breed isometry invariant representations [22].

For non-spectral approaches, feature correspondences in the input domain allow to define how the weights are tied across the layer, such as for images or manifolds. For graphs and graph signals, such correspondences doesn't necessarily exist. For example, in [16] (where the convolution is based on multiplications with powers of the probability transition matrix) weights are tied according to the power to which they are attached, in [15] an ordering of the nodes is used, in [13] an embedding is learned from the degrees of the nodes. These choices are arbitrary and unsimilar to what is done by regular convolutions. On the contrary, we propose a generic layer formulation that allows to also learn how the weights are linearly distributed over the local receptive field.

Our model is first designed for the task of graph signal classification, but another common task is the problem of node classification such as in [16], [23], [24], [13]. Models learning part of their structures have also been proposed, such as in [25], [26]. Moreover, because our model strongly resembles regular convolutions, it can also resemble some of their variants, such as group equivariant convolutions [27].

### III. METHODOLOGY

We first recall the basic principles of Deep Neural Networks (DNNs) and CNNs, then introduce our proposed graph layer.

#### A. Background

DNNs [28] consist of a composition of layers, each one parametrized by a learnable weight kernel  $W$  and a nonlinear function  $f : \mathbb{R} \rightarrow \mathbb{R}$ . Providing the input of such a layer is  $\mathbf{x}$ , the corresponding output is then:

$$\mathbf{y} = f(W \cdot \mathbf{x} + \mathbf{b}),$$

where  $\cdot$  is the matrix product operator,  $f$  is typically applied component-wise and  $\mathbf{b}$  is a learnable bias vector.

The weight kernels are learned using an optimization routine usually based on gradient descent, so that the DNN is able to approximate an objective function. A DNN containing only this type of layer is called Multi-Layer Perceptron (MLP).

In the case of CNNs [29], some of the layers have the particular form of convolution filters. In this case, the convolutional operation can also be written as the product of the input signal with a matrix  $W$ , where  $W$  is a Toeplitz matrix. Previous works [30], [31], [32], [33] have shown that to obtain the best accuracy in vision challenges, it is usually better to use very small kernels, resulting in a sparse  $W$ . Figure 1 depicts a convolutional layer.

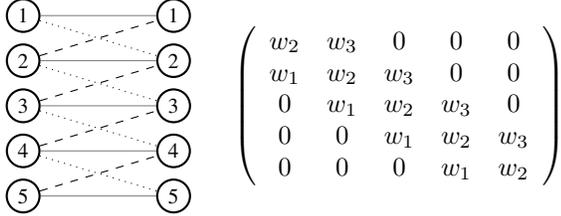


Figure 1. Depiction of a 1D-convolutional layer and its associated matrix  $W$ .

#### B. Proposed Method

We propose to introduce another type of layer, that we call *receptive graph layer*. It is based on an adjacency matrix and aims at extending the principle of convolutional layers to any domain that can be described using a graph.

Consider an adjacency matrix  $A$  that is well fitted to the signals to be learned, in the sense that it describes an underlying graph structure between the input features. We define the receptive graph layer associated with  $A$  using the product between a third rank tensor  $S$  and a weight kernel  $W$ . For now, the tensor  $W$  would be one-rank containing the weights of the layer and  $S$  is of shape  $n \times n \times \omega$ , where  $n \times n$  is the shape of the adjacency matrix and  $\omega$  is the shape of  $W$ .

On the first two ranks, the support of  $S$  must not exceed that of  $A$ , such that  $A_{ij} = 0 \Rightarrow \forall k, S_{ijk} = 0$ .

Overall, we obtain:

$$\mathbf{y} = f(W \cdot S \cdot \mathbf{x} + \mathbf{b}),$$

where here  $\cdot$  denotes the tensor product.

Intuitively, the values of the weight kernel  $W$  are linearly distributed to pairs of neighbours in  $A$  with respect to the values of  $S$ . For this reason, we call  $S$  the *scheme* (or *weight sharing scheme*) of the receptive graph. In a sense, this scheme tensor is to the receptive graph what the adjacency matrix is to the graph. An example is depicted in Figure 2.

