

GRASSMANN GRAPH EMBEDDING

Bingxin Zhou*

The University of Sydney
NSW 2006, Australia
bzho3923@uni.sydney.edu.au

Xuebin Zheng*

The University of Sydney
NSW 2006, Australia
xuebin.zheng@sydney.edu.au

Yu Guang Wang

Max Planck Institute for Mathematics in the Sciences
Leipzig 04103, Germany
yuguang.wang@mis.mpg.de

Ming Li

Zhejiang Normal University
Zhejiang 321004, China
mingli@zjnu.edu.cn

Junbin Gao

The University of Sydney
NSW 2006, Australia
junbin.gao@sydney.edu.au

ABSTRACT

Geometric deep learning that employs the geometric and topological features of data has attracted increasing attention in deep neural networks. Learning the intrinsic structure property of data is a crucial step for dimensionality reduction and effective feature extraction. This paper develops *Grassmann graph embedding*, which combines graph convolutions to capture the main components within graphs' hidden representations. Each set of featured graph nodes is mapped to a point on a Grassmann matrix manifold through Singular Value Decomposition, which is then embedded into a symmetric matrix space that approximates denoised second-order feature information. The view of treating nodes as a set could inspire many potential applications. In particular, we propose Grassmann (global graph) pooling that can connect with any graph convolution for graph neural networks. The Grassmann pooling achieves state-of-the-art performance on a variety of graph prediction benchmarks.

1 INTRODUCTION

Graphs are prevalent objects in machine learning and data science. Graph representation learning has recently attracted increasing attention with graph neural networks (GNNs) emerging as one of the most prominent avenues in deep learning (Bronstein et al., 2017; Hamilton, 2020). One critical ingredient of graph representation like GNNs is to effectively distill key features of graph data. Graph convolution (Bruna et al., 2014) to GNN, or more generally, graph neural message passing (Gilmer et al., 2017), provides an efficient way to extract core information of a graph. In node prediction tasks, graph convolution is an excellent learner that aggregates graph features and uses node relation to predict node labels. On the contrary, graph-level learning tasks are usually defined on graph data with varying size and structure, in which case dimensionality reduction (for graph size) is often an indispensable necessity. To this end, a graph pooling scheme is critical to obtaining scaled-down graphs or graph-level embeddings. For example, TOPKPOOL (Cătălina et al., 2018) exploits a node selection criterion by a learned score and drops all but the top-ranked nodes; SAGPOOL (Lee et al., 2019) calculates a score with an extra MLP in the TOPK framework to enhance the selection efficacy; and ATTENTIONPOOL (Li et al., 2016) defines a soft scoring mechanism that concatenates node state and annotation.

One basic property of a pooling or embedding operation is node permutation invariance, as the node order of an undirected graph must not incline the network to focus on specific information. The

*equal contribution first authors.

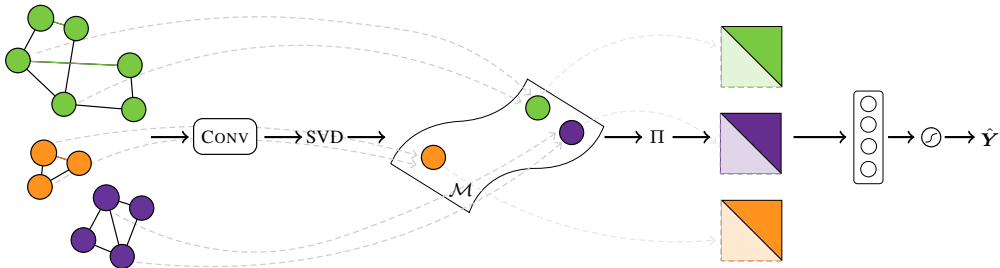


Figure 1: Computational principle of the proposed *Grassmann Graph Embedding*. A complete description is provided in Section 1. This work implements the strategy and constructs Grassmann pooling (GRPOOL).

permutation invariance is intuitive in heuristic pooling operations like summation, averaging, and maximization, where the aggregation rules do not rely on the order information of nodes. Indeed we can define such operations on subsets of a graph. The view of treating nodes as a set in pooling can be seen in the recent work of Kolouri et al. (2021), where a node set finds its Wasserstein Embedding within a linear optimal transport framework, and match each pair of nodes by a proper Wasserstein metric (Mémoli, 2011).

When graph nodes come with features, we may view the set of graph nodes from the perspective of manifold learning. The set of feature vectors naturally generates a subspace as a new representation. In other words, the node information of a graph can be abstracted into a lower-dimensional subspace, and a learning task over graphs of different size and structure can thus be converted into a learning task on a Grassmann manifold (Absil et al., 2008). We name this process of mapping a graph onto the Grassmann manifold as *Grassmann (graph) embedding*. The acquired Grassmann embeddings from graphs can be pipelined into a Grassmann learning algorithm, such as Deep Grassmann Networks (Huang et al., 2018) and Grassmann clustering (Wang et al., 2014; 2017).

