

A UNIFIED VIEW ON THE REPRESENTATIONAL POWER OF GRAPH CONVOLUTION, MESSAGE-PASSING NNS AND GRAPH TRANSFORMERS

Anonymous authors

Paper under double-blind review

ABSTRACT

TODO

Guaranteed representational power + incorporate proper graph inductive bias.

Motivation: Xxx

Outline:

- A new framework
 - Unifying existings with this new framework;
- Representational power analysis
 - The upper bound of \mathcal{L} ;
 - The Representational Power of Polynomial Filters;
 - The Representational Power of Truncate Filters;
 - The complementary between polynomial filter and truncate filter
 - * The trade-off between representational power and complexity.
- Designing powerful and computational efficient \mathbf{M} ;

Misc:

Predict (row-wise transformation) then propagate (column-wise transformation) paradigm

A more well-motivated explanation of required filters based on labels over existng low/high pass filter arguements. Lemma shows that \mathcal{L} relates to graph structure, node features and labels.

1 INTRODUCTION

2 PRELIMINARIES

3 GRAPH LEARNING MODULE UNIFICATION

Given a graph G and node features $X \in \mathbb{R}^{n \times d}$, the operations in various graph learning modules can be unified into row-wise transformation f_W^{row} and column-wise transformation f_M^{col} on X ,

$$Z = f_M^{\text{col}} \circ f_W^{\text{row}}(X). \quad (1)$$

Here, regarding to specific graph learning modesl, $f_W^{\text{row}}(\mathcal{O}) = \sigma(\mathcal{O}W)$ can be a single- or multi-layer perceptron applied on each row of X by right multiplying W , where $W \in \mathbb{R}^{d \times d'}$ is the learnable weights. $f_M^{\text{col}}(\mathcal{O}) = \sigma(\mathbf{M}\mathcal{O})$ is applied on each column of X by left multiplying \mathbf{M} , where $\mathbf{M} \in \mathbb{R}^{n \times n}$ is the matrix representation of G . \mathbf{M} can also be learnable. Most importantly, \mathbf{M} incorporates *graph inductive biases* to leverage topology information in learning node representations. For example, in GCN, $\mathbf{M} = \hat{A}$. In GAT, to utilize topology information, only entries corresponding to edges are learnable. The most fundamental difference of various graph learning modules is their

ways to incorporate graph inductive biases into \mathbf{M} . In some literatures, f_W^{row} and f_M^{col} refer to *predict* and *propagate* operations respectively (Klicpera et al., 2019), and different \mathbf{M} result in different propagation scheme.¹

3.1 STRATEGIES OF INCORPORATING GRAPH INDUCTIVE BIASES INTO \mathbf{M}

\mathbf{M} is learnable. Encoding topology information into \mathbf{M} with various graph inductive biases.

Graph Convolution.

$$\mathbf{M} = g_{\alpha}^{(k)}(\hat{L}) = U \text{diag} \left(g_{\alpha}^{(k)}(\lambda) \right) U^{\top} = U \sum_{i=0}^k \alpha_i \lambda^i U^{\top}, \quad (2)$$

where $g_{\alpha}^{(k)}$ denotes a k -degree polynomial with the polynomial coefficient $\alpha \in \mathbb{R}^k$.

Message-passing NNs. MPNNs involve fixed \mathbf{M} , such as different aggregators SUM, MEAN, MAX/MIN, or learnable \mathbf{M} such as GAT.

$$\mathbf{M} = \hat{A} \quad (3)$$

Neighbor-sampling, DropEdge.

Graph Transformers. Graph transformers incorporate graph inductive biases into \mathbf{M} by proposing a series of positional encoding techniques.

$$\mathbf{M} = QK^{\top} + \text{Bias} \quad (4)$$

There is a lack of a unified understanding of the effectiveness of various graph learning modules. In this work, we fill this part by analyzing their representational power.

4 REPRESENTATIONAL POWER ANALYSIS

Following the setting in Xu et al. (2021) and Wang & Zhang (2022), we study the representational power, i.e. fitting ability, of the linear case of Equation 1 trained with the squared loss. For the standard semi-supervised node classification, it corresponds to

$$\begin{aligned} \mathcal{L} &= \|\mathbf{T}(\mathbf{Y} - \mathbf{MH})\|_F \\ &= \|\mathbf{TY} - \mathbf{TMH}\|_F \\ &= \|\mathbf{Y}^{\text{train}} - \mathbf{TMH}\|_F, \end{aligned} \quad (5)$$

where $\mathbf{Y} \in \{0, 1\}^{n \times c}$ is one-hot encoding node labels with the number of classes c , $\mathbf{H} = f_W(\mathbf{X})$ and $\mathbf{T} \in \{0, 1\}^{m \times n}$ is the mask of training set with the size m .² We use \mathbf{Y}^* to denote valid or test set labels, then $\mathbf{Y} = [\mathbf{Y}^{\text{train}} \| \mathbf{Y}^*]$.

