

Appendix: Four-set Hypergraphlets for Characterization of Directed Hypergraphs

A. Datasets

We describe the representation, sources, and preprocessing steps of the datasets used in this work. As a default preprocessing step, we remove all duplicate hyperarcs and self-loops.

- **Metabolic datasets:** We use two metabolic datasets, metabolic-iAF1260b, and metabolic-iJO1366. Each node represents a gene, and each hyperarc represents a metabolic reaction, where each head and tail set indicates a set of genes. When the genes in the tail set participate in a metabolic reaction, they become the genes in the head set of the corresponding hyperarc. They are provided in the form of directed hypergraphs [1] which do not require any preprocessing step.
- **Email datasets:** We use two email datasets, email-enron [2], and email-eu [3]. Each node represents an account, and each hyperarc represents an email from a sender to one or more recipients, where the tail set consists of a node representing the sender, and the head set consists of nodes representing the recipients. We transformed the original pairwise graph into a directed hypergraph by considering all edges occurring at the same timestamp from the same sender as a single email (hyperarc). Note that the size of tail sets is always 1 in these datasets. (i.e., $|T_i| = 1, \forall i = \{1, \dots, |E|\}$.)
- **Citation datasets:** We use two citation datasets, citation-data-science, and citation-software, which are from the DBLP citation dataset [4], [5]. We extracted papers in the fields of data science or software from the dataset. Nodes represent authors and the (head and tail) sets indicate co-authors of each publication. Hyperarcs indicate citation relationships, with the tail set representing the paper that cites the head0set paper. We discarded duplicate hyperarcs, keeping only one. The publication year of the tail set paper was recorded as the timestamp of each hyperarc for temporal analysis.
- **Question & Answering datasets:** We use two question & answering datasets, qna-math and qna-server. Following [6], we created a directed hypergraph from the log data of the two question-answering sites [7]: Math Exchange and Server Fault. Each node represents a user, and each hyperarc represents a post, with the tail set consisting of the answerers and the head set consisting of the questioner. Note that the size of head sets is always 1 in these datasets. (i.e., $|H_i| = 1, \forall i = \{1, \dots, |E|\}$.)
- **Bitcoin transaction datasets:** We use three bitcoin transaction datasets, bitcoin-2014, bitcoin-2015, and bitcoin-2016, created from the original datasets [8], as suggested in [6]. They contain the first 1.5 million transactions in

Nov 2014, Jun 2015, and Jan 2016, respectively. Each node represents an individual account, and each hyperarc represents a cryptocurrency transaction. The tail set of a hyperarc corresponds to the accounts selling the cryptocurrency, while the head set corresponds to the accounts buying the corresponding cryptocurrency.

B. Sample Concentration Bound

Lemma 1 (Hoeffding's inequality [9]). *Let X_1, X_2, \dots, X_n be independent random variables with $a_j \leq X_j \leq b_j$ for all $j \in [n]$. Then for any $t > 0$, we have*

$$\Pr[|X - \mu| \geq t] \leq 2 \exp \left(- \frac{2t^2}{\sum_{j=1}^n (b_j - a_j)^2} \right).$$

Proposition 1 (Sample Concentration Bound of D-MoCHY). *For any $\epsilon > 0$, if $n \geq \frac{1}{2\epsilon^2} \ln(\frac{2}{\delta})$ and $|\Omega| > 0$, $\Pr(|C[i] - |\Omega_i|| \geq |\Omega| \cdot \epsilon) \leq \delta, \forall i \in [m]$.*

Proof. Let $t := |\Omega| \cdot \epsilon$. Since $\mathbb{E}[C[i]] = |\Omega_i|$ and $X_1^i, X_2^i, \dots, X_n^i$ are independent random variables such that $0 \leq X_j^i \leq \frac{1}{np(e,e')} = \frac{|\Omega|}{n}$ where $j \in [n]$, we can apply Hoeffding's inequality (Lemma 1):

$$\begin{aligned} \Pr[|C[i] - |\Omega_i|| \geq |\Omega| \cdot \epsilon] &\leq 2 \exp \left(- \frac{2\epsilon^2 |\Omega|^2}{n(|\Omega|/n)^2} \right) \\ &\leq 2 \exp(-2\epsilon^2 n) \leq \delta. \end{aligned}$$

□

Proposition 2 (Sample Concentration Bound of CODA-A). *Let $W = \sum_{v \in V} w[v]$ and $\bar{h} = HM(|\bar{e} \cap \bar{e}'|^2)$ be a harmonic mean of $|\bar{e} \cap \bar{e}'|^2$ for all $(e, e') \in \Omega$, i.e., $HM(|\bar{e} \cap \bar{e}'|^2) = n / \sum_{(e,e') \in \Omega} \frac{1}{|\bar{e} \cap \bar{e}'|^2}$. Then for any $\epsilon > 0$, if $n \geq \frac{1}{2\epsilon^2 \bar{h}} \ln(\frac{2}{\delta})$ and $W > 0$, $\Pr(|C[i] - |\Omega_i|| \geq W \cdot \epsilon) \leq \delta, \forall i \in [m]$.*

Proof. Let $t := W \cdot \epsilon$. Since $\mathbb{E}[C[i]] = |\Omega_i|$ and $X_1^i, X_2^i, \dots, X_n^i$ are independent random variables such that $0 \leq X_j^i \leq \frac{1}{np(e,e')} = \frac{W}{n|\bar{e} \cap \bar{e}'|}$ where $j \in [n]$, we can apply Hoeffding's inequality (Lemma 1):

$$\begin{aligned} \Pr[|C[i] - |\Omega_i|| \geq W \cdot \epsilon] &\leq 2 \exp \left(- \frac{2\epsilon^2 W^2}{\sum_{(e,e') \in \Omega} \left(\frac{W}{n|\bar{e} \cap \bar{e}'|} \right)^2} \right) \\ &\leq 2 \exp(-2\epsilon^2 n \bar{h}) \leq \delta. \end{aligned}$$

□

C. A2A sampling Sampling

The detailed process is presented in Algorithm 1. Here, $p(e, e') = \left(\frac{1}{|N_e|} + \frac{1}{|N_{e'}|}\right) \cdot \frac{1}{|E_{\geq 1}|}$ where $E_{\geq 1} = \{e : N_e \geq 1\}$. Assume $E_{\geq 1}$ is given at first. Also, for space efficiency, we assume N_e is maintained (Line 4). $HM(A, B)$ on Line 5 denotes the harmonic mean of A and B . Following the same flow of proofs of Proposition ?? and ??, the unbiasedness and variance analyses are immediate.

