

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

- [jqassistant](#)

- [Neo4j Python Driver](#)
- [Tutorial: Applied Graph Embeddings](#)
- [Visualizing the embeddings in 2D](#)
- [scikit-learn TSNE](#)
- [AttributeError: 'list' object has no attribute 'shape'](#)
- [Fast Random Projection \(neo4j\)](#)
- [HashGNN \(neo4j\)](#)
- [node2vec \(neo4j\)](#) computes a vector representation of a node based on second order random walks in the graph.
- [Complete guide to understanding Node2Vec algorithm](#)

The openTSNE version is: 1.0.1
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://opentsne.readthedocs.io>)

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.

👉 **Hint:** To skip existing node embeddings and always calculate them based on the parameters below edit `Node_EMBEDDINGS_0a_Query_Calculated` so that it won't return any results.

The results have been provided by the query filename: ../cypher/Node_EMBEDDINGS/Node_EMBEDDINGS_0a_Query_Calculated.cypher

	codeUnitName	shortCodeUnitName	projectName	communityId	centrality	embedding
0	org.axonframework.disruptor.commandhandling	commandhandling	axon-disruptor-4.10.0	0	0.012594	[0.5492294430732727, 0.027224401012063026, -0....]
1	org.axonframework.commandhandling	commandhandling	axon-messaging-4.10.0	0	0.073080	[0.2871908247470857, 0.24274806678295135, -0....]
2	org.axonframework.commandhandling.callbacks	callbacks	axon-messaging-4.10.0	0	0.015707	[0.3374820053577423, 0.4208607077598572, -0.02...]
3	org.axonframework.commandhandling.distributed	distributed	axon-messaging-4.10.0	0	0.023106	[0.39360445737838745, 0.2639845609664917, -0.0...]
4	org.axonframework.commandhandling.distributed....	commandfilter	axon-messaging-4.10.0	0	0.013920	[0.07936500012874603, 0.439640998840033203, -0....]

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

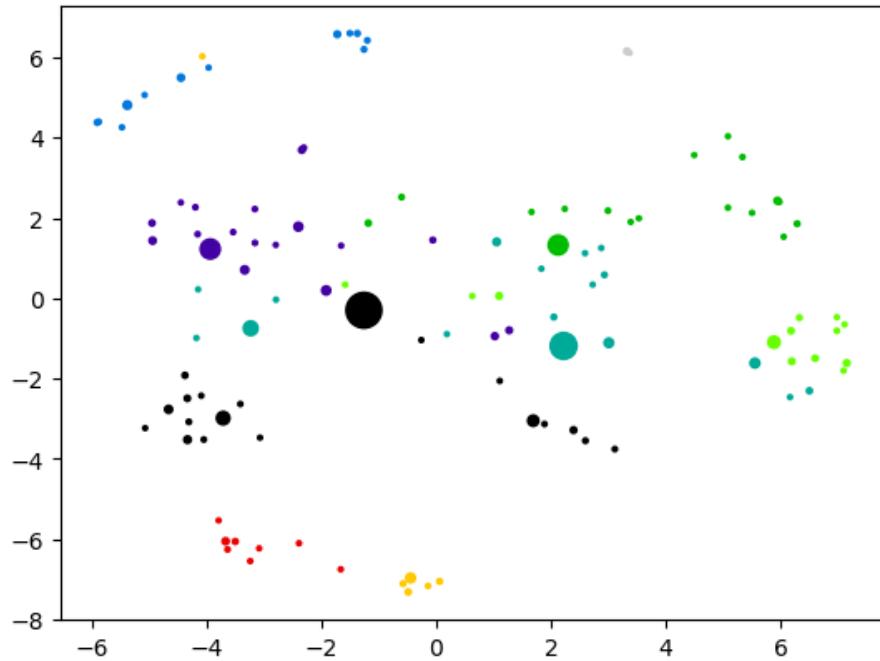
This step takes the original node embeddings with a higher dimensionality, e.g. 64 floating point numbers, and reduces them into a two dimensional array for visualization. For more details look up the function declaration for "prepare_node_embeddings_for_2d_visualization".

```
-----
TSNE(early_exaggeration=12, random_state=47, verbose=1)
-----
====> Finding 90 nearest neighbors using exact search using euclidean distance...
    --> Time elapsed: 0.03 seconds
====> Calculating affinity matrix...
    --> Time elapsed: 0.00 seconds
====> Calculating PCA-based initialization...
    --> Time elapsed: 0.00 seconds
====> Running optimization with exaggeration=12.00, lr=9.50 for 250 iterations...
Iteration 50, KL divergence -0.4140, 50 iterations in 0.0578 sec
Iteration 100, KL divergence 1.2312, 50 iterations in 0.0160 sec
Iteration 150, KL divergence 1.2312, 50 iterations in 0.0147 sec
Iteration 200, KL divergence 1.2312, 50 iterations in 0.0146 sec
Iteration 250, KL divergence 1.2312, 50 iterations in 0.0150 sec
    --> Time elapsed: 0.12 seconds
====> Running optimization with exaggeration=1.00, lr=114.00 for 500 iterations...
Iteration 50, KL divergence 0.1842, 50 iterations in 0.0530 sec
Iteration 100, KL divergence 0.1677, 50 iterations in 0.0438 sec
Iteration 150, KL divergence 0.1543, 50 iterations in 0.0412 sec
Iteration 200, KL divergence 0.1541, 50 iterations in 0.0412 sec
Iteration 250, KL divergence 0.1540, 50 iterations in 0.0407 sec
Iteration 300, KL divergence 0.1539, 50 iterations in 0.0415 sec
Iteration 350, KL divergence 0.1540, 50 iterations in 0.0410 sec
Iteration 400, KL divergence 0.1540, 50 iterations in 0.0411 sec
Iteration 450, KL divergence 0.1540, 50 iterations in 0.0408 sec
Iteration 500, KL divergence 0.1540, 50 iterations in 0.0411 sec
    --> Time elapsed: 0.43 seconds
(114, 2)
```

