

Correlation Clustering with Noisy Input

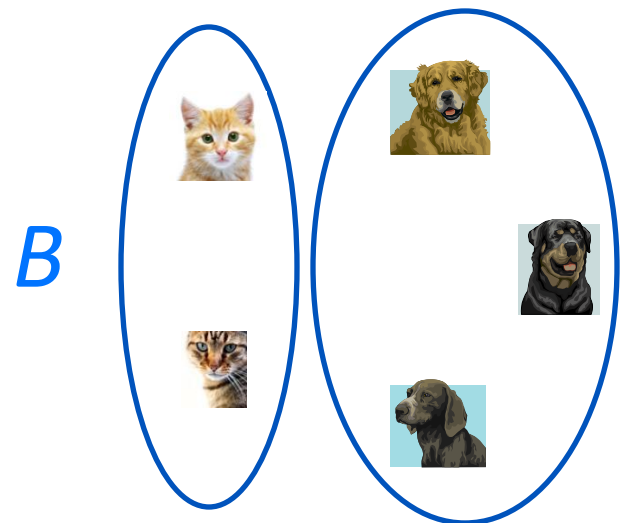
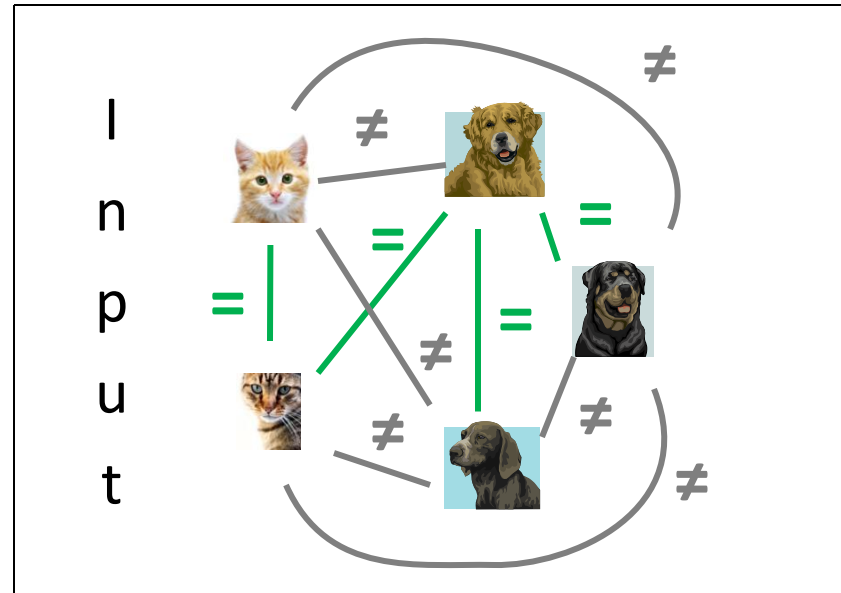
Claire Mathieu
Warren Schudy

Brown University

SODA 2010

Noisy Correlation Clustering Model

- Unknown **base clustering** B of n objects
- Noise: each edge is controlled by an adversary with probability p and “tells the truth” otherwise
- Problem: reconstruct B from the edge labels



One of our results

- **Theorem:** assume $p \leq 1/3$. If all clusters have size at least $\alpha_1 \sqrt{n}$ then the natural **semi-definite program (SDP)** recovers B **exactly** with high probability.
- **Previous best:** $\alpha_2 \sqrt{n \log n}$ [Bansal, Blum, Chawla '04, Shamir and Tsur '07], **combinatorial**.
- See paper for other results (including approximation algorithms)


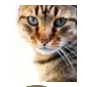



Plan

- The semi-definite program
- Its dual
- Using the dual

Clusterings



- Clusterings are represented by 0/1 matrices:
 $X_{ij}=1$: i and j in same cluster

	1	1	0	0	0
	1	1	0	0	0
	0	0	1	1	1
	0	0	1	1	1
	0	0	1	1	1

- In general a clustering satisfies:

$$X = \sum_k v_k v_k^T \text{ for some 0/1 orthogonal vectors } v_1, v_2, \dots, v_m, \\ \text{one per cluster}$$

