

project_2

March 10, 2023

In this project, you will be working with the protein interaction network (PIN) of *Saccharomyces cerevisiae* from the BioGRID Multi-Validated (MV) Datasets. The nodes and edges of this PIN denote proteins and their physical interactions, respectively. These interactions have been identified from large-scale protein interaction experiments such as yeast two-hybrid screenings and affinity purification with mass spectrometry. Additionally, they pass a specific set of criteria as defined [here](#) (BioGRID MV Datasets). See the following [link](#) to familiarize yourself with the contents of this BioGRID file. You will be analyzing the yeast PIN using k -core analysis, and further use the network as a scaffold for incorporating gene expression data, allowing you to perform s -core analysis using pairwise gene expression correlations as weights.

0.0.1 1.1 Yeast protein interaction network

Download the file containing all multi-validated protein-protein interactions from BioGRID found on Blackboard BioGRID_PIN.txt. Create a network with NetworkX using the systematic names to set source and target nodes. These corresponds to the columns named Systematic Name Interactor A and Systematic Name Interactor B. Remove all self-edges from the network.

```
[ ]: # Imports
import pandas as pd
import numpy as np
import networkx as nx
from pyvis.network import Network
```

```
[ ]: import sys
import os
path_ = sys.path[0].split("/")
os.chdir("/")
os.chdir(path_[:-1])
os.getcwd()
from SysBio_fun import *
os.chdir("/")
os.chdir(path_)
```

```
[ ]: # Read file:
df = pd.read_table("BioGRID_PIN.txt", sep="\t", low_memory=False)

# Define source and target column:
source = "Official Symbol Interactor A"
target = "Official Symbol Interactor B"

# Rename columns:
```

```

df.rename(columns={source:"source",
                  target:"target"},
          inplace=True)
# Drop rows where column "source" value is identical to column "target" value.
# I.e., drop self loops:
# **NB!** This the number of rows left in the df may vary depending on which
# column-pairs were treated as source-target. E.g. source:"- " and target:"- "
# will be removed, even though they may have different Entrez ID.
df = df[df["source"] != df["target"]]

# Generate the graph:
# **NB!** I treat the graph as multi directed, i.e. each row in the df
# amounts to a single edge in the network!
PIN = nx.from_pandas_edgelist(df, edge_attr=True, create_using=nx.
    ↪MultiDiGraph())

# Make a copy
PIN_copy = PIN.copy()

```

- (i) How large is the network (i.e. number of nodes and edges), and what is the average node degree?

```

[ ]: N = PIN.number_of_nodes()
E = PIN.number_of_edges()
mean_k = (E*2)/N # Average node degree. Each link contributes a degree "point" ↵
    ↪for each node it connects. Therefore 2*E

print(f"""
The number of nodes are: {N}
The number of edges are: {E}
The average node degree is: {mean_k}
""")

```

```

The number of nodes are: 20550
The number of edges are: 245288
The average node degree is: 23.872311435523116

```

- (ii) Look up *S. cerevisiae* S288C in the NCBI Taxonomy Database [here](#) and take note of the taxonomy ID. Using the Organism Interactor A and Organism Interactor B edge attributes in the network, use this taxonomy ID to filter out all non-*S. cerevisiae* S288C nodes from the network. Also, remove all self-loops and isolated nodes. How many nodes and edges remain? Create a plot of its degree distribution. How would you characterize this network and its degree distribution?

```

[ ]: PIN = PIN_copy.copy()

```

```

# The taxonomic ID of S. cerevisiae S288C is: 559292

# Generate a list of edges that does not fit:
r1 = [edge for edge,tax \
      in nx.get_edge_attributes(PIN, "Organism Interactor A").items() \
      if tax != 559292] # The edge information is a tuple and the key
r2 = [edge for edge,tax \
      in nx.get_edge_attributes(PIN, "Organism Interactor B").items() \
      if tax != 559292]
rem = set(r1+r2) # set() is used to remove potential duplicates
# Use the list to remove all edges
PIN.remove_edges_from(rem)
# In pandas:
# len(df[(df["Organism Interactor A"] == 559292) & (df["Organism Interactor B"]
↳ == 559292)])

# Remove nodes with a node degree of 0 (i.e. isolated nodes):
PIN.remove_nodes_from( # Remove nodes from list ...
    [nodes for nodes,degree \
     in dict(PIN.degree()).items() if degree == 0] # nodes with degree 0
)

N = PIN.number_of_nodes() # New number of nodes
E = PIN.number_of_edges() # New number of edges

print(f"""
The number of nodes is now {N},
while number of edges is {E}
""")

