# Towards the Universal Learning Principle for Graph Neural Networks

Anonymous Author(s) Affiliation Address email

# Abstract

Graph neural networks (GNNs) are currently highly regarded in graph repre-1 sentation learning tasks due to their significant performance. Although various 2 3 propagation mechanisms and graph filters were proposed, few works have investigated their rationale from the perspective of learning. In this paper, we elucidate 4 the criterion for the graph filter formed by power series, and further establish a 5 scalable regularized learning framework that theoretically realizes very deep GNN. 6 Following the framework, we introduce Adaptive Power GNN (APGNN), a deep 7 GNN that employs exponentially decaying weights to aggregate graph information 8 9 of varying orders, thus facilitating more effective mining of deeper neighbor infor-10 mation. Moreover, the multiple P-hop message passing strategy is proposed to efficiently perceive the higher-order neighborhoods. Different from other GNNs, 11 the proposed APGNN can be seamlessly extended to an infinite-depth network. To 12 clarify the learning guarantee, we theoretically analyze the generalization of the 13 proposed learning framework via uniform convergence. Experimental results show 14 that APGNN obtains superior performance compared to state-of-the-art GNNs, 15 highlighting the effectiveness of our framework. 16

# 17 **1 Introduction**

Recently, Graph Neural Networks (GNNs) have shown commendable performance on numerous graph 18 19 representation learning tasks. In addition, GNNs have been introduced in a variety of application tasks, 20 such as recommendation systems [7, 11, 37], computer vision [4, 13, 23], and traffic forecasting [8, 9]. The fundamental part of GNN is the design of the propagation mechanism or the graph 21 filter [5, 12, 27, 32, 34]. GNNs can be categorized into two groups based on the approach of 22 formulation. Spatial-based GNN formulates propagation mechanisms through the direct aggregation 23 of spatial features. As one of the most simple GNNs, Graph Convolutional Network (GCN) [15] 24 designs graph convolutional layer via aggregating one-hop information on the graph. Graph Attention 25 Network (GAT) [30] learns node relationships using an attention mechanism, enhancing the scalability 26 of the network. For extension of inductive learning, GraphSAGE [10] employs various pooling 27 operations as aggregation functions. Liu et. al proposed DAGNN, which integrates information from 28 multiple receptive fields for adaptive propagation [19]. Spectral-based GNN designs graph filters by 29 constructing filter functions in the graph Fourier domain, which aims to find a proper transformation 30 of the graph spectrum. Chebynet constructs the localized graph filter with Chebyshev polynomial [3]. 31 From the view of the spectrum, GCN could be seen as a Chebyshev filter with first-order truncation 32 [15]. To construct deeper GNN, Personalize PageRank method is employed to design graph filter 33 [16]. GNN-LF/HF [36] concludes various designs of graph filters and constructs the graph filter 34 through a graph optimization framework. 35

<sup>36</sup> Despite their success, few studies have explored the general rule for devising GNNs from the <sup>37</sup> perspective of learning. In this paper, we start from the graph filter formed by power series and

Submitted to 37th Conference on Neural Information Processing Systems (NeurIPS 2023). Do not distribute.



Figure 1: An illustration of the proposed APGNN that adheres to the learning principle. The model incorporates the decay rate  $\alpha$  to suppress the information from high-order neighbors while adaptively learning bounded coefficients  $\beta$ . Furthermore, it aggregates information with P-hop to enlarge the receptive field. This design enables the seamless extension of APGNN to an extremely deep network.

discuss what makes a legitimate graph filter for the construction of deep GNN. A learning principle 38 is then proposed to summarize the rule of formulating a graph filter. Following this, we propose 39 Adaptive Power Graph Neural Network (APGNN), which adaptively learns the task-specific graph 40 filter for node representation learning. The main idea of APGNN is depicted in Figure 1. The 41 parameterized graph filter is designed with regularization of the exponential decay rate. A multiple 42 *P*-hop strategy is applied to enhance the capacity of perceiving the higher-order neighborhoods. 43 Furthermore, the generalization bound of APGNN is presented with the setting of the continuous 44 graph, which provides a learning guarantee for the proposed principle theoretically. 45

We conduct evaluations on five benchmark datasets on node classification tasks. The experimental 46 47 results suggest the superiority of the proposed method over the existing GNNs. The theoretical analysis is also validated via the empirical study. 48

#### 2 **Preliminaries** 49

**Notations.** Suppose we have an undirected graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$  with node set  $\mathcal{V}$  and  $|\mathcal{V}| = n$ . 50  $\mathbf{A} \in \mathbb{R}^{n \times n}$  denotes the adjacency matrix indicating the edges in  $\mathcal{E}$ . Assuming that the self-loops are contained in the graph, i.e.,  $a_{ii} = 1$ . Let  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \cdots, \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}$  be the graph signals (or 51 52 features) of the nodes. We use notation  $[n] \triangleq \{1, 2, \dots, n\}$  for  $n \in \mathbb{N}_+$ . Assume that the label of  $\mathbf{x}_i$ 53 is  $y_i \in \mathcal{Y}$  for all  $i = [n_l]$ , where  $n_l \leq n$  is the number of labeled samples. 54

**Graph Neural Networks.** We introduce some essential concepts in GNNs. Let  $d_i = \sum_{j=1}^n A_{ij}$ 55 be the degree of *i*-th node, so the degree matrix of **A** can be defined as  $\mathbf{D} = \text{diag}(d_1, d_2, \cdots, d_n)$ . 56

- The symmetrically normalized Laplacian is  $\mathbf{L} = \mathbf{I} \tilde{\mathbf{A}}$ , where  $\tilde{\mathbf{A}} \triangleq \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$  is normalized 57
- adjacency matrix. Consider the eigen-decomposition  $\mathbf{L} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^{\top}$ , where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n)$ 58
- is the diagonal matrix of eigenvalues, and  $\mathbf{U} = [\mathbf{u}_1, \cdots, \mathbf{u}_n]$  represents the eigenvectors associated 59
- with the eigenvalues. Note that A shares the same eigenvectors with L. 60
- Spectral convolution on graphs is defined as the following transformation [15, 28]: 61

$$g * \mathbf{X} = \mathbf{U}g(\mathbf{\Lambda})\mathbf{U}^{\top}\mathbf{X},\tag{1}$$

where  $g(\cdot) : [0,2] \to \mathbb{R}$  is called filter function and  $g(\Lambda) = \text{diag}(g(\lambda_1), \cdots, g(\lambda_n))$ . The com-62 mon approach in GNNs is to apply polynomial functions as the filters [3, 12, 15], which leads to 63  $\mathbf{U}q(\mathbf{\Lambda})\mathbf{U}^{\top} = q(\mathbf{L})$ . Therefore, spectral convolution is usually written as  $q * \mathbf{X} = q(\mathbf{L})\mathbf{X}$ . The graph 64 representation paradigm in GNN is generally expressed as follows: 65

$$\mathbf{Z} = g(\mathbf{L})f(\mathbf{X}), \quad g(\mathbf{L}) = \sum_{k=0}^{K} \theta_k \tilde{\mathbf{A}}^k, \tag{2}$$

where  $\mathbf{Z} \in \mathbb{R}^{n \times c}$  denotes the node representation, and  $f(\cdot)$  represents a feature extractor such as 66 multi-layer perceptions (MLPs). 67

# 68 **3** Learning Principle for GNNs

### 69 3.1 The principle of devising graph filters

<sup>70</sup> Current studies suggest a significant relationship between the performance of GNN and its graph filter <sup>71</sup> [16, 19]. Predominantly, the general graph filters are characterized by polynomials associated with <sup>72</sup> the adjacency matrix  $\tilde{\mathbf{A}}$  (or Laplacian matrix L), i.e.,  $g(\mathbf{L}) = \sum_{k=0}^{K} \theta_k \tilde{\mathbf{A}}^k$ . However, the existing <sup>73</sup> methods still meet the issue that the depth of GNN is limited. The reason for this phenomenon <sup>74</sup> is that these GNNs are inconsistent with their "infinite-depth" version. That is, the corresponding <sup>75</sup> models lose some essential properties as the depth  $K \to \infty$ . Consequently, the depth of the models <sup>76</sup> is restricted. To address this issue, it is necessary to study the properties of GNNs with infinite depth. <sup>77</sup> Therefore, we explore the graph filter reformulated as power series:

$$g(\tilde{\mathbf{A}}) = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k = \sum_{k=0}^{\infty} \theta_k (\mathbf{I} - \mathbf{L})^k.$$
 (3)

<sup>78</sup> First of all, a well-defined graph filter represented as equation 3 must be convergent. Consequently, it

<sup>79</sup> becomes essential to investigate the properties that the coefficients  $\theta_k$  should exhibit. The following <sup>80</sup> lemma provides appropriate constraints for the coefficients of the graph filter.

