

Impact of Censorship in Social Media

Kacper Krasowski

June 30, 2025

Outline

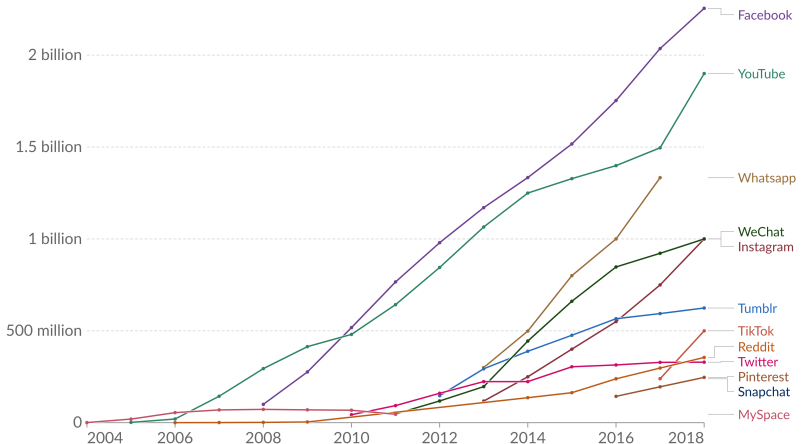
- 1 Introduction
- 2 Methodology
- 3 Results

Introduction

Our World
in Data

Number of people using social media platforms, 2004 to 2018

Estimates correspond to monthly active users (MAUs). Facebook, for example, measures MAUs as users that have logged in during the past 30 days. See source for more details.



Source: Statista and TNW (2019)

OurWorldInData.org/internet • CC BY



3/27

Introduction

- **Research Question**
 - What are the implications of limiting freedom of speech on social media for overall welfare?

Introduction

- **Research Question**

- What are the implications of limiting freedom of speech on social media for overall welfare?

- **Motivation**

- Great and growing impact of social media on opinion dynamics
- Increasing interest of regulators on social media (e.g. The Digital Services Act)
- Existing gap in the literature

Introduction

- **The model:**

- Agents' **opinions** are **repeated weighted averages** of their neighbours' opinions and their initial opinion.
- **Censorship** introduced by **banning** agents with **extreme** enough **opinions**.
- Welfare investigated in terms of:
 - **Polarization** - How much opinions differ in a network.
 - **Disagreement** - How much opinions differ among connected agents.
 - **Internal Conflict** - How much opinions have evolved.
 - **Mix** of the above indices.

Introduction

- **The model:**

- Agents' **opinions** are **repeated weighted averages** of their neighbours' opinions and their initial opinion.
- **Censorship** introduced by **banning** agents with **extreme** enough **opinions**.
- Welfare investigated in terms of:
 - **Polarization** - How much opinions differ in a network.
 - **Disagreement** - How much opinions differ among connected agents.
 - **Internal Conflict** - How much opinions have evolved.
 - **Mix** of the above indices.

- **Preview of results:**

- Higher **censorship** increases **polarization** and **reduces internal conflict**.
- The impact on **disagreement** **depends** on the network size and connectivity.
- There **exists** a **optimal censorship**.

Literature Review

- **Opinion Dynamics**

- DeGroot (1974)
- Friedkin and Johnsen (1997)
- Golub and Jackson (2010)
- Cameron Musco, Christopher Musco, and Tsourakakis (2018)

Literature Review

- **Opinion Dynamics**

- DeGroot (1974)
- Friedkin and Johnsen (1997)
- Golub and Jackson (2010)
- Cameron Musco, Christopher Musco, and Tsourakakis (2018)

- **Freedom of speech in social media**

- Dehghan (2018)
- Segado-Boj and Campo Lozano (2020)

Literature Review

- **Opinion Dynamics**

- DeGroot (1974)
- Friedkin and Johnsen (1997)
- Golub and Jackson (2010)
- Cameron Musco, Christopher Musco, and Tsourakakis (2018)

