
G²CN: Graph Gaussian Convolution Networks with Concentrated Graph Filters

Mingjie Li¹ Xiaojun Guo¹ Yifei Wang² Yisen Wang^{1 3} Zhouchen Lin^{1 3 4}

Abstract

Recently, linear GCNs have shown competitive performance against non-linear ones with less computation cost, and the key lies in their propagation layers. Spectral analysis has been widely adopted in designing and analyzing existing graph propagations. Nevertheless, we notice that existing spectral analysis fails to explain why existing graph propagations with the same global tendency, such as low-pass or high-pass, still yield very different results. Motivated by this situation, in this paper, we develop a new framework for spectral analysis from a detailed spectral analysis called concentration analysis. In particular, we propose three attributes: concentration centre, maximum response, and bandwidth for our analysis. Through a dissection of the limitations of existing graph propagations via the above analysis, we proposed a new kind of propagation layer, Graph Gaussian Convolution Networks (G²CN), in which the three properties are decoupled and the whole structure becomes more flexible and applicable on different kinds of graphs. Extensive experiments show that we can obtain state-of-the-art performance on various benchmark datasets with our proposed Graph Gaussian Convolution Networks.

1. Introduction

Graph neural networks have achieved excellent performance on various graph-related tasks on social relationship (Li & Goldwasser, 2019; Qiu et al., 2018; Tong et al., 2019), traffic (Bogaerts et al., 2020; Cui et al., 2019; Li et al., 2018b), recommendation (Ying et al., 2018; Feng et al., 2022), medical researches and others (Zhao et al., 2019; Rathi et al.,

2020; Jiang et al., 2021; Chen et al., 2022). GCNs mainly consist of the graph propagation part and feature transformation part. The graph propagation part first processes the graph signals to obtain the features with global or local graph properties. Then the signals are fed to a neural network with several or only one layer. Several works (Wu et al., 2019; Zhu & Koniusz, 2020) have shown that even a linear transformation layer can achieve better performance on many benchmark datasets. Therefore, graph propagation is the key part of a GCN model and attracts a lot of attention these days. Among the researches on graph propagation, graph spectral analysis has been widely adopted in designing and analyzing the existing graph propagation (Wu et al., 2019; Balcilar et al., 2020). It treats the graph propagations as graph filters on the spectrum of the graph Laplacian and classifies various graph propagation via the spectral response. For example, a graph filter that responds to amplify eigenvalues and diminish large eigenvalues will be classified as a low-pass graph filter.

However, the former spectral analysis mainly focuses on the global tendency. These analyses cannot completely explain why GCNs with different graph propagation will show different results even if their corresponding graph filter share the same low-pass or high-pass tendency. Motivated by this, we'd analyze the filters' pass band since it determines the output of the graph propagation. Since the filters' responses are concentrated on the passband, we called the analysis as **spectral concentration analysis** to distinguish from the previous study. And we propose three attributes called the **concentration attributes** for spectral analysis on the graph propagation contrast to the filter's global tendency: maximum response \mathcal{R} , bandwidth \mathcal{BW} and band concentration centre b in the following papers.

We analyze existing graph propagation with these attributes and find that their different performance can be attributed to their various concentration flexibility. Models with high concentration flexibility can apply different graph filters on various graphs, leading to better performance. With these findings, we also propose our Gaussian Graph Convolution, whose graph filters can be assumed as the combination of several Gaussian filters. Our Gaussian graph filter decouples the above attributes and can be easily adapted to perform as different kinds of filters to suit the graph desire. Then we implement our Gaussian graph filter in a linear graph model

¹Key Lab. of Machine Perception (MoE), School of Artificial Intelligence, Peking University. ²School of Mathematical Sciences, Peking University. ³Institute for Artificial Intelligence, Peking University. ⁴Pazhou Lab, Guangzhou, China. Correspondence to: Zhouchen Lin <zlin@pku.edu.cn>.

called Graph Gaussian Convolution Networks (G^2CN). The experiments show that our G^2CN can obtain state-of-art performance on various benchmark datasets with efficient computation cost. We summarize the contributions of our work as below:

1. We extend the former spectral analysis to our concentration spectral analysis with our proposed concentration attributes. Compared with the former study, our proposed one not only focuses on the graph filters global tendency. And we can explain the different performances of various graph propagation with the same global trend.
2. We proposed Gaussian Graph Convolution and its graph propagation formulation from the above analysis. Our proposed filter enjoys sufficient concentration flexibility and can apply to different graphs with better performance, especially the heterophily graphs.
3. We propose the Graph Gaussian Convolution Networks (G^2CN) with our Gaussian Graph Convolution, an efficient linear graph model that can achieve state-of-art performance on the node-classification tasks.

2. Related Works

Graph Neural Networks can be generally classified into spectral-based models (Bruna et al., 2013; Henaff et al., 2015; Li et al., 2018a), which use the eigendecomposition of the graph Laplacian, and spatial based model, which directly uses the graph Laplacian for graph propagation without the decomposition, which is widely used due to its efficiency. In this paper, we only discuss the latter propagation type. And we can divide the graph propagation for the current spatial-based GCNs from the spectral view into two kinds: GCNs with designed graph filters or learnable graph filters.

2.1. GCNs with Graph Propagation Constructed by designed Graph Filters

As for the graph propagation part in Graph Convolution Networks, most models utilize the designed graph filter based on their priors on the graph information they need. For example, GCNs (Kipf & Welling, 2017) can be regarded as a low-pass filter (Balcilar et al., 2020; Wu et al., 2019; Xu et al., 2019a). Furthermore, APPNP (Klicpera et al., 2019a) utilizes the Personalized PageRank and achieves a low-pass filter with different concentration properties compared to GCN as discussed in the following. Due to such limitations, these models show unstable performance on some real-world datasets. These models mainly focus on the local information on the graph and motivated many researchers to try to combine low-pass and high-pass filters (Zhu & Koniusz, 2020; Zhu et al., 2021) to balance the local and

global information for graph data.

2.2. GCNs with the Graph Propagation Constructed by learnable Filters

Without considering the priors on the possible graphs filters, GCNs with learnable filters try to approximate any filter function for different graph datasets using the generalization ability of the neural network, such as GPR-GNN (Chien et al., 2020), ChebNet (Defferrard et al., 2016), BernNet (He et al., 2021) and others (Bianchi et al., 2021). These models use the combination of a series of polynomial bases they choose to approximate an arbitrary filter function, like Chebyshev polynomials and Bernstein polynomials. For this count, these models can obtain better results compared with traditional GCNs on many graph datasets, especially the heterophily ones.

3. On the Spectral Concentration Analysis of Graph Propagation

In this section, we first propose three concentration attributes: maximum response, concentration centre and bandwidth for spectral analysis on graph propagation. Our spectral analysis with these attributes discusses the weaknesses and advantages of different graph propagation in the following. Our analysis can explain why models like GCN, APPNP and Heat Diffusion’s graph filters with the same spectral tendency still show different performance on node classification tasks.

3.1. Notations and Preliminaries on Graph Propagation

For a graph defined as $\mathcal{G} = (\mathcal{V}, \mathbf{A})$, where $\mathcal{V} = \{v_1, \dots, v_n\}$ denotes the vertex set of n nodes which contains a subset labeled nodes $\mathcal{V}_l \subset \mathbf{V}$ for training and unlabeled nodes $\mathcal{V}_u = \mathcal{V} \setminus \mathcal{V}_l$ to predict in the node classification task. $\mathbf{A} \in \mathbb{R}^{n \times n}$ is the adjacent matrix where a_{ij} denotes the edge weight between node v_i and v_j . The degree matrix $\mathbf{D} = \text{Diag}(d_1, \dots, d_n)$ of \mathbf{A} is a diagonal matrix with its i -th diagonal entry as $d_i = \sum_j a_{ij}$. Each node v_i is represented by a d -dimensional feature vector $\mathbf{x}_i \in \mathbb{R}^d$ and we use $\mathbf{X} \in \mathbb{R}^{n \times d} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$ to denote the feature matrix and use a one-hot vector $\mathbf{y}_i \in \{0, 1\}^C$ to denote the label of i -th node of C classes. In our paper, we only discuss the normalized graph Laplacian matrix defined as $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2} \in \mathcal{S}_+^n$ for convenience. The results can easily be extended to other types of Laplacian. Its eigendecomposition can be depicted as $\mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top$, where $\mathbf{\Lambda}$ is the diagonal matrix consists of \mathbf{L} ’s eigenvalues, and $\mathbf{U} \in \mathbb{R}^{n \times n}$ is made by the eigenvectors of the Laplacian matrix.