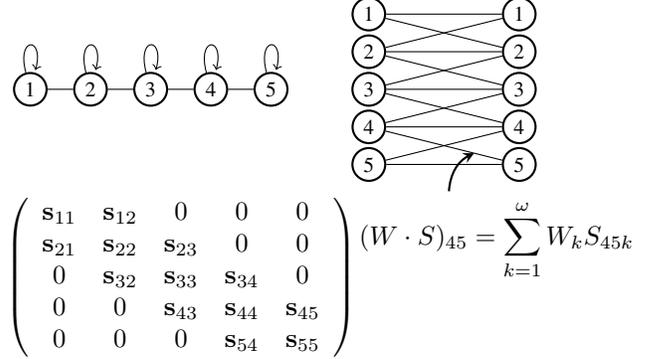


Figure 2. Depiction of a graph, the corresponding receptive graph of the propagation and its associated weight sharing scheme  $S$ . Note that  $s_{ij}$  are vector slices of  $S$  along the first two ranks,  $s_{ij}$  determines how much of each weight in  $W$  is allocated for the edge linking vertex  $i$  to vertex  $j$ .

Alike convolution on images,  $W$  is extended as a third-rank tensor to include multiple input and output channels (also known as feature maps). It is worth mentioning that an implementation must be memory efficient to take care of a possibly large sparse  $S$ .

#### C. Training

The proposed formulation allows to learn both  $S$  and  $W$ . We perform the two jointly. Learning  $W$  amounts to learning weights as in regular CNNs, whereas learning  $S$  amounts to learning how these weights are tied over the receptive fields. We also experiment a fine-tuning step, which consists in freezing  $S$  in the last epochs. Indeed, when a weight sharing scheme can be decided directly from the underlying structure, it is not necessary to train  $S$ .

Because of our inspiration from CNNs, we propose constraints on the parameters of  $S$ . Namely, we impose them to be between 0 and 1, and to sum to 1 along the third dimension. Therefore, the vectors on the third rank of  $S$  can be interpreted as performing a weighted average of the parameters in  $W$ .

We test two types of initialization for  $S$ . The first one consists in distributing one-hot-bit vectors along the third rank. We impose that for each receptive field, a particular one-hot-bit vector can only be distributed at most once more than any other. We refer to it as one-hot-bit initialization. The second one consists in using a uniform random distribution with limits as described in [34].

#### D. Genericity

For simplicity we restricted our explanation to square adjacency matrices. In the case of oriented graphs, one could remove the rows and columns of zeros and obtain a receptive

graph with a distinct number of neurons in the input ( $n$ ) than in the output ( $m$ ). As a result, receptive graph layers extend usual ones, as explained here:

- 1) To obtain a fully connected layer, one can choose  $\omega$  to be of size  $nm$  and  $S$  the matrix of vectors that contains all possible one-hot-bit vectors.
- 2) To obtain a convolutional layer, one can choose  $\omega$  to be the size of the kernel.  $S$  would be one-hot-bit encoded along its third rank and circulant along the first two ranks. A stride  $> 1$  can be obtained by removing the corresponding rows.
- 3) Similarly, most of the layers presented in related works can be obtained for an appropriate definition of  $S$ .

In our case,  $S$  is more similar to that obtained when considering convolutional layers, with the noticeable differences that we do not force which weight to allocate for which neighbor along its third rank and it is not necessarily circulant along the first two ranks.

### E. Discussion

Although we train  $S$  and  $W$ , the layer propagation is ultimately handled by their tensor product. That is, its output is determined by  $\Theta \cdot \mathbf{x}$  where  $\Theta = S \cdot W$ . For the weight sharing to make sense, we must then not over-parameterize  $S$  and  $W$  over  $\Theta$ . If we call  $l$  the number of non-zeros in  $A$  and  $w \times p \times q$  the shape of  $W$ , then the former assumption requires  $lw + wpq \leq lpq$  or equivalently  $\frac{1}{w} \geq \frac{1}{pq} + \frac{1}{l}$ . It implies that the number of weights per filter  $w$  must be lower than the total number of filters  $pq$  and than the number of edges  $l$ .

