

Bipartite correlation clustering: Maximizing agreements

Megasthenis Asteris, Anastasios Kyrillidis, Dimitris Papailiopoulos, Alexandros Dimakis

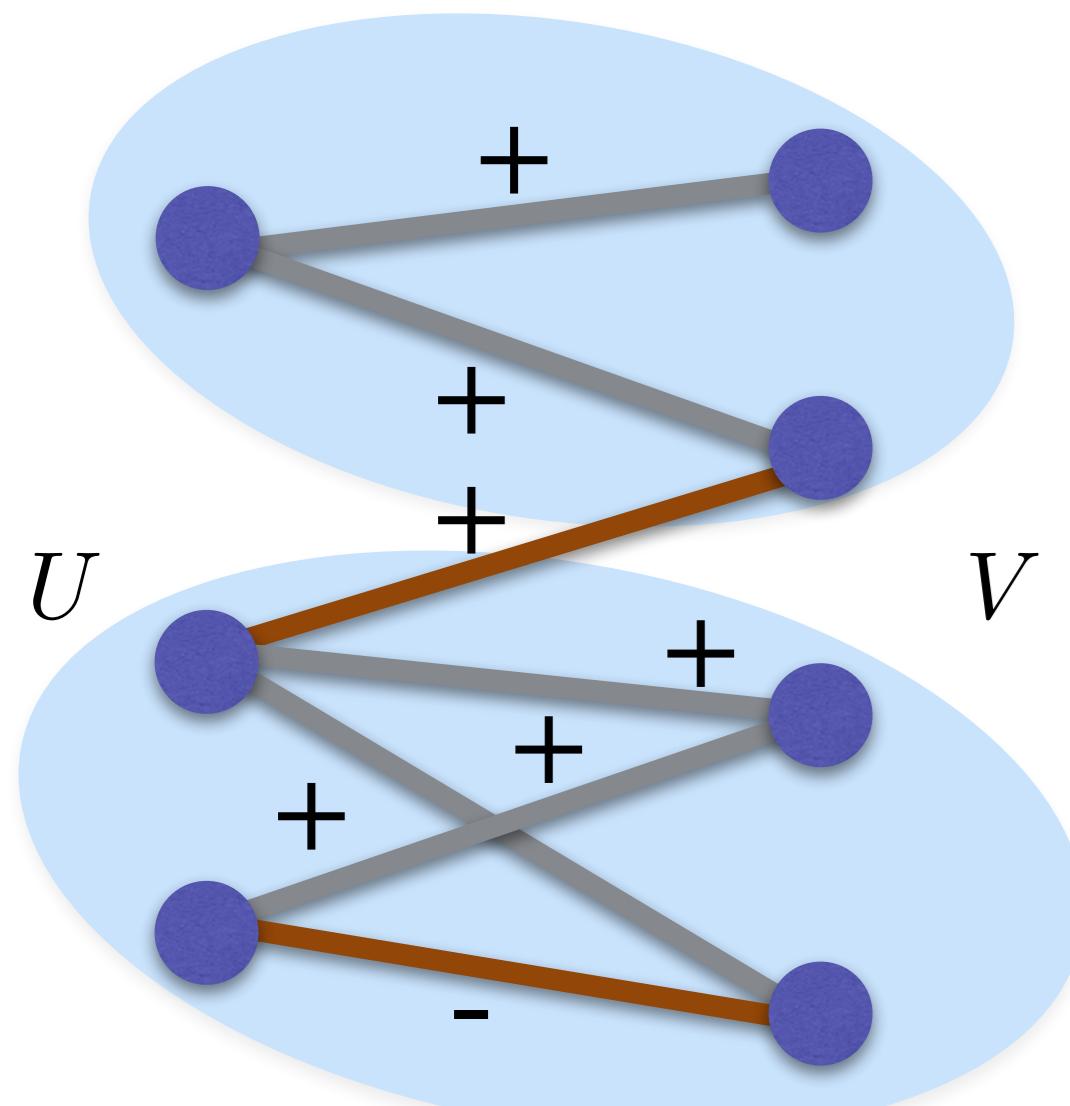
[What is bipartite correlation clustering (BCC)?]

Task: given a bipartite graph $G = (U, V, E)$, partition objects into clusters, based on pairwise relationships, in order to maximize the number of **agreements**.

