

KAISTAI Kim Jaechul Graduate School

Are Weights Useful in Graph Summarization? – A Comparative Study

Shinhwan Kang Kyuhan Lee Kijung Shin

KAIST AI





Graphs are everywhere!

• Graphs represent relationships such as

Purchase history

- Friends in social networks
- Purchase history

Social network

Hyperlinks between web pages



FAQ



Graphs become large!

• Graphs grow rapidly at an unprecedented pace





Graphs become large!

• Graphs grow rapidly at an unprecedented pace





How do we efficiently utilize such large graphs?



Graph compression

- Useful for *efficient utilization of large graphs*
- To find a compact representation that exactly or approximately describes an input graph





Graph summarization: overview

- Promising graph compression technique
- Summary graph is in the form of a graph
 - Can process graph queries directly without restoration
 - Can apply other graph compression techniques





Graph summarization: overview

 We can categorize graph summarization methods *depending* on loss of information during summarization





Graph summarization: overview

- We can categorize graph summarization methods depending on loss of information during summarization
- We focus on lossy graph summarization





Graph summarization models

- We also can categorize graph summarization methods *depending* on summarization models
- One of the two representative models allows edge weights in summary graphs; the other does not

Weighted Unweighted



Lossy graph summarization algorithm



Weighted graph summarization model

- Supernode: a node in a summary graph
- Superedge: an edge in a summary graph





Weighted graph summarization model

 Summary graphs with edge weights contain information about the number of edges on each superedge





Unweighted graph summarization model

 Summary graphs without edge weights *do not retain* information about the number of edges on each superedge





Which one is better between two models?

 There was no systematic comparison between two extensivelystudied graph summarization models



Weighted graph summarization model



Unweighted graph summarization model



Which one is better between two models?

- There was no systematic comparison between two extensivelystudied graph summarization models
- We conduct *a systematic comparison* in five aspects
 - For example, reconstruction and compression ratios
- To this end, we extend three algorithms to both models



Which one is better?



Unweighted graph summarization model

Weighted graph summarization model



Road map

- Introduction
- Notations & Problem formulation <
- Considered algorithms
- Experiments & Theoretical analysis
- Conclusion



Notations: input graph

- Input graph G = (V, E)
 - Set of *subnodes* $V = \{1, 2, ..., |V|\}$
 - Set of *subedges* $E \subseteq \binom{V}{2}$

Input graph G = (V, E)







Notations: input graph

- Input graph G = (V, E)
 - Set of subnodes $V = \{1, 2, ..., |V|\}$
 - Set of subedges $E \subseteq \binom{V}{2}$
- Adjacency matrix $A \in \mathbb{R}^{|V| \times |V|}$ of G
 - $A_{ij} = 1$ if $\{i, j\} \in E$ and $A_{ij} = 0$ otherwise

Input graph G = (V, E)



		d	D	J	u	е
	а	0	1	0	1	0
	b	1	0	1	0	0
	С	0	1	0	1	0
	d	1	0	1	0	1
	е	0	0	0	1	0

hada

←Adjacency matrix A



Notations: summary graph

- Summary graph G' = (S, P) of G = (V, E)
 - Set of supernodes S is a partition of V
 - Set of *superedges* $P \subseteq \binom{S}{2}$





(Un)weighted summary graph

• Summary graph G' is either **weighted** or **unweighted**



Superedge weight function

- Summary graph G' is either weighted or unweighted
- Weighted summary graph additionally has a superedge weight function ω
 - $\hfill \omega$ takes each superedge and returns the its weights

$$\blacktriangleright$$
(e.g.,) $\omega_{AB} = 4$, $\omega_{BC} = 1$





(Un)weighted reconstructed graph

- **Reconstructed** graph \hat{G} is obtained from a summary graph G'
- \widehat{A} denotes an *adjacency matrix* of a reconstructed graph \widehat{G}





Reconstructed adjacency matrix

• With *unweighted* summary graphs, entries of \hat{A} are defined as



Reconstructed adjacency matrix

• With *unweighted* summary graphs, entries of \hat{A} are defined as



Reconstructed adjacency matrix

d

• With weighted summary graphs, entries of \hat{A} are defined as

$$\hat{A}_{ij} = \begin{cases} \omega_{S_i S_j} \\ \pi_{S_i S_j} \\ 0, & \text{otherwise} \end{cases} \quad \begin{bmatrix} \pi_{AB} : \# \text{ of possible} \\ \text{subedges between} \\ \text{subedges between} \\ \text{supernodes } A \text{ and } B \\ \hline \\ \mathbf{Weighted} \\ \text{summary graph } G' \\ \mathbf{M}_{a} \\ \mathbf{M}_{b} \\ \mathbf{M}_{b} \\ \mathbf{C} \\ \mathbf{Reconstruct} \\ \mathbf{Reconstruct} \\ \mathbf{M}_{a} \\ \mathbf{M}_{b} \\$$

