

rdf-network-analysis

May 15, 2020

1 Network Analysis of RDF Graphs

In this notebook we provide basic facilities for performing network analyses of RDF graphs easily with Python [rdflib](#) and [networkx](#)

We do this in 4 steps: 1. Load an arbitrary RDF graph into rdflib 2. Get a subgraph of relevance (optional) 3. Convert the rdflib Graph into an networkx Graph, as shown [here](#) 4. Get an network analysis report by running networkx's algorithms on that data structure

1.1 0. Preparation

```
[13]: # Install required packages in the current Jupyter kernel
# Uncomment the following lines if you need to install these libraries
# If you run into permission issues, try with the --user option
import sys
# !pip install -q rdflib networkx matplotlib scipy
!{sys.executable} -m pip install rdflib networkx matplotlib scipy --user

# Imports
from rdflib import Graph as RDFGraph
from rdflib.extras.external_graph_libs import rdflib_to_networkx_graph
import networkx as nx
from networkx import Graph as NXGraph
import matplotlib.pyplot as plt
import statistics
import collections
```

Requirement already satisfied: rdflib in /home/amp/.local/lib/python3.8/site-packages (5.0.0)

Requirement already satisfied: networkx in /home/amp/.local/lib/python3.8/site-packages (2.4)

Requirement already satisfied: matplotlib in /home/amp/.local/lib/python3.8/site-packages (3.2.1)

Requirement already satisfied: scipy in /home/amp/.local/lib/python3.8/site-packages (1.4.1)

Requirement already satisfied: pyparsing in /usr/lib/python3/dist-packages (from rdflib) (2.4.6)

Requirement already satisfied: six in /usr/lib/python3/dist-packages (from rdflib) (1.14.0)

Requirement already satisfied: isodate in /home/amp/.local/lib/python3.8/site-packages (from rdflib) (0.6.0)

Requirement already satisfied: decorator>=4.3.0 in /usr/lib/python3/dist-packages (from networkx) (4.4.2)

Requirement already satisfied: kiwisolver>=1.0.1 in /home/amp/.local/lib/python3.8/site-packages (from matplotlib) (1.2.0)

Requirement already satisfied: numpy>=1.11 in /usr/lib/python3/dist-packages (from matplotlib) (1.17.4)

Requirement already satisfied: cyclor>=0.10 in /home/amp/.local/lib/python3.8/site-packages (from matplotlib) (0.10.0)

Requirement already satisfied: python-dateutil>=2.1 in /usr/lib/python3/dist-packages (from matplotlib) (2.7.3)

1.2 1. Loading RDF

The first thing to do is to load the RDF graph we want to perform the network analysis on. By executing the next cell, we'll be asked to fill in the path to an RDF graph. This can be any path, local or online, that we can look up.

Any of the Turtle (`ttl`) files that we include with this notebook will do; for example, `bsbm-sample.ttl`. But any Web location that leads to an RDF file (for example, the GitHub copy of that same file at <https://raw.githubusercontent.com/albertmeronyo/rdf-network-analysis/master/bsbm-sample.ttl>; or any other RDF file on the Web like <https://raw.githubusercontent.com/albertmeronyo/lodapi/master/ghostbusters.ttl>) will work too.

```
[18]: # RDF graph loading
path = input("Path or URI of the RDF graph to load: ")
rg = RDFGraph()
rg.parse(path, format='turtle')
print("rdflib Graph loaded successfully with {} triples".format(len(rg)))
```

```
Path or URI of the RDF graph to load: wechanged-german.ttl
rdflib Graph loaded successfully with 53 triples
```

1.3 2. Get a subgraph out of the loaded RDF graph (optional)

This cell can be skipped altogether without affecting the rest of the notebook; but it will be useful if instead of using the whole RDF graph of the previous step, we just want to use a subgraph that's included in it.

