

VII. SUPPLEMENT

A. Datasets

In the tax heterogeneous graph, the initial feature selection of the company node is shown in Table V. We select two public datasets from a heterogeneous graph benchmark [40] to verify the effectiveness of our proposed algorithm. The detailed statistics of the four datasets are shown in Table IV, and the degree distribution is shown in Figure 9.

- **PubMed.** We use PubMed provided by a heterogeneous graph benchmark [40], which contains four node types(i.e., 13561 genes (G), 20163 diseases (D), 26522 chemicals (C), 2863 species (S)) and ten edge types. Word2Vec [44] is used to aggregate the word embedding to obtain 200 dimensional features of all node types. The label information divide the diseases into 8 categories, and each labeled disease has only one label.

- **DBLP.** We sample DBLP from the heterogeneous graph benchmark [40], which contains four node types(i.e., 16254 phrases (P), 185048 authors (A), 5076 venues (V), 83 years (Y)) and six edge types. The initial features of papers and phrases are subjected to Word2Vec [44] to aggregate all word embeddings, and the initial features of authors and venues are aggregated according to corresponding paper features. The label information marks a part of the authors from the four research fields into 13 research groups, and each labeled author has only one label.

B. Classification

PubMed includes 8 classification tasks and DBLP includes 13 classification tasks. We use the Macro-F1 and Micro-F1 as the evaluation metrics. The results of the node classification task on the PubMed and DBLP data sets are shown in Table VI.

TED still has the best classification effect in PubMed and DBLP. Compared with the sub-optimal, in the PubMed and DBLP datasets, the Macro-F1 increased by 5.17% and 1.72%, and the Micro-F1 increased by 3.49% and 1.85%, respectively.

C. Clustering

We also conduct clustering experiments to evaluate the performance of different models on clustering tasks. We input the disease node embedding in the PubMed test set and the author node embedding in the DBLP test set into the KMeans algorithm for clustering. Then, we set the number of clusters to the number of label categories, i.e., the number of PubMed clusters is 8, and the number of DBLP clusters is 13. We use the NMI and ARI as the evaluation metrics. Since the performance of KMeans is affected by the initialization, we iteratively repeated the process 10 times and take the average value. The results are shown in Table VII.

As shown in Table VII, on the PubMed dataset, TED outperforms all other baselines. Compared with the sub-optimal, in the PubMed and DBLP datasets, the NMI increased by 3.65%, and the ARI increased by 4.67%. On the DBLP dataset, TED is suboptimal in the clustering task.

D. Visualization

To visualize the node embeddings learned by the model more intuitively, we use t-SNE [41] to project the disease embedding in the PubMed dataset into a two-dimensional space, and we use different colors to represent the corresponding disease categories.

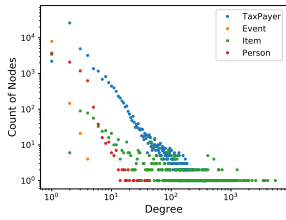
The visualization results are shown in Figure 10, where we show the embedding results of the four models of metapath2vec, R-GCN, HAN and our proposed TED. Figure 10(a) shows the embedding result of metapath2vec. We can see that the embedding is clustered into several clusters, but each cluster contains multiple types of nodes, which does not distinguish different types of nodes. Figure 10(b) is the embedding result of R-GCN. The figure is much better than Figure 10(a). R-GCN basically distinguishes each category, but there is also a cluster containing nodes of various categories, which also explains why R-GCN, a spectral domain method, may have over-smoothing problem. Figure 10(c) is the embedding result of HAN. The classification result of HAN is very good. There are obvious boundaries between categories, but the distance within the cluster is relatively large. Figure 10(d) is our model, which distinguishes different categories well and has a good degree of similarity within the categories. This also reflects that TED can learn better node embedding.

TABLE IV
STATISTICS OF THE DATASETS.

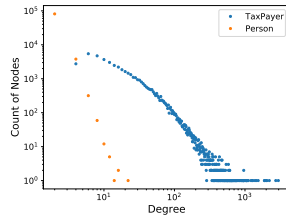
Dataset	node type	node	edge type	edge	attribute	label type	label node
T20H	4	112015	6	198903	300	2	1770
T15S	2	132522	2	467273	300	2	2072
PubMed	4	63109	10	244986	200	8	454
DBLP	4	206461	6	288959	300	13	618

TABLE V
COMPANY FEATURE SELECTION.