As a proof of concept, this work develops a simple manifold embedding strategy (Figure 1). An input graph is first passed through Graph Convolutional layers to extract the hidden representation of its node information, which is then embedded as a Grassmann point with a subspace dimension by Singular Value Decomposition (SVD). All Grassmann points that represent corresponding graphs are on the same Grassmann manifold, and they can be projected to Euclidean space in the form of symmetric matrices (Absil et al., 2008), where all graphs acquire the same size representation. The vectorized representations can be sent to a conventional learning model for any learning task such as graph classification or clustering. The architecture design is one practice of exploiting the Grassmann representation as a pooling or embedding component. Based on this, we propose Grassmann Graph Pooling (GRPOOL), a feasible graph pooling mechanism for graph data in varying size and structure. The strategy can also inspire other learning tasks, such as graph coarsening and graph clustering.

2 GRASSMANN GRAPH EMBEDDING

We can embed multiple graphs to one Grassmann manifold using their graph convolutional network representation. For a given set of N graphs $\mathbb{G} = \{\mathcal{G}_1, \dots, \mathcal{G}_N\}$, a sequence of graph convolutional layers like GCN (Kipf & Welling, 2017) and GIN (Xu et al., 2019) maps them to a set of hidden representations $\mathbb{H} = \{\mathbf{H}_1, \dots, \mathbf{H}_N\}$ where $\mathbf{H}_i \in \mathbb{R}^{n_i \times m}$ is with respect to n_i nodes in \mathcal{G}_i , and m denotes the number of hidden units in the last graph convolutional layer. To conduct graph-level learning tasks, one needs to find a uni-dimensional graph representation that is independent of the node size. We accomplish the mission by transforming each graph representation of \mathbb{H} to a Grassmann point, which can later be projected back to its Euclidean form for regular output layers.

Grassmann manifold A Grassmann manifold $\mathcal{M}(d, m)$ is the space of all d -dimensional linear subspaces of \mathbb{R}^m for $0 \leq d \leq m$. A point on $\mathcal{M}(d, m)$ is a d -dimensional subspace of \mathbb{R}^m , which can be represented by any orthonormal basis $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d] \in \mathbb{R}^{m \times d}$. The chosen orthonormal basis is called a representative of its subspace. The Grassmann manifold is an abstract

Table 1: Performance comparison for graph property prediction. **QM7** is a regression task in MSE; **ogbg-molhiv** is a classification task in ROC-AUC in percentage; others are for classification in test accuracy in percentage. The value after \pm is standard deviation.

		PROTEINS	D&D	NCII	Mutagenicity	COLLAB	ogbg-molhiv	QM7
GCN Conv	TopKPool	73.48 \pm 3.57	74.87 \pm 4.12	75.11 \pm 3.45	79.84 \pm 2.46	81.18 \pm 0.89	77.11 \pm 1.27	175.41 \pm 3.16
	Attention	73.93 \pm 5.37	77.48 \pm 2.65	74.04 \pm 1.27	80.25 \pm 2.22	81.58 \pm 1.72	77.44 \pm 1.27	177.99 \pm 2.22
	SAGPool	75.89 \pm 2.91	74.96 \pm 3.60	76.30 \pm 1.53	79.86 \pm 2.36	79.26 \pm 5.37	75.36 \pm 1.82	41.93 \pm 1.14
	SUM	74.91 \pm 4.08	78.91 \pm 3.37	76.96 \pm 1.70	80.69 \pm 3.26	80.76 \pm 1.56	74.88 \pm 2.64	42.09 \pm 0.91
	AVG	73.13 \pm 3.18	76.89 \pm 2.23	73.70 \pm 2.55	80.37 \pm 2.44	81.24 \pm 1.34	77.69 \pm 1.17	177.49 \pm 4.69
	MAX	73.57 \pm 3.94	75.80 \pm 4.11	75.96 \pm 1.82	78.83 \pm 1.70	82.28 \pm 2.10	76.95 \pm 0.94	177.48 \pm 4.70
	GrPool	77.79 \pm 2.16	79.10 \pm 2.98	77.80 \pm 2.01	81.01 \pm 1.28	82.94 \pm 1.06	76.60 \pm 1.10	172.35 \pm 3.50
GIN Conv	TopKPool	73.66 \pm 6.00	76.40 \pm 2.32	77.06 \pm 0.90	78.30 \pm 1.36	81.40 \pm 0.94	78.14 \pm 0.62	42.68 \pm 0.50
	Attention	75.63 \pm 1.13	71.76 \pm 3.26	78.22 \pm 1.32	78.54 \pm 5.37	83.22 \pm 0.30	74.44 \pm 2.12	42.67 \pm 0.80
	SAGPool	75.95 \pm 4.52	68.94 \pm 7.62	76.97 \pm 2.94	78.86 \pm 1.58	81.76 \pm 1.57	75.26 \pm 2.29	43.02 \pm 1.54
	SUM	78.04 \pm 2.30	78.57 \pm 1.26	78.83 \pm 1.49	81.31 \pm 1.10	82.64 \pm 0.85	77.41 \pm 1.16	42.41 \pm 3.36
	AVG	71.70 \pm 2.08	74.37 \pm 1.32	76.55 \pm 1.72	80.97 \pm 1.18	83.30 \pm 0.77	78.21 \pm 0.90	43.16 \pm 4.77
	MAX	76.70 \pm 1.57	77.31 \pm 2.06	79.27 \pm 1.38	80.28 \pm 0.83	80.94 \pm 0.72	78.16 \pm 1.33	42.54 \pm 5.17
	GrPool	79.80 \pm 1.09	81.18 \pm 1.14	81.31 \pm 1.55	82.53 \pm 0.72	81.32 \pm 0.68	77.82 \pm 0.90	42.36 \pm 3.39