- **Freedom of speech in social media**

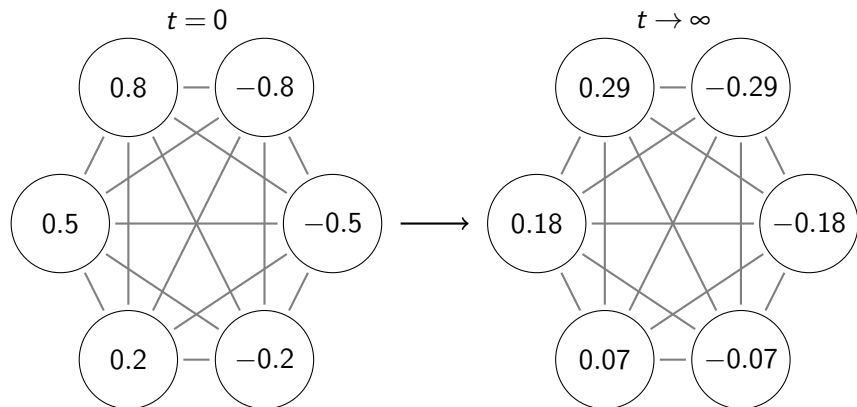
- Dehghan (2018)
- Segado-Boj and Campo Lozano (2020)

- **Impact of Social Media on Welfare**

- Allcott et al. (2020)

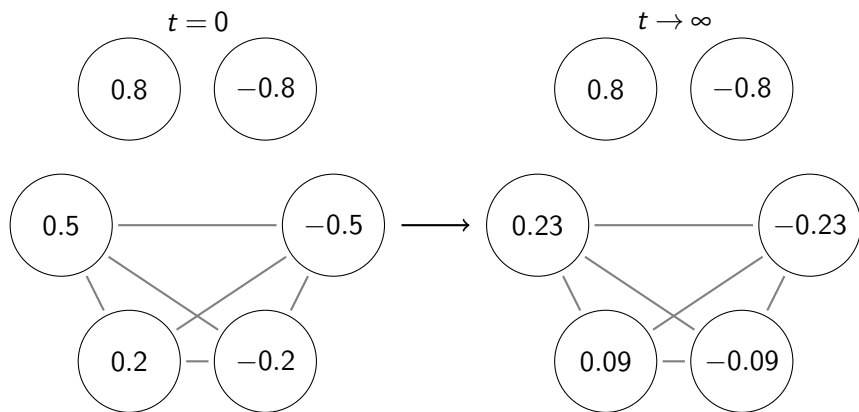
Example

Is climate change real?



Example

Is climate change real?



Set up

- Opinions $\bar{z}_i \in [-1, 1]$
 - $\bar{z}_i = 1$ - totally agree
 - $\bar{z}_i = 0$ - neutral
 - $\bar{z}_i = -1$ - totally disagree

Set up

- Opinions $\bar{z}_i \in [-1, 1]$
 - $\bar{z}_i = 1$ - totally agree
 - $\bar{z}_i = 0$ - neutral
 - $\bar{z}_i = -1$ - totally disagree
- Agents interact in a social network, represented by $G(V, E)$, where:
 - V is the set of nodes (agents)
 - E is set of edges (connections)

Set up

- Opinions $\bar{z}_i \in [-1, 1]$
 - $\bar{z}_i = 1$ - totally agree
 - $\bar{z}_i = 0$ - neutral
 - $\bar{z}_i = -1$ - totally disagree
- Agents interact in a social network, represented by $G(V, E)$, where:
 - V is the set of nodes (agents)
 - E is set of edges (connections)
- Two types of n agents:
 - Type 1 draws innate opinion \bar{s}_i from distribution with mean $-\mu$ and variance σ^2 .
 - Type 2 draws innate opinion \bar{s}_i from distribution with mean μ and variance σ^2 .

