

# Introduction to Graph Neural Networks

Graph Deep Learning 2022

Daniele Grattarola February 28, 2022 Things we are going to cover:

- Practical introduction to GNNs
- Message passing
- Anisotropic GNNs (attention, edge attributes)
- $\cdot$  Spectral graph theory and GNNs

#### From CNNs to GNNs



- The receptive field of a CNN reflects the **underlying grid structure**.
- The CNN has an inductive bias on how to process the individual pixels/timesteps/nodes.

#### From CNNs to GNNs



- Drop assumptions about underlying structure: it is now an **input of the problem**.
- The only thing we know: the representation of a node depends on its **neighbors**.



## Discrete convolution





#### Discrete convolution:

$$(f \star g)[n] = \sum_{m=-M}^{M} f[n-m]g[m]$$

Problems:

- Variable degree of nodes
- $\cdot$  Orientation



- Graph: nodes connected by edges;
- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N], \mathbf{x}_i \in \mathbb{R}^F$ , node attributes or "graph signal";
- $\mathbf{e}_{ij} \in \mathbb{R}^{S}$ , edge attribute for edge  $i \rightarrow j$ ;
- + A,  $N \times N$  adjacency matrix;
- $\mathbf{D} = \text{diag}([d_1, \ldots, d_N])$ , diagonal degree matrix;
- $\cdot \, \mathbf{L} = \mathbf{D} \mathbf{A}$ , Laplacian;
- Reference operator R:  $\mathbf{r}_{ij} \neq 0$  if  $\exists i \rightarrow j$ Note: for us, **R** is symmetric.



### A quick recipe for a local learnable filter

Applying **R** to node attributes **X** is a **local** action:

$$(\mathsf{RX})_i = \sum_{j=1}^N \mathsf{r}_{ij} \cdot \mathsf{x}_j = \sum_{j \in \mathcal{N}(i)} \mathsf{r}_{ij} \cdot \mathsf{x}_j$$

Instead of having a different weight for each neighbor, we share weights among nodes in the same neighborhood:

$$X' = RX\Theta$$

 $x_4$  $r_{14}$  $r_{14}$  $r_{14}$  $r_{14}$  $r_{13}$  $r_{13}$  $r_{13}$  $r_{13}$ 

where  $\Theta \in \mathbb{R}^{F \times F'}$ .

#### Powers of R

Let's consider the effect of applying  $R^2$  to X:

$$(\mathsf{RRX})_i = \sum_{j \in \mathcal{N}(i)} \mathsf{r}_{ij}(\mathsf{RX})_j = \sum_{j \in \mathcal{N}(i)} \sum_{k \in \mathcal{N}(j)} \mathsf{r}_{ij} \cdot \mathsf{r}_{jk} \cdot \mathsf{x}_k$$

Key idea: by applying R<sup>K</sup> we aggregate information from the *K*-th order neighborhood of a node.





To cover all neighbors of order 0 to *K*, we can just take a polynomial with weights  $\Theta^{(k)}$ :

$$\mathsf{X}' = \sum_{k=0}^{K} \mathsf{R}^k \mathsf{X} \Theta^{(k)}$$



## Chebyshev polynomials [1]

A recursive definition using Chebyshev polynomials:

$$T^{(0)} = I$$

$$T^{(1)} = \tilde{L}$$

$$T^{(k)} = 2 \cdot \tilde{L} \cdot T^{(k-1)} - T^{(k-2)}$$
Where  $\tilde{L} = \frac{2L_n}{\lambda_{max}} - I$  and  $L_n = I - \underbrace{D^{-1/2}AD^{-1/2}}_{A_n}$ 
Layer:  $X' = \sum_{k=0}^{K} T^{(k)}X\Theta^{(k)}$ 

1 00

<sup>[1]</sup> M. Defferrard et al., "Convolutional neural networks on graphs with fast localized spectral filtering," 2016.

## Graph convolutional networks [2]

Polynomial of order  $K \rightarrow K$  layers of order 1;

Three simplifications:

1. 
$$\lambda_{\max} = 2 \rightarrow \tilde{L} = \frac{2L_n}{\lambda_{\max}} - I = -D^{-1/2}AD^{-1/2} = -A_n$$
  
2.  $K = 1 \rightarrow X' = X\Theta^{(0)} - A_n X\Theta^{(1)}$   
3.  $\Theta = \Theta^{(0)} = -\Theta^{(1)}$ 

Layer:  $X' = (\underbrace{I + A_n}_{\tilde{A}}) X \Theta = \tilde{A} X \Theta$ For stability:  $(I + A_n) \rightarrow D^{-1/2} (I + A) D^{-1/2}$ 



<sup>[2]</sup> T. N. Kipf et al., "Semi-supervised classification with graph convolutional networks," 2016.

