# In and Out: Optimizing Overall Interaction in Probabilistic Graphs under Clustering Constraints

Domenico Mandaglio, Andrea Tagarelli, Francesco Gullo*

DIMES – Univ. Calabria
Rende (CS), Italy

DIMES – Univ. Calabria
Rende (CS), Italy

UniCredit
Rome, Italy

# Outline

- General setting

- Related Work

- Problem formulations:
  - Maximize overall interaction
  - Minimize overall interaction loss

- Algorithms

- Experimental results

- Conclusions & Future Work

**General context**: maximize interactions (user engagement) in social network system

Focus on two properties of a social network system:

1. Uncertainty in user behaviors

2. Exogenous conditions can affect the users' interaction behaviors

Our design choice: ***clustering constraint assumption***

- ***the (uncertain) interaction behaviors depend on a clustering of the set of users in a graph.***
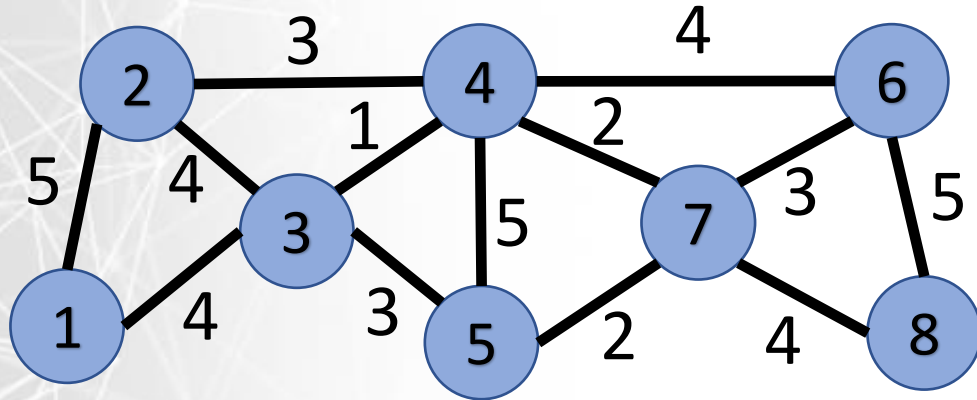
# Applications (I)

**Online social platforms**

- User personal home pages show contents produced by other users

- A cluster of users corresponds to a set of users which are interested in contents generated by users belonging to the cluster

- **Goal**: exploit a clustering of the user to drive the delivering of contents to homepages such as to maximize the overall interaction between users

# Applications (II)

**Team formation**

- Users should be grouped into teams to contribute to a common global task

- The likelihood of collaboration between any pair of users will vary in relation to their assignment to the same team

- **Goal**: partition the set of users into teams in order to maximize the total collaboration
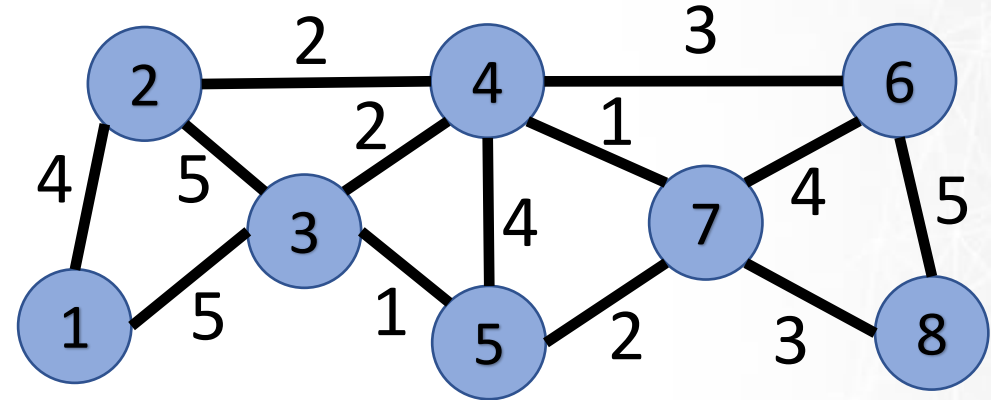
# Overall Interaction
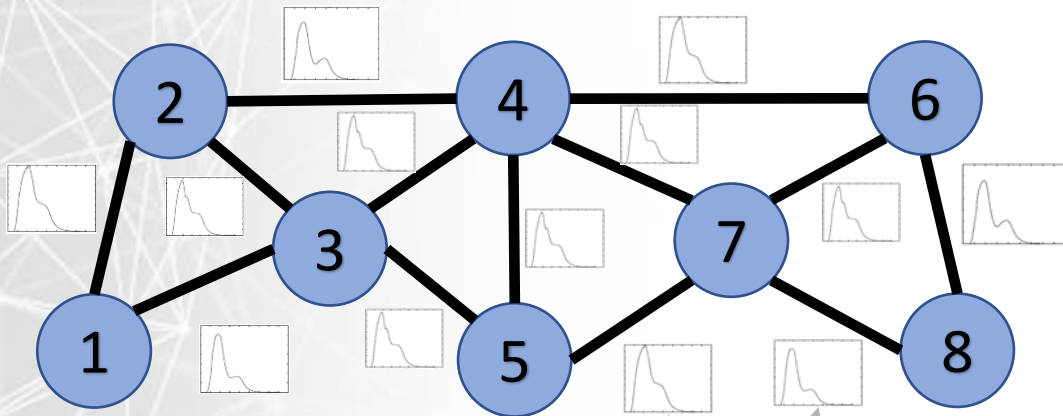