Set up

- Opinions $\bar{z}_i \in [-1, 1]$
 - $\bar{z}_i = 1$ - totally agree
 - $\bar{z}_i = 0$ - neutral
 - $\bar{z}_i = -1$ - totally disagree
- Agents interact in a social network, represented by $G(V, E)$, where:
 - V is the set of nodes (agents)
 - E is set of edges (connections)
- Two types of n agents:
 - Type 1 draws innate opinion \bar{s}_i from distribution with mean $-\mu$ and variance σ^2 .
 - Type 2 draws innate opinion \bar{s}_i from distribution with mean μ and variance σ^2 .
- Agents opinions evolve according to:

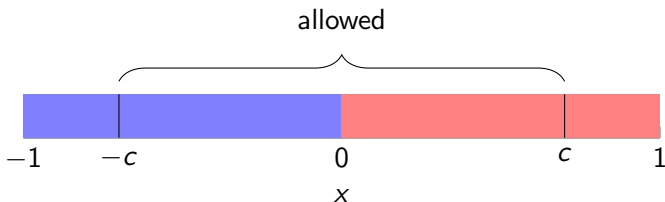
$$\bar{z}_i^{(t)} = \frac{\bar{s}_i + \sum_j a_{ij} \bar{z}_j^{(t-1)}}{1 + d_i}, \quad \bar{\mathbf{z}} = (\mathbf{I} + \mathbf{D} - \mathbf{A})^{-1} \bar{\mathbf{s}}$$

Censorship and Dynamics

- **Network Administrator** decides a threshold $[-c, c]$ of allowed opinions.
 - $c = 1$ - full freedom of speech
 - $c = 0$ - full censorship

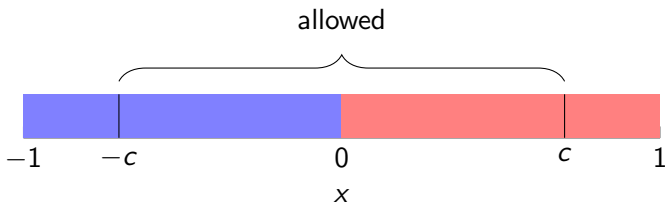
Censorship and Dynamics

- **Network Administrator** decides a threshold $[-c, c]$ of allowed opinions.
 - $c = 1$ - full freedom of speech
 - $c = 0$ - full censorship



Censorship and Dynamics

- **Network Administrator** decides a threshold $[-c, c]$ of allowed opinions.
 - $c = 1$ - full freedom of speech
 - $c = 0$ - full censorship



- The timings are as follows:
 - 1 Network Administrator learns the types, and intensity.
 - 2 Network Administrator decides the censor point.
 - 3 Opinions get drawn and evolve until they reach equilibrium.

Welfare Measures

Polarization - Variance of a set of opinions.

$$\mathcal{P} = \sum_{i=1}^n \bar{z}_i^2$$

Welfare Measures

Polarization - Variance of a set of opinions.

$$\mathcal{P} = \sum_{i=1}^n \bar{z}_i^2$$

Disagreement - How much opinions differ in the network.

$$\mathcal{D} = \sum_{i=1}^n \sum_{j \in N(i)} a_{ij} (\bar{z}_i - \bar{z}_j)^2$$

Welfare Measures

Polarization - Variance of a set of opinions.

$$\mathcal{P} = \sum_{i=1}^n \bar{z}_i^2$$

Disagreement - How much opinions differ in the network.

$$\mathcal{D} = \sum_{i=1}^n \sum_{j \in N(i)} a_{ij} (\bar{z}_i - \bar{z}_j)^2$$

Internal Conflict - How much opinions differ from the innate ones.

$$\mathcal{IC} = \sum_{i=1}^n (\bar{z}_i - \bar{s}_i)^2$$

Welfare Measures

Polarization - Variance of a set of opinions.

$$\mathcal{P} = \sum_{i=1}^n \bar{z}_i^2$$

Disagreement - How much opinions differ in the network.

$$\mathcal{D} = \sum_{i=1}^n \sum_{j \in N(i)} a_{ij} (\bar{z}_i - \bar{z}_j)^2$$

Internal Conflict - How much opinions differ from the innate ones.

$$\mathcal{IC} = \sum_{i=1}^n (\bar{z}_i - \bar{s}_i)^2$$

Objective - Convex combination of the above indices.

