

# Deep Networks on Graph-Structured Data

Anonymous Author(s)

Affiliation

Address

email

## Abstract

## 1 Introduction

In recent times, deep learning models have proven extremely successful in a wide variety of tasks, from computer vision and acoustic modeling, to natural language processing [8]. At the core of their success lies an important assumption on the statistical properties of the data, namely the *stationarity* and the *compositionality* through local statistics, which is present in natural images, video, or speech. These properties are exploited efficiently by ConvNets [7, 6], which are designed to extract local features that are shared across the signal domain. Thanks to this, they are able to greatly reduce the number of parameters in the network with respect to generic Deep architectures, without sacrificing the capacity to extract informative statistics from the data. Similarly, Recurrent Neural Nets (RNNs) trained on temporal data implicitly assume a stationary distribution.

One can think such data as being signals defined on a low-dimensional grid, where the stationarity is well defined via the natural translation operator on the grid, locality is defined via the metric of the grid, and compositionality is obtained from downsampling. However, many types of data are defined on more complex graphs. For example, text documents represented as bags of words can be thought of as signals defined on a graph whose nodes are vocabulary terms and whose weights represent some similarity measure between terms, such as co-occurrence statistics. In medicine, a patient's gene expression data can be viewed as a signal defined on the graph imposed by the regulatory network. In fact, computer vision and audio, which are the main focus of research efforts in deep learning, only represent a special case of data defined on an extremely simple low-dimensional graph. For those type of data of dimension  $n$ , deep learning strategies are reduced to learning with fully-connected layers, which have  $o(n^2)$  parameters, and regularization is carried out via weight decay and dropout [14].

When the graph structure of the input is known, [2] introduced a model to generalize ConvNets using low learning complexity, similar to what a convnet uses, that was shown on simple low-dimensional graphs. In this work, we are interested in generalizing ConvNets to more general data distributions, and, most importantly, to the setting where the graph structure is not known a priori. In this context, learning the graph structure amounts to estimating the similarity matrix, which has complexity  $o(n^2)$ . One may therefore wonder whether the graph estimation followed by graph convolutions offers advantages with respect to learning directly from the data with fully connected layers. We attempt to answer this question experimentally as well as mathematically.

We explore these approaches in two areas of application for which it has not been possible to apply convolutional networks before: text categorization and bioinformatics. Our results show that our method is capable of matching or outperforming large, fully-connected networks trained with dropout using fewer parameters. Our main contributions can be summarized as follows:

- We extend the ideas from [2] to large scale classification problems, on Object Recognition, text categorization and bioinformatics.

- We consider the most general setting where no prior information on the graph structure is available, and propose unsupervised and supervised graph estimation strategies.
- Finally, we introduce an alternative formulation for efficient learning in graph-structured domains which works directly in the feature domain.

The rest of the paper is structured as follows. Section 2 reviews similar works in the literature. Section 3 discusses generalizations of convolutions on graphs, and Section 4 addresses the question of graph estimation. Finally, Section ?? shows numerical experiments on large scale object recognition, text categorization and bioinformatics.

## 2 Related Work

There have been several works which have explored architectures using the so-called local receptive fields [5, 4, 11], mostly with applications to image recognition. In particular, [4] proposes a scheme to learn how to group together features based upon a measure of similarity that is obtained in an unsupervised fashion. However, it does not attempt to exploit any weight-sharing strategy.

Recently, [2] proposed a generalization of convolutions to graphs via the Graph Laplacian. By identifying a linear, translation invariant operator in the grid (the Laplacian operator), with its counterpart in a general graph (the Graph Laplacian), one can view convolutions as the family of linear transforms commuting with the Laplacian. By combining this commutation property with a rule to find localized filters, the model requires only  $o(1)$  parameters per “feature map”. However, this construction requires prior knowledge of the graph structure, and was shown only on simple, low-dimensional graphs. More recently, [10] introduced Shapenet, another generalization of Convolutions on non-Euclidean domains, based on geodesic polar coordinates, which was successfully applied to shape analysis, and allows comparison across different manifolds. However, it also requires prior knowledge of the manifolds.

The graph or similarity estimation aspects have also been extensively studied in the past. For instance, [12] studies the estimation of the graph from a statistical point of view, through the identification of a certain graphical model using  $\ell_1$  penalized logistic regression. Also, [3] considers the problem of learning a deep architecture through a series of Haar contractions, which are learnt using an unsupervised pairing criteria over the features.

