



A New Perspective on Low-Rank Optimization

Ryan Cory-Wright

MIT, Operations Research Center

Paper available at: ryancorywright.github.io

Questions? Email ryancw@mit.edu

Joint work with

Dimitris Bertsimas (MIT)

Jean Pauphilet (LBS)

Motivating Example: A Tale of Two Problems

Sparse Linear Regression:

$$\min_{\mathbf{w} \in \mathbb{R}^p} \quad \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2\gamma} \|\mathbf{w}\|_2^2 + \mu \|\mathbf{w}\|_0$$

Routinely solved to optimality at scale, because convex relaxation well-studied.

Reduced Rank Regression:

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{p \times n}} \quad \frac{1}{2m} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_F^2 + \frac{1}{2\gamma} \|\boldsymbol{\beta}\|_F^2 + \mu \cdot \text{Rank}(\boldsymbol{\beta})$$

Almost the same problem, but intractable, because no known relaxations give valid bounds.

We show: How to get strong & scalable relaxations.

Motivating Example: Linear Regression and Relaxations

Sparse Linear Regression: Fit **interpretable** model using small number of features :

$$\min_{\mathbf{w} \in \mathbb{R}^p} \quad \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2\gamma} \|\mathbf{w}\|_2^2 + \mu \|\mathbf{w}\|_0$$

Perspective Relaxation:

$$\min_{\mathbf{w}, \boldsymbol{\rho} \in \mathbb{R}^p, \mathbf{z} \in \{0,1\}^p} \quad \frac{1}{2n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{1}{2\gamma} \mathbf{e}^\top \boldsymbol{\rho} + \mu \cdot \mathbf{e}^\top \mathbf{z} \quad \text{s.t.} \quad z_i \rho_i \geq w_i^2 \quad \forall i \in [p].$$

Relaxation due to Frangioni and Gentile 2006, Günlük and Linderoth 2010: **strong & scalable**.

Allows exact solutions with $p = 10^7$ features (Bertsimas and van Parys 2020, Hazimeh and Mazumder 2021)

Motivating Example: Rank Regression and Relaxations

Reduced Rank Regression: Fit **interpretable** model using small number of singular values

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{p \times n}} \quad \frac{1}{2m} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_F^2 + \frac{1}{2\gamma} \|\boldsymbol{\beta}\|_F^2 + \mu \cdot \text{Rank}(\boldsymbol{\beta})$$

Matrix Perspective Relaxation (new):

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{p \times n}, \mathbf{W} \in \mathcal{S}_+^n, \boldsymbol{\theta} \in \mathcal{S}_+^p} \quad \frac{1}{2m} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_F^2 + \frac{1}{2\gamma} \text{tr}(\boldsymbol{\theta}) + \mu \cdot \text{tr}(\mathbf{W}) \quad \text{s.t.} \quad \mathbf{W} \preceq \mathbb{I}, \begin{pmatrix} \boldsymbol{\theta} & \boldsymbol{\beta} \\ \boldsymbol{\beta}^\top & \mathbf{W} \end{pmatrix} \succeq \mathbf{0}.$$

In this talk, we derive this relaxation.

A generalization of perspective reformulation technique due to Frangioni and Gentile '06.

Even Stronger Relaxations

Saddle-Point Sparsity Relaxation (Dong et al. '15): Also apply to diagonal “extracted” from $X^T X$