Then we define the graph Fourier transform of a graph signal \mathbf{x} as $\hat{\mathbf{x}} = \mathbf{U}^\top \mathbf{x}$ and its inverse is $\mathbf{x} = \mathbf{U} \hat{\mathbf{x}}$. Thus the graph

propagation for signal \mathbf{x} with kernel \mathbf{g} can be defined as:

$$\mathbf{H} = \mathbf{g} *_G \mathbf{x} = \mathbf{U} \left((\mathbf{U}^\top \mathbf{g}) \odot \mathbf{U} \mathbf{x} \right) = \mathbf{U} \hat{\mathbf{G}} \mathbf{U}^\top \mathbf{X}, \quad (1)$$

where $\hat{\mathbf{G}} = \text{diag}(\hat{g}_1, \dots, \hat{g}_2)$ denotes the spectral kernel coefficients. Exactly computing the graph convolution on a random kernel \mathbf{g} is complicated due to the eigendecomposition for graph's Laplacian consumes unaffordable computation complexity. For this count, as we illustrate in the following, current works with spatial convolution usually use the polynomial functions $g(\mathbf{L})$ to approximate different kernels:

$$\mathbf{H} = g(\mathbf{L})\mathbf{x} = \mathbf{U}g(\mathbf{\Lambda})\mathbf{U}^\top \mathbf{X}, \quad (2)$$

where $g(\mathbf{\Lambda}) = \text{Diag}(g(\lambda_1), \dots, g(\lambda_n))$ and we call g the **graph filter** function for above graph propagation Eqn. (2). Then GCNs with such graph propagation can be considered as filtering the graph signals by the corresponding eigenvalues for Laplacian. Those g preserve small λ can be regarded as low-pass filters while those who keep large λ are high-pass filters. Apart from the graph filters global tendency, we also analyzed the concentrated properties via the concentration attributes defined as follows.

3.2. Concentration Attributes on Graph Propagation

First, we'd like to define the graph filters maximum response and the concentration centre:

Definition 3.1. The maximum response \mathcal{R}_g for graph filter g is denoted as:

$$\mathcal{R}_g = \max_{\lambda \in [0, 2]} |g(\lambda)|, \quad (3)$$

since the eigenvalues of the graph Laplacian \mathbf{L} lies in $[0, 2]$.

And the concentration centre can be defined as:

Definition 3.2. Then we define b as the **centre for graph filter** g if,

$$g(b) = \mathcal{R}_g, \quad (4)$$

Then we give the definition for the graph filter's bandwidth.

Definition 3.3. The bandwidth \mathcal{BW}_g for a given graph filter $g(\lambda)$ can be defined as:

$$\mathcal{BW}_g = \int_0^2 \mathbb{I}(g(\lambda), \frac{\mathcal{R}}{\sqrt{2}}) d\lambda,$$

where $\mathbb{I}(g(\lambda), \frac{\mathcal{R}_g}{\sqrt{2}})$ is a indicator function:

$$\mathbb{I}(\gamma_1, \gamma_2) = \begin{cases} 1, & \text{if } \gamma_1 \geq \gamma_2 \\ 0, & \text{if } \gamma_1 \leq \gamma_2 \end{cases},$$

for $\gamma_1, \gamma_2 \in \mathbb{R}$.

The above attributes tell which spectrum the graph filter will amplify to diminish. And different graph filters enjoy unique concentration attributes and lead to different performances, as we discussed. The bandwidth and concentration centre should be flexible to ensure the proper graph propagation for a different dataset, while the maximum response should be bounded for stable training. Then we try to dissect existing graph propagations via our spectral concentration analysis.

3.3. Dissecting Existing Designed Graph Propagations

Graph Convolutional Networks (GCN). Firstly, we dissect GCN (Kipf & Welling, 2017) via our concentration spectral analysis. The graph propagation for GCN can be formulated as:

$$\mathbf{H}_{gcn} = (\mathbf{2I} - \mathbf{L})^K \mathbf{X} \quad (5)$$

where $K \in \mathbb{Z}^+$ is the propagation steps for SGC. Its graph filter can be formulated as $g_{gcn}(\lambda) = (2 - \lambda)^K$, which is also a low-pass filter, whose concentration centre is 0 with $\mathcal{R}_{gcn} = 2^K$. And its bandwidth can be formulated as:

$$\mathcal{BW}_{gcn} = 2 - 2^{1 - \frac{1}{2K}}. \quad (6)$$

One can see that the bandwidth of GCN and its latter works (Wu et al., 2019; Xu et al., 2018) are determined by its propagation steps. Therefore, the bandwidth for GCN's graph filter is limited and such defects restrict its performance since it cannot flexibly implement low-pass filters with different bandwidths as the graph desires.

Personalized PageRank (PPR). Then Klicpera et al. (2019a) proposed the personalized PageRank (PPR) graph propagation, trying to make the graph filter more flexible. Its graph propagation tries to approximate:

$$\mathbf{H}_{ppr} = (\mathbf{I} - (1 - \alpha)(\mathbf{I} - \mathbf{L}))^{-1} \mathbf{X}, \quad (7)$$

with $\alpha \in (0, 1)$ to ensure the inverse exists. Its graph filter is $g_{ppr} = \frac{1}{1 - (1 - \alpha)(1 - \lambda)}$ and its concentration centres are both fixed at $\lambda = 0$ with $\mathcal{R}_{ppr} = \frac{1}{\alpha}$ and also a low pass filter. The bandwidth for PPR can be formulated as:

$$\mathcal{BW}_{ppr} = \frac{(\sqrt{2} - 1)\alpha}{1 - \alpha}. \quad (8)$$

Since the bandwidth for PPR can continuously change with different α , models with PPR propagation are more flexible than original GCN and enjoy better performance. Whereas when $\mathbf{BW}_{ppr} \rightarrow 0$, $\alpha \rightarrow 0$ and its maximum response will explode $\mathcal{R}_{ppr} = \frac{1}{\alpha} \rightarrow \infty$. Therefore, we cannot choose a narrow PPR filter in practice, although the bandwidth for PPR can vary from 0 to ∞ in theory.

Frequency Adaptation GCN (FAGCN). The spectral filter for FAGCN (Bo et al., 2021) is the combination of:

$$\begin{aligned} h_{FAGCN_L}(\lambda) &= (1 - \lambda + \epsilon)^2, \\ h_{FAGCN_H}(\lambda) &= (\lambda - 1 + \epsilon)^2, \end{aligned} \quad (9)$$

with $\epsilon \in [0, 1]$. From the above equation, one can see the center of FAGCN’s base filters are fixed at 0 or 2 like GCN. Their bandwidth are $\mathcal{BW}_{FAGCN} = (2^{1/4} - 1)(1 + \epsilon) \in [0.19, 0.38]$ and maximum responses are $(1 + \epsilon)^2$ which are less flexible than ours.

Auto-Regressive Moving Average (ARMA). The spectral filter for ARMA (Bianchi et al., 2021) is the combination of:

$$h_{ARMA_K}(\mu) = \frac{b_K}{1 - a_K\mu}, \quad (10)$$

with $\mu = (\lambda_{\max} - \lambda_{\min})/2 - \lambda$, λ is eigenvalue of graph Laplacian, $|a_K| < 1$ and b_K . It is a low-pass filter when $a_K > 0$ with fixed center $\lambda = 0$ and it is a high-pass filter when $a_K < 0$ with fixed center $\lambda = 2$. In most cases, $\lambda_{\max} = 2$ and $\lambda_{\min} = 0$, then maximum response is $\frac{b_K}{1 - |a_K|}$ with bandwidth $\mathcal{BW}_{ARMA} = \frac{(\sqrt{2}-1)(1-|a_K|)}{|a_K|}$ which performs like APPNP but less flexible than ours.