quotient manifold that can be represented in many ways. In this paper, we embed the Grassmann manifold into the space of symmetric matrices $\text{Sym}(m)$ by

$$\Pi : \mathcal{M}(d, m) \longrightarrow \text{Sym}(m), \quad \Pi(\mathbf{U}) = \mathbf{U}\mathbf{U}^\top. \quad (1)$$

This method is called *Grassmann embedding*, which we now introduce.

Rectified manifold embedding Suppose we have obtained a graph hidden representation $\mathbf{H} \in \mathbb{H} \subset \mathbb{R}^{n \times m}$ and its row-generated subspace $\text{span}(\mathbf{H}^\top)$, which as mentioned can be achieved by applying one or multiple layers of graph convolution like GCN. We then go to find a better Grassmann representation for the subspace to extract the most important features of the data. Such representation is characterized by an orthogonal basis of the subspace, which can be rectified in several ways. A typical method as demonstrated in Huang et al. (2018) is QR decomposition of matrix \mathbf{H}^\top . Our preliminary goal is to bring representation towards a unique subspace dimension k . That is, we need to find out the most representative basis $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d]$, thus $\mathcal{U} = [\mathbf{U}]$, an equivalence class of \mathbf{U} , is a point of an underlying Grassmann manifold \mathcal{M} . To this end, we employ singular value decomposition (SVD) for the hidden feature

$$\mathbf{H}^\top = \mathbf{U}\mathbf{S}\mathbf{V}^\top, \quad (2)$$

where $\mathbf{U} \in \mathbb{R}^{m \times k}$ with $\text{rank}(\mathbf{H}^\top) = k, 1 \leq k \leq \min\{m, n\}$ is an orthonormal basis that spans the column space of \mathbf{H}^\top . The diagonal $\mathbf{S} := \text{Diag}([\sigma_1, \dots, \sigma_k]) \in \mathbb{R}^{k \times k}$ contains k singular values sorted in the descending order, where the l th singular value indicates the percentage importance of the l th basis vector. The corresponding singular vectors constitute $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_k] \in \mathbb{R}^{n \times k}$.

The first d -columns of \mathbf{U} ($d \leq k$) from SVD include the d most important components in the original space of \mathbf{H} , and the subspace it spans naturally composes a Grassmann point $\mathcal{U} := [\mathbf{U}_d]$ (the equivalent class under the orthogonal group $\mathcal{O}(d)$) on the Grassmann manifold $\mathcal{M}(d, M)$. This embedding operation can be applied to any hidden representation from \mathbb{H} , which intuitively regards a graph with a set of n featured nodes as a d -dimensional subspace of the embedding space \mathbb{R}^m . Each graph feature is transformed to a hidden feature, and different graphs (possibly with varying node size) are taken as points on the (same) Grassmann manifold, where similar instances are also geodesically close to each other. Such similarity provides a criterion for clustering or classification. At this point, we have finished embedding graphs to the Grassmann manifold where each graph is represented by a subspace of orthogonal basis $\mathbf{U}_d \in \mathbb{R}^{m \times d}$.

Projection mapping With the above Grassmann embedding, we have a set of embedded Grassmann points $\{\mathcal{U}_1, \dots, \mathcal{U}_N\} \subset \mathcal{M}(d, M)$ with their relevant orthogonal matrix basis as representatives. Any follow-up learning can then be built on $\{\mathcal{U}_1, \dots, \mathcal{U}_N\}$ by the Grassmann geometry, see Absil et al. (2008). As pointed out in (1), the entire manifold $\mathcal{M}(d, M)$ can be embedded (mapped) into the Euclidean space $\text{Sym}(m)$. This embedding naturally transforms graph features to their Euclidean representation and thus conventional deep learning can be applied.