С

d

e

0

1

0

1

0

0.5

1

0

0.5

0

0.5

0

0

1

0



Reconstructed adjacency matrix

• With weighted summary graphs, entries of \hat{A} are defined as





Optimization problem formulation

- *Given*: input graph *G* and a size budget *k*
- Find: summary graph G'
 - G': (Un)weighted summary graph
- To minimize: the L_p reconstruction error $||A \hat{A}||_p$
- Subject to: the size of summary graph G', Size(G') ≤ k
 (e.g.,) # of supernodes in G', # of bits to encode G'



• The size of an unweighted summary graph in bits is defined as

$$Size_{bits}(G') = 2|P| \log_2 |S| + |V| \log_2 |S|$$





• The size of an unweighted summary graph in bits is defined as

 $Size_{bits}(G') = 2|P| \log_2 |S| + |V| \log_2 |S|$

• 2 | *P* | *log*₂ | *S* | corresponds to | *P* | *superedges in bits*





• The size of an unweighted summary graph in bits is defined as

 $Size_{bits}(G') = 2|P| \log_2 |S| + |V| \log_2 |S|$

- $2|P|\log_2 |S|$ corresponds to |P| superedges in bits
- |V| log₂ |S| corresponds to supernodes membership of |V| subnodes in bits





• The size of a *weighted* summary graph in bits is defined as

 $Size_{bits}(G') = 2|P|\log_2|S| + |V|\log_2|S| + |P|\log_2\omega_{max}$

- $2|P|\log_2 |S|$ corresponds to |P| superedges in bits
- $|V| \log_2 |S|$ corresponds to supernodes membership of |V| subnodes in bits
- For a weighted summary graph, |P| log₂ ω_{max} corresponds to |P| superedge weights in bits



Road map

- Introduction
- Notations & Problem formulation
- Considered algorithms <<
- Experiments & Theoretical analysis
- Conclusion



Optimization algorithms

- We introduce three optimization algorithms (W) for finding a weighted summary graph
 - k-Grass (W) [1], SSumM (W) [2], MoSSo-Lossy (W) [3]





Extending optimization algorithms

- We introduce three optimization algorithms (W) for finding a weighted summary graph
 - k-Grass (W) [1], SSumM (W) [2], MoSSo-Lossy (W) [3]
- We extend each algorithm to provide an unweighted summary graph (U) by modifying their objective function

Algorithm	Algorithm	
k-Grass (W)	k-Grass (U)	
SSumM (W)	<i>SSumM</i> (U)	
MoSSo-Lossy (W)	MoSSo-Lossy (U)	



Extending optimization algorithms

- We introduce three optimization algorithms (W) for finding a weighted summary graph
 - k-Grass (W) [1], SSumM (W) [2], MoSSo-Lossy (W) [3]
- We extend each algorithm to provide an unweighted summary graph (U) by modifying their objective function

Algorithm	Algorithm
k-Grass (W)	k-Grass (U)
SSumM (W)	SSumM (U)
MoSSo-Lossy (W)	MoSSo-Lossy (U)





Inputs & output

• Given an input graph *G*, *k*-Grass produces a summary graph whose size is smaller than the size budget

Input: (1) *input graph G*, and (2) *size budget k* Output: *summary graph G*'

- 1. initialize G'
- 2. while Size(G') > k do
- merge a supernode pair {A, B} whose merger increases Loss() least
 return G'



Initialization

• It *first initializes the set of supernodes* so that each subnode forms a singleton supernode

Input: (1) input graph G, and (2) size budget kOutput: summary graph G'

- 1. *initialize* G'
- 2. while Size(G') > k do
- merge a supernode pair {A, B} whose merger increases Loss() least
 return G'



Merging

 After that, it greedily merges a supernode pair whose merger increases the objective function Loss() least

Input: (1) input graph G, and (2) size budget kOutput: summary graph G'

- 1. initialize G'
- 2. while Size(G') > k do

merge a supernode pair {A, B} whose merger increases Loss() least
 return G'



Loss function of k-Grass

• Loss function of k-Grass (U, W) is the L_p reconstruction error

 $\left\| A - \widehat{A} \right\|_{n}$

where \hat{A} is an adjacency matrix of a reconstructed graph

 Note that the adjacency matrix is *reconstructed from an unweighted summary graph (U) or from a weighted summary graph (U)*





