

Non-convex relaxations for the densest submatrix problem

Brendan Ames

Department of Mathematics
The University of Alabama

UA Applied Math Seminar

Friday August 26, 2022

Agenda

Consider convex and non-convex relaxations for the **maximum clique** and **densest submatrix** problems.

Give a probabilistic model for **“clusterable”** data and graphs, and theoretical recovery guarantees.

Propose efficient first-order methods for solving these relaxations.

Joint work with **Polina Bombina, UA**.

Clustering

Clustering: partition data so that items in each cluster are similar to each other and items not in the same cluster are dissimilar.

Fundamental problem in statistics and machine learning:

- pattern recognition, computational biology, image processing/computer vision, network analysis.

No consensus on what constitutes a **good** clustering; depends heavily on application.

Intractable: usually modeled as some NP-hard problem (e.g., clique, normalized cut, k-means).

A sanity check

Clustering seems to be a very difficult/ill-posed problem.

Many heuristics seem to work well in practice.

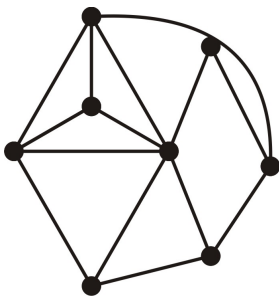
Question: can we show that we can cluster “clusterable” data?
How do we model **clusterable** data?

Cliques of a graph

Given graph $G = (V, E)$, a **clique** of G is a pairwise adjacent subset of V .

The vertex set $C \subseteq V$ is a clique of G if $uv \in E$ for all $u, v \in C$.

The subgraph $G(C)$ induced by C is **complete**.

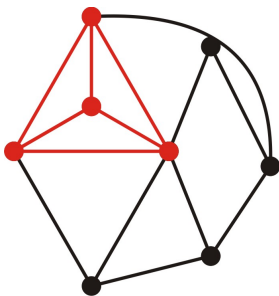


Cliques of a graph

Given graph $G = (V, E)$, a **clique** of G is a pairwise adjacent subset of V .

The vertex set $C \subseteq V$ is a clique of G if $uv \in E$ for all $u, v \in C$.

The subgraph $G(C)$ induced by C is **complete**.



The Clique problem

Optimization version: Find the clique of G of maximum size. Size of the largest clique is the **clique number** $\omega(G)$.

Decision version: Given graph G , integer k : does G contain a clique of cardinality at least k .

Complexity: NP-complete, cannot approximate within a ratio of $N^{1-\epsilon}$ for any $\epsilon > 0$.

Many applications: communication, biological, and social networks. Find large group of related objects.

The planted case

Hardness results are **worst** case.

There should be instances we should be able to solve efficiently.

In particular, if G has a clique of size k , we should be able to find it if k is large.

Alon et al. 1998, Feige and Krauthgamer 2000, Ames and Vavasis 2011: if $k \geq \Omega(\sqrt{N})$ and all other edges are added independently at random then we can find the maximum clique in polynomial time.

A more general model?

These recovery guarantees rely heavily on the fact that G is an undirected graph:

- e.g., symmetry of A_G , the fact that a stable set of \bar{G} is a clique of G , etc.

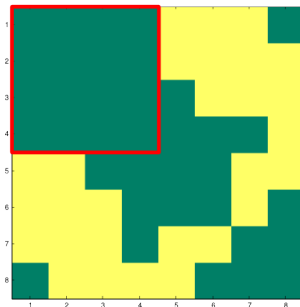
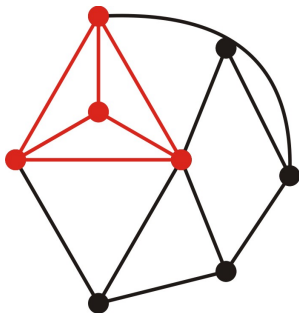
Would like an approach that translates to finding other “clique-like” objects with minimal effort.

e.g., the maximum biclique of a bipartite graph, fully dense block in a matrix.

Cliques and low-rank matrices

Every clique C (with characteristic vector \mathbf{v}) of the graph $G = (V, E)$ defines a rank-one matrix by $\mathbf{X} = \mathbf{v}\mathbf{v}^T$.

Moreover, nonzero entries of \mathbf{X} form a $|C| \times |C|$ rank-one block in $\mathbf{A}_G + \mathbf{I}$.