$G_t = (V, E, w_t)$

$G_{t+1} = (V, E, w_{t+1})$

**Overall interaction**

$$f(G_t) = \sum_{(u,v) \in E} w_t(u,v) = 45$$

$$f(G_{t+1}) = \sum_{(u,v) \in E} w_{t+1}(u,v) = 42$$

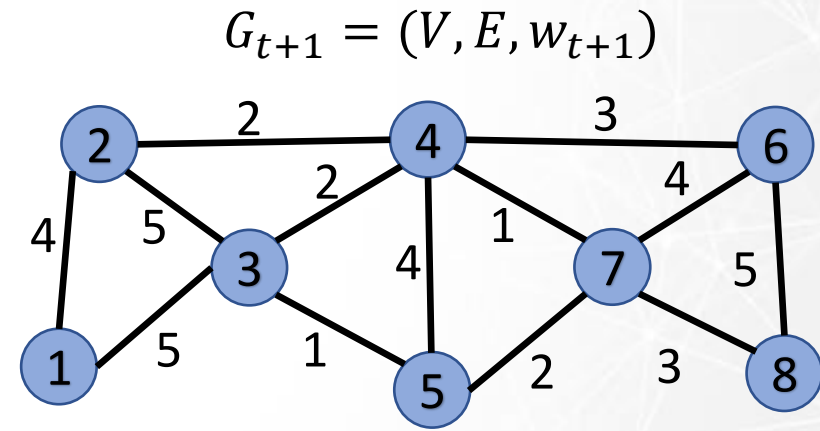# Probabilistic Interaction Graph



...

$G_t = (V, E, w_t)$

$G_t, G_{t+1} \sqsubseteq \mathcal{G}$

Possible worlds

$\mathcal{G} = (V, E, P)$    $P = \{p_{uv}\}_{(u,v) \in E}$

**Probabilistic interaction graph**

$G_{t+1} = (V, E, w_{t+1})$

...

# Clustering-Conditional Probabilistic Graph



$$\mathcal{G}^+ = (V, E, P^+)$$

$$\mathcal{G}^- = (V, E, P^-)$$

$\mathcal{C}$ Clustering of $V$

$$P_{\mathcal{C}} = \{p_{uv} \in P^+ | \mathcal{C}(u) = \mathcal{C}(v)\} \cup \{p_{uv} \in P^- | \mathcal{C}(u) \neq \mathcal{C}(v)\}$$

$\mathcal{G}_{\mathcal{C}} = (V, E, P_{\mathcal{C}})$ **Clustering-conditional probabilistic graph**

# Problem formulations

Given two interaction graphs $\mathcal{G}^+ = (V, E, P^+), \mathcal{G}^- = (V, E, P^-)$ find a clustering $\mathcal{C}^*: V \rightarrow \mathbb{N}$ to:

1) **Maximize (expected) overall interaction**

2) **Minimize (expected) overall interaction loss**

# Related Work

- Clustering uncertain graphs:
  - Interactions are binary
  - Maximize the intra-cluster connectivity and minimize the inter-cluster connectivity
  - Exogenous conditioning factors are not considered

- Community detection in signed graphs:
  - Edges with a sign and a weight
  - Maximize positive (resp. negative) links within (resp. between) communities

- Correlaton clustering:
  - Advice on whether two nodes should be clustered together or not

G. Kollios, M. Potamias, and E. Terzi. 2013. Clustering Large Probabilistic Graphs. IEEE TKDE 25,2(2013),325–336.

M. Ceccarello, C. Fantozzi, A. Pietracaprina, G. Pucci, and F. Vandin. 2017. Clustering Uncertain Graphs. PVLDB 11,4(2017),472–484.

Yu G., Chunpeng G., Gao C., and Ge Y. 2014. Effective and Efficient Clustering Methods for Correlated Probabilistic Graphs. IEEE TKDE 26, 5(2014),1117–1130.

A. Khan, F. Bonchi, F. Gullo, and A. Nufer. 2018. Conditional Reliability in Uncertain Graphs. IEEE TKDE 30,11(2018),2078–2092.