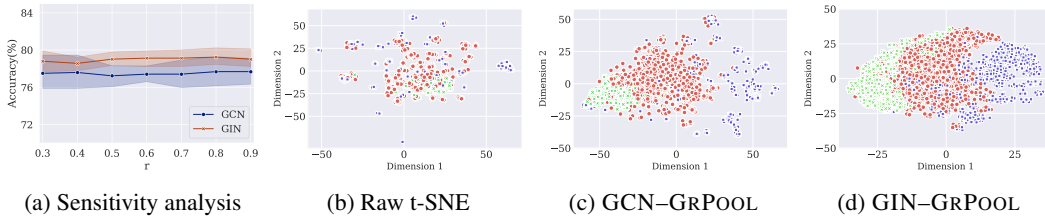


Figure 2: (a) Sensitivity analysis for the threshold information ratio r on **PROTEINS**. The plots (b)-(d) present the t-SNE visualizations of the vectorized graph representations produced by our GRPOOL on **COLLAB** with GCN before training; with GCN after training; and with GIN after training, respectively.

To be precise, we project the graph latent representation to the symmetric positive definite (SPD) matrix $\Pi(U_d) = U_d U_d^\top$. The resulting SPD matrix representation, an analog to a bilinear mapping, captures a second-order statistics that can better reflect regional features (Tuzel et al., 2006). Moreover, the rectified representation U from (2) can reduce noise, and the approximated covariance matrix $U_d U_d^\top$ is thus more expressive. This projected Euclidean representation can be used for various tasks. For example, in graph property prediction, the Grassmann embedding can be vectorized and output as the readout for graph prediction. In subspace clustering, $\Pi(U_d)$ with respect to H is analogous to the affinity matrix and it can be used by clustering algorithms for graph segmentation.

Remark 1. *The SPD matrices from the projection mapping are points on the positive semidefinite (PSD) cone. However, we do not leverage this special topology here when calculating the pair-wise distance. We measure the distance between Grassmann points with Euclidean metrics, which only requires the vector representation of $m(m+1)/2$ dimensions.*

3 EXPERIMENTS FOR GRASSMANN POOLING

We evaluate the proposed GRPOOL on a set of graph classification tasks with different graph sizes, volume, and density. We also include a graph regression task with moderate data size.

We benchmark the performance on five binary classification, one multi-class classification and one regression tasks. Except for **ogbg-molhiv** (Hu et al., 2020), all datasets are provided by **TUdataset benchmarks** (Morris et al., 2020) with default graph attributes. The performance is compared against two hierarchical (TOPKPOOL (Gao & Ji, 2019; Cătălina et al., 2018), SAGPOOL (Lee et al., 2019)) and four global pooling methods (ATTENTIONPOOL (Li et al., 2016), SUM, AVG and MAX pooling) with two types of network architectures. For completeness, we report prediction results of all models with GCN (Kipf & Welling, 2017) and GIN (Xu et al., 2019) (under JKNET structure (Xu et al., 2018)) convolutional layers. We detail the model setup in Appendix C.

Table 1 compares the prediction performance of GRPOOL with baselines. We report mean test accuracy for the first five classification tasks; ROC-AUC score for **ogbg-molhiv** classification; and mean square error (MSE) for **QM7** regression. All the mean scores are averaged over 10 repetitions. Our GRPOOL achieves the top score with generally a lower volatility on all tasks with primary feature matrices. The advantage is more salient on the bottom half of the table, where we apply JKNET structure with GIN convolution for GRPOOL. It provides a global pooling method and can express more information from its compressed graph representation. The performance of GRPOOL on **COLLAB** and **QM7** is less promising, which is likely due to the lack of non-structural feature information from the dataset. In general, the abstracted subspaces provide more information compared to the baseline methods with common graph convolutional layers especially when there are feature matrices on graphs or nodes.

4 FURTHER INVESTIGATION ON GRASSMANN EMBEDDING

Real datasets could have a great number of graphs with a large variation on node sizes from a few to thousands. Consequently, defining an identical relatively small subspace dimension d for all graphs could potentially hurt the expressiveness of graph embedding. Instead, we let d of a Grassmann point $[U_d] \in \mathcal{M}(m, d)$ be determined by $d = \sum_{i=1}^k \mathbb{1}\{\sigma_i > r\}$. The σ_i corresponds to the singular

values from (2), and r denotes the global threshold of the percentage importance. However, all these Grassmannian points $\{\mathcal{U}\}$ can be naturally mapped to the embedded Euclidean space $\text{Sym}(m)$ of the manifold $\mathcal{M}(m, d_{\max})$, with d_{\max} the highest d over all Grassmann points.