Two vertices are similar: '+' edge / otherwise: '-' edge.

Number of agreements: # of '+' within resulting clusters
+ # of '-' across resulting clusters.

Example:



of output clusters: *not part of the input*

Might be useful to restrict number of output clusters:

k-BCC

- Applications:
1. Recommendation systems
 2. Gene expression data analysis
 3. Graph partitioning

Number of agreements: 5 '+'
Number of disagreements: 1 '+' and 1 '-'

[Contributions]

We develop a novel approximation algorithm for k-BCC with provable guarantees for the "maximizing agreements" problem.

In particular:

- **k-BCC:** given $G = (U, V, E)$, a number of clusters k and any accuracy parameter $\delta \in (0, 1)$, our algorithm computes a clustering of at most k clusters and within $(1 - \delta)$ -factor from the optimal.
- **BCC:** We show that, in order to achieve a $(1 - \delta)$ -approximation in the "maximizing agreements" objective, it suffices to use at most $O(\delta^{-1})$ clusters.
- **Proposed scheme:** We propose a new algorithm, based on constrained bilinear maximization techniques, that yields an EPTAS for the unconstrained BCC problem

[Related work]

Ref.	Min./Max.	Guarantee	Complexity	D/R	Setting
[13]	MAXAGREE	(E)PTAS	$n/\delta \cdot k^{O(\delta^{-2} \log k \log(1/\delta))}$	R	k -CC
[13]	MINDISAGREE	PTAS	$n^{O(100k/\delta^2)} \cdot \log n$	R	k -CC
[15]	MINDISAGREE	PTAS	$n^{O(9k/\delta^2)} \cdot \log n$	R	k -CC
[11, 9]	MINDISAGREE	$O(\log n)$ -OPT	LP	D	Inc. CC
[17]	MAXAGREE	0.766-OPT	SDP	D	Inc. CC
[5]	MINDISAGREE	11-OPT	LP	D	BCC
[2]	MINDISAGREE	4-OPT	LP	D	BCC
[2]	MINDISAGREE	4-OPT	$ E $	R	BCC
Ours	MAXAGREE	(E)PTAS	$2^{O(k/\delta^2 \cdot \log \sqrt{k}/\delta)} \cdot (\delta^{-2} + k) \cdot n + T_{SVD}(\delta^{-2})$	R	k -BCC
Ours	MAXAGREE	(E)PTAS	$2^{O(\delta^{-3} \cdot \log \delta^{-3})} \cdot O(\delta^{-2}) \cdot n + T_{SVD}(\delta^{-2})$	R	BCC

Table 1: Summary of results on BCC and related problems. For each scheme, we indicate the problem setting, the objective (MAXAGREE/MINDISAGREE), the guarantees (c -OPT implies a multiplicative factor approximation), and its computational complexity (n denotes the total number of vertices, LP/SDP denotes the complexity of a linear/semidefinite program, and $T_{SVD}(r)$ the time required to compute a rank- r truncated SVD of a $n \times n$ matrix). The D/R column indicates whether the scheme is deterministic or randomized.

[k-BCC as a bilinear maximization problem]

Setting: For any instance $G = (U, V, E)$ of k-BCC problem, let

$$\begin{aligned} \mathbf{B} &\in \{\pm 1\}^{|U| \times |V|} && : \text{Bi-adjacency matrix} \\ \mathbf{X} &\in \{0, 1\}^{|U| \times k} && : \text{Cluster assignment matrices} \\ \mathbf{Y} &\in \{0, 1\}^{|V| \times k} \end{aligned} \quad \left. \begin{array}{l} \text{(i.e., } X_{ij} = 1 \text{ iff vertex } i \in U \text{ is} \\ \text{assigned to cluster } C_j.\text{)} \end{array} \right\}$$

Lemma 3.2 : For any clustering \mathcal{C} of $U \cup V$ into k clusters, the number of agreements achieved by \mathcal{C} is:

$$\text{AGREE}(\mathcal{C}) = \text{TR}(\mathbf{X}^\top \mathbf{B} \mathbf{Y}) + |E^-|$$

where E^- denotes the set of negative edges in the graph.