- E.g.

$$v_{cat} : \begin{bmatrix} 1 & 1 & 0 & 0 & 0 \end{bmatrix}$$

$$v_{dog} : \begin{bmatrix} 0 & 0 & 1 & 1 & 1 \end{bmatrix}$$

Relaxation of clusterings

Relaxation







- ~~Clustering~~

- $X = \sum_k v_k v_k^T$ for some ~~0/1~~ vectors v_1, v_2, \dots, v_m

- $X_{ij} = 1$ for all i

- $X_{ij} \geq 0$ for all i, j

X

			
	1	0.7	0
	0.7	1	0.7
	0	0.7	1

- The following are equivalent (X symmetric):

- $X = \sum_k v_k v_k^T$ for some vectors v_1, v_2, \dots, v_m

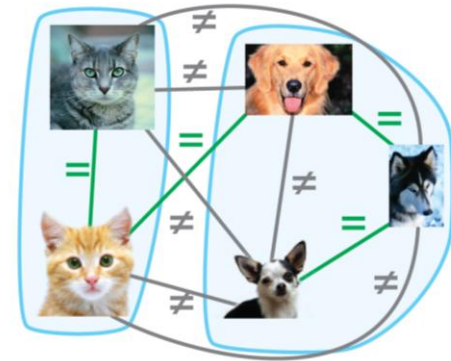
- X is *positive semi-definite (p.s.d.)*

Objective

- Maximize number of agreements:

$$\max \sum_{i < j} \begin{cases} X_{ij} & \text{if } i \text{ —} \underline{=} \text{— } j \\ 1 - X_{ij} & \text{if } i \text{ —} \underline{\neq} \text{— } j \end{cases}$$

Drop the constant



- i.e. $\max \sum_{i < j} X_{ij} \bar{E}_{ij}$

where

$$\bar{E}_{ij} = \begin{cases} 1 & \text{if } i \text{ —} \underline{=} \text{— } j \\ -1 & \text{if } i \text{ —} \underline{\neq} \text{— } j \end{cases}$$

Summary of SDP

$$\max \sum_{i < j} X_{ij} \bar{E}_{ij} \text{ s.t.}$$

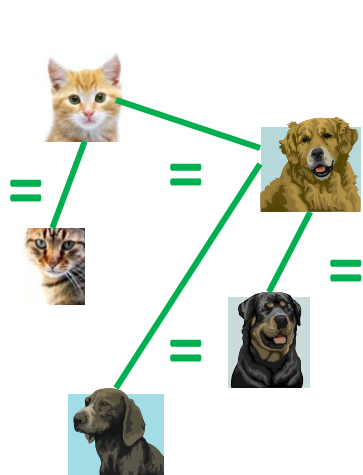
X p.s.d.


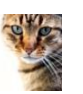




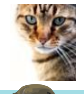



$$X_{ii} = 1$$







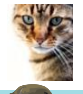



$$X_{ij} \geq 0$$

This SDP was previously used by:

- [Charikar, Guruswami, Wirth '05]
- [Swamy '04]