Clique as rank minimization

G has a clique of cardinality at least k if and only if there exists rank-one symmetric binary matrix \mathbf{X} such that

$$\begin{aligned}\sum \sum x_{ij} &\geq k^2 \\ x_{ij} &= 0 \quad \forall ij \notin E, i \neq j.\end{aligned}$$

Otherwise $\omega(G) < k$.

Therefore **Clique** is equivalent to the rank minimization problem:

$$\min_{\substack{\mathbf{X} \in \{0,1\}^{V \times V} \\ \mathbf{X} \in \Sigma^V}} \left\{ \text{rank}(\mathbf{X}) : \mathbf{e}^T \mathbf{X} \mathbf{e} \geq k^2, x_{ij} = 0 \text{ if } (i,j) \in \tilde{E} \right\}$$

where $\tilde{E} = V \times V - \{E \cup \{(u, u) : u \in V\}\}$.

Rank minimization

Affine rank minimization problem: find matrix with minimum rank satisfying linear constraints:

$$\min\{\text{rank}(\mathbf{X}) : \mathcal{A}(\mathbf{X}) = \mathbf{b}\}.$$

Well-known to be NP-hard.

Relax $\text{rank}(\mathbf{X})$ with nuclear norm $\|\mathbf{X}\|_* = \sigma_1(\mathbf{X}) + \cdots + \sigma_N(\mathbf{X})$:

$$\text{rank}(\mathbf{X}) = \text{card } \sigma(\mathbf{X}), \quad \|\mathbf{X}\|_* = \|\sigma(\mathbf{X})\|_1.$$

If \mathcal{A} satisfies certain “niceness” conditions then the minimum nuclear norm solution is the minimum rank solution.

The densest (m,n) -submatrix problem

We want to find a **dense $k \times k$ submatrix** in $A_G + I$, not necessarily a clique.

Densest $m \times n$ -submatrix problem (DSM): Given a matrix $A \in \mathbf{R}^{M \times N}$, find submatrix with m rows and n columns with maximum number of nonzero entries.

NP-hard: proof is by reduction to **Clique**; hard to approximate.

Duality of density and number of missing edges / zero entries

Let U and V be subsets of $\{1, 2, \dots, M\}$ and $\{1, 2, \dots, N\}$ with characteristic vectors \mathbf{u} and \mathbf{v} respectively.

Introduce a new variable \mathbf{Y} to act as a **correction** for entries of $\mathbf{X} = \mathbf{u}\mathbf{v}^T$ that should be 0:

$$y_{ij} = \begin{cases} -x_{ij}, & \text{if } a_{ij} = 0 \\ 0, & \text{otherwise.} \end{cases}$$

Cardinality of \mathbf{Y} acts as a dual of density of $\mathbf{A}(U, V)$:

$$\text{card}(\mathbf{A}(U, V)) = mn - \sum_{i=1}^M \sum_{j=1}^N y_{ij}$$

Formulation as sparse plus low-rank decomposition

Can formulate **(DSM)** as

$$\begin{aligned} \min \quad & \text{rank } \mathbf{X} + \gamma \text{card } \mathbf{Y} \\ \text{s. t.} \quad & \mathbf{e}^T \mathbf{X} \mathbf{e} = mn \\ & x_{ij} + y_{ij} = 0 \text{ if } a_{ij} = 0 \\ & x_{ij} \in \{0, 1\} \end{aligned}$$

where γ is a regularization parameter.

Formulation as sparse plus low-rank decomposition

Can formulate (**DSM**) as

$$\begin{aligned} \min \quad & \|X\|_* + \gamma \|Y\|_1 \\ \text{s. t.} \quad & \mathbf{e}^T X \mathbf{e} = mn \\ & x_{ij} + y_{ij} = 0 \text{ if } a_{ij} = 0 \\ & 0 \leq x_{ij} \leq 1 \end{aligned}$$

where γ is a regularization parameter.

Relax **card** \mathbf{Y} using the ℓ_1 -norm $\|\mathbf{Y}\|_1$, and **rank** \mathbf{X} with the nuclear norm $\|\mathbf{X}\|_*$.