Thus, computing a k -clustering that achieves the **maximum number of agreements**, builds down to a **constrained bilinear maximization problem**:

$$\text{MAXAGREE}[k] = \max_{\mathbf{X} \in \mathcal{X}, \mathbf{Y} \in \mathcal{Y}} \text{TR}(\mathbf{X}^\top \mathbf{B} \mathbf{Y}) + |E^-|,$$

where

$$\begin{aligned} \mathcal{X} &\triangleq \{\mathbf{X} \in \{0, 1\}^{m \times k} : \|\mathbf{X}\|_{\infty, 1} = 1\}, \\ \mathcal{Y} &\triangleq \{\mathbf{Y} \in \{0, 1\}^{n \times k} : \|\mathbf{Y}\|_{\infty, 1} = 1\}. \end{aligned}$$

[Bilinear maximization framework]

Consider the following maximization problem:

$$\max_{\mathbf{X} \in \mathcal{X}, \mathbf{Y} \in \mathcal{Y}} \text{TR}(\mathbf{X}^\top \mathbf{A} \mathbf{Y})$$

where \mathbf{A} is a real $m \times n$ input matrix.

In general **NP-hard**; we approximate the above problem as:

$$\max_{\mathbf{X} \in \mathcal{X}, \mathbf{Y} \in \mathcal{Y}} \text{TR}(\mathbf{X}^\top \tilde{\mathbf{A}} \mathbf{Y})$$

where $\tilde{\mathbf{A}}$ is the rank- r approximation of \mathbf{A} .

[Algorithm for bilinear maximization]

Algorithm 1 BiLinearLowRankSolver

```

input :  $m \times n$  real matrix  $\tilde{\mathbf{A}}$  of rank  $r$ ,  $\epsilon \in (0, 1)$ 
output  $\tilde{\mathbf{X}} \in \mathcal{X}$ ,  $\tilde{\mathbf{Y}} \in \mathcal{Y}$ 
1:  $\mathcal{C} \leftarrow \{\}$  {Candidate solutions}
2:  $\tilde{\mathbf{U}}, \tilde{\Sigma}, \tilde{\mathbf{V}} \leftarrow \text{SVD}(\tilde{\mathbf{A}})$  { $\tilde{\Sigma} \in \mathbb{R}^{r \times r}$ }
3: for each  $\mathbf{C} \in (\epsilon\text{-net of } \mathbb{B}_2^{r-1})^{\otimes k}$  do
4:    $\mathbf{L} \leftarrow \tilde{\mathbf{U}} \tilde{\Sigma} \mathbf{C}$  { $\mathbf{L} \in \mathbb{R}^{m \times k}$ }
5:    $\mathbf{X} \leftarrow P_{\mathcal{X}}(\mathbf{L})$ 
6:    $\mathbf{R} \leftarrow \mathbf{X}^\top \tilde{\mathbf{A}}$ 
7:    $\mathbf{Y} \leftarrow P_{\mathcal{Y}}(\mathbf{R})$ 
8:    $\mathcal{C} \leftarrow \mathcal{C} \cup \{(\mathbf{X}, \mathbf{Y})\}$ 
9: end for
10:  $(\tilde{\mathbf{X}}, \tilde{\mathbf{Y}}) \leftarrow \arg \max_{(\mathbf{X}, \mathbf{Y}) \in \mathcal{C}} \text{TR}(\mathbf{X}^\top \tilde{\mathbf{U}} \tilde{\Sigma} \tilde{\mathbf{V}}^\top \mathbf{Y})$ 

```

Lemma 3.3.

$$(\mathbf{X}_*, \mathbf{Y}_*) \triangleq \arg \max_{\mathbf{X} \in \mathcal{X}, \mathbf{Y} \in \mathcal{Y}} \text{TR}(\mathbf{X}^\top \mathbf{A} \mathbf{Y}),$$

where \mathcal{X} and \mathcal{Y} satisfy the conditions of Lemma 3.2. Let $\tilde{\mathbf{A}}$ be a rank- r approximation of \mathbf{A} , and $\tilde{\mathbf{X}} \in \mathcal{X}$, $\tilde{\mathbf{Y}} \in \mathcal{Y}$ be the output of Alg. 1 with input $\tilde{\mathbf{A}}$ and accuracy ϵ . Then,

$$\begin{aligned} \text{TR}(\mathbf{X}_*^\top \mathbf{A} \mathbf{Y}_*) - \text{TR}(\tilde{\mathbf{X}}^\top \mathbf{A} \tilde{\mathbf{Y}}) \\ \leq 2 \cdot (\epsilon \sqrt{k} \cdot \|\tilde{\mathbf{A}}\|_2 + \|\mathbf{A} - \tilde{\mathbf{A}}\|_2) \cdot \mu_{\mathcal{X}} \cdot \mu_{\mathcal{Y}}. \end{aligned}$$