```

The number of nodes is now 4104,
while number of edges is 68704

- (iii) Compare the yeast PIN with Barabási–Albert (BA) and Erdős–Rényi (ER) networks with approximately the same number of nodes and edges using a few selected network measurements and the degree distributions. Do these models do a good job in describing the yeast PIN? Discuss why/why not.

```

[ ]: l_comp = nx.Graph(PIN.subgraph(
    sorted(
        nx.connected_components(
            nx.to_undirected(PIN.copy())
            #PIN.copy().to_undirected(as_view=True)
        ),
        reverse=True
    )
)

```

```

    )[0]
).copy())

n = l_comp.number_of_nodes()
e = l_comp.number_of_edges()

ER = connected_ER(n, e)
BA = connected_BA(n, round(((e*2)/n)))

```

```

[ ]: #type(PIN.subgraph(nx.connected_components(PIN.to_undirected())[0]))
type(
    nx.Graph(PIN.subgraph(
        sorted(
            nx.connected_components(
                nx.Graph(PIN.copy())
            ), reverse=True)[0]
        ))
    )
type(PIN)

# PIN2 = nx.Graph(PIN.copy())
# print(type(PIN2))
# print(PIN == PIN2)
# PIN2 = PIN2.to_undirected()
# print(type(PIN2))

```

```

[ ]: networkx.classes.multidigraph.MultiDiGraph

```

```

[ ]: graph_info(l_comp)
graph_info(ER)
graph_info(BA)

```

```

Number of nodes: 4022
Number of edges: 16700
Is connected:    True
Number of nodes: 4022
Number of edges: 16700
Is connected:    True
Number of nodes: 4022
Number of edges: 32112
Is connected:    True

```

```

[ ]: def degree_distribution(graph):
    """
    Function that returns a scatter plot of the degree
    distribution of a graph.
    """
    N = graph.number_of_nodes() # Total number of nodes

```

```

    y = np.array(nx.degree_histogram(graph))/N # All occurrences of "n"
    ↪ degree, divided by total number of nodes
    x = np.arange(len(y))[y != 0] # x values
    y = y[y != 0] # Remove the values = 0 from the array

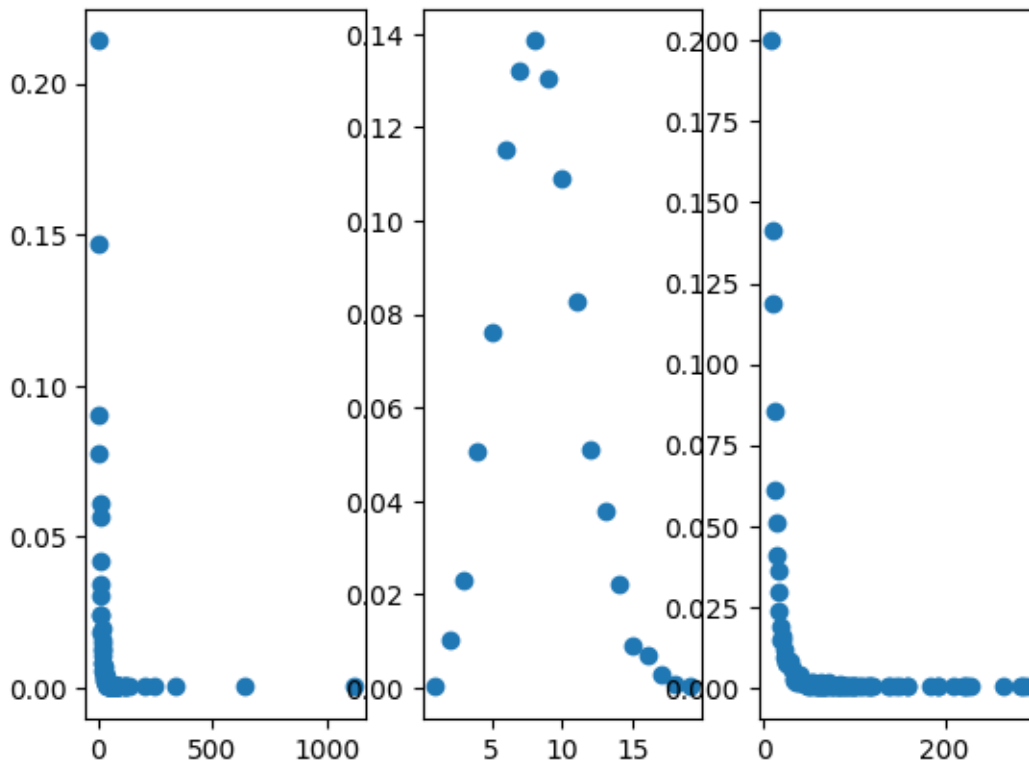
    return x, y

x_l_comp, y_l_comp = degree_distribution(l_comp)
x_ER, y_ER = degree_distribution(ER)
x_BA, y_BA = degree_distribution(BA)

figure, axis = plt.subplots(1, 3)
axis[0].scatter(x_l_comp, y_l_comp)
axis[1].scatter(x_ER, y_ER)
axis[2].scatter(x_BA, y_BA)