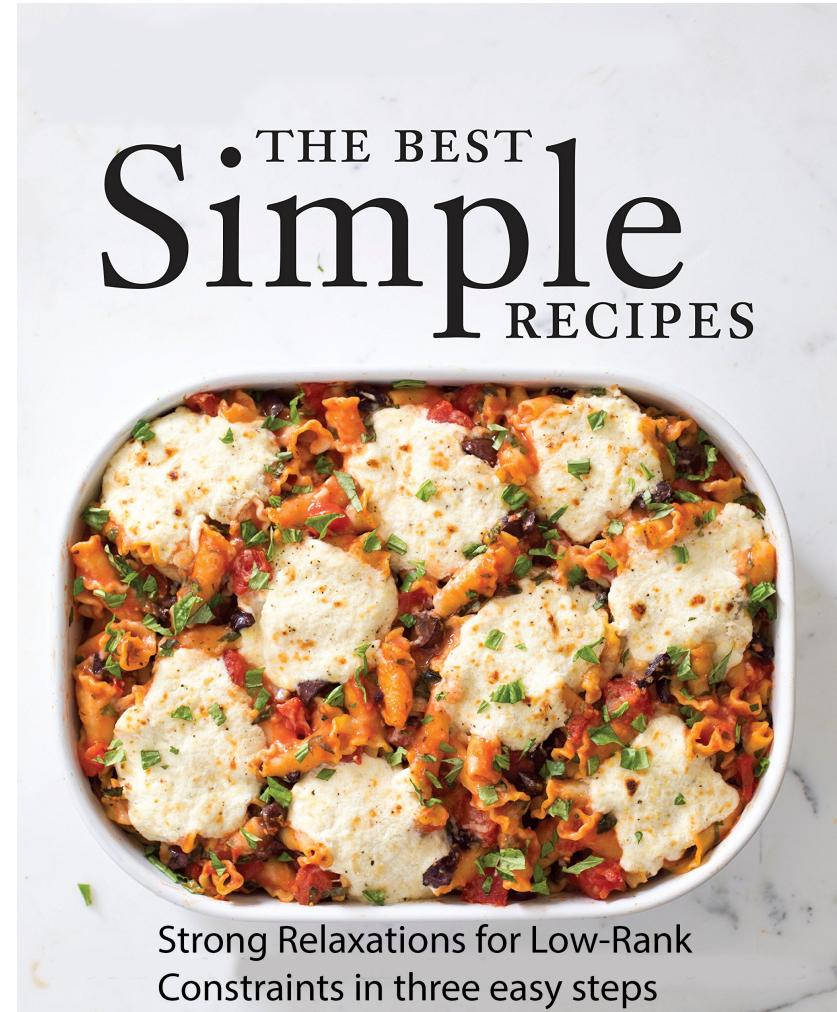
$$\begin{aligned} \min_{\boldsymbol{w} \in \mathbb{R}^p, \boldsymbol{z} \in [0,1]^p, \boldsymbol{W} \in S_+^p} \quad & \frac{1}{2n} \|\boldsymbol{y}\|_2^2 - \frac{1}{n} \langle \boldsymbol{y}, \boldsymbol{X} \boldsymbol{w} \rangle + \frac{1}{2} \langle \boldsymbol{W}, \frac{1}{\gamma} \mathbb{I} + \frac{1}{n} \boldsymbol{X}^\top \boldsymbol{X} \rangle + \mu \boldsymbol{e}^\top \boldsymbol{z} \\ \text{s.t.} \quad & \boldsymbol{W} \succeq \boldsymbol{w} \boldsymbol{w}^\top, \ z_i W_{i,i} \geq w_i^2 \ \forall i \in [p]. \end{aligned}$$

Saddle-Point Rank Relaxation (new): Play same game in low-rank case

$$\begin{aligned} \min_{\boldsymbol{\theta} \in \mathcal{S}_+^n, \boldsymbol{\beta} \in \mathbb{R}^{p \times n}, \boldsymbol{B} \in \mathcal{S}_+^n, \boldsymbol{W} \in \mathcal{S}_+^n} \quad & \frac{1}{2m} \|\boldsymbol{Y}\|_F^2 - \frac{1}{m} \langle \boldsymbol{Y}, \boldsymbol{X} \boldsymbol{\beta} \rangle + \frac{1}{2} \langle \boldsymbol{B}, \frac{1}{\gamma} \mathbb{I} + \frac{1}{m} \boldsymbol{X}^\top \boldsymbol{X} \rangle + \mu \cdot \text{tr}(\boldsymbol{W}) \\ \text{s.t.} \quad & \begin{pmatrix} \boldsymbol{B} & \boldsymbol{\beta} \\ \boldsymbol{\beta} & \boldsymbol{W} \end{pmatrix} \succeq \mathbf{0}, \boldsymbol{W} \preceq \mathbb{I}. \end{aligned}$$

Matrix Perspective Reformulation Technique: Recipe

1. Consider low-rank problem with spectral regularization
2. Formulate as mixed-projection optimization problem
3. Take matrix perspective of spectral regularizer.



