

$$\frac{16.25}{17}$$

CS-E5740 Complex Networks, Answers to exercise set 8

Marco Di Francesco, Student number: 100632815

November 15, 2022

Problem 1

a) Look at table 1.

samp.	triangles	two-stars	transit.	triang.frac.	two-st.frac.
node	345	472	0.7309	0.0135	0.0128
edge	258	1802	0.1432	0.0101	0.0491
star	3174	7653	0.4147	0.1244	0.2083
orig.	25524	36732	0.6949	1.0000	1.0000

Figure 1: Table exercise 1 section a

*How do sampling schemes compare in the fraction of triangles/two-stars they preserve?
Do sampling schemes affect two-stars/triangles in the same way?*

Sampling shemes given that they work in different ways have very different results for the number of two-starts and triangles. Explanation for each sampling sheme:

- Node sampling: we have very little number of two-stars and triangles (less than 2%) because in general we have very little edges. This is given because we have an edge only in the case we select both nodes it connects, thus the probability an edge exists is very small.
- Edge sampling: we have an higher number of edges compared to node sampling, thus we have a higher number of two-stars fraction, but because the selecting the edges is random (not like in star sampling) the number of triangles is still very small.
- Star sampling: we have way more edges compared to the other sampling techniques, and they are not selected randomly but each node selectes also its neighbors, thus we have not only a high number of two-stars included, but also a high number of triangles.

Transitivity via node sampling should be similar to the real value, what could be the reason for this?

As explained above, the edges are sampled randomly in all cases, and it's way more likely to get two-stars than triangles. This is because getting 3 nodes that constitute a triangle is way more unlikely than getting 2 nodes that constitute a two-stars.

- b)
- Node sampling two-stars: p^3 . This is because to get an edge all 3 nodes must be selected.
 - Node sampling triangles: p^3 . As for two-stars, because to get an edge all 3 nodes must be selected.
 - Star sampling two-star: p . This is because for each node we get we also have the neighbors. (Assumption: we start from 55 cliques of 12 nodes and rewiring only 10% of the nodes we are not considering the degree of each node)
 - Star sampling triangles: $\binom{3}{2}p^2(1-p) + p^3$. This is because we need a binomial to select 2 out of 3 neighbor nodes to create a triangle, or we are just choosing all 3 nodes (p^3).

- c) *Why the HT estimator corrects for sampling bias?*

The HT estimator corrects for sampling bias because it normalizes the count of two-stars or three stars by the probability of observing them after sampling. This probability is distorted and will not always be p .

HT estimators for the number of two-stars, triangles and transitivity for the three sampling schemes we explored

Two-stars

- Nodes $\hat{\tau}_{\angle}^{HT} = \hat{\tau}_{\angle}^n / p^3$
- Edges $\hat{\tau}_{\angle}^{HT} = \hat{\tau}_{\angle}^e / p^2$
- Stars $\hat{\tau}_{\angle}^{HT} = \hat{\tau}_{\angle}^s / p$

Triangles

- Nodes $\hat{\tau}_{\Delta}^{HT} = \hat{\tau}_{\Delta}^n / p^3$
- Edges $\hat{\tau}_{\Delta}^{HT} = \hat{\tau}_{\Delta}^e / p^3$
- Stars $\hat{\tau}_{\Delta}^{HT} = \hat{\tau}_{\Delta}^s / (\hat{\tau}_{\angle}^s \binom{3}{2} p^2 (1-p) + p^3)$

Transitivity

- Nodes $\hat{\tau}_C^{HT} = \hat{\tau}_{\Delta}^n / \hat{\tau}_{\angle}^n$
- Edges $\hat{\tau}_C^{HT} = \hat{\tau}_{\Delta}^e / (p \hat{\tau}_{\angle}^e)$
- Stars $\hat{\tau}_C^{HT} = \hat{\tau}_{\Delta}^s / (p \hat{\tau}_{\angle}^s \binom{3}{2} (1-p) + p)$

d) Looks plots in figure 2.

What is the effect of the sampling probability p on your estimators?

Higher probability raises variance in the prediction. We can see that the distributions with $p = 0.35$ have higher variance than with $p = 0.5$.

How the HT distributions differ between sampling schemes? In all 3 cases the HT distributions were able to align the mean value the estimation, the difference lies on the variance these distribution have. Having lower probability of having more two-stars and triangles (lower bias) led to having less variance.

In the last plot, the empirical estimators (without HT correction) should be centered around the true value, why could this be?

This happens because there is no bias for $p = 0.5$, thus the distributions are centered around the true value.

