
Scattering GCN: Overcoming Oversmoothness in Graph Conv. Networks

Yimeng Min^{*1} Frederik Wenkel^{*21} Guy Wolf²¹

Abstract

Graph convolutional networks (GCNs) are widely used for semi-supervised node classification on graphs today. The graph structure is however only accounted for by considering the similarity of activations between adjacent nodes, in turn degrading the results. In this work, we augment GCN models by incorporating richer notions of regularity by leveraging cascades of band-pass filters, known as geometric scatterings. We introduce a new hybrid architecture for the task and demonstrate its potential on multiple graph datasets, where it outperforms leading GCN models.

1. Introduction

Geometric deep learning approaches typically use graphs to model geometries, either by constructing them from input data (e.g., via similarity kernels) or directly given as quantified interactions between data points (Bronstein et al., 2017). Using such models, recent works have shown that graph neural networks (GNNs) perform well in multiple fields, including biology, chemistry and social networks (Gilmer et al., 2017; Hamilton et al., 2017; Kipf & Welling, 2016). Most GNNs consider each graph together with given node features as a generalization of images or audio signals, aiming to compute whole-graph representations. These can be applied to graph classification, for example, when each graph represents the molecular structure of proteins or enzymes classified by their chemical properties (Fout et al., 2017; De Cao & Kipf, 2018; Knyazev et al., 2018).

On the other hand, methods such as graph convolutional networks (GCNs) presented by Kipf & Welling (2016) consider node-level tasks, and in particular node classification. As explained in Kipf & Welling (2016), such tasks are of-

ten considered in the context of semi-supervised learning, as typically only a small portion of nodes on the graph possesses labels. In these settings, the entire dataset is considered as one graph and the network is tasked with learning node representations that infer information from node features as well as the graph structure. However, most state-of-the-art approaches for incorporating graph structure information in neural network operations aim to enforce similarity between representations of adjacent nodes, which essentially implements local smoothing of neuron activations over the graph (Li et al., 2018). Such smoothing operations often cause degradation of results in node processing tasks due to oversmoothing (Li et al., 2018; NT & Mahara, 2019), as nodes become indistinguishable with deeper and increasingly complex network architectures. Graph attention networks (Veličković et al., 2017) have shown promising results in overcoming such limitations by introducing adaptive weights for graph smoothing via message passing operations, using attention mechanisms computed from node features and masked by graph edges. However, these networks still essentially rely on enforcing similarity (although adaptive) between neighboring nodes, while also requiring more intricate training as their attention mechanism requires gradient computations driven not only by graph nodes, but also by graph edges.

In this paper, we propose a new approach for node-level processing in GNNs by introducing neural pathways that encode higher-order forms of regularity in graphs. Our construction is inspired by recently proposed geometric scattering networks (Gama et al., 2018; Gao et al., 2019; Zou & Lerman, 2019), which leverage deep cascades of graph wavelets (Hammond et al., 2011; Coifman & Maggioni, 2006) and pointwise nonlinearities to capture multiple modes of variation from node features or labels. Using the terminology of graph signal processing, these can be considered as generalized band-pass filtering operations, while GCNs (and many other GNNs) can be considered as relying on low-pass filters only. Our approach combines together the merits of both GCN and geometric scattering architectures to enable the learning of node-level features that encode geometric information beyond smoothed activation signals, thus alleviating oversmoothing concerns.

^{*}Equal contribution ¹Mila – Quebec AI Institute, Montréal, QC, Canada ²Department of Mathematics & Statistics, Université de Montréal, Montréal, QC, Canada. Correspondence to: Yimeng Min <yimeng.min@mail.utoronto.ca>, Guy Wolf <guy.wolf@umontreal.ca>.

Preliminaries: We refer the reader to the supplement for notations and graph signal processing preliminaries.

2. Graph Convolutional Network

Convolutional filters can be parametrized using the terminology introduced in the supplement. As for graph convolutional networks (GCNs) introduced in Kipf & Welling (2016), we choose just one learnable parameter θ (to avoid overfitting). This is done by setting $\hat{g}[i] := \theta(2 - \lambda_i)$. The resulting convolutional filtering operation is given by

$$g_\theta \star x = \theta \left(I_n + D^{-1/2} W D^{-1/2} \right) x. \quad (1)$$