Heat Kernels (HK). Apart from the above models, another widespread designed graph propagation is the heat kernel propagation (Klicpera et al., 2019b; Wang et al., 2021) which is inspired by the heat diffusion. Such propagation tries to approximate the following graph propagation:

$$\mathbf{H}_{HK} = e^{-T\mathbf{L}}\mathbf{X}, \quad (11)$$

where $T \in \mathbb{R}^+$ is the diffusion factor and its graph filter is an exponential function $g_{dgc}(\lambda) = e^{-T\lambda}$. Its maximum response is fixed to be 1 with fixed centre $\lambda = 0$. And the bandwidth for HK is:

$$\mathcal{BW}_{HK} = \frac{\log(\sqrt{2})}{T}. \quad (12)$$

The above formula shows that the bandwidth of DGC is flexible and can be easily controlled by T , and the fixed maximum response can also stabilize the training procedure. Thereby, as illustrated in (Klicpera et al., 2019b; Wang et al., 2021), heat kernels can show better performance against the above graph propagations on most datasets. However, the unflexible concentration centre for the above models also hinders their performance on many datasets since these low-pass models will permanently cut off the graph information corresponding to large eigenvalues of the graph Laplacian.

3.4. Dissecting Graph Propagation with Learnable Graph Propagations

To approximate the graph filters with different concentration centres, many graph models combine a series of polynomial bases with learnable weights as their graph propagation. Although these models can theoretically approximate every filter with proper weight for each basis, different properties for the basis will infect the difficulties of weight training and lead to different performances. Thereby, this section discusses the concentration attributes for a single filter basis.

Chebyshev Polynomials. ChebNet (Defferrard et al., 2016) which utilizes the Chebyshev Polynomials to approximate various graph filters, whose i -th basis can be formulated as:

$$\mathbf{C}^{(i)} = 2\mathbf{C}^{(2)}\mathbf{C}^{(s-1)} - \mathbf{C}^{s-2} \quad (13)$$

with $\mathbf{C}^{(0)} = 0$ and $\mathbf{C}^{(1)} = \frac{2\mathbf{L}}{\lambda_{\max}}$ and its propagation for i -th basis is $\mathbf{H}_{Cheb}^i = \mathbf{C}^{(i)}\mathbf{X}$. Although we cannot give the closed-form formula for the concentration centre and bandwidth for its i -th basis, its concentration centre and bandwidth are determined by its index i . We left the detailed calculation of these attributes for ChebNet filters in the Appendix E. The bases for ChebNet have more than one concentration centre lies in $[0, 2]$ and relies on the complicated weights to approximate the usually used filters like Comb Filters as experiment shows in (He et al., 2021).

Bernstein Polynomials. BernNet (He et al., 2021) which sum the Bernstein bases with weights for their network, the graph propagation for i -th basis for an order- K BernNet graph propagation can be formulated as follows:

$$\mathbf{H}_{Bern}^{K,i} = \frac{1}{2^K} \binom{K}{k} (2\mathbf{I} - \mathbf{L})^{K-k} \mathbf{L}^k, \quad (14)$$

whose corresponding graph filter can be simplified as $(2 - \lambda)^{K-i} \lambda^i$. Its concentration centre of with order K is $\frac{2i}{K}$, which cannot continuously change since $K \in \mathbb{Z}^+$. But when comparing with ChebNet, each Bernstein filters only have one concentration centre, which makes the weights’ training easier and leads to better performance.

The concentration centre and bandwidth for each Bernstein basis is fixed by its index i and order K as we show in the Appendix. When the objective filter is similar to their bases, the approximation is easy. And single concentration centre for each basis make it easier to approximate different filters compared with ChebNet. Since real-world graphs usually needs passbands of low or high frequency with various bandwidths, BernNet needs a complicated linear transformation layer to train a complicated weight and then they can obtain a complicated weights for bases during training. Thereby, the BernNet is not efficient and cannot easily adapt on different graphs. We summarized the analysis above as a Table in Appendix.

4. The proposed G²CN models.

4.1. Gaussian Graph Filters and their proposed

From the above analysis, one can see that models with better flexibility like Heat Kernel propagation and BernNet perform better than the others since they can easily approximate unique filters adapted to different graphs. However, as we stated above, these models are not flexible enough as the concentration center and bandwidth for Bernstein basis are unchangeable and Heat Kernel based propagation always

performs as a low-pass filter. Therefore, we'd like to construct a series of bases with better concentration flexibility and combine them to approximate unique filters for different graphs.

A proper graph filter should have flexible concentration attributes which means that its concentration centre and bandwidth can change smoothly and its maximum response is bounded, then it can easily adapted to various graphs. To build such a graph filter, we use the a series of Gaussian bases as our graph filter. We propose our Gaussian Graph Convolution inspired from the close form for graph heat diffusion. The graph propagation for a single Gaussian basis concentrated at b can be formulated as follows:

$$\mathbf{H}_G = e^{-T(\mathbf{L}-b\mathbf{I})^2} \mathbf{X}, \quad (15)$$

with $b \in [0, 2]$ and its maximum response 1. Then for its corresponding Gaussian filter, we can obtain its bandwidth as follows:

Proposition 4.1. *For the Gaussian filter denoted as:*

$$g_{G_{T,b}}(\lambda) = e^{-T(\lambda-b)^2},$$

with $b \in [0, 2]$ and T is the hyper-parameter related the bandwidth. It's bandwidth is:

$$\mathcal{BW}_{g_{T,b}} = \left\{ \begin{array}{ll} b, & b < \sqrt{\frac{\log\sqrt{2}}{T}} \\ 2\sqrt{\frac{\log\sqrt{2}}{T}}, & \sqrt{\frac{\log\sqrt{2}}{T}} \leq b \leq 2 - \sqrt{\frac{\log\sqrt{2}}{T}} \\ 2 - b, & b > 2 - \sqrt{\frac{\log\sqrt{2}}{T}} \end{array} \right\}$$

The proposition can be quickly proved by calculating the $g_{G_{T,b}}(\lambda) \leq \frac{1}{\sqrt{2}}$ with $\lambda \in [0, 2]$. From the above results, we can demonstrate the advantages for the Gaussian filter. Like heat kernel filters, Gaussian filters decoupled the three concentration attributes. For this count, we can flexibly choose different Gaussian filters and combine them to achieve our graph filter with the desired pass-band. We can choose T for different graph propagation to different approximated filters with various bandwidths and choose b to change the concentration centre smoothly. Empirical results show that such flexibility enables our G^2CN to perform as different graph filters and show satisfying results on real-world datasets.

Analysis via the function approximation perspective. As the following lemma shows, all square-integrable functions can approximate f with a series of Gaussian kernels.

Lemma 4.2. (Calcaterra & Boldt, 2008) *For any $f \in L^2\mathbb{R}$ and any $\epsilon > 0$ there exists $T > 0$ and $b > 0$ and $N \in \mathbb{N}$ and $\theta_i \in \mathbb{R}$ such that:*

$$\left\| f(x) - \sum_{i=0}^N \theta_i e^{-T(x-nb)^2} \right\| \leq \epsilon.$$

We can obtain the above conclusion based on the findings proposed in (Calcaterra & Boldt, 2008). Thereby, we can approximate the most practical filters with proper series Gaussian filter. The above results demonstrate the expressive power of our proposed Gaussian Graph Convolution.

