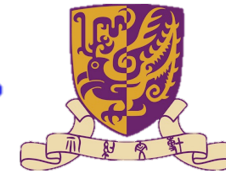




清华大学
Tsinghua University

idea

Carnegie
Mellon
University



FEATURE EXPANSION FOR GRAPH NEURAL NETWORKS

1. Tsinghua University 2. International Digital Economic Academy 3. Carnegie Mellon University 4. Mohamed bin Zayed University of Artificial Intelligence
5. Chinese University of Hong Kong

JIAQI SUN¹, LIN ZHANG², GUANGYI CHEN^{3,4}, PENG XU⁵, KUN ZHANG^{3,4}, YUJIU YANG¹

SPEAKER: JIAQI SUN

2023/07/27



CONTENT

- Background
- Motivation
- Analysis
- Method
- Experiments
- Conclusions & open questions

BACKGROUND



- Graph neural networks are widely adopted on a wide range of applications, e.g., bioinformatics and social study, especially popular for node classification task.
- Recently, many methods have studied the representations of GNNs from the perspective of optimization goals and spectral graph theory.
- However, there still no perspective could explain the existing problems integratively, e.g., over-smoothing problem and inferior performance for heterophilic graphs.
- GNNs is actually representation learning methods, while its representation space is not studied yet – **feature space**.

MOTIVATION



- We propose to analysis the feature space for existing GNNs with the help of a linear approximation.
 - Reasonableness will be indirectly verified in the post-experiments.

$$H = \overline{\text{GNN}}(X, \hat{A}) = \sum_{t=0}^{T-1} \Phi_t(X, \hat{A}) \Theta_t$$

- where $\Phi_t(X, \hat{A}) \in \mathbb{R}^{n \times d_t}$ is the non-parametric feature space constructing function that inputs the graph data (e.g., node attributes and graph structure) and outputs a feature subspace, and
- $\Theta \in \mathbb{R}^{d_t \times c}$ is the parameter space to reweight the corresponding feature subspace for each class c , and T is a hyper-parameter of the number of the feature subspaces that the GNN contains.

- We can summarize existing spatial and spectral GNN approaches as follows:

Table 1. Feature space and parameters for GNN models (better viewed in color)

	Original formula*	Linear approximation formulations
GCN (Kipf & Welling, 2017)	$H^{(k+1)} = \sigma(\hat{A}H^{(k)}W^{(k)})$	$H^{(K)} = \hat{A}^K X \prod_{i=0}^{K-1} W^{(i)}$
GIN (Xu et al., 2018)	$H^{(k+1)} = \sigma\left(\left(\epsilon^{(k)}I + \hat{A}\right)H^{(k)}W_0^{(k)}\right)W_1^{(k)}$	$H^{(K)} = \sum_{t=0}^K \hat{A}^t X \sum_{\{q_0, \dots, q_{K-t-1}\} \subseteq \{\epsilon^{(0)}, \dots, \epsilon^{(K-1)}\}} \prod_i q_i \cdot \prod_{j=0}^{K-1} W_0^{(j)} W_1^{(j)}$
GCNII (Chen et al., 2020)	$H^{(l+1)} = \sigma\left(\left(\left(1 - \alpha^{(l)}\right)\hat{A}H^{(l)} + \alpha^{(l)}H^{(0)}\right)\left(\left(1 - \beta^{(l)}\right)I + \beta^{(l)}W^{(l)}\right)\right)$	$H^{(K)} = \sum_{l=0}^{K-1} \hat{A}^l X \prod_{i=L-l}^{L-1} \left(1 - \alpha^{(i)}\right) \alpha^{(L-l-1)} \prod_{j=L-l-1}^{L-1} W^{(j)}\} + \hat{A}^K \prod_{h=0}^{K-1} \left(1 - \alpha^{(h)}\right) W^{(h)}$
ARMA (Bianchi et al., 2021)	$H^{(K)} = \sigma(\tilde{L}H^{(K-1)}W_1 + XW_2)$	$H^{(K)} = \sum_{t=0}^K \tilde{L}XW_2^t W_1^{K-t}$
APPNP (Klicpera et al., 2019)	$H^{(k+1)} = (1 - \alpha)\hat{A}H^{(k)} + \alpha H^{(0)}; H^{(0)} = \sigma(XW_1)W_2$	$H^{(K)} = \sum_{t=0}^K (1 - \alpha)^t \hat{A}^t H^{(0)} + \sum_{i=0}^{t-1} \alpha(1 - \alpha)^i \hat{A}^i H^{(0)} W_1 W_2$
ChebyNet** (Defferrard et al., 2016)	$H = \sum_{k=0}^K P_k(\hat{L})XW^{(k)}$	$H^{(K)} = \sum_{t=0}^K P_t(\hat{L})XW^{(t)}$
GPRGNN (Chien et al., 2021)	$H = \sum_{k=0}^K \gamma^{(k)} \hat{L}^k \sigma(XW_1)W_2$	$H^{(K)} = \sum_{t=0}^K \hat{L}^t X \gamma^{(t)} W_1 W_2$
BernNet (He et al., 2021)	$H = \sum_{k=0}^K \frac{1}{2^k} \binom{K}{k} \gamma^{(k)} (2I - \hat{L})^{K-k} \hat{L}^k \sigma(XW_1)W_2$	$H^{(K)} = \sum_{t=0}^K \sum_{j=0}^t \frac{1}{2^j} \binom{K}{j} \hat{L}^t \sum_{j=0}^t \gamma^{(j)} W_1 W_2$