	codeUnit	artifact	communityId	centrality	x	y
0	org.axonframework.disruptor.commandhandling	axon-disruptor-4.10.0	0	0.012594	-4.100150	-2.427136
1	org.axonframework.commandhandling	axon-messaging-4.10.0	0	0.073080	1.689716	-3.050854
2	org.axonframework.commandhandling.callbacks	axon-messaging-4.10.0	0	0.015707	2.600156	-3.548316
3	org.axonframework.commandhandling.distributed	axon-messaging-4.10.0	0	0.023106	2.393333	-3.284461
4	org.axonframework.commandhandling.distributed....	axon-messaging-4.10.0	0	0.013920	3.112906	-3.757783

1.3 Visualization of the node embeddings reduced to two dimensions

Java Package positioned by their dependency relationships (FastRP node embeddings + t-SNE)



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.

The results have been provided by the query filename: `.../cypher/Node_EMBEDDINGS/Node_EMBEDDINGS_0a_Query_Calculated.cypher`

	codeUnitName	shortCodeUnitName	projectName	communityId	centrality	embedding
0	org.axonframework.disruptor.commandhandling	commandhandling	axon-disruptor-4.10.0	0	0.012594	[0.6495190411806107, 1.5155444294214249, -1.29...
1	org.axonframework.commandhandling	commandhandling	axon-messaging-4.10.0	0	0.073080	[0.21650634706020355, 1.0825317353010178, -1.9...
2	org.axonframework.commandhandling.callbacks	callbacks	axon-messaging-4.10.0	0	0.015707	[-0.4330126941204071, 0.4330126941204071, -1.2...
3	org.axonframework.commandhandling.distributed	distributed	axon-messaging-4.10.0	0	0.023106	[0.21650634706020355, -0.21650634706020355, -1...
4	org.axonframework.commandhandling.distributed....	commandfilter	axon-messaging-4.10.0	0	0.013920	[0.21650634706020355, -1.0825317353010178, -1....]

