

Exercise set #7 (18 pts)

- The deadline for handing in your solutions is November 9th 2020 23:55.
- Return your solutions (one `.pdf` file and one `.zip` file containing Python code) in MyCourses (Assignments tab). Additionally, submit your pdf file also to the Turnitin plagiarism checker in MyCourses.
- Check also the course practicalities page in MyCourses for more details on writing your report.

1. Weight–topology correlations in social networks (12 pts)

In this exercise, we will do some weighted network analysis using a social network data set describing private messaging in a Facebook-like web-page¹. In the network, each node corresponds to a user of the website and link weights describe the total number of messages exchanged between users.

In the file `0Clinks_w_undir.edg`, the three entries of each row describe one link:

`(node_i node_j w_ij)`,

where the last entry `w_ij` is the weight of the link between nodes `node_i` and `node_j`.

You can use the accompanying Python template (`weight_topology_correlations.py`) to get started. `scipy.stats.binned_statistic` function is especially useful throughout this exercise.

- a) (3 pts) Before performing more sophisticated analysis, it is always good to get some idea on what the network is like. To this end, plot the complementary cumulative distribution (1-CDF) for node degree k , node strength s and link weight w .
- **Show** all three distributions **in one plot** using loglog-scale and verify that all the distributions have heavier tails than Gaussian.
 - Based on the plots, roughly **estimate** the 90th percentiles of the degree, strength, and weight distributions.

Hints:

- For reading in the network from the given datafile, use `net = nx.read_weighted_edgelist`
- See the binning tutorial for help on computing the 1-CDFs.
- For getting node strengths, use `strengths = nx.degree(net, weight="weight")`
- As a check, calculate the maximum degree k_{\max} and maximum strength s_{\max} . If everything is running okay so far, you should be getting $k_{\max} = 255$ and $s_{\max} = 1546.0$.

- b) (5 pts) Next, we will study how the average link weight per node $\langle w \rangle = \frac{s}{k}$ behaves as a function of the node degree k .

¹Data originally from <http://toreopsahl.com/datasets/>

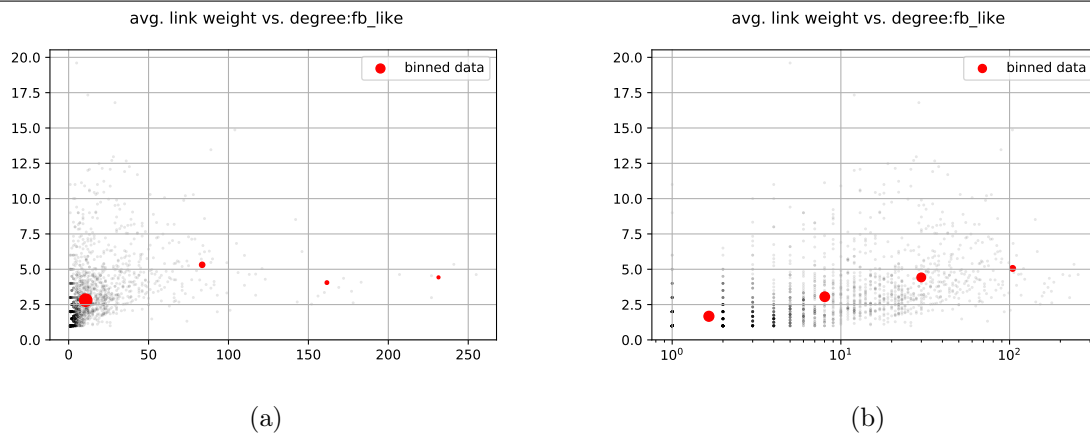


Figure 1: Example of $\langle w \rangle$ as a function of k using another Facebook-like social network data. 1a: linear axes. 1b: logarithmic axes.

- Compute s , k , and $\langle w \rangle = \frac{s}{k}$ for each node.
- **Make scatter plots** of all the data points of $\langle w \rangle$ as a function of k . Create two versions of the plots: one with linear and one with logarithmic x -axes.
- The large variance of the data can make the scatter plots a bit messy. To make the relationship between $\langle w \rangle$ and k more visible, **create bin-averaged versions** of the plots, *i.e.* divide nodes into bins based on their degree and calculate the average $\langle w \rangle$ in each bin. Plot the bin-averaged versions on top of the scatter plots.
- Based on the plots, which of the two approaches (linear or logarithmic x -axes) suits better for presenting $\langle w \rangle$ as a function of k ? Why?