V. A. Traag and J. Bruggeman. 2009. Community detection in networks with positive and negative links. Physical Review E 80,3(2009),036115.

S. Gómez, P. Jensen, and A. Arenas. 2009. Analysis of community structure in networks of correlated data. Physical Review E 80,1(2009),016114.

P. Esmailian and M. Jalili. 2015. Community detection in signed networks: the role of negative ties in different scales. Scientific reports 5(2015),14339.

N. Bansal, A. Blum, and S. Chawla. 2004. Correlation Clustering. Machine Learning 56,1(2004),89–113.

# Background: (Weighted) Correlation Clustering

Given a set $\Omega$ of entities, and weights $\omega_{xy}^+, \omega_{xy}^- \in \mathbb{R}_0^+$ for all unordered pairs $x, y \in \Omega$ find a clustering $\mathcal{C}: \Omega \to \mathbb{N}$ that:

- **Maximize Agreements (Max-CC)**

$$\sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)=\mathcal{C}(y)}} \omega_{xy}^+ + \sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)\neq\mathcal{C}(y)}} \omega_{xy}^-$$

- **Minimize Disagreements (Min-CC)**

$$\sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)=\mathcal{C}(y)}} \omega_{xy}^- + \sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)\neq\mathcal{C}(y)}} \omega_{xy}^+$$

# Maximizing Interaction

**MAX-INTERACTION-CLUSTERING.** Given two interaction graphs $\mathcal{G}^+ = (V, E, P^+)$, $\mathcal{G}^- = (V, E, P^-)$ find a clustering $\mathcal{C}^* : V \to \mathbb{N}$ such that:

Overall interaction

$$\mathcal{C}^* = \underset{\mathcal{C}}{\operatorname{argmax}} \; \overline{f}(\mathcal{G}_\mathcal{C}) = \underset{\mathcal{C}}{\operatorname{argmax}} \; \mathbb{E}_{G \sqsubseteq \mathcal{G}_\mathcal{C}} \left[ f(G) \right]$$

$$\overline{f}(\mathcal{G}_\mathcal{C}) = \sum_{\substack{u,v \in V \\ \mathcal{C}(u) = \mathcal{C}(v)}} \mathbb{E}[p_{uv}^+] + \sum_{\substack{u,v \in V \\ \mathcal{C}(u) \neq \mathcal{C}(v)}} \mathbb{E}[p_{uv}^-]$$

**Max-CC instance**

$$\sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x) = \mathcal{C}(y)}} \omega_{xy}^+ + \sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x) \neq \mathcal{C}(y)}} \omega_{xy}^-$$

$$\Omega = V, \omega_{xy}^+ = \mathbb{E}[p_{uv}^+], \omega_{xy}^- = \mathbb{E}[p_{uv}^-]$$

# Maximizing Interaction

- MAX-INTERACTION-CLUSTERING is **NP**-Hard

- Approximation algorithms designed for Max-CC keep their guarantees on MAX-INTERACTION-CLUSTERING too

- State-of-the-art approximation algorithms for Max-CC (on general, weighted graphs) are inefficient and impractical (output at most a small, fixed number of clusters)

# Minimizing Interaction Loss

**MIN-INTERACTION-LOSS-CLUSTERING**. Given two interaction graphs $\mathcal{G}^+ = (V, E, P^+)$, $\mathcal{G}^- = (V, E, P^-)$ find a clustering $\mathcal{C}^*: V \rightarrow \mathbb{N}$ such that:

Overall interaction loss

$$\mathcal{C}^* = \underset{\mathcal{C}}{\operatorname{argmin}} \, \overline{\ell}(\mathcal{G}_\mathcal{C}) = \underset{\mathcal{C}}{\operatorname{argmin}} \, \mathbb{E}_{G \sqsubseteq \mathcal{G}_\mathcal{C}}[\ell(G)]$$

$$\overline{\ell}(\mathcal{G}_\mathcal{C}) = \sum_{\substack{u,v \in V \\ \mathcal{C}(u)=\mathcal{C}(v)}} M(\mathcal{G}^+, \mathcal{G}^-) - \mathbb{E}[p_{uv}^+] + \sum_{\substack{u,v \in V \\ \mathcal{C}(u)\neq\mathcal{C}(v)}} M(\mathcal{G}^+, \mathcal{G}^-) - \mathbb{E}[p_{uv}^-]$$

### Min-CC instance

$$\sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)=\mathcal{C}(y)}} \omega_{xy}^- + \sum_{\substack{x,y \in \Omega \\ \mathcal{C}(x)\neq\mathcal{C}(y)}} \omega_{xy}^+$$

$$\Omega = V$$

$$\omega_{xy}^+ = M(\mathcal{G}^+, \mathcal{G}^-) - \mathbb{E}[p_{uv}^-]$$
$$\omega_{xy}^- = M(\mathcal{G}^+, \mathcal{G}^-) - \mathbb{E}[p_{uv}^+]$$

