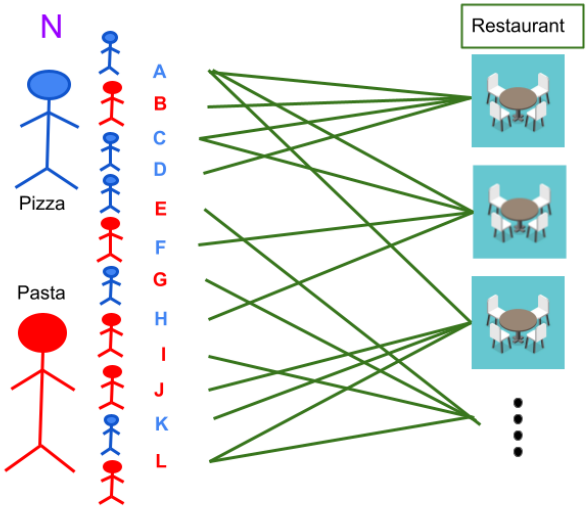# Community detection on multilayer hypergraphs using the aggregate similarity matrix
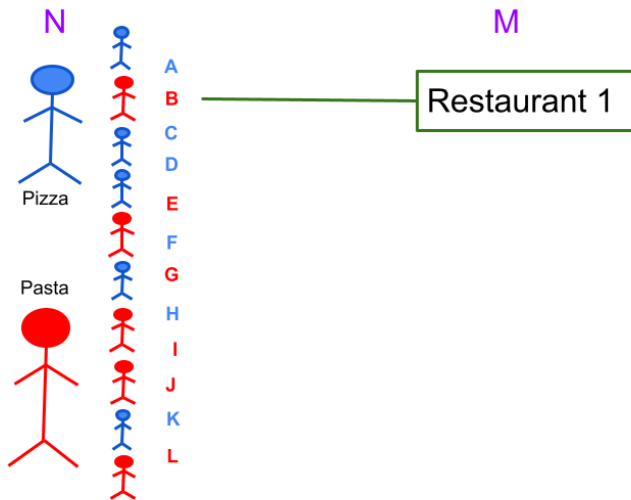
Kalle Alaluusua, Konstantin Avrachenkov, **B R Vinay Kumar**, Lasse Leskelä

# Motivating example:



N

Pizza

Pasta

A
B
C
D
E
F
G
H
I
J
K
L

Restaurant
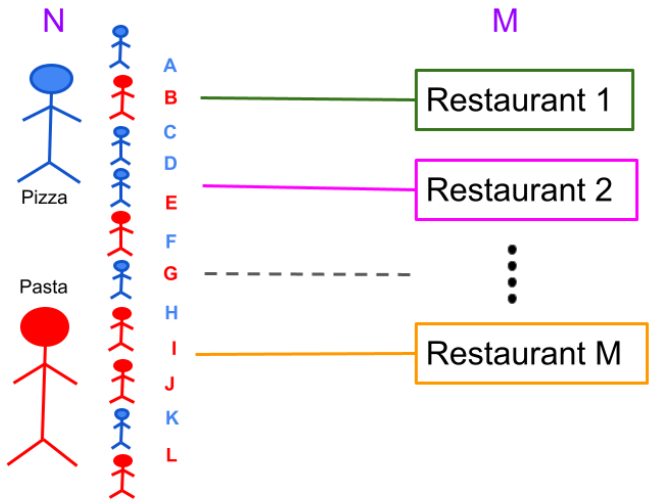
# Motivating example

# Motivating example



Multilayer hypergraph

# Multilayer HSBM

- $N$ Vertices - $\{1, \cdots, N\} =: [N]$. Two communities $\{-1, +1\}$.
- $M$ Layers - $\{1, \cdots, M\}$ indexed by $m$
- $d$ vertices in every hyperedge

# Multilayer HSBM

- $N$ Vertices - $\{1, \cdots, N\} =: [N]$. Two communities $\{-1, +1\}$.
- $M$ Layers - $\{1, \cdots, M\}$ indexed by $m$
- $d$ vertices in every hyperedge

Step 1: Sample the communities
$$\boldsymbol{\sigma} \sim \text{Unif}\left(\left\{\sigma \in \{\pm 1\}^N \mid \text{equal number of } +1 \text{ and } -1\right\}\right)$$