The r is a hyperparameter of the GRPOOL. However, the choice of r has little impact on the performance of the Grassmann embedding, as demonstrated by our sensitivity analysis. The experiment is conducted on **PROTEINS** with both GCN and GIN network architectures from Section 3. We report the mean test accuracy over 10 repetitions with seven different values of the threshold ratio r , ranging from 0.3 to 0.9 with step size 0.1. All other hyperparameters are tuned in the same way as Section 3. Figure 2a shows that the mean test accuracy is stable over all values of the threshold ratio r . This suggests that the hyperparameter r (within a wide range) has a negligible impact on the performance of our proposed GRPOOL, which indicates that a moderately high value of r (e.g., $r > 0.5$) is sufficient to retain the essential information of a graph for Grassmann embedding.

The rest three plots in Figure 2 exploit two-dimensional t-distributed Stochastic Neighbor Embedding (t-SNE) of flattened Euclidean graph embeddings on **COLLAB**, which is a 3-classification task. Each point denotes a graph representation after GRPOOL, and their color indicates the true labels. We randomly sampled a sufficient amount of 4,000 instances. As the GIN-GRPOOL model employs JKNET structure, we aggregate all four pooling results as an approximated graph representation. Both Figures 2c and 2d suggest a clear clustering pattern of the pooled graphs.

5 DISCUSSION AND CONCLUSION

In this paper, we propose GRPOOL, a new graph pooling strategy that connects with the Grassmann geometry. The underlying method treats the hidden feature subspace of graphs as Grassmann points, and make analysis on them. The process of extracting subspace of graph features is a non-linear transformation that outputs the most valuable feature information, and the projection embedding makes second-order approximation on feature relationship. The result graph representation supports Euclidean metrics for loss design. Such strategy can be used for various scenarios. In specific, we propose GRPOOL that pools graph hidden representations to identical vectors which we employ for label prediction. In this way, non-linear analysis on the graph features is defined, which is useful for working on complex or massive attributes. Moreover, the way of treating regional graphs as Grassmann points might motivate more connections of mutual learning schemes between the two non-Euclidean spaces of graph and manifold.

REFERENCES

- P.-A. Absil, R. Mahony, and R. Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- Michael M Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann Lecun. Spectral networks and locally connected networks on graphs. In *ICLR*, 2014.
- Cangea Cătălina, Petar Veličković, Nikola Jovanović, Thomas Kipf, and Pietro Liò. Towards sparse hierarchical graph classifiers. In *NeurIPS Workshop on Relational Representation Learning*, 2018.
- Hongyang Gao and Shuiwang Ji. Graph U-nets. In *ICML*, 2019.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In *NeurIPS*, volume 70, pp. 1263–1272, 2017.
- William L Hamilton. *Graph Representation Learning*. Morgan & Claypool Publishers, 2020.
- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. In *NeurIPS*, 2020.

- Zhiwu Huang and Luc Van Gool. A Riemannian network for SPD matrix learning. In *AAAI*, volume 31, 2017.
- Zhiwu Huang, Jiqing Wu, and Luc Van Gool. Building deep networks on grassmann manifolds. In *AAAI*, volume 32, 2018.
- Katsuhiko Ishiguro, Shin-ichi Maeda, and Masanori Koyama. Graph warp module: an auxiliary module for boosting the power of graph neural networks. *arXiv:1902.01020*, 2019.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- Soheil Kolouri, Navid Naderializadeh, Gustavo K. Rohde, and Heiko Hoffmann. Wasserstein embedding for graph learning. In *ICLR*, 2021.
- Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. In *ICML*, 2019.
- Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. In *ICLR*, 2016.
- Facundo Mémoli. Gromov-Wasserstein distances and the metric approach to object matching. *Foundations of Computational Mathematics*, 11:417–487, 2011.
- Christopher Morris, Nils M Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion Neumann. TUDataset: A collection of benchmark datasets for learning with graphs. In *ICML Workshop Graph Representation Learning and Beyond*, 2020.
- Shayle R Searle and Andre I Khuri. *Matrix algebra useful for statistics*. John Wiley & Sons, 2017.
- Oncel Tuzel, Fatih Porikli, and Peter Meer. Region covariance: A fast descriptor for detection and classification. In *ECCV*, pp. 589–600. Springer, 2006.
- Boyue Wang, Yongli Hu, Junbin Gao, Yanfeng Sun, and Baocai Yin. Low rank representation on Grassmann manifolds. In *Asian Conference on Computer Vision (ACCV)*, 2014.
- Boyue Wang, Yongli Hu, Junbin Gao, Yanfeng Sun, Haoran Chen, Muhammad Ali, and Baocai Yin. Locality preserving projections for Grassmann manifold. In *IJCAI*, pp. 2893–2900, 2017.
- Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *ICML*, 2018.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *ICLR*, 2019.
- Xuebin Zheng, Bingxin Zhou, Ming Li, Yu Guang Wang, and Junbin Gao. MathNet: Haar-like wavelet multiresolution-analysis for graph representation and learning. *arXiv:2007.11202*, 2020.
- Xuebin Zheng, Bingxin Zhou, Junbin Gao, Yu Guang Wang, Pietro Liò, Ming Li, and Guido Montúfar. How framelets enhance graph neural networks. *arXiv:2102.06986*, 2021.