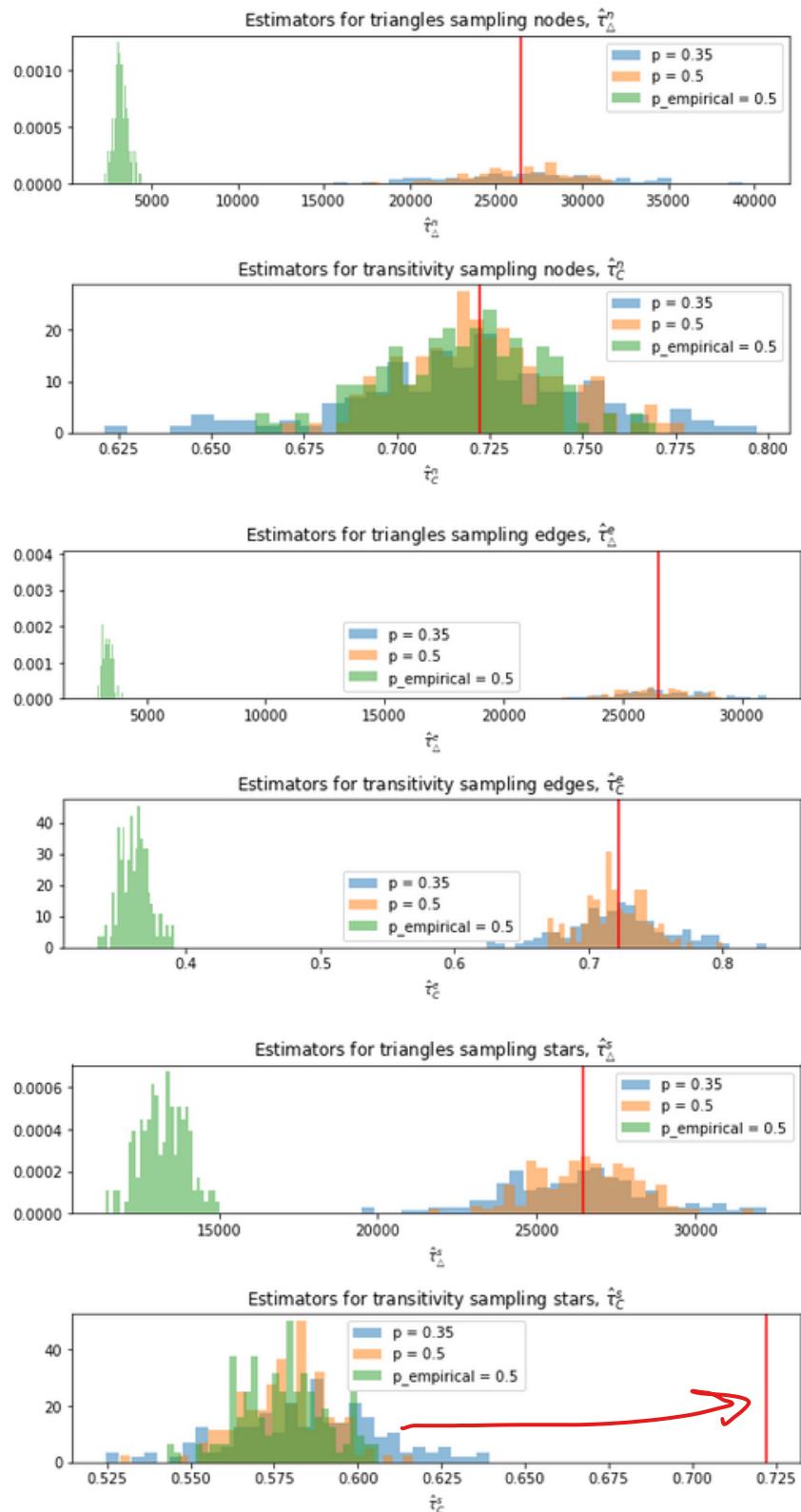


Figure 2: Exercise 1 part d

Problem 2

a) Look at image 3.

$$\text{i) yellow} \quad \text{pink}$$

$$Q = \sum \left(\frac{6}{13} - \left(\frac{13}{2 \cdot 13} \right)^2 \right) + \left(\frac{6}{13} - \left(\frac{13}{2 \cdot 13} \right)^2 \right) = \frac{11}{26}$$

$$\text{ii) green}$$

$$Q = \left(\frac{13}{13} - \left(\frac{26}{2 \cdot 13} \right)^2 \right) = 0$$

Figure 3: Exercise 2 part a

b) Look at images 4 and 5.

$$\Delta Q = Q_2 - Q_1 = Q_{ab} + Q_b - Q_a - Q_a - Q_b =$$

$$= Q_{ab} - Q_a - Q_b$$

Figure 4: Exercise 2 part b image 1

$$\Delta Q = \left(\frac{l_{ab}}{L} - \left(\frac{d_{ab}}{2L} \right)^2 \right) - \left(\frac{l_a}{L} - \left(\frac{d_a}{2L} \right)^2 \right) - \left(\frac{l_b}{L} - \left(\frac{d_b}{2L} \right)^2 \right)$$

$$\xrightarrow{\text{similar terms}}$$

$$= \frac{d_a^2 + d_b^2 - d_{ab}^2 - 4L(l_a + l_b - l_{ab})}{4L^2} =$$

$$= \frac{1}{4L^2} (d_a^2 + d_b^2 - d_{ab}^2) - \frac{4L}{4L^2} (l_a + l_b - l_{ab})$$

Figure 5: Exercise 2 part b image 2

c) Look at image 6.