Algorithm 1: A2A Sampling

Input: (1) a directed hypergraph: $G = (V, E)$
(2) # of samples $n = q \cdot |E|$ for a given fraction $q \in (0, 1)$
Output: $C[i]$ for every $i \in [m]$

```

1  $C[i] \leftarrow 0, \forall i \in [m]$ 
2 for  $1 : n$  do
3   Choose  $e \in E_{\geq 1}$  uniformly at random
4    $N_e \leftarrow \{e' \in E \setminus \{e\} : e \cap e' \neq \emptyset\}$ 
5    $C[f(e, e')] \leftarrow C[f(e, e')] + \frac{|E_{\geq 1}|}{2 \cdot n} \cdot HM(|N_e|, |N_{e'}|)$ 
6 return  $C$ 

```

Proposition 3 (Unbiasedness of A2A sampling). *Algorithm 1 is unbiased, i.e., $\mathbb{E}[C[i]] = |\Omega_i|$.*

Proposition 4 (Variance of A2A sampling). *The variance of $C[i]$ obtained by Algorithm 1 is*

$$\begin{aligned}
Var[C[i]] &= \sum_{(e, e') \in \Omega_i} \frac{1}{n} \left(\frac{1}{p(e, e')} - 1 \right) \\
&= \sum_{(e, e') \in \Omega_i} \frac{1}{n} \left(\frac{|E_{\geq 1}|}{2} \cdot HM(|N_e|, |N_{e'}|) - 1 \right).
\end{aligned}$$

Proposition 5 (Time & Space complexity of A2A sampling). *The time complexity of Algorithm 1 is $O(n \cdot (\max_{e \in E} |\bar{e}| \cdot \max_{e \in E} |N_e|))$. Its space complexity is $O(\sum_{e \in E} |\bar{e}|)$.*

Proof. The information of a given directed graph is stored in $O(\sum_{e \in E} |\bar{e}|)$ space at first. For time complexity, $O(\max_{e \in E} |\bar{e}| \cdot \max_{e \in E} |N_e|)$ time is required assuming $O(p \cdot q)$ time is taken for set union when there are p sets, of which size bounded by q . For space complexity, $O(|N_e|) \in O(\sum_{e \in E} |\bar{e}|)$ space is needed. Checking $f(e, e')$ requires $O(\max_{(e, e') \in \Omega} \min(|\bar{e}|, |\bar{e}'|))$ -time, which is bounded by $O(\max_{e \in E} |\bar{e}| \cdot \max_{e \in E} |N_e|)$. \square

D. Count Distributions

We analyze the occurrence distributions of DHGs in real-world and randomized directed hypergraphs (DHs). To ensure statistical significance, we generate ten randomized DHs and report the average counts. As shown in Figure 1-5, the counts of DHGs in real-world directed hypergraphs are distinct from those in randomized directed hypergraphs.

E. Error measure for counts of DHGs

As shown in Figure 6, the error values of CODA-A are significantly lower than those of D-MoCHy and A2A sampling on all datasets, when their running time is similar. When

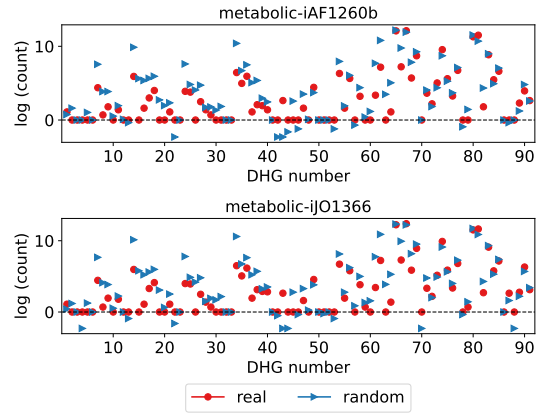


Fig. 1. [Metabolic datasets] Log counts of DHGs in both real-world and randomized directed hypergraphs.

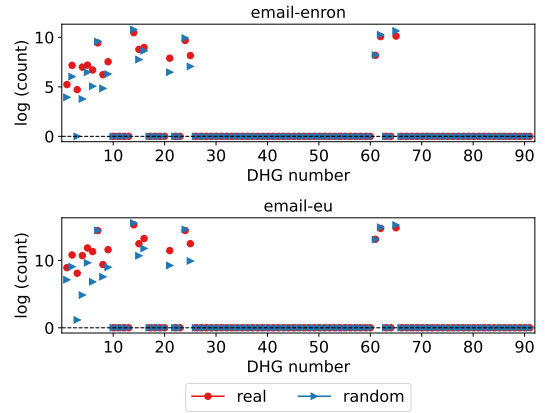


Fig. 2. [Email datasets] Log counts of DHGs in both real-world and randomized directed hypergraphs.

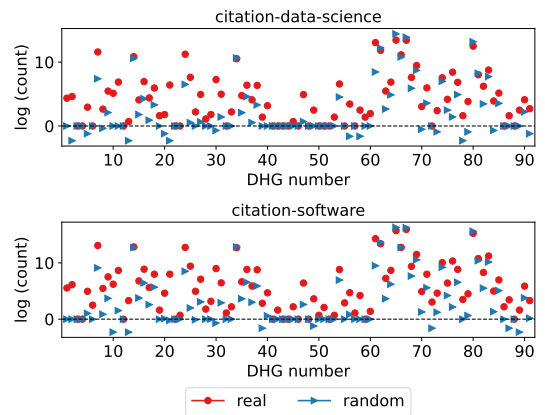


Fig. 3. [Citation datasets] Log counts of DHGs in both real-world and randomized directed hypergraphs.

the error values are similar in different methods, CODA-A is up to 40× faster than D-MoCHy and A2A sampling. In all methods, the error tends to be smaller as the number of samples increases.

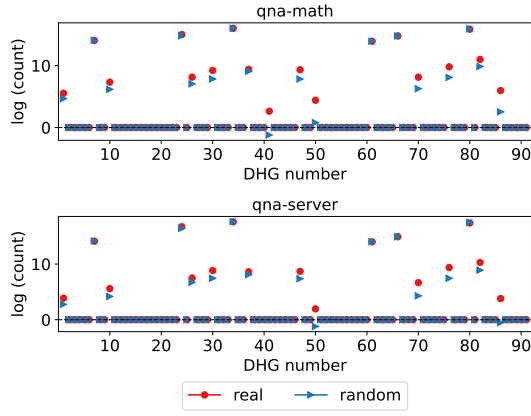


Fig. 4. [Qna datasets] Log counts of DHGs in both real-world and randomized directed hypergraphs.

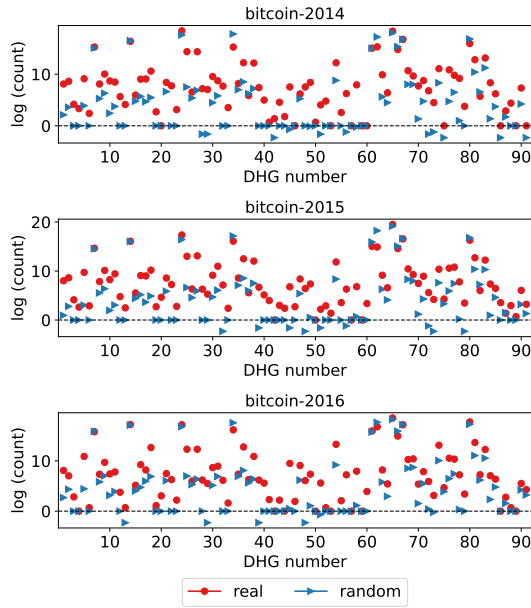


Fig. 5. [Bitcoin datasets] Log counts of DHGs in both real-world and randomized directed hypergraphs.