Lemma 1. Given $\mathcal{L} = \|\mathbf{T}(\mathbf{Y} - \mathbf{MH})\|_F$, suppose \mathbf{H} has no missing frequency components over all channels and frequency profiles, i.e. $\mathbf{u}_i^{\top} \mathbf{h}_j \neq 0$, and $\|\mathbf{h}_j\|_2 \leq \eta$ for all $i \in [n], j \in [c]$, then:

$$\mathcal{L} \leq \sqrt{mn}\eta \sum_{i=1}^c \left\| \text{vec}_j \left(\frac{\mathbf{u}_j^{\top} [\mathbf{y}_i^{\text{train}} \| \mathbf{y}_i^*]}{\mathbf{u}_j^{\top} \mathbf{h}_i} \right) - \mathbf{t} \right\|_2, \quad (6)$$

where $\mathbf{M} = U \text{diag}(\mathbf{t}) U^{\top}$.

¹In the multi-layer settings, they apply multi-layer f_W^{row} and f_M^{col} in an alternate manner, i.e. $Z^{(k)} = f_M^{\text{col}(k)}(f_W^{\text{row}(k)}(Z^{(k-1)}))$. There are also GNNs applying decouple architectures where node features are processed by the multi-layer f_W^{row} and f_M^{col} individually (Klicpera et al., 2019; Wu et al., 2019; Klicpera et al., 2019; Liu et al., 2020; Zhu & Koniusz, 2020; Zhang et al., 2021). σ in f_M^{col} and f_W^{row} share the same role that introduces nonlinearity to the transformation. Some works study linear GNNs which correspond to removing σ in Equ. 1 (Wu et al., 2019; Xu et al., 2021; Wang & Zhang, 2022; Liu et al., 2021; 2022).

²Note that although \mathbf{Y} in Equation 5 involves valid or test set labels, it does not lead to label leakage as these labels are masked by \mathbf{T} and do not affect \mathcal{L} .

Due to semi-supervised settings, the derived upper bound of \mathcal{L} in Equation 6 indicates a connection with valid or test set labels \mathbf{Y}^* , which makes it inconsistent with a practical scenario. However, it provides valuable insights by showing that no matter the assignments of \mathbf{Y}^* , an expressive filter \mathbf{t} can always achieve a lower upper bound than an inexpressive filter by approximating $\text{vec}_j \left(\frac{\mathbf{u}_j^\top [\mathbf{y}_i^{\text{train}} \parallel \mathbf{y}_i^*]}{\mathbf{u}_j^\top \mathbf{h}_i} \right)$ with smaller errors.

We prove Lemma 1 in Appendix A. As $H = f_W(X)$ and $\mathbf{M} = U \text{diag}(\mathbf{t}) U^\top$, Lemma 1 indicates the upper bound of \mathcal{L} is controlled by both f_M^{col} and f_W^{row} . This is consistent with our empirical understanding that the predictions of nodes should consider both the node features and topology information.

f_W^{row} is generally implemented as single- or multi-layer perceptron whose representational power is guaranteed by universal approximation theorem (Hornik et al., 1989; Cybenko, 1989). In contrast, the design of the transformation f_M^{col} is more challenging since apart from maintaining the representational power, it also needs to effectively encode topology information into the learned node representations, which is also mentioned as graph inductive biases (Ma et al., 2023). This highlights the difficulty of graph learning. Next, we study the effects of f_M^{col} to the fitting ability.

4.1 THE REPRESENTATIONAL POWER OF POLYNOMIAL FILTERS

Filter expressiveness in graph convolution has drawn extensive research interests. A common understanding is that a more expressive filter can better filtering signal patterns over different frequency profiles thus achieves better performance (Chien et al., 2021; He et al., 2021; Wang & Zhang, 2022; Yang et al., 2022; Bo et al., 2022). Here, we study how filter expressiveness relates to the representational power of models.