## 3 Generalizing Convolutions in Graphs

### 3.1 Spectral Networks

Our work builds upon [2] which introduced spectral networks. We recall the definition here and its main properties.

A spectral network generalizes a convolutional network through the Graph Fourier Transform, which is in turn defined via a generalization of the Laplacian operator on the grid to the graph Laplacian. An input vector  $x \in \mathbb{R}^N$  is seen as a signal defined on a graph  $G$  with  $N$  nodes.

**Definition 1.** Let  $W$  be a  $N \times N$  similarity matrix representing an undirected graph  $G$ , and let  $L = D - W$  be its graph Laplacian with  $D = W \cdot \mathbf{1}$ , with eigenvectors  $U = (u_1, \dots, u_N)$ . Then a graph convolution of input signals  $x$  with filters  $g$  on  $G$  is defined by  $x *_G g = U^T (Ux \odot Ug)$ , where  $\odot$  represents a point-wise product.

Here, the unitary matrix  $U$  plays the role of the Fourier Transform in  $\mathbb{R}^d$ . There are several ways of computing the graph Laplacian  $L$  [1]. In this paper, we choose the normalized version  $L = I - D^{-1/2} W D^{-1/2}$ , where  $D$  is a diagonal matrix with entries  $D_{ii} = \sum_j W_{ij}$ . Note that in the case where  $W$  represents the lattice, from the definition of  $L$  we recover the discrete Laplacian operator  $\Delta$ . Also note that the Laplacian commutes with the translation operator, which is diagonalized in the Fourier basis. It follows that the eigenvectors of  $\Delta$  are given by the Discrete Fourier Transform (DFT) matrix. We then recover a classical convolution operator that noting that convolutions are by definition linear operators that diagonalize in the Fourier domain (also known as the Convolution Theorem [9]).

Learning filters in a Graph thus amounts to learning spectral multipliers  $w_g = (w_1, \dots, w_N)$

$$x *_G g := U^T (\text{diag}(w_g) U x) .$$

Extending the convolution to inputs  $x$  with multiple input channels is straightforward. If  $x$  is a signal with  $M$  input channels and  $N$  locations, we apply the transformation  $U$  on each channel, and then use multipliers  $w_g = (w_{i,j}; i \leq N, j \leq M)$ .

However, for each feature map  $g$  we need convolutional kernels are typically restricted to have small spatial support, independent of the number of input pixels  $N$ , which enables the model to learn a number of parameters independent of  $N$ . In order to recover a similar learning complexity in the spectral domain, it is thus necessary to restrict the class of spectral multipliers to those corresponding to localized filters.

For that purpose, we seek to express spatial localization of filters in terms of their spectral multipliers. In the grid, smoothness in the frequency domain corresponds to the spatial decay, since

$$\left| \frac{\partial^k \hat{x}(\xi)}{\partial \xi^k} \right| \leq C \int u^k |x(u)| du ,$$

where  $\hat{x}(\xi)$  is the Fourier transform of  $x$ . In [2] it was suggested to use the same principle in a general graph, by considering a smoothing kernel  $\mathcal{K} \in \mathbb{R}^{N \times N_0}$ , such as splines, and searching for spectral multipliers of the form

$$w_g = \mathcal{K} \tilde{w}_g .$$

The algorithm which implements the graph convolution is described in 1.

---

**Algorithm 1** Train Graph Convolution Layer

---

- 1: Given GFT matrix  $U$ , interpolation kernel  $K$ , weights  $w$ .
  - 2: **Forward Pass:**
  - 3: Fetch input batch  $x$  and gradients w.r.t outputs  $\nabla y$ .
  - 4: Compute interpolated weights:  $w_{f'f} = \mathcal{K} w_{f'f}$ .
  - 5: Compute output:  $y_{sf'} = U^T \left( \sum_f U x_{sf} \odot w_{f'f} \right)$ .
  - 6: **Backward Pass:**
  - 7: Compute gradient w.r.t input:  $\nabla x_{sf} = U^T \left( \sum_{f'} \nabla y_{sf'} \odot w_{f'f} \right)$
  - 8: Compute gradient w.r.t interpolated weights:  $\nabla w_{f'f} = U^T \left( \sum_s \nabla y_{sf'} \odot x_{sf} \right)$
  - 9: Compute gradient w.r.t weights  $\nabla w_{\tilde{f}f} = K^T \nabla w_{f'f}$ .
- 

### 3.2 Pooling with Hierarchical Graph Clustering

In image and speech applications, and in order to reduce the complexity of the model, it is often useful to trade-off spatial resolution with feature resolution as the representation becomes deeper. For that purpose, pooling layers compute statistics in local neighborhoods, such as the average amplitude, energy or maximum activation.