```

[]: <matplotlib.collections.PathCollection at 0x7c30ebb3e700>



[]: `#plt.scatter(nx.clustering(BA).keys(), nx.clustering(BA).values())`

(iv) You will now start peeling away layers of the yeast PIN by k -core analysis. In your own words,

describe/define k -core analysis, and explain how it works on a network. What network does the 1-core correspond to?

k -core analysis is a method to

- (v) Calculate the 2-core of the yeast PIN. How large is this network? Plot and characterize its degree distribution and compare it to the original yeast PIN.

```
[ ]: l_comp_2_core = nx.k_core(l_comp, k=2)
      graph_info(l_comp_2_core)

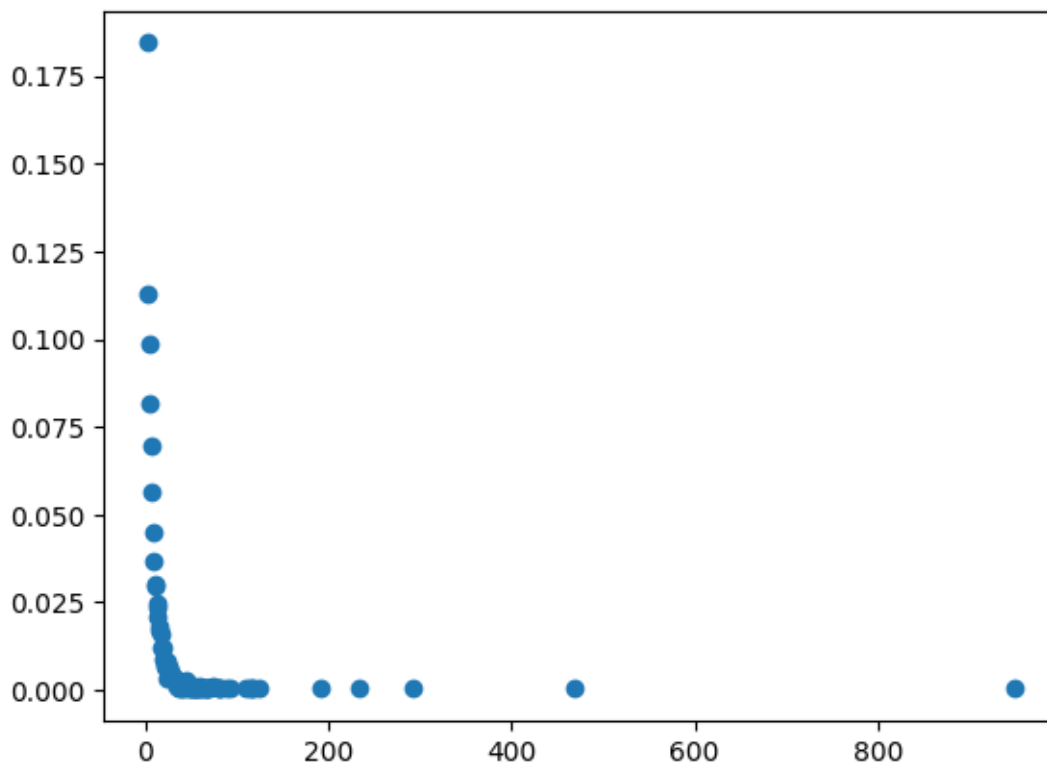
      x, y = degree_distribution(l_comp_2_core)
      plt.scatter(x,y)
```

