Beyond Neighbors: Distance-Generalized Graphlets for Enhanced Graph Characterization



Yeongho Kim, Yuyeong Kim, Geon Lee, Kijung Shin {yeongho, youyoung, geonlee0325, kijungs} @ kaist.ac.kr https://github.com/thisis05/EDGE



1. What is Graphlet?

[Graphlet]

- Recurring small subgraph patterns in graphs with N Nodes.
- A Key to analyze local structural patterns in complex networks.

[Examples]



3-size Graphlets

4-size Graphlets

[How to use Graphlets]

2. How to Count Graphlets ?





3-size Graphlet

 G_1

Base Approach: Enumerate all Connected Subgraphs

 $\{A, B, C\} \longrightarrow G_1$ $\{A, B, D\} \longrightarrow G_2 \qquad \{B, C, D\} \longrightarrow G_2 \\ \{B, C, E\} \longrightarrow G_2 \\ \{B, C, E\} \longrightarrow G_2 \end{cases}$ $\{A, C, E\} \longrightarrow G_2$



Efficient Approach: Combinatorial Techniques

Deduce the count of some graphlets based on the counts of others.

(1) Get the count of G_1 with enumeration.

(2) Get the count of G_2 with a pre-defined combinatorial equation.

Count each graphlet instances in graphs and normalize the counts.

• Use the normalized values in downstream tasks. (e.g. Graph Characterization, Graph Clustering)

$\bullet \bullet \bullet \bullet$ Count(G_2) = $\sum_{u \in V} {\binom{|N_u|}{2}} - 3 \times \text{Count}(G_1) = 7 - 3 \times 1 = 4$ G_2

 $* N_u$: The set of neighbor nodes connected to node u

3. Limitation and Motivation to Distance-Generalized Graphlets

[Limitations of Original Graphlets]

 Graphlets capture only direct connections between nodes. Original graphlets miss some local structures.

[Need for More Expressive Graphlets]

• We propose (*d*, *s*)-graphlets, which incorporate indirect **connections up to distance** *d* for deeper analysis.





Size-4 Graphlets

Graphlets cannot distinguish & *identify* these local structures...



(2, 4)-Graphlets

Considering distance up to 2 makes them distinguishable and identifiable!

4. Proposed Concept: (*d*, *s*)-Graphlets

[Definition]

(1) *d*-edges: Pairs of nodes whose distance is exactly *d*.

(2) *d*-induced subgraphs: Subgraphs induced by a set of nodes that consist of 1-edge, 2-edge, ..., *d*-edges (up to *d*-edges).

(3) (*d*, *s*)-graphlets: *d*-induced subgraphs patterns with *s* nodes.

(3, 3)-Graphlets (2, 3)-Graphlets

5. Proposed Counting Algorithm: EDGE

EDGE: A New Counting Algorithm for (*d*, *s*)-Graphlets Using Combinatorial Techniques

[1. Preprocessing]

Construct *d*-edges to capture both direct and indirect (up to distance *d*) connections.





• Size-3 graphlets T(d): 6 types for (2, 3)-graphlets / 13 types for (3, 3)-graphlets.

• Size-4 graphlets Q(d): 36 types for (2, 4)-graphlets.

6. Experimental Results



Require explicit enumeration Partially require enumeration

No enumeration! Only need combinatorial equations

[2. Enumeration : Non-deducible / Semi-deducible (*d*, *s*)-Graphlets]

Count non-deducible graphlets with enumeration and save the counts.

• Enumerate to obtain the candidate counts of semi-deducible graphlets and save the counts.

[3. Combinatorial Equations : Compute the counts of Deducible (d, s)-Graphlets]

• Apply pre-defined combinatorial equations to obtain the counts of deducible graphlets.

Example Count(
$$\bullet_{T_5^{(2)}}$$
) = $\sum_{u \in V} \begin{pmatrix} |N_u^{(2)}| \\ 2 \end{pmatrix}$ - 3×Count($\bullet_{T_1^{(2)}}$) - Count($\bullet_{T_2^{(2)}}$)

 $* N_{u}^{(2)}$: The set of nodes connected to node u by a 2-edge



[R3. Speed and Scalability]



• **Domain-based similarity is much higher** for (*d*, *s*)-graphlets than for original graphlets.

(*d*, *s*)-graphlets accurately characterize real-world graphs.

EDGE achieves up to 14.86× speedup over

EDGE is competitive with conventional graphlet counting methods in **runtime per graphlet instance**.