Contributions

- **Methodological:** We propose a **simple preprocessing technique** which gives **strong & scalable** bounds for low-rank problems. Generalizes perspective reformulation technique from MIO.
- **Theoretical:** We invoke technique to **explicitly characterize** convex hulls of simple low-rank sets
- **Algorithmic:** We demonstrate technique's **efficacy** across diverse range of low-rank problems

Matrix Perspective Reformulation Technique I

Consider low-rank problem with spectral regularization

$$\min_{\mathbf{X} \in \mathcal{S}_+^n} \langle \mathbf{C}, \mathbf{X} \rangle + \Omega(\mathbf{X}) + \boxed{\mu \cdot \text{Rank}(\mathbf{X})} \text{ s.t. } \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} \in \mathcal{K}, \quad \boxed{\text{Rank}(\mathbf{X}) \leq k},$$

Where:

- Rank \rightarrow \mathbf{x} of minimal complexity,

Matrix Perspective Reformulation Technique I

Consider low-rank problem with spectral regularization

$$\min_{\mathbf{X} \in \mathcal{S}_+^n} \langle \mathbf{C}, \mathbf{X} \rangle + \boxed{\Omega(\mathbf{X})} + \mu \cdot \text{Rank}(\mathbf{X}) \text{ s.t. } \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} \in \mathcal{K}, \quad \text{Rank}(\mathbf{X}) \leq k,$$

Where:

- $\Omega(\mathbf{X}) := \sum_{i=1}^n \omega(\lambda_i(\mathbf{X}))$ convex spectral function with $0 \in \omega$; ω closed convex.
- Example: $\Omega(\mathbf{X}) = \|\mathbf{X}\|_F^2$; like ridge regularization in regression.
- Example: $\Omega(\mathbf{X}) = 0$ if $\|\mathbf{X}\|_* \leq M$, $+\infty$ otherwise; like big-M penalty in regression.

Matrix Perspective Reformulation Technique II: Formulation

Low-rank problem

$$\min_{\mathbf{X} \in \mathcal{S}_+^n} \langle \mathbf{C}, \mathbf{X} \rangle + \Omega(\mathbf{X}) + \mu \cdot \text{Rank}(\mathbf{X}) \text{ s.t. } \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} \in \mathcal{K}, \quad \text{Rank}(\mathbf{X}) \leq k,$$

can be expressed as [Mixed-Projection Optimization](#) problem

$$\begin{aligned} & \min_{\mathbf{Y} \in \mathcal{Y}_n^k} \min_{\mathbf{X} \in \mathcal{S}_+^n} \quad \langle \mathbf{C}, \mathbf{X} \rangle + \mu \cdot \text{tr}(\mathbf{Y}) + \text{tr}(f(\mathbf{X})) \\ & \text{s.t.} \quad \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} = \mathbf{Y}\mathbf{X}, \quad \mathbf{X} \in \mathcal{K} \end{aligned}$$

where \mathbf{Y} is a projection matrix, f is a spectral function such that $\text{tr}(f(\mathbf{X})) = \Omega(\mathbf{X})$.