Message passing

## Message passing neural networks [3]



<sup>[3]</sup> J. Gilmer et al., "Neural message passing for quantum chemistry," 2017.

A general scheme for message-passing networks:

$$\mathbf{x}_{i}^{\prime} = \gamma \left( \mathbf{x}_{i}, \Box_{j \in \mathcal{N}(i)} \phi \left( \mathbf{x}_{i}, \mathbf{x}_{j}, \mathbf{e}_{ji} \right) \right),$$

- φ: message function, depends on x<sub>i</sub>, x<sub>j</sub> and possibly the edge attribute e<sub>ji</sub>;
- □<sub>j∈N(i)</sub>: aggregation function (sum, average, max, or something else...);
- $\gamma$ : **update function**, final transformation to obtain new attributes after aggregating messages.



<sup>[3]</sup> J. Gilmer et al., "Neural message passing for quantum chemistry," 2017.

# Anisotropic GNNs

## Graph attention networks [4]

- 1. Messages:  $\mathbf{h}_i = \Theta_f \mathbf{x}_i$  with  $\Theta_f \in \mathbb{R}^{F' \times F}$ .
- 2. Attention between neighbors:

2.1 Score: 
$$\mathbf{e}_{ij} = \sigma \left( \theta_a^\top [\mathbf{h}_i \parallel \mathbf{h}_j] \right)$$
, with  $\theta_a \in \mathbb{R}^{2F'}$ .  
2.2 Normalize with Softmax:  $\mathbf{a}_{ij} = \frac{\exp \left( \mathbf{e}_{ij} \right)}{\sum\limits_{k \in \mathcal{N}(i)} \exp \left( e_{ik} \right)}$ 

3. Aggregate using the attention coefficients:

$$\mathbf{x}'_i = \sum_{j \in \mathcal{N}(i)} \mathbf{a}_{ij} \mathbf{h}_j$$

|| indicates concatenation



<sup>[4]</sup> P. Velickovic et al., "Graph attention networks," 2017.

## Edge-conditioned convolution [5]

**Key idea**: incorporate edge attributes into the messages.

Use a MLP  $\phi : \mathbb{R}^{S} \to \mathbb{R}^{FF'}$  to generate weights:

 $\Theta^{(ji)} = \operatorname{reshape}(\phi(\mathbf{e}_{ji}))$ 

Use the edge-dependent weights to compute messages:

$$\mathbf{x}'_i = \Theta^{(i)} \mathbf{x}_i + \sum_{j \in \mathcal{N}(i)} \Theta^{(ji)} \mathbf{x}_j + \mathbf{b}$$



<sup>[5]</sup> M. Simonovsky et al., "Dynamic edge-conditioned filters in convolutional neural networks on graphs," 2017.

**GCNConv** Kipf & Welling

ECCConv Simonovsky & Komodakis **ChebConv** Defferrard et al.

**GATConv** Velickovic et al. GraphSageConv

Hamilton et al.

**GCSConv** Bianchi et al. ARMAConv Bianchi et al.

APPNPConv Klicpera et al.

**GINConv** Xu et al.

**TAGConv** Du et al. DiffusionConv

CrystalConv Xie & Grossman GatedGraphConv

EdgeConv Wang et al. **AGNNConv** Thekumparampil et al.

MessagePassing Gilmer et al. Message passing scheme:

- Message:  $\mathbf{m}_{ji} = \mathsf{PReLU} \left( \mathsf{BatchNorm} \left( \Theta \mathbf{x}_j + \mathbf{b} \right) \right)$
- Aggregate:  $\mathbf{m}^{agg} = \sum_{j \in \mathcal{N}(i)} \mathbf{m}_{ji}$
- Update:  $\mathbf{x}' = \mathbf{x} \mid\mid \mathbf{m}^{\text{agg}}$ ;

Architecture:

- Pre- and post-process node features using 2-layer MLPs;
- 4-6 message passing steps;



<sup>[6]</sup> J. You et al., "Design space for graph neural networks," 2020.