# Minimizing Interaction Loss

- MIN-INTERACTION-LOSS-CLUSTERING is **NP**-Hard and it is equivalent to its maximization formulation counterpart

- Approximation algorithms designed for Min-CC keep their guarantees (under certain conditions) on MIN-INTERACTION-LOSS-CLUSTERING too

- More practical and efficient algorithms available for Min-CC

# Pivot algorithm

- Pick a node $u$ uniformly at random
- Build a cluster upon $u$ together with its similar nodes that are still unclustered
- Remove the built cluster from the graph
- Repeat until the graph is empty

**Properties of Pivot**:
- (expected) 5-approximation guarantee
- Can be improved to 2-approximation guarantee provided that weights satisfy the triangle inequality property

# Theoretical basis

It does not hold for our instances 😫

**Probability Constraint:** $\omega_{xy}^+ + \omega_{xy}^- = 1$ for all $x, y \in \Omega$ ⬅

**Question**: is solving our problem equivalent to solve Min-CC instance where the probability constraint is satisfied?

$$\overline{\ell}(\mathcal{G}_C) = g(\mathcal{G}_C) \times M(\mathcal{G}^+, \mathcal{G}^-) + K(\mathcal{G}^+, \mathcal{G}^-)$$

🙂

Min-CC objective over an instance where probability constraint holds

Constant depending on the input graphs

# Theoretical basis

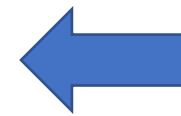**Theorem**. If $K(\mathcal{G}^+, \mathcal{G}^-) \geq 0,$ then Pivot is a randomized expected 5-approximation algorithm for MIN-INTERACTION-LOSS-CLUSTERING

**Condition for approximation guarantees**

$$K(\mathcal{G}^+, \mathcal{G}^-) \geq 0$$

$$\Updownarrow$$

$$\sum_{(u,v)\in E} \mathbb{E}[p_{uv}^+] + \mathbb{E}[p_{uv}^-] \leq M(\mathcal{G}^+, \mathcal{G}^-) \times \binom{|V|}{2}$$

It holds for sparse graphs

# Pivot for minimizing interaction loss (*MIL*)

**Algorithm 1** MIL

**Input:** Interaction graphs $\mathcal{G}^+ = (V, E, P^+)$, $\mathcal{G}^- = (V, E, P^-)$

**Output:** A clustering $C$ of $V$

1: compute $\tau_{uv}^+$, $\tau_{uv}^-$ for all $(u, v) \in E$     ⬅ Build the Min-CC instance where the probability constraint holds

2: $C \leftarrow \emptyset$, $V' \leftarrow V$

3: **while** $V' \neq \emptyset$ **do**

4:     pick a pivot vertex $u \in V'$ uniformly at random

5:     $C_u \leftarrow \{u\} \cup \{v \in V' \mid (u, v) \in E, \tau_{uv}^+ > \tau_{uv}^-\}$

6:     add cluster $C_u$ to $C$ and remove all vertices in $C_u$ from $V'$

Pivot algorithm

MIL runs in $O(|V| + |E|)$ time

# MIL algorithm: effect of sampling pivots uniformly at random in general graphs



$$sign((u,v)) = \begin{cases} + & if \ \ \mathbb{E}[p_{uv}^+] > \mathbb{E}[p_{uv}^-] \\ - & otherwise \end{cases}$$