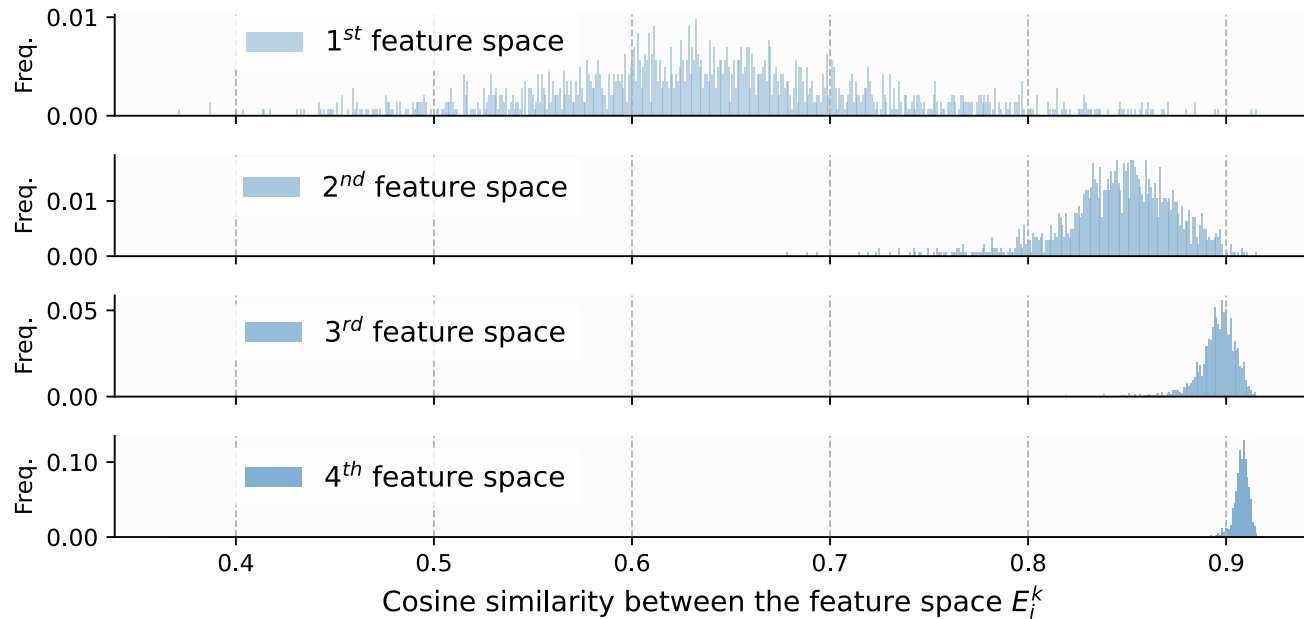
* Without specification, $H^{(0)} = X$.