#### How do we use this?



Node-level learning. (e.g., social networks)



Graph-level learning. (e.g., molecules) Graph convolution





**Recall**: CNNs compute a discrete convolution

$$(f \star g)[n] = \sum_{m=-M}^{M} f[n-m]g[m]$$
 (1)

Given two functions f and g, their convolution  $f \star g$  can be expressed as:

$$f \star g = \mathcal{F}^{-1} \left\{ \mathcal{F} \left\{ f \right\} \cdot \mathcal{F} \left\{ g \right\} \right\}$$
(2)

Where  $\mathcal{F}$  is the Fourier transform and  $\mathcal{F}^{-1}$  its inverse. Can we use this major property? Key intuition – we are representing a function in a different basis.



The eigenvectors of the Laplacian for a path graph can be obtained analytically:

$$\mathbf{u}_{k}[n] = \begin{cases} 1, & \text{for } k = 0\\ e^{i\pi(k+1)n/N}, & \text{for odd } k, k < N-1\\ e^{-i\pi kn/N}, & \text{for even } k, k > 0\\ \cos(\pi n), & \text{for odd } k, k = N-1 \end{cases}$$

Looks familiar?





- Drop the "grid" assumption
- Replace  $e^{-i\frac{2\pi}{N}kn}$  with generic  $\mathbf{u}_k[n]$ :

$$\mathcal{F}_{G}\left\{f\right\}\left[k\right] = \sum_{n=0}^{N-1} f[n] \mathbf{u}_{k}[n]$$

• GFT: 
$$\mathcal{F}_G{f} = \hat{f} = \mathbf{U}^\top f;$$
  
• IGFT:  $\mathcal{F}_G^{-1}{\hat{f}} = f = \mathbf{U}\hat{f}$ 

## Graph convolution

#### Recall:

• Convolution theorem:  $f \star g = \mathcal{F}^{-1} \{ \mathcal{F} \{ f \} \cdot \mathcal{F} \{ g \} \}$ 

• Spectral theorem: 
$$\mathbf{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{\top} = \sum_{i=0}^{N-1} \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

Graph signals:  

$$f,g$$
 $\mathsf{GFT:}$ 
 $\mathsf{U}^{\top}f, \mathsf{U}^{\top}g$ 
 $\mathsf{U}^{\top}f \odot \mathsf{U}^{\top}g$ 
 $\mathsf{U}^{\top}f \odot \mathsf{U}^{\top}g$ 
 $\mathsf{U}^{\top}f \odot \mathsf{U}^{\top}g$ 
 $\mathsf{U}^{\top}f \odot \mathsf{U}^{\top}g$ 
 $\mathsf{Graph}$  filter:  $\mathsf{U}\left(\mathsf{U}^{\top}f \odot \mathsf{U}^{\top}g\right) = \mathsf{U} \cdot \underbrace{\operatorname{diag}(\mathsf{U}^{\top}g)}_{g(\Lambda)} \cdot \mathsf{U}^{\top}f = \underbrace{\mathsf{U}} \cdot g(\Lambda) \cdot \mathsf{U}^{\top}f = g(\mathsf{L})f$ 

 $<sup>^{1}\</sup>odot$  indicates element-wise multiplication

Spectral GCNs

A first idea [7]: transformation of **each individual eigenvalue** is learned with a free parameter  $\theta_i$ .

Problems:

- O(N) parameters;
- not localized in node space (the only thing that we want);
- $\mathbf{U} \cdot g(\mathbf{\Lambda}) \cdot \mathbf{U}^{\top}$  costs  $O(N^2)$ ;



<sup>[7]</sup> J. Bruna et al., "Spectral networks and locally connected networks on graphs," 2013.

Better idea [7]:

- Localized in node domain ↔ smooth in spectral domain;
- Learn only a few parameters  $\theta_i$ ;
- Interpolate the other eigenvalues using a smooth cubic spline;

Localized and O(1) parameters, but multiplying by U twice is still expensive.



<sup>[7]</sup> J. Bruna et al., "Spectral networks and locally connected networks on graphs," 2013.

The same recursion is used to filter eigenvalues:

$$\begin{split} \mathbf{T}^{(0)} &= \mathbf{I} \\ \mathbf{T}^{(1)} &= \tilde{\Lambda} \\ \mathbf{T}^{(k)} &= 2 \cdot \tilde{\Lambda} \cdot \mathbf{T}^{(k-1)} - \mathbf{T}^{(k-2)} \end{split}$$



<sup>[1]</sup> M. Defferrard et al., "Convolutional neural networks on graphs with fast localized spectral filtering," 2016.