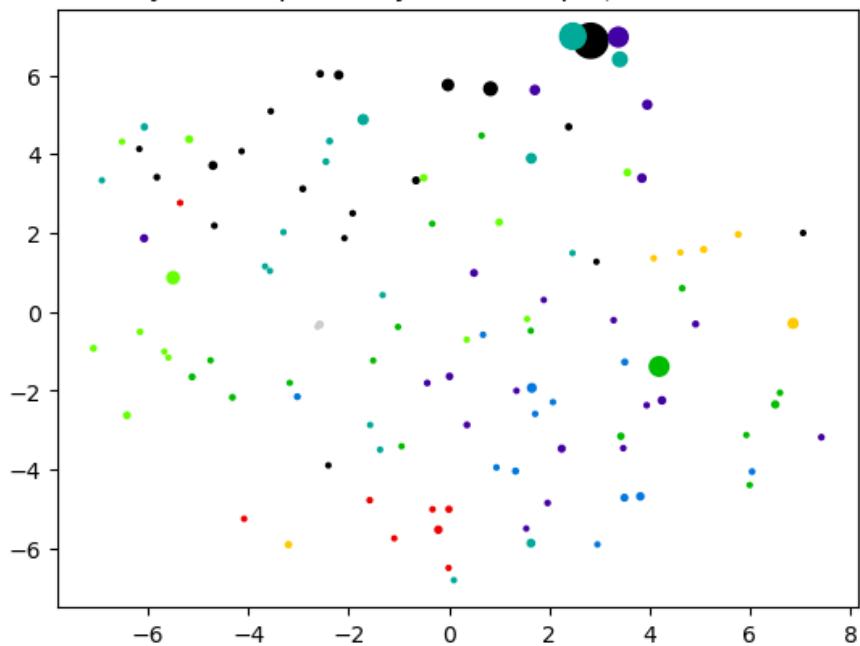
```

-----
TSNE(early_exaggeration=12, random_state=47, verbose=1)
-----
====> Finding 90 nearest neighbors using exact search using euclidean distance...
    --> Time elapsed: 0.00 seconds
====> Calculating affinity matrix...
    --> Time elapsed: 0.00 seconds
====> Calculating PCA-based initialization...
    --> Time elapsed: 0.00 seconds
====> Running optimization with exaggeration=12.00, lr=9.50 for 250 iterations...
Iteration 50, KL divergence -0.5039, 50 iterations in 0.0670 sec
Iteration 100, KL divergence 1.2392, 50 iterations in 0.0178 sec
Iteration 150, KL divergence 1.2392, 50 iterations in 0.0154 sec
Iteration 200, KL divergence 1.2392, 50 iterations in 0.0154 sec
Iteration 250, KL divergence 1.2392, 50 iterations in 0.0154 sec
    --> Time elapsed: 0.13 seconds
====> Running optimization with exaggeration=1.00, lr=114.00 for 500 iterations...
Iteration 50, KL divergence 0.6574, 50 iterations in 0.0491 sec
Iteration 100, KL divergence 0.6427, 50 iterations in 0.0472 sec
Iteration 150, KL divergence 0.6298, 50 iterations in 0.0459 sec
Iteration 200, KL divergence 0.6298, 50 iterations in 0.0468 sec
Iteration 250, KL divergence 0.6300, 50 iterations in 0.0458 sec
Iteration 300, KL divergence 0.6297, 50 iterations in 0.0452 sec
Iteration 350, KL divergence 0.6302, 50 iterations in 0.0455 sec
Iteration 400, KL divergence 0.6304, 50 iterations in 0.0455 sec
Iteration 450, KL divergence 0.6307, 50 iterations in 0.0454 sec
Iteration 500, KL divergence 0.6290, 50 iterations in 0.0455 sec
    --> Time elapsed: 0.46 seconds
(114, 2)

```

	codeUnit	artifact	communityId	centrality	x	y
0	org.axonframework.disruptor.commandhandling	axon-disruptor-4.10.0	0	0.012594	-4.133647	4.073027
1	org.axonframework.commandhandling	axon-messaging-4.10.0	0	0.073080	-0.018907	5.755876
2	org.axonframework.commandhandling.callbacks	axon-messaging-4.10.0	0	0.015707	-2.914155	3.119171
3	org.axonframework.commandhandling.distributed	axon-messaging-4.10.0	0	0.023106	-0.656205	3.332650
4	org.axonframework.commandhandling.distributed....	axon-messaging-4.10.0	0	0.013920	-1.919528	2.499617