Step 2: For each layer $m \in \{1, \cdots, M\}$ and for each hyperedge $e \subset [N]$ with $|e| = d$, set

$$A_e^{(m)} = \begin{cases} 1 & \text{with prob. } p_e^{(m)} \quad (e \text{ is present in layer } m) \\ 0 & \text{with prob. } 1 - p_e^{(m)} \quad (e \text{ is not present in layer } m) \end{cases},$$

Hypergraph incidence matrix - $\mathbf{A} = (A_e^{(m)})$

# Multilayer HSBM

- $N$ Vertices - $\{1, \cdots, N\} =: [N]$. Two communities $\{-1, +1\}$.
- $M$ Layers - $\{1, \cdots, M\}$ indexed by $m$
- $d$ vertices in every hyperedge

Step 1: Sample the communities

$$\boldsymbol{\sigma} \sim \mathsf{Unif}\left(\left\{\sigma \in \{\pm 1\}^N \;\middle|\; \text{equal number of } +1 \text{ and } -1\right\}\right)$$

Step 2: For each layer $m \in \{1, \cdots, M\}$ and for each hyperedge $e \subset [N]$ with $|e| = d$, set

$$A_e^{(m)} = \begin{cases} 1 & \text{with prob. } p_e^{(m)} & (e \text{ is present in layer } m) \\ 0 & \text{with prob. } 1 - p_e^{(m)} & (e \text{ is not present in layer } m) \end{cases},$$

Hypergraph incidence matrix - $\mathbf{A} = (A_e^{(m)})$

$$(\boldsymbol{\sigma}, \mathbf{A}) \sim \mathsf{HSBM}(N, M, d, (p_e^{(m)}))$$

# Multilayer HSBM: Specifications

$$(\boldsymbol{\sigma}, \mathbf{A}) \sim \text{HSBM}(N, M, d, (p_e^{(m)}))$$

1. **Community profile** of hyperedge $e$ denoted $\boldsymbol{\tau} \equiv (\boldsymbol{\tau}(e))$

$$\boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\end{array}\raisebox{-0.5em}{\includegraphics{spider}}\right) = (3, 2), \quad \boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\\\bullet\\\bullet\end{array}\raisebox{-0.5em}{\includegraphics{spider}}\right) = (2, 5), \quad \boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\end{array}\raisebox{-0.5em}{\includegraphics{spider}}\right) = (3, 0)$$

$$p_e^{(m)} = p_{\tau(e)}^{(m)}$$

For two communities, $\boldsymbol{\tau}(e) \in \{(0, d), (1, d-1), \cdots, (d, 0)\}$.

2. **Scaling regime**: For an edge with community profile $\tau(e)$,

$$p_{\tau(e)}^{(m)} = \alpha_{\tau(e)}^{(m)} \frac{\log N}{\binom{N-1}{d-1}}.$$

3. **Symmetricity**:

$$\alpha_{(r, d-r)}^{(m)} = \alpha_{(d-r, r)}^{(m)}$$

---

Homogeneous: $\alpha_{\boldsymbol{\tau}} = \alpha$ if $\boldsymbol{\tau} \in \{(d, 0), (0, d)\}$; else $\alpha_{\boldsymbol{\tau}} = \beta$.

Gaudio, J. and Joshi, N., 2022. Community detection in the hypergraph sbm: Optimal recovery given the similarity matrix. arXiv preprint arXiv:2208.12227.

# Multilayer HSBM: Specifications

$$(\boldsymbol{\sigma}, \mathbf{A}) \sim \text{HSBM}(N, M, d, (p_e^{(m)}))$$

1. **Community profile** of hyperedge $e$ denoted $\boldsymbol{\tau} \equiv (\boldsymbol{\tau}(e))$

$$\boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\end{array}\!\!\!\rightarrow\!\!\blacksquare\right) = (3, 2), \quad \boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\\\bullet\\\bullet\end{array}\!\!\!\rightarrow\!\!\blacksquare\right) = (2, 5), \quad \boldsymbol{\tau}\left(\begin{array}{c}\bullet\\\bullet\\\bullet\end{array}\!\!\!\rightarrow\!\!\blacksquare\right) = (3, 0)$$

$$p_e^{(m)} = p_{\tau(e)}^{(m)}$$

For two communities, $\boldsymbol{\tau}(e) \in \{(0, d), (1, d-1), \cdots, (d, 0)\}$.