As shown in Equ. 6, \mathbf{t} is shared over all c channels, for a given channel $\mathbf{y}^* \in \mathbb{R}^n$ and the corresponding $\mathbf{h} \in \mathbb{R}^n$, we have

$$\begin{aligned} \mathcal{L} &= \|T(\mathbf{y}^* - \mathbf{M}\mathbf{h})\|_2 \\ &\leq \sqrt{mn\eta} \left\| (U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*) - \mathbf{t} \right\|_2. \end{aligned} \quad (7)$$

In graph convolution, \mathbf{t} is interpreted as a filter and is generally approximated by polynomials. Although there are extensive studies on approximation abilities, there is a lack of the understanding of the approximation objective. Equ. 7 fills this gap from the perspective of the representational power. Specifically, we use \mathbf{t} to approximate $(U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*)$, and $\|(U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*) - \mathbf{t}\|_2$ is the approximation error. A smaller \mathcal{L} requires a smaller approximation error. In the context of filter studies, the desired filter relates to graph topology U , input signals \mathbf{h} and the label distribution \mathbf{y}^* . U and \mathbf{y}^* together reflect the homophily of a graph. So the desired filter also relates to the graph homophily, which is consistent with our intuition.

Let $\epsilon = \|(U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*) - \mathbf{t}\|_2$ denote the approximation error of polynomials. Then according to Equ. 7, we have

$$\mathcal{L} \leq \sqrt{mn\eta} \epsilon. \quad (8)$$

A high-degree polynomial achieves smaller approximation error, i.e. $\epsilon_{g^{(k+1)}} \leq \epsilon_{g^{(k)}}$, thus a high-degree polynomial corresponds to higher representational power. The perfect fitting requires a degree- $(n-1)$ polynomial which is also known as the universal filter approximator (He et al., 2021; Yang et al., 2022; Wang & Zhang, 2022; Bo et al., 2022). Now, we align filter expressiveness with the representational power. Although theoretically we can increase the degree of polynomials to consistently improve the expressiveness, unfortunately, due to the limited computational precision and numerical instability, the high-degree polynomial are intractable in implementations (Yang et al., 2022).

4.2 THE REPRESENTATIONAL POWER OF TRUNCATE FILTER

Due to the limited expressiveness of polynomials, some recent work choose to learn the filter with more expressive MLP or Transformer (Lingam et al., 2022; Yang et al., 2022; Bo et al., 2022). Correspondingly, following the formulation of Equ. 1, $\mathbf{M} = U \text{diag}(\text{MLP}(\boldsymbol{\lambda})) U^\top$ or

$\mathbf{M} = U \text{diag}(\text{Transformer}(\boldsymbol{\lambda})) U^\top$ respectively. The limitation is that it requires eigendecomposition computation to obtain U and $\boldsymbol{\lambda}$ first, which is intractable on large-scale graphs. Luckily, the results show that utilizing partial spectrum information of the graph matrix can still obtain competitive prediction performance. Thus they apply truncate eigendecomposition which achieves a good balance between prediction performance and computation efficiency on large graphs. We call such partial eigendecomposition-based methods truncate filter.

We use $\tilde{\boldsymbol{\lambda}}^{(k)} \in \mathbb{R}^k$ and $\tilde{U}^{(k)} \in \mathbb{R}^{n \times k}$ to denote the top k spectrum information obtained by truncate eigendecomposition, and $\tilde{U}^{(n-k)} \in \mathbb{R}^{n \times (n-k)}$ to denote the remaining $n - k$ eigenvectors.

Proposition 1. When $\mathbf{M} = \tilde{U}^{(k)} \text{diag}(\text{MLP}(\tilde{\boldsymbol{\lambda}}^{(k)})) \tilde{U}^{(k)\top}$, we have

$$\mathcal{L} \leq \sqrt{mn}\eta \left(\epsilon + \left\| \left(\tilde{U}^{(n-k)\top} \mathbf{h} \right)^{-1} \odot \left(\tilde{U}^{(n-k)\top} \mathbf{y}^* \right) \right\|_2 \right), \quad (9)$$

where $\epsilon > 0$ denotes the approximation error of MLP.

We prove Proposition 1 in Appendix B. Proposition 1 shows that a larger k achieves better fitting ability. When $k = n$, $\mathcal{L} \leq \sqrt{mn}\eta\epsilon$, where ϵ relates to the expressiveness of MLP.