4.2. Formulation of the Gaussian Filter in Graph and our G^2CN

To achieve the Gaussian filter in the graph propagation, we propose two formulations for the approximated graph propagation as listed below:

G^2CN -Taylor Firstly, we can obtain dissecting the Gaussian filter as graph propagation by K -order Taylor expansion:

$$\mathbf{H} = g_{G_{T,b}}(\mathbf{L})\mathbf{X} \approx \sum_{k=0}^K \frac{T^k}{k!} (\mathbf{L} - b\mathbf{I})^{2k} \mathbf{X}, \quad (16)$$

where K here denotes the propagation times, which only influences the accuracy of the approximation.

G^2CN -Euler Since the equivalent graph filter $g_{G_{T,b}}$ can also be regarded as the close form solution for the following differential system at $t = 1$:

$$\frac{\partial \mathbf{X}(t)}{\partial t} = -T(\mathbf{L} - b\mathbf{I})^2 \mathbf{X}(t), \quad (17)$$

with $\mathbf{X}(0) = \mathbf{X}$. Then we can use the forward Euler method to solve the graph propagation with the step size equal to $1/K$.

$$\mathbf{X}\left(\frac{i}{K}\right) = \mathbf{X}\left(\frac{i-1}{K}\right) - \frac{T}{K}(\mathbf{L} - b\mathbf{I})^2 \mathbf{X}\left(\frac{i-1}{K}\right),$$

with $i = 0, \dots, K$ to obtain the final state $\mathbf{H} \approx \mathbf{X}(1)$. We only need to store the state i/K after the i -th iteration with such a scheme. And since the matrix-matrix multiplication is too expansive, we iteratively use matrix-vector multiplication in practice.

Using the above methods, we can obtain the equivalent graph propagation for the graph filter concentrated at b . Since one Gaussian filter can only respond to the eigenvalues around the basis, we can combine features obtained by Gaussian filters with different concentration centre as follow:

$$\mathbf{H} = \sum_{i=0}^N \theta_i g_{G_{T_i,b_i}}(\mathbf{L})\mathbf{X}, \quad (18)$$

where N denotes the bases number we used, and θ_i are learnable weights balancing Gaussian filters with different concentration centres and bandwidth.

The whole structure of G^2CN . After combining our Gaussian Graph Convolution, we can obtain the graph information filtered by the graph spectral filter we desired. As

shown in (He et al., 2021), the node classification mainly needs graph information centred on the high-frequency or low-frequency part in practice, we only adopt two Gaussian filters for our G^2CN with fixed centre $b = 0$ and $b = 2$ and different T for the following equations. Since the high flexibility makes the weights training much easier, we only use a linear layer for classification after the propagation for computational efficiency. Then we finally get our G^2CN .

4.3. Numerical Analysis on above two G^2CN type

In this section, we'd like to compare to type of G^2CN from the perspective of their approximated error from their closed-form solution which can be defined as:

$$e_{\mathcal{M}}^{T,K,b} = \|\mathbf{H}_{\mathcal{M}} - e^{-T(b\mathbf{I}-\mathbf{L})^2} \mathbf{X}\|, \quad (19)$$

with \mathcal{M} is "Eu" or "Ta" to denote G^2CN -Euler and G^2CN -Taylor and $\|\cdot\|$ denotes the operator norm. Then the numerical error for the Taylor method is:

Theorem 4.3. *For a single G^2CN -Taylor's graph filter with order K with width factor T and bias b , the approximated error $e_{Ta}^{T,K,b}$ compared with its closed-form $e^{-T(b\mathbf{I}-\mathbf{L})^2} \mathbf{X}$ is:*

$$e_{Ta}^{T,K,b} \leq \frac{(T\|b\mathbf{I} - \mathbf{L}\|^2)^{K+1}}{(K+1)!} \|\mathbf{X}\|. \quad (20)$$

And the numerical error for Euler method is:

Theorem 4.4. *For a single G^2CN -Euler's graph filter with order K with width factor T and bias b , if $2T < K$ the approximated error $e_{Eu}^{T,K,b}$ compared with its closed-form $e^{-T(b\mathbf{I}-\mathbf{L})^2} \mathbf{X}$ can be bounded by:*

$$e_{Eu}^{T,K,b} \leq \frac{T^2\|b\mathbf{I} - \mathbf{L}\|_2^4}{K} \|\mathbf{X}\|. \quad (21)$$

From the above results, one can see that the increase of K can reduce the numerical error for the Euler method. However, when $K < T\|b\mathbf{I} - \mathbf{L}\|_2^2$ the error for the Taylor expansion will increase first as K increases and then decrease when K is larger enough and will return a smaller error than the Euler method when K is large enough. Therefore, we recommend the G^2CN Euler when K is small due to the limited computation complexity. And as the results demonstrate, we only recommend the G^2CN -Taylor when K is significantly large compared with T , G^2CN -Taylor will return a better approximation in this scenario.

5. Experiments

5.1. Homophily and Heterophily Datasets

The graph node classification datasets can be divided into homophily graphs and heterophily graphs with the homophily principle (McPherson et al., 2001). Homophily principle

states that contact between similar nodes occurs at a higher rate than among dis-similar nodes. In recent work, Pei et al. (2020) designs an index denoted by \mathcal{H} to measure the homophily in a graph:

$$\mathcal{H}(G) = \frac{1}{|V|} \sum_{v \in V} \frac{\text{Number of neighbors of } v \in V \text{ that have the same label as } v}{\text{Number of neighbors of } v}.$$

Lower \mathcal{H} indicates a strong heterophilic graph, where nodes with distinct labels are more likely to link together. As ? states, the homophilic graphs need low-pass filters for a satisfactory performance while heterophilic graphs need other passbands in graph filters for generalization.

In our work, we conduct experiments on five widely used benchmark homophylic graphs as others (Wang et al., 2021) to demonstrate our Gaussian filter achieve the comparable state-of-the-art performance compared with other low-pass graph filters, including the standard citation graphs Cora, CiteSeer and PubMed (Sen et al., 2008; Yang et al., 2016) and the Amazon co-purchase graphs Computers and Photo (McAuley et al., 2015; Shchur et al., 2018).

Then we finish the experiments for five benchmark heterophilic graphs used in (Pei et al., 2020), including Wikipedia graphs Chameleon and Squirrel (Rožemberczki et al., 2021), the Actor co-occurrence graph and the WebKB graphs Texas and Cornell (Pei et al., 2020). Experiments on heteophilic graphs demonstrate the superiority of our G^2CN models compared with other linear or non-linear graph models due to the flexibility of the Gaussian Graph Filters we used. More rigorously, we also conduct the comparison on a large scale node classification dataset, OGB-Arxiv (?). An overview summary of characteristics of the datasets is given in Table 1.

5.2. Experimental Settings

Firstly, we experiment on the semi-supervised node classification task for standard citation networks with the standard data split as in (Wang et al., 2021; Kipf & Welling, 2017; Veličković et al., 2018; Xu et al., 2019a;b; Poli et al., 2019) with the same setting as DGC. Furthermore, we also evaluate the performance of the homophilic Amazon co-purchase graph Computers and Photo for fully supervised node classification tasks with the same setting as BernNet, who randomly split the node sets into train/validation/test set ratio 60%/20%/20%.

Then for the heterophilic graphs, we use the same setting as BernNet and finish the experiments on the fully-supervised node classification tasks for these graphs. We compare G^2CN with three baseline models: GCN (Kipf & Welling, 2017), GAT (Veličković et al., 2018), APPNP (Klicpera et al., 2019a), BernNet (He et al., 2021), SGC (Wu et al., 2019), DGC (Wang et al., 2021). We directly tuned the

Table 1. Datasets statistics

Dataset	Cora	Citeseer	Pubmed	Computers	Photo	Chameleon	Squirrel	Actor	Texas	Cornell	OGB-Arxiv
#Nodes	2708	3327	19717	13752	7650	2277	5201	7600	183	183	169343
#Edges	5278	4552	44324	245861	119081	36101	217073	26659	295	309	1166243
#Features	1433	3703	500	767	745	2325	2089	932	1703	1703	128
#Classes	7	6	3	10	8	6	5	5	5	5	40
\mathcal{H}	0.83	0.71	0.79	0.47	0.59	0.25	0.22	0.24	0.06	0.11	

hyper-parameters for non-linear models follows the same steps as He et al. (2021) stated and the channel numbers hidden layers for non-linear models are 64. As for linear GCNs, we performed a hyper-parameter search for linear models’ hyper-parameter (T and weight decay) on the validation set as used in DGC for 1000 trials. We choose $K = 200$ for DGC and $K = 100$ for our G^2CN to ensure the same propagation steps with graph Laplacian for fairness.