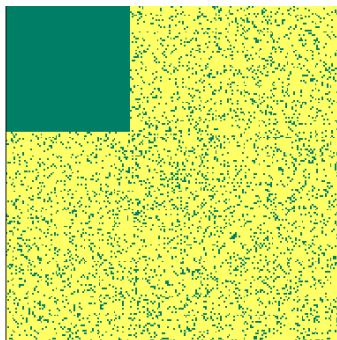
Planted case

Start with $M \times N$ all-zeros matrix \mathbf{A} .

Set all entries in $m \times n$ block equal to 1.

Add noise:

- Add some of the remaining potential entries with probability p .
- Delete some entries in $m \times n$ block with probability $1 - q$, $q > p$.



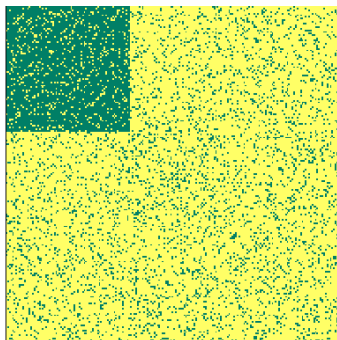
Planted case

Start with $M \times N$ all-zeros matrix \mathbf{A} .

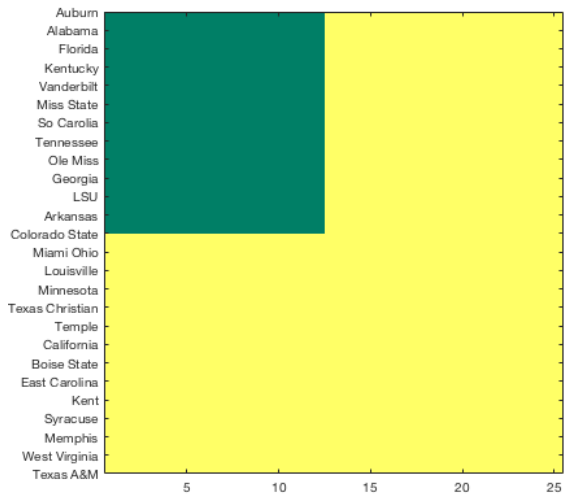
Set all entries in $m \times n$ block equal to 1.

Add noise:

- Add some of the remaining potential entries with probability p .
- Delete some entries in $m \times n$ block with probability $1 - q$, $q > p$.



Back to the SEC Example



Recovery Guarantee

Theorem (Bombina-Ames 2020)

Suppose that \mathbf{A} is sampled from the planted dense $m \times n$ -submatrix model with edge probabilities q and p .

Let $(\mathbf{X}^*, \mathbf{Y}^*)$ denote the matrix representation of the planted submatrix and assume $m \leq n$, $M \leq N$.

Then there exists constants $c_1, c_2, c_3 > 0$ such that if

$$q - p \geq c_1 \max \left\{ \sqrt{\max\{\sigma_q^2, \sigma_p^2\} \frac{\log N}{m}}, \frac{\log N}{m} \sqrt{\sigma_p^2 N}, \frac{(\log N)^{3/2}}{m} \right\}$$

then $(\mathbf{X}^*, \mathbf{Y}^*)$ is the **unique optimal solution** of (DSM) for regularization parameter

$$\gamma = \frac{t}{(q - p)m}, \quad c_2 \leq t \leq c_3$$

with high probability.

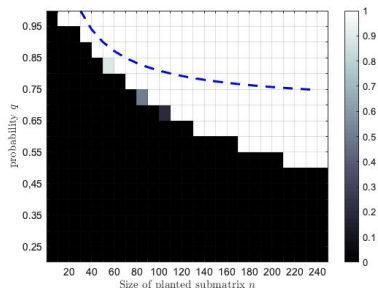
Example: Dense Case

Suppose that p, q are fixed or shrink very slowly, i.e.,
 $p, 1 - q > 1/\log k$.

Then we can recover the planted submatrix with high probability provided that

$$m \geq C\sqrt{N \log N}.$$

Ignoring log-term, we have the same results as before.



Sparse Graphs

In most practical examples, the following are not necessarily true:

- 1 $m = \Omega(\sqrt{N})$.
- 2 The noise probabilities p, q are not fixed.

Example: Community Detection. In most real-world social networks, community size does not grow as the number of users increases. (Seems to be capped at a very small fraction of the total population.)