$$\mathcal{W} = -[\alpha \mathcal{P} + \beta \mathcal{D} + (1 - \alpha - \beta) \mathcal{IC}]$$

Methodology

- Each agent talk with one another with intensity p .
- The point of interest is the perspective of the Network administrator.
- As the network administrator does not know which agents have what opinions, while deciding on c , from hers perspective each agent has a probability $F(c)$ to stay in the social network.
- Then the equilibrium opinions are given by:

$$\bar{z} = (\mathbf{I} + \mathbf{D} - \mathbf{A})^{-1} \bar{s} = (\mathbf{I} + F(c)n\mathbf{p}\mathbf{I} - F(c)p\mathbf{J})^{-1} \bar{s}$$

Results ▶

$$\mathbb{E}[\mathcal{P}] = \sigma^2 + \frac{1}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{D}] = \frac{F(c)np}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{SC}] = \frac{(F(c)np)^2}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

Results ▶

$$\mathbb{E}[\mathcal{P}] = \sigma^2 + \frac{1}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{D}] = \frac{F(c)np}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{SC}] = \frac{(F(c)np)^2}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

- ❶ **Polarization is increasing** in censorship

Results ▶

$$\mathbb{E}[\mathcal{P}] = \sigma^2 + \frac{1}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{D}] = \frac{F(c)np}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{I}\mathcal{C}] = \frac{(F(c)np)^2}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

- 1 **Polarization** is **increasing** in censorship
- 2 **Disagreement** is:
 - **decreasing** in censorship if $\frac{1}{np} < 1$
 - initially **increasing** and then **decreasing** otherwise

Results ▶

$$\mathbb{E}[\mathcal{P}] = \sigma^2 + \frac{1}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

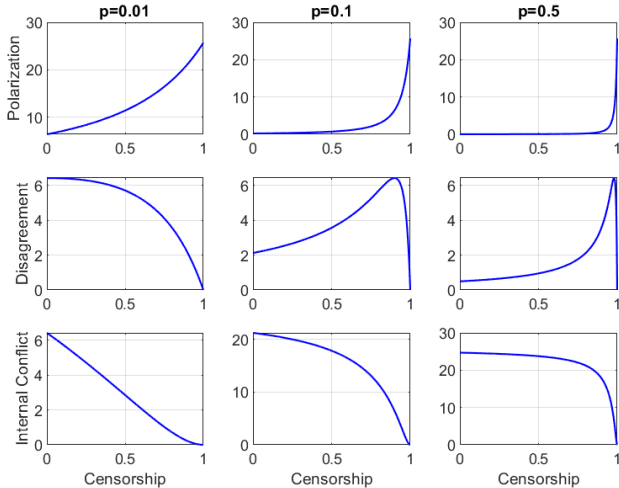
$$\mathbb{E}[\mathcal{D}] = \frac{F(c)np}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

$$\mathbb{E}[\mathcal{IC}] = \frac{(F(c)np)^2}{(1 + F(c)np)^2} (n\mu^2 + (n-1)\sigma^2)$$

- ❶ **Polarization** is **increasing** in censorship
- ❷ **Disagreement** is:
 - **decreasing** in censorship if $\frac{1}{np} < 1$
 - initially **increasing** and then **decreasing** otherwise
- ❸ **Internal conflict** is **decreasing** in censorship