The matrix $I_n + D^{-1/2} W D^{-1/2}$ has eigenvalues in $[0, 2]$, which could lead to vanishing or exploding gradients. This issue is addressed by the following renormalization trick (Kipf & Welling, 2016): $I_n + D^{-1/2} W D^{-1/2} \rightarrow \tilde{D}^{-1/2} \tilde{W} \tilde{D}^{-1/2}$, where $\tilde{W} := I_n + W$ and \tilde{D} is a diagonal matrix with $\tilde{D}[i, i] := \sum_{j=1}^n \tilde{W}[i, j]$ for $i \in [n]$. This operation replaces the features of the nodes by a weighted average of itself and its neighbors. Note that the repeated execution of graph convolutions will enforce similarity throughout higher-order neighborhoods with order equal to the number of stacked layers. Setting $A := \tilde{D}^{-1/2} \tilde{W} \tilde{D}^{-1/2}$, the complete layer-wise propagation rule takes the form $h_j^\ell = \sigma \left(\sum_{i=1}^{N_{\ell-1}} \theta_i^\ell A h_i^{\ell-1} \right)$, where $\theta_i^\ell A h_i^{\ell-1} = g_{\theta_i^\ell} \star h_i^{\ell-1}$ is the convolutional filter operation. Here, ℓ indicates the layer with N_ℓ neurons, $h_j^\ell \in \mathbb{R}^n$ the activation vector of the j^{th} neuron, θ_i^ℓ the learned parameter of the convolution with the i^{th} incoming activation vector from the preceding layer and $\sigma(\cdot)$ an element-wise applied activation function. Written in matrix notation, this gives

$$H^\ell = \sigma \left(A H^{\ell-1} \Theta^\ell \right), \quad (2)$$

where $\Theta^\ell \in \mathbb{R}^{N_{\ell-1} \times N_\ell}$ is the weight-matrix of the ℓ^{th} layer and $H^\ell \in \mathbb{R}^{n \times N_\ell}$ contains the activations of the ℓ^{th} layer.

We remark that the above explained GCN model can be interpreted as a low-pass operation. For the sake of simplicity, let us consider the convolutional operation (Eq. 1) before the renormalization trick. If we observe the convolution operation as the summation $g_\theta \star x = \sum_{i=1}^n \gamma_i \hat{x}[i] q_i$, we clearly see that higher weights $\gamma_i = \theta(2 - \lambda_i)$ are put on the low-frequency harmonics, while high-frequency harmonics are progressively less involved as $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq 2$. This indicates that the model can only access a diminishing portion of the original information contained in the input signal the more graph convolutions are stacked. This observation is in line with the well-known oversmoothing problem (Li et al., 2018) related to GCN models. The repeated application of graph convolutions will successively smooth the signals of the graph such that nodes cannot be distinguished anymore.

3. Geometric Scattering

The construction of geometric scattering on graphs is based on the *lazy random walk* matrix $P := \frac{1}{2}(I_n + W D^{-1})$, which is closely related to the *graph random walk* defined as a Markov process with transition matrix $R := W D^{-1}$. The matrix P however allows self loops while renormalizing to retain a Markov process. Therefore, considering an initial distribution $\mu_0 \in \mathbb{R}^n$, its positional distribution after t steps is encoded by $\mu_t = P^t \mu_0$.

As discussed in Gao et al. (2019), the propagation of a graph signal $x \in \mathbb{R}^n$ by $x_t = P^t x$ performs a low-pass operation that suppresses high frequencies. In geometric scattering, this low-pass information is augmented by introducing the *wavelet* matrices $\Psi_k \in \mathbb{R}^{n \times n}$ of scales 2^k , $k \in \mathbb{N}_0$,

$$\begin{cases} \Psi_0 := I_n - P, \\ \Psi_k := P^{2^{k-1}} - P^{2^k} = P^{2^{k-1}} (I_n - P^{2^k}), \end{cases} \quad (3)$$

for $k \geq 1$. This leverages the fact that high frequencies can be recovered with multiscale wavelet transforms, e.g., by decomposing nonzero frequencies into dyadic frequency bands. The operation $(\Psi_k x)[v_i]$ collects signals from a neighborhood of order 2^k , but extracts multiscale differences rather than averaging over them. The wavelets in Eq. 3 can be organized in a filter bank $\{\Psi_k, \tilde{\Psi}_K\}_{0 \leq k \leq K}$, where $\tilde{\Psi}_K := P^{2^K}$ is a pure low-pass filter. The telescoping sum of matrices in this filter bank constitutes the identity matrix, thus enabling to reconstruct processed signals from their filter responses. See, e.g., Perlmutter et al. (2019) for further discussion of this construction and its properties.