**Lemma 1.** Let  $\{a_k\}$  and  $\{\gamma^k\}$  be the real number sequences, where  $\gamma \in (-1, 1]$  and  $k \in \mathbb{N}$ . Then  $\sum_{k=0}^{\infty} a_k \gamma^k$  converges uniformly and absolutely if and only if the series  $\sum_{k=0}^{\infty} a_k$  converges absolutely.

As a direct corollary, the weights of the graph filter (i.e.,  $\theta_k$ ) should satisfy the following theorem.

**Theorem 1.** Let  $\tilde{\mathbf{A}} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$  be the normalized adjacency matrix of a graph  $\mathcal{G}$  with spectral radius  $\rho(\tilde{\mathbf{A}}) \leq 1$ . The matrix series  $\sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k$  converges uniformly and absolutely if and only if the series  $\sum_{k=0}^{\infty} \theta_k$  converges absolutely.

The proofs are shown in Appendix. Theorem 1 offers a sufficient and necessary condition for the convergence of graph filters formed by power series. Specifically, the condition requires the existence of a finite real number  $M \ge 0$ ,

$$\|\boldsymbol{\theta}\|_1 \triangleq \sum_{k=0}^{\infty} |\theta_k| \le M.$$
<sup>(4)</sup>

<sup>90</sup> Therefore, an arbitrary graph filter formed by power series should satisfy the above convergence <sup>91</sup> condition, which gives the first requirement while designing GNN. Apart from convergence, we <sup>92</sup> expect the graph filter to possess good analytic properties such as smoothness. To this end, Lipschitz <sup>93</sup> continuity should be considered the second requirement of the graph filter. Let  $g(\cdot)$  be an *L*-Lipschitz <sup>94</sup> continuous function, meaning that

$$|g(\lambda) - g(\lambda')| \le L|\lambda - \lambda'|, \quad \forall \lambda, \lambda' \in [0, 2).$$
(5)

<sup>95</sup> This property indicates the stability or robustness of the model [6, 24]. If the graph is contaminated <sup>96</sup> and its eigenvalues are perturbed by at most  $\epsilon$ , Lipschitz continuity ensures the perturbation of the <sup>97</sup> graph-filtered result is at most  $L\epsilon$ . For instance, considering  $g(\lambda) = \sum_{k=0}^{\infty} (1-\lambda)^k / k^2$ , which is <sup>98</sup> convergent, yet the Lipschitz condition does not satisfy for  $\lambda$  closed to zero. Therefore, this graph

<sup>99</sup> filter might be sensitive to the input graph. Subsequently, we conclude the following criterion.

$$\mathbf{Z} = g_{\boldsymbol{\theta}}(\mathbf{L}) f(\mathbf{X}), \text{ with } \|\boldsymbol{\theta}\|_1 \le M, \ g_{\boldsymbol{\theta}}(\cdot) \text{ is a Lipschitz function.}$$
(6)

To enhance the scalability of the model, we define  $\theta$  as a learnable parameter (though its dimension is infinite). In this way, (6) gives a regularized learning framework for GNN. Therefore, for a *K*-order polynomial graph filter  $g_{\theta}^{K}(\lambda) = \sum_{k=0}^{K} \theta_{k}(1-\lambda)^{k}$ , which is what we can implement in practice, the condition (6) should be satisfied to keep the consistency with its infinitely deep version  $g_{\theta}^{\infty}(\lambda) = \sum_{k=0}^{\infty} \theta_{k}(1-\lambda)^{k}$ . We will present the applications of this criterion in this section, and further analyze the learning guarantee with generalization in section 4.

#### 106 3.2 Related works

In this subsection, we investigate the relationship between our learning framework and several wellknown Graph Neural Networks (GNNs), focusing on the design of graph filters. Our findings indicate that these GNNs are all special cases of our learning framework, which are summarized in Table 1. **GCN/SGC [15, 33].** Graph convolutional network (GCN) aims to learn a node representation by stacking multiple graph convolutional layers. In each layer, GCN applies first-order Chebyshev approximation as the graph filter followed by a fully connected layer. For simplicity, we analyze one-layer GCN, which is formulated as  $\mathbf{Z} = \sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W})$ , where  $\mathbf{W}$  is a learnable weight matrix for linear transformation. Therefore, the graph filter of one-layer GCN is  $g_{\text{GCN}}(L) = \mathbf{I} - \mathbf{L} = \tilde{\mathbf{A}}$ , or a trivial matrix power series:

$$g_{\rm GCN}(\mathbf{L}) = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k, \quad \text{where } \theta_k = \begin{cases} 1, & \text{if } k = 1, \\ 0, & \text{otherwise.} \end{cases}$$
(7)

116 It should be noted that this equation satisfies the condition described in (6).

SGC is a simplified version of GCN that eliminates the activation function and applies a single linear projection to extract features. This simplification reduces the multiple-layer GCN into a more concise model as  $\mathbf{Z} = \tilde{\mathbf{A}}^{K} \mathbf{X} \mathbf{W}$ . Similarly, the graph filter of SGC can be represented as:

$$g_{\rm SGC}(\mathbf{L}) = \tilde{\mathbf{A}}^K = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k, \quad \text{where } \theta_k = \begin{cases} 1, & \text{if } k = K, \\ 0, & \text{otherwise.} \end{cases}$$
(8)

Both GCN and SGC use a monomial to construct the graph filter. Therefore, in the viewpoint of spectral-GNN, their graph filters are too simple to capture the spectral characteristic. Besides, the small eigenvalue vanishes when K becomes very large, leaving only the largest eigenvalue, which leads to the well-known over-smoothing problem [17].

124 **PPNP [16].** PPNP uses Personalized PageRank as the graph filter, which balances local information 125 preservation and the utilization of high-order neighbor information. The model of PPNP is  $\mathbf{Z} =$ 126  $\alpha (\mathbf{I} - (1 - \alpha)\tilde{\mathbf{A}})^{-1}\mathbf{H} = (\mathbf{I} + \beta \mathbf{L})^{-1}\mathbf{H}$ , where  $\mathbf{H} = f(\mathbf{X})$  is a two-layer MLPs and  $\beta = 1/\alpha - 1$ . 127 Hence, the graph filter of PPNP is  $g_{\text{PPNP}}(\mathbf{L}) = (\mathbf{I} + \beta \mathbf{L})^{-1}$ . Considering its Taylor series, we have

$$g_{\text{PPNP}}(\mathbf{L}) = (\mathbf{I} + \beta \mathbf{L})^{-1} = \frac{1}{1+\beta} \sum_{k=0}^{\infty} \left(\frac{\beta}{1+\beta}\right)^k \tilde{\mathbf{A}}^k = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}},\tag{9}$$

where  $\theta_k = \beta^k / (1 + \beta)^{k+1}$ . It is straightforward to validate that  $\sum_{k=0}^{\infty} \theta_k = 1$ , and thus the convergence requirement (4) holds. Moreover, the Lipschitz condition is easily verified. Thus PPNP satisfies the criterion of (6). However, the performance of PPNP is heavily dependent on the hyperparameter  $\beta$ , which must be carefully tuned to achieve optimal performance.

**DAGNN [19].** DAGNN adaptively adjusts the weight of information aggregation from different neighbors to solve the over-smoothing problem. It designs a parameterized graph filter formulated as a *K*-order polynomial:

$$g_{\text{DAGNN}}(\mathbf{L}) = \sum_{k=0}^{K} \theta_k \tilde{\mathbf{A}}^k, \quad \text{s.t. } 0 \le \theta_k \le 1,$$
(10)

where  $\theta_k$  is the learnable parameter with bounded constraint. Due to this adaptive learning strategy, DAGGN is able to learn a graph filter more suitable for node classification. The empirical studies suggest DAGNN works well with a proper K. However, as  $K \to \infty$ , the constraint  $0 \le \theta_k \le 1$ cannot guarantee the convergence of the graph filter. It indicates that DAGNN is "inconsistent" with its infinitely deep version. Therefore, it can not be naturally extended to significantly deep GNN.