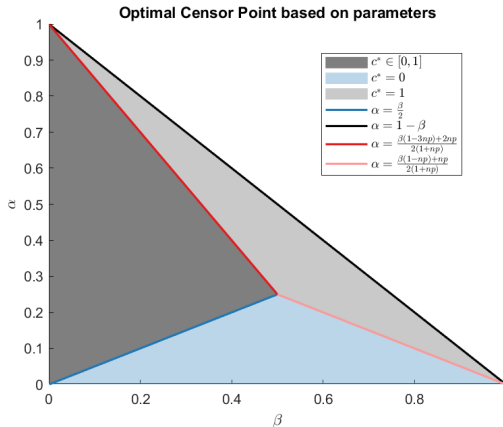
Results

Polarization, Disagreement, and Internal Conflict



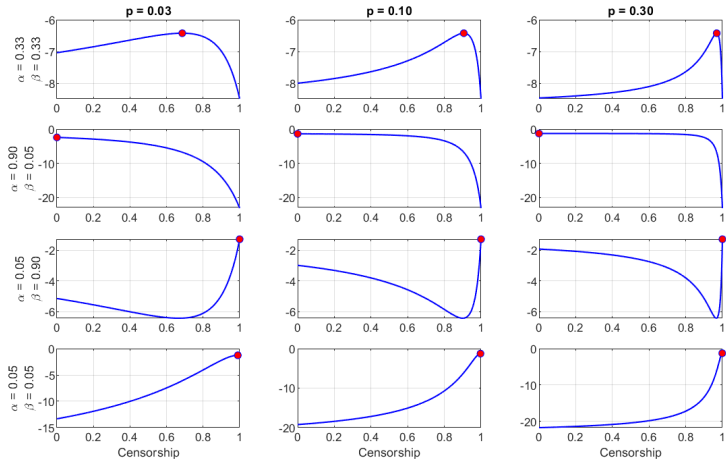
Results

If the none of the indexes have too much weight on it in the welfare function, then $c^* = Q(\frac{1}{np} \frac{2\alpha - \beta}{2 - 2\alpha - 3\beta})$ is the optimal censoring point. Otherwise its either 0 or 1.



Results

Welfare



Extensions

No cross types communication

- **Polarization** is **increasing** in censorship
- **Disagreement** is:
 - **decreasing** in censorship if $\frac{2}{np} < 1$
 - initially **increasing** and then **decreasing** otherwise
- **Internal conflict** is decreasing in **censorship**
- $c^* = Q\left(\frac{2}{np} \frac{2\alpha - \beta}{2 - 2\alpha - 3\beta}\right)$
- Minor changes in indices caused by censorship

Extensions

No cross types communication

- **Polarization** is **increasing** in censorship
- **Disagreement** is:
 - **decreasing** in censorship if $\frac{2}{np} < 1$
 - initially **increasing** and then **decreasing** otherwise
- **Internal conflict** is decreasing in **censorship**
- $c^* = Q(\frac{2}{np} \frac{2\alpha - \beta}{2 - 2\alpha - 3\beta})$
- Minor changes in indices caused by censorship

Different intensities across types

- **Polarization** is **increasing** in censorship
- **Internal conflict** is decreasing in **censorship**
- **Optimal censoring point** could be found by **numerical methods**.

Conclusions

- **Objective:** Investigate the effects of censorship in social media

Conclusions

- **Objective:** Investigate the effects of censorship in social media
- **Model:** Two types of agents with opposing opinions; censorship introduced by banning agents with opinions outside a set threshold. Opinions evolve according to Friedkin-Johansen model
- **Results:**
 - Higher censorship levels generally increase polarization and reduce internal conflict.
 - Disagreement initially rises with censorship but eventually decreases to zero.
 - Optimal censorship thresholds identified for specific parameters.

Avenue for Future work

- Investigate the incentives of network administrators, focusing on profit maximization rather than welfare.
- Incorporate censorship costs and user activity into profit functions.
- Explore empirical validation using data from platforms like Twitter or Reddit.

Thank You!

$$\bar{z}_i^{(t)} = \frac{\bar{s}_i + \sum_j a_{ij} \bar{z}_j^{(t-1)}}{1 + d_i} \quad (1)$$

$$\bar{\mathbf{z}}^{(t)} = (\mathbf{I} + \mathbf{D})^{-1}(\bar{\mathbf{s}} + \mathbf{A}\bar{\mathbf{z}}^{(t-1)}) \quad (2)$$

$$\bar{\mathbf{z}} = (\mathbf{I} + \mathbf{D})^{-1}(\bar{\mathbf{s}} + \mathbf{A}\bar{\mathbf{z}}) \quad (3)$$

$$(\mathbf{I} + \mathbf{D})\bar{\mathbf{z}} = (\bar{\mathbf{s}} + \mathbf{A}\bar{\mathbf{z}}) \quad (4)$$

$$(\mathbf{I} + \mathbf{D} - \mathbf{A})\bar{\mathbf{z}} = \bar{\mathbf{s}} \quad (5)$$

$$\bar{\mathbf{z}} = (\mathbf{I} + \mathbf{D} - \mathbf{A})^{-1}\bar{\mathbf{s}} \quad (6)$$