Geometric scattering was originally introduced in the context of whole-graph classification and consisted of aggregating scattering features. These are stacked wavelet transforms parameterized via tuples $J := (k_1, \dots, k_m) \in \cup_{m \in \mathbb{N}} \mathbb{N}_0^m$ containing the bandwidth scale parameters separated by element-wise absolute value nonlinearities, i.e.,

$$\Phi_J x := \Psi_{k_m} |\Psi_{k_{m-1}} \dots |\Psi_{k_2} |\Psi_{k_1} x|| \dots|,$$

where m corresponds to the length of the tuple J . The scattering features are aggregated over the whole graph by taking q^{th} -order moments over the set of nodes,

$$S(J, q) x := \sum_{i=1}^n |\Phi_J x[v_i]|^q. \quad (4)$$

Here, we modify this construction to keep the scattering transform Φ_J at the node-level by dismissing the aggregation step in Eq. 4. For each tuple J , we define the scattering propagation rule

$$H^\ell = \sigma \left(\Phi_J H^{\ell-1} \Theta^\ell \right), \quad (5)$$

which mirrors the GCN one while replacing the low-pass filter by a geometric scattering operation. We note that in practice, we only use a subset of tuples, selected as part of our network design explained in the following section.

4. Combining GCN and Scattering Models

To combine the benefits of GCNs and geometric scattering adapted to the node level, we propose a hybrid architecture that combines low-pass operations based on the GCN model with band-pass operations from geometric scattering. To define the layer-wise propagation rule, we introduce

$$\mathbf{H}_{gcn}^\ell := \left[\mathbf{H}_{gcn,1}^\ell \parallel \dots \parallel \mathbf{H}_{gcn,C_{gcn}}^\ell \right],$$

and

$$\mathbf{H}_{sct}^\ell := \left[\mathbf{H}_{sct,1}^\ell \parallel \dots \parallel \mathbf{H}_{sct,C_{sct}}^\ell \right],$$

which are the concatenations of channels $\{\mathbf{H}_{gcn,k}^\ell\}_{k=1}^{C_{gcn}}$ and $\{\mathbf{H}_{sct,k}^\ell\}_{k=1}^{C_{sct}}$, respectively. Every $\mathbf{H}_{gcn,k}^\ell$ is defined according to Eq. 2 with the slight modification of added biases and powers of \mathbf{A} , i.e.,

$$\mathbf{H}_{gcn,k}^\ell := \sigma \left(\mathbf{A}^k \mathbf{H}^{\ell-1} \Theta_{gcn,k}^\ell + \mathbf{B}_{gcn,k}^\ell \right).$$

Note that every GCN filter uses a different propagation matrix \mathbf{A}^k and therefore aggregates information from k -step neighborhoods. Similarly, we proceed with $\mathbf{H}_{sct,k}^\ell$ according to Eq. 5 and calculate

$$\mathbf{H}_{sct,k}^\ell := \sigma \left(\Phi_{J_k} \mathbf{H}^{\ell-1} \Theta_{sct,k}^\ell + \mathbf{B}_{sct,k}^\ell \right),$$

where $J_k \in \bigcup_{m \in \mathbb{N}} \mathbb{N}_0^m$, $k = 1, \dots, C_{sct}$ enables scatterings of different orders and scales. Finally, the GCN components and scattering components get concatenated to

$$\mathbf{H}^\ell := \left[\mathbf{H}_{gcn}^\ell \parallel \mathbf{H}_{sct}^\ell \right]. \quad (6)$$

The learned parameters are the weight matrices $\Theta_{gcn,k}^\ell, \Theta_{sct,k}^\ell \in \mathbb{R}^{N_{\ell-1} \times N_\ell}$ coming from the convolutional and scattering layers. These are complemented by vectors of the biases $\mathbf{b}_{gcn,k}^\ell, \mathbf{b}_{sct,k}^\ell \in \mathbb{R}^{N_\ell}$, which are transposed and vertically concatenated n times to the matrices $\mathbf{B}_{gcn,k}, \mathbf{B}_{sct,k} \in \mathbb{R}^{n \times N_\ell}$. To simplify notation, we assume here that all channels use the same number of neurons (N_ℓ). Waiving this assumption would slightly complicate the notation but works perfectly fine in practice.