** $T_k(x)$ denotes Chebyshev polynomial $P_0(x) = 1, P_1(x) = x, P_k(x) = 2xP_{k-1} - P_{k-2}$.

ANALYSIS



Proposition 3.1. Suppose the feature subspaces are constructed sequentially by $\{\Phi_t = \hat{A}^t X\}_t$. As $i \in \mathbb{Z}$ increases, the subspace Φ_{t+i} gradually tends to be linearly correlated with Φ_t .



E.g., GPRGNN,

$$E_i^k = \max_{j=0, \dots, k-1} \mu(\hat{L}^j X, \hat{L}^k X.i),$$

Figure 1. Distribution of the mutual correlation values between the later feature (sub)spaces to the previous total ones.



■ Issue I: Constraint from the weight sharing mechanism.

From Table I. we see that existing GNNs usually share parameter weights between different subspaces. Under linear correlation condition, using the weight-sharing method limits the expressiveness of the feature space.

➤ Modification I: Feature Subspaces Flattening

Theorem 3.2. *(its proof can be found in Appendix A.4)*
Suppose $\Phi_a, \Phi_b \in \mathbb{R}^{n \times d}$ are two linearly correlated feature subspaces, i.e. there exists $W_a \in \mathbb{R}^{d \times d}$ such that $\Phi_a W_a = \Phi_b$, and suppose a matrix $B \in \mathbb{R}^{n \times c}$, $c \ll d$. If B can be represented by using both subspaces with a common weight W_B , i.e., $\gamma_a \Phi_a W_B + \gamma_b \Phi_b W_B = B$ and $\gamma_a, \gamma_b \in \mathbb{R}$, then B can always be represented by only one subspace Φ_a , i.e., $\Phi_a W'_B = B$ and $W'_B \in \mathbb{R}^{d \times c}$.

Theorem 4.1. *(its proof can be found in Appendix A.5)*
Suppose $\Phi_a, \Phi_b \in \mathbb{R}^{n \times d}$ are two linearly correlated feature subspaces, i.e., there exists $W_a \in \mathbb{R}^{d \times d}$ such that $\Phi_a W_a = \Phi_b$, and a matrix $B \in \mathbb{R}^{n \times c}$, $c \ll d$. If B can be expressed by Φ_a , i.e., $\Phi_a W_B = B$, then using both subspaces Φ_a and Φ_b independently, i.e., $\Phi_a W_a + \Phi_b W_b = B$, the optimum is more easily achieved than a weight-sharing style.



- **Issue 2: Constraint from limited dimensionality of node attributes.**

When the node attributes X have a thin shape, is far from a sufficient feature space, making the regression system strictly over-determined.

In addition, without any assumption about the feature space construction, there is hardly an exact solution, e.g. heterophilic condition.

- **Modification 2: Structural Principal Components**

Theorem 4.2. *(its proof can be found in Appendix A.6)*
 Suppose the dimensionality of the node attributes is much smaller than the number of nodes, i.e., $d \ll n$, $X \in \mathbb{R}^{n \times d}$, and a z -truncated SVD of \hat{L} , which satisfies $\|U_z S_z - \hat{L}\|_2 < \epsilon$, where ϵ is a sufficiently small constant. Then the linear system $(\Phi_k, U_z S_z) W_B' = B$ can achieve a minor error than the linear system $\Phi_k W_B = B$.