2. **Scaling regime**: For an edge with community profile $\boldsymbol{\tau}(e)$,

$$p_{\tau(e)}^{(m)} = \alpha_{\tau(e)}^{(m)} \frac{\log N}{\binom{N-1}{d-1}}.$$

3. **Symmetricity**:

$$\alpha_{(r, d-r)}^{(m)} = \alpha_{(d-r, r)}^{(m)}$$

$$(\boldsymbol{\sigma}, \mathbf{A}) \sim \text{HSBM}(N, M, d, (\alpha_{\tau}^{(m)}))$$

Homogeneous: $\alpha_{\boldsymbol{\tau}} = \alpha$ if $\boldsymbol{\tau} \in \{(d, 0), (0, d)\}$; else $\alpha_{\boldsymbol{\tau}} = \beta$.

Gaudio, J. and Joshi, N., 2022. Community detection in the hypergraph sbm: Optimal recovery given the similarity matrix. arXiv preprint arXiv:2208.12227.

$\alpha_{(\bullet\bullet\bullet\bullet\bullet)}$ $\qquad$ $\alpha_{(\bullet\bullet\bullet\bullet\bullet)}$ $\qquad$ $\alpha_{(\bullet\bullet\bullet\bullet\bullet)}$

# Assortativity

Assortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

Assortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

$$\parallel \qquad\qquad\qquad \parallel \qquad\qquad\qquad \parallel$$

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

# Assortativity

Assortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

$$\| \qquad\qquad\qquad \| \qquad\qquad\qquad \|$$

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

Disassortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad < \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad < \quad \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$$

# Assortativity

Assortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{(\bullet\bullet\bullet\bullet{\color{blue}\bullet})} \quad > \quad \alpha_{({\color{blue}\bullet}\bullet\bullet\bullet{\color{blue}\bullet})}$$

$$\shortparallel \qquad\qquad \shortparallel \qquad\qquad \shortparallel$$

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{({\color{blue}\bullet}\bullet\bullet\bullet\bullet)} \quad > \quad \alpha_{({\color{blue}\bullet}\bullet\bullet\bullet{\color{blue}\bullet})}$$

Disassortative:

$$\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \quad < \quad \alpha_{(\bullet\bullet\bullet\bullet{\color{blue}\bullet})} \quad < \quad \alpha_{({\color{blue}\bullet}\bullet\bullet\bullet{\color{blue}\bullet})}$$

▶ Each layer could be either assortative or disassortative.
▶ Define assortativity

$$\xi := \sum_{m=1}^{M}\sum_{r=0}^{d-1}\binom{d-1}{r}(d-1-2r)\alpha_{(r,d-r)}^{(m)}.$$

▶ Assortative: $\xi > 0$ and Disassortative: $\xi < 0$.
▶ For $d = 5$, $\quad \xi = 4\alpha_{(0,5)} + 4\alpha_{(1,4)} - 8\alpha_{(2,3)}$.

# Motivation

# Motivation



N

Pizza

Pasta

A B C D E F G H I J K L

Restaurant

N

Pizza

Pasta

A
B
C
D
E
F
G
H
I
J
K
L

Restaurant

Recover communities while maintaining privacy.

# Similarity matrix

# Similarity matrix



|   | A | B | C | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 0 | 1 | 2 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 1 |
| B |   | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C |   |   | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| D |   |   |   | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| E |   |   |   |   | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| F |   |   |   |   |   | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| G |   |   |   |   |   |   | 0 | 0 | 1 | 0 | 0 | 1 |
| H |   |   |   |   |   |   |   | 0 | 0 | 0 | 0 | 0 |
| I |   |   |   |   |   |   |   |   | 0 | 0 | 0 | 1 |
| J |   |   |   |   |   |   |   |   |   | 0 | 0 | 0 |
| K |   |   |   |   |   |   |   |   |   |   | 0 | 0 |
| L |   |   |   |   |   |   |   |   |   |   |   | 0 |

## Problem formulation

Data: $(\sigma, \mathbf{A}) \sim \mathsf{HSBM}(N, M, d, (\alpha_\tau^{(m)}))$

Given: $N \times N$ aggregate similarity matrix $\mathbf{W} = (\mathbf{W}_{ij})$, such that

$$\mathbf{W}_{ij} = \sum_{m=1}^{M} W_{ij}^{(m)},$$

where

$$
\begin{aligned}
W_{ij}^{(m)} &= \text{\# of hyperedges that contain both } i \text{ and } j \text{ in layer } m \\
&= \sum_{e:e \ni i,j} A_e^{(m)}
\end{aligned}
$$