Appendix ▶

$$\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} p & \dots & p \\ \vdots & \ddots & \vdots \\ p & \dots & p \end{bmatrix} = p\mathbf{J}$$

$$\mathbf{D} = \begin{bmatrix} d_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & d_n \end{bmatrix} = \begin{bmatrix} \sum_i^n a_{1i} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sum_i^n a_{ni} \end{bmatrix} = \begin{bmatrix} np & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & np \end{bmatrix} = np\mathbf{I}$$

Appendix

$$\mathcal{P} = \sum_{i=1}^n \bar{z}_i^2 = \bar{\mathbf{z}}^T \bar{\mathbf{z}} = \bar{\mathbf{s}}^T (\mathbf{I} + \mathbf{L})^{-2} \bar{\mathbf{s}}$$

$$\begin{aligned} \mathcal{D} &= \sum_{i=1}^n \sum_{j \in N(i)} a_{ij} (\bar{z}_i - \bar{z}_j)^2 = \bar{\mathbf{z}}^T \mathbf{L} \bar{\mathbf{z}} \\ &= \bar{\mathbf{s}}^T (\mathbf{I} + \mathbf{L})^{-1} \mathbf{L} (\mathbf{I} + \mathbf{L})^{-1} \bar{\mathbf{s}} \end{aligned}$$

$$\mathcal{IC} = \sum_{i=1}^n (\bar{z}_i - \bar{s}_i)^2 = \bar{\mathbf{s}}^T (\mathbf{I} + \mathbf{L})^{-1} \mathbf{L}^2 (\mathbf{I} + \mathbf{L})^{-1} \bar{\mathbf{s}}$$

$$\mathcal{P} = \bar{\mathbf{s}}^T f_{\mathcal{P}}(\mathbf{L}) \bar{\mathbf{s}},$$

$$\mathcal{D} = \bar{\mathbf{s}}^T f_{\mathcal{D}}(\mathbf{L}) \bar{\mathbf{s}},$$

$$\mathcal{IC} = \bar{\mathbf{s}}^T f_{\mathcal{IC}}(\mathbf{L}) \bar{\mathbf{s}}$$

Appendix

$$\begin{aligned}\mathbb{E}[\mathcal{P}] &= \mathbb{E}[\bar{\mathbf{s}}^T f_{\mathcal{P}}(\mathbf{L}) \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}^T \mathbf{U} \Lambda_{\mathcal{P}} \mathbf{U}^T \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}_{\mathbf{U}}^T \Lambda_{\mathcal{P}} \bar{\mathbf{s}}_{\mathbf{U}}] \\ &= \sum_i^n \lambda_{\mathcal{P}i} \mathbb{E}[\bar{s}_{Ui}^2] = \sum_i^n \frac{1}{(1 + \lambda_{Li})^2} \mathbb{E}[\bar{s}_{Ui}^2]\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathcal{D}] &= \mathbb{E}[\bar{\mathbf{s}}^T f_{\mathcal{D}}(\mathbf{L}) \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}^T \mathbf{U} \Lambda_{\mathcal{D}} \mathbf{U}^T \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}_{\mathbf{U}}^T \Lambda_{\mathcal{D}} \bar{\mathbf{s}}_{\mathbf{U}}] \\ &= \sum_i^n \lambda_{\mathcal{D}i} \mathbb{E}[\bar{s}_{Ui}^2] = \sum_i^n \frac{\lambda_{Li}}{(1 + \lambda_{Li})^2} \mathbb{E}[\bar{s}_{Ui}^2]\end{aligned}$$