4.3 POLYNOMIAL FILTER VS TRUNCATE FILTER

On the one hand, polynomial filters have drawn extensive research interests, and many spectral GNNs have been proposed based on different polynomials. On the other hand, truncate filters show its competitive performance in recent work. However, there is a lack of a detailed understanding of the difference between them. The representational power analysis provides a way to fill this gap. Specifically, we can compare the power of two filters by comparing which one achieves smaller loss bound. We use $\mathcal{L}_{\text{poly.}}$ to denote the loss of polynomial filter as shown in Equ. 8, and $\mathcal{L}_{\text{trunc.}}$ to denote the loss of polynomial filter as shown in Equ. 9. Both $\mathcal{L}_{\text{poly.}}$ and $\mathcal{L}_{\text{trunc.}}$ are related to the approximation abilities of approximators and the availability of frequency profiles (FPs) as shown in Fig. 1.

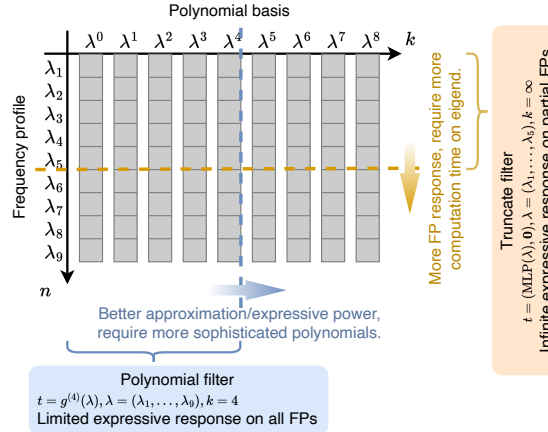


Figure 1: The comparison between polynomial filter and truncate filter.

Polynomial filters are implemented as the polynomial of the entire spectrum, thus they can operate all frequency profiles. While the limitation is that the polynomials as approximators whose approximation abilities are bounded by the degree of polynomials. Generally, the applicable degree of polynomials is far from the desired infinite approximation power, i.e. $k = n$. Therefore, although polynomial filters have all profiles response, the approximation error can be large, and this leads to a larger bound of $\mathcal{L}_{\text{poly.}}$. In contrast, truncate filters conduct truncate eigendecomposition and can only

operate partial spectrum, thus they can only operate partial frequency profiles. For the missing frequency profiles, it corresponds to $\|(\tilde{U}^{(n-k)\top} \mathbf{h})^{-1} \odot (\tilde{U}^{(n-k)\top} \mathbf{y}^*)\|_2$ as in Equ. 9. The positive side is that for the available frequency profiles, truncate filter can achieve infinite approximation power (which corresponds to $k = n$) with powerful approximators like MLP or Transformer. Therefore, truncate filters are complementary to polynomial filters that well handle the approximation error but only on the partial profiles. If the label information has no frequency components in the missing profiles, i.e. $(\tilde{U}^{(n-k)\top} \mathbf{y}^*) = \mathbf{0}$, truncate filters can achieve much smaller loss bound.

5 DESIGNING POWERFUL AND COMPUTATIONAL EFFICIENT M

Following the principle of graph convolution, the most powerful model requires the filter approximator to have infinite approximation power over all frequency profiles. This objective is much challenging as discussed in Sec. 4.3. Also, we have shown that various efforts in improving graph convolutions can be unified into improving the representational power. Naturally, this motivates us to alter the objective to improving the representational power. Let $\tilde{\mathcal{L}} = \sqrt{m}\|\mathbf{y}^* - \mathbf{M}\mathbf{h}\|_2$ denote the upper bound of \mathcal{L} .

Proposition 2. *Given a group of $k + 1$ linearly independent vectors, let $U^{(k)} = (\mathbf{u}_1, \dots, \mathbf{u}_k) \in \mathbb{R}^{n \times k}$, $U^{(k+1)} = (\mathbf{u}_1, \dots, \mathbf{u}_{k+1}) \in \mathbb{R}^{n \times (k+1)}$ and $\mathbf{M}^{(k)} = U^{(k)} \text{diag}(\boldsymbol{\alpha}) U^{(k)\top}$. Then we have*

$$\min_{\boldsymbol{\alpha}} \tilde{\mathcal{L}}_{\mathbf{M}^{(k+1)}} \leq \min_{\boldsymbol{\alpha}} \tilde{\mathcal{L}}_{\mathbf{M}^{(k)}}. \quad (10)$$

When $k = n$, we have $\min_{\boldsymbol{\alpha}} \tilde{\mathcal{L}}_{\mathbf{M}^{(k)}} = 0$.