In this work, for simplicity, and because it is sufficient to establish our claim, we limit our architecture to three GCN channels and two scattering channels. Inspired by the aggregation step in classical geometric scattering, we use $\sigma(\cdot) := |\cdot|^q$ as nonlinearity. However, unlike in Eq. 4, the q^{th} power is applied at the node-level instead of being aggregated as moments over the entire graph, thus retaining the distinction between node-wise activations. We set the input of the first layer \mathbf{H}^0 to be the original node features. Each subchannel transforms the original feature space to a new hidden space with the dimension determined by the number

of neurons encoded in the columns of the corresponding submatrix of \mathbf{H}^ℓ . These transformations are learned via the weights and biases. Larger matrices \mathbf{H}^ℓ (i.e., more columns, as the number of nodes in the graph is fixed) indicate that the weight matrices have more parameters to learn. Thus, the information in these channels can be propagated well and will be sufficiently represented. The channel width reflects the importance of the captured regularities. Wider channels suggest that their frequency components are more critical and need to be sufficiently learned.

We note that in sparsely labelled graphs, high-frequency noise (e.g., differences between labelled and unlabelled nodes) can unintentionally be captured by scattering features. To mitigate such artifacts, we introduce the *graph residual convolution*, an adjustable lowpass filter inspired by skip connections in residual networks. Governed by the hyperparameter α , it is defined via the matrix $\mathbf{A}_{res}(\alpha) = \frac{1}{\alpha+1}(\mathbf{I}_n + \alpha \mathbf{W} \mathbf{D}^{-1})$ and applied after the hybrid layer of GCN and scattering filters. We have $\mathbf{A}_{res}(0) = \mathbf{I}_n$, while $\alpha \rightarrow \infty$ gives $\mathbf{R} = \mathbf{W} \mathbf{D}^{-1}$, which is an interpolation between the completely lazy random walk \mathbf{I}_n and the non-resting random walk \mathbf{R} . We apply the graph residual layer on the output \mathbf{H}^ℓ of the scattering GCN layer (Eq. 6). The update rule for this step is expressed by $\mathbf{H}^{\ell+1} = \mathbf{A}_{res}(\alpha) \mathbf{H}^\ell \Theta_{res} + \mathbf{B}_{res}$, where $\Theta_{res} \in \mathbb{R}^{N \times N_{\ell+1}}$ are learned weights, $\mathbf{B}_{res} \in \mathbb{R}^{n \times N_{\ell+1}}$ are learned biases (similar to the notations used previously), and N is the number of features of the concatenated layer \mathbf{H}^ℓ in Eq. 6. If $\mathbf{H}^{\ell+1}$ is the final layer, we set $N_{\ell+1}$ equal to the number of classes.

5. Additional Information Introduced by Node-level Scattering Features

Before empirically verifying the viability of the proposed architecture in node classification tasks, we first discuss a property demonstrating the additional information (in particular, carried by node features) provided by scattering channels beyond that provided by traditional GCN channels. The following lemma shows a particular type of such information in the form of two-coloring features provided on bipartite graphs. We note that such features are highly regular and completely determined by the graph structure, but the smoothing in GCN channels would eliminate their information while the addition of certain scattering channels would retain it for further downstream processing.

Lemma. *Consider a bipartite graph on $n \in \mathbb{N}$ nodes with constant node degree β . Let $\mathbf{x} \in \mathbb{R}^n$ be a 2-coloring signal (i.e., one part assigned constant a and the other b , where $a \neq b \in \mathbb{R}$). Then, for any $\theta \in \mathbb{R}$, the GCN filtering $\mathbf{g}_\theta \star \mathbf{x}$ from Eq. 1 yields a constant signal, while the scattering filter $\Psi_0 \mathbf{x}$ from Eq. 3 still produces a (non-constant) 2-coloring. This result extends to any finite linear filter cascade (i.e., $\mathbf{g}_\theta \star \dots \star \mathbf{g}_\theta \star \mathbf{x}$ or $\Psi_0^k \mathbf{x}$ with $k \in \mathbb{N}$ filter applications).*

Table 1. Dataset characteristics

Dataset	Nodes	Edges	Features	Degrees	$\frac{\text{Edges}}{\text{Nodes}}$
Citeseer	3,327	4,732	3,703	3.77 ± 3.38	1.42
Cora	2,708	5,429	1,433	4.90 ± 5.22	2.00
Pubmed	19,717	44,338	500	5.50 ± 7.43	2.25
DBLP	17,716	105,734	1639	6.97 ± 9.35	5.97