F. Domain-specific patterns

We compare the CPs and CP distances from the same and different domain datasets. As shown in Figure 7, CPs of DHs from the same domain tend to be similar, while those from the different domains tend to be different. When we define the CP distance using the Manhattan distance, the within CP distance is the CP distance between the same domain datasets, and the across CP distance is the CP distance between the different domain datasets. The within CP distances, across CP distances, and the ratio of these two are in Table I. These numerical results indicate that CPs from the same domain datasets are closer to each other.

G. Discoveries

To examine time-evolving patterns, we conduct a case study, excluding the metabolic datasets which lack temporal information. We consider a temporal DH $G = (V, E)$ where E

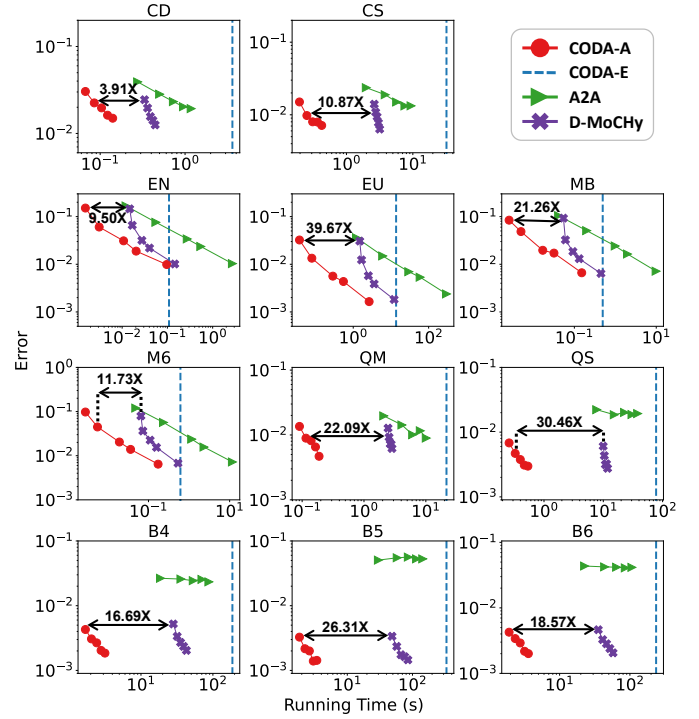


Fig. 6. Estimation error err^G (where a lower error indicates better performance) of each algorithm on all datasets. Note that CODA-A consistently provides the best trade-off among the approximate algorithms.

TABLE I
THE WITHIN- AND ACROSS-DOMAIN CP DISTANCES FROM EACH DATASET, AND THE RATIO OF THESE TWO DISTANCES. DATASETS FROM THE SAME DOMAIN TEND TO HAVE SMALLER CP DISTANCES.

Dataset	within	across	ratio
MB	0.0142	0.1420	9.988
M6	0.0142	0.1434	10.084
EN	0.0091	0.1001	11.002
EU	0.0091	0.1008	11.078
CD	0.0244	0.0971	3.977
CS	0.0244	0.0955	3.911
QM	0.0067	0.1011	15.204
QS	0.0067	0.1012	15.212
B4	0.0196	0.1054	5.381
B5	0.0196	0.1045	5.334
B6	0.0219	0.1047	4.772

has timestamp τ_e for each $e \in E$, i.e., $e = \langle H, T, \tau_e \rangle$. For the citation datasets, each hyperedge representing a publication is assigned a timestamp equal to the publication year of its tail set. Also, in the email, qna, and bitcoin datasets, each hyperarc is assigned a unique timestamp.

For each DH, we consider 10 timestamps with equal intervals $\{t_1, t_2, \dots, t_{10}\}$, where $t_1 = \min_{e \in E} \tau_e$ and $t_{10} = \max_{e \in E} \tau_e$. We create 10 snapshots (i.e., sub-DHs) where the edge set of each i -th sub-DH is $E_i = \{e : \tau_e \leq t_i\}$ with the node set $V_i = \bigcup_{e \in E_i} \bar{e}$. Then, we calculate the ratio of every DHG count of a sub-DH.

In Figure 8, we visualize DHGs whose ratio is greater than specific thresholds while aggregating the rest as Others. In addition, we summarize the time-evolving patterns of the top

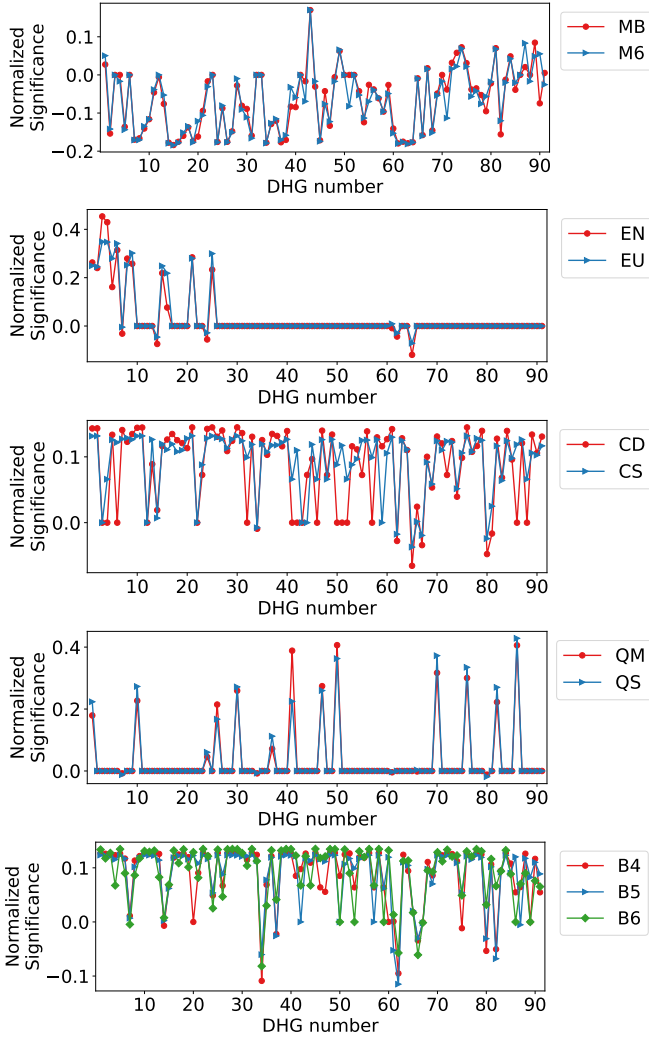


Fig. 7. Characteristic profile (CP) of metabolic, email, citation, qna, and bitcoin datasets with ten randomized directed hypergraphs (DHGs). The CPs of DHs from the same domain tend to be remarkably similar, while those from different domains tend to be different.