#### 140 3.3 Instantiation: Adaptive Power Graph Neural Network

We now introduce a novel GNN following the framework in section 3.1, called Adaptive Power GNN (APGNN). We first consider the following graph filter parameterized by  $\beta$  with the form:

$$g^{\infty}_{\beta}(\lambda) = \sum_{k=0}^{\infty} \beta_k \alpha^k (1-\lambda)^k, \quad \text{where } |\beta_k| \le 1, \ 0 < \alpha < 1, \tag{11}$$

where the coefficient of the power series  $\theta_k = \beta_k \alpha^k$ , with hyper-parameter  $\alpha \in (0, 1)$  ensuring the convergence. Immediately, we check the condition of Lemma 1.

$$\|\boldsymbol{\theta}\|_1 = \sum_{k=0}^{\infty} \left|\beta_k \alpha^k\right| \le \sum_{k=0}^{\infty} \alpha^k \le \frac{1}{1-\alpha}.$$
(12)

Table 1: Graph filter for various GNNs

Model	Filter function	Setting of $\theta$	Learnable $g(\cdot)$
1-layer GCN	$g(\mathbf{L}) = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k$	$\theta_k = \begin{cases} 1, & \text{if } k = 1 \\ 0, & \text{otherwise} \end{cases}$	No
SGC	$g(\mathbf{L}) = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k$	$\theta_k = \begin{cases} 1, & \text{if } k = K \\ 0, & \text{otherwise} \end{cases}$	No
PPNP	$g(\mathbf{L}) = \sum_{k=0}^{\infty} \theta_k \tilde{\mathbf{A}}^k$	$\theta_k = \frac{\beta^k}{(1+\beta)^{k+1}}, \ \beta > 0$	No
DAGNN	$g(\mathbf{L}) = \sum_{k=0}^{K}  heta_k  ilde{\mathbf{A}}^k$	$0 \leq \theta_k \leq 1$	Yes

Hence, the power series converges on [0, 2] absolutely and uniformly. Similarly, the associated matrix series  $g^{\infty}_{\beta}(\mathbf{L}) = \sum_{k=0}^{\infty} \beta_k \alpha^k \tilde{\mathbf{A}}^k$  also converges uniformly and absolutely by Theorem 1. Moreover,  $g^{\infty}_{\beta}(\cdot)$  is  $\alpha(1-\alpha)^{-2}$ -Lipschitz. To see this, for any  $|\beta_k| \leq 1$  and  $1-\lambda \in (-1, 1]$ , we have

$$\left|\nabla g^{\infty}_{\beta}(\lambda)\right| = \left|\sum_{k=1}^{\infty} (-1)^k k \beta_k \alpha^k (1-\lambda)^{k-1}\right| \le \sum_{k=1}^{\infty} k \alpha^k = \frac{\alpha}{(1-\alpha)^2},\tag{13}$$

which implies the Lipschitz continuous property. Thus, this graph filter fits the requirement of the proposed criterion. However, the model with this graph filter is unavailable in practice as the number of parameters to be learned is infinite. The *K*-order truncated polynomial is utilized for substitution, i.e.,  $g_{\beta}^{K}(\mathbf{L}) = \sum_{k=0}^{K} \beta_{k} \alpha^{k} \tilde{\mathbf{A}}^{k}$ . We evaluate the approximation via the upper bound of *K*-order truncation error:

$$|g^{\infty}_{\beta}(\lambda) - g^{K}_{\beta}(\lambda)| \le \sum_{k=K+1}^{\infty} \left|\beta_{k}\alpha^{k}(1-\lambda)^{k}\right| \le \sum_{k=K+1}^{\infty} \alpha^{k} = \frac{\alpha^{K+1}}{1-\alpha},\tag{14}$$

which uniformly holds for  $\forall \lambda \in [0, 2]$ . Likewise, the approximation error of matrix series is given by

$$\left\|g_{\boldsymbol{\beta}}^{\infty}(\mathbf{L}) - g_{\boldsymbol{\beta}}^{K}(\mathbf{L})\right\|_{2} = \left\|\mathbf{U}\left(g_{\boldsymbol{\beta}}^{\infty}(\boldsymbol{\Lambda}) - g_{\boldsymbol{\beta}}^{K}(\boldsymbol{\Lambda})\right)\mathbf{U}^{\top}\right\|_{2} = \sup_{i \in [n]} \left|g_{\boldsymbol{\beta}}^{\infty}(\lambda_{i}) - g_{\boldsymbol{\beta}}^{K}(\lambda_{i})\right| \le \frac{\alpha^{K+1}}{1 - \alpha},$$
(15)

where  $\lambda_i$  denotes the *i*-th eigenvalue of **L**. This upper bound is independent of the given graph, which can be controlled via tuning  $\alpha$  and K. The higher K and smaller  $\alpha$  yield a better approximation to the exact graph filter  $g^{\infty}_{\beta}(\cdot)$ . Nevertheless, the small  $\alpha$  tends to limit the capability of the graph filter. Extremely,  $\alpha \to 0$  gives a trivial function  $g^K_{\beta}(\lambda) = \beta_0$ . This suggests that  $\alpha$  should be elaborately tuned to improve the performance.

Though the aforementioned graph filter is primarily motivated via spectral analysis, we can still 159 present the spatial perspective explanation for its design. Existing GNNs aggregate the neighbor 160 information of different hops with certain weights, which could be either manually assigned or learned 161 adaptively. Typically, methods like GPR-GNN [2] and DAGNN [19] that learn the aggregation weight, 162 tend to treat the neighbor's information of different hops equally. That is, the k-hop's weight are 163 assigned with  $\theta_k = \mathcal{O}(1)$  for each  $k \in [K]$ . However, it is shown in the previous research that the 164 propagation with the very high-order neighbor potentially leads to the over-smoothing issue [25, 33]. 165 The current methods magnify this flaw of the high-order graph since they cannot distinguish the 166 significance of the information of different hops. This motivates the design of the decay rate in 167 APGNN, i.e., we employ weights with exponential decaying rate by assigning  $\theta_k = \mathcal{O}(\alpha^k)$  for some 168  $0 < \alpha < 1$ . This approach emphasizes the contribution of lower-order neighbors and restricts the 169 over-weighting of the information from high-order neighbors due to  $\theta_k \to 0$  with  $k \to \infty$ . Therefore, 170 it provides more effective aggregation and thus enhances the model's scalability. 171

To take a further step in the construction of a deep GNN, we introduce a multiple P-hop strategy for the graph filter of (11), which effectively extends the utmost neighborhood range that the graph filter can perceive by P times. Consider a different perspective regarding the construction of a filter with the utmost order T = KP. The previous methods can be viewed as a one-hop graph filter by setting

P = 1. For P > 1, the graph filter is able to aggregate information from a larger neighborhood in

the same order. In addition, we will illustrate the advantages of this strategy from the perspective of generalization in the following section.

<sup>179</sup> Summarizing the above analysis, we present the following comprehensive architecture of APGNN:

$$\mathbf{Z} = g_{\boldsymbol{\beta}}(\mathbf{L})f(\mathbf{X}), \quad f(\mathbf{X}) = \mathrm{MLP}(\mathbf{X}), \quad g_{\boldsymbol{\beta}}(\mathbf{L}) = \sum_{k=0}^{K} \beta_k \alpha^k \tilde{\mathbf{A}}^{kP}.$$
(16)

180

In short, APGNN incorporates the benefits from the decay rate  $\alpha$  that exponentially suppresses the information of extremely high-order neighbors and the multiple *P*-hop strategy to enlarge receptive fields. These approaches make it possible to realize a sufficiently deep GNN.

# **184 Generalization analysis**

The theoretical analysis of GNN's generalization is widely studied. [31] provides the generalization result of the algorithmic stability of GCN in the discrete graph setting. In contrast, [14] shows the convergence and stability guarantee over the random and continuous graph. In this section, we will present the uniform generalization bound of the proposed GNN learning framework under the continuous setup.