The same layers can be defined in a Graph by providing the equivalent notion of neighborhood. In this work, we construct such neighborhoods at different scales using multi-resolution spectral clustering [15], and consider both average and max-pooling as in standard convolutional network architectures.

## 4 Graph Construction

### 4.1 Using Prior Knowledge

Whereas some recognition tasks in non-Euclidean domains, such as those considered in [2] or [10], might have a prior knowledge of the graph structure of the input data, many other real-world applications do not have such knowledge. It is thus necessary to estimate a similarity matrix  $W$  from the data before constructing the spectral network.

We consider in this paper two possible graph constructions, one unsupervised by measuring joint feature statistics, and another one supervised using an initial network as a proxy for the estimation.

## 4.2 Unsupervised Graph Estimation

Given data  $X \in \mathbb{R}^{L \times N}$ , where  $L$  is the number of samples and  $N$  the number of features, the simplest approach to estimating a graph structure from the data is to consider a distance between features  $i$  and  $j$  given by

$$d(i, j) = \|X_i - X_j\|^2,$$

where  $X_i$  is the  $i$ -th column of  $X$ . While correlations are typically sufficient to reveal the intrinsic geometrical structure of images [13], the effects of higher-order statistics might be non-negligible in other contexts, especially in presence of sparsity. Indeed, in many situations the pairwise Euclidean distances might suffer from unnormalized measurements. Several strategies and variants exist to gain some robustness, for instance replacing the Euclidean distance by the  $Z$ -score (thus renormalizing each feature by its standard deviation), the “square-correlation” (computing the correlation of squares of previously whitened features), or the mutual information.

This distance is then used to build a Gaussian diffusion Kernel [1]

$$\omega(i, j) = \exp^{-\frac{d(i, j)}{\sigma^2}}. \quad (1)$$

In our experiments, consider the variant of self-tuning diffusion kernel [16]

$$\omega(i, j) = \exp^{-\frac{d(i, j)}{\sigma_i \sigma_j}},$$

where  $\sigma_i$  is computed as the distance  $d(i, i_k)$  corresponding to the  $k$ -th nearest neighbor  $i_k$  of feature  $i$ .

The main advantage of (1) is that it does not require labeled data. Therefore, it is possible to estimate the similarity using several datasets that share the same features, for example in text classification.

## 4.3 Supervised Graph Estimation

As discussed in the previous section, the notion of feature similarity is not well defined, as it depends on our choice of kernel and criteria. Therefore, in the context of supervised learning, the relevant statistics from the input signals might not correspond to our imposed similarity criteria. It may thus be interesting to ask for the feature similarity that best suits a particular classification task.

A particularly simple approach is to use a fully connected network to determine the feature similarity. Given a training set with normalized<sup>1</sup> features  $X \in \mathbb{R}^{L \times N}$  and labels  $y \in \{1, \dots, C\}^L$ , we initially train a fully connected network  $\phi$  with  $K$  layers of weights  $W_1, \dots, W_K$ , using standard ReLU activations and dropout. We then extract the first layer features  $W_1 \in \mathbb{R}^{N \times M_1}$ , where  $M_1$  is the number of first-layer hidden features, and consider the distance

$$d_{sup}(i, j) = \|W_{1,i} - W_{1,j}\|^2, \quad (2)$$

that is then fed into the Gaussian kernel as in (1). The interpretation is that the supervised criterion will extract through  $W_1$  a collection of linear measurements that best serve the classification task. Thus two features are similar if the network decides to extract similar collection of linear measurements.

Another supervised strategy that further exploits the dependencies learnt at the deeper layers of the network consists in using the gradients of the supervised loss function with respect to the inputs as discriminative information. More concretely, if  $\mathcal{L}(X_l, y_l, (W_k)_{k \leq K})$  is the loss function associated to training example  $l$  and  $Z_{l,i} = \nabla_{x_i} \mathcal{L}(X_l, y_l, (W_k)_{k \leq K})$  is the input gradient corresponding to that example with respect to the  $i$ -th feature, then we consider the distance

$$d_{gsup}(i, j) = \|Z_{\cdot,i} - Z_{\cdot,j}\|^2, \quad (3)$$

that is, we measure the average similarity between the gradients for features  $i$  and  $j$ . That is, the similarity of two features depends on their respective influence (averaged over the training set) to the outcome of the network.