6. Empirical Results

To evaluate our Scattering GCN approach, we compare it to several established methods for semi-supervised node classification, including the original GCN (Kipf & Welling, 2016) and the GAT network (Veličković et al., 2017), which indirectly addresses oversmoothing by training adaptive node-wise weighting of the smoothing operation via an attention mechanism, as well as the methods from Li et al. (2018, Partially Absorbing), Defferrard et al. (2016, Chebyshev), and Zhu et al. (2003, Label Propagation). Our comparisons are based on four popular graph datasets¹ with varying sizes and connectivity structures summarized in Tab. 1. We order the datasets here by increasing connectivity structure, reflected by node degrees and edges-to-nodes ratios. As discussed in Li et al. (2018), increased connectivity leads to faster mixing of node features in GCN, exacerbating the oversmoothing problem (as nodes quickly become indistinguishable) and degrading classification performance. Therefore, we expect the relative improvement achieved by Scattering GCN to correspond to the increasing connectivity order in Tab. 1, which is maintained in Tab. 2 and Fig. 1.

Table 2. Classification accuracy comparison (top two methods marked in bold; best one underlined) on four benchmark datasets.

MODEL	CITSEER	CORA	PUBMED	DBLP
SCATTERING GCN (OURS)	71.7	84.2	79.4	81.5
GAT	72.5	83.0	79.0	66.1
PARTIALLY ABSORBING	71.2	81.7	79.2	56.9
GCN	70.3	81.5	79.0	59.3
CHEBYSHEV	69.8	78.1	74.4	57.3
LABEL PROPAGATION	58.2	77.3	71.0	53.0
NODE FEATURES (SVM)	61.1	58.0	49.9	48.2

We first consider test classification accuracy reported in Tab. 2, which shows that our approach outperforms other methods on three out of the four considered datasets. On the remaining one (namely Citeseer) we are only outperformed by GAT. However, we note that this dataset has the weakest connectivity structure (see Tab. 1) and the most informative node features (e.g., achieving 61.1% accuracy via linear SVM without considering any graph information). In contrast, on DBLP, which has the richest connectivity structure and least informative features (only 48.2% SVM accuracy), we significantly outperform GAT (over 15% improvement),

¹See, e.g., Yang et al. (2016) for Citeseer, Cora, and Pubmed, and Pan et al. (2016) for DBLP

which in itself significantly outperforms all other methods.

Next, we consider the impact of training size on classification performance, as we are interested in semi-supervised settings where only a small portion of nodes in the graph are labelled. Fig. 1 presents the classification accuracy (on validation set) for the training size reduced to 20%, 40%, 60%, 80% and 100% of the original training size available for each dataset. These results indicate that generally Scattering GCN exhibits greater stability to sparse training conditions compared to other methods. Importantly, we note that on Citeseer, while GAT outperforms our method for the original training size, its performance degrades rapidly when the training size is reduced below 60% of the original one, at which point Scattering GCN outperforms all other methods. We also note that on Pubmed, even a small decrease in training size creates a significant performance gap between Scattering GCN and GAT, which we believe is due to node features being less independently informative in this case (see baseline in Tab. 2) compared to Citeseer and Cora.

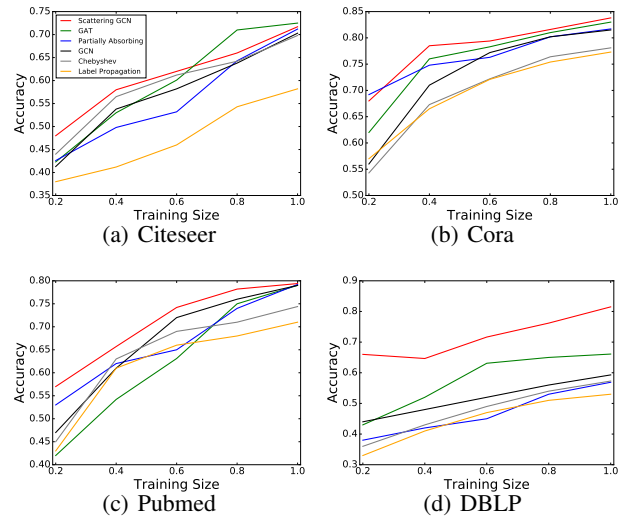


Figure 1. Impact of training set size on classification accuracy.

7. Conclusion

Our study of semi-supervised node classification in graphs presents a new approach to address some of the main concerns and limitations of GCN models. Our construction is inspired by geometric scattering, which has mainly been used for whole-graph classification so far. Our results demonstrate several benefits of incorporating the elements presented here (i.e., scattering channels and residual convolution) in GCN architectures. We expect future work looking to incorporate these elements together in more intricate architectures (e.g., with attention mechanisms) to provide promising new capabilities of pattern recognition and local information extraction in graphs.