$$\mathcal{C}^* = \{\{1, 5, 8, 9\}, \{3, 7, 10, 11\}, \{2, 4, 6, 12\}\}$$

No matter to put non-linked nodes in the same cluster or in different clusters

**Idea**: sample pivots by degree!

# D-MIL algorithm

**Algorithm 2** D-MIL

**Input:** Interaction graphs $\mathcal{G}^+ = (V, E, P^+)$, $\mathcal{G}^- = (V, E, P^-)$
**Output:** A clustering $C$ of $V$
1: compute $\tau_{uv}^+$, $\tau_{uv}^-$ for all $(u, v) \in E$
2: $C \leftarrow \emptyset$, $V' \leftarrow V$
3: **while** $V' \neq \emptyset$ **do**
4:    compute $d_{V'}(u) = |\{v \in V' \mid (u, v) \in E\}|$, for all $u \in V'$
5:    sample a pivot vertex $u \in V'$ with probability proportional to $d_{V'}(u)$
6:    $C_u \leftarrow \{u\} \cup \{v \in V' \mid (u, v) \in E, \tau_{uv}^+ > \tau_{uv}^-\}$
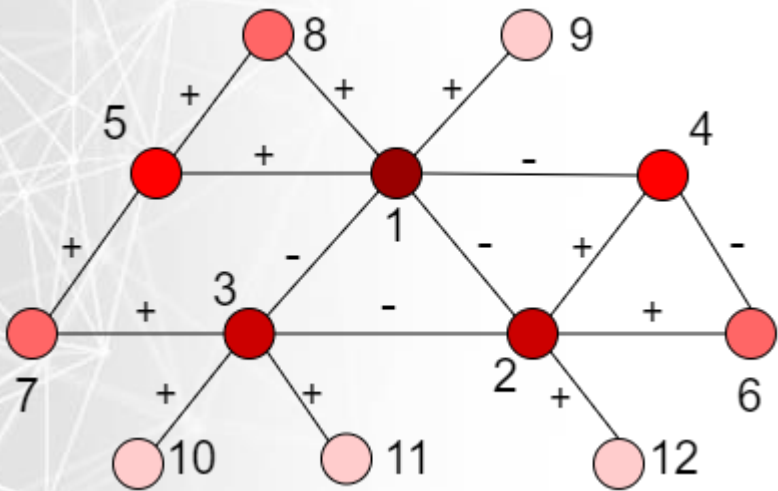7:    add cluster $C_u$ to $C$ and remove all vertices in $C_u$ from $V'$

Build the Min-CC instance where the probability constraint holds

Pivot algorithm with sampling by degree (no approximation guarantee)

D-MIL runs in $O(|E| \log |V|)$ time

# Hill Climbing

**Algorithm 3** HillClimbing

**Input:** Interaction graphs $\mathcal{G}^+ = (V, E, P^+)$, $\mathcal{G}^- = (V, E, P^-)$;  A clustering $C$ of $V$;  An integer $I > 0$
**Output:** A clustering $C'$ of $V$
1: $C' \leftarrow C$
2: **for all** $i = 1, \ldots, I$ **do**
3:     for every $u \in V$ let $C_u \in C'$ the cluster of $C'$ where $u$ belongs to
4:     pick $u \in V$ and cluster $C'_u \in C'$ ($C'_u \neq C_u$) that minimize Eq. (16)
5:     $C'' \leftarrow$ clustering obtained from $C'$ by moving $u$ from $C_u$ to $C'_u$
6:     **if** $\bar{\ell}(\mathcal{G}_{C''}) < \bar{\ell}(\mathcal{G}_{C'})$ **then**
7:         $C' \leftarrow C''$

Hill Climbing runs in
$O(I \times (|V| + |E|))$ time

- MIL + Hill Climbing = MIL_R          →     - Approximation guarantees of MIL

- D-MIL + Hill Climbing = D-MIL_R   →     - No approximation guarantees

# Evaluation

**Data**

- <u>Real</u> network data
- <u>Syntethic</u> network data: Barabasi-Albert (BA) and Watts-Strogatz (WS) random graph models

**Evaluation goals**

- Interaction Loss
- Clustering size
- Efficiency evaluation
- Comparison with competing methods (CPM[1], CJA[2], CPMap[3])

1. V. A. Traag and J. Bruggeman. 2009. Community detection in networks with positive and negative links. Physical Review E 80,3(2009),036115.
2. S. Gómez, P. Jensen, and A. Arenas. 2009. Analysis of community structure in networks of correlated data. Physical Review E 80,1(2009),016114.
3. P. Esmailian and M. Jalili. 2015. Community detection in signed networks: the role of negative ties in different scales. Scientific reports 5(2015),14339.

# Real network data

Preprocessing of timestamped networks:
- The edge set $E$ of the probabilistic graph are obtained by "flattening" the temporal network
- The distributions $p_{uv}^+$ and $p_{uv}^-$ are estimated based on the fraction of clusters shared by $u, v$ over all graphs

**Table 1: Summary of real networks**

|  | $|V|$ | $\sum_{t=1}^{T} |E_t|$ | $T$ | edge semantics | $|E|$ |
|---|---|---|---|---|---|
| *Amazon* | 2 146 057 | 22 728 036 | 115 | co-rating | 22 507 680 |
| *DBLP* | 1 824 701 | 11 865 584 | 80 | co-authorship | 8 344 615 |
| *Epinions* | 120 492 | 33 412 111 | 25 | co-rating | 24 994 363 |
| *HighSchool* | 327 | 47 589 | 1212 | face-to-face | 5 818 |
| *Last.fm* | 992 | 4 342 951 | 77 | co-listening | 369 973 |
| *PrimarySchool* | 242 | 55 043 | 390 | face-to-face | 8 317 |
| *ProsperLoans* | 89 269 | 3 343 271 | 307 | economic | 3 330 022 |
| *StackOverflow* | 2 433 067 | 16 200 209 | 51 | Q/A | 15 786 816 |
| *Wikipedia* | 343 860 | 18 086 734 | 101 | co-editing | 10 519 921 |
| *WikiTalk* | 2 863 439 | 10 335 318 | 192 | communication | 8 146 544 |