Number of nodes: 3102

Number of edges: 15780

Is connected: True

```
[ ]: <matplotlib.collections.PathCollection at 0x7c30ee367190>
```



- (vi) Find the innermost k -core (the last k before the network is empty). How many edges and nodes are there in the innermost k -core? What is the absolute maximal amount of edges that there could be between this number of nodes, and how does the innermost core compare to that?

```
[ ]: current_k = 0
g = l_comp.copy()
while nx.number_of_nodes(g) > 0:
    current_k += 1
    g = nx.k_core(g, k=current_k)
innermost_k_core = current_k-1

print(f"The innermost k-core is: {innermost_k_core}")

graph_info(nx.k_core(l_comp.copy(), k=innermost_k_core))
```

The innermost k-core is: 18

Number of nodes: 19

Number of edges: 171

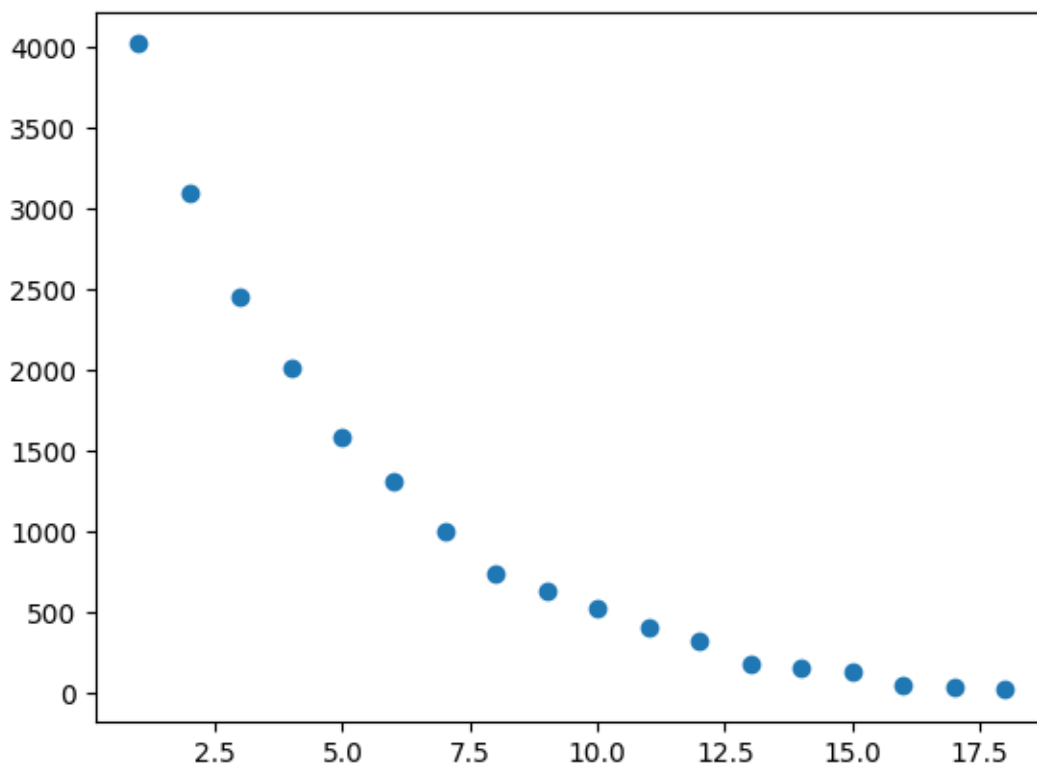
Is connected: True

(vii) Create a plot of the number of nodes and the number of edges in a k -core against k . How would you characterize this plot? What does it tell you?

```
[ ]: y = [nx.number_of_nodes(nx.k_core(l_comp.copy(), k=k_)) for k_ in range(1,
↪innermost_k_core+1)]
x = range(1, innermost_k_core+1)
```

```
[ ]: plt.scatter(x, y)
```

```
[ ]: <matplotlib.collections.PathCollection at 0x7c30e126ea60>
```