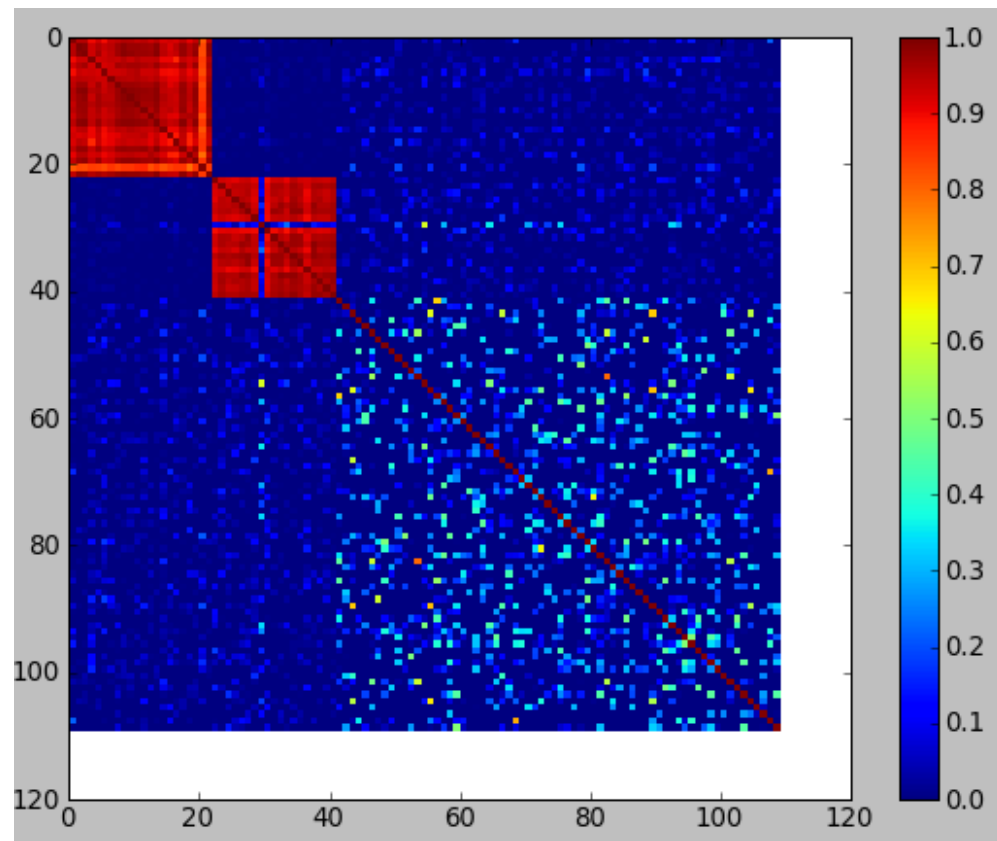


\bar{E}					
	0	1	-1	-1	1
	1	0	-1	-1	-1
	-1	-1	0	-1	1
	-1	-1	-1	0	1
	1	-1	1	1	0

X					
	1	1	0	0	0
	1	1	0	0	0
	0	0	1	0	0.7
	0	0	0	1	0.7
	0	0	0.7	0.7	1

Discussion

- Algorithm:
 - Solve SDP
 - If integral, output it. Otherwise fail.
- **Thm:** assume $p \leq 1/3$. If all clusters have size at least $\alpha_1 \sqrt{n}$ then the SDP recovers B exactly with high probability.



An example X matrix from solver in [Elsner and Schudy '09]. That solver scales to a few thousand objects.

Plan

 The semi-definite program

- Its dual
- Using the dual

Translate SDP into LP

The following are equivalent (X symmetric):

– X positive semi-definite

– $u^T X u \geq 0$ for all vectors u ←———— Linear in X for fixed u

SDP again:

$$\max \sum_{i < j} X_{ij} \bar{E}_{ij} \text{ s.t.}$$

$$X_{ii} = 1$$

$$X_{ij} \geq 0$$

X p.s.d.

LP form:

$$\max \sum_{i < j} X_{ij} \bar{E}_{ij} \text{ s.t.}$$

$$X_{ii} = 1$$

$$X_{ij} \geq 0$$

$$u^T X u \geq 0 \text{ for all vectors } u$$

SDP Dual

Primal:

$$\max \sum_{i < j} X_{ij} \bar{E}_{ij} \text{ s.t.}$$

$$X_{ii} = 1 \text{ for all } i \quad (d_i)$$

$$-X_{ij} \leq 0 \text{ for all } i, j \quad (h_{ij})$$

$$-u^T X u \leq 0 \text{ for all } u \quad (a_u)$$

Dual:

$$\min \sum_i d_i \text{ s.t.}$$

$$-\sum_u a_u u_i u_j + d_i 1(i=j) - h_{ij} = \bar{E}_{ij} \text{ for all } i \leq j$$

$$a_u, h_{ij} \geq 0$$

Translate dual LP into SDP

The following are equivalent (X symmetric):

- $X = \sum_k a_u u_k u_k^T$ with $a_u \geq 0$
- X positive semi-definite

Dual again:

$$\min \sum_i d_i \text{ s.t.}$$

$$-\sum_u a_u u_i u_j + d_i 1(i=j) - h_{ij} = \bar{E}_{ij}$$

$$a_u, h_{ij} \geq 0$$

Matrix form:

Arbitrary
positive semi-
definite matrix

$$\min \text{Trace}(D) \text{ s.t.}$$

$$-\bar{E} + D - H = \sum_u a_u u u^T$$

D diagonal

$$a_u \geq 0$$

$$H \geq 0$$

The Dual SDP

$\min \text{Trace}(D)$ s.t.