# Avg. loss values and clustering sizes on real data

| | MIL | | MIL_R | | D-MIL | | D-MIL_R | | CPM [21] | | GJA [9] | | CPMap [7] | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | loss | #clusters | loss | #clusters | loss | #clusters | loss | #clusters | loss | #clusters | loss | #clusters | loss | #clusters |
| Amazon | $4.80 \times 10^6$ | $1.51 \times 10^6$ | $3.82 \times 10^6$ | $1.36 \times 10^6$ | $4.49 \times 10^6$ | $1.47 \times 10^6$ | $3.69 \times 10^6$ | $1.34 \times 10^6$ | $4.38 \times 10^6$ | $1.17 \times 10^6$ | $4.33 \times 10^6$ | $1.03 \times 10^6$ | $\mathbf{3.66 \times 10^6}$ | $1.34 \times 10^6$ |
| DBLP | $3.94 \times 10^6$ | $986.02 \times 10^3$ | $3.17 \times 10^6$ | $614.86 \times 10^3$ | $3.70 \times 10^6$ | $858.93 \times 10^3$ | $3.01 \times 10^6$ | $557.90 \times 10^3$ | $\mathbf{2.55 \times 10^6}$ | $354.03 \times 10^3$ | $2.89 \times 10^6$ | $506.72 \times 10^3$ | $2.81 \times 10^6$ | $393.38 \times 10^3$ |
| Epinions | $12.92 \times 10^6$ | $76.81 \times 10^3$ | $4.71 \times 10^6$ | $47.54 \times 10^3$ | $9.06 \times 10^6$ | $65.59 \times 10^3$ | $\mathbf{4.70 \times 10^6}$ | $47.51 \times 10^3$ | $9.80 \times 10^6$ | $16.73 \times 10^3$ | $8.82 \times 10^6$ | $16.68 \times 10^3$ | $5.06 \times 10^6$ | $65.31 \times 10^3$ |
| HighSchool | $4.59 \times 10^3$ | $45.26$ | $3.50 \times 10^3$ | $8.16$ | $4.44 \times 10^3$ | $37.66$ | $3.35 \times 10^3$ | $6.38$ | $4.29 \times 10^3$ | $9.00$ | $3.43 \times 10^3$ | $7.00$ | $\mathbf{3.29 \times 10^3}$ | $8.00$ |
| Last.fm | $164.67 \times 10^3$ | $57.04$ | $\mathbf{150.25 \times 10^3}$ | $37.64$ | $163.35 \times 10^3$ | $42.10$ | $\mathbf{150.25 \times 10^3}$ | $36.94$ | $161.53 \times 10^3$ | $3.00$ | $160.66 \times 10^3$ | $4.00$ | $151.60 \times 10^3$ | $37.00$ |
| PrimarySchool | $6.95 \times 10^3$ | $16.44$ | $5.01 \times 10^3$ | $1.20$ | $6.80 \times 10^3$ | $15.12$ | $\mathbf{4.92 \times 10^3}$ | $1.04$ | $6.48 \times 10^3$ | $5.00$ | $6.27 \times 10^3$ | $5.00$ | $5.46 \times 10^3$ | $2.00$ |
| ProsperLoans | $1.82 \times 10^6$ | $39.60 \times 10^3$ | $1.30 \times 10^6$ | $3.75 \times 10^3$ | $1.81 \times 10^6$ | $26.06 \times 10^3$ | $\mathbf{1.28 \times 10^6}$ | $3.70 \times 10^3$ | $\mathbf{1.28 \times 10^6}$ | $1.54 \times 10^3$ | $1.30 \times 10^6$ | $1.13 \times 10^3$ | $1.39 \times 10^6$ | $7.49 \times 10^3$ |
| StackOverflow | $12.39 \times 10^6$ | $1.74 \times 10^6$ | $8.83 \times 10^6$ | $308.66 \times 10^3$ | $11.91 \times 10^6$ | $1.27 \times 10^6$ | $\mathbf{8.65 \times 10^6}$ | $237.36 \times 10^3$ | $9.90 \times 10^6$ | $106.58 \times 10^3$ | $9.26 \times 10^6$ | $13.78 \times 10^3$ | $10.81 \times 10^6$ | $188.44 \times 10^3$ |
| Wikipedia | $6.74 \times 10^6$ | $276.14 \times 10^3$ | $5.31 \times 10^6$ | $157.77 \times 10^3$ | $6.44 \times 10^6$ | $246.29 \times 10^3$ | $\mathbf{5.26 \times 10^6}$ | $168.80 \times 10^3$ | $5.84 \times 10^6$ | $113.64 \times 10^3$ | $5.84 \times 10^6$ | $108.82 \times 10^3$ | $5.83 \times 10^6$ | $209.68 \times 10^3$ |
| WikiTalk | $6.29 \times 10^6$ | $2.77 \times 10^6$ | $3.72 \times 10^6$ | $381.88 \times 10^3$ | $5.41 \times 10^6$ | $1.99 \times 10^6$ | $\mathbf{3.38 \times 10^6}$ | $485.66 \times 10^3$ | $3.68 \times 10^6$ | $351.73 \times 10^3$ | $5.13 \times 10^6$ | $1.69 \times 10^6$ | NA | NA |
| **tot. average** | $4.91 \times 10^6$ | $7.40 \times 10^5$ | $3.10 \times 10^6$ | $2.87 \times 10^5$ | $4.30 \times 10^6$ | $5.93 \times 10^5$ | $3.01 \times 10^6$ | $2.84 \times 10^5$ | $3.76 \times 10^6$ | $2.11 \times 10^5$ | $3.77 \times 10^6$ | $3.37 \times 10^5$ | $3.30 \times 10^6$ | $2.45 \times 10^5$ |