A STABLE BACK-PROPAGATION FOR SVD

As Singular Value Decomposition (SVD) is one of the essential parts of our proposed *Grassmann graph embedding*, we present in this section the derivation for the back-propagation for SVD. This back-propagation is numerically stable, especially in the case when the input matrix has some extremely small singular values. The following materials are mainly based on the note of Townsend¹.

For a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with rank $k \leq \min(m, n)$, we decompose $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$ by using SVD, where \mathbf{U} is a $m \times k$ matrix, \mathbf{V} is in $n \times k$, and $\mathbf{S} = \text{diag}(s_1, s_2, \dots, s_k)$ is a $k \times k$ diagonal matrix which stores the singular values of \mathbf{A} . Here, we have the constraints that the matrices \mathbf{U} and \mathbf{V} are orthogonal. We denote the objective function by $f(\mathbf{A})$ whose argument is a square matrix \mathbf{A} . We aim to find the gradient of the objective function: $\bar{\mathbf{A}} := \nabla_{\mathbf{A}} f$. For simplicity, we employ the shorthand $\bar{\mathbf{A}} := \nabla_{\mathbf{A}} f$ for the derivation below. By the property of matrix algebra (Searle & Khuri, 2017), we reach the relation

$$df(\mathbf{A}) = \text{tr}(\bar{\mathbf{A}}^\top d\mathbf{A}) = \text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) + \text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S}) + \text{tr}(\bar{\mathbf{V}}^\top d\mathbf{V}). \quad (3)$$

To calculate $\bar{\mathbf{A}}$, we need the following three differentials

$$d\mathbf{U} = \mathbf{U} (\mathbf{F} \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}]) + (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1} \quad (4)$$

$$d\mathbf{S} = \mathbf{I}_k \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}] \quad (5)$$

$$d\mathbf{V} = \mathbf{V} (\mathbf{F} \circ [\mathbf{S}\mathbf{U}^\top d\mathbf{A}\mathbf{V} + \mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}\mathbf{S}]) + (\mathbf{I}_n - \mathbf{V}\mathbf{V}^\top) d\mathbf{A}^\top \mathbf{U}\mathbf{S}^{-1}, \quad (6)$$

where $\mathbf{F}_{ij} = \frac{1}{s_j^2 - s_i^2} \cdot \mathbb{1}\{i \neq j\}$ which satisfies the identity $\mathbf{F}^\top = -\mathbf{F}$. The calculations of \mathbf{F} and \mathbf{S}^{-1} are often numerically unstable due to the possible near-zero singular values. To circumvent this difficulty, we employ the following improved strategy, as similarly used by Huang & Van Gool (2017):

$$\mathbf{S}_{i,i}^{\text{new}} = \mathbf{S}_{i,i} \cdot \mathbb{1}\{\mathbf{S}_{i,i} > \epsilon\} + \epsilon \cdot \mathbb{1}\{\mathbf{S}_{i,i} \leq \epsilon\},$$

where ϵ is a small number and can usually be set to 10^{-12} . In practice, we can replace \mathbf{S} by the modified matrix \mathbf{S}^{new} in the back-propagation formulas.

We evaluate the three terms in (3) respectively, as follows. By (5), we can write $\text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S})$ as

$$\begin{aligned} \text{tr}(\bar{\mathbf{S}}^\top d\mathbf{S}) &= \text{tr}(\bar{\mathbf{S}}^\top (\mathbf{I}_k \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}])) \\ &= \text{tr}(\mathbf{U}^\top d\mathbf{A}\mathbf{V} (\mathbf{I}_k \circ \bar{\mathbf{S}})) \\ &= \text{tr}(\mathbf{V} (\mathbf{I}_k \circ \bar{\mathbf{S}}) \mathbf{U}^\top d\mathbf{A}). \end{aligned} \quad (7)$$

Using (4) for the term $\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U})$ yields

$$\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) = \text{tr}(\bar{\mathbf{U}}^\top [\mathbf{U} (\mathbf{F} \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}]) + (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1}]).$$