Acknowledgements

The authors would like to thank Dongmian Zou for fruitful discussions. This work was partially funded by IVADO (l'institut de valorisation des données) and NIH grant R01GM135929. The content provided here is solely the responsibility of the authors and does not necessarily represent the official views of the funding agencies.

References

- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34 (4):18–42, 2017.
- Chung, F. R. K. *Spectral Graph Theory*. American Mathematical Society, 1997.
- Coifman, R. R. and Maggioni, M. Diffusion wavelets. *Applied and Computational Harmonic Analysis*, 21(1):53 – 94, 2006.
- De Cao, N. and Kipf, T. Molgan: An implicit generative model for small molecular graphs. *arXiv preprint arXiv:1805.11973*, 2018.
- Defferrard, M., Bresson, X., and Vandergheynst, P. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in neural information processing systems*, pp. 3844–3852, 2016.
- Fout, A., Byrd, J., Shariat, B., and Ben-Hur, A. Protein interface prediction using graph convolutional networks. In *Advances in neural information processing systems*, pp. 6530–6539, 2017.
- Gama, F., Ribeiro, A., and Bruna, J. Diffusion scattering transforms on graphs. *arXiv preprint arXiv:1806.08829*, 2018.
- Gao, F., Wolf, G., and Hirn, M. Geometric scattering for graph data analysis. In *International Conference on Machine Learning*, pp. 2122–2131, 2019.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1263–1272. JMLR, 2017.
- Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pp. 1024–1034, 2017.
- Hammond, D. K., Vandergheynst, P., and Gribonval, R. Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis*, 30(2):129 – 150, 2011.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International conference on learning representations*, 2016.
- Knyazev, B., Lin, X., Amer, M. R., and Taylor, G. W. Spectral multigraph networks for discovering and fusing relationships in molecules. *arXiv preprint arXiv:1811.09595*, 2018.
- Li, Q., Han, Z., and Wu, X.-M. Deeper insights into graph convolutional networks for semi-supervised learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Liao, R., Zhao, Z., Urtasun, R., and Zemel, R. S. Lanczosnet: Multi-scale deep graph convolutional networks. *arXiv preprint arXiv:1901.01484*, 2019.
- NT, H. and Maehara, T. Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550*, 2019.
- Pan, S., Wu, J., Zhu, X., Zhang, C., and Wang, Y. Tri-party deep network representation. *Network*, 11(9):12, 2016.
- Perlmutter, M., Gao, F., Wolf, G., and Hirn, M. Understanding graph neural networks with asymmetric geometric scattering transforms. *arXiv preprint arXiv:1911.06253*, 2019.
- Shuman, D. I., Ricaud, B., and Vandergheynst, P. Vertex-frequency analysis on graphs. *Applied and Computational Harmonic Analysis*, 40(2):260–291, 2016.
- Susnjara, A., Perraudin, N., Kressner, D., and Vandergheynst, P. Accelerated filtering on graphs using lanczos method. *arXiv preprint arXiv:1509.04537*, 2015.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Yang, Z., Cohen, W. W., and Salakhutdinov, R. Revisiting semi-supervised learning with graph embeddings. *arXiv preprint arXiv:1603.08861*, 2016.
- Zhu, X., Ghahramani, Z., and Lafferty, J. D. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pp. 912–919, 2003.
- Zou, D. and Lerman, G. Graph convolutional neural networks via scattering. *Applied and Computational Harmonic Analysis*, 2019.

Supplement

A. Notations

We denote matrices and vectors with bold letters with uppercase letters representing matrices and lowercase letters representing vectors. In particular, $\mathbf{I}_n \in \mathbb{R}^{n \times n}$ is used for the identity matrix and $\mathbf{1}_n \in \mathbb{R}^n$ denotes the vector with ones in every component. We write $\langle \cdot, \cdot \rangle$ for the standard scalar product in \mathbb{R}^n . We will interchangeably consider functions of graph nodes as vectors indexed by the nodes, implicitly assuming a correspondence between a node and a specific index. This carries over to matrices, where we relate nodes to column or row indices. We further use the abbreviation $[n] := \{1, \dots, n\}$ where $n \in \mathbb{N}$ and write $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$.