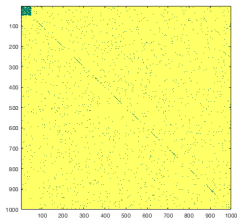
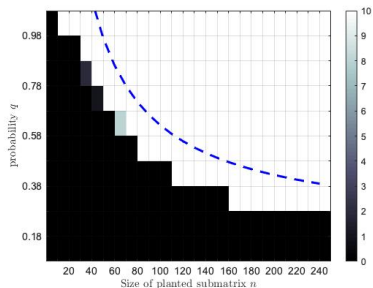
Need to modify model to use **sparse** noise: p and/or q tend to zero as $N \rightarrow \infty$.

Example: Sparse Case

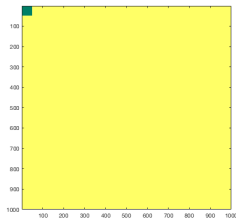
Suppose that noise is **sparse**.

Suppose q is fixed and
 $p \leq \log N/N$.

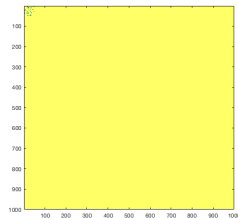
Then we have exact recovery
w.h.p. if $m \geq C(\log N)^{3/2}$



\Rightarrow



$+$



Proof Idea

Apply KKT conditions and SDP duality to derive conditions ensuring optimality and uniqueness of \mathbf{X}^* .

Propose a choice of Lagrange multipliers corresponding to \mathbf{X}^* .

Use bounds on concentration of norms of random matrices to establish that these multipliers satisfy the optimality and uniqueness conditions (with high probability).

ADMM Approach

Introduce artificial variables \mathbf{Q} , \mathbf{W} , \mathbf{Z} to obtain the equivalent convex optimization problem

$$\begin{aligned} \min \quad & \|\mathbf{X}\|_* + \gamma \|\mathbf{Y}\|_1 + \mathbf{1}_{\Omega_Q}(\mathbf{Q}) + \mathbf{1}_{\Omega_W}(\mathbf{W}) + \mathbf{1}_{\Omega_Z}(\mathbf{Z}) \\ & \mathbf{X} = \mathbf{Y} = \mathbf{Q}, \mathbf{X} - \mathbf{W} = \mathbf{0}, \mathbf{X} - \mathbf{Z} = \mathbf{0}, \end{aligned}$$

where $\Omega_Q, \Omega_W, \Omega_Z$ denote the constraint sets

$$\begin{aligned} \Omega_Q &:= \{\mathbf{Q} : P_{\tilde{N}}(\mathbf{Q}) = \mathbf{0}\}, \\ \Omega_W &:= \{\mathbf{W} : \mathbf{e}^T \mathbf{W} \mathbf{e} = mn\}, \\ \Omega_Z &= \{\mathbf{Z} : Z_{ij} \leq 1 \forall (i,j) \in M \times N\}, \end{aligned}$$

and $\mathbf{1}_S : \mathbf{R}^{M \times M} \rightarrow \{0, +\infty\}$ is the indicator function of the set $S \subseteq \mathbf{R}^{M \times N}$ ($\mathbf{1}_S(\mathbf{X}) = 0$ if $\mathbf{X} \in S$, and $+\infty$ otherwise).

ADMM Idea

We solve using the **Alternating Direction Method of Multipliers (ADMM)**.

We update each primal variable by minimizing the augmented Lagrangian in Gauss-Seidel fashion with respect to each primal variable. Then the dual variables are updated using approximate gradient ascent.

ADMM Update Steps

The augmented Lagrangian is given by

$$\begin{aligned} L_\tau = & \| \mathbf{X} \|_* + \gamma \| \mathbf{Y} \|_1 + \mathbf{1}_{\Omega_Q}(\mathbf{Q}) + \mathbf{1}_{\Omega_W}(\mathbf{W}) + \mathbf{1}_{\Omega_Z}(\mathbf{Z}) \\ & + \text{tr}(\boldsymbol{\Lambda}_Q(\mathbf{X} - \mathbf{Y} - \mathbf{Q})) + \text{tr}(\boldsymbol{\Lambda}_W(\mathbf{X} - \mathbf{W})) + \text{tr}(\boldsymbol{\Lambda}_Z(\mathbf{X} - \mathbf{Z})) \\ & + \frac{\tau}{2} (\| \mathbf{X} - \mathbf{Y} - \mathbf{Q} \|_F^2 + \| \mathbf{X} - \mathbf{W} \|_F^2 + \| \mathbf{X} - \mathbf{Z} \|_F^2), \end{aligned}$$

where τ is a regularization parameter chosen so that L_τ is strongly convex in each primal variable.