## Problem formulation

Data: $(\boldsymbol{\sigma}, \mathbf{A}) \sim \mathsf{HSBM}(N, M, d, (\alpha_\tau^{(m)}))$

Given: $N \times N$ aggregate similarity matrix $\mathbf{W} = (\mathbf{W}_{ij})$, such that

$$\mathbf{W}_{ij} = \sum_{m=1}^{M} W_{ij}^{(m)},$$

where

$$
\begin{aligned}
W_{ij}^{(m)} &= \# \text{ of hyperedges that contain both } i \text{ and } j \text{ in layer } m \\
&= \sum_{e: e \ni i, j} A_e^{(m)}
\end{aligned}
$$

Want to find an estimate $\hat{\boldsymbol{\sigma}}^{(N)}$ of $\boldsymbol{\sigma}$ ($\equiv \boldsymbol{\sigma}^{(N)}$) that exactly recovers the communities,

$$\lim_{N \to \infty} \mathbb{P}\left(\hat{\boldsymbol{\sigma}}^{(N)} \in \{\pm \boldsymbol{\sigma}^{(N)}\}\right) = 1.$$

# A first approach

In the assortative case:

▶ Solve the min-bisection problem:

$$\max \sum_{i,j} W_{ij} x_i x_j \qquad \text{subject to } \mathbf{x} \in \{\pm 1\}^N, \mathbf{1}^T \mathbf{x} = 0. \qquad (1)$$

[1] Kim, C., Bandeira, A.S. and Goemans, M.X., 2018. Stochastic block model for hypergraphs: Statistical limits and a semidefinite programming approach. arXiv preprint arXiv:1807.02884.

# A first approach

In the assortative case:

▶ Solve the min-bisection problem:

$$\max \sum_{i,j} W_{ij} x_i x_j \quad \text{subject to } \mathbf{x} \in \{\pm 1\}^N, \mathbf{1}^T \mathbf{x} = 0. \quad (1)$$

▶ SDP relaxation[1]:

$$
\begin{aligned}
\text{maximize} \quad & \sum_{1 \le i < j \le N} W_{ij} X_{ij} \\
\text{subject to} \quad & \sum_{1 \le i < j \le N} X_{ij} = 0, \\
& X_{ii} = 1 \text{ for all } i \in [N] \\
& \mathbf{X} \succeq 0.
\end{aligned}
\quad (2)
$$

▶ Any solution $\mathbf{x}$ of (1) is a solution of (2) by taking $\mathbf{X} = \mathbf{xx}^T$.

[1] Kim, C., Bandeira, A.S. and Goemans, M.X., 2018. Stochastic block model for hypergraphs: Statistical limits and a semidefinite programming approach. arXiv preprint arXiv:1807.02884.

# Algorithm and main result

### Step 1

Given $s \in \{\pm 1\}$ and $\mathbf{W}$, solve:

$$\text{maximize} \sum_{1 \le i < j \le N} s W_{ij} X_{ij}$$

$$\text{subject to} \sum_{1 \le i < j \le N} X_{ij} = 0,$$

$$X_{ii} = 1 \text{ for all } i \in [N]$$

$$\mathbf{X} \succeq 0.$$

### Step 2

The optimal solution $\mathbf{X}^* = \sum_{i=1}^{N} \lambda_i \mathbf{v_i} \mathbf{v_i}^{\mathsf{T}}$ with $\lambda_1 \ge \cdots \ge \lambda_N$.

### Step 3

Output $\hat{\boldsymbol{\sigma}} = \text{sgn}(\mathbf{v_1})$

# Algorithm and main result

**Step 1**
Given $s \in \{\pm 1\}$ and $\mathbf{W}$, solve:

$$\text{maximize} \sum_{1 \leq i < j \leq N} s W_{ij} X_{ij}$$

$$\text{subject to} \sum_{1 \leq i < j \leq N} X_{ij} = 0,$$

$$X_{ii} = 1 \text{ for all } i \in [N]$$

$$\mathbf{X} \succeq 0.$$

**Step 2**
The optimal solution $\mathbf{X}^* = \sum_{i=1}^{N} \lambda_i \mathbf{v_i} \mathbf{v_i}^{\mathsf{T}}$ with $\lambda_1 \geq \cdots \geq \lambda_N$.

