

Learning Graph Representations of Biochemical Networks and its Application to Enzymatic Link Prediction



Julie Jiang¹, Li-Ping Liu¹, and Soha Hassoun^{1,2}

¹Department of Computer Science, Tufts University, Medford, MA

²Department of Chemical and Biological Engineering, Tufts University, Medford, MA



Motivation and Contribution

Problem: No complete characterization of enzymatic reactions

The curation of enzyme functions and the reactions they catalyze remains elusive, hindering biological engineering and discovery.

Goal: Predict new enzymatic transformation of molecules

- Exploit existing biochemical databases (e.g. KEGG [1]) to better understand their underlying enzymatic transformation relationship
- Enhance biological discovery of undocumented molecular reactions

Contributions: Utilize graph embedding techniques to model molecular reactions

- Apply graph embedding methods to learn latent representations of molecules, capturing both the structural properties of molecules and connectivity among molecules
- Develop a practical and accurate machine learning framework to predict new enzymatic reactions
- Derive meaningful visualizations of pathway metabolites

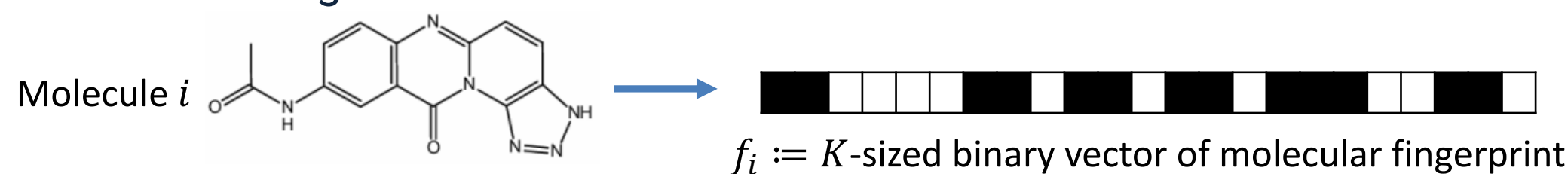
KEGG Database as Graph

The KEGG Dataset

- KEGG is a large database of catalogued enzymatic reactions
- We assume all reactions are reversible, as most are in the database
- We remove cofactors as they are high-connectivity hub nodes

Graph Construction

- Every molecule is a node
- Each substrate-product pair within a reaction is an undirected edge
- Molecular fingerprints (MACCS [2] or PubChem [3]) are used as node attributes
 - A fingerprint is a binary, fixed size vector
 - Each element indicates the presence or absence of pre-defined structural molecular fragments

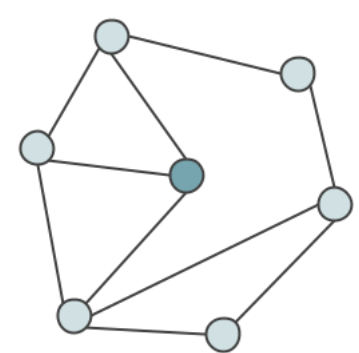


- Enzyme classes are used as edge attributes – enzyme commission (EC) number or KEGG reaction class (RC)

Method: Enzymatic Link Prediction (ELP)

1 Embedding Propagation on Graph

We use Embedding Propagation [4], a graph embedding method, to learn embedding vectors of nodes



All embeddings are randomly initialized:

- Connectivity-based node embeddings $\{u_i\}$,
- Fingerprint embeddings $\{v_k\}$, one for each fingerprint entry
- Enzyme embeddings $\{z_r\}$, one for each enzyme label

- Fingerprint-based node embeddings $\{u_i^{fp}\}$ are constructed from fingerprint embeddings

$$u_i^{fp} = \frac{1}{\sum_{k=1}^K f_{ik}} \sum_{k=1}^K f_{ik} v_k$$

k^{th} value of node i 's fingerprint Fingerprint embedding of entry k

- Reconstruct node embedding (\tilde{u}_i) from the embeddings of its neighbors

$$\tilde{u}_i = \frac{1}{|\mathcal{N}(i)|} \sum_{j \in \mathcal{N}(i)} u_j + \alpha z_{r(i,j)}$$

Neighbors of node i Node embedding of node j Enzyme embedding of the edge (i,j)