By executing the next cell, we'll be asked two things:

- The URI of the "entity" type we are interested in (e.g. <http://dbpedia.org/ontology/Band>)
- The URI of the "relation" connecting entities we are interested in (e.g. <http://dbpedia.org/ontology/influencedBy>)

Using these two, the notebook will replace the original graph with the subgraph that's constructed by those entity types and relations only.

```
[ ]: # Subgraph construction (optional)
entity = input("Entity type to build nodes of the subgraph with: ")
relation = input("Relation type to build edges of the subgraph with: ")

# TODO: Use entity and relation as parameters of a CONSTRUCT query
query = """
PREFIX bsbm: <http://www4.wiwiss.fu-berlin.de/bizer/bsbm/v01/vocabulary/>
CONSTRUCT {{ ?u a {} . ?u {} ?v }} WHERE {{ ?u a {} . ?u {} ?v }}"""
    ↪format(entity, relation, entity, relation)
# print(query)
subg = rg.query(query)

rg = subg
```

1.4 3. Converting rdflib.Graph to networkx.Graph

Thanks to [the great work done by the rdflib developers](#) this step, which converts the basic graph data structure of rdflib into its equivalent in networkx, is straightforward. Just run the next cell to make our RDF dataset ready for network analysis!

```
[19]: # Conversion of rdflib.Graph to networkx.Graph
G = rdflib_to_networkx_graph(rg)
print("networkx Graph loaded successfully with length {}".format(len(G)))
```

networkx Graph loaded successfully with length 65

1.5 4. Network analysis

At this point we can run the network analysis on our RDF graph by using the networkx algorithms. Executing the next cell will output a full network analysis report, with the following parts:

- General network metrics (network size, pendants, density)
- Node centrality metrics (degree, eigenvector, betweenness). For these, averages, stdevs, maximum, minimum and distribution histograms are given
- Clustering metrics (connected components, clustering)
- Overall network plot

The report can be easily selected and copy-pasted for further use in other tools.

```
[20]: # Analysis

def mean(numbers):
    return float(sum(numbers)) / max(len(numbers), 1)

def number_of_pendants(g):
    """
    Equals the number of nodes with degree 1
    """
    pendants = 0
```

```

for u in g:
    if g.degree[u] == 1:
        pendants += 1
return pendants

def histogram(l):
    degree_sequence = sorted([d for n, d in list(l.items())], reverse=True)
    degreeCount = collections.Counter(degree_sequence)
    deg, cnt = zip(*degreeCount.items())
    print(deg, cnt)

    fig, ax = plt.subplots()
    plt.bar(deg, cnt, width=0.80, color='b')

    plt.title("Histogram")
    plt.ylabel("Count")
    plt.xlabel("Value")
    ax.set_xticks([d + 0.4 for d in deg])
    ax.set_xticklabels(deg)

    plt.show()

# Network size
print("NETWORK SIZE")
print("=====")
print("The network has {} nodes and {} edges".format(G.number_of_nodes(), G.
    ↳number_of_edges()))
print()

# Network size
print("PENDANTS")
print("=====")
print("The network has {} pendants".format(number_of_pendants(G)))
print()

# Density
print("DENSITY")
print("=====")
print("The network density is {}".format(nx.density(G)))
print()

# Degree centrality -- mean and stdev
dc = nx.degree_centrality(G)
degrees = []
for k,v in dc.items():
    degrees.append(v)

```

```

print("DEGREE CENTRALITY")
print("=====")
print("The mean degree centrality is {}, with stdev {}".format(mean(degrees),
↳statistics.stdev(degrees)))
print("The maximum node is {}, with value {}".format(max(dc, key=dc.get),
↳max(dc.values())))
print("The minimum node is {}, with value {}".format(min(dc, key=dc.get),
↳min(dc.values())))
histogram(dc)
print()

# Eigenvector centrality -- mean and stdev
ec = nx.eigenvector_centrality_numpy(G)
degrees = []
for k,v in ec.items():
    degrees.append(v)

print("EIGENVECTOR CENTRALITY")
print("=====")
print("The mean network eigenvector centrality is {}, with stdev {}".
↳format(mean(degrees), statistics.stdev(degrees)))
print("The maximum node is {}, with value {}".format(max(ec, key=ec.get),
↳max(ec.values())))
print("The minimum node is {}, with value {}".format(min(ec, key=ec.get),
↳min(ec.values())))
histogram(ec)
print()

# Betweenness centrality -- mean and stdev
# bc = nx.betweenness_centrality(G)
# degrees = []
# for k,v in bc.items():
#     degrees.append(v)
# print("BETWEENNESS CENTRALITY")
# print("=====")
# print("The mean betweenness centrality is {}, with stdev {}".
↳format(mean(degrees), statistics.stdev(degrees)))
# print("The maximum node is {}, with value {}".format(max(bc, key=bc.get),
↳max(bc.values())))
# print("The minimum node is {}, with value {}".format(min(bc, key=bc.get),
↳min(bc.values())))
# histogram(bc)
# print()