Comments:

- Our approach has complexity **exponential** in the rank of $\tilde{\mathbf{A}}$
- The quality of the output depends on the spectrum of \mathbf{A} and the level of rank approximation r .

[Our k-BCC algorithm]

Algorithm 4 k-BCC/MAXAGREE

```

input : Bi-adjacency matrix  $\mathbf{B} \in \{\pm 1\}^{m,n}$ ,  
Target number of  $k$ ,  
Accuracy  $\delta \in (0, 1)$ .
output Clustering  $\tilde{\mathcal{C}}^{(k)}$  of  $U \cup V$  such that  
 $\text{AGREE}(\tilde{\mathcal{C}}^{(k)}) \geq (1 - \delta) \cdot \text{AGREE}(\mathcal{C}_*^{(k)})$ .
1: Set up parameters:  
 $\epsilon \leftarrow 2^{-3} \cdot \delta \cdot k^{-1/2}$ ,  $r \leftarrow 2^6 \cdot \delta^{-2} - 1$ .
2: Return output of Alg. 3 for input  $(\mathbf{B}, k, r, \epsilon)$ .

```

Theorem 1. (k-BCC.) For any instance $(G = (U, V, E), k)$ of the k -BCC problem with bi-adjacency matrix \mathbf{B} and for any desired accuracy parameter $\delta \in (0, 1)$, Algorithm 4 computes a clustering $\tilde{\mathcal{C}}^{(k)}$ of $U \cup V$ into at most k clusters, such that

$$\text{AGREE}(\tilde{\mathcal{C}}^{(k)}) \geq (1 - \delta) \cdot \text{AGREE}(\mathcal{C}_*^{(k)}),$$

where $\mathcal{C}_*^{(k)}$ is the optimal k -clustering, in time $2^{O(k/\delta^2 \cdot \log \sqrt{k}/\delta)} \cdot (\delta^{-2} + k) \cdot (|U| + |V|) + T_{SVD}(\delta^{-2})$.

[A (E)PTAS algorithm for BCC]

Algorithm 5 A PTAS for BCC/MAXAGREE

```

input : Bi-adjacency matrix  $\mathbf{B} \in \{\pm 1\}^{m,n}$ ,  
Accuracy  $\delta \in (0, 1)$ .
output Clustering  $\tilde{\mathcal{C}}$  of  $U \cup V$  (into at most  
 $2^3 \cdot \delta^{-1}$  clusters) such that  
 $\text{AGREE}(\tilde{\mathcal{C}}) \geq (1 - \delta) \cdot \text{AGREE}(\mathcal{C}_*)$ .
1: Set up parameters:  
 $k \leftarrow 2^3 \cdot \delta^{-1}$ ,  $\epsilon \leftarrow 2^{-6} \cdot \delta^2$ ,  $r \leftarrow 2^8 \cdot \delta^{-2} - 1$ .
2: Return output of Alg. 3 for input  $(\mathbf{B}, k, r, \epsilon)$ .