Note that without the constraint that the support of  $S$  must not exceed that of  $A$  (or if the used graph is complete), the proposed formulation could also be applied to structure learning of the input features space [35], [36]. That is, operations on  $S$  along the third rank might be exploitable in some way, e.g. dropping connections during training [37] or discovering some sort of structural correlations. However, even if this can be done for toy image datasets, such  $S$  wouldn't be sparse and would lead to memory issues in higher dimensions. So we didn't include these avenues in the scope of this paper.

## IV. EXPERIMENTS

### A. Description

We are interested in comparing various receptive graph layers with convolutional ones. For this purpose, we use image datasets, but restrain priors about the underlying structure.

We first present experiments on MNIST [38]. It contains 10 classes of gray levels images (28x28 pixels) with 60'000 examples for training, 10'000 for testing. We also do experiments on a scrambled version to hide the underlying structure, as done in previous work [39]. Then we present experiments on Cifar10 [40]. It contains 10 classes of RGB images (32x32 pixels) with 50'000 examples for training, 10'000 for testing.

Because receptive graph layers are wider than their convolutional counterparts ( $lw$  more parameters from  $S$ ), experiments are done on shallow (but wide) networks for this introductory paper. Also note that they require  $w + 1$  times more multiply

operations than a convolution lowered to a matrix multiplication [41]. In practice, they roughly took 2 to 2.5 more time.

### B. Experiments with grid graphs on MNIST

Here we use models composed of a single receptive graph (or convolutional) layer made of 50 feature maps, without pooling, followed by a fully connected layer of 300 neurons, and terminated by a softmax layer of 10 neurons. Rectified Linear Units [42] are used for the activations and a dropout [43] of 0.5 is applied on the fully-connected layer. Input layers are regularized by a factor weight of  $10^{-5}$  [44]. We optimize with ADAM [45] up to 100 epochs and fine-tune (while  $S$  is frozen) for up to 50 additional epochs.

We consider a grid graph that connects each pixel to itself and its 4 nearest neighbors (or less on the borders). We also use the square of this graph (pixels are connected to their 13 nearest neighbors, including themselves), the cube of this graph (25 nearest neighbors), up to 10 powers (211 nearest neighbors). Here we use one-hot-bit initialization. We test the model under two setups: either the ordering of the node is unknown, and then the one-hot-bit vectors are distributed randomly and modified upon training ; either an ordering of the node is known, and then the one-hot-bit vectors are distributed in a circulant fashion in the third rank of  $S$  which is freed in this state. We use the number of nearest neighbors as for the dimension of the third rank of  $S$ . We also compare with a convolutional layer of size 5x5, thus containing as many weights as the cube of the grid graph. Table I summarizes the obtained results. The ordering is unknown for the first result given, and known for the second result between parenthesis.

Table I  
ERROR RATES ON POWERS OF THE GRID GRAPHS ON MNIST.

Conv5x5	Grid <sup>1</sup>	Grid <sup>2</sup>	Grid <sup>3</sup>
(0.87%)	1.24% (1.21%)	1.02% (0.91%)	0.93% (0.91%)
Grid <sup>4</sup>	Grid <sup>5</sup>	Grid <sup>6</sup>	Grid <sup>10</sup>
0.90% (0.87%)	0.93% (0.80%)	1.00% (0.74%)	0.93% (0.84%)

We observe that even without knowledge of the underlying euclidean structure, receptive grid graph layers obtain comparable performances as convolutional ones, and when the ordering is known, they match convolutions. We also noticed that after training, even though the one-hot-bit vectors used for initialization had changed to floating point values, their most significant dimension was always the same. That suggests there is room to improve the initialization and the optimization.