$-\bar{E} + D - H$ positive semi-definite

D diagonal

$H \geq 0$

Plan

- ✓ The semi-definite program
- ✓ Its dual
 - Using the dual

This proof is inspired by a similar result for the planted clique problem [Feige and Krauthgamer '00].

Using the dual - overview

- Prove optimality of the base clustering by presenting dual solution (D, H) whose value matches value of base clustering B (see paper)
- Difficult part: proving that $-\bar{E} + D - H$ is p.s.d.
- The following are equivalent (Y symmetric):
 - Y positive semi-definite
 - All eigenvalues of Y are ≥ 0
- We present b eigenvectors with eigenvalue 0 (see paper), where b is the number of clusters in B
- We prove that the $(b+1)^{th}$ smallest eigenvalue, denoted $\lambda_{b+1}(-\bar{E} + D - H)$, is positive (sketched next)
- Hence all eigenvalues of $-\bar{E} + D - H$ are ≥ 0

Eigenvalue analysis

$$-\bar{E} + D - H = M_1 + M_2 + M_3 + M_4$$

$$\lambda_{b+1} = \theta(\text{min cluster size})$$

(see paper)

$$\lambda_1 \geq -\theta(\sqrt{n})$$

(next)

We apply the following:

Theorem [Weyl]: If M and N are symmetric matrices then

$$\lambda_{b+1}(M + N) \geq \lambda_1(M) + \lambda_{b+1}(N)$$

Hence for sufficiently large min cluster size

$$\lambda_{b+1}(-\bar{E} + D - H) > 0.$$

Random matrices

Theorem [Füredi and Komlós '81]: Let M be a random symmetric matrix with independent entries of mean zero. Then with high probability

$$|\lambda_i(M)| = O(\sqrt{n}) \text{ for all } i.$$

Application:

$$\lambda_1(M_2) = \lambda_1\left(-\bar{E} - \mathbf{Expectation}\left[-\bar{E}\right]\right) \geq -\theta(\sqrt{n})$$

To analyze M_3 we developed a **generalization** of this theorem that handles **limited dependence** between the entries.