These constructions can be seen as “distilling” the information learnt by a first network into a kernel. In the general case where no assumptions are made on the dimension of the graph, it amounts

<sup>1</sup>In our experiments we simply normalized each feature by its standard deviation, but one could also whiten completely the data.

to extracting  $N^2$  parameters from the first learning stage (which typically involves a much larger number of parameters). If, moreover, we assume a low-dimensional graph structure of dimension  $m$ , then  $mN$  parameters are extracted by projecting the resulting kernel into its leading  $m$  directions.

Finally, observe that one could simply replace the eigen-basis  $U$  obtained by diagonalizing the graph Laplacian by an arbitrary unitary matrix, which is then optimized by back-propagation together with the rest of the parameters of the model. We do not report results on this strategy, although we point out that it was the same learning complexity as the Fully Connected network (requiring  $O(KN^2)$  parameters, where  $K$  is the number of layers and  $N$  is the input dimension).

## 5 Experiments

Word datasets. Mention that Using Word embeddings is also a possibility for future work.

## 6 Discussion

The graph construction using the Supervised Graph Estimation. Related to discovering latent graphical models

Limitations: When is weight sharing appropriate/ not.

When is locality appropriate/not appropriate.

Complexity.

Alternatives we are

## References

- [1] Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*, volume 14, pages 585–591, 2001.
- [2] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and deep locally connected networks on graphs. In *Proceedings of the 2nd International Conference on Learning Representations*, 2013.
- [3] Xu Chen, Xiuyuan Cheng, and Stéphane Mallat. Unsupervised deep haar scattering on graphs. In *Advances in Neural Information Processing Systems*, pages 1709–1717, 2014.
- [4] Adam Coates and Andrew Y Ng. Selecting receptive fields in deep networks. In *Advances in Neural Information Processing Systems*, pages 2528–2536, 2011.
- [5] Karo Gregor and Yann LeCun. Emergence of complex-like cells in a temporal product network with local receptive fields. *arXiv preprint arXiv:1006.0448*, 2010.
- [6] Geoffrey Hinton, Li Deng, Dong Yu, George Dahl, Abdel rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara Sainath, and Brian Kingsbury. Deep neural networks for acoustic modeling in speech recognition. *Signal Processing Magazine*, 2012.
- [7] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. Imagenet classification with deep convolutional neural networks. In F. Pereira, C.J.C. Burges, L. Bottou, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 1097–1105. Curran Associates, Inc., 2012.
- [8] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 05 2015.
- [9] Stéphane Mallat. *A wavelet tour of signal processing*. Academic press, 1999.
- [10] Jonathan Masci, Davide Boscaini, Michael M. Bronstein, and Pierre Vandergheynst. Shapenet: Convolutional neural networks on non-euclidean manifolds. *CoRR*, abs/1501.06297, 2015.
- [11] Jiquan Ngiam, Zhenghao Chen, Daniel Chia, Pang W Koh, Quoc V Le, and Andrew Y Ng. Tiled convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1279–1287, 2010.

270 [12] Pradeep Ravikumar, Martin J Wainwright, John D Lafferty, et al. High-dimensional ising  
271 model selection using  $\ell_1$ -regularized logistic regression. *The Annals of Statistics*, 38(3):1287–  
272 1319, 2010.

273 [13] Nicolas L Roux, Yoshua Bengio, Pascal Lamblin, Marc Joliveau, and Balázs Kégl. Learning  
274 the 2-d topology of images. In *Advances in Neural Information Processing Systems*, pages  
275 841–848, 2008.

276 [14] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov.  
277 Dropout: A simple way to prevent neural networks from overfitting. *The Journal of*  
278 *Machine Learning Research*, 15(1):1929–1958, 2014.

279 [15] Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–  
280 416, 2007.

281 [16] Lihi Zelnik-Manor and Pietro Perona. Self-tuning spectral clustering. In *Advances in neural*  
282 *information processing systems*, pages 1601–1608, 2004.

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323