**Step 3**
Output $\hat{\boldsymbol{\sigma}} = \mathrm{sgn}(\mathbf{v_1})$

## Theorem

*Suppose $(\boldsymbol{\sigma}, \mathbf{A}) \sim \mathrm{HSBM}(N, M, d, (\alpha_\tau^{(m)}))$, and let $\mathbf{W}$ be the aggregate similarity matrix of $\mathbf{A}$. When $I > 1$, the above algorithm with $\mathbf{W}$ and $s = \mathrm{sgn}(\xi)$ as inputs, exactly recovers $\sigma$.*

$$I = \sup_{\lambda \in \mathbb{R}} \sum_{m=1}^{M} \sum_{r=0}^{d-1} 2^{-(d-1)} \binom{d-1}{r} \alpha_{(r,d-r)}^{(m)} \left(1 - e^{-\lambda(d-1-2r)}\right)$$

## Dual formulation

### Primal problem:

$$\begin{aligned}
\max \quad & \sum_{1 \le i < j \le N} W_{ij} X_{ij} \\
\text{subject to} \quad & X_{ii} = 1, \ \forall i \in [N] \\
& \langle \mathbf{X}, \mathbf{1}\mathbf{1}^T \rangle = 0, \\
& \mathbf{X} \succeq 0.
\end{aligned}
\quad \equiv \quad
\begin{aligned}
\min \quad & \langle \mathbf{W}', \mathbf{X} \rangle \\
\text{subject to} \quad & \langle \mathbf{A}_i, \mathbf{X} \rangle = 1, \ \forall i \in [N] \\
& \langle \mathbf{X}, \mathbf{J} \rangle = 0, \\
& - \mathbf{X} \preceq 0.
\end{aligned}$$

where

$$(\mathbf{A}_i)_{jk} = 0 \text{ for } j \ne k, \text{ and } (\mathbf{A}_i)_{jj} = \begin{cases} 1 & i = j \\ 0 & i \ne j \end{cases}$$

### Lagrangian:

$$\mathcal{L}(\mathbf{X}, \mathbf{S}, \nu, \boldsymbol{d}) = \langle \mathbf{W}', \mathbf{X} \rangle - \langle \mathbf{S}, \mathbf{X} \rangle + \nu \langle \mathbf{X}, \mathbf{J} \rangle + \sum_{i=1}^{N} d_i \left( \langle \mathbf{A}_i, \mathbf{X} \rangle - 1 \right).$$

where $\mathbf{S} \succeq 0$.

## Dual formulation

Dual objective:

$$
\begin{aligned}
g(\mathbf{S}, \nu, \boldsymbol{d}) &= \inf_{\mathbf{X}} \ \mathcal{L}(\mathbf{X}, \mathbf{S}, \nu, \boldsymbol{d}) \\
&= \inf_{\mathbf{X}} \ \langle \mathbf{W}' - \mathbf{S} + \nu \mathbf{J} + \text{ diag } (\boldsymbol{d}), \mathbf{X} \rangle - \sum_{i=1}^{N} d_i \\
&= \inf_{\mathbf{X}} \ \langle \mathbf{W}' - \mathbf{S} + \nu \mathbf{J} + \mathbf{D}, \mathbf{X} \rangle - \text{ trace } (\mathbf{D}) \\
&= \begin{cases} - \text{ trace } (\mathbf{D}) & \text{if } \mathbf{W}' - \mathbf{S} + \nu \mathbf{J} + \mathbf{D} = 0 \\ -\infty & \text{o.w.} \end{cases}
\end{aligned}
$$

Dual problem:

$$
\begin{array}{ll}
\max & - \text{ trace } (\mathbf{D}) \\
\text{subject to} & \mathbf{W}' - \mathbf{S} + \nu \mathbf{J} + \mathbf{D} = 0, \\
& \mathbf{S} \succeq 0
\end{array}
\quad \equiv \quad
\begin{array}{ll}
\min & \text{trace } (\mathbf{D}) \\
\text{subject to} & \mathbf{D} + \nu \mathbf{J} - \mathbf{W} \succeq 0.
\end{array}
$$