Matrix Perspective Reformulation Technique III: Reformulation

Mixed-Projection Conic Optimization problem

$$\begin{aligned} \min_{\mathbf{Y} \in \mathcal{Y}_n^k} \min_{\mathbf{X} \in \mathcal{S}_+^n} \quad & \langle \mathbf{C}, \mathbf{X} \rangle + \mu \cdot \text{tr}(\mathbf{Y}) + \text{tr}(f(\mathbf{X})) \\ \text{s.t.} \quad & \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} = \mathbf{Y}\mathbf{X}, \quad \mathbf{X} \in \mathcal{K} \end{aligned}$$

Rewrite as equivalent problem which gives stronger relaxations:

$$\begin{aligned} \min_{\mathbf{Y} \in \mathcal{Y}_n^k} \min_{\mathbf{X} \in \mathcal{S}_+^n} \quad & \langle \mathbf{C}, \mathbf{X} \rangle + \mu \cdot \text{tr}(\mathbf{Y}) + \text{tr}(g_f(\mathbf{X}, \mathbf{Y})) + (n - \text{tr}(\mathbf{Y}))\omega(0) \\ \text{s.t.} \quad & \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} \in \mathcal{K}, \end{aligned}$$

where g_f is the matrix perspective of f (Effros, PNAS, 2009; Ebadian et al., PNAS, 2011)

$$g_f(\mathbf{X}, \mathbf{Y}) = \begin{cases} \mathbf{Y}^{\frac{1}{2}} f(\mathbf{Y}^\dagger \mathbf{X}) \mathbf{Y}^{\frac{1}{2}} & \text{if } \mathbf{Y}^\dagger \mathbf{X} \in \mathcal{X}, \mathbf{Y} \succeq \mathbf{0}, \\ +\infty & \text{otherwise,} \end{cases}$$

Matrix Perspective Reformulation Technique III: Reformulation

Mixed-Projection Conic Optimization problem

$$\begin{aligned} \min_{\mathbf{Y} \in \mathcal{Y}_n^k} \min_{\mathbf{X} \in \mathcal{S}_+^n} \quad & \langle \mathbf{C}, \mathbf{X} \rangle + \mu \cdot \text{tr}(\mathbf{Y}) + \text{tr}(f(\mathbf{X})) \\ \text{s.t.} \quad & \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} = \mathbf{Y}\mathbf{X}, \quad \mathbf{X} \in \mathcal{K}. \end{aligned}$$

Rewrite as equivalent problem which gives stronger relaxations:

$$\begin{aligned} \min_{\mathbf{Y} \in \mathcal{Y}_n^k} \min_{\mathbf{X} \in \mathcal{S}_+^n} \quad & \langle \mathbf{C}, \mathbf{X} \rangle + \mu \cdot \text{tr}(\mathbf{Y}) + \text{tr}(g_f(\mathbf{X}, \mathbf{Y})) + (n - \text{tr}(\mathbf{Y}))\omega(0) \\ \text{s.t.} \quad & \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \forall i \in [m], \quad \mathbf{X} \in \mathcal{K}, \end{aligned}$$

equivalent to applying perspective reformulation to eigenvalues of X, Y.

Can use relaxation to get feasible solutions!

Feasible Solutions via Greedily Rounding Semidefinite Relaxation

Take dual with respect to \mathbf{X} , obtain saddle-point problem:

$$\min_{\mathbf{Y} \in \text{Conv}(\mathcal{Y}_n^k)} \max_{\boldsymbol{\alpha}, \mathbf{V}_{11}, \mathbf{V}_{22} \in S^m} \lambda \cdot \text{tr}(\mathbf{Y}) + h(\boldsymbol{\alpha}) - \Omega^*(\boldsymbol{\alpha}, \mathbf{Y}, \mathbf{V}_{11}, \mathbf{V}_{22}).$$

- Greedily rounding relaxation \mathbf{Y}^* using SVD gives projection matrix \mathbf{Y} . **How good is it?**