$$S = \tilde{Q} \tilde{V}; \hat{A} = Q V R^T$$

$$H = \sum_{k=0}^K P_k(\hat{L}) X W^{(k)} + S W_s.$$

METHOD

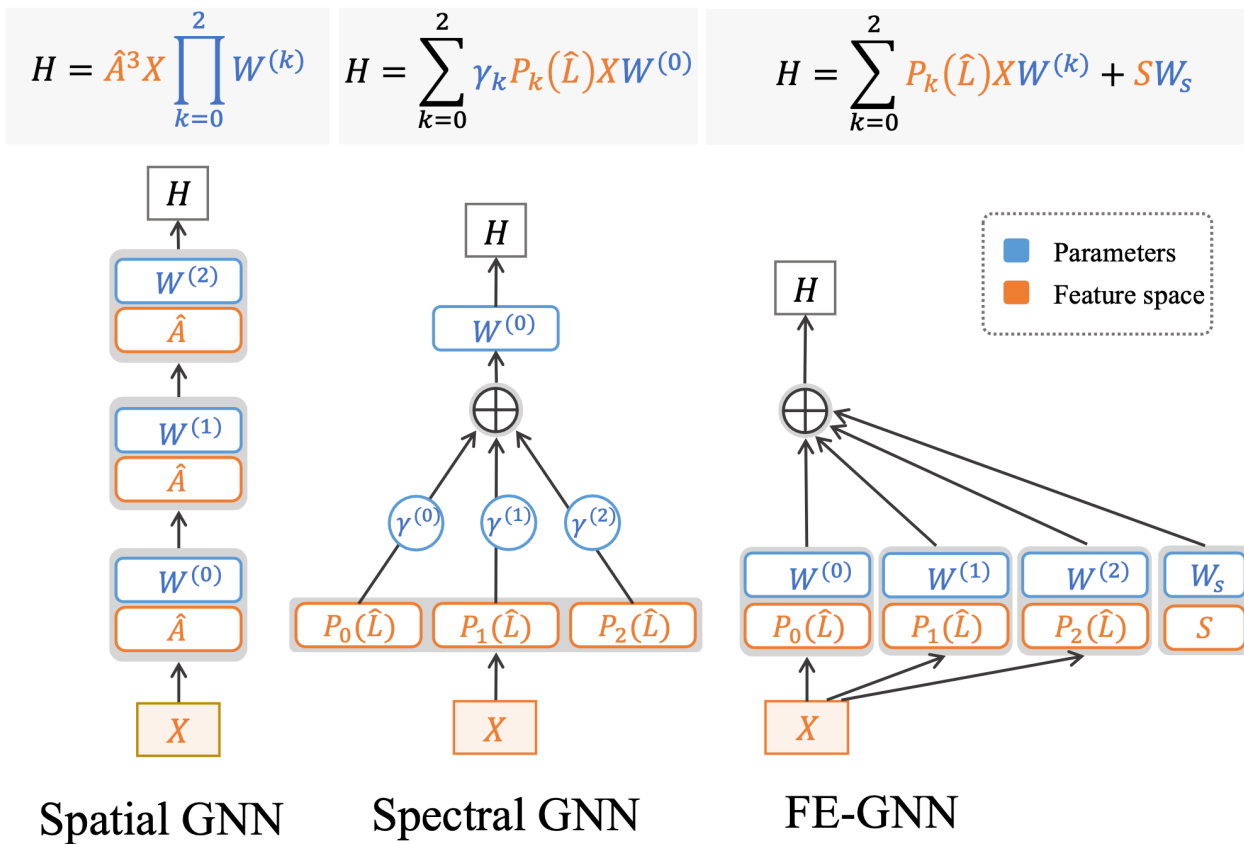


Figure 2. Architecture of our proposal

EXPERIMENTS



- Node classification

Table 2. Overall performance of FE-GNN in node classification

Type	Baseline	Time (ms)	Homophilic graphs					Heterophilic graphs		
			Cora	CiteSeer	PubMed	Computers	Photo	Squirrel	Chameleon	Actor
Spatial	MLP	-	76.70 \pm 0.15	76.67 \pm 0.26	85.11 \pm 0.26	82.62 \pm 0.21	84.16 \pm 0.13	37.86 \pm 0.39	57.83 \pm 0.31	38.99 \pm 0.17
	GCN	17.42 \pm 1.64	87.69 \pm 0.40	79.31 \pm 0.46	86.71 \pm 0.18	83.24 \pm 0.11	88.61 \pm 0.36	47.21 \pm 0.59	61.85 \pm 0.38	28.61 \pm 0.39
	GAT	18.06 \pm 1.18	88.07 \pm 0.41	80.80 \pm 0.26	86.69 \pm 0.14	82.86 \pm 0.35	90.84 \pm 0.32	33.40 \pm 0.16	51.82 \pm 1.33	33.48 \pm 0.35
	GraphSAGE	10.72 \pm 0.25	87.74 \pm 0.41	79.20 \pm 0.42	87.65 \pm 0.14	87.38 \pm 0.15	93.59 \pm 0.13	48.15 \pm 0.45	62.45 \pm 0.48	36.39 \pm 0.35
	GCNII	8.48 \pm 0.24	87.46 \pm 0.31	80.76 \pm 0.30	88.82 \pm 0.08	84.75 \pm 0.22	93.21 \pm 0.25	43.28 \pm 0.35	61.80 \pm 0.44	38.61 \pm 0.26
	APPNP	23.74 \pm 2.08	87.92 \pm 0.20	81.42 \pm 0.26	88.16 \pm 0.14	85.88 \pm 0.13	90.40 \pm 0.34	39.63 \pm 0.49	59.01 \pm 0.48	39.90 \pm 0.25