- D-MIL outperforms MIL
- MIL_R and D-MIL_R produce better solutions than MIL and D-MIL
- D-MIL_R is the best performing method and outperforms competing methods

- D-MIL yelds a smaller number of clusters than MIL
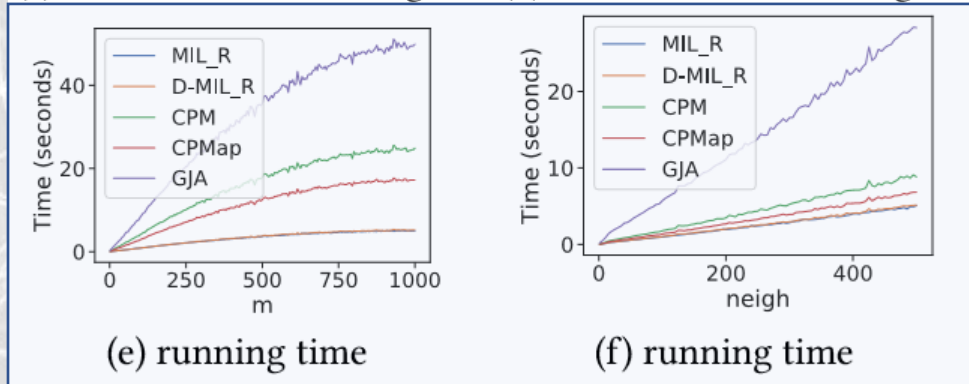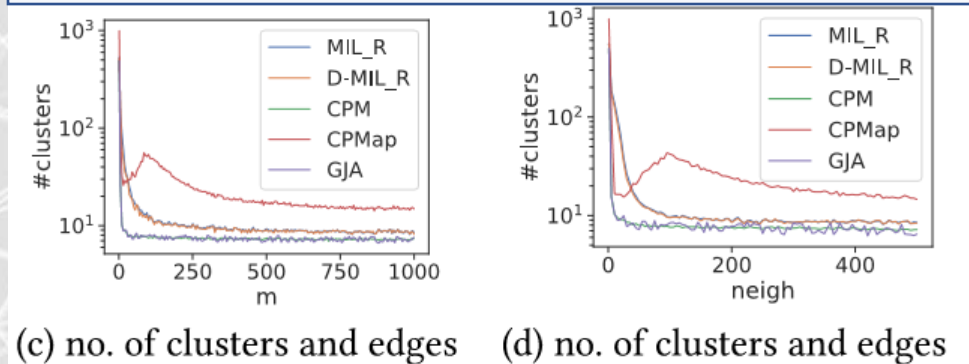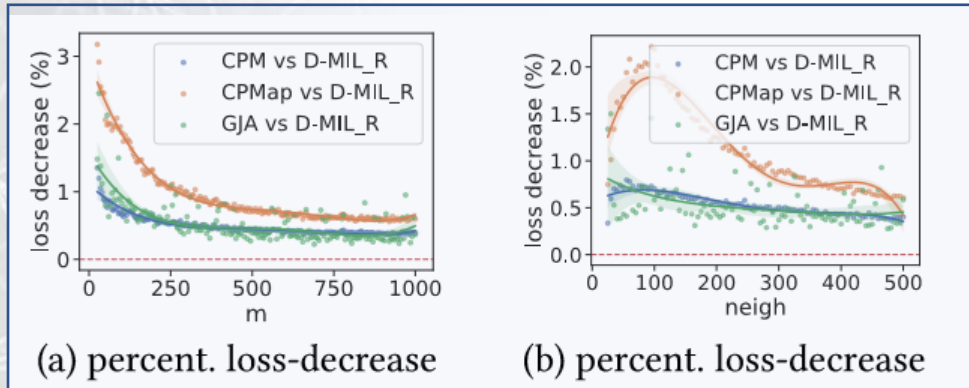- MIL and D-MIL produce more clusters than MIL_R and D-MIL_R