The above expression is a sum of two terms. By $\mathbf{F}^\top = -\mathbf{F}$, these two terms can be simplified as

$$\begin{aligned} &\text{tr}(\bar{\mathbf{U}}^\top \mathbf{U} (\mathbf{F} \circ [\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}])) \\ &= \text{tr}([\mathbf{U}^\top d\mathbf{A}\mathbf{V}\mathbf{S} + \mathbf{S}\mathbf{V}^\top d\mathbf{A}^\top \mathbf{U}] (\mathbf{F} \circ \mathbf{U}^\top \bar{\mathbf{U}})) \\ &= \text{tr}(\mathbf{V}\mathbf{S} (\mathbf{F} \circ \mathbf{U}^\top \bar{\mathbf{U}}) \mathbf{U}^\top d\mathbf{A} - \mathbf{V}\mathbf{S} (\mathbf{F} \circ \bar{\mathbf{U}}^\top \mathbf{U}) \mathbf{U}^\top d\mathbf{A}) \\ &= \text{tr}(\mathbf{V}\mathbf{S} (\mathbf{F} \circ [\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U}]) \mathbf{U}^\top d\mathbf{A}) \end{aligned}$$

and

$$\text{tr}(\bar{\mathbf{U}}^\top (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top) d\mathbf{A}\mathbf{V}\mathbf{S}^{-1}) = \text{tr}(\mathbf{V}\mathbf{S}^{-1} \bar{\mathbf{U}}^\top (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top) d\mathbf{A}).$$

Thus, an equivalent expression for the term $\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U})$ reads

$$\text{tr}(\bar{\mathbf{U}}^\top d\mathbf{U}) = \text{tr}(\mathbf{V} [\mathbf{S} (\mathbf{F} \circ [\mathbf{U}^\top \bar{\mathbf{U}} - \bar{\mathbf{U}}^\top \mathbf{U}]) \mathbf{U}^\top + \mathbf{S}^{-1} \bar{\mathbf{U}}^\top (\mathbf{I}_m - \mathbf{U}\mathbf{U}^\top)] d\mathbf{A}). \quad (8)$$

¹<https://j-towns.github.io/papers/svd-derivative.pdf>

Table 2: Summary of the datasets for the graph property prediction tasks.

Datasets	PROTEINS	D&D	NCI1	Mutagenicity	COLLAB	ogbg-molhiv	QM7
# graphs	1,113	1,178	4,110	4,337	5,000	41,127	7,165
# classes	2	2	2	2	3	2	1 (R)
Min # nodes	4	30	3	30	32	2	4
Max # Nodes	620	5,748	111	417	492	222	23
Avg # nodes	39	284	30	30	74	26	15
Avg # edges	73	716	32	31	2,458	28	123
# Features	3	89	37	14	0	9	0

Similarly, the term $\text{tr}(\bar{\mathbf{V}}^\top d\mathbf{V})$ has the following form after some simple algebraic manipulations

$$\text{tr}(\bar{\mathbf{V}}^\top d\mathbf{V}) = \text{tr}\left(\left[\mathbf{V}\left(\mathbf{F}\circ\left[\mathbf{V}^\top\bar{\mathbf{V}}-\bar{\mathbf{V}}^\top\mathbf{V}\right]\right)\mathbf{S}+\left(\mathbf{I}_m-\mathbf{V}\mathbf{V}^\top\right)\bar{\mathbf{V}}\mathbf{S}^{-1}\right]\mathbf{U}^\top d\mathbf{A}\right). \quad (9)$$

Finally, we substitute (7)–(9) into (3) and then retrieve the expression for \mathbf{B} to obtain the solution

$$\begin{aligned} \mathbf{B} = & \left[\mathbf{U}\left(\mathbf{F}\circ\left[\mathbf{U}^\top\bar{\mathbf{U}}-\bar{\mathbf{U}}^\top\mathbf{U}\right]\right)\mathbf{S}+\left(\mathbf{I}_m-\mathbf{U}\mathbf{U}^\top\right)\bar{\mathbf{U}}\mathbf{S}^{-1}\right]\mathbf{V}^\top + \\ & \mathbf{U}\left(\mathbf{I}_k\circ\bar{\mathbf{S}}\right)\mathbf{V}^\top + \mathbf{U}\left[\mathbf{S}\left(\mathbf{F}\circ\left[\mathbf{V}^\top\bar{\mathbf{V}}-\bar{\mathbf{V}}^\top\mathbf{V}\right]\right)\mathbf{V}^\top + \mathbf{S}^{-1}\bar{\mathbf{V}}^\top\left(\mathbf{I}_n-\mathbf{V}\mathbf{V}^\top\right)\right] = \bar{\mathbf{A}}. \end{aligned}$$

B ESSENTIAL REQUIREMENTS FOR GRAPH POOLING

In this section, we check the two essential requirements of a valid graph pooling method, which are *graph embedding with a unified size* and *permutation invariance*.

Proposition 1. GRPOOL always produces a graph embedding $g \in \mathbb{R}^{\frac{m(m+1)}{2}}$ for the node representation matrix $\mathbf{H} \in \mathbb{R}^{n \times m}$, regardless of the graph size n .