Recap

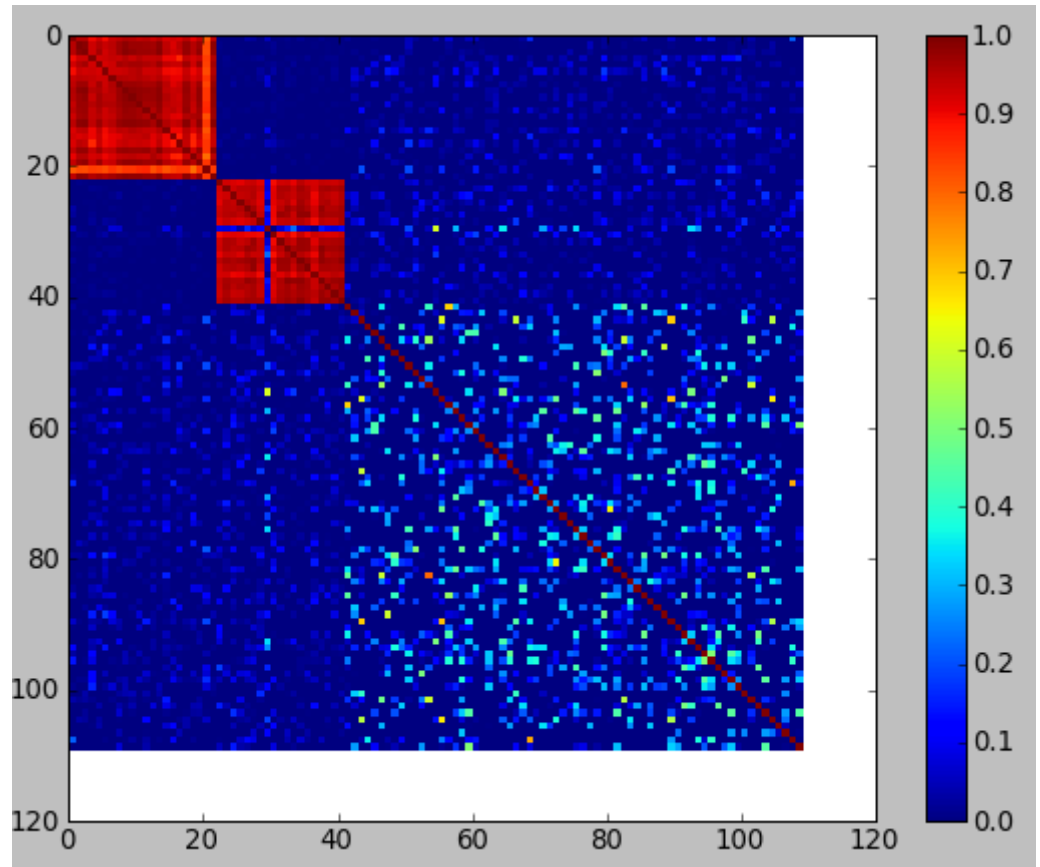
Theorem: assume $p \leq 1/3$. If all clusters have size at least $\alpha_1 \sqrt{n}$ then the SDP recovers B exactly with high probability.

Proof:

- We wrote a dual solution matrix as a sum of 4 random matrices, used Füredi-Komlós variants to bound their eigenvalues, used Weyl to infer bound on eigenvalues of the matrix, hence p.s.d., hence solution is feasible.
- That solution has value equal to the value of B , hence by duality B is primal optimal
- B is the *unique* primal optimum (see paper), hence SDP will exactly return B
- Hence algorithm reconstructs B exactly when all clusters have size at least $\alpha_1 \sqrt{n}$.

Open Question 1

- Suppose some clusters are size $c_3\sqrt{n}$ and others are size 1. Can the SDP be used to reconstruct the large clusters?



Software: [Elsner and Schudy '09].

Open Question 2

- Planted clique problem = correlation clustering with only one non-singleton and no corruption of within-cluster edges
- Exist polynomial-time algorithm when clique size = $c_1 \sqrt{n}$
- Exists $n^{O(\log n)}$ -time algorithm when clique size = $c_1 \log n$
- Can polynomial-time algorithms beat the $c_1 \sqrt{n}$ barrier?

Clustering References

- Nir Ailon, Moses Charikar, and Alantha Newman. Aggregating inconsistent information: ranking and clustering. In *STOC '05*, pages 684–693, 2005.
- Nikhil Bansal, Avrim Blum, and Shuchi Chawla. Correlation clustering. *Mach. Learn.*, 56(1-3):89–113, 2004.
- Moses Charikar, Venkatesan Guruswami, and Anthony Wirth. Clustering with qualitative information. *J. Comput. Syst. Sci.*, 71(3):360–383, 2005.
- M. Elsner and W. Schudy. Bounding and Comparing Methods for Correlation Clustering Beyond ILP. In *ILP-NLP '09: Proc. NAACL/HLT 2009 Workshop on Integer Linear Programming for Natural Language Processing*, pages 19–27, 2009.
- Ron Shamir and Dekel Tsur. Improved algorithms for the random cluster graph model. *Random Structures and Algorithms*, 31(4):418–449, 2007.

Other References

- F. Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM Journal on Optimization*, 5(1):13–51, 1995
- Uriel Feige and Robert Krauthgamer. Finding and certifying a large hidden clique in a semirandom graph. *Random Struct. Algorithms*, 16(2):195–208, 2000
- Zoltán Füredi and János Komlós. The eigenvalues of random symmetric matrices. *Combinatorica*, 1(3):233–241, 1981