We use Adam optimizer (Kingma & Ba, 2014) for training with 100 epochs for semi-supervised tasks and 1000 epochs for fully-supervised tasks as BernNet for fair with lr=0.5 for our model and report results averaged over ten random runs. And we use G^2CN -Euler with $2T < K$ in the following since its performance is more stable when T is large.

5.3. Experimental Results on Homophily Datasets.

The test results for homophily datasets are shown in Table 2. The results show that our G^2CN can show comparable performance on the semi-supervised citation graphs. Furthermore, one can see that our G^2CN enjoys an evident advantage on the fully-supervised graph tasks for large graphs. The empirical results demonstrate that our G^2CN can return better graph signals output by our graph filters on homophilic graphs, which mainly depend on the graph’s low-frequent composition.

We notice that BernNet perform worse than our G^2CN on the homophilic graphs. The proper passband for these graphs shown in Table 3 can explain such phenomenon.

The table shows that The bandwidth for each bernstein basis in BernNet is not flexible. For example, the pass band for the last three filters of a BernNet ($K = 10$) is $[1.3, 1.92]$, $[1.61, 1.92]$, $[1.93, 2.0]$. Due to such reason, BernNet needs to rely on complicated weighted parameter to approximate the narrow high-pass filter stated in G^2CN and it will make the learning task much harder. Thereby, it cannot utilize proper high-frequency graph information as ours which leads to worse results than ours.

5.4. Experimental Results on Heterophily Datasets

Apart from the homophilic graphs, we also conduct experiments on the heterophilic graphs. The results are shown

Table 2. Test accuracy (%) comparison on common used citation graphs with semi-supervised scheme (above) and Amazon co-purchase graphs with fully-supervised scheme (blow). The results are averaged over 10 runs.

Type	Dataset	Cora	CiteSeer	Pubmed
Nonlinear	GCN	81.5	70.3	79.0
	GAT	83.0	72.5	79.0
	APPNP	83.3	71.8	80.1
Linear	SGC	81.1	71.9	78.9
	DGC	83.3	73.3	80.3
	G^2CN	82.7	73.8	80.4

Type	Dataset	Computers	Photo
Nonlinear	GCN	83.32 ± 0.33	88.26 ± 0.73
	GAT	83.32 ± 0.39	90.94 ± 0.68
	APPNP	85.32 ± 0.32	88.51 ± 0.31
	BernNet	87.64 ± 0.44	94.50 ± 0.40
Linear	SGC	84.74 ± 0.27	86.48 ± 0.41
	DGC	89.54 ± 0.64	94.89 ± 0.50
	G^2CN	91.46 ± 0.35	95.29 ± 0.48

Table 3. Bandwidth for DGC and G^2CN (Bandwidth for two Gaussian filters with $b = 0, 2$) Computers and Photo.

Datasets	DGC	G^2CN
Computers	0.21	(0.35, 0.15)
Photo	0.23	(0.37, 0.12)

in Table 4. Compared with the homophilic graphs, the advantages for our G^2CN over the low-pass models is more evident since the heterophilic graphs need the graph signals with different frequency for their tasks. Compared with BernNet, our G^2CN can still show better performance on these models except Actor. The above results can also demonstrate the superiority of our proposed model.

Understandings via our Concentration Analysis From the experimental results, one can see that our G^2CN or BernNet shows significant performance on Chameleon and

Table 4. Test accuracy (%) comparison on heterophily datasets. Reported results are averaged over 10 runs.

Dataset	GCN	GAT	APPNP	FAGCN	ARMA	DGC	BernNet	G ² CN
Chameleon	42.98 ± 1.41	47.62 ± 0.83	43.07 ± 1.55	63.3 ± 1.41	55.1 ± 1.32	67.08 ± 1.58	54.11 ± 1.73	73.61 ± 0.89
Squirrel	28.41 ± 0.57	27.33 ± 0.81	31.71 ± 0.47	39.7 ± 0.67	35.5 ± 0.81	50.51 ± 0.81	41.35 ± 0.81	66.91 ± 1.13
Actor	33.23 ± 1.16	33.93 ± 2.47	39.66 ± 0.55	40.11 ± 0.77	40.79 ± 0.89	40.35 ± 0.73	41.79 ± 1.01	41.44 ± 0.76
Texas	77.38 ± 3.28	80.82 ± 2.13	90.98 ± 1.64	96.5 ± 0.47	92.3 ± 0.66	94.74 ± 0.33	96.22 ± 0.79	96.72 ± 0.73
Cornell	65.90 ± 4.43	78.21 ± 2.95	91.81 ± 1.96	93.3 ± 1.21	93.4 ± 1.13	92.45 ± 1.21	92.29 ± 2.74	94.11 ± 1.81

Squirrel, while on Cornell, Actor and Texas, the difference is not significant compared with low-pass graph models. To explain such a phenomenon, we calculate the bandwidth of DGC and G²CN in Table 5.

Table 5. Bandwidth for DGC and G²CN (Bandwidth for two Gaussian filters with $b = 0, 2$) for heterophilic Squirrel and Cornell.

Datasets	DGC	G ² CN
Texas	2	(1.13, 0.16)
Actor	2	(2.0, 0.27)
Cornell	2	(2.0, 0.08)
Squirrel	0.09	(0.12, 0.37)
Chameleon	0.07	(0.11, 0.27)

The table shows that the Texas, Actor and Cornell mainly need to apply an all-pass filter. For this count, APPNP and DGC can also perform better than GCN and GAT since the bandwidth for APPNP and DGC is more flexible than GCN and GAT. As for Chameleon and Squirrel datasets, the best bandwidth adapted for DGC and G²CN demonstrates that they need a comb graph filter with narrow pass-bands. Therefore, graph models cannot perform well since they lack high-frequent graph signals. As for BernNet, as illustrated in the former, BernNet needs complicated weight to approximate graph filters, making its training procedure much harder. Empirical results demonstrate the effectiveness of our G²CN and the rationality of our concentration spectral analysis since we can use it to understand different models’ performance on various datasets.

5.5. Performance on Large Scale Datasets

Apart from the above datasets, we also finish experiments for our G²CN on OGB-Arxiv. Since the OGB-Arxiv dataset is large and hard, we replace the linear layer with 2-layer MLP after the output of our Gaussian graph filters. The results are listed in Table 6. From the table, one can see that our G²CN+MLP can outperform the DGC+MLP and SGC on the large scale dataset. The results on the large graph OGB-Arxiv demonstrate the effectiveness of our G²CN.

Table 6. Test accuracy (%) comparison on the large scale dataset OGB-Arxiv. OOM: out of memory.

Method	Accuracy
GCN	OOM
BernNet	OOM
SGC	68.9
DGC+MLP	71.2
G ² CN+MLP	71.6

5.6. Computation Time for Different Models

We list the computation time for different models trained on large graphs for 1000 epochs in Table 7. The models are evaluated on a single RTX-3070 and Intel I5-10400 (2.90GHz). As a linear GCN, DGC and G²CN is much faster than non-linear GCNs like GCN and BernNet. Our G²CN is only slightly slower compared with DGC. Since our G²CN needs an additional weighted summation for the outputs of two filters, our model is slightly slower than DGC.

Table 7. Comparison of total training time for different models trained for 1000 epochs. The figures in the brackets are preprocessing times (s) and those outside the brackets are total times (s). The K for DGC and G²CN are 200 and 100, and G²CN _{n} denotes our G²CN with n Gaussian graph filters.