We prove Proposition 2 in Appendix C. In graph convlution, the columns of U is considered as the bases of graph Fourier transform and are orthogonal to each other. However, from the perspective of improving the representational power, the orthogonality is unnecessary, and one can apply a larger k to consistently improve the representational power as long as all \mathbf{u}_i are linearly independent. Benefit from relaxing the orthogonality requirement, we can make the highest representational power implementable.

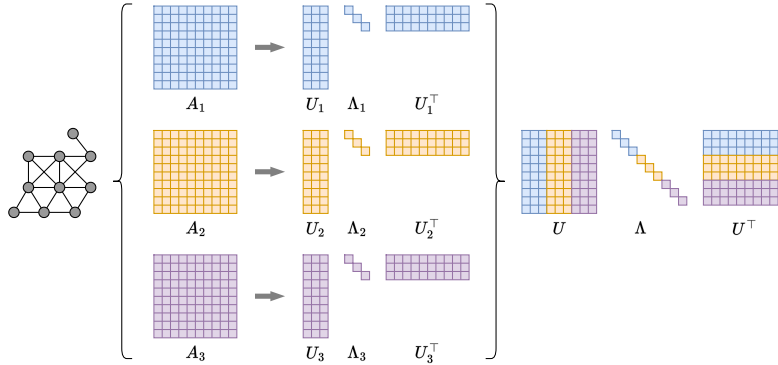


Figure 2: Xxx.

6 EXPERIMENTS

7 CONCLUSION

REFERENCES

Deyu Bo, Chuan Shi, Lele Wang, and Renjie Liao. Specformer: Spectral graph neural networks meet transformers. In *The Eleventh International Conference on Learning Representations*, 2022.

- Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized pagerank graph neural network. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=n6jl7fLxrP>.
- George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of control, signals and systems*, 2(4):303–314, 1989.
- Mingguo He, Zhewei Wei, Zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary graph spectral filters via bernstein approximation. In *NeurIPS*, 2021.
- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations (ICLR)*, 2019.
- Vijay Lingam, Manan Sharma, Chanakya Ekbote, Rahul Ragesh, Arun Iyer, and Sundararajan Selamanickam. A piece-wise polynomial filtering approach for graph neural networks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 412–452. Springer, 2022.
- Juncheng Liu, Kenji Kawaguchi, Bryan Hooi, Yiwei Wang, and Xiaokui Xiao. Eignn: Efficient infinite-depth graph neural networks. *Advances in Neural Information Processing Systems*, 34: 18762–18773, 2021.
- Juncheng Liu, Bryan Hooi, Kenji Kawaguchi, and Xiaokui Xiao. Mgnni: Multiscale graph neural networks with implicit layers. In *Advances in Neural Information Processing Systems*, 2022.
- Meng Liu, Hongyang Gao, and Shuiwang Ji. Towards deeper graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 338–348, 2020.
- Liheng Ma, Chen Lin, Derek Lim, Adriana Romero-Soriano, Puneet K Dokania, Mark Coates, Philip Torr, and Ser-Nam Lim. Graph inductive biases in transformers without message passing. *arXiv preprint arXiv:2305.17589*, 2023.
- Xiyuan Wang and Muhan Zhang. How powerful are spectral graph neural networks. *ICML*, 2022.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6861–6871. PMLR, 2019.
- Keyulu Xu, Mozhi Zhang, Stefanie Jegelka, and Kenji Kawaguchi. Optimization of graph neural networks: Implicit acceleration by skip connections and more depth. In *International Conference on Machine Learning*, pp. 11592–11602. PMLR, 2021.
- Mingqi Yang, Yanming Shen, Rui Li, Heng Qi, Qiang Zhang, and Baocai Yin. A new perspective on the effects of spectrum in graph neural networks. In *Proceedings of the 39th International Conference on Machine Learning*, 2022.
- Shanzhuo Zhang, Lihang Liu, Sheng Gao, Donglong He, Xiaomin Fang, Weibin Li, Zhengjie Huang, Weiyue Su, and Wenjin Wang. Litegem: Lite geometry enhanced molecular representation learning for quantum property prediction, 2021.
- Hao Zhu and Piotr Koniusz. Simple spectral graph convolution. In *International Conference on Learning Representations*, 2020.