We first introduce some notations for later discussion. Denote  $\mathbf{x} \in \mathcal{X}$  as any samples from the input space  $\mathcal{X}$  (we generally set  $\mathcal{X}$  as a subset of  $\mathbb{R}^d$ ). Let  $\rho(\cdot)$  be a probability measure defined over  $\mathcal{X}$ . Assume  $x_j$  is the *j*-th coordinate of  $\mathbf{x} \in \mathcal{X}$  and  $\mathbb{E}[x_j^2] \leq c_{\mathcal{X}}^2$  for any  $j \in [d]$ . To describe the graph relation between each pair  $(\mathbf{x}, \mathbf{x}')$  over  $\mathcal{X} \times \mathcal{X}$ , we define a continuous graph function  $A(\cdot, \cdot) : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R}_+$ , and its corresponding degree function is

$$d(\mathbf{x}') = \int_{\mathcal{X}} A(\mathbf{x}, \mathbf{x}') d\rho(\mathbf{x}').$$
(17)

Different from the setting of [18, 26], we assume  $0 \le A(\mathbf{x}, \mathbf{x}') \le c_U$ , and  $0 < c_L \le d(\mathbf{x})$  for any  $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$ . Therefore, we can define the symmetric normalized graph:

$$\tilde{A}(\mathbf{x}, \mathbf{x}') = \frac{A(\mathbf{x}, \mathbf{x}')}{\sqrt{d(\mathbf{x})d(\mathbf{x}')}}.$$
(18)

Then the corresponding normalized Laplacian is  $L = I - \tilde{A}$ , where *I* indicates the identity operator over  $\mathcal{X}$ . For a graph filter function  $g_{\theta}(\lambda) = \sum_{k=0}^{K} \theta_k (1-\lambda)^k$ , graph convolution of the continuous graph is defined as the following integral operator:

$$g_{\boldsymbol{\theta}}Lf = \sum_{k=0}^{K} \theta_k \tilde{A}^k f, \quad \tilde{A}f = \int_{\mathcal{X}} \tilde{A}(\cdot, \mathbf{x}) f(\mathbf{x}) \mathrm{d}\rho(\mathbf{x}), \tag{19}$$

where  $\tilde{A}^k = \tilde{A}^{k-1} \circ \tilde{A}$  denotes k-order composition of integral operator with  $\tilde{A}^0 = I$ . Note we have  $\sum_{k=0}^{K} \theta_k \|\tilde{A}\| \leq \|\boldsymbol{\theta}\|_1 \leq M$  for any  $K \in \mathbb{N}$ , indicating  $\sum_{k=0}^{\infty} \theta_k \tilde{A}$  is absolutely summable. This guarantees the existence of graph filter on the continuous graph when  $K \to \infty$ . For convenience in understanding, we provide the analysis on a simplified GNN, where we consider a semi-supervised learning task with two classes, i.e.,  $y_i \in \mathcal{Y} \triangleq \{-1, 1\}$ , and utilize linear feature extractor  $f(\mathbf{X}) =$   $\mathbf{w}^{\top} \mathbf{X}$ . Note that we can still extend our result for  $f(\mathbf{X}) = MLP(\mathbf{X})$  and multi-class cases using the techniques proposed in [1]. With the above setting, the hypothesis set over is described as

$$\mathcal{H}_{\mathcal{X}} = \{h : h(\mathbf{x}) = g_{\boldsymbol{\theta}} L f(\mathbf{x}), \ f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle, \ \|\mathbf{w}\|_2 \le B, \ \|\boldsymbol{\theta}\|_1 \le M\}.$$
(20)

However, the integral in each hypothesis  $h \in \mathcal{H}_{\mathcal{X}}$  is intractable since the underlying graph function and the data distribution are unknown. Therefore, we should use the "empirical version" of the hypothesis to estimate  $h \in \mathcal{H}_{\mathcal{X}}$ . For this reason, we introduce the hypothesis set defined over the observed samples S and graph  $\mathcal{G}$ :

$$\mathcal{H}_{S} = \left\{ h : h(\mathbf{x}_{i}) = \sum_{j=1}^{n} g_{\boldsymbol{\theta}}(\mathbf{L})_{ij} \mathbf{x}_{j}^{\top} \mathbf{w}, \quad \|\mathbf{w}\|_{2} \le B, \ \|\boldsymbol{\theta}\|_{1} \le M \right\}.$$
 (21)

<sup>211</sup> Define the generalization error and the empirical error [21] as follows

$$R(h) = \mathbb{E}_{(\mathbf{x},y)}[1_{yh(\mathbf{x}) \le 0}], \quad \hat{R}(h) = \frac{1}{n_l} \sum_{i=1}^{n_l} \min(1, \max(0, 1 - y_i h(\mathbf{x}_i))).$$
(22)

<sup>212</sup> We have the following theorem on the generalization of the proposed learning paradigm.

**Theorem 2.** Suppose  $g_{\theta}(\cdot)$  is  $L_M$ -Lipschitz. Let  $h_{\mathbf{w},\theta} \in \mathcal{H}_X$  and  $h_{\mathbf{w},\theta} \in \mathcal{H}_S$  share the same parameter  $(\mathbf{w}, \theta)$ . Then there exists a constant C > 0 related to the graph function, with the probability at least  $1 - \delta$ , the following inequality holds.

$$R(h_{\mathbf{w},\boldsymbol{\theta}}) \lesssim \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}}) + 2BMc_{\mathcal{X}}\sqrt{\frac{2d\log(2K+2)}{n_l}} + BCL_M dc_{\mathcal{X}}\sqrt{\frac{\log(2/\delta)}{n}}.$$
 (23)

The proof is given by excess risk decomposition, shown in Appendix. The notation "<" denotes 216 "less than or approximately equal to the right-hand side" and guarantees an approximation error of 217 at most  $\mathcal{O}(\sqrt{\frac{\log(1/\tau)}{n_l}})$  with a probability of at least  $1 - \mathcal{O}(\tau)$ . We remind readers the important 218 difference between  $R(h_{\mathbf{w},\theta})$  and  $\hat{R}(\hat{h}_{\mathbf{w},\theta})$ . The former term measures the population error over the 219 whole input space with the **continuous** graph filter  $g_{\theta}L$ . In contrast,  $\hat{R}(\hat{h}_{w,\theta})$  is the empirical risk 220 (i.e., training risk) on the sample set S with the **discrete** graph filter  $g_{\theta}(\mathbf{L})$ .  $h_{\mathbf{w},\theta}$  shares the same 221 learning parameter with  $\hat{h}_{\mathbf{w},\theta}$ . Therefore, the minimization of the right-hand-side of (23) w.r.t ( $\mathbf{w}, \theta$ ) 222 reduces the upper bound of the population error. 223

We observe the first term of generalization bound is of order  $\mathcal{O}((dn_l^{-1}\log K)^{1/2})$ , which outlines the 224 model's complexity. Although it becomes infinity when  $K \to \infty$ , the growth of this term is extremely 225 slow as K increases. In practice, we generally set K < n since the neighbor information beyond 226 *n*-hops is redundant, restricting the complexity away from infinity. Therefore, the generalization of 227 the model is rigorously guaranteed for sufficiently large K, which allows us to construct significantly 228 deep GNN in the proposed framework. We can obtain a more precise estimation for a certain model. 229 In the following proposition, we unveil the generalization of APGNN as a direct application of 230 Theorem 2. 231

**Proposition 1.** Let  $\beta \in \mathbb{R}^{K}$  and  $g_{\beta}^{K}(\lambda) = \sum_{k=0}^{K} \beta_{k} \alpha^{k} (1-\lambda)^{k}$  where  $0 < \alpha < 1$  and  $\|\beta\|_{\infty} \leq 1$ . with the probability at least  $1 - \delta$ , the following inequality holds.

$$R(h_{\mathbf{w},\boldsymbol{\beta}}) \lesssim \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\beta}}) + \frac{2Bc_{\mathcal{X}}(1-\alpha^{K})}{1-\alpha} \sqrt{\frac{2d\log(2K+2)}{n_{l}} + \frac{BCdc_{\mathcal{X}}\alpha}{(1-\alpha)^{2}} \sqrt{\frac{\log(2/\delta)}{n}}}.$$
 (24)

234 *Proof.* This is a direct result with  $M = (1 - \alpha^K)/(1 - \alpha)$  and  $L_M = \alpha/(1 - \alpha)^2$  in Theorem 2.  $\Box$ 

In (24), the complexity term becomes  $\mathcal{O}(\sqrt{\log K}(1-\alpha^K))$  with K = |T/P|, which is relatively 235 tighter than  $\mathcal{O}(\sqrt{\log K})$ . For this term, we promote further discussion with P-hop. Since it takes 236 |T/P| steps to reach the T-order graph, the term becomes  $\mathcal{O}(\sqrt{\log |T/P|}(1-\alpha^{\lfloor T/P \rfloor}))$ . It is 237 observed that the term decreases as P increases. Therefore, the appropriate P reduces the bound, 238 explaining the mechanism of the P-hop method. On the other hand, larger  $\alpha$  leads to a higher bound. 239 From the point of spatial view, the information from high-order neighbors is underused, which limits 240 the range of the graph filter. Thus  $\alpha$  should be moderate to leverage the generalization and the 241 capability of the model. 242

# 243 5 Experiment

In this section, we conduct node classification experiments on various benchmark datasets to evaluate the performance of APGNN. Specifically, we compare our method with state-of-the-art methods and display the corresponding learned graph filter on different data sets. Moreover, to validate the theoretical analysis, the influence of parameters K,  $\alpha$ , and P is also investigated in experiments.