$$\begin{aligned}\mathbb{E}[\mathcal{IC}] &= \mathbb{E}[\bar{\mathbf{s}}^T f_{\mathcal{IC}}(\mathbf{L}) \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}^T \mathbf{U} \Lambda_{\mathcal{IC}} \mathbf{U}^T \bar{\mathbf{s}}] = \mathbb{E}[\bar{\mathbf{s}}_{\mathbf{U}}^T \Lambda_{\mathcal{IC}} \bar{\mathbf{s}}_{\mathbf{U}}] \\ &= \sum_i^n \lambda_{\mathcal{IC}i} \mathbb{E}[\bar{s}_{Ui}^2] = \sum_i^n \frac{\lambda_{Li}^2}{(1 + \lambda_{Li})^2} \mathbb{E}[\bar{s}_{Ui}^2]\end{aligned}$$

Appendix

It is common knowledge in the SBM literature that the eigenvalues of **A** are:

$$\lambda_{A1} = F(c)n\frac{p+q}{2}, \lambda_{A2} = F(c)n\frac{p-q}{2}, \lambda_{A3} = \dots \lambda_{An} = 0$$

And the eigenvalues of **D** are:

$$\lambda_{D1} = \lambda_{D2} = \dots \lambda_{Dn} = F(c)n\frac{p+q}{2}$$

, then the eigenvalues of **L** are:

$$\lambda_{L1} = 0, \lambda_{L2} = F(c)qn, \lambda_{L3} = \dots = \lambda_{Ln} = F(c)n\frac{p+q}{2}$$

Appendix

And from Chen, Lijffijt, and De Bie (2018) we know the mappings from eigenvalues of \mathbf{L} to eigenvalues of each index:

$$\textcircled{1} \lambda_{\mathcal{P}} = \frac{1}{(1+\lambda_L)^2}$$

$$\textcircled{2} \lambda_{\mathcal{D}} = \frac{\lambda_L}{(1+\lambda_L)^2}$$

$$\textcircled{3} \lambda_{\mathcal{J}\mathcal{C}} = \frac{\lambda_L^2}{(1+\lambda_L)^2}$$

,thus:

$$\lambda_{\mathcal{P}1} = 1, \lambda_{\mathcal{P}2} = \frac{1}{(1+F(c)qn)^2}, \lambda_{\mathcal{P}3} = \dots = \lambda_{\mathcal{P}n} = \frac{1}{(1+F(c)\frac{p+q}{2}n)^2},$$

$$\lambda_{\mathcal{D}1} = 0, \lambda_{\mathcal{D}2} = \frac{F(c)qn}{(1+F(c)qn)^2}, \lambda_{\mathcal{D}3} = \dots = \lambda_{\mathcal{D}n} = \frac{F(c)\frac{p+q}{2}n}{(1+F(c)\frac{p+q}{2}n)^2},$$

$$\lambda_{\mathcal{J}\mathcal{C}1} = 0, \lambda_{\mathcal{J}\mathcal{C}2} = \frac{(F(c)qn)^2}{(1+F(c)qn)^2}, \lambda_{\mathcal{J}\mathcal{C}3} = \dots = \lambda_{\mathcal{J}\mathcal{C}n} = \frac{(F(c)\frac{p+q}{2}n)^2}{(1+F(c)\frac{p+q}{2}n)^2},$$

Appendix

$$\mathbb{E}[\mathcal{P}] = \sigma^2 + \frac{1}{(1 + F(c)nq)^2}(n\mu^2 + \sigma^2) + (n-2)\frac{1}{(1 + F(c)n\frac{p+q}{2})^2}\sigma^2$$

$$\mathbb{E}[\mathcal{D}] = \frac{F(c)nq}{(1 + F(c)nq)^2}(n\mu^2 + \sigma^2) + (n-2)\frac{F(c)n\frac{p+q}{2}}{(1 + F(c)n\frac{p+q}{2})^2}\sigma^2$$

$$\mathbb{E}[\mathcal{IC}] = \frac{(F(c)nq)^2}{(1 + F(c)nq)^2}(n\mu^2 + \sigma^2) + (n-2)\frac{(F(c)n\frac{p+q}{2})^2}{(1 + F(c)n\frac{p+q}{2})^2}\sigma^2$$