In Figure 3, we plot the test error rate for various normalizations when using the square of the grid graph, as a function of the number of epochs of training. We observe that they have little influence on the performance and sometimes improve it a bit. Thus, we use them as optional hyperparameters.

### C. Experiments with covariance graphs on Scrambled MNIST

We use a thresholded covariance matrix obtained by using all the training examples. We choose the threshold so that the

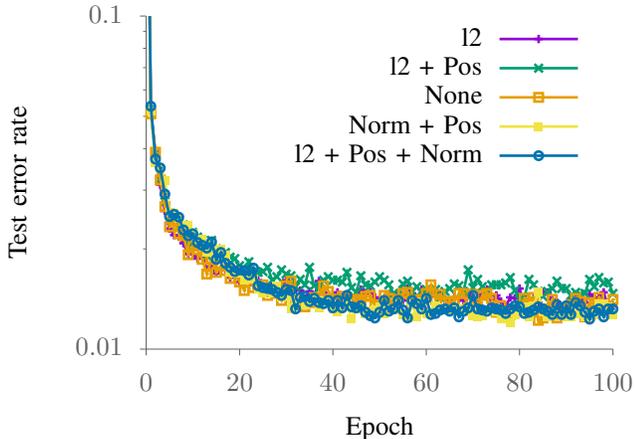


Figure 3. Evolution of the test error rate when learning MNIST using the square of a grid graph and for various normalizations, as a function of the epoch of training. The legend reads: “l2” means  $\ell_2$  normalization of weights is used (with weights  $10^{-5}$ ), “Pos” means parameters in  $S$  are forced to being positive, and “Norm” means that the  $\ell_1$  norm of each vector in the third dimension of  $S$  is forced to 1.

number of remaining edges corresponds to a certain density  $p$  ( $5 \times 5$  convolutions correspond approximately to a density of  $p = 3\%$ ). We also infer a graph based on the  $k$  nearest neighbors of the inverse of the values of this covariance matrix ( $k$ -NN). The latter two are using no prior about the signal underlying structure. The pixels of the input images are shuffled and the same re-ordering of the pixels is used for every image. Dimension of the third rank of  $S$  is chosen equal to  $k$  and its weights are initialized random uniformly [34]. The receptive graph layers are also compared with models obtained when replacing the first layer by a fully connected or convolutional one. Architecture used is the same as in the previous section. Results are reported on table II.

Table II  
ERROR RATES WHEN TOPOLOGY IS UNKNOWN ON SCRAMBLED MNIST.

MLP	Conv5x5	Thresholded ( $p = 3\%$ )	$k$ -NN ( $k = 25$ )
1.44%	1.39%	1.06%	0.96%

We observe that the receptive graph layers outperforms the CNN and the MLP on scrambled MNIST. This is remarkable because that suggests it has been able to exploit information about the underlying structure thanks to its graph.

#### D. Experiments with shallow architectures on Cifar10

On Cifar10, we made experiments on shallow CNN architectures and replaced convolutions by receptive graphs. We report results on a variant of AlexNet [3] using little distortion on the input that we borrowed from a tutorial of tensorflow [46]. It is composed of two  $5 \times 5$  convolutional layers of 64 feature maps, with max pooling and local response normalization, followed by two fully connected layers of 384 and

192 neurons. We compare two different graph supports: the one obtained by using the underlying graph of a regular  $5 \times 5$  convolution, and the support of the square of the grid graph. Optimization is done with stochastic gradient descent on 375 epochs where  $S$  is frozen on the 125 last ones. Circulant one-hot-bit initialization is used. These are weak classifiers for Cifar10 but they are enough to analyse the usefulness of the proposed layer. Exploring deeper architectures is left for further work. Experiments are run five times each. Means and standard deviations of accuracies are reported in table III. “Pos” means parameters in  $S$  are forced to being positive, “Norm” means that the  $\ell_1$  norm of each vector in the third dimension of  $S$  is forced to 1, “Both” means both constraints are applied, and “None” means none are used.