10 most frequent DHGs in Figure 9. As shown in the figures, datasets from the same domain have the same set of frequent DHGs and exhibit similar time-evolving tendencies for each DHG.

Citation datasets: In citation datasets, the ratios of DHGs 65, 67, and 80 increase, while those of DHGs 7, 61, and 62 decrease. In DHGs 65, 67, and 80, all non-intersecting regions (i.e., regions 1, 4, 5, and 8 in Figure ??) contain at least one node, while in DHGs 7, 61, and 62, some non-intersecting regions are empty. That is, the number of non-empty non-intersecting regions increases over time.

Email datasets: Among the frequent DHGs, only the ratio of DHG 24 decreases, whereas the ratios of DHGs 7, 14, 62, and 65 increase. DHG 24 is distinct from the other DHGs in that two tail sets intersect in it. Given that the size of tail sets in email datasets is always 1, the probability of two hyperarcs sharing the same tail set decreases over time, resulting in the decrease of the DHG 24 ratio.

Qna datasets: In qna datasets, there is a dramatic decrease in

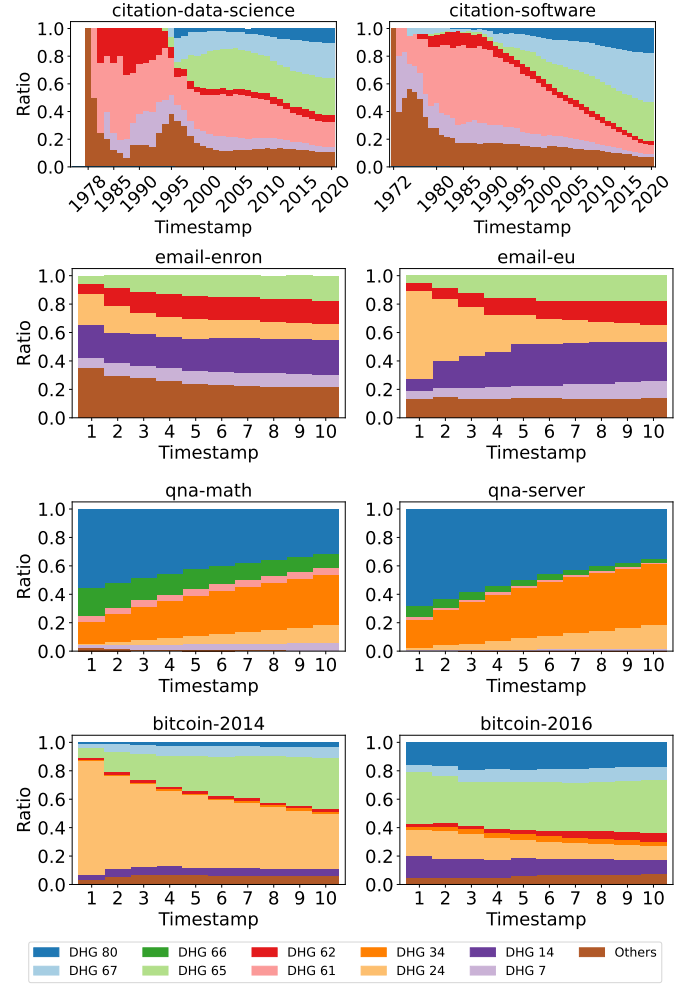


Fig. 8. DHs from the same domain have the same set of frequent DHGs and exhibit similar time-evolving tendencies for each DHG. We visualize DHGs whose ratio is greater than specific thresholds while aggregating the rest as Others. The thresholds for the citation, email, qna, and bitcoin datasets are 0.03, 0.1, 0.01, and 0.03, respectively.

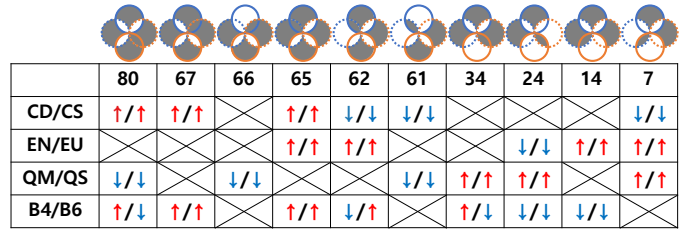


Fig. 9. Time-evolving trend. Each row represents a domain of datasets and each column represents a DHG. DHGs marked with 'X' have ratios less than a predefined threshold. The ratios of DHGs marked with ↑ or ↓ tend to increase or decrease, respectively.

the ratio of DHG 80 and a dramatic increase in the ratio of DHG 34. The ratios of DHGs 61, 66, and 80 decrease, while the ratios of DHGs 7, 24, and 34 increase. The increasing DHGs have empty non-intersecting areas in their tail sets, while the declining DHGs have no such areas. This indicates that the number of users who frequently answer questions increases.

Bitcoin datasets: Although the two bitcoin datasets share the

same set of frequent DHGs, their tendencies to increase or decrease over time differ in them. For example, DHG 65 becomes dominant in the bitcoin-2016 dataset, with a ratio of 0.37, while it has a ratio of only 0.07 in the other dataset. Conversely, DHG 24 becomes dominant in the bitcoin-2014 dataset, with a ratio of 0.81, but only has a ratio of 0.10 in the other dataset. The main difference between DHG 65 and DHG 24 is the presence or absence of an intersection between tail sets. DHG 24 has only an intersection between tail sets, but DHG 65 does not. This indicates that the diversity of accounts participating in transactions increases over time.

H. Application results

In this section, we report the average results of the hyperarc prediction problem in Table II, III. The best performances are in bold. The use of DHG vectors results in the best performance on most datasets, indicating that informative feature vectors can be obtained from DHGs. In particular, DHG shows a 56% improvement in performance on the **citation-software (CS)** dataset compared to the competitors.

REFERENCES

- [1] N. Yadati, V. Nitin, M. Nimishakavi, P. Yadav, A. Louis, and P. Talukdar, "Nhp: Neural hypergraph link prediction," in *CIKM*, 2020.
- [2] P. Chodrow and A. Mellor, "Annotated hypergraphs: models and applications," *Applied Network Science*, vol. 5, no. 1, pp. 1–25, 2020.
- [3] J. Leskovec and A. Krevl, "SNAP Datasets: Stanford large network dataset collection," <http://snap.stanford.edu/data>, Jun. 2014.
- [4] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su, "Arnetminer: Extraction and mining of academic social networks," in *KDD*, 2008.
- [5] A. Sinha, Z. Shen, Y. Song, H. Ma, D. Eide, B.-J. Hsu, and K. Wang, "An overview of microsoft academic service (mas) and applications," in *WWW*, 2015.
- [6] S. Kim, M. Choe, J. Yoo, and K. Shin, "Reciprocity in directed hypergraphs: Measures, findings, and generators," in *ICDM*, 2022.
- [7] S. Exchange, "Stack exchange data dump," <https://archive.org/details/stackexchange>, 2020.
- [8] J. Wu, J. Liu, W. Chen, H. Huang, Z. Zheng, and Y. Zhang, "Detecting mixing services via mining bitcoin transaction network with hybrid motifs," *IEEE Trans. Syst. Man Cybern.: Syst.*, vol. 52, no. 4, pp. 2237–2249, 2021.
- [9] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *The collected works of Wassily Hoeffding*, pp. 409–426, 1994.