Hints:

- For the bin-averaged plots, use bins that are consistent with the scale of the x -axis: bins with constant width for the linear scale and logarithmic bins for the logarithmic scale. If in trouble, see the binning tutorial for help.
 - An example of how the scatter and bin-averaged plots may look like is shown in Fig. 1. Obviously, the number of bins are too few in these plots. Typically, it is better to use too many than too few bins. A good choice for the number of bins in this case would be 20.
 - Check if each degree bin includes roughly the same number of nodes. Nonequal distribution of observations may obscure the results.
- c) (4 pts) Lets consider a link between nodes i and j . For this link, *link neighborhood overlap* O_{ij} is defined as the fraction of common neighbors of i and j out of all their neighbors:

$$O_{ij} = \frac{n_{ij}}{(k_i - 1) + (k_j - 1) - n_{ij}}, \quad (1)$$

where n_{ij} denotes the number of common neighbors.

According to the Granovetter hypothesis, link neighborhood overlap is an increasing function of link weight in social networks. Your task is to find out whether this is the case also for the present data set by visualizing it in an appropriate way. Use the binning strategy (linear or logarithmic) that is most suitable for this case.

- Calculate the link neighborhood overlap for each link.
- **Create a scatter plot** showing the overlap for every link as a function of link weight.
- As in b), **produce also a bin-averaged version** of the plot on top of the scatter plot. You should find a reasonable number of bins, which, in this case, will be between 10 and 30. Specify how many bins you used.
- In the end, you should be able to spot a subtle trend in the data. Based on your plot, show that the trend is in accordance with the Granovetter hypothesis.

2. Network thresholding and spanning trees: the case of US air traffic (6 pts)

In this exercise, we will get familiar with different approaches to thresholding networks, and also learn how they can be used for efficiently visualizing networks. Now, you are given an undirected network describing the US Air Traffic between 14th and 23rd December 2008². In the network, each node corresponds to an airport and link weights describe the number of flights between the airports during the time period.

The data and some code for visualizing the network are provided at the course web-page. The network is given in the file `aggregated_US_air_traffic_network_undir.edg`, and `us_airport_id_info.csv` contains information about names and locations of the airports. In this exercise, you may also freely use all available `networkx` functions.

- a) (1 pt) When facing a new network, it is always good to first get some idea on what the network is like. Thus, **compute** and list the following basic network properties:
- Number of network nodes N and density D
 - Network diameter d
 - Average clustering coefficient C

Hints:

- To make sure that the network is loaded correctly, calculate the number of links L . You should be getting $L = 2088$.
 - For the clustering coefficient, consider the unweighted version of the network, where two airports are linked if there is a flight between them.
- b) (1 pt) **Visualize** the full network with all links on top of the map of USA. The resulting figure is somewhat messy due to the large number of visible links.
- c) (2 pts) In order to reduce the number of plotted links, **compute** both the *maximal* and *minimal spanning trees* of the network and **visualize** them. Then, answer the following question:
- If you would like to understand the overall organization of the air traffic in the US, would you use the minimal or maximal spanning tree? Why?

²Data from <http://www.rita.dot.gov/bts/>

Hint: For computing minimum spanning trees, use `nx.minimum_spanning_tree`. For computing maximum spanning trees, use `nx.maximum_spanning_tree`.

- d) (2 pts) **Threshold and visualize** the network by taking only the strongest M links into account, where $M = N - 1$ is the number of links in the maximal spanning tree. Then, answer following questions.
- How many of the links in the thresholded network are the same as those in the maximal spanning tree?
 - Given this number and the visualizations, does simple thresholding yield a similar network as the maximal spanning tree?

Feedback (1 pt)

To earn one bonus point, give feedback on this exercise set and the corresponding lecture latest two days after the report's submission deadline.

Link to the feedback form: <https://forms.gle/ZhZ4zhAPKv6bUGPQA>.

References