Node Embeddings

This notebook demonstrates different methods for node embeddings and how to further reduce their dimensionality to be able to visualize them in a 2D plot.

Node embeddings are essentially an array of floating point numbers (length = embedding dimension) that can be used as "features" in machine learning. These numbers approximate the relationship and similarity information of each node and can also be seen as a way to encode the topology of the graph.

Considerations

Due to dimensionality reduction some information gets lost, especially when visualizing node embeddings in two dimensions. Nevertheless, it helps to get an intuition on what node embeddings are and how much of the similarity and neighborhood information is retained. The latter can be observed by how well nodes of the same color and therefore same community are placed together and how much bigger nodes with a high centrality score influence them.

If the visualization doesn't show a somehow clear separation between the communities (colors) here are some ideas for tuning:

- Clean the data, e.g. filter out very few nodes with extremely high degree that aren't actually that important
- Try directed vs. undirected projections
- Tune the embedding algorithm, e.g. use a higher dimensionality
- Tune t-SNE that is used to reduce the node embeddings dimension to two dimensions for visualization.

It could also be the case that the node embeddings are good enough and well suited the way they are despite their visualization for the down stream task like node classification or link prediction. In that case it makes sense to see how the whole pipeline performs before tuning the node embeddings in detail.

Note about data dependencies

PageRank centrality and Leiden community are also fetched from the Graph and need to be calculated first. This makes it easier to see if the embeddings approximate the structural information of the graph in the plot. If these properties are missing you will only see black dots all of the same size.

References

igassistant

- Neo4j Python Driver
- Tutorial: Applied Graph Embeddings
- · Visualizing the embeddings in 2D
- Fast Random Projection
- scikit-learn TSNE
- AttributeError: 'list' object has no attribute 'shape'

The scikit-learn version is 1.3.0. The pandas version is 1.5.1.

Dimensionality reduction with t-distributed stochastic neighbor embedding (t-SNE)

The following function takes the original node embeddings with a higher dimensionality, e.g. 64 floating point numbers, and reduces them into a two dimensional array for visualization.

It converts similarities between data points to joint probabilities and tries to minimize the Kullback-Leibler divergence between the joint probabilities of the low-dimensional embedding and the high-dimensional data.

(see https://scikit-

learn.org/stable/modules/generated/sklearn.manifold.TSNE.html#sklearn.manifold.TSNE)

1. Java Packages

1.1 Generate Node Embeddings using Fast Random Projection (Fast RP) for Java Packages

Fast Random Projection is used to reduce the dimensionality of the node feature space while preserving most of the distance information. Nodes with similar neighborhood result in node embedding with similar vectors.

No projected data for node embeddings calculation available

1.2 Dimensionality reduction with t-distributed stochastic neighbor embedding (t-SNE)

This step takes the original node embeddings with a higher dimensionality, e.g. 64 floating point numbers, and reduces them into a two dimensional array for visualization. For more details look up the

function declaration for "prepare node embeddings for 2d visualization".

No projected data for node embeddings dimensionality reduction available

1.3 Visualization of the node embeddings reduced to two dimensions

No projected data to plot available

1.4 Node Embeddings for Java Packages using HashGNN

HashGNN resembles Graph Neural Networks (GNN) but does not include a model or require training. It combines ideas of GNNs and fast randomized algorithms. For more details see HashGNN. Here, the latter 3 steps are combined into one for HashGNN.

No projected data for node embeddings calculation available No projected data for node embeddings dimensionality reduction available No projected data to plot available

2.5 Node Embeddings for Java Packages using node2vec

No projected data for node embeddings calculation available No projected data for node embeddings dimensionality reduction available No projected data to plot available

2. Typescript Modules

2.1 Generate Node Embeddings for Typescript Modules using Fast Random Projection (Fast RP)

See section 1.1 for some background about node embeddings.

No projected data for node embeddings calculation available

2.2 Dimensionality reduction with t-distributed stochastic neighbor embedding (t-SNE)

See section 1.2 for some background about t-SNE.

No projected data for node embeddings dimensionality reduction available

2.3 Plot the node embeddings reduced to two dimensions for Typescript

No projected data to plot available

2.4 Node Embeddings for Typescript Modules using HashGNN

HashGNN resembles Graph Neural Networks (GNN) but does not include a model or require training. It combines ideas of GNNs and fast randomized algorithms. For more details see HashGNN. Here, the latter 3 steps are combined into one for HashGNN.

No projected data for node embeddings calculation available No projected data for node embeddings dimensionality reduction available No projected data to plot available

2.5 Node Embeddings for Typescript Modules using node2vec

No projected data for node embeddings calculation available No projected data for node embeddings dimensionality reduction available No projected data to plot available