A PROOF OF LEMMA 1

Proof.

$$\begin{aligned}
\mathcal{L} &= \|\mathbf{T}(\mathbf{Y} - \mathbf{MH})\|_F \\
&\leq \|\mathbf{T}\|_F \|\mathbf{Y} - \mathbf{MH}\|_F \\
&= \sqrt{m} \|\mathbf{U}\mathbf{U}^\top \mathbf{Y} - \mathbf{U} \text{diag}(\mathbf{t}) \mathbf{U}^\top \mathbf{H}\|_F \\
&\leq \sqrt{m} \|\mathbf{U}\|_F \|\mathbf{U}^\top \mathbf{Y} - \text{diag}(\mathbf{t}) \mathbf{U}^\top \mathbf{H}\|_F \\
&= \sqrt{m} \sqrt{n} \sqrt{\sum_{i=1}^c \|\mathbf{U}^\top \mathbf{y}_i - \text{diag}(\mathbf{t}) \mathbf{U}^\top \mathbf{h}_i\|_2^2} \\
&\leq \sqrt{m} \sqrt{n} \sum_{i=1}^c \|\mathbf{U}^\top \mathbf{y}_i - \text{diag}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{t}\|_2 \\
&= \sqrt{m} \sqrt{n} \sum_{i=1}^c \|\text{diag}(\mathbf{U}^\top \mathbf{h}_i) (\text{diag}^{-1}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{U}^\top \mathbf{y}_i - \mathbf{t})\|_2 \\
&\leq \sqrt{m} \sqrt{n} \sum_{i=1}^c \|\mathbf{U}^\top \mathbf{h}_i\|_2 \|\text{diag}^{-1}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{U}^\top \mathbf{y}_i - \mathbf{t}\|_2 \\
&\leq \sqrt{m} \sqrt{n} \sum_{i=1}^c \|\mathbf{U}^\top\|_F \|\mathbf{h}_i\|_2 \|\text{diag}^{-1}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{U}^\top \mathbf{y}_i - \mathbf{t}\|_2 \\
&= \sqrt{mn} \sum_{i=1}^c \|\mathbf{h}_i\|_2 \|\text{diag}^{-1}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{U}^\top \mathbf{y}_i - \mathbf{t}\|_2 \\
&\leq \sqrt{mn} \eta \sum_{i=1}^c \|\text{diag}^{-1}(\mathbf{U}^\top \mathbf{h}_i) \mathbf{U}^\top \mathbf{y}_i - \mathbf{t}\|_2 \\
&= \sqrt{mn} \eta \sum_{i=1}^c \left\| \text{vec}_j \left(\frac{\mathbf{u}_j^\top \mathbf{y}_i}{\mathbf{u}_j^\top \mathbf{h}_i} \right) - \mathbf{t} \right\|_2 \\
&= \sqrt{mn} \eta \sum_{i=1}^c \left\| \text{vec}_j \left(\frac{\mathbf{u}_j^\top [\mathbf{y}_i^{\text{train}} \|\mathbf{y}_i^*\]}{\mathbf{u}_j^\top \mathbf{h}_i} \right) - \mathbf{t} \right\|_2
\end{aligned} \tag{11}$$

□

B PROOF OF PROPOSITION 1

Assume the k -truncate eigendecomposition is applied. Correspondingly,

$$\begin{aligned}
\mathbf{M} &= \tilde{\mathbf{U}}^{(k)} \text{diag} \left(\text{MLP} \left(\tilde{\boldsymbol{\lambda}}^{(k)} \right) \right) \tilde{\mathbf{U}}^{(k)\top} + \tilde{\mathbf{U}}^{(n-k)} \text{diag}(\mathbf{0}) \tilde{\mathbf{U}}^{(n-k)\top} \\
&= \mathbf{U} \text{diag} \left(\left(\text{MLP} \left(\tilde{\boldsymbol{\lambda}}^{(k)} \right), \mathbf{0} \right) \right) \mathbf{U}^\top \\
&= \mathbf{U} \text{diag}(\mathbf{t}) \mathbf{U}^\top,
\end{aligned} \tag{12}$$

where $\mathbf{t} = (\text{MLP}(\tilde{\lambda}), \mathbf{0}) \in \mathbb{R}^n$. Let $\epsilon > 0$ denotes the approximation error of MLP. Then