```

```

# Connected components
cc = list(nx.connected_components(G))
print("CONNECTED COMPONENTS")
print("=====")
print("The graph has {} connected components".format(len(cc)))
for i,c in enumerate(cc):
    print("Connected component {} has {} nodes".format(i,len(c)))
print()

# Clusters
cl = nx.clustering(G)
print("CLUSTERS")
print("=====")
print("The graph has {} clusters".format(len(cl)))
for i,c in enumerate(cl):
    print("Cluster {} has {} nodes".format(i,len(c)))
print()

# Plot
print("Visualizing the graph:")
plt.plot()
plt.figure(1)
nx.draw(G, with_labels=False, font_weight='normal', node_size=60, font_size=8)
plt.figure(1,figsize=(120,120))
plt.savefig('example.png', dpi=1000)

```

NETWORK SIZE

=====

The network has 65 nodes and 53 edges

PENDANTS

=====

The network has 51 pendants

DENSITY

=====

The network density is 0.02548076923076923

DEGREE CENTRALITY

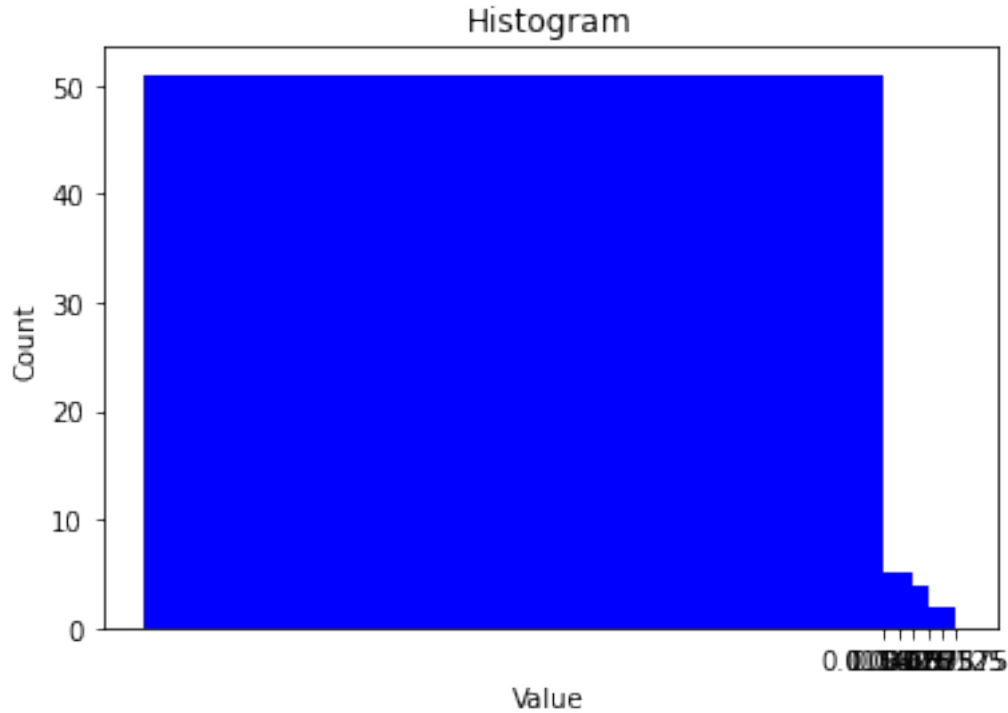
=====

The mean degree centrality is 0.02548076923076923, with stdev
0.020774719730186218