TABLE II

HYPERARC PREDICTION ACCURACY RESULTS OF EACH DATASET. THE BEST PERFORMANCES ARE IN BOLD. USING DHG VECTORS OUTPERFORMS ALL BASELINE METHODS IN ALMOST ALL CASES, INDICATING THAT INFORMATIVE FEATURE VECTORS CAN BE OBTAINED BY DHGs.

Model	Dataset	DHG	h-motif	triad	n2v	h2v	deep-h	Dataset	DHG	h-motif	triad	h2v	h2v	deep-h
LR	MB	0.656±0.023	0.649±0.015	0.549±0.018	0.504±0.017	0.511±0.012	0.701±0.014	M6	0.656±0.016	0.648±0.019	0.584±0.011	0.506±0.011	0.505±0.011	0.713±0.012
RF		0.690±0.011	0.681±0.021	0.650±0.009	0.531±0.014	0.518±0.012	0.714±0.015		0.698±0.014	0.704±0.017	0.666±0.010	0.532±0.010	0.527±0.015	0.723±0.017
DT		0.651±0.009	0.632±0.019	0.597±0.016	0.512±0.014	0.506±0.011	0.583±0.016		0.656±0.018	0.641±0.015	0.612±0.009	0.509±0.014	0.517±0.016	0.582±0.014
KNN		0.696±0.014	0.696±0.018	0.625±0.014	0.537±0.014	0.534±0.016	0.704±0.014		0.682±0.014	0.691±0.020	0.628±0.019	0.520±0.016	0.539±0.011	0.725±0.013
MLP		0.654±0.011	0.646±0.014	0.603±0.019	0.533±0.008	0.537±0.014	0.705±0.014		0.656±0.012	0.653±0.008	0.618±0.015	0.537±0.009	0.539±0.011	0.710±0.012
XGB		0.708±0.016	0.666±0.024	0.625±0.015	0.514±0.016	0.519±0.015	0.695±0.015		0.709±0.022	0.681±0.017	0.657±0.013	0.519±0.017	0.522±0.013	0.703±0.014
LGBM		0.697±0.062	0.658±0.053	0.626±0.063	0.527±0.054	0.516±0.012	0.698±0.012		0.720±0.057	0.663±0.056	0.632±0.058	0.509±0.050	0.530±0.011	0.702±0.014
HGNN		0.549±0.055	0.543±0.040	0.538±0.063	0.483±0.045	0.498±0.067	0.526±0.057		0.569±0.055	0.553±0.036	0.550±0.043	0.499±0.056	0.522±0.037	0.545±0.036
HGCN		0.666±0.067	0.535±0.051	0.555±0.055	0.507±0.045	0.491±0.059	0.538±0.050		0.653±0.071	0.532±0.053	0.566±0.059	0.504±0.045	0.498±0.048	0.548±0.060
UGCNII		0.618±0.061	0.625±0.038	0.596±0.069	0.491±0.050	0.494±0.054	0.597±0.070		0.621±0.051	0.616±0.048	0.610±0.050	0.518±0.047	0.514±0.034	0.633±0.035
Max		0.708±0.016	0.696±0.018	0.650±0.009	0.537±0.014	0.537±0.014	0.714±0.015		0.720±0.057	0.704±0.017	0.666±0.010	0.537±0.009	0.539±0.011	0.725±0.013
Avg.	0.659±0.047	0.633±0.054	0.596±0.038	0.514±0.018	0.512±0.016	0.646±0.076	0.662±0.044	0.638±0.057	0.612±0.037	0.515±0.012	0.521±0.013	0.658±0.074		
Rank Avg.	1.600±0.516	2.700±0.823	3.600±0.699	5.500±0.527	5.500±0.527	2.100±1.287		1.700±0.823	2.700±0.675	3.600±0.699	5.600±0.516	5.400±0.516	2.000±1.247	
LR	EN	0.804±0.014	0.752±0.011	0.732±0.017	0.578±0.017	0.492±0.012	0.590±0.018	EU	0.869±0.001	0.776±0.005	0.837±0.004	0.618±0.008	0.496±0.003	0.659±0.006
RF		0.796±0.013	0.773±0.017	0.712±0.023	0.626±0.024	0.562±0.024	0.592±0.022		0.907±0.003	0.838±0.004	0.839±0.003	0.652±0.003	0.515±0.002	0.668±0.005
DT		0.705±0.011	0.689±0.018	0.654±0.020	0.551±0.022	0.528±0.018	0.542±0.021		0.849±0.003	0.761±0.005	0.787±0.005	0.546±0.004	0.504±0.007	0.564±0.007
KNN		0.778±0.014	0.737±0.016	0.694±0.019	0.636±0.020	0.571±0.017	0.567±0.022		0.875±0.002	0.780±0.005	0.838±0.005	0.573±0.002	0.556±0.012	0.677±0.003
MLP		0.805±0.014	0.751±0.011	0.731±0.013	0.639±0.021	0.551±0.018	0.588±0.023		0.906±0.003	0.821±0.007	0.857±0.005	0.660±0.016	0.507±0.005	0.675±0.011
XGB		0.775±0.018	0.763±0.014	0.709±0.026	0.614±0.020	0.579±0.018	0.577±0.020		0.903±0.003	0.831±0.005	0.854±0.005	0.654±0.005	0.522±0.005	0.656±0.003
LGBM		0.756±0.059	0.763±0.056	0.709±0.060	0.609±0.064	0.580±0.018	0.581±0.019		0.906±0.010	0.839±0.010	0.856±0.011	0.645±0.027	0.512±0.003	0.645±0.005
HGNN		0.499±0.049	0.543±0.063	0.538±0.074	0.513±0.055	0.526±0.048	0.512±0.058		0.529±0.020	0.523±0.020	0.520±0.015	0.512±0.018	0.505±0.014	0.513±0.017
HGCN		0.693±0.117	0.651±0.090	0.703±0.101	0.536±0.061	0.566±0.076	0.550±0.069		0.742±0.072	0.638±0.060	0.790±0.136	0.512±0.020	0.519±0.038	0.547±0.054
UGCNII		0.710±0.050	0.708±0.065	0.727±0.046	0.673±0.045	0.689±0.055	0.582±0.051		0.783±0.013	0.726±0.009	0.859±0.008	0.724±0.014	0.740±0.020	0.706±0.013
Max		0.805±0.014	0.773±0.017	0.732±0.017	0.673±0.045	0.689±0.055	0.592±0.022		0.907±0.003	0.839±0.010	0.859±0.008	0.724±0.014	0.740±0.020	0.706±0.013
Avg.	0.732±0.092	0.713±0.071	0.691±0.058	0.598±0.051	0.564±0.052	0.568±0.026	0.827±0.119	0.753±0.102	0.804±0.103	0.610±0.071	0.538±0.073	0.631±0.065		
Rank Avg.	1.800±1.549	2.000±0.667	2.500±0.850	4.400±0.699	5.100±1.101	5.200±0.632		1.200±0.422	3.000±0.471	1.900±0.568	5.100±0.316	5.600±0.966	4.200±0.632	