**GNN** libraries



#### References i

- [1] 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. 3844–3852.
- [2] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations (ICLR)*, 2016.
- [3] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," *arXiv preprint arXiv:1704.01212*, 2017.
- [4] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.
- [5] M. Simonovsky and N. Komodakis, "Dynamic edge-conditioned filters in convolutional neural networks on graphs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

- [6] J. You, R. Ying, and J. Leskovec, "Design space for graph neural networks," *arXiv* preprint arXiv:2011.08843, 2020.
- [7] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv preprint arXiv:1312.6203*, 2013.



# Pooling in Graph Neural Networks

Graph Deep Learning 2022

Daniele Grattarola February 28, 2022

# Pooling in CNNs





Things we are going to cover:

- $\cdot$  A "message passing" for pooling
- Methods
- Global pooling
- $\cdot$  Open questions

**Source**: "Understanding pooling in graph neural networks", Grattarola et al., 2021 https://arxiv.org/abs/2110.05292

- Graph: nodes connected by edges;
- A, adjacency matrix of shape  $N \times N$ ;
- $\mathbf{D} = \text{diag}([d_1, \dots, d_N])$ , diagonal degree matrix;
- · L = D A, Laplacian matrix;
- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]^{\top}$ ,  $\mathbf{x}_i \in \mathbb{R}^F$ , node attributes or "graph signal";
- $\mathbf{e}_{ij} \in \mathbb{R}^{S}$ , edge attribute for edge  $i \rightarrow j$ ;



## Graph pooling by example

Strategy 1: aggregate same attributes (Candy Crush pooling). Strategy 2: aggregate cliques. Strategy 3: keep only some types/colors.



- 1. How to identify groups of related nodes?
- 2. How to get new node attributes from the groups?
- 3. How to **connect** the new nodes?

<sup>[1]</sup> D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," 2021 (In preparation).

# Step 1: Select



Example 1: partition. $\{ \bigcirc \bigcirc \bigcirc \}$  $\{ \bigcirc \bigcirc \bigcirc \bigcirc \}$  $\{ \bigcirc \bigcirc \bigcirc \}$ 

Example 2: cover (possible overlaps). $\{ \bigcirc \bigcirc \bigcirc \bigcirc \}$  $\{ \bigcirc \bigcirc \bigcirc \}$  $\{ \bigcirc \bigcirc \bigcirc \}$  $\{ \bigcirc \bigcirc \bigcirc \}$ 

Example 3: sparse.  $\left\{ \bigcirc \right\} \left\{ \bigcirc \right\} \left\{ \bigcirc \right\} \left\{ \bigcirc \right\}$ 

The **selection** stage computes *K* **supernodes**:

 $SEL: \mathcal{G} \mapsto \mathcal{S} = \{\mathcal{S}_1, \dots, \mathcal{S}_K\}.$ 

Each supernode is a set of nodes associated with a score:

 $\mathcal{S}_k = \{ (\mathsf{x}_i, \mathsf{s}_i) \mid \mathsf{s}_i \in \mathbb{R}_{>0} \},\$ 



## Spectral clustering [3]



[2] J. Shi *et al.*, "Normalized cuts and image segmentation," 2000.[3] U. Von Luxburg, "A tutorial on spectral clustering," 2007.

The low-frequency eigenvectors naturally cluster the nodes.



**Idea:** run k-means clustering (or similar) using the first few eigenvectors.

## Node decimation [5]



Alternative: use the highest-frequency eigenvector to do something similar to a regular subsampling.



[5] F. M. Bianchi et al., Hierarchical Representation Learning in Graph Neural Networks with Node Decimation Pooling, 2019.

<sup>[4]</sup> L. Palagi et al., "Computational approaches to max-cut," 2012.

## Some problems



Problems with spectral methods:

- Computing eigenvectors is **expensive** (*O*(*N*<sup>3</sup>));
- They do not consider **attributes**.

But we get the general idea...