Spectral	ChebyNet	20.26 \pm 1.03	87.17 \pm 0.19	77.97 \pm 0.36	89.04 \pm 0.08	87.92 \pm 0.13	94.58 \pm 0.11	44.55 \pm 0.28	64.06 \pm 0.47	25.55 \pm 1.67
	GPRGNN	23.55 \pm 1.26	87.97 \pm 0.24	78.57 \pm 0.31	89.11 \pm 0.08	86.07 \pm 0.14	93.99 \pm 0.11	43.66 \pm 0.22	63.67 \pm 0.34	36.93 \pm 0.26
	BernNet	36.88 \pm 0.84	87.66 \pm 0.26	79.34 \pm 0.32	89.33 \pm 0.07	88.66 \pm 0.08	94.03 \pm 0.08	44.57 \pm 0.33	63.07 \pm 0.43	36.89 \pm 0.30
Unified	GNN-LF	52.77 \pm 4.50	88.12 \pm 0.06	83.66\pm0.06	87.79 \pm 0.05	87.63 \pm 0.05	93.79 \pm 0.06	39.03 \pm 0.08	59.84 \pm 0.09	41.97 \pm 0.06
	GNN-HF	53.28 \pm 4.51	88.47 \pm 0.09	83.56 \pm 0.10	87.83 \pm 0.10	86.94 \pm 0.06	93.89 \pm 0.10	39.01 \pm 0.51	63.90 \pm 0.11	42.47\pm0.07
	ADA-UGNN	14.36 \pm 0.21	88.92 \pm 0.11	79.34 \pm 0.09	90.08 \pm 0.05	89.56 \pm 0.09	94.66 \pm 0.07	44.58 \pm 0.16	59.25 \pm 0.16	41.38 \pm 0.12
	FE-GNN (C)	15.8 \pm 0.11	89.45\pm0.22	81.96 \pm 0.23	90.27\pm0.49	90.79\pm0.08	95.36 \pm 0.14	67.82 \pm 0.26	73.33\pm0.35	40.54 \pm 0.15
	FE-GNN (M)	14.6 \pm 0.32	89.09 \pm 0.22	81.76 \pm 0.23	89.93 \pm 0.23	90.60 \pm 0.11	95.45\pm0.15	67.90\pm0.23	73.26 \pm 0.38	40.91 \pm 0.22