LR	CD	0.921±0.004	0.751±0.009	0.602±0.004	0.527±0.004	0.504±0.002	0.593±0.007	CS	0.919±0.002	0.767±0.002	0.662±0.005	0.541±0.005	0.508±0.005	0.625±0.007
RF		0.977±0.001	0.855±0.004	0.644±0.007	0.548±0.003	0.500±0.003	0.599±0.005		0.984±0.001	0.866±0.006	0.702±0.006	0.568±0.005	0.501±0.003	0.621±0.005
DT		0.963±0.001	0.777±0.007	0.583±0.006	0.511±0.004	0.497±0.003	0.539±0.004		0.974±0.001	0.783±0.011	0.623±0.004	0.519±0.003	0.498±0.004	0.548±0.002
KNN		0.917±0.003	0.787±0.005	0.594±0.007	0.558±0.004	0.514±0.006	0.592±0.003		0.941±0.001	0.834±0.002	0.652±0.004	0.578±0.003	0.521±0.008	0.606±0.001
MLP		0.969±0.001	0.822±0.008	0.637±0.008	0.545±0.006	0.502±0.004	0.633±0.008		0.980±0.001	0.844±0.008	0.699±0.006	0.575±0.009	0.512±0.006	0.652±0.012
XGB		0.975±0.001	0.843±0.005	0.639±0.007	0.547±0.004	0.505±0.004	0.618±0.006		0.984±0.001	0.850±0.009	0.704±0.006	0.562±0.006	0.506±0.003	0.636±0.006
LGBM		0.977±0.004	0.849±0.020	0.652±0.022	0.546±0.020	0.504±0.003	0.601±0.006		0.984±0.002	0.860±0.032	0.713±0.016	0.560±0.028	0.502±0.002	0.622±0.007
HGNN		0.595±0.009	0.543±0.015	0.534±0.013	0.542±0.012	0.535±0.017	0.519±0.013		0.555±0.007	0.534±0.011	0.529±0.008	0.553±0.009	0.568±0.009	0.521±0.006
HGCN		0.754±0.091	0.597±0.087	0.556±0.067	0.505±0.015	0.504±0.013	0.502±0.008		0.738±0.098	0.609±0.086	0.594±0.089	0.503±0.013	0.502±0.009	0.503±0.009
UGCNII		0.932±0.005	0.798±0.013	0.657±0.012	0.769±0.016	0.630±0.028	0.541±0.011		0.917±0.006	0.739±0.010	0.718±0.010	0.827±0.016	0.823±0.009	0.578±0.012
Max		0.977±0.001	0.855±0.004	0.657±0.012	0.769±0.016	0.630±0.028	0.633±0.008		0.984±0.002	0.866±0.006	0.718±0.010	0.827±0.016	0.823±0.009	0.652±0.012
Avg.	0.898±0.126	0.762±0.107	0.610±0.043	0.560±0.075	0.520±0.040	0.574±0.045	0.898±0.142	0.769±0.114	0.660±0.062	0.579±0.091	0.544±0.100	0.591±0.051		
Rank Avg.	1.000±0.000	2.000±0.000	3.300±0.675	4.500±0.850	5.600±0.699	4.600±0.966		1.100±0.316	2.400±0.843	3.400±0.843	4.400±1.075	5.200±1.751	4.500±0.850	
LR	QM	0.604±0.003	0.579±0.004	0.553±0.003	0.500±0.001	0.504±0.001	0.566±0.002	QS	0.561±0.002	0.533±0.003	0.530±0.003	0.502±0.001	0.509±0.002	0.556±0.002
RF		0.673±0.003	0.613±0.005	0.620±0.004	0.503±0.003	0.502±0.003	0.581±0.005		0.661±0.001	0.565±0.002	0.590±0.002	0.500±0.002	0.503±0.001	0.565±0.004
DT		0.617±0.003	0.547±0.006	0.572±0.004	0.502±0.001	0.502±0.004	0.547±0.004		0.626±0.002	0.528±0.002	0.572±0.003	0.500±0.002	0.501±0.002	0.545±0.002
KNN		0.601±0.002	0.576±0.003	0.553±0.003	0.504±0.002	0.506±0.004	0.523±0.001		0.619±0.002	0.538±0.002	0.566±0.002	0.504±0.002	0.502±0.003	0.576±0.001
MLP		0.679±0.002	0.598±0.004	0.590±0.002	0.505±0.002	0.503±0.002	0.611±0.011		0.668±0.001	0.544±0.003	0.603±0.003	0.509±0.002	0.506±0.003	0.598±0.013
XGB		0.681±0.002	0.607±0.005	0.629±0.003	0.502±0.002	0.508±0.002	0.601±0.004		0.679±0.001	0.575±0.002	0.612±0.003	0.503±0.001	0.503±0.001	0.585±0.005
LGBM		0.690±0.009	0.618±0.019	0.641±0.012	0.504±0.009	0.506±0.002	0.577±0.005		0.688±0.004	0.585±0.011	0.624±0.013	0.504±0.006	0.503±0.002	0.567±0.005
HGNN		0.549±0.011	0.525±0.011	0.529±0.008	0.520±0.011	0.546±0.007	0.545±0.008		0.579±0.006	0.539±0.005	0.546±0.005	0.535±0.007	0.572±0.006	0.571±0.005
HGCN		0.512±0.022	0.503±0.007	0.516±0.029	0.500±0.003	0.502±0.008	0.507±0.016		0.530±0.041	0.505±0.014	0.515±0.027	0.501±0.003	0.501±0.005	0.502±0.009
UGCNII		0.607±0.006	0.583±0.009	0.599±0.010	0.615±0.010	0.637±0.008	0.743±0.013		0.645±0.005	0.595±0.007	0.605±0.005	0.563±0.005	0.653±0.006	0.752±0.007
Max		0.690±0.009	0.618±0.019	0.641±0.012	0.615±0.010	0.637±0.008	0.743±0.013		0.645±0.005	0.595±0.007	0.624±0.013	0.563±0.005	0.653±0.006	0.752±0.007
Avg.	0.621±0.060	0.575±0.039	0.580±0.043	0.516±0.035	0.522±0.043	0.580±0.066	0.626±0.053	0.551±0.028	0.576±0.037	0.512±0.021	0.525±0.050	0.582±0.065		
Rank Avg.	1.400±0.966	3.400±1.265	2.900±1.287	5.500±0.972	4.600±1.430	3.200±1.033		1.100±0.316	3.900±0.738	2.700±0.949	5.600±0.516	4.900±1.370	2.800±0	

TABLE III

HYPERARC PREDICTION AUROC RESULTS OF EACH DATASET. THE BEST PERFORMANCES ARE IN BOLD. USING DHG VECTORS OUTPERFORMS ALL BASELINE METHODS IN ALMOST ALL CASES, INDICATING THAT INFORMATIVE FEATURE VECTORS CAN BE OBTAINED BY DHGs.