(viii) Generate the second innermost k -core and visualize it using pyvis. Describe the network.

```
[ ]: # def show_html(graph, name="nx", show=True, size="small"):
#     """
#     Generate and display the graph through pyvis
#     """
#
#     if size == "small":
#         nt = Network("500px", "500px")
#     else:
#         nt = Network("1080px", "1920px")
#
#     nt.from_nx(graph)
#     if show == True:
#         nt.show(f"{name}.html")
#
s = nx.k_core(l_comp.copy(), k=innermost_k_core-1)
show_html(s, name=f"{innermost_k_core-1}-core_network", show=True, size="large")
```

(ix) We will now load gene expression data onto our yeast PIN. Go to this [url](#) containing the gene expression data set from *Transcription profiling by array of yeast to investigate expression of the beta-subunit of Snf1 kinase*. For those interested, you can find detailed information about the dataset in the original publication [here](#). Download `normalized_expressions.tsv` from Blackboard and import the list of gene names, as well as corresponding expression data from 24 experimental conditions. For all pairwise gene combinations, calculate the Pearson correlation coefficients and associated p-values for the N genes present in the yeast PIN (both as $N \times N$ Numpy arrays). Out of the gene pairs with positive correlation values, how many of them have a significance value below 0.0001?

```
[ ]: norm_ex = pd.read_table("normalized_expressions.tsv", sep="\t")
norm_ex.head(10)
```

```
[ ]:
Gene ID  GSM517188  GSM517189  GSM517170  GSM517171  GSM517180  GSM517186  \
0  YAL001C  10.862551  10.787718  10.609314  10.533807  10.683814  10.884063
1  YAL002W  10.318229  10.376402  10.307895  10.231124  10.217388  10.289048
2  YAL003W  14.116945  14.110254  14.038039  14.042464  14.116978  14.100915
3  YAL005C  14.372108  14.419899  14.327958  14.283262  14.436713  14.418395
4  YAL007C  11.717200  11.861313  11.594489  11.750768  11.574141  11.824636
5  YAL008W  11.275044  11.248086  11.370819  11.108830  11.510864  11.330115
6  YAL009W   9.971481   9.935119   9.988509   9.993210   9.915578   9.936548
7  YAL010C   9.788592   9.882031   9.969160   9.935823   9.869899   9.830950
8  YAL011W   9.434274   9.396159   9.500599   9.525675   9.558400   9.427803
9  YAL012W  13.515144  13.565094  13.187785  13.172306  13.216801  13.380632

GSM517187  GSM517183  GSM517172  ...  GSM517184  GSM517168  GSM517181  \
```


0	10.991316	10.793569	10.641122	...	10.776250	10.829294	10.706164
1	10.571356	10.192751	10.307293	...	10.166869	10.360606	10.282979
2	14.101845	14.105961	13.999617	...	14.085988	14.127210	14.048146
3	14.354286	14.369767	14.285516	...	14.243826	14.349704	14.361417
4	11.856330	11.650583	11.548949	...	11.612014	11.475502	11.608367
5	11.461117	11.363152	11.246455	...	11.431879	11.554264	11.363908
6	10.152794	9.824603	10.018877	...	9.809905	9.740123	9.854044
7	10.069362	9.904191	9.780692	...	9.966418	9.828681	9.752460
8	9.519610	9.578430	9.450371	...	9.538738	9.512650	9.614558
9	13.419438	13.374047	13.050084	...	13.449695	13.217452	13.372705

	GSM517185	GSM517167	GSM517169	GSM517190	GSM517173	GSM517174	GSM517175
0	11.036464	10.876691	10.778201	10.860265	10.581020	10.698511	10.561507
1	10.015871	10.210560	10.433530	10.334844	10.246341	10.279370	10.347147
2	14.031750	14.115299	14.085665	14.117709	13.981556	14.010271	13.993024
3	14.238842	14.298123	14.325895	14.371525	14.310817	14.346674	14.345590
4	11.670285	11.479741	11.490664	11.975316	11.606178	11.605336	11.535457
5	11.351615	11.426774	11.533108	11.299099	11.358341	11.310492	11.474645
6	10.033609	9.888328	9.792665	10.048278	9.968829	10.066921	9.871540
7	10.116446	9.899894	9.906184	9.823420	9.742537	9.947832	9.807137
8	9.257099	9.566624	9.492080	9.458073	9.336601	9.489274	9.422245
9	13.160807	13.213176	13.258713	13.537055	13.238104	13.163476	13.089239