# Time performances (in seconds) on real data

| | MIL | MIL_R | opt. time | D-MIL | D-MIL_R | opt. time | CPM [21] | GJA [9] | CPMap [7] |
|---|---|---|---|---|---|---|---|---|---|
| *Amazon* | 8.63 | 347.77 | 339.14 | 97.40 | 427.28 | 329.88 | 2 248.9 | 1 020 122.23 | 669.114 |
| *DBLP* | 6.11 | 189.63 | 183.52 | 71.15 | 251.24 | 180.09 | 1 570.41 | 147 159.68 | 601.044 |
| *Epinions* | 5.90 | 327.27 | 321.38 | 18.90 | 348.11 | 329.21 | 797.71 | 34 998.9 | 592.901 |
| *High School* | 0.00 | 0.04 | 0.04 | 0.01 | 0.04 | 0.03 | 0.2 | 0.19 | 2.716 |
| *Last.fm* | 0.03 | 3.48 | 3.45 | 0.14 | 3.72 | 3.58 | 7.73 | 21.54 | 10.467 |
| *PrimarySchool* | 0.00 | 0.06 | 0.05 | 0.01 | 0.05 | 0.04 | 0.125 | 0.1 | 3.698 |
| *ProsperLoans* | 0.70 | 48.74 | 48.04 | 4.31 | 52.06 | 47.75 | 179.78 | 30 152.47 | 116.59 |
| *StackOverflow* | 7.88 | 319.67 | 311.79 | 105.68 | 397.39 | 291.72 | 2 465.76 | 1 140 054.23 | 1519.943 |
| *Wikipedia* | 2.41 | 150.03 | 147.62 | 19.05 | 160.69 | 141.64 | 826.93 | 189 345.74 | 316.438 |
| *WikiTalk* | 13.92 | 203.49 | 189.56 | 129.10 | 300.68 | 171.58 | 1 165.01 | 650 282.4 | NA |

- MIL is faster than D-MIL
- Running times for MIL_R and D-MIL_R are dominated by optimization time (Hill Climbing)
- MIL_R is faster than D-MIL_R
- Competing methods are always outperformed by our algorithms

# Comparison with competitors on synthetic data



(a) percent. loss-decrease

(b) percent. loss-decrease

(c) no. of clusters and edges

(d) no. of clusters and edges

(e) running time

(f) running time

For the BA model, $m$ is the number of edges to attach with a new vertex

For the WS model, $neigh$ is the distance within which two vertices will be connected

- Superiority of D-MIL_R w.r.t. competing methods in terms of effectiveness
- Our algorithms are faster than competing methods

# Conclusions & Future Work

**Summary**:

- We introduced the problem of optimizing the overall interaction among a set of entities whose interaction patterns depend on their cluster memberships

- We exploit the connection with correlation clustering to develop both approximation algorithms and heuristics (for the minimization formulation)

- Experimental evaluation of our algorithms on both synthetic and real network datasets:
  - Better effectiveness and efficiency than competing methods

**Future Work**:

- Clustering properties:
  - Overlapping
  - Size bounds

- Probability distributions of interactions are not given

# Thank you!
# Questions?

## In and Out: Optimizing Overall Interaction in Probabilistic Graphs under Clustering Constraints

Domenico Mandaglio
DIMES Dept., University of Calabria
Rende (CS), Italy
d.mandaglio@dimes.unical.it

Andrea Tagarelli
DIMES Dept., University of Calabria
Rende (CS), Italy
andrea.tagarelli@unical.it

Francesco Gullo
UniCredit, R&D Dept.
Rome, Italy
gullof@acm.org

### ABSTRACT

We study two novel clustering problems in which the pairwise interactions between entities are characterized by probability distributions and conditioned by external factors within the environment where the entities interact. This covers any scenario where a set of actions can alter the entities' interaction behavior. In particular, we consider the case where the interaction conditioning factors can be modeled as cluster memberships of entities in a graph and the goal is to partition a set of entities such as to maximize the overall vertex interactions or, equivalently, minimize the loss of

behaviors into a representation of user interaction patterns [14]. A common way of modeling uncertainty in a graph, which we refer to in this work, is to associate each pair of (linked) users with a probability value that expresses the likelihood of *observing* and *quantifying* an interaction between the two users. In this regard, one important aspect is that the modeling of user interactions should also account for exogenous conditions or events that occur within the social environment where the users belong to, which indeed can significantly affect the users' interaction behaviors. For example, delivering a post on a user's page (e.g., Facebook wall) that contains