THEOREM 7. Let \mathbf{Y}^* denote a solution to the semidefinite relaxation (19), $\mathbf{Y}^* = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^\top$ be a singular value decomposition of \mathbf{Y}^* , \mathcal{R} denote the indices of strictly fractional diagonal entries in $\boldsymbol{\Lambda}$, and $\boldsymbol{\alpha}^*(\mathbf{Y})$ denote a best choice of $\boldsymbol{\alpha}$ for a given \mathbf{Y} . Suppose that for any $\mathbf{Y} \in \mathcal{Y}_k^n$, we have $\sigma_{\max}(\boldsymbol{\alpha}^*(\mathbf{Y})) \leq L$. Then, a greedy rounding of \mathbf{Y}^* , i.e., $\mathbf{Y}_{\text{greedy}} = \mathbf{U}\boldsymbol{\Lambda}_{\text{greedy}}\mathbf{U}^\top$, where $\boldsymbol{\Lambda}_{\text{greedy}}$ is a diagonal matrix such that $\Lambda_{i,i} = 1$ for the k highest diagonal coefficients in $\boldsymbol{\Lambda}^*$, satisfies $0 \leq f(\mathbf{Y}_{\text{greedy}}) - f(\mathbf{Y}^*) \leq \epsilon$, where

- $\epsilon = ML \min(|\mathcal{R}|, n - k)$ for the spectral penalty.
- $\epsilon = \frac{\gamma}{2} \min(|\mathcal{R}|, n - k)L^2$ for the Frobenius penalty.

Feasible Solutions via Greedily Rounding Semidefinite Relaxation

Take dual with respect to X , obtain saddle-point problem:

$$\min_{\mathbf{Y} \in \text{Conv}(\mathcal{Y}_n^k)} \max_{\boldsymbol{\alpha}, \mathbf{V}_{11}, \mathbf{V}_{22} \in S^m} \lambda \cdot \text{tr}(\mathbf{Y}) + h(\boldsymbol{\alpha}) - \Omega^*(\boldsymbol{\alpha}, \mathbf{Y}, \mathbf{V}_{11}, \mathbf{V}_{22}).$$

- Greedily rounding relaxation \mathbf{Y}^* using SVD gives projection matrix \mathbf{Y} . **How good is it?**
- **Answer:** Near-exact in theory, often exact in practice.

Matrix Perspective Reformulation: Worked Example

Reduced Rank Regression: Fit **interpretable** model using small number of singular values

Step 1: Consider problem with spectral regularization:

$$\min_{\beta \in \mathbb{R}^{p \times n}} \quad \frac{1}{2m} \|\mathbf{Y} - \mathbf{X}\beta\|_F^2 + \boxed{\frac{1}{2\gamma} \|\beta\|_F^2 + \mu \cdot \text{Rank}(\beta)}$$

Where $\Omega(\mathbf{X}) = \frac{1}{2\gamma} \sum_{i=1}^n \lambda_i(\beta)^2 = \frac{1}{2\gamma} \|\beta\|_F^2$

Matrix Perspective Reformulation: Worked Example

Reduced Rank Regression: Fit **interpretable** model using small number of singular values

Step 2: Formulate as Mixed-Projection problem

$$\min_{\beta \in \mathbb{R}^{p \times n}, W \in \mathcal{Y}_n^n} \quad \frac{1}{2m} \|Y - X\beta\|_F^2 + \frac{1}{2\gamma} \|\beta\|_F^2 + \mu \cdot \text{tr}(W), W = \beta W$$

Where $\mathcal{Y}_n := \{\mathbf{P} \in S^n : \mathbf{P}^2 = \mathbf{P}\}$ is set of $n \times n$ projection matrices

Matrix Perspective Reformulation: Worked Example

Reduced Rank Regression: Fit **interpretable** model using small number of singular values

Step 3: Reformulate by taking matrix perspective

$$\min_{\boldsymbol{\beta} \in \mathbb{R}^{p \times n}, \mathbf{W} \in \mathcal{S}_+^n, \boldsymbol{\theta} \in S_+^p} \quad \frac{1}{2m} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|_F^2 + \boxed{\frac{1}{2\gamma} \text{tr}(\boldsymbol{\theta})} + \mu \cdot \text{tr}(\mathbf{W}) \quad \text{s.t.} \quad \mathbf{W} \preceq \mathbb{I}, \boxed{\begin{pmatrix} \boldsymbol{\theta} & \boldsymbol{\beta} \\ \boldsymbol{\beta}^\top & \mathbf{W} \end{pmatrix} \succeq \mathbf{0}}$$