B. Graph Signal Processing Preliminaries

Let $G = (V, E, w)$ be a weighted graph with $V := \{v_1, \dots, v_n\}$ the set of nodes, $E \subset \{\{v_i, v_j\} \in V \times V, i \neq j\}$ the set of (undirected) edges and $w : E \rightarrow (0, \infty)$ assigning (positive) edge weights to the graph edges. We note that w can equivalently be considered as a function of $V \times V$, where we set the weights of non-adjacent node pairs to zero. We define a *graph signal* as a function $x : V \rightarrow \mathbb{R}$ on the nodes of G and aggregate them in a signal vector $\mathbf{x} \in \mathbb{R}^n$ with the i^{th} entry being $x(v_i)$.

We define the (combinatorial) *graph Laplacian* matrix $\mathbf{L} := \mathbf{D} - \mathbf{W}$, where $\mathbf{W} \in \mathbb{R}^{n \times n}$ is the *weighted adjacency matrix* of the graph G given by

$$\mathbf{W}[v_i, v_j] := \begin{cases} w(v_i, v_j) & \text{if } \{v_i, v_j\} \in E \\ 0 & \text{otherwise,} \end{cases}$$

and $\mathbf{D} \in \mathbb{R}^{n \times n}$ is the *degree matrix* of G defined by $\mathbf{D} := \text{diag}(d_1, \dots, d_n)$ with $d_i := \deg(v_i) := \sum_{j=1}^n \mathbf{W}[v_i, v_j]$ being the *degree* of the node v_i . In practice, we work with the (symmetric) *normalized Laplacian* matrix $\mathbf{L} := \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2} = \mathbf{I}_n - \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$. It can be verified that \mathbf{L} is symmetric and positive semi-definite and can thus be orthogonally diagonalized as $\mathbf{L} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^T = \sum_{i=1}^n \lambda_i \mathbf{q}_i \mathbf{q}_i^T$, where $\mathbf{\Lambda} := \text{diag}(\lambda_1, \dots, \lambda_n)$ is a diagonal matrix with the eigenvalues on the main diagonal and \mathbf{Q} is an orthogonal matrix containing the corresponding normalized eigenvectors $\mathbf{q}_1, \dots, \mathbf{q}_n \in \mathbb{R}^n$ as its columns.

A detailed study (see, e.g., (Chung, 1997)) of the eigenvalues reveals that $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n \leq 2$. We can interpret the $\lambda_i, i \in [n]$ as the frequency magnitudes and the \mathbf{q}_i as the corresponding Fourier modes. We accordingly define the *Fourier transform* of a signal vector $\mathbf{x} \in \mathbb{R}^n$ by $\hat{\mathbf{x}}[i] = \langle \mathbf{x}, \mathbf{q}_i \rangle$ for $i \in [n]$. The corresponding inverse Fourier transform is given by $\mathbf{x} = \sum_{i=1}^n \hat{\mathbf{x}}[i] \mathbf{q}_i$. Note that this can be written compactly as $\hat{\mathbf{x}} = \mathbf{Q}^T \mathbf{x}$ and

$\mathbf{x} = \mathbf{Q} \hat{\mathbf{x}}$. Finally, we introduce the concept of *graph convolutions*. We define a filter $g : V \rightarrow \mathbb{R}$ defined on the set of nodes and want to convolve the corresponding filter vector $\mathbf{g} \in \mathbb{R}^n$ with a signal vector $\mathbf{x} \in \mathbb{R}^n$, i.e. $\mathbf{g} \star \mathbf{x}$. To explicitly compute this convolution, we recall that in the Euclidean setting, the Fourier transform of the convolution of two signals equals the product of the corresponding Fourier transforms. This property generalizes to graphs (Shuman et al., 2016) in the sense that $(\widehat{\mathbf{g} \star \mathbf{x}})[i] = \hat{\mathbf{g}}[i] \hat{\mathbf{x}}[i]$ for $i \in [n]$. Applying the inverse Fourier transform yields $\mathbf{g} \star \mathbf{x} = \sum_{i=1}^n \hat{\mathbf{g}}[i] \hat{\mathbf{x}}[i] \mathbf{q}_i = \sum_{i=1}^n \hat{\mathbf{g}}[i] \langle \mathbf{q}_i, \mathbf{x} \rangle \mathbf{q}_i = \mathbf{Q} \hat{\mathbf{G}} \mathbf{Q}^T \mathbf{x}$, where $\hat{\mathbf{G}} := \text{diag}(\hat{\mathbf{g}}) = \text{diag}(\hat{\mathbf{g}}[1], \dots, \hat{\mathbf{g}}[n])$. Hence, convolutional graph filters can be parameterized by using the Fourier coefficients in $\hat{\mathbf{G}}$.