Model	Dataset	DHG	h-motif	triad	n2v	h2v	deep-h	Dataset	DHG	h-motif	triad	n2v	h2v	deep-h
LR	MB	0.728±0.022	0.724±0.015	0.580±0.017	0.506±0.021	0.509±0.013	0.785±0.017	M6	0.721±0.016	0.727±0.020	0.622±0.012	0.501±0.015	0.504±0.013	0.795±0.010
RF		0.826±0.015	0.784±0.020	0.703±0.011	0.539±0.018	0.533±0.013	0.793±0.018		0.834±0.011	0.794±0.018	0.730±0.013	0.541±0.010	0.546±0.027	0.803±0.012
DT		0.651±0.009	0.632±0.019	0.599±0.015	0.512±0.014	0.506±0.011	0.583±0.016		0.656±0.018	0.641±0.015	0.611±0.009	0.509±0.014	0.517±0.016	0.582±0.014
KNN		0.762±0.013	0.755±0.018	0.673±0.011	0.550±0.018	0.552±0.020	0.760±0.015		0.746±0.012	0.752±0.017	0.673±0.018	0.535±0.015	0.560±0.018	0.778±0.013
MLP		0.696±0.016	0.694±0.021	0.622±0.021	0.543±0.016	0.559±0.020	0.800±0.009		0.687±0.016	0.705±0.011	0.646±0.018	0.548±0.012	0.559±0.017	0.804±0.011
XGB		0.812±0.013	0.752±0.025	0.684±0.019	0.518±0.018	0.526±0.018	0.775±0.010		0.819±0.020	0.765±0.016	0.710±0.014	0.530±0.014	0.537±0.019	0.792±0.011
LGBM		0.799±0.057	0.736±0.062	0.680±0.072	0.540±0.059	0.529±0.021	0.785±0.012		0.822±0.047	0.747±0.055	0.687±0.063	0.515±0.057	0.541±0.016	0.794±0.015
HGNN		0.594±0.060	0.554±0.059	0.541±0.067	0.479±0.058	0.506±0.086	0.522±0.060		0.601±0.062	0.551±0.049	0.552±0.052	0.509±0.051	0.526±0.027	0.534±0.043
HGCN		0.729±0.065	0.567±0.072	0.612±0.056	0.506±0.059	0.490±0.067	0.561±0.060		0.720±0.072	0.553±0.070	0.615±0.065	0.506±0.059	0.484±0.059	0.566±0.065
UGCNII		0.680±0.060	0.674±0.048	0.641±0.097	0.488±0.049	0.489±0.068	0.607±0.092		0.663±0.058	0.663±0.044	0.647±0.056	0.502±0.045	0.503±0.035	0.642±0.043
Max		0.826±0.015	0.784±0.020	0.703±0.011	0.550±0.018	0.559±0.020	0.800±0.009		0.834±0.011	0.794±0.018	0.730±0.013	0.548±0.012	0.560±0.018	0.804±0.011
Avg.	0.728±0.075	0.687±0.080	0.634±0.052	0.518±0.024	0.520±0.024	0.697±0.113	0.727±0.079	0.690±0.086	0.649±0.053	0.520±0.017	0.528±0.025	0.709±0.113		
Rank Avg.	1.200±0.422	2.700±0.483	3.500±0.707	5.600±0.516	5.400±0.516	2.600±1.265	1.600±0.966	2.600±0.699	3.400±0.843	5.900±0.316	5.100±0.316	2.400±1.265		
LR	EN	0.883±0.014	0.826±0.015	0.783±0.016	0.627±0.017	0.480±0.020	0.634±0.023	EU	0.933±0.002	0.838±0.004	0.876±0.003	0.691±0.004	0.494±0.003	0.722±0.002
RF		0.880±0.016	0.856±0.010	0.773±0.021	0.684±0.024	0.624±0.029	0.623±0.024		0.960±0.002	0.921±0.003	0.901±0.003	0.737±0.003	0.529±0.004	0.770±0.004
DT		0.707±0.010	0.690±0.019	0.652±0.019	0.551±0.022	0.529±0.017	0.542±0.021		0.852±0.003	0.762±0.004	0.785±0.005	0.546±0.004	0.504±0.007	0.564±0.007
KNN		0.846±0.013	0.810±0.015	0.745±0.020	0.685±0.019	0.597±0.026	0.591±0.024		0.920±0.002	0.854±0.005	0.881±0.003	0.701±0.003	0.586±0.019	0.749±0.003
MLP		0.883±0.013	0.825±0.017	0.780±0.013	0.697±0.019	0.618±0.024	0.636±0.022		0.962±0.001	0.909±0.004	0.911±0.004	0.790±0.005	0.539±0.015	0.802±0.003
XGB		0.863±0.017	0.847±0.014	0.765±0.022	0.666±0.022	0.632±0.025	0.610±0.019		0.961±0.001	0.917±0.002	0.905±0.003	0.747±0.005	0.557±0.011	0.777±0.002
LGBM		0.842±0.056	0.847±0.060	0.769±0.063	0.655±0.081	0.635±0.027	0.612±0.027		0.963±0.005	0.923±0.006	0.908±0.008	0.761±0.014	0.550±0.007	0.796±0.002
HGNN		0.505±0.052	0.543±0.070	0.554±0.070	0.532±0.070	0.543±0.062	0.532±0.080		0.529±0.028	0.520±0.020	0.516±0.017	0.504±0.018	0.502±0.014	0.500±0.017
HGCN		0.804±0.102	0.738±0.099	0.737±0.073	0.548±0.075	0.619±0.084	0.570±0.073		0.849±0.054	0.724±0.051	0.888±0.041	0.550±0.029	0.550±0.057	0.614±0.070
UGCNII		0.787±0.048	0.767±0.065	0.788±0.052	0.706±0.045	0.739±0.056	0.606±0.059		0.874±0.010	0.805±0.010	0.912±0.007	0.793±0.014	0.812±0.020	0.784±0.012
Max		0.883±0.014	0.856±0.010	0.788±0.052	0.706±0.045	0.739±0.056	0.636±0.022		0.963±0.005	0.923±0.006	0.788±0.052	0.706±0.045	0.739±0.056	0.802±0.003
Avg.	0.800±0.117	0.775±0.098	0.738±0.076	0.635±0.067	0.602±0.071	0.596±0.037	0.880±0.132	0.817±0.126	0.848±0.123	0.679±0.113	0.563±0.092	0.708±0.108		
Rank Avg.	1.700±1.567	2.200±0.632	2.500±0.850	4.500±0.707	4.800±1.229	5.300±0.823	1.200±0.422	2.700±0.675	2.200±0.789	5.000±0.471	5.500±0.972	4.400±0.843		
LR	CD	0.969±0.002	0.857±0.003	0.644±0.003	0.564±0.004	0.512±0.003	0.653±0.004	CS	0.980±0.001	0.890±0.002	0.722±0.004	0.584±0.002	0.516±0.002	0.688±0.002