The maximum node is <http://www.wikidata.org/entity/Q165824>, with value 0.09375

The minimum node is wcd_00814_id, with value 0.015625

(0.09375, 0.078125, 0.0625, 0.046875, 0.03125, 0.015625) (2, 2, 4, 5, 1, 51)



EIGENVECTOR CENTRALITY

=====

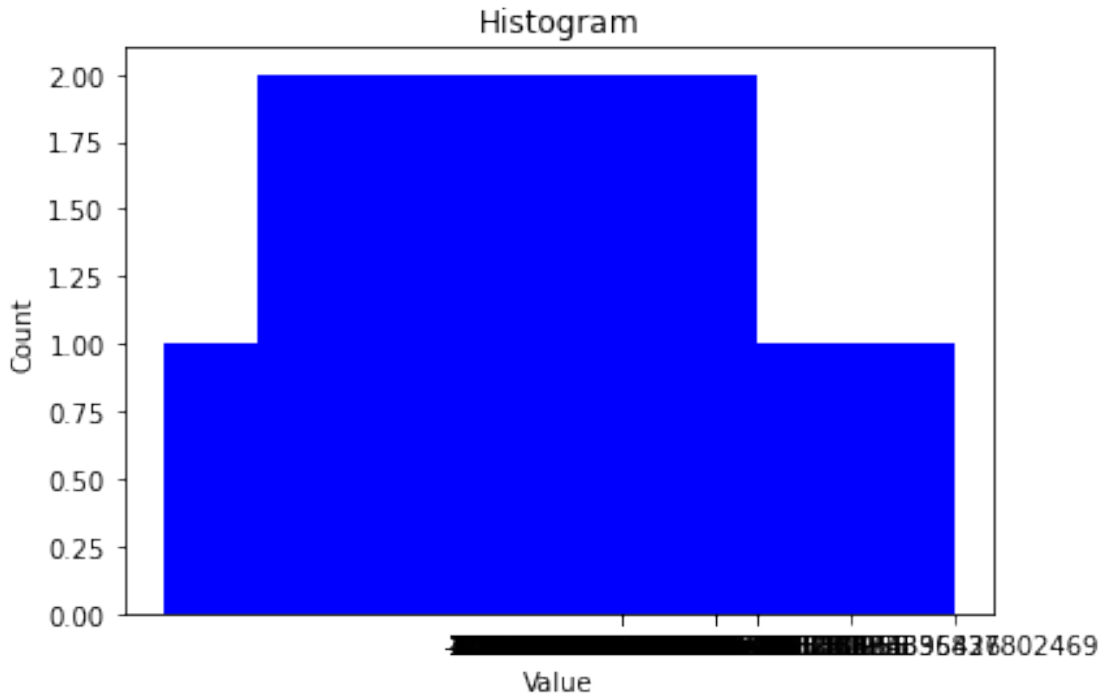
The mean network eigenvector centrality is 0.052190328754421186, with stdev 0.11339581003839311

The maximum node is <http://www.wikidata.org/entity/Q165824>, with value 0.5823336837802469

The minimum node is <http://www.wikidata.org/entity/Q2653682>, with value -4.221865487293616e-16

(0.5823336837802469, 0.40110781684595426, 0.23773673088280958, 0.23773673088280955, 0.23773673088280953, 0.2377367308828095, 0.16375158051905314, 0.1637515805190531, 0.16375158051905309, 0.16375158051905306, 0.16375158051905303, 4.0375682011403296e-16, 3.867620361637153e-16, 2.423140802377901e-16, 2.1493997183573874e-16, 2.0826656295842276e-16, 1.8925354328681477e-16, 1.860937486981313e-16, 1.7617115305582745e-16, 1.667799235809163e-16, 1.5837406027033936e-16, 1.3001507409184495e-16, 1.2555668835368535e-16, 1.1354908529513395e-16, 1.0195879145351594e-16, 9.659051261599858e-17, 8.249547678468506e-17, 7.150789936020287e-17, 6.477197106861316e-17, 6.34748456895395e-17, 6.255159781101699e-17, 5.748386506242491e-17, 4.5070823959909377e-17, 4.028075437461217e-17, 3.191177899331292e-17, 1.7134628208983114e-17, 1.504830421598233e-17, 1.5222146334878828e-18, -1.2422309969230075e-18, -8.56840964600123e-18, -1.0566991786028656e-17, -1.3380020371347866e-17, -1.82290962330364e-17, -2.5459354322188953e-17, -2.9381498680853597e-17,