## Dual certificate

**W**: Observed aggregate similarity matrix

### Lemma

*Suppose there is a $N \times N$ diagonal matrix **D** such that
$\mathbf{S} := \mathbf{D} + \mathbf{1}\mathbf{1}^T - \mathbf{W}$ satisfies:*

$$\mathbf{S} \succeq 0, \quad \lambda_{N-1}(\mathbf{S}) > 0, \quad and \quad \mathbf{S}\boldsymbol{\sigma} = 0,$$

*then $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$ is the unique optimal solution to the SDP.*

# Dual certificate

**W**: Observed aggregate similarity matrix

> ### Lemma
>
> *Suppose there is a $N \times N$ diagonal matrix* **D** *such that*
> $\mathbf{S} := \mathbf{D} + \mathbf{1}\mathbf{1}^T - \mathbf{W}$ *satisfies:*
>
> $$\mathbf{S} \succeq 0, \quad \lambda_{N-1}(\mathbf{S}) > 0, \quad and \quad \mathbf{S}\boldsymbol{\sigma} = 0,$$
>
> *then* $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$ *is the unique optimal solution to the SDP.*

Weak duality: **Optimality**
Let $X$ be primal feasible and $X^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$. Then

$$
\begin{aligned}
\langle \mathbf{W}, \mathbf{X} \rangle &\leq \ \text{trace}\,(\mathbf{D}) \\
&= \langle \mathbf{D}, \mathbf{X} \rangle \ = \ \langle \mathbf{D}, \mathbf{X}^* \rangle && \text{(since } X_{ii} = X_{ii}^* = 1) \\
&= \langle \mathbf{W} + \mathbf{S} - \nu \mathbf{J}, \mathbf{X}^* \rangle = \langle \mathbf{W}, \mathbf{X}^* \rangle && \text{(since } \langle \mathbf{S}, \mathbf{X}^* \rangle = \boldsymbol{\sigma}^T(S\boldsymbol{\sigma}) = 0)
\end{aligned}
$$

# Dual certificate

**W**: Observed aggregate similarity matrix

> ### Lemma
>
> *Suppose there is a $N \times N$ diagonal matrix **D** such that*
> $\mathbf{S} := \mathbf{D} + \mathbf{1}\mathbf{1}^T - \mathbf{W}$ *satisfies:*
>
> $$\mathbf{S} \succeq 0, \quad \lambda_{N-1}(\mathbf{S}) > 0, \quad and \quad \mathbf{S}\boldsymbol{\sigma} = 0,$$
>
> *then $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$ is the unique optimal solution to the SDP.*

**Strong duality: Uniqueness**
Let $\tilde{X}$ be an optimal solution and $X^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$. Then

$$
\begin{aligned}
\langle \mathbf{S}, \tilde{\mathbf{X}} \rangle &= \langle \mathbf{D} + \mathbf{J} - \mathbf{W}, \tilde{\mathbf{X}} \rangle = \langle \mathbf{D} - \mathbf{W}, \tilde{\mathbf{X}} \rangle \\
&= \langle \mathbf{D} - \mathbf{W}, \mathbf{X}^* \rangle \qquad (\langle \mathbf{W}, \tilde{\mathbf{X}} \rangle = \langle \mathbf{W}, \mathbf{X}^* \rangle \text{ and } \tilde{X}_{ii} = X_{ii}^* = 1) \\
&= \langle \mathbf{S}, \mathbf{X}^* \rangle = 0 \qquad (\text{since } \langle \mathbf{S}, \mathbf{X}^* \rangle = \boldsymbol{\sigma}^T(S\boldsymbol{\sigma}) = 0)
\end{aligned}
$$

Since $\mathbf{S} \succeq 0$ and $\lambda_{N-1} > 0$, the Null($\mathbf{S}$) is spanned by $\boldsymbol{\sigma}$ only.
$\tilde{\mathbf{X}} \succeq 0$ now implies that it should be a multiple of $\boldsymbol{\sigma}\boldsymbol{\sigma}^T$ as well.