Theoretical Contribution: Convex Hulls of Simple Low-Rank Sets

Bertsimas, Cory-Wright, and Pauphilet (21+): Theorem 2

Let T denote epigraph of spectral function under rank constraints:

$$\mathcal{T} = \{\mathbf{X} \in \mathcal{S}_+^n : \text{tr}(f(\mathbf{X})) + \mu \cdot \text{Rank}(\mathbf{X}) \leq t, \text{Rank}(\mathbf{X}) \leq k\}$$

$\omega(\cdot)$ scalar convex function such that $\text{tr}(f(X)) = \sum_{i=1}^n \omega(\lambda_i(X))$

Then, convex hull of T given by:

$$\mathcal{T}^c = \{(\mathbf{X}, \mathbf{Y}) \in \mathcal{S}_+^n \times \text{Conv}(\mathcal{Y}_n^k) : \text{tr}(g_f(\mathbf{X}, \mathbf{Y})) + \mu \cdot \text{tr}(\mathbf{Y}) + (n - \text{tr}(\mathbf{Y}))\omega(0) \leq t\}$$

Where:

- g_f matrix perspective of f
- $\text{Conv}(\mathcal{Y}_n^k) = \{\mathbf{Y} \in \mathcal{S}_+^n : \mathbf{Y} \preceq \mathbb{I}, \text{tr}(\mathbf{Y}) \leq k\}$ is convex hull of rank- k projection matrices.

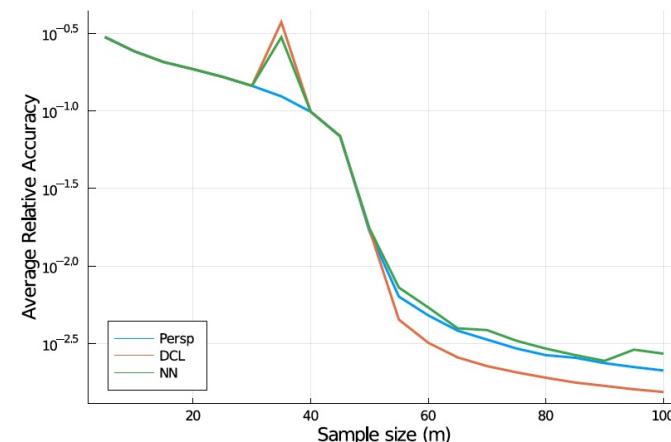
Matrix perspective reformulation gives convex hull of simple low-rank sets

Application I: Reduced Rank Regression

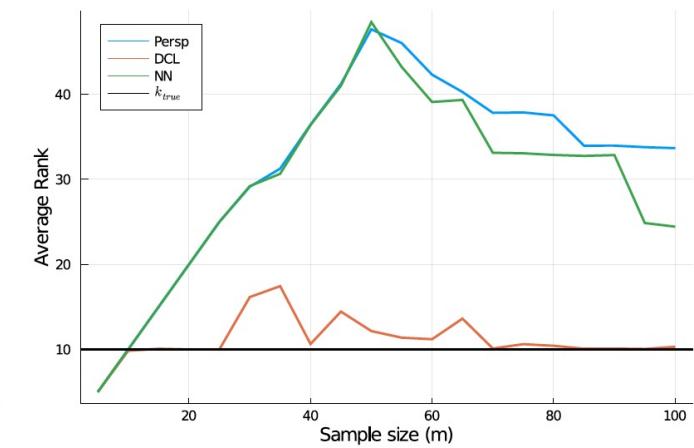
Example:

Recover rank-10 $50 \times m$ matrix:

- Vary m
- Measure MSE, rank from relaxations
- Compare against nuclear norm
- Matrix perspective dominates nuclear norm
- Saddle-point much more accurate and much faster than matrix perspective or NN, recovers true rank rapidly.
- Saddle-point w. Mosek solves for 300×300 matrices on Macbook Pro in minutes, nuclear norm takes hours for 150×150 .
- Code available on GitHub:
[ryancorywright/MatrixPerspectiveSoftware](https://github.com/ryancorywright/MatrixPerspectiveSoftware)



(a) Accuracy



(b) Rank

In practice, saddle-point relaxation *right way* to solve reduced rank regression.