RF		0.997±0.000	0.939±0.001	0.703±0.007	0.573±0.006	0.498±0.004	0.707±0.005		0.999±0.000	0.945±0.003	0.777±0.004	0.611±0.004	0.502±0.005	0.739±0.004
DT		0.963±0.001	0.777±0.007	0.583±0.006	0.511±0.004	0.497±0.003	0.539±0.004		0.974±0.001	0.783±0.011	0.623±0.004	0.519±0.003	0.498±0.004	0.548±0.002
KNN		0.962±0.002	0.857±0.004	0.629±0.008	0.595±0.005	0.520±0.010	0.633±0.005		0.974±0.001	0.899±0.002	0.702±0.004	0.632±0.003	0.531±0.011	0.659±0.002
MLP		0.990±0.001	0.914±0.007	0.693±0.009	0.597±0.006	0.510±0.005	0.776±0.003		0.996±0.000	0.930±0.004	0.775±0.005	0.671±0.006	0.544±0.014	0.806±0.002
XGB		0.997±0.000	0.937±0.002	0.703±0.008	0.579±0.004	0.507±0.007	0.747±0.004		0.999±0.000	0.939±0.005	0.781±0.004	0.620±0.004	0.522±0.007	0.786±0.002
LGBM		0.998±0.001	0.941±0.008	0.719±0.021	0.587±0.014	0.511±0.004	0.782±0.003		0.999±0.000	0.946±0.017	0.792±0.012	0.626±0.019	0.525±0.005	0.807±0.001
HGNN		0.629±0.008	0.528±0.017	0.523±0.013	0.537±0.014	0.540±0.022	0.510±0.013		0.587±0.006	0.509±0.023	0.510±0.012	0.562±0.009	0.590±0.012	0.508±0.007
HGCN		0.861±0.059	0.700±0.053	0.637±0.060	0.516±0.028	0.513±0.021	0.505±0.011		0.852±0.044	0.718±0.050	0.714±0.054	0.513±0.025	0.511±0.024	0.505±0.013
UGCNII		0.975±0.005	0.874±0.010	0.714±0.013	0.851±0.016	0.672±0.042	0.547±0.012		0.971±0.003	0.812±0.010	0.792±0.011	0.899±0.016	0.895±0.007	0.609±0.015
Max		0.998±0.001	0.941±0.008	0.719±0.021	0.851±0.016	0.672±0.042	0.782±0.003		0.999±0.000	0.946±0.017	0.792±0.012	0.899±0.016	0.895±0.007	0.807±0.001
Avg.	0.934±0.115	0.832±0.132	0.655±0.064	0.591±0.096	0.528±0.052	0.640±0.110	0.933±0.129	0.837±0.139	0.719±0.091	0.624±0.109	0.563±0.119	0.666±0.119		
Rank Avg.	1.000±0.000	2.200±0.632	3.900±0.568	4.500±0.850	5.400±1.265	4.000±1.414	1.100±0.316	2.500±1.080	3.600±0.699	4.400±1.075	5.100±1.729	4.300±1.252		
LR	QM	0.652±0.003	0.620±0.004	0.580±0.004	0.499±0.003	0.514±0.003	0.600±0.002	QS	0.598±0.001	0.553±0.002	0.558±0.003	0.512±0.001	0.528±0.002	0.586±0.002
RF		0.734±0.003	0.657±0.006	0.663±0.006	0.505±0.003	0.504±0.004	0.767±0.005		0.728±0.001	0.595±0.002	0.637±0.003	0.501±0.003	0.504±0.003	0.748±0.002
DT		0.621±0.003	0.547±0.006	0.569±0.004	0.502±0.001	0.502±0.004	0.547±0.004		0.630±0.002	0.524±0.002	0.571±0.003	0.500±0.002	0.501±0.002	0.545±0.002
KNN		0.638±0.002	0.601±0.004	0.570±0.003	0.506±0.003	0.511±0.006	0.561±0.002		0.669±0.002	0.552±0.003	0.599±0.002	0.505±0.003	0.503±0.004	0.610±0.002
MLP		0.737±0.004	0.649±0.005	0.634±0.003	0.514±0.004	0.509±0.002	0.834±0.004		0.717±0.002	0.573±0.002	0.630±0.002	0.511±0.002	0.525±0.004	0.815±0.001
XGB		0.744±0.003	0.660±0.005	0.677±0.005	0.504±0.003	0.513±0.004	0.823±0.002		0.745±0.001	0.612±0.002	0.653±0.003	0.504±0.002	0.507±0.002	0.801±0.002
LGBM		0.755±0.010	0.679±0.016	0.694±0.015	0.505±0.014	0.513±0.003	0.844±0.002		0.753±0.005	0.628±0.014	0.669±0.018	0.506±0.009	0.513±0.003	0.813±0.002
HGNN		0.570±0.013	0.520±0.011	0.535±0.012	0.519±0.012	0.566±0.009	0.551±0.010		0.619±0.008	0.543±0.006	0.563±0.006	0.534±0.008	0.599±0.007	0.585±0.009
HGCN		0.538±0.031	0.507±0.010	0.545±0.044	0.501±0.005	0.507±0.020	0.521±0.028		0.593±0.041	0.528±0.025	0.546±0.049	0.502±0.007	0.503±0.012	0.511±0.018
UGCNII		0.658±0.007	0.620±0.011	0.642±0.011	0.661±0.012	0.689±0.007	0.817±0.014		0.713±0.005	0.636±0.008	0.649±0.007	0.594±0.010	0.713±0.006	0.824±0.008
Max		0.755±0.010	0.679±0.016	0.694±0.015	0.661±0.012	0.689±0.007	0.844±0.002		0.753±0.005	0.636±0.008	0.669±0.018	0.594±0.010	0.713±0.006	0.824±0.008
Avg.	0.665±0.076	0.606±0.061	0.611±0.058	0.522±0.049	0.533±0.058	0.687±0.140	0.677±0.062	0.574±0.041	0.608±0.045	0.517±0.029	0.540±0.068	0.684±0.127		
Rank Avg.	1.800±0.919	3.800±1.317	3.200±1.135	5.500±0.972	4.500±1.434	2.200±1.317	1.700±0.949	4.100±0.568	2.900±0.568	5.900±0.316	4.500±1.354	1.900±1.101		
LR	B4	0.693±0.003	O.O.T.*	0.616±0.001	0.569±0.001	0.640±0.000	0.696±0.002	B5	0.696±0.002	O.O.T.*	0.612±0.003	0.592±0.001		0.627±0.001
RF		0.976±0.001		0.799±0.009	0.562±0.003	0.729±0.008	0.977±0.001			0.758±0.004		0.715±0.005		0.748±0.002
DT		0.900±0.002		0.704±0.008	0.507±0.003	0.549±0.4								