$$\begin{aligned}
\mathcal{L} &\leq \sqrt{mn}\eta \left\| (U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*) - \mathbf{t} \right\|_2 \\
&= \sqrt{mn}\eta \left\| (U^\top \mathbf{h})^{-1} \odot (U^\top \mathbf{y}^*) - (\text{MLP}(\tilde{\lambda}), \mathbf{0}) \right\|_2 \\
&= \sqrt{mn}\eta \sqrt{\left\| (\tilde{U}^\top \mathbf{h})^{-1} \odot (\tilde{U}^\top \mathbf{y}^*) - \text{MLP}(\tilde{\lambda}) \right\|_2^2 + \left\| (\tilde{U}'^\top \mathbf{h})^{-1} \odot (\tilde{U}'^\top \mathbf{y}^*) - \mathbf{0} \right\|_2^2} \\
&\leq \sqrt{mn}\eta \sqrt{\epsilon^2 + \left\| (\tilde{U}'^\top \mathbf{h})^{-1} \odot (\tilde{U}'^\top \mathbf{y}^*) \right\|_2^2} \\
&\leq \sqrt{mn}\eta \left(\epsilon + \left\| (\tilde{U}'^\top \mathbf{h})^{-1} \odot (\tilde{U}'^\top \mathbf{y}^*) \right\|_2 \right).
\end{aligned} \tag{13}$$

C PROOF OF PROPOSITION 2

$$\begin{aligned}
\min_{\alpha} \bar{\mathcal{L}}_{\mathbf{M}^{(k+1)}} &= \min_{\alpha} \sqrt{m} \left\| \mathbf{y}^* - \mathbf{M}^{(k+1)} \mathbf{h} \right\|_2 \\
&= \min_{\alpha_0, \dots, \alpha_{k+1}} \sqrt{m} \left\| \mathbf{y}^* - \sum_{i=1}^{k+1} \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&\leq \min_{\alpha_0, \dots, \alpha_k | \alpha_{k+1}=0} \sqrt{m} \left\| \mathbf{y}^* - \sum_{i=1}^{k+1} \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&= \min_{\alpha_0, \dots, \alpha_k} \sqrt{m} \left\| \mathbf{y}^* - \sum_{i=1}^k \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&= \min_{\alpha} \sqrt{m} \left\| \mathbf{y}^* - \mathbf{M}^{(k)} \mathbf{h} \right\|_2 \\
&= \min_{\alpha} \bar{\mathcal{L}}_{\mathbf{M}^{(k)}}
\end{aligned} \tag{14}$$

When $k = n$, $\mathbf{M}^{(n)} \mathbf{h} = U \text{diag}(\alpha) U^\top \mathbf{h} = U(\alpha \odot U^\top \mathbf{h})$. Let $\mathbf{x} = \alpha \odot U^\top \mathbf{h}$. Consider the linear system $U\mathbf{x} = \mathbf{y}^*$. As $\text{rank}(U) = n$, a solution \mathbf{x}_0 always exists. Suppose $U_i^\top \mathbf{h} \neq 0$ for all $i \in [n]$. Correspondingly, $\alpha_0 = \mathbf{x}_0 \odot (U^\top \mathbf{h})^{-1}$. Therefore, $\min_{\alpha} \bar{\mathcal{L}}_{\mathbf{M}^{(n)}} = \bar{\mathcal{L}}_{\mathbf{M}^{(n)}}|_{\alpha_0} = 0$.

D PROOF OF PROPOSITION 2

For any $\alpha \in \mathbb{R}^{k+1}$, we have

$$\begin{aligned}
\bar{\mathcal{L}}_{\mathbf{M}^{(k+1)}} &= \sup \left\| \mathbf{y}^* - \mathbf{M}^{(k+1)} \mathbf{h} \right\|_2 \\
&= \left\| \mathbf{y}^* - \sum_{i=1}^{k+1} \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&\leq \left\| \mathbf{y}^* - \sum_{i=1}^{k+1} \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&= \left\| \mathbf{y}^* - \sum_{i=1}^k \alpha_i U_i U_i^\top \mathbf{h} \right\|_2 \\
&= \left\| \mathbf{y}^* - \mathbf{M}^{(k)} \mathbf{h} \right\|_2 \\
&= \bar{\mathcal{L}}_{\mathbf{M}^{(k)}}
\end{aligned} \tag{15}$$

E MISC

H is channel-independent. \mathbf{t} is shared over different channels. There are also channel-independent filter design where each signal channel learns an individual filter (Yang et al., 2022; Wang & Zhang, 2022; Bo et al., 2022). Correspondingly, the loss is $\mathcal{L} \leq \sqrt{mn}\eta \sum_{i=1}^c \left\| (U^\top H_i)^{-1} \odot (U^\top Y_i^*) - \mathbf{T}_i \right\|_F$ with $\mathbf{T} \in \mathbb{R}^{n \times c}$.