Java Package positioned by their dependency relationships (HashGNN node embeddings + t-SNE)



2.5 Node Embeddings for Java Packages using node2vec

The results have been provided by the query filename: `../cypher/Node_EMBEDDINGS/Node_EMBEDDINGS_0a_Query_Calculated.cypher`

	codeUnitName	shortCodeUnitName	projectName	communityId	centrality	embedding
0	org.axonframework.disruptor.commandhandling	commandhandling	axon-disruptor-4.10.0	0	0.012594	<code>[-0.11765751242637634, -0.16672852635383606, 0.1...</code>
1	org.axonframework.commandhandling	commandhandling	axon-messaging-4.10.0	0	0.073080	<code>[0.2573402523994446, -0.29474595189094543, 0.1...</code>
2	org.axonframework.commandhandling.callbacks	callbacks	axon-messaging-4.10.0	0	0.015707	<code>[0.3123716413974762, -0.5747050642967224, -0.0...</code>
3	org.axonframework.commandhandling.distributed	distributed	axon-messaging-4.10.0	0	0.023106	<code>[-0.22914829950196838, -0.3431680202484131, 0....</code>
4	org.axonframework.commandhandling.distributed....	commandfilter	axon-messaging-4.10.0	0	0.013920	<code>[-0.3670963943004608, -0.3763198256492615, 0.3...</code>

```

-----
TSNE(early_exaggeration=12, random_state=47, verbose=1)
-----
====> Finding 90 nearest neighbors using exact search using euclidean distance...
    --> Time elapsed: 0.00 seconds
====> Calculating affinity matrix...
    --> Time elapsed: 0.00 seconds
====> Calculating PCA-based initialization...
    --> Time elapsed: 0.00 seconds
====> Running optimization with exaggeration=12.00, lr=9.50 for 250 iterations...
Iteration 50, KL divergence -0.9374, 50 iterations in 0.0662 sec
Iteration 100, KL divergence 1.1597, 50 iterations in 0.0172 sec
Iteration 150, KL divergence 1.1597, 50 iterations in 0.0150 sec
Iteration 200, KL divergence 1.1597, 50 iterations in 0.0149 sec
Iteration 250, KL divergence 1.1597, 50 iterations in 0.0149 sec
    --> Time elapsed: 0.13 seconds
====> Running optimization with exaggeration=1.00, lr=114.00 for 500 iterations...
Iteration 50, KL divergence 0.3132, 50 iterations in 0.0531 sec
Iteration 100, KL divergence 0.2635, 50 iterations in 0.0507 sec
Iteration 150, KL divergence 0.2546, 50 iterations in 0.0474 sec
Iteration 200, KL divergence 0.2541, 50 iterations in 0.0472 sec
Iteration 250, KL divergence 0.2542, 50 iterations in 0.0466 sec
Iteration 300, KL divergence 0.2538, 50 iterations in 0.0470 sec
Iteration 350, KL divergence 0.2536, 50 iterations in 0.0472 sec
Iteration 400, KL divergence 0.2537, 50 iterations in 0.0480 sec
Iteration 450, KL divergence 0.2538, 50 iterations in 0.0467 sec
Iteration 500, KL divergence 0.2537, 50 iterations in 0.0468 sec
    --> Time elapsed: 0.48 seconds

```

(114, 2)

	codeUnit	artifact	communityId	centrality	x	y
0	org.axonframework.disruptor.commandhandling	axon-disruptor-4.10.0	0	0.012594	-2.465379	1.104010
1	org.axonframework.commandhandling	axon-messaging-4.10.0	0	0.073080	-3.746857	-1.501235
2	org.axonframework.commandhandling.callbacks	axon-messaging-4.10.0	0	0.015707	-4.715600	-2.404075
3	org.axonframework.commandhandling.distributed	axon-messaging-4.10.0	0	0.023106	-4.900034	-3.777513
4	org.axonframework.commandhandling.distributed....	axon-messaging-4.10.0	0	0.013920	-4.944299	-3.894449

Java Package positioned by their dependency relationships (node2vec node embeddings + t-SNE)