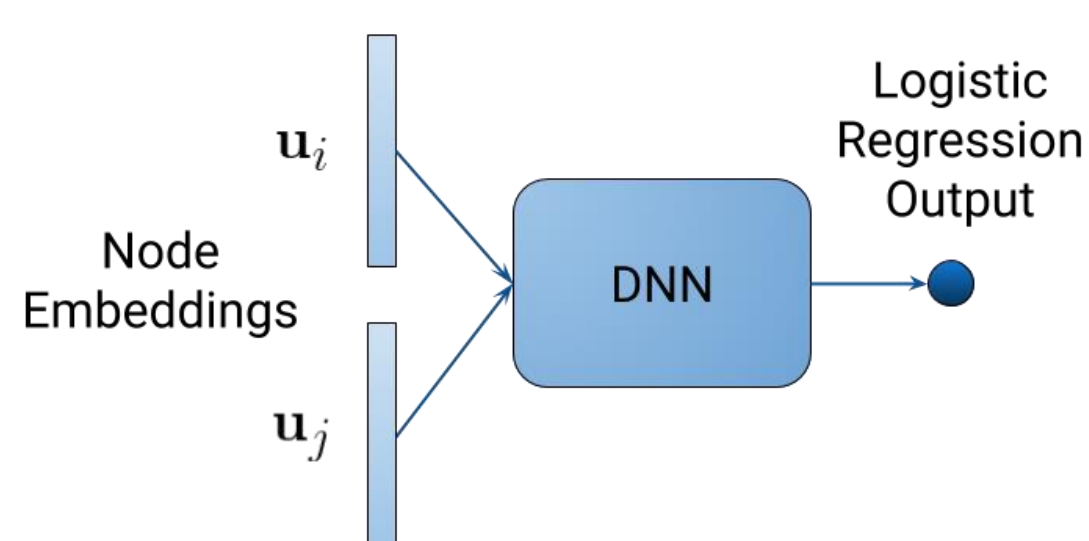
- Margin-based ranking loss.

- Aim to maximize the similarity between the reconstruction of node embedding \tilde{u}_i with node embedding u_i

$$\mathcal{L} = \sum_{i \in V} \sum_{j \in V, j \neq i} \max\{\gamma - \tilde{u}_i^\top u_i + \tilde{u}_i^\top u_j, 0\}$$

Random node j as the negative example for each node in every iteration

- Concatenate u_i and u_i^{fp} to form final node embedding vectors



2 Link Prediction Using Embedding Vectors

We train a logistic regression model using deep neural nets to predict the likelihood of an edge between two nodes

Experiments

Transductive Learning

- Model is trained on all nodes and evaluated for edge recovery for a held out set of test edges.
- The training graph must be connected
- Different combinations of fingerprints and enzyme labels are explored

Inductive Learning

- Model is trained to predict possible interactions for *out-of-sample* nodes excluded from training
- This type of prediction is made possible by only using fingerprint-based node embeddings

Results

Baselines

- Other graph embedding methods (node2vec and Deepwalk), but which do not have a way of utilizing node and edge attributes
- Jaccard similarity score of fingerprints as a way to predict links

The ELP model, usually with MACCS fingerprints, yields the best performance across all test scenarios

Transductive Learning Results

Model		AUC				
Method	Connectivity Embedding	Node Attribute	Edge Attribute	Test Ratios		
				0.1	0.3	0.5
<i>A. Connectivity-based embeddings only</i>						
ELP	Yes	–	–	0.801	0.789	0.761
node2vec	Yes	–	–	0.824	0.736	0.776
DeepWalk	Yes	–	–	0.847	0.763	0.749
<i>B. Connectivity and one additional attribute</i>						
ELP	Yes	MACCS	–	0.953*	0.935*	0.900
ELP	Yes	PubChem	–	0.891	0.882	0.864
ELP	Yes	–	EC	0.795	0.808	0.810
ELP	Yes	–	RC	0.810	0.798	0.810
<i>C. Connectivity with one node and one edge attribute</i>						
ELP	Yes	MACCS	EC	0.941	0.933	0.922*
ELP	Yes	MACCS	RC	0.942	0.929	0.895
ELP	Yes	PubChem	EC	0.892	0.879	0.867
ELP	Yes	PubChem	RC	0.892	0.876	0.859
<i>D. Embedding based on MACCS fingerprints</i>						
ELP	No	MACCS	–	0.931	0.916	0.898
ELP	No	MACCS	EC	0.940	0.925	0.913
ELP	No	MACCS	RC	0.939	0.904	0.896
<i>E. Embeddings based on PubChem fingerprints</i>						
ELP	No	PubChem	–	0.665	0.709	0.682
ELP	No	PubChem	EC	0.745	0.707	0.728
ELP	No	PubChem	RC	0.728	0.706	0.720
<i>F. Jaccard index similarity scoring; no embeddings</i>						
Jaccard	No	MACCS	–	0.808	0.778	0.767
Jaccard	No	PubChem	–	0.542	0.526	0.535