Model	Dataset									
Widdei	Cora	Citeseer	Pubmed	Wiki-CS	MS-Academic					
MLP	$57.79_{\pm 0.11}$	$61.20_{\pm 0.08}$	$73.23_{\pm 0.05}$	$65.66_{\pm 0.20}$	$87.79_{\pm 0.42}$					
ChebNet	$79.92{\scriptstyle \pm 0.18}$	$70.90{\scriptstyle \pm 0.37}$	$76.98{\scriptstyle \pm 0.16}$	$63.24_{\pm 1.43}$	$90.76 \pm 0.73$					
GCN	$82.03 \pm 0.27$	$71.05{\scriptstyle \pm 0.33}$	$79.26{\scriptstyle \pm 0.18}$	$72.05 _{\pm 0.45}$	$92.07_{\pm 0.13}$					
SGC	$81.89_{\pm 0.26}$	$\underline{72.18}_{\pm 0.24}$	$78.58_{\pm 0.15}$	$72.76_{\pm 0.35}$	$89.01_{\pm 0.40}$					
GAT	$82.82_{\pm 0.36}$	$71.96_{\pm 0.39}$	$79.15_{\pm 0.34}$	$74.36_{\pm0.58}$	$91.86_{\pm 0.27}$					
GraphSage	$82.14_{\pm 0.25}$	$71.80_{\pm 0.36}$	$79.20_{\pm 0.27}$	$73.17_{\pm 0.41}$	$91.53_{\pm 0.15}$					
PPNP	$83.73{\scriptstyle \pm 0.31}$	$71.74_{\pm 0.44}$	$80.28 _{\pm 0.22}$	$74.69{\scriptstyle \pm 0.53}$	$92.58{\scriptstyle \pm 0.06}$					
APPNP	$83.73_{\pm 0.21}$	$71.70_{\pm 0.21}$	$80.07_{\pm 0.21}$	$74.91_{\pm 0.61}$	$92.81_{\pm 0.12}$					
GNN-LF(iter)	$83.83 \pm 0.36$	$71.44_{\pm 0.42}$	$\underline{80.31}_{\pm 0.16}$	$75.19_{\pm 0.49}$	$92.78_{\pm 0.22}$					
GNN-HF(iter)	$83.68_{\pm 0.31}$	$71.58_{\pm 0.36}$	$79.99_{\pm 0.22}$	$74.71_{\pm 0.55}$	$92.72_{\pm 0.31}$					
DAGNN	$82.70_{\pm 0.17}$	$71.90{\scriptstyle \pm 0.06}$	$80.06 \pm 0.30$	$\underline{75.63}_{\pm 0.48}$	$93.24 \pm 0.21$					
Ours	$84.15_{\pm 0.23}$	<b>72.44</b> $_{\pm 0.56}$	$80.74_{\pm 0.24}$	<b>76.03</b> $_{\pm 0.51}$	<b>93.69</b> ±0.20					

Table 2: The average accuracy (%) and standard deviation (%) on five benchmark datasets. The highest accuracy in each column is shown in bold, while the second-best result is underlined.

#### 248 5.1 Experiment Setup

**Datasets.** We perform experiments on five benchmark datasets commonly used in node classification tasks. 1). **Cora, Citeseer, Pubmed**[29, 35]: These are three standard citation networks where each node is a paper and each edge is a citation link. 2).**Wiki-CS**[20]: This dataset defines the computer science articles as nodes, while the hyperlinks are edges. 3). **MS Acadamic**[16]: The nodes represent the author and the edges represent the co-authorships. A co-authorship Microsoft Academic Graph, where the nodes are the bag-of-words representation of the papers' abstract and edges are co-authorship. The data statistics and their partitions are presented in Appendix.



Figure 2: The graph filters learned on different data sets, with the parameter P being odd in subfigure (a) and even in subfigure (b).

Baselines. To evaluate the effectiveness of APGNN, we compare it with the following baseline
models: 1) MLP [22], a traditional method that does not use graphs, 2) GAT [30] and GraphSAGE
[10], spatial methods that aggregate neighborhoods' information, and 3) ChebNet [3], GCN [15],
SGC [33], PPNP, APPNP[16], GNN-LF (iteration form), GNN-HF (iteration form) [36], and DAGNN
[19], spectral methods analyzing GNNs with graph Fourier transform.

Settings. We conducted 10 runs for each method on each dataset, with a hidden dimension of 64. For all compared methods, their parameter settings follow the previous practices [19, 36]: the dropout rate is 0.5 except for Cora, which had a rate of 0.8. Furthermore, the learning rate is 0.01 for Cora, Citeseer, and Pubmed, but 0.03 for Wiki-CS and 0.02 for MS-Academic, while the weight decay is 0.005 for Cora and Pubmed, 0.02 for Citeseer, 0.0005 for Wiki-CS, and 0.00525 for MS-Academic. We fix the polynomial order K to 10 in ChebNet, APPNP, GNN-LF, GNN-HF, DAGNN, and APPNP. The best hyperparameters we choose for APGNN are presented in Appendix. To ensure a fair comparison with the compared methods, we also applied our optimal hyperparameters to them,
 selecting the maximum value to display.



Figure 3: Accuracy with different (a) K. (b)  $\alpha$ . (c) P (for fixing T).

#### 270 5.2 Analysis

Node Classification. As the metric for evaluation, the mean accuracy of 10 runs is used. We compare the performance of APGNN with other methods on five benchmark datasets. Experiment results are reported in Table 2. We can observe that APGNN achieves the highest accuracy across all five datasets, demonstrating its superior performance.

**Learnable Graph Filters.** Figure 2 shows the graph filters learned on various datasets via APGNN. When the parity of P varies, the graph filter has a distinctive shape. However, their shapes exhibit minimal impact on their accuracy regardless of the parity of P according to the experiment results. Moreover, the graph filters of each dataset are plotted in Appendix, more details are included in Appendix. Our results show that the graph filters learned from different datasets vary in detail, even when their parameters have similar parity, demonstrating the efficacy of APGNN in learning task-specific graph filters.

**Polynomial Order** K. To gain insight into the role of polynomial order K, we conduct the experiment tuning K in  $\{1, 2, ..., 20\}$  on Cora, Citeseer, and Pubmed dataset. Our theoretical analysis supports the observation that a small K can result in a large truncation error, leading to a low accuracy rate. It can be observed that the accuracy rate has little promotion when K is larger than 10, although at the cost of high computational resources.

**Decay Rate**  $\alpha$ . Figure 3 (b) depicts the accuracy curve corresponding to various  $\alpha$  values ranging from 0.1 to 0.9 and 0.99 on Cora, Citeseer and Pubmed datasets. As  $\alpha$  decreases, the classification accuracy initially increases and then declines sharply. This phenomenon verifies the theory that the truncation error decreases as  $\alpha$  decreases, but it leads to a trivial function when  $\alpha$  is extremely small.

**P-hop strategy.** We investigate the accuracy associated with varying parameters P taken from the set {1, 2, 3, 4, 5, 6} when fixing T = KP = 60. As we can see in Figure 3 (c), the accuracy increase when P > 1. This phenomenon can be attributed to the fact that the generalization bounding decreases when P increases, which suggests that the P-hop strategy can effectively explore deeper information with the same computational complexity.