```

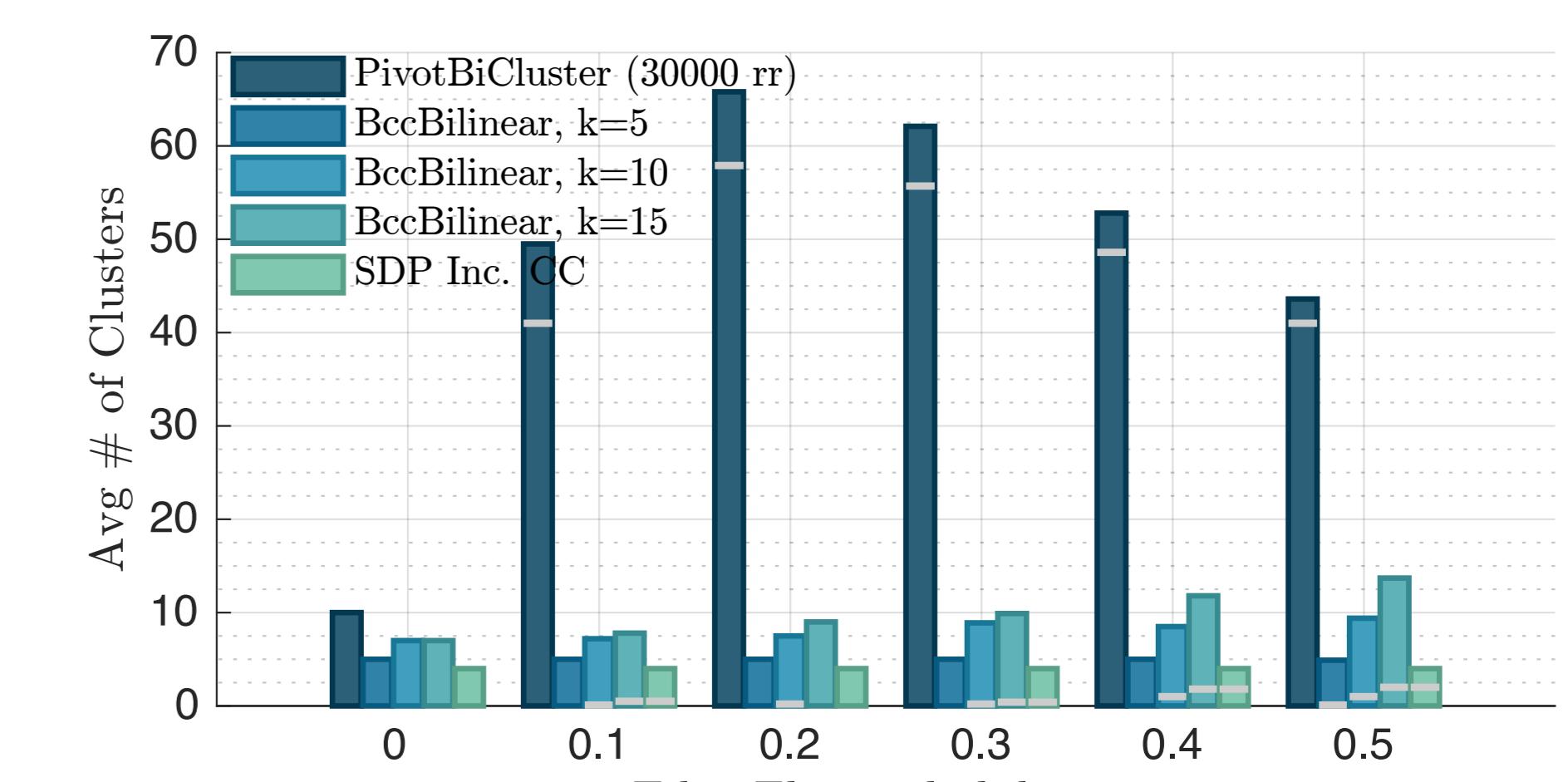
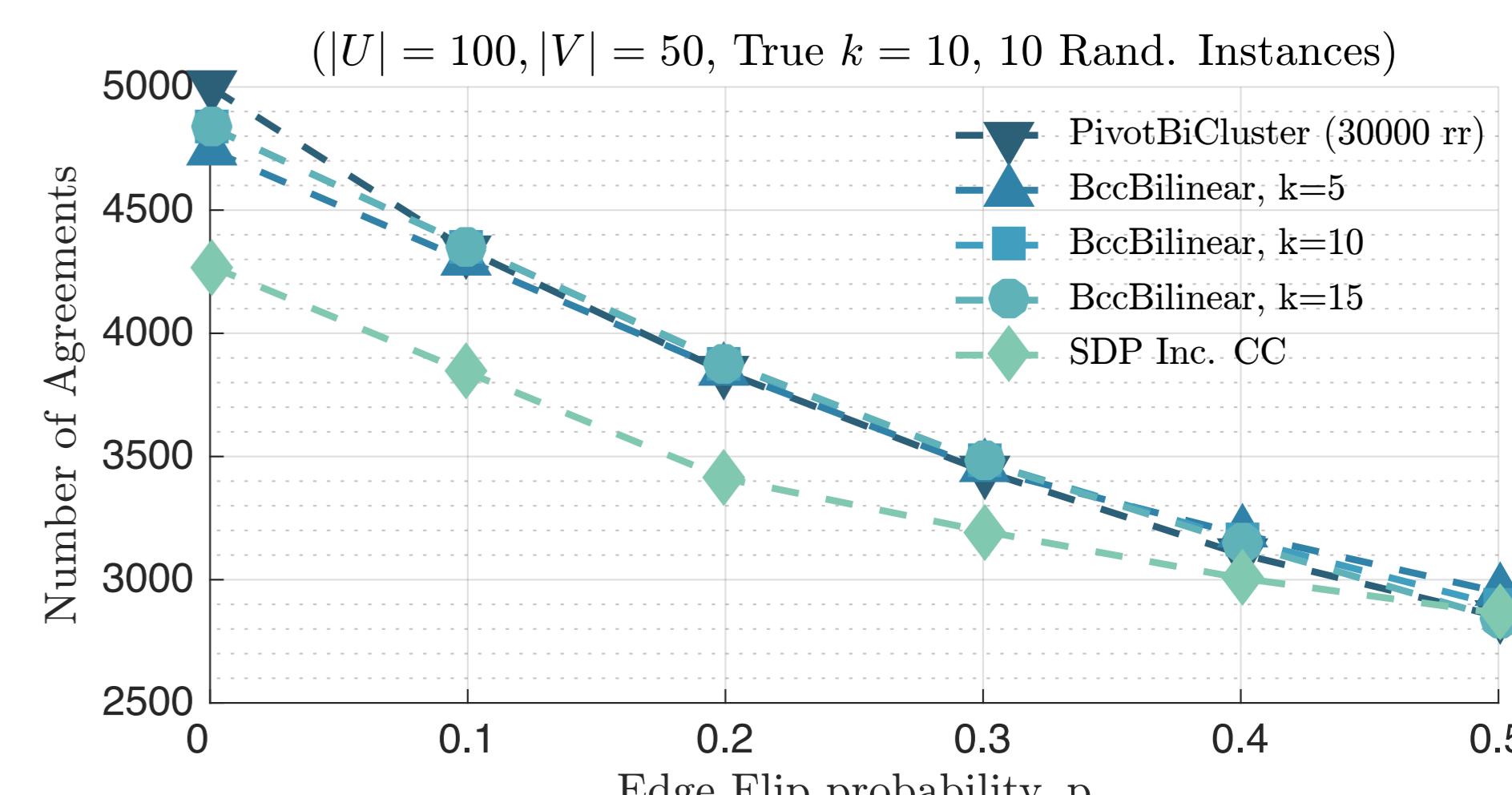
Theorem 2. (PTAS for BCC.) For any instance $G = (U, V, E)$ of the BCC problem with bi-adjacency matrix \mathbf{B} and for any desired accuracy parameter $\delta \in (0, 1)$, Algorithm 5 computes a clustering $\tilde{\mathcal{C}}$ of $U \cup V$ into (at most) $2^3 \cdot \delta^{-1}$ clusters, such that

$$\text{AGREE}(\tilde{\mathcal{C}}) \geq (1 - \delta) \cdot \text{AGREE}(\mathcal{C}_*),$$

[In practice...]

Generate BCC instances as follows:

- Generate complete bipartite graphs with $k=5$ clusters
- Perturb label signs with probability p
- Generate 10 random instances and compute vertex clustering using PivotBiCluster (Ailon et al 2012), SDP approach (Swamy 2004) and ours.



	$p = 0.0$	0.1	0.2	0.3	0.4	0.5
PivotBiCluster [2]	124	83	66	54	45	39
BccBil. (k=5)	69	70	71	71	71	71
BccBil. (k=10)	73	74	75	75	75	75
BccBil. (k=15)	78	79	79	80	80	80
SDP Inc CC [17]	198	329	257	258	266	314

Average runtimes/instance (in seconds).