EXPERIMENTS



■ Ablation study

Table 3. Ablation study of 1) flattening feature subspaces

	Cora	CiteSeer	Chameleon	Squirrel
k=2	89.15±0.86	81.97±1.10	73.41±1.40	67.37±0.83
k=2 (WS)	87.21±0.83	78.39±0.72	72.94±1.22	66.01±1.31
k=4	88.56±2.01	80.19±0.84	73.27±1.56	67.40±0.90
k=4 (WS)	87.28±1.38	77.72±0.86	73.23±1.72	66.43±1.72
k=8	88.92±0.88	81.11±0.89	73.85±1.52	67.93±2.04
k=8 (WS)	86.92±1.66	77.32±0.33	73.15±1.83	66.63±2.38
k=16	88.26±0.14	80.54±1.03	73.88±1.53	67.82±1.54
k=16 (WS)	87.34±1.98	78.60±0.70	72.94±1.79	66.65±2.15

Table 5. Ablation study of 2) structural principal components

	Cora	CiteSeer	PubMed	Squirrel	Chameleon
FE-GNN(C)	89.45 ±0.22	81.96 ±0.23	89.87±0.49	67.82±0.26	73.33 ±0.35
FE-GNN(M)	89.09±0.22	81.76±0.23	89.93±0.23	67.90 ±0.23	73.26±0.38
w/o norm	86.23±1.43	79.32±0.59	90.27 ±0.49	64.70±1.10	68.25±1.64
w/o S	89.20±0.93	81.95±0.87	89.76±0.46	43.21±0.99	61.54±1.52
w/o $P_k(\hat{L})X_{k>0}$	71.10±1.72	74.38±1.01	86.61±0.54	67.90±0.96	73.35±1.21
w/o $P_k(\hat{L})X_{k=0}$	84.70±1.05	58.60±2.19	85.84±0.45	65.75±0.63	72.61±1.60