Step 2: Reduce

The **reduction** stage aggregates the supernodes in a **permutation-invariant** way:

 $\mathsf{Red}:\mathcal{G},\mathcal{S}_k\mapsto \mathbf{x}'_k$ 



Typical approach is to take a **weighted sum** (weights given by the scores in the supernodes):

$$X' = SX \ (\in \mathbb{R}^{K \times F})$$

Step 3: Connect

The **connection** function decides whether two supernodes are connected (and, in case, computes the associated attributes):

 $CON : \mathcal{G}, \mathcal{S}_k, \mathcal{S}_l \mapsto \mathbf{e}'_{kl}$ 



Typical approach is again to take a **weighted sum** of edges between two supernodes:

 $\mathsf{A}' = \mathsf{S}\mathsf{A}\mathsf{S}^{ op} \ \ (\in \mathbb{R}^{K imes K})$ 

## Select, Reduce, Connect [1]

Putting everything together:



<sup>[1]</sup> D. Grattarola et al., "Understanding Pooling in Graph Neural Networks," 2021 (In preparation).

Methods

A few ideas:

- 1. **Graclus** [6]: visit nodes randomly, merge pairs that maximize  $\frac{\mathbf{a}_{ij}}{w_i} + \frac{\mathbf{a}_{ij}}{w_j}$ , <sup>1</sup> In [7], they reduce supernodes with element-wise max.
- 2. Clique Pooling [8]: merge together cliques.
- 3. LaPool [9]: select "leaders" that have higher local variation ||LX|| w.r.t. all their neighbors. Create clusters by assigning nodes to nearest leader.

<sup>[6]</sup> I. S. Dhillon et al., "Weighted graph cuts without eigenvectors a multilevel approach," 2007.

<sup>[7]</sup> M. Defferrard et al., "Convolutional neural networks on graphs with fast localized spectral filtering," 2016.

<sup>[8]</sup> E. Luzhnica et al., "Clique pooling for graph classification," 2019.

<sup>[9]</sup> E. Noutahi et al., "Towards Interpretable Sparse Graph Representation Learning with Laplacian Pooling," 2019.

Key idea: learn to output  $S^{\top}$  by giving node features X as input to a neural network.

- **DiffPool** [10]: GNN for  $S^{\top}$ , regularize with "link prediction" loss;
- **MinCutPool** [11]: MLP for **S**<sup>⊤</sup>, regularize with "minimum cut" loss (same objective as spectral clustering);
- **Deep Graph Mapper** [12]: combine Mapper [13] and GCN [14] to compute clusters.



<sup>[10]</sup> R. Ying et al., "Hierarchical Graph Representation Learning with Differentiable Pooling," 2018.

<sup>[11]</sup> F. M. Bianchi *et al.,* "Mincut pooling in Graph Neural Networks," 2019.

<sup>[12]</sup> C. Bodnar et al., "Deep Graph Mapper: Seeing Graphs through the Neural Lens," 2020.

- Select:  $S^{\top} = MLP(X)$
- Reduce: X' = SX
- · Connect:  $\mathbf{A}' = \mathbf{S}\mathbf{A}\mathbf{S}^{\top}$
- MinCut loss:  $\mathcal{L}_c = -\frac{Tr(SAS^{\top})}{Tr(SDS^{\top})}$
- Orthogonality loss:

$$\mathcal{L}_{o} = \left\| \frac{\mathsf{S}\mathsf{S}^{\top}}{\|\mathsf{S}\mathsf{S}^{\top}\|_{F}} - \frac{\mathsf{I}_{K}}{\sqrt{K}} \right\|_{F}$$



<sup>[11]</sup> F. M. Bianchi et al., "Mincut pooling in Graph Neural Networks," 2019.

**Problem**: computing **S** with neural network is likely to yield a very **dense** matrix.

Can we learn a sparse selection?



## Top-K methods



Different ways of computing the selection indices i:

- Select with a simple linear projection  $\theta \in \mathbb{R}^{F}$  [15];
- Select with a GNN [16];
- Train the selection with a supervised objective (needs ground truth for which nodes to keep) [17].

<sup>[15]</sup> S. J. Hongyang Gao, "Graph U-Net," 2019.

<sup>[16]</sup> J. Lee et al., "Self-Attention Graph Pooling," 2019.

<sup>[17]</sup> B. Knyazev et al., "Understanding attention in graph neural networks," 2019.

Reduce:  $X' = X_i$  - Connect:  $A' = A_{i,i}$ 

Problems:

Top-k selection is non-differentiable (no way of training φ).
 Solved by gating (multiplying) the node attributes with

the scores.

 Graph is likely to be disconnected or simply cut off (like in the image on the right).
 Not really solvable...