Update \mathbf{Q} , \mathbf{W} and \mathbf{Z} by projection onto each of the sets Ω_Q, Ω_W and Ω_Z .

Update \mathbf{X} and \mathbf{Y} using proximal operators of $\| \cdot \|_*$ and $\| \cdot \|_1$ respectively.

The Algorithm

```
while convergence==0 % Repeat until converged.
    % Update Q. Project onto support of A.
    Q = (X - Y + mu*LambdaQ).*A;

    % Update X by singular value shrinkage.
    X = mat_shrink(1/3*(Y + Q + Z + W
        - mu*(LambdaQ + LambdaW + LambdaZ)), 1/(3*tau));

    % Update Y as projection of residual onto nonnegative cone.
    Y = max(X-Q-gamma*ones(M,N)*mu + LambdaQ*mu, zeros(M,N));

    % Scale/shift W so that entries sum to m*n.
    newW = X + mu*LambdaW;
    alfa = (m*n-sum(newW(:)))/(M*N);
    W = newW + alfa*ones(M,N);

    % Update Z.
    Z = X+ mu*LambdaZ; Z = min(max(Z,0),1);

    % Update dual variables by approximate gradient ascent.
    LambdaQ = LambdaQ + tau*(X-Y-Q);
    LambdaW = LambdaW + tau*(X-W);
    LambdaZ = LambdaZ + tau*(X-Z);
end
```

A Problem

ADMM algorithm requires $O(N^3)$ floating point operations for **singular value decomposition** each iteration; algorithm converges linearly.

Cannot solve large-scale problem instances.

Limited to graphs/matrices with $N = O(1000)$.

Quadratic Programming Relaxation

If $\text{rank } \mathbf{X} = 1$ then $\mathbf{X} = \mathbf{u}\mathbf{v}^T \in \mathbf{R}^{M \times N}$ for some $\mathbf{u} \in \mathbf{R}^M$, $\mathbf{v} \in \mathbf{R}^N$.

(DSM) can be relaxed as

$$\begin{aligned} \min \quad & \frac{\lambda}{2} \left(\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 \right) + \mathbf{u}^T \bar{\mathbf{A}} \mathbf{v} \\ \text{s. t.} \quad & \sum u_i = m, \quad \sum v_i = n \\ & 0 \leq u_i \leq 1, \quad 0 \leq v_i \leq 1 \end{aligned}$$

This is a **non-convex** quadratic program in \mathbf{u} and \mathbf{v} .

A Translation of Recovery Guarantees

Theorem

Suppose that the **nuclear norm relaxation is exact**.

That is $\mathbf{X}^* = \mathbf{u}^*(\mathbf{v}^*)^T$, is the optimal solution for **(DSM)** and the nuclear norm relaxation with regularization parameter γ .

Then $(\mathbf{u}^*, \mathbf{v}^*)$ is the optimal solution of the non-convex QP relaxation with

$$\lambda \leq \frac{1}{2\gamma} \min \left\{ \sqrt{\frac{m}{n}}, \sqrt{\frac{n}{m}} \right\}.$$

Proof Idea: Use optimality of \mathbf{X}^* to establish that

$$\frac{\lambda}{2} \left(\|\mathbf{u}\|_2^2 + \|\mathbf{v}\|_2^2 \right) + \mathbf{u}^T \bar{\mathbf{A}} \mathbf{v} \geq \frac{\lambda}{2} \left(\|\mathbf{u}^*\|_2^2 + \|\mathbf{v}^*\|_2^2 \right) + (\mathbf{u}^*)^T \bar{\mathbf{A}} \mathbf{v}^*$$

for every feasible \mathbf{u} and \mathbf{v} for this choice of γ and λ .