Model	Pubmed	Computers	Photo
BernNet	59.33s(-)	107.1s(-)	53.33s(-)
DGC	1.04s(19ms)	2.01s(20ms)	1.40s(19ms)
G ² CN	1.24s(18ms)	2.27s(21ms)	1.56s(19ms)

5.7. Ablation Studies on Spectral Response Curve For Selected Dataset

In this section, we draw the approximated spectral response curve of our G²CN’s Gaussian graph filter for selected datasets as shown in Figure 1. The curve depicts whether our Gaussian graph filter amplifies or diminishes the signals corresponding to eigenvalue λ . From the above figures, one

can see that same as our former analysis. Pubmed needs a low-pass filter while Computers and Squirrel needs comb filters with narrow bands and filters. And Squirrel mostly rely on its graph signals with high frequency for classification which is the same as Chameleon listed in the Appendix. Since our G^2CN can easily approximate such filter, our outperforms other datasets with significant advantages. Cornell is more likely to be an all-pass filter. We left the curves for other graphs in the Appendix F.

62006153), Project 2020BD006 supported by PKU-Baidu Fund, and Open Research Projects of Zhejiang Lab (No. 2022RC0AB05).

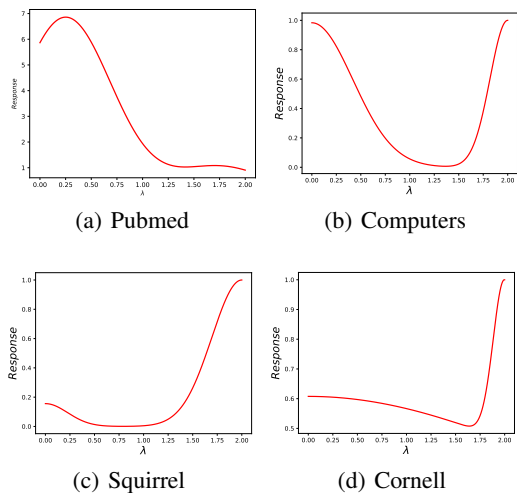


Figure 1. G^2CN 's equivalent graph filters for selected real-world datasets.

6. Conclusions

In this work, we first proposed three attributes called concentration attributes to analyze the spectral properties for different graph models called concentration spectral analysis. Our analysis can explain why various graph models will show different results even if their global tendency is the same. Furthermore, we analyzed the graph filters' weaknesses and advantages for different models. We notice that these models' concentration attributes are not flexible, making it hard to apply proper graph filters for different graphs. We then proposed our Gaussian graph filters and our G^2CN using this filter. The concentration attributes for our G^2CN is much more flexible and can easily apply proper filters for different graphs. Empirical results also demonstrate the superiority of our model.

ACKNOWLEDGMENTS

Zhouchen Lin is supported by the NSF China (No. 61731018), NSFC Tianyuan Fund for Mathematics (No. 12026606), and Project 2020BD006 supported by PKU-Baidu Fund.

Yisen Wang is partially supported by the NSF China (No.

References

- Balcilar, M., Renton, G., Héroux, P., Gaüzère, B., Adam, S., and Honeine, P. Analyzing the expressive power of graph neural networks in a spectral perspective. In *ICLR*, 2020.
- Bianchi, F. M., Grattarola, D., Livi, L., and Alippi, C. Graph neural networks with convolutional arma filters. In *TPAMI*, 2021.
- Bo, D., Wang, X., Shi, C., and Shen, H. Beyond low-frequency information in graph convolutional networks. In *arXiv preprint arXiv:2101.00797*, 2021.
- Bogaerts, T., Masegosa, A. D., Angarita-Zapata, J. S., Onieva, E., and Hellinckx, P. A graph cnn-lstm neural network for short and long-term traffic forecasting based on trajectory data. In *Elsevier*, 2020.
- Bruna, J., Zaremba, W., Szlam, A., and LeCun, Y. Spectral networks and locally connected networks on graphs. In *arXiv preprint arXiv:1312.6203*, 2013.
- Calcaterra, C. and Boldt, A. Approximating with gaussians. In *arXiv preprint arXiv:0805.3795*, 2008.
- Chen, Q., Wang, Y., Wang, Y., Yang, j., and Lin, z. Optimization-induced graph implicit nonlinear diffusion. In *ICML*, 2022.
- Chien, E., Peng, J., Li, P., and Milenkovic, O. Adaptive universal generalized pagerank graph neural network. In *ICLR*, 2020.
- Cui, Z., Henrickson, K., Ke, R., Pu, Z., and Wang, Y. Traffic graph convolutional recurrent neural network: A deep learning framework for network-scale traffic learning and forecasting. In *TITS*, 2019.
- Defferrard, M., Bresson, X., and Vandergheynst, P. Convolutional neural networks on graphs with fast localized spectral filtering. In *NeurIPS*, 2016.
- Feng, L., Cai, Y., Wei, E., and Li, J. Graph neural networks with global noise filtering for session-based recommendation. In *Neurocomputing*, 2022.
- He, M., Wei, Z., Huang, Z., and Xu, H. Bernnet: Learning arbitrary graph spectral filters via bernstein approximation. In *NeurIPS*, 2021.
- Henaff, M., Bruna, J., and LeCun, Y. Deep convolutional networks on graph-structured data. In *arXiv preprint arXiv:1506.05163*, 2015.
- Jiang, D., Wu, Z., Hsieh, C.-Y., Chen, G., Liao, B., Wang, Z., Shen, C., Cao, D., Wu, J., and Hou, T. Could graph neural networks learn better molecular representation for drug discovery? a comparison study of descriptor-based and graph-based models. In *booktitle of cheminformatics*, 2021.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *arXiv preprint arXiv:1412.6980*, 2014.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- Klicpera, J., Bojchevski, A., and Günnemann, S. Predict then propagate: Graph neural networks meet personalized pagerank. In *ICLR*, 2019a.
- Klicpera, J., Weißenberger, S., and Günnemann, S. Diffusion improves graph learning. In *NeurIPS*, 2019b.
- Li, C. and Goldwasser, D. Encoding social information with graph convolutional networks for political perspective detection in news media. In *ACL*, 2019.
- Li, R., Wang, S., Zhu, F., and Huang, J. Adaptive graph convolutional neural networks. In *AAAI*, 2018a.
- Li, Y., Yu, R., Shahabi, C., and Liu, Y. Diffusion convolutional recurrent neural network: Data-driven traffic forecasting. In *ICLR*, 2018b.
- McAuley, J., Targett, C., Shi, Q., and Van Den Hengel, A. Image-based recommendations on styles and substitutes. In *SIGIR*, 2015.
- McPherson, M., Smith-Lovin, L., and Cook, J. M. Birds of a feather: Homophily in social networks. In *Annual review of sociology*, 2001.
- Pei, H., Wei, B., Chang, K. C.-C., Lei, Y., and Yang, B. Geom-gcn: Geometric graph convolutional networks. In *arXiv preprint arXiv:2002.05287*, 2020.
- Poli, M., Massaroli, S., Park, J., Yamashita, A., Asama, H., and Park, J. Graph neural ordinary differential equations. In *NeurIPS*, 2019.
- Qiu, J., Tang, J., Ma, H., Dong, Y., Wang, K., and Tang, J. Deepinf: Social influence prediction with deep learning. In *SIGKDD*, 2018.
- Rathi, P. C., Ludlow, R. F., and Verdonk, M. L. Practical high-quality electrostatic potential surfaces for drug discovery using a graph-convolutional deep neural network. In *booktitle of Medicinal Chemistry*, 2020.
- Rozemberczki, B., Allen, C., and Sarkar, R. Multi-scale attributed node embedding. In *JCN*, 2021.
- Sen, P., Namata, G., Bilgic, M., Getoor, L., Galligher, B., and Eliassi-Rad, T. Collective classification in network data. In *AI magazine*, 2008.