Table III  
ACCURACIES (IN %) OF SHALLOW NETWORKS ON CIFAR10.

Support	Learn $S$	None	Pos	Norm	Both
Conv5x5	No	/	/	/	$86.8 \pm 0.2$
Conv5x5	Yes	$87.4 \pm 0.1$	$87.1 \pm 0.2$	$87.1 \pm 0.2$	$87.2 \pm 0.3$
Grid <sup>2</sup>	Yes	$87.3 \pm 0.2$	$87.3 \pm 0.1$	$87.5 \pm 0.1$	$87.4 \pm 0.1$

The receptive graph layers are able to outperform the corresponding CNNs by a small amount in the tested configurations, opening the way for more complex architectures.

## V. CONCLUSION

We introduced a new class of layers for deep neural networks which consists in using the support of a graph operator and linearly distributing a pool of weights over the defined edges. The linear distribution is learned jointly with the pool of weights. Thanks to these structural dependencies, we showed it is possible to share weights in a fashion similar to Convolutional Neural Networks (CNNs).

We performed experiments on vision datasets where the receptive graph layer obtains similar performance as convolutional ones, even when the underlying image structure is hidden. We believe that with further work, the proposed layer could fully extend the performance of CNNs to many other domains described by a graph.

Future works will also include exploration of more advanced graph inference techniques. One example is using gradient descent from the supervised task at hand [19]. We can also notice that in our case, this amounts to select receptive fields, breeding another avenue [47].

## ACKNOWLEDGMENTS

This work was funded in part by the CominLabs project Neural Communications, and by the ANRT (Agence Nationale de la Recherche et de la Technologie) through a CIFRE (Convention Industrielle de Formation par la REcherche).