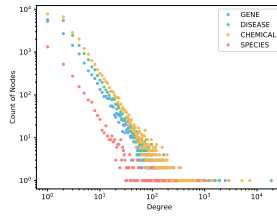
Feature Category	Attribute Symbols	Attribute Type	Attribute Meaning
Registration Information	HYMC	text	industry
	SCCYDM	category	industry code
	DJZCLXMC	text	taxpayer registration type
	NSRZTMC	text	taxpayer current status
	FDDBRNL	number	age of legal representative
	CWFZRNL	number	age of financial officer
	BSRNL	number	age of tax preparers
	JYFW	text	business scope
	CYRS	number	number of people engaged
	ZCZB	number	registration capital
Business Information	TZZE	number	total investment
	XFKPZB	number	sales side invoicing percentage
	GFKPZB	number	purchaser invoicing percentage
	XGFPSB	number	the ratio of the number of invoices from the seller to the buyer
	FPSYSZB	number	the proportion of upstream companies
	FPXYSZB	number	the proportion of downstream companies
	XSYSB	number	ratio of downstream companies to upstream companies
	XFJEZB	number	seller's amount percentage
	GFJEZB	number	purchaser's amount percentage
	XGFPJEB	number	sales to purchase amount ratio



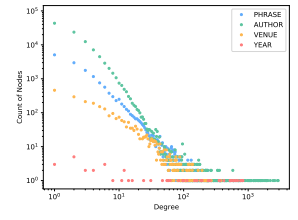
(a) T20H



(b) T15S



(c) PubMed



(d) DBLP

Fig. 9. Degree distribution in four real-world tax datasets.

TABLE VI
EXPERIMENTAL RESULTS (%) OF NODE CLASSIFICATION TASK ON PUBMED AND DBLP DATASETS.

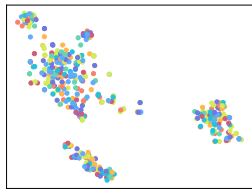
Datasets	Metrics	H2V	PTE	m2v	TransE	ConvE	DistMult	ComplEx	HAN	MAGNN	R-GCN	HGT	TED
PubMed	Macro-F1	11.41	9.25	12.48	12.84	10.31	10.02	8.26	52.03	33.98	46.94	36.38	57.20
	Micro-F1	16.73	12.11	14.53	15.87	13.00	14.97	14.09	55.81	37.21	48.84	38.37	59.30
DBLP	Macro-F1	15.37	40.89	42.91	36.46	40.71	18.37	32.87	42.35	28.18	25.44	32.79	47.30
	Micro-F1	32.86	50.98	55.50	50.49	52.44	33.50	47.08	56.41	38.46	36.75	44.44	58.97

*Note: H2V is HIN2Vec, m2v is metapath2vec.

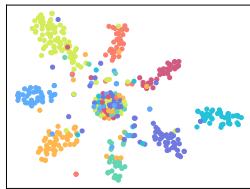
TABLE VII
EXPERIMENTAL RESULTS (%) OF CLUSTERING TASK ON PUBMED AND DBLP DATASETS.

Datasets	Metrics	H2V	PTE	m2v	TransE	ConvE	DistMult	ComplEx	HAN	MAGNN	R-GCN	HGT	TED
PubMed	NMI	16.96	13.21	16.01	14.87	13.34	15.43	15.24	43.15	20.20	26.97	27.41	46.80
	ARI	0.21	-0.13	-1.74	-1.00	-1.16	-0.78	-1.27	22.46	0.31	0.24	4.48	27.13
DBLP	NMI	29.71	34.63	41.32	42.39	43.30	36.00	43.23	54.25	33.90	28.26	44.60	51.91
	ARI	2.11	6.99	18.58	16.79	15.78	9.90	21.23	26.68	9.55	4.09	15.84	23.51

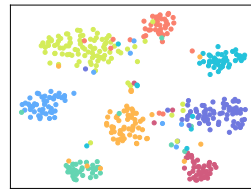
*Note: H2V is HIN2Vec, m2v is metapath2vec.



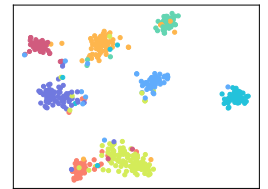
(a) metapath2vec



(b) R-GCN



(c) HAN



(d) TED

Fig. 10. Visualization on the PubMed dataset. Each node represents a disease, and the color of the node represents the type of disease.