- Shchur, O., Mumme, M., Bojchevski, A., and Günnemann, S. Pitfalls of graph neural network evaluation. In *arXiv preprint arXiv:1811.05868*, 2018.
- Tong, P., Zhang, Q., and Yao, J. Leveraging domain context for question answering over knowledge graph. In *Springer*, 2019.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., and Bengio, Y. Graph attention networks. In *ICLR*, 2018.
- Wang, Y., Wang, Y., Yang, J., and Lin, Z. Dissecting the diffusion process in linear graph convolutional networks. In *NeurIPS*, 2021.
- Wu, F., Souza, A., Zhang, T., Fifty, C., Yu, T., and Weinberger, K. Simplifying graph convolutional networks. In *ICML*, 2019.
- Xu, B., Shen, H., Cao, Q., Cen, K., and Cheng, X. Graph convolutional networks using heat kernel for semi-supervised learning. In *IJCAI*, 2019a.
- Xu, B., Shen, H., Cao, Q., Qiu, Y., and Cheng, X. Graph wavelet neural network. In *ICLR*, 2019b.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *arXiv preprint arXiv:1810.00826*, 2018.
- Yang, Z., Cohen, W., and Salakhudinov, R. Revisiting semi-supervised learning with graph embeddings. In *ICML*, 2016.
- Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W. L., and Leskovec, J. Graph convolutional neural networks for web-scale recommender systems. In *SIGKDD*, 2018.
- Zhao, L., Peng, X., Tian, Y., Kapadia, M., and Metaxas, D. N. Semantic graph convolutional networks for 3d human pose regression. In *CVPR*, 2019.
- Zhu, H. and Koniusz, P. Simple spectral graph convolution. In *ICLR*, 2020.
- Zhu, M., Wang, X., Shi, C., Ji, H., and Cui, P. Interpreting and unifying graph neural networks with an optimization framework. In *WebConf*, 2021.

A. Proof for Theorem 4.3

The proof for Theorem 4.3 is as follows:

Proof. From the definition, the approximated error for G²CN-Taylor is:

$$e_{Ta}^{T,K,b} = \left\| \mathbf{H}_{Ta} - e^{T(b\mathbf{I}-\mathbf{L})^2} \mathbf{X} \right\|_2 \quad (22)$$

$$= \left\| \sum_{i=0}^K \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{i!} \mathbf{X} - \sum_{i=0}^{\infty} \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{i!} \mathbf{X} \right\|_2 \quad (23)$$

$$= \left\| \sum_{i=K+1}^{\infty} \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{i!} \mathbf{X} \right\|_2 \quad (24)$$

$$= \left\| (-T(b\mathbf{I}-\mathbf{L})^2)^{K+1} \sum_{i=0}^{\infty} \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{(i+K+1)!} \mathbf{X} \right\|_2 \quad (25)$$

$$= \left\| (-T(b\mathbf{I}-\mathbf{L})^2)^{K+1} \right\|_2 \left\| \sum_{i=0}^{\infty} \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{(i+K+1)!} \right\|_2 \|\mathbf{X}\|_2 \quad (26)$$

$$\leq \left\| \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^{K+1}}{(K+1)!} \right\|_2 \left\| \sum_{i=0}^{\infty} \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^i}{i!} \right\|_2 \|\mathbf{X}\|_2 \quad (27)$$

$$= \left\| \frac{(-T(b\mathbf{I}-\mathbf{L})^2)^{K+1}}{(K+1)!} \right\|_2 \left\| e^{-T(b\mathbf{I}-\mathbf{L})^2} \right\|_2 \|\mathbf{X}\|_2 \quad (28)$$

$$\leq \frac{(T\|(b\mathbf{I}-\mathbf{L})\|^2)^{K+1}}{(K+1)!} \|\mathbf{X}\|_2 \quad (29)$$

Then Theorem 4.3 is proved. \square

B. Proof for Theorem 4.4

The proof for Theorem 4.4 is as follows:

Proof. The Euler Scheme is to solve the following differential equation at $t = 1$:

$$\begin{cases} \frac{d\mathbf{X}_t}{dt} = -T(b\mathbf{I}-\mathbf{L})^2 \mathbf{X}_t, \\ \mathbf{X}_0 = \mathbf{X}, \end{cases} \quad (30)$$

with step size $1/K$ for K steps. Then consider a general Euler forward scheme for our problem:

$$\hat{\mathbf{X}}(k+1) = \hat{\mathbf{X}}(k) - h\mathbf{L}_{b,T}^2 \hat{\mathbf{X}}(k), \quad k = 0, \dots, K-1, \quad \mathbf{X}^{(0)} = \mathbf{X}, \quad (31)$$

with $\hat{\mathbf{X}}(k)$ denotes approximated $\mathbf{X}(k/K)$ by the forward Euler, $h = \frac{1}{K}$ and $\mathbf{L}_{b,T} = \sqrt{T}(b\mathbf{I}-\mathbf{L})$. Then we denote the error at step k as:

$$\mathbf{e}_k = \mathbf{X}^{(k)} - \hat{\mathbf{X}}(k), \quad (32)$$

then the approximated error is $e_{Eu}^{T,K,b} = \mathbf{e}_K$, and the truncation error of the Euler forward finite difference (Eqn (31)) at step k as:

$$\mathbf{T}^{(k)} = \frac{\mathbf{X}^{(k+1)} - \mathbf{X}^{(k)}}{h} + \mathbf{L}_{b,T}^2 \mathbf{X}^{(k)}, \quad (33)$$

which can be reformulated as follows:

$$\mathbf{X}^{(k+1)} = \mathbf{X}^{(k)} + h(\mathbf{T}^{(k)} - \mathbf{L}_{b,T}^2 \mathbf{X}^{(k)}). \quad (34)$$

Then subtract Eqn 34 by Eqn 31, we can obtain:

$$\mathbf{e}^{(k+1)} = (\mathbf{I} - h\mathbf{L}_{b,T}^2) \mathbf{e}^{(k)} + h\mathbf{T}^{(k)}, \quad (35)$$

whose norm can be upper bounded as:

$$\left\| \mathbf{e}^{(k+1)} \right\|_2 \leq \left\| \mathbf{I} - h\mathbf{L}_{b,T}^2 \right\|_2 \left\| \mathbf{e}^{(k)} \right\|_2 + h \left\| \mathbf{T}^{(k)} \right\|. \quad (36)$$

Since $\mathbf{L}_{b,T} = (b\mathbf{I} - \mathbf{L})$ with $b \in [0, 2]$ and the eigenvalue of the symmetric matrix \mathbf{L} also lies in $[0, 2]$, Then the eigenvalue of $\mathbf{L}_{b,T}^2$ lies in $[0, 4T]$ and eigenvalue of $h\mathbf{L}_{b,T}^2$ lies in $[0, 2]$ since $h = T/K < 0.5$. Thereby, the spectral norm of $\mathbf{I} - h\mathbf{L}_{b,T}^2$ is upper bounded by 1. Then with $M = \max_{0 \leq k \leq K-1} \left\| \mathbf{T}^{(k)} \right\|_2$ be the upper bound on the truncation error, we have:

$$\left\| \mathbf{e}^{k+1} \right\|_2 \leq \left\| \mathbf{e}^{(k)} \right\|_2 + hM. \quad (37)$$

Since $\mathbf{e}^{(0)} = \mathbf{X}^{(0)} - \hat{\mathbf{X}}^{(0)} = 0$, we have the following formula by induction:

$$\left\| \mathbf{e}^{(K)} \right\|_2 \leq KhM = M. \quad (38)$$

Then since $\frac{d\mathbf{X}^{(k)}}{dt} = -L\mathbf{X}^{(k)}$ and applying Taylor's theorem, there exists $\delta \in [kh, (k+1)h]$ such that the truncation error $\mathbf{T}^{(k)}$ in Eqn 33 follows:

$$\mathbf{T}^{(k)} = \frac{h}{2} \mathbf{L}_{b,T}^4 \mathbf{X}_\delta. \quad (39)$$

Furthermore, since

$$\left\| \mathbf{X}_\delta \right\|_2 = \left\| e^{-T\delta \mathbf{L}_b^2} \mathbf{X}_0 \right\| \leq \left\| \mathbf{X} \right\|_2. \quad (40)$$

Thereby, we can bound the truncated error:

$$\left\| \mathbf{T}^{(k)} \right\|_2 = \frac{h}{2} \left\| \mathbf{L}_b \right\|^4 \left\| \mathbf{X}_\delta \right\|_2 \leq \frac{h}{2} \left\| \mathbf{L}_b \right\|^4 \left\| \mathbf{X} \right\|_2, \quad (41)$$

Finally, we have

$$\left\| \mathbf{e}_{Eu}^{T,K,b} \right\|_2 = \left\| \mathbf{e}^{(K)} \right\|_2 \leq M \leq \frac{T^2}{2K} \left\| b\mathbf{I} - \mathbf{L} \right\|^4 \left\| \mathbf{X} \right\|_2. \quad (42)$$

Then we complete the proof. \square

C. Bandwidth and Concentration Centers for Graph filters approximated by Polynomial Bases.

We list the bandwidth of each bernstein basis for BernNet with $K = 10$ listed in Table 8.