## REFERENCES

- [1] Y. LeCun, K. Kavukcuoglu, and C. Farabet, "Convolutional networks and applications in vision," in *Circuits and Systems (ISCAS), Proceedings of 2010 IEEE International Symposium on*. IEEE, 2010, pp. 253–256.
- [2] D. Cireřan, U. Meier, J. Masci, and J. Schmidhuber, "A committee of neural networks for traffic sign classification," in *Neural Networks (IJCNN), The 2011 International Joint Conference on*. IEEE, 2011, pp. 1918–1921.
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [4] T. N. Sainath, A.-r. Mohamed, B. Kingsbury, and B. Ramabhadran, "Deep convolutional neural networks for lvcsr," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 8614–8618.
- [5] C. Farabet, C. Couprie, L. Najman, and Y. LeCun, "Learning hierarchical features for scene labeling," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1915–1929, 2013.
- [6] J. J. Tompson, A. Jain, Y. LeCun, and C. Bregler, "Joint training of a convolutional network and a graphical model for human pose estimation," in *Advances in neural information processing systems*, 2014, pp. 1799–1807.
- [7] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [8] S. Mallat, "Understanding deep convolutional networks," *Philosophical Transactions of the Royal Society A*, vol. 374, no. 2065, 2016.
- [9] J. Moody and C. Darken, *Learning with localized receptive fields*. Yale Univ., Department of Computer Science, 1988.
- [10] N. Grelier, B. Pasdeloup, J.-C. Vialatte, and V. Gripon, "Neighborhood-preserving translations on graphs," in *Proceedings of IEEE GlobalSIP*, 2016, pp. 410–414.
- [11] J.-C. Vialatte, V. Gripon, and G. Mercier, "Generalizing cnns for data structured on locations irregularly spaced out," *arXiv preprint arXiv:1606.01166*, 2016.
- [12] J. Masci, D. Boscaini, M. Bronstein, and P. Vandergheynst, "Shapenet: Convolutional neural networks on non-euclidean manifolds," 2015.
- [13] F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. M. Bronstein, "Geometric deep learning on graphs and manifolds using mixture model cnns," *arXiv preprint arXiv:1611.08402*, 2016.
- [14] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," in *Advances in Neural Information Processing Systems*, 2015, pp. 2215–2223.
- [15] M. Niepert, M. Ahmed, and K. Kutzkov, "Learning convolutional neural networks for graphs," in *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, 2016, pp. 2014–2023.
- [16] J. Atwood and D. Towsley, "Diffusion-convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 1993–2001.
- [17] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, "Geometric deep learning: going beyond euclidean data," *arXiv preprint arXiv:1611.08097*, 2016.
- [18] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.
- [19] M. Henaff, J. Bruna, and Y. LeCun, "Deep convolutional networks on graph-structured data," *arXiv preprint arXiv:1506.05163*, 2015.
- [20] D. I. Shuman, S. K. Narang, P. Frossard, A. Ortega, and P. Vandergheynst, "The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains," *IEEE Signal Processing Magazine*, vol. 30, pp. 83–98, 2013.
- [21] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," in *Advances in Neural Information Processing Systems*, 2016, pp. 3837–3845.
- [22] R. Khasanova and P. Frossard, "Graph-based isometry invariant representation learning," *arXiv preprint arXiv:1703.00356*, 2017.
- [23] Z. Yang, W. W. Cohen, and R. Salakhutdinov, "Revisiting semi-supervised learning with graph embeddings," *arXiv preprint arXiv:1603.08861*, 2016.
- [24] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [25] J. Feng and T. Darrell, "Learning the structure of deep convolutional networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 2749–2757.
- [26] C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, and S. Yang, "Adanet: Adaptive structural learning of artificial neural networks," *arXiv preprint arXiv:1607.01097*, 2016.
- [27] T. Cohen and M. Welling, "Group equivariant convolutional networks," in *International Conference on Machine Learning*, 2016, pp. 2990–2999.
- [28] K. Hornik, M. Stinchcombe, and H. White, "Multilayer feedforward networks are universal approximators," *Neural Networks*, vol. 2, no. 5, pp. 359–366, 1989.
- [29] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [30] B. Graham, "Spatially-sparse convolutional neural networks," *arXiv preprint arXiv:1409.6070*, 2014.
- [31] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [32] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014.
- [33] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [34] X. Glorot and Y. Bengio, "Understanding the difficulty of training deep feedforward neural networks." 2010.
- [35] T. Richardson, "A discovery algorithm for directed cyclic graphs," in *Proceedings of the Twelfth international conference on Uncertainty in artificial intelligence*. Morgan Kaufmann Publishers Inc., 1996, pp. 454–461.
- [36] T.-Y. Kwok and D.-Y. Yeung, "Constructive algorithms for structure learning in feedforward neural networks for regression problems," *IEEE Transactions on Neural Networks*, vol. 8, no. 3, pp. 630–645, 1997.
- [37] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Advances in Neural Information Processing Systems*, 2015, pp. 1135–1143.
- [38] Y. LeCun, C. Cortes, and C. J. Burges, "The mnist database of handwritten digits," 1998.
- [39] X. Chen, X. Cheng, and S. Mallat, "Unsupervised deep haar scattering on graphs," in *Advances in Neural Information Processing Systems*, 2014, pp. 1709–1717.
- [40] A. Krizhevsky, "Learning multiple layers of features from tiny images," 2009.
- [41] S. Chetlur, C. Woolley, P. Vandermersch, J. Cohen, J. Tran, B. Catanzaro, and E. Shelhamer, "cudnn: Efficient primitives for deep learning," *arXiv preprint arXiv:1410.0759*, 2014.
- [42] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *International Conference on Artificial Intelligence and Statistics*, 2011, pp. 315–323.
- [43] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [44] A. Y. Ng, "Feature selection, l1 vs. l2 regularization, and rotational invariance," in *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004, p. 78.
- [45] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [46] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, "TensorFlow: Large-scale machine learning on heterogeneous systems," 2015. [Online]. Available: <http://tensorflow.org/>
- [47] A. Coates and A. Y. Ng, "Selecting receptive fields in deep networks," in *Advances in Neural Information Processing Systems*, 2011, pp. 2528–2536.