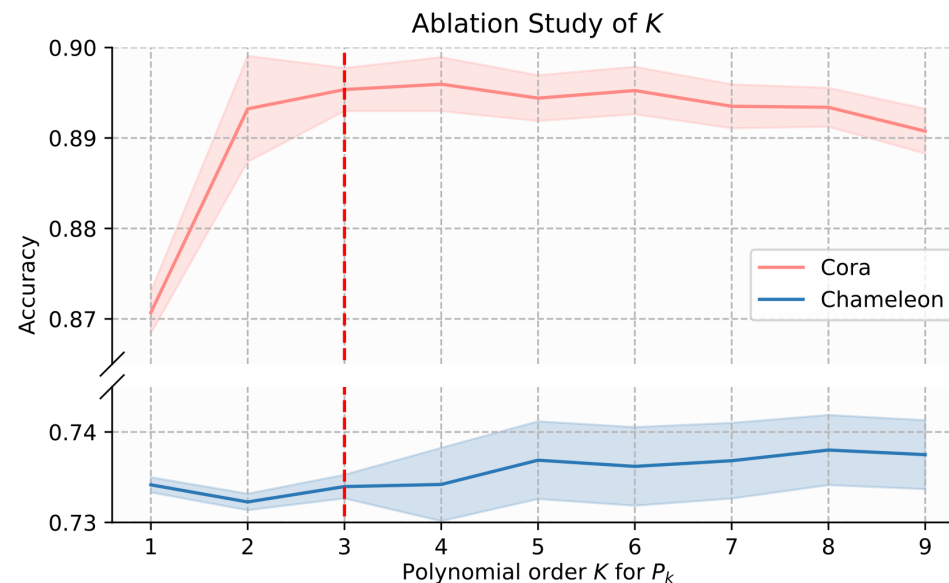


Figure 4. Ablation study of K

EXPERIMENTS



■ SVD analysis

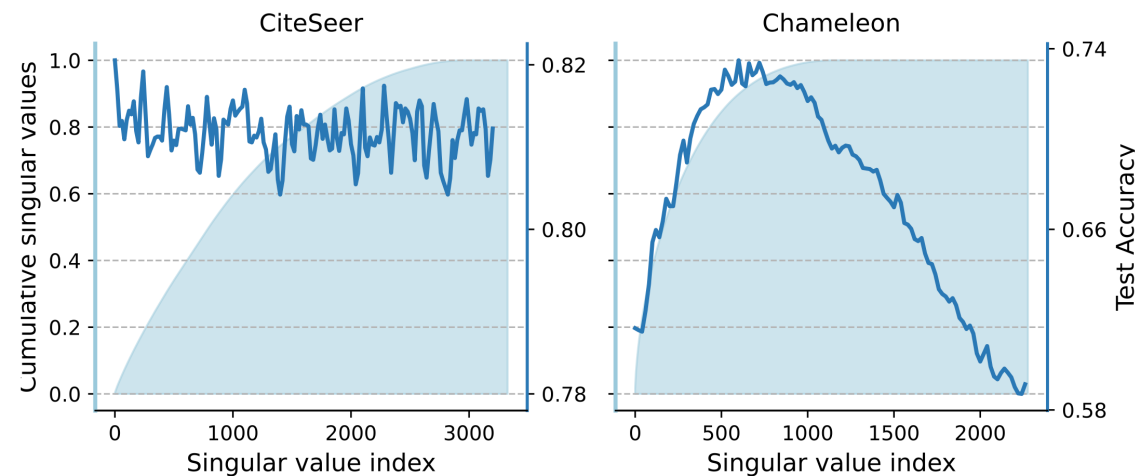


Figure 6. Sensitivity study of truncated SVD

Table 4. Time consumption of SVD

	Cora	CiteSeer	Chameleon	Squirrel	Actor
Training time (ms)	4000.01 \pm 52.23	4103.46 \pm 133.77	2818.21 \pm 81.87	6096.43 \pm 403.95	6074.39 \pm 547.34
SVD time (ms)	3.88 \pm 0.08	8.76 \pm 0.05	61.80 \pm 0.24	432.43 \pm 0.70	3.99 \pm 0.02
# of epochs	252	252	252	271	252
SVD time / Training time (%)	0.097	0.21	2.2	7.1	0.066