LADMM setup

We can write the QP relaxation as

$$\begin{aligned} \min \quad & \frac{\lambda}{2} (\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2) + \mathbf{u}^T \bar{\mathbf{A}} \mathbf{v} + \mathbf{1}_{\Omega_1}(\mathbf{x}) + \mathbf{1}_{\Omega_2}(\mathbf{w}) \\ \text{s. t.} \quad & \mathbf{u} = \mathbf{x}, \mathbf{v} = \mathbf{w}, \end{aligned}$$

where

$$\begin{aligned} \Omega_1 &= \{\mathbf{x} : \mathbf{0} \leq \mathbf{x} \leq \mathbf{e}, \mathbf{x}^T \mathbf{e} = m\}, \\ \Omega_2 &= \{\mathbf{w} : \mathbf{0} \leq \mathbf{w} \leq \mathbf{e}, \mathbf{w}^T \mathbf{e} = n\}. \end{aligned}$$

The augmented Lagrangian is given by:

$$\begin{aligned} L_\tau &= \frac{\lambda}{2} (\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2) + \mathbf{u}^T \bar{\mathbf{A}} \mathbf{v} + \mathbf{1}_{\Omega_1}(\mathbf{x}) + \mathbf{1}_{\Omega_2}(\mathbf{w}) \\ &\quad + \boldsymbol{\Lambda}_1^T (\mathbf{u} - \mathbf{x}) + \boldsymbol{\Lambda}_2^T (\mathbf{v} - \mathbf{w}) + \frac{\tau}{2} (\|\mathbf{u} - \mathbf{x}\|^2 + \|\mathbf{v} - \mathbf{w}\|^2) \end{aligned}$$

Outline of the Algorithm

Minimization of the augmented Lagrangian with respect to each of the artificial primal variables \mathbf{x} and \mathbf{w} is equivalent to projection onto the capped simplex.

To update \mathbf{u} , we replace $\mathbf{u}^T \bar{\mathbf{A}} \mathbf{v}^i + \frac{\lambda}{2} \|\mathbf{u}\|^2$ by

$$\langle \mathbf{u} - \mathbf{u}^i, \bar{\mathbf{A}} \mathbf{v}^i + \lambda \mathbf{u}^i \rangle + \frac{\ell_u}{2} \|\mathbf{u} - \mathbf{u}^i\|^2,$$

where ℓ_u is a regularization term.

Similarly for \mathbf{v} : we replace $\mathbf{u}^T \bar{\mathbf{A}} \mathbf{v} + \frac{\lambda}{2} \|\mathbf{v}\|^2$ by

$$\langle \mathbf{v} - \mathbf{v}^i, \bar{\mathbf{A}}^T \mathbf{u}^{i+1} + \lambda \mathbf{v}^i \rangle + \frac{\ell_v}{2} \|\mathbf{v} - \mathbf{v}^i\|^2,$$

where ℓ_v is a regularization term.

The LADMM Algorithm

```
while convergence==0
    %update x
    y0 = u + 1/tau*Lambda_x;
    x = projection(y0,m,tau);

    % Update u
    u = 1/(L_v+tau)*(tau*x-Lambda_x-A_bar*v+L_v*u_old-lambda*u_old);

    %update w
    y1 = v + 1/tau*Lambda_w;
    w = projection(y1,n,tau);

    % Update v
    v = 1/(L_v+tau)*(tau*w-Lambda_w-A_bar'*u+L_v*v_old-lambda*v_old);

    % Update dual variables
    Lambda_x_old = Lambda_x;
    Lambda_x = Lambda_x_old+tau*(u-k);

    Lambda_w_old = Lambda_w;
    Lambda_w = Lambda_w_old + tau*(v-w);
end
```

Remarks

The sequences of iterates $\{\mathbf{u}^k\}, \{\mathbf{v}^k\}, \{\mathbf{x}^k\}, \{\mathbf{w}^k\}$ are convergent if we choose regularization parameter τ and linearization parameters ℓ_u, ℓ_v in a certain range.

The QP relaxation is **degenerate** (i.e., doesn't satisfy usual constraint qualifications) at binary feasible solutions.

Can show that there is a non-zero duality gap between the QP relaxation and its dual for modestly large planted solutions, even when we have perfect recovery.

In practice, method converges quickly with initial solution $\mathbf{u}^0 = \mathbf{e}/m \in \mathbf{R}^M$ and $\mathbf{v}^0 = \mathbf{e}/n \in \mathbf{R}^N$.