Furthermore, it can be verified (Defferrard et al., 2016) that when these coefficients are defined as polynomials $\hat{\mathbf{g}}[i] := \sum_k \gamma_k \lambda_i^k$ for $i \in \mathbb{N}$ of the Laplacian eigenvalues in $\mathbf{\Lambda}$ (i.e. $\hat{\mathbf{G}} = \sum_k \gamma_k \mathbf{\Lambda}^k$), the resulting filter convolution are localized in space and can be written in terms of \mathcal{L} as $\mathbf{g} \star \mathbf{x} = \sum_k \gamma_k \mathcal{L}^k \mathbf{x}$ without requiring spectral decomposition of the normalized Laplacian. This motivates the standard practice (Kipf & Welling, 2016; Defferrard et al., 2016; Susnjara et al., 2015; Liao et al., 2019) of using filters that have polynomial forms, which we follow here as well.

C. Proof of the Lemma

Proof. We first notice that if $\mathbf{g}_{\frac{1}{2}} \star \mathbf{x}$ is constant, then $\mathbf{g}_{\theta} \star \mathbf{x} = 2\theta(\mathbf{g}_{\frac{1}{2}} \star \mathbf{x})$ is constant for any θ . Furthermore, for the considered class of graphs, $\mathbf{D} = \beta \mathbf{I}_n$ with $\beta > 0$, implying that $\mathbf{D}^{-1} = \frac{1}{\beta} \mathbf{I}_n$ and $\mathbf{D}^{-1/2} = \frac{1}{\sqrt{\beta}} \mathbf{I}_n$. Therefore, as a direct result of Eq. 1 in the main paper, it holds that

$$\mathbf{g}_{\frac{1}{2}} \star \mathbf{x} = \left(\frac{1}{2} \mathbf{I}_n + \frac{1}{2\beta} \mathbf{W} \right) \mathbf{x} = \mathbf{P} \mathbf{x}. \quad (7)$$

Similarly, it is easily verified that any $k \in \mathbb{N}$ applications of the convolution with \mathbf{g}_{θ} (for any $\theta \in \mathbb{R}$) can be written as $2^k \theta^k \mathbf{P}^k \mathbf{x}$. Furthermore, since \mathbf{P} is column-stochastic and (here) symmetric (thus also row-stochastic), we have $\mathbf{P} \mathbf{c} = \mathbf{c}$ for any constant signal $\mathbf{c} = c \mathbf{1}_{2n}$. Thus, it is sufficient to show that $\mathbf{P} \mathbf{x}$ is a constant signal to verify the first claim of the lemma.

We consider the set of nodes V . For any node $v \in V$, according to Eq. 7, we can write

$$\begin{aligned} (\mathbf{P} \mathbf{x})[v] &= \overbrace{\mathbf{P}[v, v]}^{=\frac{1}{2}} \mathbf{x}[v] + \sum_{u \in \mathcal{N}(v)} \overbrace{\mathbf{P}[v, u]}^{=\frac{1}{2\beta}} \mathbf{x}[u] \\ &\quad + \sum_{w \in V^v} \underbrace{\mathbf{P}[v, w]}_{=0} \mathbf{x}[w], \end{aligned}$$

where we denote by $\mathcal{N}(v)$ the neighborhood of the node v

and set $V^v := V \setminus (\{v\} \cup \mathcal{N}(v))$. This implies that

$$(\mathbf{P}\mathbf{x})[v] = \frac{\mathbf{x}[v]}{2} + \sum_{u \in \mathcal{N}(v)} \frac{\mathbf{x}[u]}{2\beta}. \quad (8)$$

We now consider a 2-coloring signal $\mathbf{x} \in \mathbb{R}^n$. W.l.o.g., let $\mathbf{x}[v] = a$, which implies $\mathbf{x}[u] = b$ for all $u \in \mathcal{N}(v)$. Now, since $|\mathcal{N}(v)| = \beta$, it holds

$$(\mathbf{P}\mathbf{x})[v] = \frac{a+b}{2},$$

thus verifying the first claim of the lemma as the choice of v was arbitrary. Finally, it is now straightforward to verify the second claim as well, since the operation $\Psi_0\mathbf{x} = (\mathbf{I}_n - \mathbf{P})\mathbf{x} = \mathbf{x} - \frac{a+b}{2} \mathbf{1}_n$ retains a 2-coloring signal (the original colors are shifted by a constant: $-\frac{a+b}{2}$). \square