Proof. Given the node representation matrix $\mathbf{H} \in \mathbb{R}^{n \times m}$ of a graph \mathcal{G} with n nodes and m features, the Grassmann graph embedding gives the output $\mathbf{U}\mathbf{U}^\top$, where $\mathbf{U} \in \mathbb{R}^{m \times k}$ with $k = \text{rank}(\mathbf{H}^\top)$. Then, the output of GRPOOL is the flattened representation of the upper triangular matrix of $\mathbf{U}\mathbf{U}^\top \in \mathbb{R}^{m \times m}$, which is independent of the graph size n . \square

Proposition 2. GRPOOL satisfies the requirement of permutation invariance so that it produces the same Grassmann graph embedding under row permutations of the input node representation matrix.

Proof. Suppose \mathbf{H}_1 is the node representation matrix of a graph and let $\mathbf{H}_2 = \mathbf{P}\mathbf{H}_1$, where \mathbf{P} is a permutation matrix. Then,

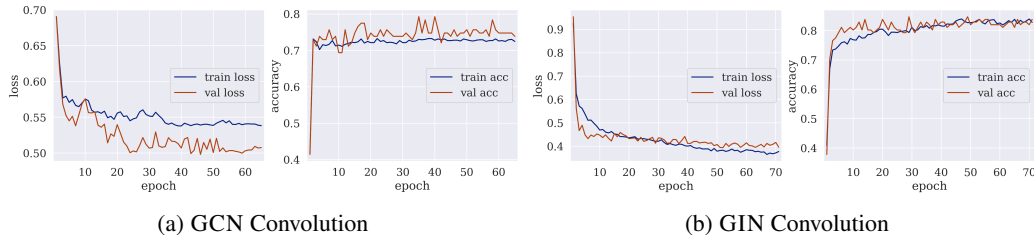
$$\mathbf{H}_1^\top = \mathbf{U}\mathbf{S}\mathbf{V}^\top, \quad \mathbf{H}_2^\top = \mathbf{H}_1^\top\mathbf{P}^\top = \mathbf{U}\mathbf{S}\mathbf{V}^\top\mathbf{P}^\top.$$

The Grassmann point for both \mathbf{H}_1 and \mathbf{H}_2 can be accessed by the same matrix \mathbf{U} . Hence, the proposed graph embedding method is permutation invariant. \square

C EXPERIMENT

C.1 MORE ON MODEL TRAINING

We supplement below a summary of the datasets for graph property prediction tasks in Section 3. The ‘R’ in the bracket of the last row of **QM7** represents the regression task. For the two feature-less datasets **COLLAB** and **QM7**, we follow Zheng et al. (2020) to create the constant uninformative node features for **QM7**; and follow Xu et al. (2019) to apply one-hot encoding of node degrees for **COLLAB**. For **ogbg-molhiv**, virtual nodes (Ishiguro et al., 2019) are used in the GNNs.

Figure 3: Training and Validation learning curves on **PROTEINS** with GRPOOL.

In the ablation, we consider two graph convolutions to be connected with pooling methods: GCN (Kipf & Welling, 2017) and GIN (Xu et al., 2019). For **TU Datasets**, the first set connects two GCN convolutional layers with one pooling layer except for **NCI1**, where we add a third convolutional layer. We also study the performance of the second GNN architecture with four GIN convolutional layers and JKNET structure (Xu et al., 2018). For the **ogbg-molhiv**, both models use four convolutional layers. The hidden units of the one-layer MLP for baselines are set identical to that of the convolutional layer. For GRPOOL, it is set to 64 for the first layer and 16 for the second.

Below we provide the searching space of the key hyperparameters that are fine-tuned with grid search, including learning rate, weight decay, hidden unit numbers of convolutional layers, and dropout ratio in MLP layers. Unless stated, models on **OGB** applies the same searching space. Other hyperparameters, if not specified, are set to default values.

For model training and evaluation, we split each dataset into 80% training, 10% validation, and 10% test samples. The stopping criterion is defined as the validation loss stops improving for 20 consecutive epochs, or the maximum 200 epochs is reached. In Table 1, except for those on **COLLAB**, the results of baselines with GCN on **TU Datasets** are retrieved from Zheng et al. (2021).

C.2 LEARNING CURVE AND CONVERGENCE

The training and validation curves for loss and accuracy are shown in Figure 3. The results are retrieved from a (random) single run on **PROTEINS**. The hyperparameters and network architectures follow exactly the same routine as Section 3. We only display results from one repetition due to the employment of early stopping criteria, where each independent run would stop at an individual epoch. The minor volatility after epoch 5 is generally brought about by stochastic gradients. All four visualizations validate an efficient convergence of our GRPOOL, where the loss curve stabilizes quickly after a few epochs. The training process is slightly longer for GIN convolution, which is partly due to the more sophisticated network architecture for the model to fit.

Table 3: Search space for hyperparameter tuning.

Hyperparameter	Searching Space
Learning rate	5e-3, 1e-3, 5e-4
Weight decay (L_2)	5e-3, 1e-3, 5e-4
Threshold ratio r	0.5, 0.8
Hidden units	32, 64
Dropout ratio	0, 0.5