# 296 6 Conclusion

This paper proposes a universal learning principle for a valid construction of GNN. An instantiation 297 named APGNN is proposed to verify the effectiveness of our framework. APGNN employs a decay 298 rate and a multiple P-hop strategy to learn the coefficients adaptively, which can efficiently aggregate 299 the information from high-order neighbors. We present a theoretical analysis of the generalization 300 capabilities of both our framework and APGNN, which provides a learning guarantee. Comprehensive 301 experiments show the superior performance of APGNN. In the future, it is worth exploring diverse 302 graph filters based on the proposed principle. As shown in the generalization analysis, the upper 303 bound of the model complexity relies on  $\mathcal{O}(\sqrt{\log K})$ . How to devise the GNN with complexity free 304 of the hyperparameter K is also a meaningful research direction. 305

# **306** References

- [1] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds
   for neural networks. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 30, page 6241–6250, 2017.
- [2] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. Adaptive universal generalized
   pagerank graph neural network. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [3] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
   on graphs with fast localized spectral filtering. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 29, 2016.
- [4] Leyan Deng, Defu Lian, Chenwang Wu, and Enhong Chen. Graph convolution network based
   recommender systems: Learning guarantee and item mixture powered strategy. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 35,
   pages 3900–3912, 2022.
- [5] Jiarui Feng, Yixin Chen, Fuhai Li, Anindya Sarkar, and Muhan Zhang. How powerful are k-hop
   message passing graph neural networks. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 35, pages 4776–4790, 2022.
- [6] Fernando Gama, Joan Bruna, and Alejandro Ribeiro. Stability properties of graph neural networks. *IEEE Transactions on Signal Processing*, 68:5680–5695, 2020.
- [7] Francesco Giuliari, Geri Skenderi, Marco Cristani, Yiming Wang, and Alessio Del Bue. Spatial
   commonsense graph for object localisation in partial scenes. In 2022 IEEE/CVF Conference on
   *Computer Vision and Pattern Recognition (CVPR)*, pages 19496–19505, 2022.
- [8] Kan Guo, Yongli Hu, Yanfeng Sun, Sean Qian, Junbin Gao, and Baocai Yin. Hierarchical graph
   convolution network for traffic forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 35, pages 151–159, 2021.
- [9] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. Attention based spatial temporal graph convolutional networks for traffic flow forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, volume 33, pages 922–929, 2019.
- [10] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
   graphs. In *Proceedings of the Annual Conference on Neural Information Processing Systems* (*NeurIPS*), volume 30, 2017.
- [11] Kai Han, Yunhe Wang, Jianyuan Guo, Yehui Tang, and Enhua Wu. Vision GNN: An image
   is worth graph of nodes. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 35, pages 8291–8303, 2022.
- [12] Mingguo He, Zhewei Wei, zengfeng Huang, and Hongteng Xu. Bernnet: Learning arbitrary
   graph spectral filters via bernstein approximation. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 34, pages 14239–14251, 2021.
- [13] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, YongDong Zhang, and Meng Wang. LightGCN:
   Simplifying and powering graph convolution network for recommendation. In *Proceedings* of the International ACM SIGIR Conference on Research and Development in Information
   Retrieval (SIGIR), page 639–648, 2020.
- [14] Nicolas Keriven, Alberto Bietti, and Samuel Vaiter. Convergence and stability of graph convolutional networks on large random graphs. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, volume 33, pages 21512–21523, 2020.
- [15] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- Iohannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:
   Graph neural networks meet personalized pagerank. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.
- [17] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks
   for semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence* (AAAI), 2018.

- [18] Shaojie Li, Sheng Ouyang, and Yong Liu. Understanding the generalization performance of
   spectral clustering algorithms. *arXiv preprint arXiv:2205.00281*, 2022.
- [19] Meng Liu, Hongyang Gao, and Shuiwang Ji. Towards deeper graph neural networks. In
   *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, page 338–348, 2020.
- Péter Mernyei and Cătălina Cangea. Wiki-CS: A wikipedia-based benchmark for graph neural
   networks. *arXiv preprint arXiv:2007.02901*, 2020.
- [21] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*.
   MIT press, 2018.
- [22] Sankar K Pal and Sushmita Mitra. Multilayer perceptron, fuzzy sets, and classification. *IEEE Transactions on neural networks*, 3(5):683–697, 1992.
- Yitong Pang, Lingfei Wu, Qi Shen, Yiming Zhang, Zhihua Wei, Fangli Xu, Ethan Chang,
   Bo Long, and Jian Pei. Heterogeneous global graph neural networks for personalized session based recommendation. In *Proceedings of the ACM International Conference on Web Search* and Data Mining (WSDM), page 775–783, 2022.
- Patricia Pauli, Anne Koch, Julian Berberich, Paul Kohler, and Frank Allgöwer. Training robust
   neural networks using lipschitz bounds. *IEEE Control Systems Letters*, 6:121–126, 2021.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Dropedge: Towards deep graph
   convolutional networks on node classification. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- [26] Lorenzo Rosasco, Mikhail Belkin, and Ernesto De Vito. On learning with integral operators.
   *Journal of Machine Learning Research (JMLR)*, 11(2), 2010.
- [27] Aliaksei Sandryhaila and José M. F. Moura. Discrete signal processing on graphs: Graph filters.
   In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing* (*ICASSP*), pages 6163–6166, 2013.
- [28] Aliaksei Sandryhaila and José M. F. Moura. Discrete signal processing on graphs: Graph fourier
   transform. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6167–6170, 2013.
- [29] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi Rad. Collective classification in network data. *AI Magazine*, 29(3):93, 2008.
- [30] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
   Bengio. Graph attention networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2018.
- [31] Saurabh Verma and Zhi-Li Zhang. Stability and generalization of graph convolutional neural networks. In *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, pages 1539–1548, 2019.
- [32] Xiyuan Wang and Muhan Zhang. How powerful are spectral graph neural networks. In
   *Proceedings of the International Conference on International Conference on Machine Learning* (*ICML*), volume 162, pages 23341–23362, 2022.
- [33] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.
   Simplifying graph convolutional networks. In *Proceedings of the International Conference on International Conference on Machine Learning (ICML)*, 2019.
- [34] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
   networks? In *Proceedings of the International Conference on Learning Representations (ICLR)*,
   2019.
- [35] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning
   with graph embeddings. In *Proceedings of the International Conference on International Conference on Machine Learning (ICML)*, page 40–48, 2016.
- [36] Meiqi Zhu, Xiao Wang, Chuan Shi, Houye Ji, and Peng Cui. Interpreting and unifying graph
   neural networks with an optimization framework. In *Proceedings of The International World Wide Web Conference (WWW)*, page 1215–1226, 2021.

- [37] Stefano Zorzi, Shabab Bazrafkan, Stefan Habenschuss, and Friedrich Fraundorfer. PolyWorld: Polygonal building extraction with graph neural networks in satellite images. In *Proceedings*
- of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages
- 1848–1857, 2022.

# 414 A Appendix

415 A.1 Data statistics

Dataset	Nodes	Edges	Features	Class	Train	Val	Test
cora	2708	5429	1433	7	140	500	1000
citeseer	3327	4732	3703	6	120	500	1000
pubmed	19717	44338	500	3	60	500	1000
wiki-cs	11701	216123	300	10	200	500	1000
ms academic	18333	81894	6805	15	300	500	1000

Table 3: Data statistics for the node classification task.

# 416 A.2 The detail of learned graph filters



Figure 4: The graph filters learned using different data sets, with parameter P being odd.



Figure 5: The graph filters learned using different data sets, with parameter P being even.

# 417 A.3 The hyperparameters settings

Table 4: The hyperparameters of APGNN on various datasets when parameter P is odd.

Dataset	K	P	α	Weight decay	Learning rate	Dropout rate
Cora	10	3	0.7	0.005	0.01	0.8
Citeseer	10	5	0.7	0.00625	0.01	0.5
Pubmed	10	5	0.9	0.005	0.01	0.5
Wiki-CS	10	1	0.9	0.000525	0.03	0.4
MS-Academic	10	3	0.9	0.00525	0.02	0.4

# 418 A.4 Proof for Lemma 1

( $\Rightarrow$ ). We show the result by contradiction. If  $\sum_{k}^{\infty} |a_k|$  is not convergent, then at  $\gamma = 1$ , we have  $\sum_{k}^{\infty} |a_k \gamma^k|$  is not convergent, which occurs a contradiction. Therefore, series  $\sum_{k}^{\infty} |a_k|$  converges.

421 ( $\Leftarrow$ ). It is obvious that for  $\forall \gamma \in (-1, 1], \sum_{k}^{\infty} |a_k \lambda^k| \le \sum_{k}^{\infty} |a_k|$ . Therefore,  $\sum_{k}^{\infty} |a_k \lambda^k|$  uniformly 422 converges in  $\lambda \in (-1, 1]$ .

Table 5: The hyperparameters of APGNN on various datasets when parameter P is even.