Improvement: Adaptive LADMM

Performance depends on augmented Lagrangian parameter τ .

Number of iterations and run-time increase significantly if τ is too small or too large.

Need to automate choice of τ :

- 1 **Residual balancing:** increment/decrement τ^i to tune between primal and dual residuals.
- 2 **Line-search** to choose τ^i ensuring sufficient decrease in residual each iteration.

Empirical Trials

We randomly generate 500×500 matrices with randomly generated planted densest $m \times n$ submatrices according to the planted submatrix model with

$$\begin{aligned}n &\in \{10, 20, 30, \dots, 250\} & m &= 2n \\p &= 0.25 & q &\in \{0.3, 0.4, 0.5, \dots, 1\}.\end{aligned}$$

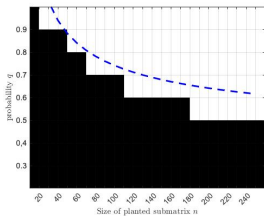
We use **ADMM**, **LADMM**, and adaptive ADMM with **line search (AdaLADMM-LS)** and **residual balancing (AdaLADMM-RB)** with $\gamma = 6/(q - p)n$ and $\lambda = (q - p)n/10$.

Augmented Lagrangian parameters and adaptation parameters are chosen to ensure convergence.

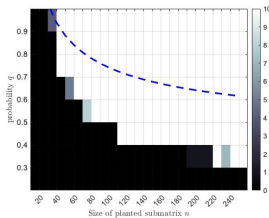
Stop each algorithm with stopping tolerance $\epsilon = 10^{-4}$ and maximum number of iterations **2000**.

Recovery Rates for Randomly Generated Matrices

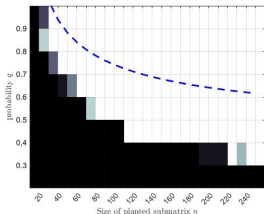
Declare DSM recovered if relative error between planted solution and calculated solution is within 10^{-2} . Repeat 10 times.



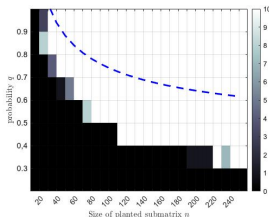
ADMM



LADMM

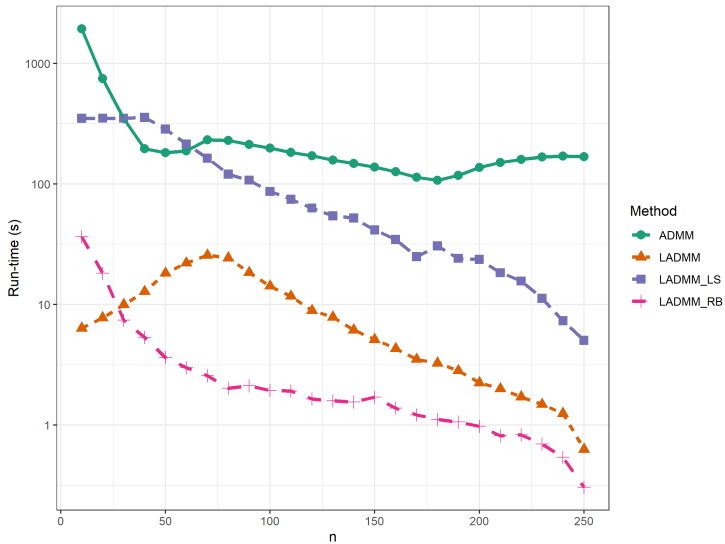


AdaLADMM:RB



AdaLADMM:LS

Run Time



Thank you!

P. Bombina and B. Ames. *Convex optimization for the densest subgraph and densest submatrix problems*. SN Operations Research Forum. Year: 2020, Vol: 1, No: 3.
<https://link.springer.com/article/10.1007/s43069-020-00020-5>

Software available from bpames.people.ua.edu/software

B. Ames supported by **NSF Grants #2012554** and **#2108645**; **UA Cyberseed Grant SP14572**; University of Alabama RGC grants **RG14678** and **RG14838**.

P. Bombina supported by **Alabama EPSCoR Graduate Research Scholars Program**.