

# Correlation Clustering with Global Weight Bounds

Domenico Mandaglio<sup>1</sup>, Andrea Tagarelli<sup>1</sup>, Francesco Gullo<sup>2</sup>

<sup>1</sup>DIMES Dept., University of Calabria, Rende (CS), Italy; <sup>2</sup>UniCredit, Rome, Italy



- We focus for the first time on **global weight bounds** for correlation clustering, focusing on its minimization objective.
- We identify a sufficient condition on input weights' aggregate functions to **extend the validity range of the approximation guarantees** of existing correlation-clustering algorithms beyond the probability constraint.
- We experimentally assess that our condition is an effective **indicator of the empirical performance** of existing probability-constraint-aware correlation-clustering algorithms.
- We showcase our results in a real-world scenario of **fair clustering**.

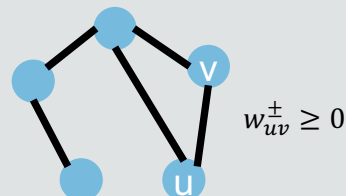
## Min-Disagreement Correlation Clustering Problem

Given an undirected graph  $G = (V, E)$ , with vertex set  $V$  and edge set  $E \subseteq V \times V$ , and weights  $w_{uv}^+, w_{uv}^- \in \mathbb{R}_0^+$  for all edges  $(u, v) \in E$ , find a clustering  $\mathcal{C}: V \rightarrow \mathbb{N}^+$  that minimizes:

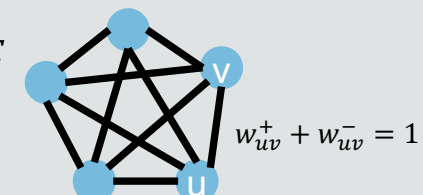
$$\sum_{\substack{(u,v) \in E \\ \mathcal{C}(u) = \mathcal{C}(v)}} w_{uv}^- + \sum_{\substack{(u,v) \in E \\ \mathcal{C}(u) \neq \mathcal{C}(v)}} w_{uv}^+$$

Any  $w_{uv}^+$  (resp.  $w_{uv}^-$ ) weight expresses the benefit of clustering  $u$  and  $v$  together (resp. separately)

1. General graph and general weights
  - Linear Programming + Rounding with  $O(\log n)$  approximation guarantees



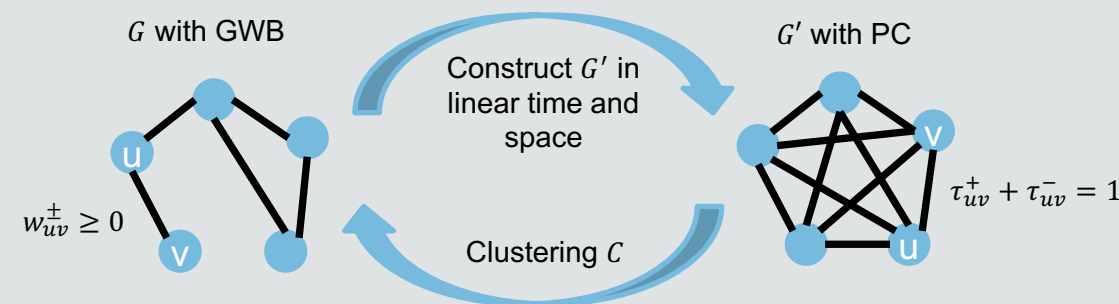
2.  $E = \binom{V}{2}$  and  $w_{uv}^+ + w_{uv}^- = 1 \forall (u, v) \in E$ 
  - Pivot algorithm with (expected) 5-approximation guarantees and  $O(|E|)$  time complexity



## Global Weight Bound (GWB) condition

$$\binom{|V|}{2}^{-1} \sum_{(u,v) \in E} w_{uv}^+ + \binom{|V|}{2}^{-1} \sum_{(u,v) \in E} w_{uv}^- \geq \max_{(u,v) \in E} |w_{uv}^+ - w_{uv}^-|$$

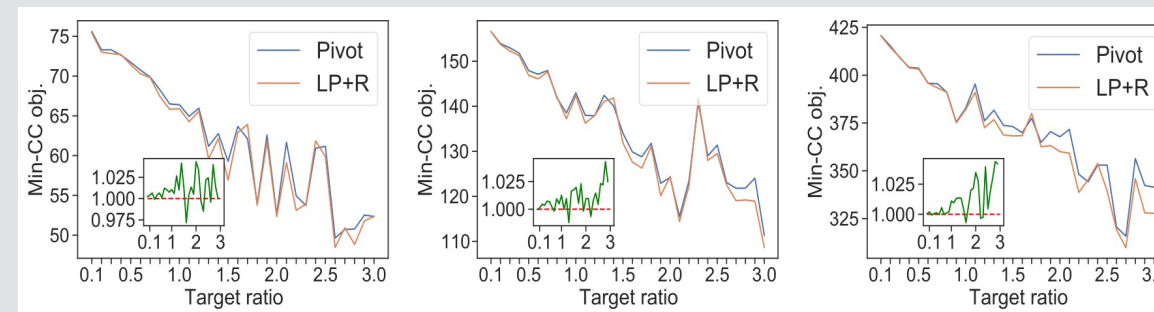
Strict approximation-preserving reduction:



An  $\alpha$ -approximate clustering on  $G'$  is also  $\alpha$ -approximate clustering on  $G$  too

- **Practical benefits:**
  - Extend the validity range of the approximation guarantees of algorithms for correlation clustering (Exp1)
  - Application to feature selection for fair clustering (Exp2)
- **Theoretical benefits:** enable better theoretical results on complex problems which exploit correlation clustering as a building block
- **Benefits for the research community:** brand new line of research

## Exp1: Analysis of the global-weight-bounds condition



A better fulfilment of our GWB leads to Pivot's performance closer to the linear programming approach's one (LP+R, for short), and vice versa.

## Exp2: Application to fair clustering

**Data.** Relational datasets describing a set of objects  $X$  defined over a set of attributes  $A$  (numerical or categorical) that can be divided into:

- *Fairness-aware* (or *sensitive*) attributes  $A^F$
- *Non-sensitive* attributes  $A^{-F}$

**Fair clustering objective:**

- **non-sensitive attributes:** minimize the inter-cluster similarities and maximize the intra-cluster similarities
- **sensitive attributes:** minimize the intra-cluster similarities and maximize the inter-cluster similarities

**Mapping to Correlation Clustering instance:**

$$w_{uv}^+ := \varphi^+ (\alpha_N^{-F} \cdot sim_{A_N^{-F}}(u, v) + (1 - \alpha_N^{-F}) \cdot sim_{A_C^{-F}}(u, v))$$

$$w_{uv}^- := \varphi^- (\alpha_N^F \cdot sim_{A_N^F}(u, v) + (1 - \alpha_N^F) \cdot sim_{A_C^F}(u, v))$$

$$\alpha_N^F = \frac{|A_N^F|}{|A_N^F| + |A_C^F|}, \alpha_N^{-F} = \frac{|A_N^{-F}|}{|A_N^{-F}| + |A_C^{-F}|}, \varphi^+ = \exp\left(\frac{|A^F|}{|A^F| + |A^{-F}| - 1}\right), \varphi^- = \exp\left(\frac{|A^{-F}|}{|A^F| + |A^{-F}| - 1}\right)$$

**Attribute selection for fair clustering.** Given a set of objects  $X$  defined over the attribute sets  $A^F$  and  $A^{-F}$ , find maximal subsets  $S_F \subseteq A^F$  and  $S_{-F} \subseteq A^{-F}$ , with  $|S_F| \geq 1$  and  $|S_{-F}| \geq 1$ , s.t. the above correlation-clustering weights satisfy the GWB condition.

	#it	target ratio	%( $w^+ > w^-$ )	orig.-weights Min-CC obj.	avg. Eucl. fairness	avg. #clusts.	intra-clust $\mathcal{A}^{-F}$	intra-clust $\mathcal{A}^F$	inter-clust $\mathcal{A}^{-F}$	inter-clust $\mathcal{A}^F$	time (seconds)
<i>Adult</i>											
initial	-	1.086	90.34	1.1915E+08	0.082	77	0.699	0.672	0.378	0.181	-
Hlv	12	0.986	93.19	1.122659E+08	<b>0.031</b>	9	0.465	0.326	0.347	0.194	545.249
Hlv_B	12	0.765	78.09	1.119757E+08	0.039	69	0.608	0.547	0.375	0.184	529.674
Hmv	5	0.974	90.83	1.21187E+08	0.094	79	0.689	0.687	0.373	<b>0.203</b>	220.056
Hmv_B	4	0.936	87.39	1.25516E+08	0.109	905	0.963	0.96	0.377	0.199	<b>178.813</b>
Hlv_BW	5	0.963	83.17	1.343503E+08	0.152	1479	<b>0.969</b>	0.964	0.384	0.199	217.333
Hmv_SW	9	0.926	91.41	1.159874E+08	0.037	5	0.451	<b>0.308</b>	<b>0.329</b>	0.195	380.875
Greedy	2	0.967	92.36	<b>1.094787E+08</b>	0.036	32	0.668	0.654	0.361	0.195	595.610

The GWB condition can help to define weights so as to account for both an effective representation of the semantics underlying objects' features, and the peculiarities that make the downstream correlation-clustering algorithm effective.