E.1 EXPLORING GRAPH REPRESENTATION SPACE

According to spectral theorem $\mathbf{M} = U \text{diag}(\mathbf{t}) U^\top$, a graph matrix representation \mathbf{M} provides a unique tuple (U, \mathbf{t}) . (U, \mathbf{t}) decides the fitting ability. Improving fitting ability while maintaining graph topology information.

$$\begin{aligned}
\mathcal{L} &= \|T(Y^* - \mathbf{M}H)\|_F \\
&\leq \|T\|_F \|Y^* - \mathbf{M}H\|_F \\
&= \|T\|_F \left\| Y^* - \sum_{i=1}^n \lambda_i U_i U_i^\top H \right\|_F \\
&= \|T\|_F \left\| \sum_{i=1}^n U_i U_i^\top Y^* - \sum_{i=1}^n \lambda_i U_i U_i^\top H \right\|_F \\
&= \|T\|_F \left\| \sum_{i=1}^n U_i U_i^\top (Y^* - \lambda_i H) \right\|_F \\
&\leq \|T\|_F \sum_{i=1}^n \|U_i U_i^\top (Y^* - \lambda_i H)\|_F \\
&= \|T\|_F \sum_{i=1}^n \sqrt{\sum_{j=1}^c \|U_i U_i^\top (Y_j^* - \lambda_i H_j)\|_F^2} \\
&= \|T\|_F \sum_{i=1}^n \sqrt{\sum_{j=1}^c (U_i^\top (Y_j^* - \lambda_i H_j))^2 \|U_i\|_F^2} \\
&= \|T\|_F \sum_{i=1}^n \sqrt{\sum_{j=1}^c (U_i^\top (Y_j^* - \lambda_i H_j))^2}
\end{aligned} \tag{16}$$

$$\begin{aligned}
\mathcal{L} &= \|T(Y^* - \mathbf{M}H)\|_F \\
&\leq \|T\|_F \|Y^* - \mathbf{M}H\|_F \\
&= \sqrt{m} \|UU^\top Y^* - U \text{diag}(\mathbf{t}) U^\top H\|_F \\
&\leq \sqrt{m} \|U\|_F \|U^\top Y^* - \text{diag}(\mathbf{t}) U^\top H\|_F \\
&= \sqrt{m} \sqrt{n} \sqrt{\sum_{j=1}^c \|U^\top Y_j^* - \text{diag}(\mathbf{t}) U^\top H_j\|_F^2} \\
&= \sqrt{m} \sqrt{n} \sqrt{\sum_{j=1}^c \|U^\top Y_j^* - \text{diag}(U^\top H_j) \boldsymbol{\lambda}\|_F^2} \\
&= \sqrt{m} \sqrt{n} \sqrt{\sum_{j=1}^c \|\text{diag}(U^\top H_j) (\text{diag}^{-1}(U^\top H_j) U^\top Y_j^* - \boldsymbol{\lambda})\|_F^2} \tag{17} \\
&\leq \sqrt{m} \sqrt{n} \sqrt{\sum_{j=1}^c \|U^\top H_j\|_F^2 \|(U^\top H_j)^{-1} \odot (U^\top Y_j^*) - \boldsymbol{\lambda}\|_F^2} \\
&\leq \sqrt{m} \sqrt{n} \sqrt{\sum_{j=1}^c \|U^\top\|_F^2 \|H_j\|_F^2 \|(U^\top H_j)^{-1} \odot (U^\top Y_j^*) - \boldsymbol{\lambda}\|_F^2} \\
&= \sqrt{mn} \sqrt{\sum_{j=1}^c \|H_j\|_F^2 \|(U^\top H_j)^{-1} \odot (U^\top Y_j^*) - \boldsymbol{\lambda}\|_F^2} \\
&\leq \sqrt{mn} \sum_{j=1}^c \|H_j\|_F \|(U^\top H_j)^{-1} \odot (U^\top Y_j^*) - \boldsymbol{\lambda}\|_F
\end{aligned}$$