Application II: Matrix Completion

Formulation:

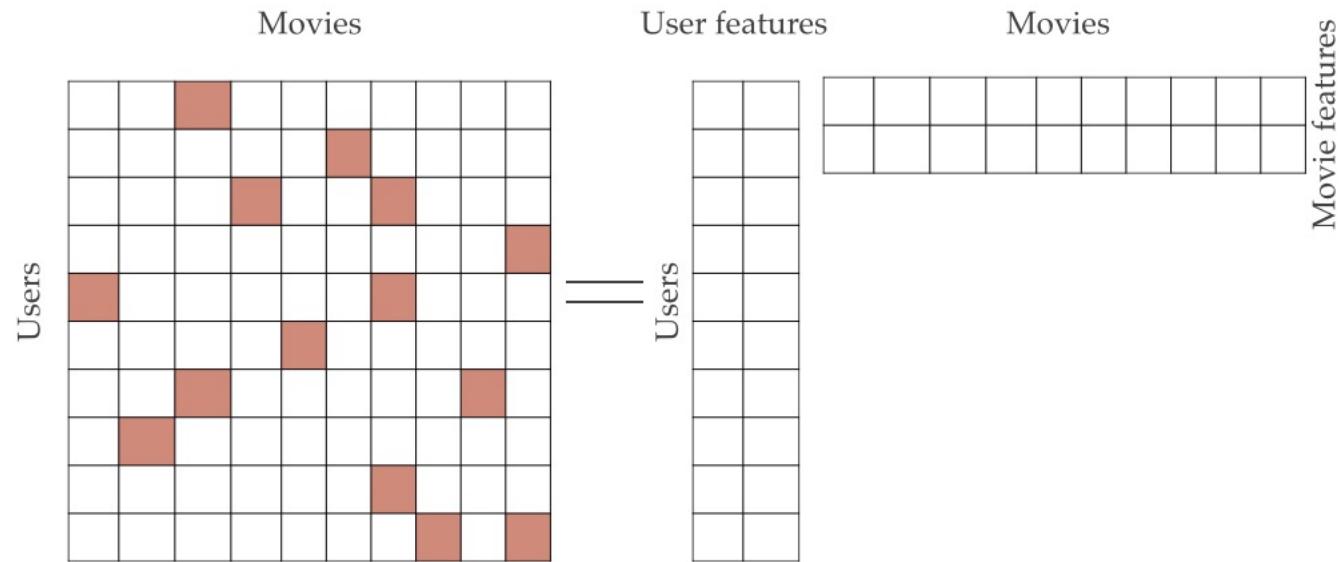
$$\min_{\mathbf{X} \in \mathbb{R}^{n \times p}} \frac{1}{2} \sum_{(i,j) \in \mathcal{I}} (X_{i,j} - A_{i,j})^2 \quad \text{s.t.} \quad \text{Rank}(\mathbf{X}) \leq k.$$

Decision variables/Problem data

$X_{i,j}$: Predicted rating movie j by user i
 $A_{i,j}$: Reported rating movie j by user i

Movie Recommendation:

- Given user movie ratings, predict ratings for unseen movies.
- To make problem tractable, assume ratings depend on k factors (lead actor, lead actress, director, genre, year, ...)

