

EmbAssi: Embedding Assignment Costs for Similarity Search in Large Graph Databases

Franka Bause¹, Erich Schubert², Nils M. Kriege¹

¹University of Vienna ²TU Dortmund University

Searching in Large Graph Databases

Goal: Find graphs similar to a query graph efficiently



Graph Edit Distance (GED)

- **Distance measure** for graphs
- Transform one graph into the other (edit path)
- **Cost of transformation** $\hat{=}$ **GED**



► **NP-hard** problem



Cost of Edit Operations:

Notation

 C_{v} : delete/insert vertex C_e : delete/insert edge C_{vl} : change label of vertex *C_{el}* : change label of edge

Our Contribution



Tight lower bounds for efficient filtering

- **Efficient** lower bounds for the graph edit distance based on **tree metrics**
- Embeddable to $\ell_1 \rightarrow index$ for vectors can be used
- Suitable for similarity search
- Lower bound for generating initial candidates

Lower Bounds Based on Optimal Assignments

- Compute optimal assignment between vertices
- \blacktriangleright Cost of assignment $\hat{=}$ lower bound for GED
- Embeddable to ℓ_1 if cost function is **tree metric**

Embedding Assignment Costs

Cost function has to be a tree metric

Representable as tree T with cost $\hat{=}$ length of path

Label Lower Bound (LLB)

- Minimum vertex relabelings/deletions
- Edge weights:

$$w_1 = c_v - 0.5 \cdot c_{vl}$$
 $w_2 = 0.5 \cdot c_{vl}$



Degree Lower Bound (DLB)

- Minimum edge insertions/deletions based on vertex degrees
- Edge weights: $w_3 = 0.5 \cdot c_e$
- Each operation influences degree of two vertices









 $4 \cdot w_3$



 $\phi(G_1)$



d





- \blacktriangleright $S_{\overline{uv}}: \#s \in S$ associated with nodes in subtree containing u when deleting *uv*
- Cost of optimal assignment: $\sum_{uv \in E(T)} |A_{\overleftarrow{uv}} B_{\overleftarrow{uv}}| \cdot w(uv)$
- $\blacktriangleright \phi(S) = [S_{\overleftarrow{uv}} \cdot w(uv)]_{uv \in E(T)}$

 \bigcirc

$2 \cdot w_2$ $1 \cdot w_2$

 $| \cdot w_2 \rangle$



Combined Lower Bound (CLB)

- **Combines** *LLB* and *DLB*
- $d_{CLB}(G_1, G_2) = d_{LLB}(G_1, G_2) + d_{DLB}(G_1, G_2)$

Comparison of Approximation Quality

 \mathbf{CLB} DLB SLF BLP LLB Δ Δ **×** CStarLB**•** CStarUB**•** CStarUBRef **×** BranchLB**•** BranchUB * BLPlb • LinD • BeamD

su

 $B \bigcirc$

3

4



Comparison to State-of-the-Art

 $\phi(G_2)$

 $0 \cdot w_1$

 $1 \cdot w_2$

 $2 \cdot w_2$



 \triangle : our approximation x: other lower bound \bigcirc : upper bound

Conclusion

- Efficiently computable tight lower bounds
- **Scalable** to databases with millions of graphs
- Optimization of runtime (choice of edge direction)
- ► *k*NN and range queries both supported
- **Speed up state-of-the-art** approaches through pre-filtering

European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases

September 19 – 23, 2022