Inductive Learning Results

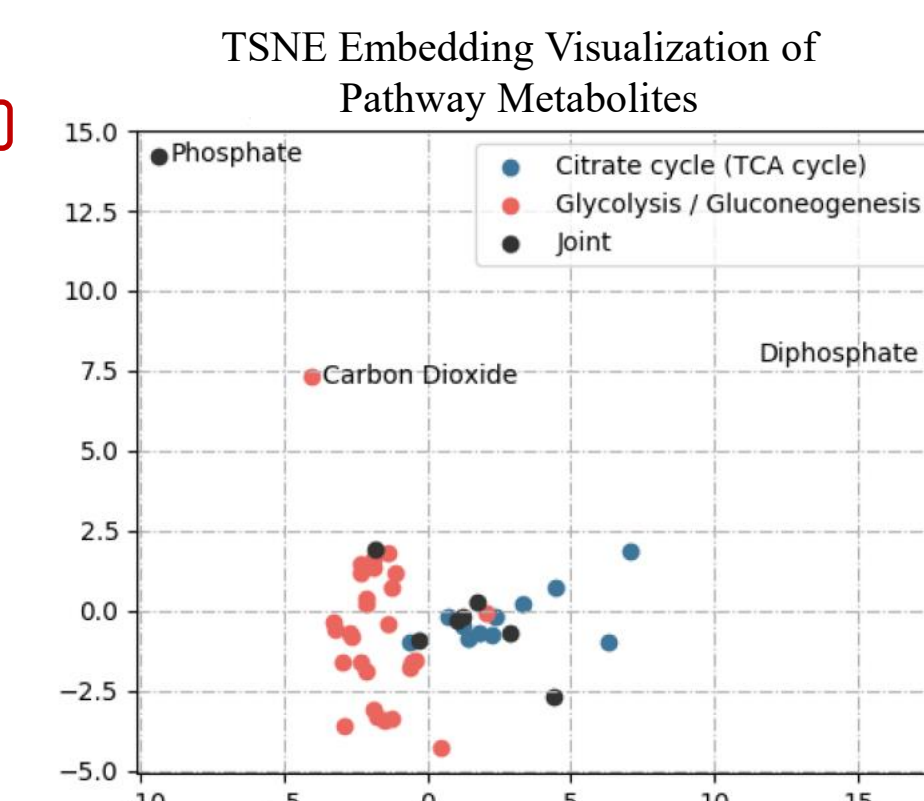
Method	Connectivity Embedding	Node Attribute	AUC
<i>A. Embeddings based on node attributes</i>			
ELP	Yes	MACCS	0.921
ELP	Yes	PubChem	0.605
<i>B. Jaccard index similarity scoring</i>			
Jaccard	No	MACCS	0.744
Jaccard	No	PubChem	0.553

Bold value indicates the best result.

TSNE Embedding Visualization of Pathway Metabolites

TSNE Embedding Visualization of Pathway Metabolites. The plot shows metabolites from three pathways: Citrate cycle (TCA cycle) in blue, Glycolysis / Gluconeogenesis in red, and Joint in black. Metabolites are clustered based on their pathway. Phosphate is at the top left, Carbon Dioxide is in the middle left, and Diphosphate is at the top right. The x-axis ranges from -10 to 15, and the y-axis ranges from -5.0 to 15.0.

Bold values in each partition indicate the best result in that train-test split, and bold values with * indicate the best overall result. The Connectivity Embedding column refers to the use of connectivity-based node embeddings.



Conclusion

This work presents ELP, a framework that learns molecular representations that capture graph connectivity, enzymatic properties, and structural molecular properties

- ELP shows high accuracy in link prediction when using both graph connectivity and molecular attributes
- ELP can be used as a guide to identifying catalyzing enzymes when constructing novel synthesis pathways or predicting interaction between microbes and human hosts
- ELP can enhance link prediction in chemical networks, where previously rule-based and path-based link prediction respectively yielded 52.7% and 67.5% prediction accuracy

References

1. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M., and Tanabe, M. (2015). Kegg as a reference resource for gene and protein annotation. *Nucleic acids research*, **44**(D1), D457–D462.
2. Durant, J. L., Leland, B. A., Henry, D. R., and Nourse, J. G. (2002). Reoptimization of mdl keys for use in drug discovery. *Journal of chemical information and computer sciences*, **42**(6), 1273–1280.
3. Kim, S., Thiessen, P. A., Bolton, E. E., Chen, J., Fu, G., Gindulyte, A., Han, L., He, J., He, S., Shoemaker, B. A., et al. (2015). Pubchem substance and compounddatabases. *Nucleic acids research*, **44**(D1), D1202–D1213.
4. Duran, A. G. and Niepert, M. (2017). Learning graph representations with embedding propagation. In *Advances in neural information processing systems*, pages 5119–5130.

Acknowledgements

This research is supported by NSF under Award Number CCF-1909536 and also by NIGMS of the National Institutes of Health under Award Number R01GM132391. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.