[10 rows x 25 columns]

- (x) Download `s_core.py`. This is similar to the NetworkX function `k_core`, but uses node strength instead of node degree. Strength is defined as the sum of the weights of all the links attached to a node (i.e. Person correlations). With negative weights, the absolute values of the weights are used to calculate the strength. In this case, however, we are only looking at positively correlated nodes. Do you think it would make sense to include negative correlations in this particular case? Discuss.
- (xi) Calculate the s -core with $s = 2.5, 3, 3.5$, and 4 , and characterize their degree distributions. Discuss how the distribution changes for increasing values of s .

```
[ ]: from s_core import *
```

- (xii) Find an integer value for s that gives as good match as possible relative to the second innermost k -core in terms of number of distinct modules and number of nodes. Visualize this s -core with `pyvis` and describe the network.

```
[ ]:
```

- (xiii) Using the full yeast PIN network (network from (ii)), extract the nodes that are either part of the second innermost k -core, the s -core in (xii), or part of both cores. How many nodes and edges are there in this new network? Create a visualization of this subnetwork using `pyvis`, and color the nodes according to the following classifications:

- Red if it was only found in the second innermost k -core.
- Green if it was only found in the s -core.

- Blue if it was found in both.

Describe the network and discuss what you observe. Is it a single connected component? What do the differences between the s - and k -cores tell you?

```
[ ]: path_ = "/" .join(sys.path[0].split("/")[:-1])
      #os.chdir("/") .join(path_[:-1]))

      #file_path = '~/TBT4165/SysBio_fun.py'
      #print(path_+"/SysBio_fun.py")
      with open(path_+"/SysBio_fun.py") as file:
          print(file.read())
```

```
import networkx as nx
import matplotlib.pyplot as plt
import numpy as np
from pyvis.network import Network
```

```
#####
#   This is a module created to hold           #
#   all custom functions used in multiple     #
#   project throughout the TBT4165 course    #
#####
```

```
def connected_ER(N, E):
    """
    N : int, Number of nodes
    E : int, Number of edges
    """
    g = nx.gnm_random_graph(N, E)
    while nx.is_connected(g) == False:
        g = nx.gnm_random_graph(N, E)
    return g
```

```
def connected_BA(N, m):
    """
    N : int, Number of nodes
    m : int, Number of edges to attach from a new node to existing nodes
    """
    g = nx.barabasi_albert_graph(n=N, m=m)
    while nx.is_connected(g) == False:
        g = nx.barabasi_albert_graph(n=N, m=m)
    return g
```

```
def get_largest_component(g, is_undirected=True):
    if is_undirected:
        l_comp = g.subgraph(
```

```

        sorted(
            nx.connected_components(
                nx.to_undirected(g.copy())
            ),
            reverse=True
        )[0]
    ).copy()
else:
    l_comp = g.subgraph(
        sorted(
            nx.connected_components(g),
            reverse=True
        )[0]
    ).copy()
return l_comp

def graph_info(g):
    print(f"\nNumber of nodes: {g.number_of_nodes()}")
    print(f"Number of edges: {g.number_of_edges()}")
    try:
        print(f"Is connected: {nx.is_connected(g)}")
    except:
        print("""\
The graph is not undirected. Therefore .is_connected() does not work.
""")

def show_html(graph, name="nx", show=True, size="small"):
    """
    Generate and display the graph through pyvis
    """

    if size == "small":
        nt = Network("500px", "500px")
    else:
        nt = Network("1080px", "1920px")

    nt.from_nx(graph)
    if show == True:
        nt.show(f"{name}.html")

def nice_plot(title, xlab, ylab, show=True):
    """
    Function to generate title and axis labels to a plot.
    The argument "show" is by default True.

```

```

    """
    plt.title(title)
    plt.xlabel(xlab)
    plt.ylabel(ylab)
    if (show == True):
        plt.show()

def plot_degree_distribution(graph):
    """
    Function that returns a scatter plot of the degree
    distribution of a graph.
    """
    N = graph.number_of_nodes() # Total number of nodes
    y = np.array(nx.degree_histogram(graph))/N # All occurrences of "n" degree,
divided by total number of nodes
    x = np.arange(len(y))[y != 0] # x values
    y = y[y != 0] # Remove the values = 0 from the array

    return [x, y]
#return plt.scatter(x,y)

```