EXPERIMENTS



Efficiency check

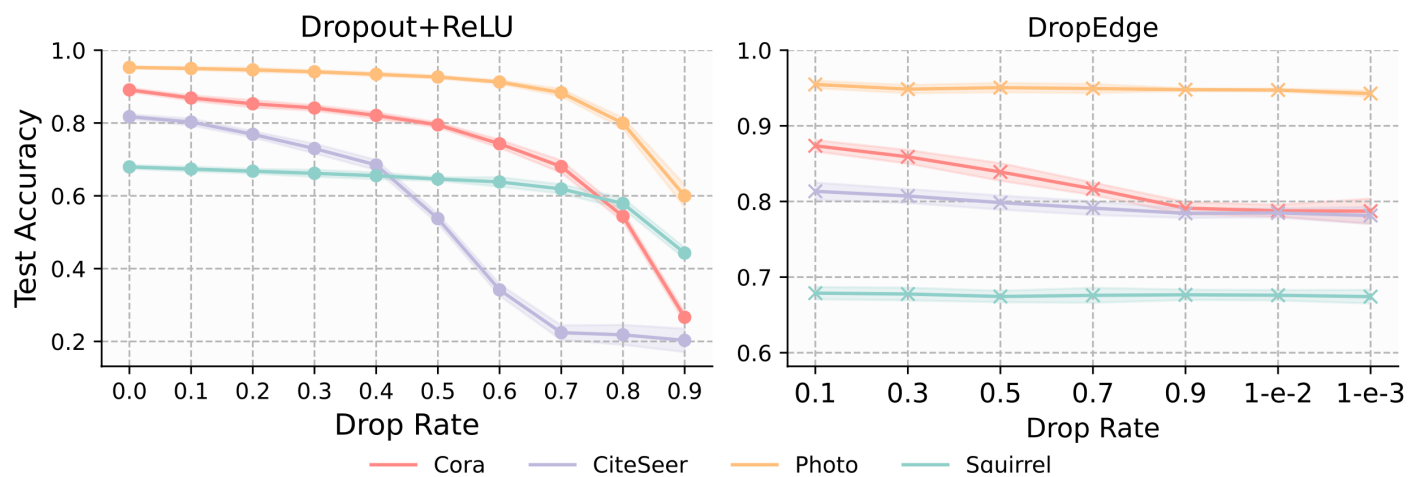


Figure 3. Analysis of the deep artifices on FE-GNN.

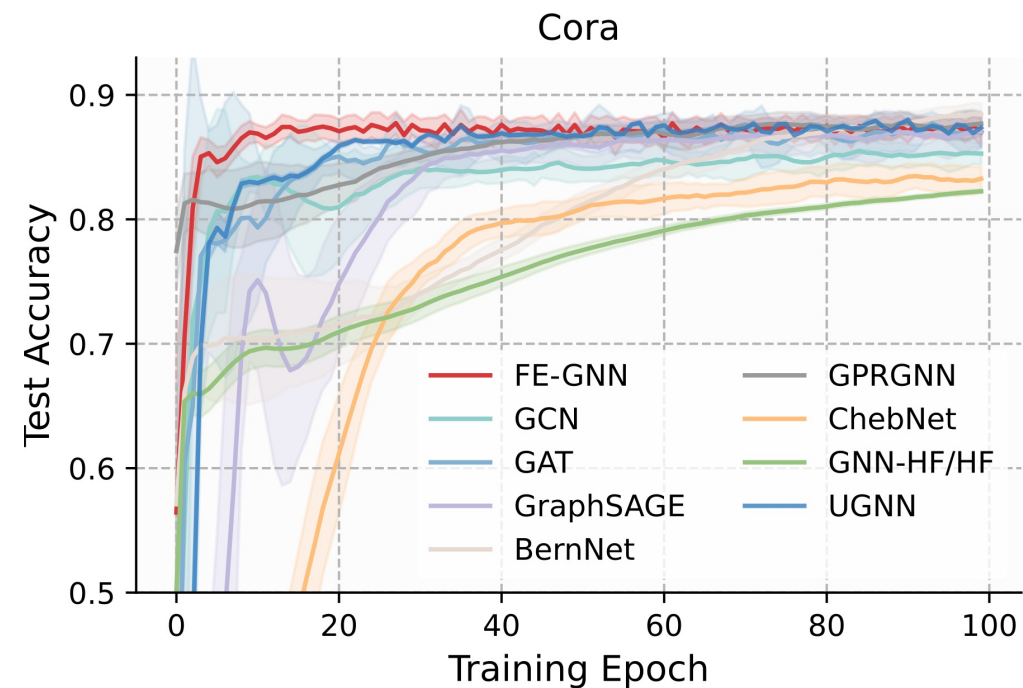


Figure 5. Convergence curve

CONCLUSIONS & OPEN QUESTIONS



- In this paper, we provide the feature space view to analyze GNNs, which separates the feature space and the parameters. Together, we provide a theoretical analysis of the existing feature space of GNNs and summarize two issues. We propose 1) feature subspace flattening and 2) structural principal components for these issues. Extensive experimental results verify their superiority.
- Limitations:
 - Nonlinear cases are not included in our work and will be considered in future work.
 - Also, the correlation between the subspaces should be studied more carefully beyond the linear correlation property; in a sense, the parameters can be further reduced by introducing reasonable constraints.
 - Finally, more feature space construction methods should be discovered for future work.



清华大学

Tsinghua University

THANK YOU!

ANY FEEDBACK OR DISCUSSIONS ARE WELCOMED!!!

CONTACT: SUNJQ20@MAILS.TSINGHUA.EDU.CN