- · Dense vs. Sparse: how many nodes are selected for the supernodes;
- · Fixed vs. Adaptive: how many supernodes does the selection compute;
- · Trainable vs. Non-trainable: learn to pool from data or not;

Global pooling

In CNNs, after convolution, we usually **flatten** out the images to give a vector as input to a MLP:

The graph equivalent must be **invariant to permutations** of the nodes:



1	2	3		
4	5	6		
7	8	9		

1	2	3	4	5	6	7	8	9

Once again, there are many ways to do this:

- Sum, average, product, max;
- Weighted sum with attention [18];
- Sum and then apply a neural network [19];

<sup>[18]</sup> Y. Li et al., "Gated graph sequence neural networks," 2015.

<sup>[19]</sup> N. Navarin et al., "Universal readout for graph convolutional neural networks," 2019.

Open questions

- Does pooling really work?
  - Dense selection + message passing + small graphs is a bad idea [20]
  - Which tasks benefit from pooling a priori?
  - Problems with inherent hierarchy?
- Can we make a pooling layer that is dense, trainable, and adaptive?

<sup>[20]</sup> D. Mesquita et al., "Rethinking pooling in graph neural networks," 2020.

#### References i

- [1] D. Grattarola, D. Zambon, F. M. Bianchi, and C. Alippi, "Understanding pooling in graph neural networks,", 2021 (In preparation).
- [2] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 8, pp. 888–905, 2000.
- [3] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, vol. 17, no. 4, pp. 395–416, 2007.
- [4] L. Palagi, V. Piccialli, F. Rendl, G. Rinaldi, and A. Wiegele, "Computational approaches to max-cut," in *Handbook on semidefinite, conic and polynomial optimization*, Springer, 2012, pp. 821–847.
- [5] F. M. Bianchi, D. Grattarola, L. Livi, and C. Alippi, *Hierarchical representation learning in graph neural networks with node decimation pooling*, 2019. arXiv: **1910.11436** [cs.LG].

### References ii

- [6] I. S. Dhillon, Y. Guan, and B. Kulis, "Weighted graph cuts without eigenvectors a multilevel approach," *IEEE transactions on pattern analysis and machine intelligence*, vol. 29, no. 11, pp. 1944–1957, 2007.
- [7] 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. 3844–3852.
- [8] E. Luzhnica, B. Day, and P. Lio, "Clique pooling for graph classification," International Conference of Learning Representations (ICLR) – Representation Learning on Graphs and Manifolds workshop, 2019.
- [9] E. Noutahi, D. Beani, J. Horwood, and P. Tossou, "Towards interpretable sparse graph representation learning with laplacian pooling," *arXiv preprint arXiv:1905.11577*, 2019.

## References iii

- [10] R. Ying, J. You, C. Morris, X. Ren, W. L. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," *arXiv preprint arXiv:1806.08804*, 2018.
- [11] F. M. Bianchi, D. Grattarola, and C. Alippi, "Mincut pooling in graph neural networks," CoRR, vol. abs/1907.00481, 2019. arXiv: 1907.00481. [Online]. Available: http://arxiv.org/abs/1907.00481.
- [12] C. Bodnar, C. Cangea, and P. Liò, "Deep graph mapper: Seeing graphs through the neural lens," *arXiv preprint arXiv:2002.03864*, 2020.
- [13] G. Singh, F. Mémoli, and G. E. Carlsson, "Topological methods for the analysis of high dimensional data sets and 3d object recognition.," in *SPBG*, 2007, pp. 91–100.
- [14] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in International Conference on Learning Representations (ICLR), 2016.

#### References iv

- [15] S. J. Hongyang Gao, "Graph u-net," Submitted to ICLR, 2019.
- [16] J. Lee, I. Lee, and J. Kang, "Self-attention graph pooling," *arXiv preprint arXiv:1904.08082*, 2019.
- [17] B. Knyazev, G. W. Taylor, and M. R. Amer, "Understanding attention in graph neural networks," CoRR, vol. abs/1905.02850, 2019. arXiv: 1905.02850. [Online]. Available: http://arxiv.org/abs/1905.02850.
- [18] Y. Li, D. Tarlow, M. Brockschmidt, and R. Zemel, "Gated graph sequence neural networks," *arXiv preprint arXiv:1511.05493*, 2015.
- [19] N. Navarin, D. Van Tran, and A. Sperduti, "Universal readout for graph convolutional neural networks," in 2019 International Joint Conference on Neural Networks (IJCNN), IEEE, 2019, pp. 1–7.

[20] D. Mesquita, A. Souza, and S. Kaski, "Rethinking pooling in graph neural networks," Advances in Neural Information Processing Systems, vol. 33, 2020.