Diagram:

Degree: is the sum
 $d_{ab} = d_a + d_b$
 Internal links: requires $+1$ for connecting edge
 $l_{ab} = l_a + l_b + 1$

$$\Rightarrow \Delta Q = \frac{1}{4L^2} (d_a^2 + d_b^2 - (d_a + d_b)^2) + \frac{1}{L} (l_a + l_b - l_a l_b + 1)$$

$$= \frac{1}{4L^2} (d_a^2 + d_b^2 - d_a^2 - d_b^2 - 2d_a d_b) + \frac{1}{L} =$$

$$= -\frac{1}{2L^2} d_a d_b + \frac{1}{L} =$$

$$= -\frac{1}{2L^2} d_a d_b + \frac{1}{L} =$$

~~Average~~

$$= \frac{1}{L} \left(1 - \frac{1}{2L} d_a d_b \right)$$

Merge the clusters with $\Delta Q > 0$

$$\frac{1}{L} \left(1 - \frac{d_a d_b}{2L} \right) > 0 \Rightarrow$$

$$\Rightarrow \frac{1}{L} - \frac{d_a d_b}{2L^2} > 0 \Rightarrow$$

$$\Rightarrow \frac{d_a d_b}{2L^2} < \frac{1}{L} \cdot L \Rightarrow$$

$$\Rightarrow \frac{d_a d_b}{2} < L \Rightarrow L > \frac{d_a d_b}{2}$$

Figure 6: Exercise 2 part c

Explain why this result indicates that modularity optimization must have a resolution limit, or a minimum size of community that can be found that depends on the size of the network

Increasing the number of links increases we need to have a resolution limit because otherwise the number of edges required to connect 2 communities drops below 1, thus connecting clusters of nodes that are in reality weakly connected.

For a large
two "small" communities
are merged

d) Look at image 7

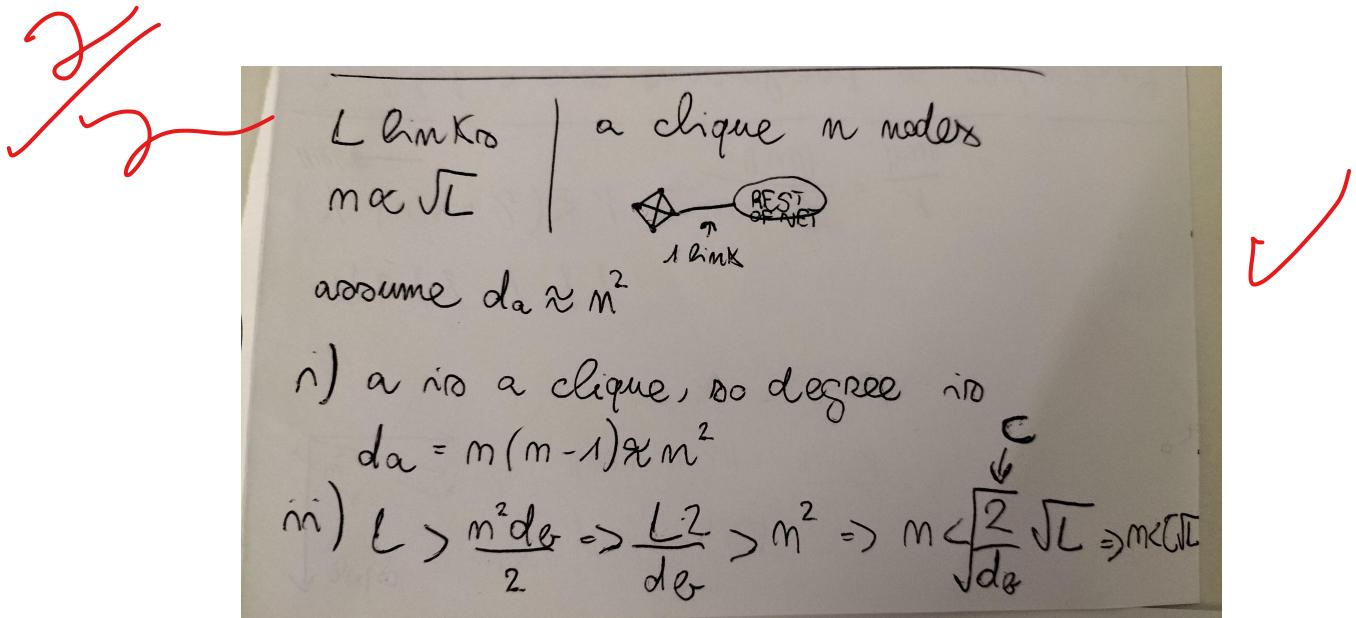


Figure 7: Exercise 2 part d