Termination

 If the size of summary graph Size(G') is smaller than the budget, it returns a summary graph

Input: (1) input graph G, and (2) size budget kOutput: summary graph G'

- 1. initialize G'
- 2. while *Size(G') > k* do

merge a supernode pair {A, B} whose merger increases Loss() least
 return G'



Road map

- Introduction
- Notations & Problem formulation
- Considered algorithms
- Experiments & Theoretical analysis <<
- Conclusion



Experiments: settings

- Datasets
 - 8 Real-world graphs (0.2M 0.2B edges)



- Considered graph summarization algorithms
 - {k-Grass, SSumM, MoSSo-Lossy} X {Weighted (W), Unweighted (U)}



• Compression ratio $\frac{Size_{bits}(G')}{2|E|\log_2|V|}$ where G = (V, E) is the input graph and G' is a summary graph





• Compression ratio $\frac{Size_{bits}(G')}{2|E|\log_2|V|}$ where G = (V, E) is the input graph and G' is a summary graph

• Quality: L_1/L_2 reconstruction error $||A - A'||_{p=\{1,2\}}$



- Compression ratio $\begin{array}{l} Size_{bits}(G') \\ \hline 2|E|\log_2|V| \\ \end{array}$ where G = (V, E) is the input graph and G' is a summary graph
- Quality: L_1/L_2 reconstruction error $||A A'||_{p=\{1,2\}}$
- Node importance: PageRank (PR) [4]
 - damping factor = 0.85



- Compression ratio $\begin{array}{l} Size_{bits}(G')\\ \hline 2|E|\log_2|V|\\ \end{array}$ where G=(V,E) is the input graph and G' is a summary graph
- Quality: L_1/L_2 reconstruction error $||A A'||_{p=\{1,2\}}$
- Node importance: PageRank (PR) [4]
 - damping factor = 0.85
- Node proximity: Random Walk with Restart (RWR) [5]
 - damping factor = 0.95 with 100 randomly-sampled seeds



0.4 0.6 0.8

LIVEJOURNAL

Compression Ratio

Experiments: quality of summary graphs

0.2 0.4 0.6 0.8

Compression Ratio

DBpedia

 Unweighted summary graphs described the input graph up to 8.2X more accurately than weighted ones



0.2 0.4 0.6 0.8

Compression Ratio

MoSSo-Lossy (W) ▼MoSSo-Lossy (U) SSumM (W) **▼**SSumM (U)



10

0.0.

8.2X

0.2 0.4 0.6 0.8



Experiments: maintain node importance

 Unweighted summary graphs *maintained the node importance up to 7.8X more accurate* than weighted ones





Experiments: preserve node proximity

• The answers from unweighted summary graphs are up to 5.9X more accurate than from weighted summary graphs





Why can edge weights be harmful?

- Consider a graph G and its weighted summary (S, P, ω)
- Assume ω is not fixed but variable
- Theorem 1. When the L1 reconstruction error is minimized, the weight of each superedge {A, B} ∈ P is set so that the weights of all reconstructed subedges are either 1 or 0, as if an unweighted summary graph is used
- That is, edge weights do not contribute to accuracy while requiring additional space



Road map

- Introduction
- Notations & Problem formulation
- Considered algorithms
- Experiments & Theoretical analysis
- Conclusion <<



Conclusion

- We *conducted a systematic comparison* between two graph summarization models with and without superedge weights
- We empirically *revealed a surprising finding* that removing superedge weights leads to significant improvements
- We *developed a theoretical analysis* to shed light on this counterintuitive observations (*Theorem 1*)

Github link: <u>https://github.com/ShinhwanKang/PAKDD22-ComparativeStudy</u>



References

[1] K. LeFevre and E. Terzi, "Grass: Graph structure summarization," in SDM, 2010

[2] K. Lee, H. Jo, J. Ko, S. Lim, and K. Shin, "Ssumm: Sparse summarization of massive graphs," in KDD, 2020

[3] Ko, Jihoon, Yunbum Kook, and Kijung Shin. "Incremental lossless graph summarization," in KDD, 2020.

[4] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. Technical report, Stanford InfoLab, 1999.

[5] H. Tong, C. Faloutsos, and J.-Y. Pan, "Random walk with restart: fast solutions and applications," KAIS, vol. 14, no. 3, pp. 327–346, 2008.



KAISTAI Kim Jaechul Graduate School

Are Weights Useful in Graph Summarization? – A Comparative Study

Shinhwan Kang Kyuhan Lee Kijung Shin

KAIST AI