Dataset	K	$\boldsymbol{P}$	$\alpha$	Weight decay	Learning rate	Dropout rate
Cora	10	4	0.7	0.005	0.01	0.8
Citeseer	10	6	0.7	0.00625	0.01	0.5
Pubmed	10	6	0.9	0.005	0.01	0.5
Wiki-CS	10	2	0.9	0.000525	0.03	0.4
MS-Academic	10	2	0.9	0.00525	0.02	0.4

### 423 A.5 Proof for Theorem 1

<sup>424</sup>  $\tilde{\mathbf{A}}$  is an adjacency matrix of a graph, which is a real symmetric matrix. Since we can decompose  $\tilde{\mathbf{A}}$ <sup>425</sup> as  $\tilde{\mathbf{A}} = \mathbf{U}\Gamma\mathbf{U}^{\top}$ , where  $\mathbf{U}$  is a matrix composed of the eigenvectors of  $\tilde{\mathbf{A}}$  and  $\Gamma = \text{diag}(\gamma_1, \dots, \gamma_n)$ <sup>426</sup> is the diagonal matrix of the corresponding eigenvalues. Therefore, we have

$$g(\mathbf{L}) = \sum_{k=1}^{\infty} \theta_k \tilde{\mathbf{A}}^k = \mathbf{U} \operatorname{diag} \left( \sum_{k=1}^{\infty} \theta_k \gamma_1^k, \ \cdots, \ \sum_{k=1}^{\infty} \theta_k \gamma_n^k \right) \mathbf{U}^{\top}$$
(25)

Therefore, the  $g(\mathbf{L})$  converges absolutely and uniformly if and only if  $\sum_{k=1}^{\infty} \theta_k \gamma_i^k$  converges absolutely and uniformly for all  $i \in [n]$ . Then apply Lemma 1 and we can obtain the result.

# 429 A.6 Proof of Theorem 2

- 430 We first introduce some definitions and Lemma for assisting with the proof.
- **Definition 1.** Consider the sample set  $S = {\mathbf{x}_1, \dots, \mathbf{x}_n}$  and function set  $\mathcal{F}$ , where  $f(\mathbf{x})$  is bounded for any  $f \in \mathcal{F}$ . Then the empirical Rademacher complexity is defined as:

$$\mathfrak{R}_{S}(\mathcal{F}) = \frac{1}{n} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{f \in \mathcal{F}} \sum_{i=1}^{n} \sigma_{i} f(\mathbf{x}_{i}) \right],$$
(26)

433 where  $\sigma_i$  is i.i.d. Rademacher random variable defined by  $Pr(\sigma_i = -1) = Pr(\sigma_i = 1) = 0.5$ .

# 434 Lemma 2. Consider the hypothesis set

$$\mathcal{H}_{\mathcal{X}} = \{h : h(\mathbf{x}) = g_{\boldsymbol{\theta}} L f(\mathbf{x}), \ f(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle, \ \|\mathbf{w}\|_2 \le B, \ \|\boldsymbol{\theta}\|_1 \le M\},$$
(27)

where  $\theta = [\theta_0, \theta_1, \dots, \theta_K]$ . Let  $x_j$  denote the *j*-th element of  $\mathbf{x} \in \mathcal{X}$ , and  $\mathbb{E}[x_j^2] \leq c_{\mathcal{X}}$  for any  $j \in [d]$ . Then for any sample set  $S = {\mathbf{x}_1, \dots, \mathbf{x}_{n_l}} \subset \mathcal{X}$  we have

$$\Re_S(\mathcal{F}_{\mathcal{X}}) \lesssim 2BMc_{\mathcal{X}} \sqrt{\frac{2\log(2K+2)}{n_l}}.$$
(28)

437 *Proof.* Based on the definition, we can write

$$\begin{aligned} \mathfrak{R}_{S}(\mathcal{H}_{\mathcal{X}}) &= \frac{1}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\substack{h_{\mathbf{w},\boldsymbol{\theta}} \in \mathcal{H}_{\mathcal{X}} \\ n_{\mathbf{w},\boldsymbol{\theta}} \in \mathcal{H}_{\mathcal{X}}} \sum_{i=1}^{n_{l}} \sigma_{i} h_{\mathbf{w},\boldsymbol{\theta}}(\mathbf{x}_{i}) \right] \\ &= \frac{1}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\substack{\|\mathbf{w}\|_{2} \leq B, \|\boldsymbol{\theta}\|_{1} \leq M}} \sum_{i=1}^{n_{l}} \sigma_{i} g_{\boldsymbol{\theta}} L f(\mathbf{x}_{i}) \right] \\ &= \frac{1}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\substack{\|\mathbf{w}\|_{2} \leq B, \|\boldsymbol{\theta}\|_{1} \leq M}} \sum_{i=1}^{n_{l}} \sigma_{i} \int_{\mathcal{X}} \sum_{k=0}^{K} \theta_{k} \tilde{A}^{k}(\mathbf{x}_{i}, \mathbf{x}) \mathbf{x}^{\top} \mathbf{w} d \rho(\mathbf{x}) \right] \\ &\leq \frac{B}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\substack{\|\boldsymbol{\theta}\|_{1} \leq M}} \left\| \sum_{i=1}^{n_{l}} \sigma_{i} \int_{\mathcal{X}} \sum_{k=0}^{K} \theta_{k} \tilde{A}^{k}(\mathbf{x}_{i}, \mathbf{x}) \mathbf{x} d \rho(\mathbf{x}) \right\|_{2} \right] \\ &= \frac{B}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\substack{\|\boldsymbol{\theta}\|_{1} \leq M}} \left\| \sum_{i=1}^{n_{l}} \sigma_{i} \mathbf{v}_{i} \right\|_{2} \right] \end{aligned}$$

where the inequality follows from the Cauchy-Schwarz inequality, and the set V is defined as

$$V \triangleq \left\{ \{\mathbf{v}_i\}_{i=1}^n : \mathbf{v}_i = \int_{\mathcal{X}} \sum_{k=0}^K \theta_k \tilde{A}^k(\mathbf{x}_i, \mathbf{x}) \mathbf{x} \mathrm{d}\rho(\mathbf{x}), \quad \|\boldsymbol{\theta}\|_1 \le M \right\}.$$
 (29)

Define  $q_j(\mathbf{x}) = x_j$  returning the *j*-th coordinate of the input. Hence, the *j*-th coordinate of  $\mathbf{v}_i$  can be rewritten as

$$v_{ij} = \int_{\mathcal{X}} \sum_{k=0}^{K} \theta_k \tilde{A}^k(\mathbf{x}_i, \mathbf{x}) x_j \mathrm{d}\rho(\mathbf{x}) = \int_{\mathcal{X}} \sum_{k=0}^{K} \theta_k \tilde{A}^k(\mathbf{x}_i, \mathbf{x}) q_j(\mathbf{x}) \mathrm{d}\rho(\mathbf{x}) = \sum_{k=0}^{K} \theta_k \tilde{A}^k q_j(\mathbf{x}_i).$$
(30)

441 Since  $\|\mathbf{u}\|_2 \leq \sqrt{d} \|\mathbf{u}\|_\infty$  for any  $\mathbf{u} \in \mathbb{R}^d$ , we have

$$\begin{aligned} \Re_{S}(\mathcal{H}_{\mathcal{X}}) &\leq \frac{B\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\{\mathbf{v}_{i}\}_{i=1}^{n} \in V} \left\| \sum_{i=1}^{n_{l}} \sigma_{i} \mathbf{v}_{i} \right\|_{\infty} \right] \\ &\leq \frac{B\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\{\mathbf{v}_{i}\}_{i=1}^{n} \in V} \max_{j \in [d]} \left| \sum_{i=1}^{n_{l}} \sigma_{i} v_{ij} \right| \right] \\ &\leq \frac{2B\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\{\mathbf{v}_{i}\}_{i=1}^{n} \in V} \max_{j \in [d]} \sum_{i=1}^{n_{l}} \sigma_{i} v_{ij} \right] \\ &\leq \frac{2B\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\|\|\boldsymbol{\theta}\|_{1} \leq M} \sum_{i=1}^{n_{l}} \sigma_{i} \sum_{k=0}^{K} \theta_{k} \tilde{A}^{k} q_{j}(\mathbf{x}_{i}) \right] \\ &= \frac{2B\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\|\|\boldsymbol{\theta}\|_{1} \leq M} \sum_{k=0}^{K} \theta_{k} \sum_{i=1}^{n_{l}} \sigma_{i} \tilde{A}^{k} q_{j}(\mathbf{x}_{i}) \right] \\ &= \frac{2BM\sqrt{d}}{n_{l}} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\boldsymbol{\theta} \in \boldsymbol{\Theta}} \sum_{k=0}^{K} \theta_{k} \sum_{i=1}^{n_{l}} \sigma_{i} \tilde{A}^{k} q_{j}(\mathbf{x}_{i}) \right] \\ &= 2BM\sqrt{d} \Re_{S}(\mathcal{H}'), \end{aligned}$$

where  $\Theta = \bigcup_{k=0}^{K} \{-e_k, e_k\}$  and  $e_k$  denote k-th vector with k-th entry as one and others are zero. The set  $\mathcal{H}' = \{h(\mathbf{x}) = \sum_{k=0}^{K} \theta_k A^k q_j(\mathbf{x}) : \boldsymbol{\theta} \in \Theta\}$  is a finite set with  $|\mathcal{H}'| = 2(K+1)$ . We bound  $\mathfrak{R}_S(\mathcal{H}')$  with Massart's Lemma:

Lemma 3 (Massart's Lemma [21]). Let  $\mathcal{X} \subset \mathbb{R}^n$  be a finite set and  $\sup_{\mathbf{x} \in \mathcal{X}} \|\mathbf{x}\|_2 \leq r\sqrt{n}$ , then the following inequality holds:

$$\mathbb{E}_{\boldsymbol{\sigma}}\left[\frac{1}{n}\left\langle\boldsymbol{\sigma},\mathbf{x}\right\rangle\right] \leq r\sqrt{\frac{2\log|\mathcal{X}|}{n}},\tag{31}$$

447 where  $\boldsymbol{\sigma} = [\sigma_1, \cdots, \sigma_n]$  denote the vector of Rademacher random variables.

448 Since  $\mathcal{H}'$  is a finite set and for any  $h \in \mathcal{H}'$ ,

$$\frac{1}{n_l} \sum_{i=1}^{n_l} h(\mathbf{x}_i)^2 = \frac{1}{n_l} \sum_{i=1}^{n_l} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}_i) \right]^2$$
$$\approx \int_{\mathcal{X}} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}) \right]^2 \mathrm{d}\rho(\mathbf{x}) \le \|q_j\|^2 \le c_{\mathcal{X}}^2$$

449 where we use  $\|\tilde{A}^k q_j\| \le \|\tilde{A}^k\| \|q_j\|$  and  $\|\tilde{A}^k\| \le 1$  for any  $k \in [n]$ , and

$$\|q_j\|^2 = \int_{\mathcal{X}} q_j(\mathbf{x})^2 \mathrm{d}\rho(\mathbf{x}) = \int_{\mathcal{X}} x_j^2 \mathrm{d}\rho(\mathbf{x}) = \mathbb{E}[x_j^2] \le c_{\mathcal{X}}^2.$$
(32)

450 Therefore we finally obtain

$$\Re_S(\mathcal{F}_{\mathcal{X}}) \lesssim 2BMc_{\mathcal{X}}\sqrt{\frac{2d\log(2K+2)}{n_l}}$$
(33)

As a remark, we can present a more precise bound through McDiarmid's inequality. consider the 451 convergence 452

$$\frac{1}{n_l} \sum_{i=1}^{n_l} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}_i) \right]^2 \to \int_{\mathcal{X}} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}) \right]^2 \mathrm{d}\rho(\mathbf{x}).$$
(34)

/

With the probability of at least  $1 - \delta$ , 453

$$\frac{1}{n_l} \sum_{i=1}^{n_l} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}_i) \right]^2 \le \int_{\mathcal{X}} \sup_{k \in [n]} \left[ \tilde{A}^k q_j(\mathbf{x}) \right]^2 \mathrm{d}\rho(\mathbf{x}) - \mathcal{O}\left(\sqrt{\frac{\log 1/\delta}{n_l}}\right).$$
(35)

- The details are omitted since it is not the major part of the analysis. 454
- *Proof.* We first write the excess risk decomposition: 455

$$R(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}}) = \underbrace{R(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(h_{\mathbf{w},\boldsymbol{\theta}})}_{A \text{ part}} + \underbrace{\hat{R}(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}})}_{B \text{ part}}$$
(36)

For the A part, we first apply Theorem 5.8 in [21]. With probability at least  $1 - \delta$ , 456

$$R(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(h_{\mathbf{w},\boldsymbol{\theta}}) \le \Re_{S}(\mathcal{H}_{\mathcal{X}}) + 3\sqrt{\frac{\log(2/\delta)}{2n_{l}}}.$$
(37)

Since the last term is of order  $\mathcal{O}(\sqrt{\log(1/\delta)n_l^{-1}})$ , which is significantly smaller than  $\mathfrak{R}_S(\mathcal{H}_{\mathcal{X}})$ , we 457 rewrite the above inequality as 458

$$R(h_{\mathbf{w},\boldsymbol{\theta}}) \lesssim \hat{R}(h_{\mathbf{w},\boldsymbol{\theta}}) + 2BMc_{\mathcal{X}}\sqrt{\frac{2d\log(2K+2)}{n_l}},\tag{38}$$

- where we replace the Rademacher complexity with its upper bound by Lemma 2. 459
- For the *B* part, we first define the empirical operator over  $S = {\mathbf{x}_1, \cdots, \mathbf{x}_n}$ , 460

$$L_n f = \frac{1}{n} \sum_{i=1}^n \frac{A(\mathbf{x}_i, \cdot)}{\sqrt{d_n(\mathbf{x}_i)d_n(\cdot)}} f(\mathbf{x}_i), \quad d_n = \frac{1}{n} \sum_{i=1}^n A(\mathbf{x}_i, \cdot)$$
(39)

461 and  $g_{\theta}L_n = \sum_{k=0}^{K} \theta_k L_n^k$ . Then we have

$$\begin{split} \hat{R}_{=0} \theta_k L_n^k. \text{ Then we have} \\ \hat{R}(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}}) &\leq \left| \hat{R}(h_{\mathbf{w},\boldsymbol{\theta}}) - \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}}) \right| \\ &\leq \frac{1}{n_l} \left| \sum_{i=1}^{n_l} h_{\mathbf{w},\boldsymbol{\theta}}(\mathbf{x}_i) - \hat{h}_{\mathbf{w},\boldsymbol{\theta}}(\mathbf{x}_i) \right| \\ &= \frac{1}{n_l} \left| \sum_{i=1}^{n_l} g_{\boldsymbol{\theta}} L f(\mathbf{x}_i) - g_{\boldsymbol{\theta}} L_n f(\mathbf{x}_i) \right| \\ &\leq \frac{1}{n_l} \left[ \sum_{i=1}^{n_l} (g_{\boldsymbol{\theta}} L f(\mathbf{x}_i) - g_{\boldsymbol{\theta}} L_n f(\mathbf{x}_i))^2 \right]^{1/2} \\ &\approx \|g_{\boldsymbol{\theta}} L f - g_{\boldsymbol{\theta}} L_n f\| \\ &\leq \|g_{\boldsymbol{\theta}} L - g_{\boldsymbol{\theta}} L_n\| \|f\|. \end{split}$$

where the second inequality follows from the Lipschitz property. With Cauchy-Schwarz inequality, 462

$$\|f\|_{2} = \int_{\mathcal{X}} \langle \mathbf{w}, \mathbf{x} \rangle \mathrm{d}\rho(\mathbf{x}) \le B \cdot \mathbb{E}_{\mathbf{x}}[\|\mathbf{x}\|_{2}] \le B dc_{\mathcal{X}}.$$
(40)

According to Theorem 15 of [26], there exists a proper constant C > 0 related to  $A(\cdot, \cdot)$ , such that 463

$$||L - L_n|| \le ||L - L_n||_{HS} \le C\sqrt{\frac{\log(2/\delta)}{n}}.$$
 (41)

with probability at least  $1 - \delta$ . Since the polynomial  $g_{\theta}$  is  $L_M$ -Lipschitz, we have

$$\|g_{\theta}L - g_{\theta}L_n\| \le L_M C \sqrt{\frac{\log(2/\delta)}{n}}$$
(42)

<sup>465</sup> Combining the above results, one can conclude that for any  $(h_{\mathbf{w},\boldsymbol{\theta}}, \hat{h}_{\mathbf{w},\boldsymbol{\theta}}) \in \mathcal{H}_{\mathcal{X}} \times \mathcal{H}_{S}$ ,

$$R(h_{\mathbf{w},\boldsymbol{\theta}}) \lesssim \hat{R}(\hat{h}_{\mathbf{w},\boldsymbol{\theta}}) + 2BMc_{\mathcal{X}}\sqrt{\frac{2d\log(2K+2)}{n_l}} + BCL_M dc_{\mathcal{X}}\sqrt{\frac{\log(2/\delta)}{n}}.$$
 (43)

466 with probability at least  $1 - \delta$ .