# Proof: Dual certificate

**W**: Observed aggregate similarity matrix

### Lemma

*Suppose there is a $N \times N$ diagonal matrix $\mathbf{D}$ such that*
$\mathbf{S} := \mathbf{D} + \mathbf{1}\mathbf{1}^T - \mathbf{W}$ *satisfies:*

$$\mathbf{S} \succeq 0, \quad \lambda_{N-1}(\mathbf{S}) > 0, \quad \text{and} \quad \mathbf{S}\boldsymbol{\sigma} = 0,$$

*then $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$ is the unique optimal solution to the SDP.*

# Proof: Dual certificate

**W**: Observed aggregate similarity matrix

---

### Lemma

*Suppose there is a $N \times N$ diagonal matrix $\mathbf{D}$ such that $\mathbf{S} := \mathbf{D} + \mathbf{1}\mathbf{1}^T - \mathbf{W}$ satisfies:*

$$\mathbf{S} \succeq 0, \quad \lambda_{N-1}(\mathbf{S}) > 0, \quad \text{and} \quad \mathbf{S}\boldsymbol{\sigma} = 0,$$

*then $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^T$ is the unique optimal solution to the SDP.*

---

Taking

$$D_{ii} := \sum_j W_{ij}\sigma_i\sigma_j,$$

easy to verify $\mathbf{S}\boldsymbol{\sigma} = 0$. Suffices to show

$$\mathbb{P}\left(\inf_{\mathbf{x} \perp \boldsymbol{\sigma}: \|\mathbf{x}\|_2 = 1} \mathbf{x}^T \mathbf{S}\mathbf{x} > 0\right) = 1 - o(1).$$

### Lemma

Let $\mathbb{E}\mathbf{W}$ is the expected aggregate similarity matrix. Then

$$\mathbf{x}^T \mathbf{S} \mathbf{x} \ \geq \ \min_i D_{ii} - \|\mathbf{W} - \mathbb{E}\mathbf{W}\|_2.$$

# Proof: Bounds

## Lemma

Let $\mathbb{E}\mathbf{W}$ is the expected aggregate similarity matrix. Then

$$\mathbf{x}^T\mathbf{S}\mathbf{x} \ \geq \ \min_i D_{ii} - \|\mathbf{W} - \mathbb{E}\mathbf{W}\|_2.$$

## Proposition

Let $I > 1$. Then there exists a constant $\epsilon > 0$ dependent on model parameters such that for all $i \in [N]$,

$$\mathbb{P}(D_{ii} > \epsilon \log N) \geq 1 - o(N^{-1}).$$

## Proposition

There exists a constant $C$ such that

$$\mathbb{P}\left(\|\mathbf{W} - \mathbb{E}\mathbf{W}\|_2 \leq CM\sqrt{\log N}\right) \geq 1 - O(N^{-11}).$$

# Proof: Main ingredients

▶ Rank-2 decomposition:

$$\mathbb{E}\mathbf{W} = \left(\frac{w_{\text{in}} + w_{\text{out}}}{2}\right)\mathbf{1}\mathbf{1}^T + \left(\frac{w_{\text{in}} - w_{\text{out}}}{2}\right)\boldsymbol{\sigma}\boldsymbol{\sigma}^T - w_{\text{in}}\mathbf{I}_N,$$

where $w_{\text{in}} = \mathbb{E}[W_{ij}|\sigma_i = \sigma_j]$ and $w_{\text{out}} = \mathbb{E}[W_{ij}|\sigma_i \neq \sigma_j]$.

▶ Assortativity:

$$w_{\text{in}} - w_{\text{out}} \approx \frac{\log N}{2^{d-2}N}\xi.$$

# Summary

- Multilayer HSBM $(\boldsymbol{\sigma}, \mathbf{A}) \sim \text{HSBM}(N, M, d, (\alpha_\tau^{(m)}))$

- Inhomogeneous hyperedge probabilities

- Assortative and disassortative cases

- Exact recovery using the similarity matrix $\mathbf{W}$

  SDP algorithm recovers the clusters exactly when $I > 1$.

# Summary

- Multilayer HSBM $(\boldsymbol{\sigma}, \mathbf{A}) \sim \text{HSBM}(N, M, d, (\alpha_\tau^{(m)}))$

- Inhomogeneous hyperedge probabilities

- Assortative and disassortative cases

- Exact recovery using the similarity matrix $\mathbf{W}$

  SDP algorithm recovers the clusters exactly when $I > 1$.

## Future work

- Asymmetric: $\alpha_{(\bullet\bullet\bullet\bullet\bullet)} \neq \alpha_{(\bullet\bullet\bullet\bullet\bullet)}$

- Necessary conditions for exact recovery from $\mathbf{W}$

- For different hyperedge sizes $d$

- $M$, $d$ depending on $N$

Thank you !!