Table 8. Concentration Attributes for Bernstein Basis.

Index	0	1	2	3	4	5	6	7	8	9	10
Concentration Center	0.0	0.2	0.4	0.6	0.8	1.0	1.2	1.4	1.6	1.8	2.0
PassBand	[0.0, 0.069]	[0.079, 0.390]	[0.219, 0.635]	[0.381, 0.855]	[0.555, 1.062]	[0.741, 1.258]	[0.938, 1.445]	[1.145, 1.619]	[1.365, 1.781]	[1.610, 1.921]	[1.931, 2.0]
Maximum Response	1	0.387	0.301	0.267	0.251	0.246	0.251	0.267	0.301	0.387	1

D. Ablation Study on the Propagation Stpes K .

We also finish the experiments on Photo with different propagation times K for our G²CN Euler to explore K 's influence. The results are shown in Figure 2. Same to our analysis stated in Theorem 4.3, as K increases, the approximated error will decrease and the performance will increase and finally converge.

E. Summary for our concentration analysis on different graph propagation.

We summarize the concentration attributes for different graph propagation as follows:

F. Response Curve for Other Graphs

We list the response curve for the rest Graphs as follows:

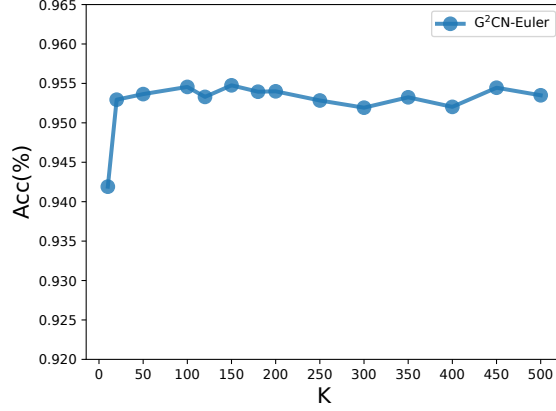

 Figure 2. Test Accuracy with respect to different propagation times for our G^2CN on Photo.

Table 9. Concentration attributes for different graph propagation.

	Graph Propagation Kernel (Or Basis)	\mathcal{R}	b	\mathcal{BW}
GCN	$(2\mathbf{I} - \mathbf{L})^K$	2^K	0	$2 - 2^{1 - \frac{1}{2K}}$
PPR	$(\mathbf{I} - (1 - \alpha)(1 - \mathbf{L}))^{-1}$	$\frac{1}{\alpha}$	0	$\frac{(\sqrt{2}-1)\alpha}{1-\alpha}$
ARMA	$\frac{b_K}{1 - a_K \mu}$	$\frac{b_K}{1 - a_K }$	0 or 2	$\frac{(\sqrt{2}-1)(1 - a_K)}{ a_K }$
FAGCN	$(1 - \lambda + \epsilon)^2$ $(\lambda - 1 + \epsilon)^2$	$(1 + \epsilon)^2$	0 or 2	$(2^{1/4} - 1)(1 + \epsilon) \in [0.19, 0.38]$
Heat Kernel	$e^{-T\mathbf{L}}$	1	0	$\frac{\log(\sqrt{2})}{T}$
ChebNet	$\mathbf{C}^{(0)} = \mathbf{I}$	1	$[0, 2]$	2
	$\mathbf{C}^{(1)} = 2\mathbf{L}/\lambda_{\max} - \mathbf{I}$		1, -1	$2 - \sqrt{2}$
	$\mathbf{C}^{(k)} = 2\mathbf{C}^{(2)}\mathbf{C}^{(s-1)} - \mathbf{C}^{(s-2)}$	
BernNet	$\frac{1}{2^K} \binom{K}{k} (2\mathbf{I} - \mathbf{L})^{K-k} \mathbf{L}^k$	Appendix C	$\frac{2}{K}$	Appendix C
Graph Gaussian	$e^{-T(\mathbf{b}-\mathbf{L})^2}$	1	b	Proposition 4.1

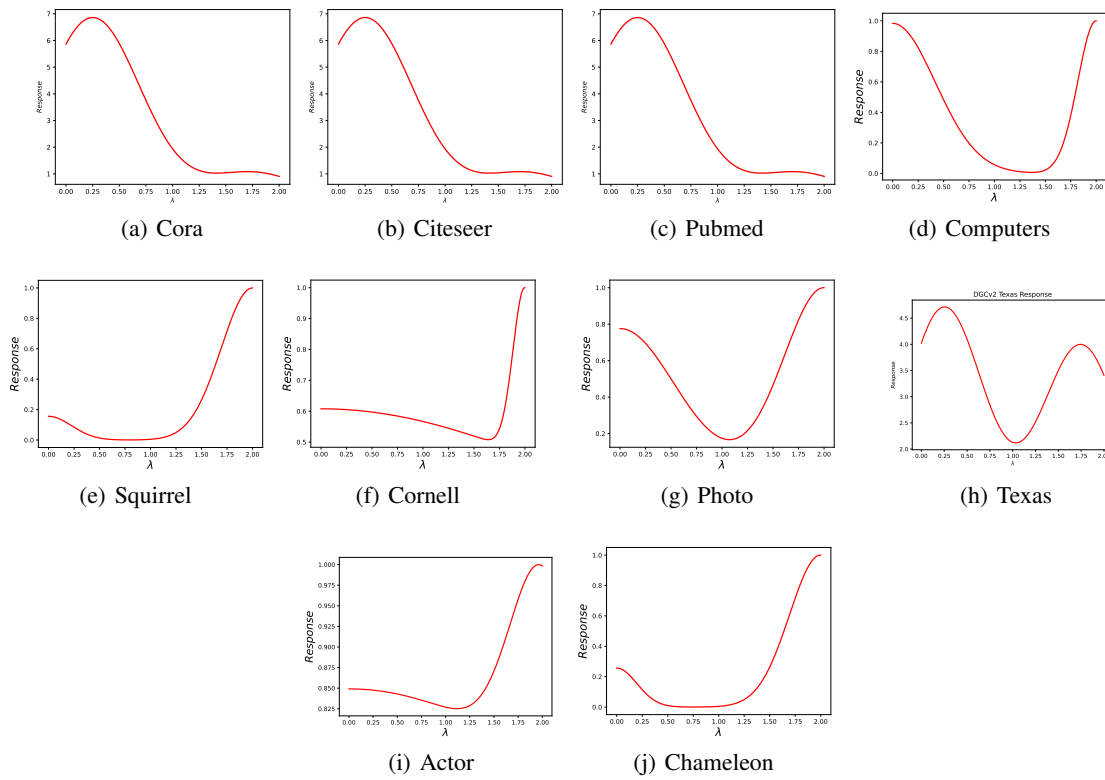


Figure 3. G^2CN 's equivalent graph filters for selected real-world datasets.