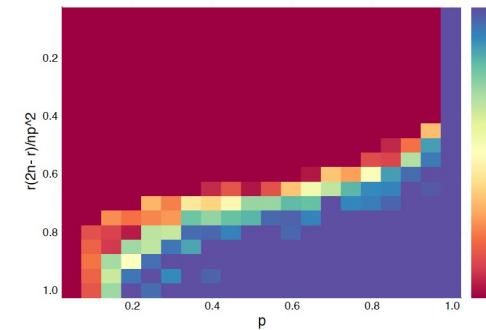


Application II: Matrix Completion

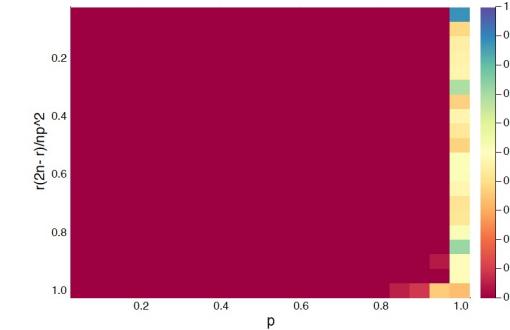
Example:

Recover low-rank 100×100 matrix:

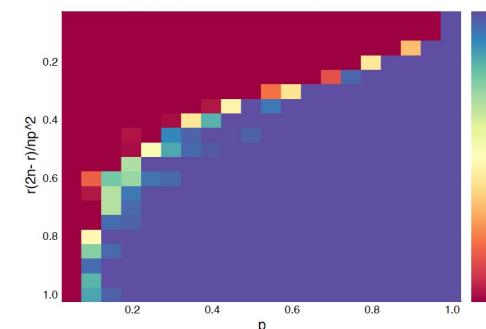
- Vary rank, proportion entries sampled
- Measure % time recover matrix to 1% MSE (more purple=better)
- Nuclear norm by far worst approach
- New penalty better, new penalty with rounding much better
- Combining greedy rounding with local search: solutions within 1% of optimality in practice, 10% in theory.
- Code available on GitHub:
[ryancorywright/MixedProjectionSoftware](https://github.com/ryancorywright/MixedProjectionSoftware)



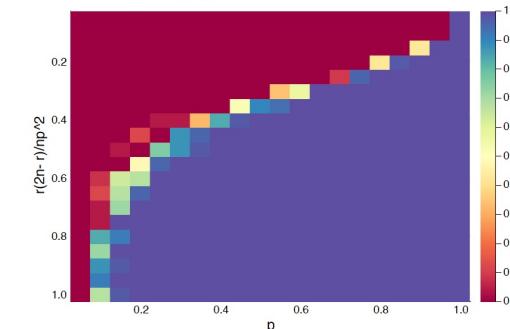
(a) New Penalty



(b) Nuclear Norm



(c) SVD+Local Improvement



(d) New Penalty+Local Improvement

In practice, new penalty is viable and *often* more accurate

Details in: [Mixed-Projection Conic Optimization: A New Paradigm for Modeling Rank Constraints](#)
D. Bertsimas, R. Cory-Wright, J. Pauphilet, Operations Research, 2021.
Winner, 2020 INFORMS George Nicholson Best Paper Competition

Technique Works for Many Other Applications!

Machine Learning

Decompose matrix into sparse plus low-rank matrix

- See Nicholas Johnson's talk later in the session

Control Theory

Identify minimum complexity realization of a system

Algebraic Geometry

Minimum degree sum-of-squares decomposition of polynomial

Computer Vision

Decompose a matrix into a product of non-negative factors

Many, many more...

Conclusion

Matrix perspective is natural generalization of perspective reformulation

- Exploit separability of eigenvalues to obtain “embarrassingly tight” formulation.
- Leads to relaxations which outperform state-of-the-art for central problems in OR/ML.
- Suggests this is a very general story, often useful to think about problems this way.

Two future directions:

1. Improve generality → extend framework to tensors, non-negative polynomials.
2. Develop rounding → find “best” rounding strategy for matrix perspective relax



Thank you for listening!
Lingering questions? Email ryancw@mit.edu