-5.607146449104495e-17, -6.208476377865393e-17, -6.278772145308986e-17,
-1.0752076302444307e-16, -1.0828082789917711e-16, -1.7349787989201016e-16,
-1.765445338168193e-16, -1.879654759105251e-16, -1.9130239507871805e-16,
-1.93439031608548e-16, -1.9861678449786301e-16, -1.9935548922841931e-16,
-2.3141199604139336e-16, -2.477318548940688e-16, -2.722433427921724e-16,
-3.7706856978891896e-16, -4.221865487293616e-16) (1, 1, 1, 1, 2, 2, 1, 1, 1, 2,
1,
1, 1)



CONNECTED COMPONENTS

=====

The graph has 12 connected components
Connected component 0 has 5 nodes
Connected component 1 has 7 nodes
Connected component 2 has 5 nodes
Connected component 3 has 4 nodes
Connected component 4 has 9 nodes
Connected component 5 has 4 nodes
Connected component 6 has 6 nodes
Connected component 7 has 6 nodes
Connected component 8 has 4 nodes
Connected component 9 has 7 nodes
Connected component 10 has 4 nodes
Connected component 11 has 4 nodes

CLUSTERS

=====

The graph has 65 clusters

Cluster 0 has 37 nodes

Cluster 1 has 12 nodes

Cluster 2 has 38 nodes

Cluster 3 has 37 nodes

Cluster 4 has 37 nodes

Cluster 5 has 25 nodes

Cluster 6 has 37 nodes

Cluster 7 has 25 nodes

Cluster 8 has 39 nodes

Cluster 9 has 40 nodes

Cluster 10 has 39 nodes

Cluster 11 has 25 nodes

Cluster 12 has 25 nodes

Cluster 13 has 37 nodes

Cluster 14 has 25 nodes

Cluster 15 has 25 nodes

Cluster 16 has 39 nodes

Cluster 17 has 25 nodes

Cluster 18 has 25 nodes

Cluster 19 has 39 nodes

Cluster 20 has 12 nodes

Cluster 21 has 37 nodes

Cluster 22 has 12 nodes

Cluster 23 has 12 nodes

Cluster 24 has 25 nodes

Cluster 25 has 37 nodes

Cluster 26 has 12 nodes

Cluster 27 has 25 nodes

Cluster 28 has 25 nodes

Cluster 29 has 37 nodes

Cluster 30 has 12 nodes

Cluster 31 has 25 nodes

Cluster 32 has 39 nodes

Cluster 33 has 12 nodes

Cluster 34 has 25 nodes

Cluster 35 has 40 nodes

Cluster 36 has 25 nodes

Cluster 37 has 12 nodes

Cluster 38 has 40 nodes

Cluster 39 has 12 nodes

Cluster 40 has 25 nodes

Cluster 41 has 39 nodes

Cluster 42 has 25 nodes

Cluster 43 has 38 nodes

Cluster 44 has 12 nodes
Cluster 45 has 25 nodes
Cluster 46 has 25 nodes
Cluster 47 has 25 nodes
Cluster 48 has 25 nodes
Cluster 49 has 25 nodes
Cluster 50 has 25 nodes
Cluster 51 has 37 nodes
Cluster 52 has 25 nodes
Cluster 53 has 12 nodes
Cluster 54 has 25 nodes
Cluster 55 has 40 nodes
Cluster 56 has 40 nodes
Cluster 57 has 40 nodes
Cluster 58 has 25 nodes
Cluster 59 has 12 nodes
Cluster 60 has 38 nodes
Cluster 61 has 25 nodes
Cluster 62 has 12 nodes
Cluster 63 has 25 nodes
Cluster 64 has 25 nodes

Visualizing the graph:

