

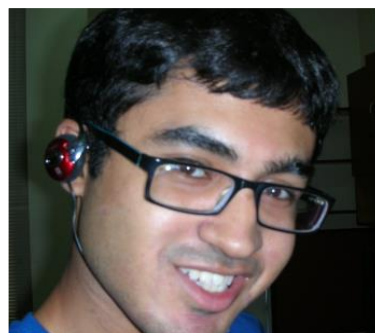
# Iterative Hard Thresholding for Sparse/Low-rank Linear Regression

Prateek Jain

Microsoft Research, India



Ambuj Tewari  
Univ of Michigan



Purushottam Kar  
MSR, India



Praneeth Netrapalli  
MSR, NE

# Microsoft Research India



## Our work

- Foundations
- Systems
- Applications
- Interplay of society and technology
- Academic and government outreach

## Our vectors of impact

- Research impact
- Company impact
- Societal impact

# Machine Learning and Optimization @ MSRI

- High-dimensional Learning & Optimization
- Extreme Classification
- Online Learning / Multi-armed Bandits

**We are Hiring!**

- Interns
- PostDocs
- Applied Researchers
- Full-time Researchers

- Ranking & Recommendation

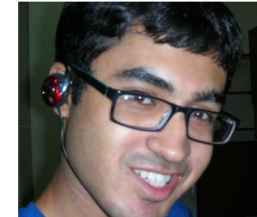
Manik Varma



Prateek Jain



Purushottam Kar



Ravi Kannan



Amit Deshpande



Navin Goyal



Sundarajan S.



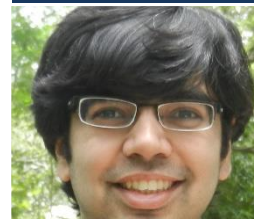
Vinod Nair



Sreangsu Acharyya



Kush Bhatia



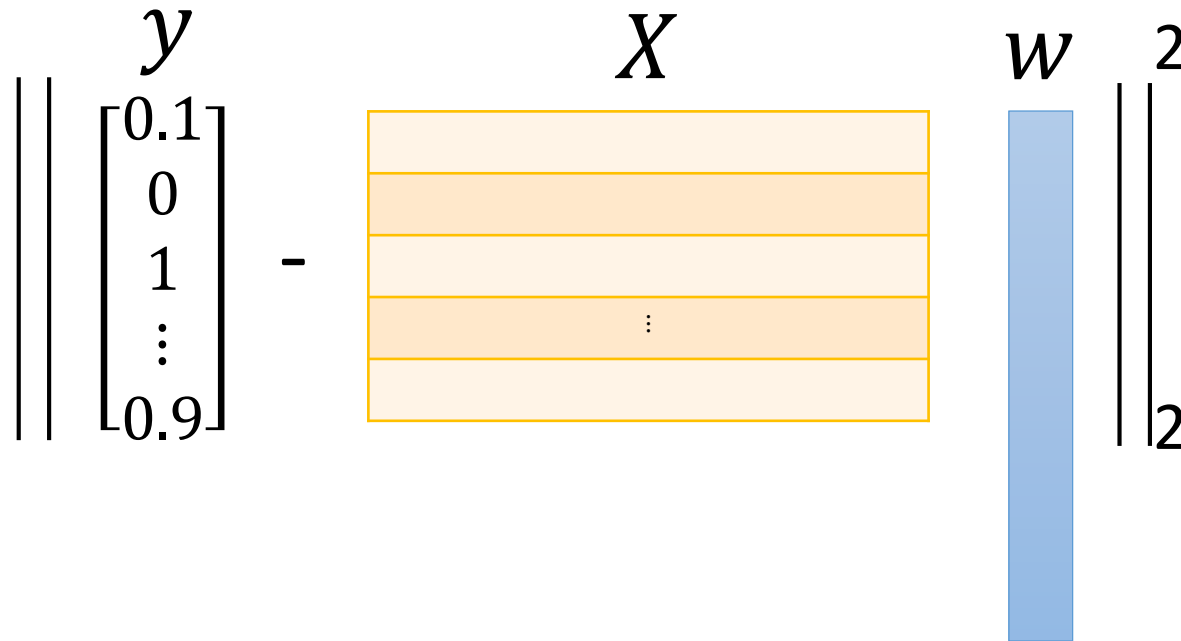
Aditya Nori



Raghavendra Udapa



# Learning in High-dimensions



• Need to solve:  $\min_w ||y - Xw||_2^2$  s. t.  $w \in \mathcal{C}$

•  $\mathcal{C}$  can be:

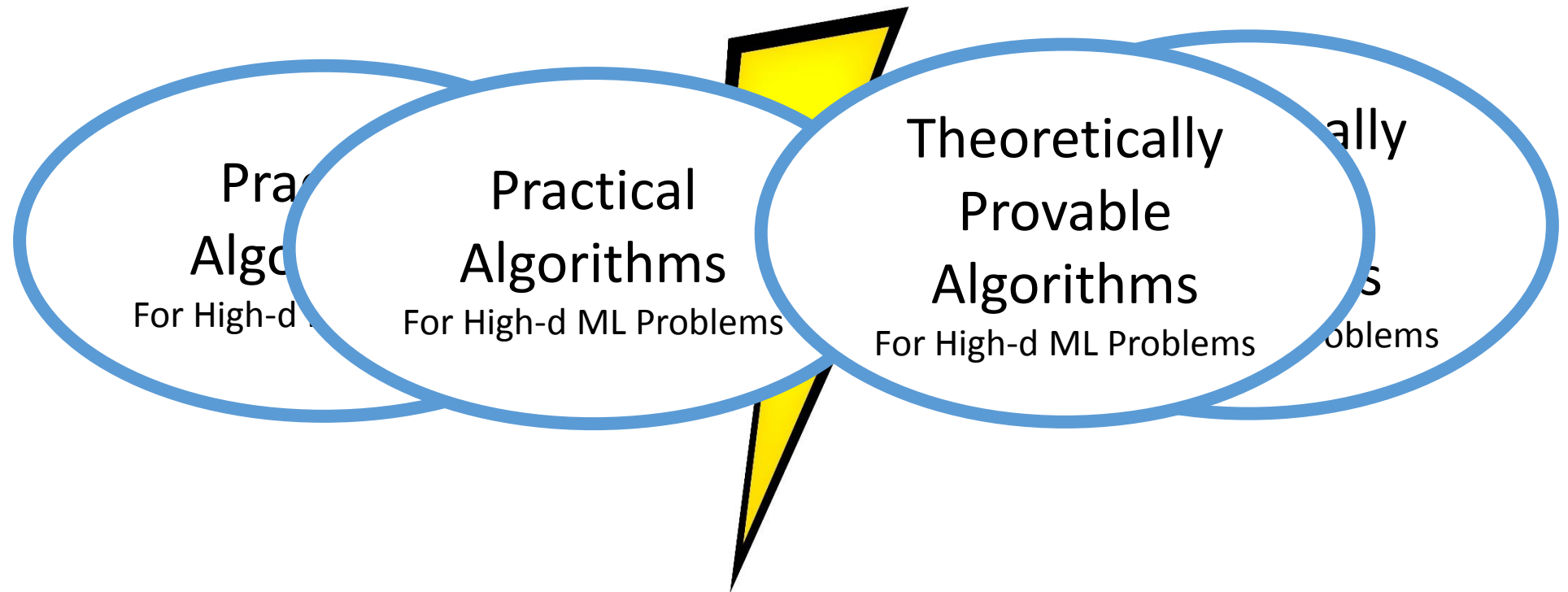
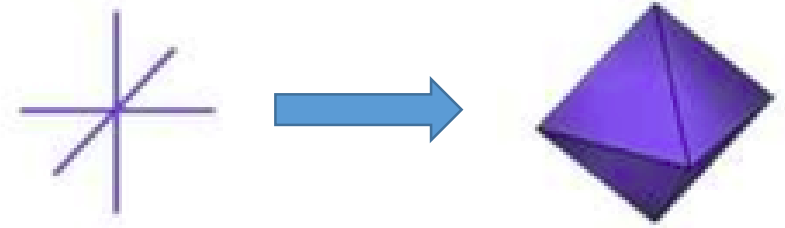
- Set of sparse vectors
- Set of group-sparse vectors
- Set of low-rank matrices



- Non-convex
- Comp. Complexity: NP-Hard

# Overview

- Most popular approach: convex relaxation
  - Solvable in poly-time
  - Guarantees under certain assumptions
  - **Slow in practice**



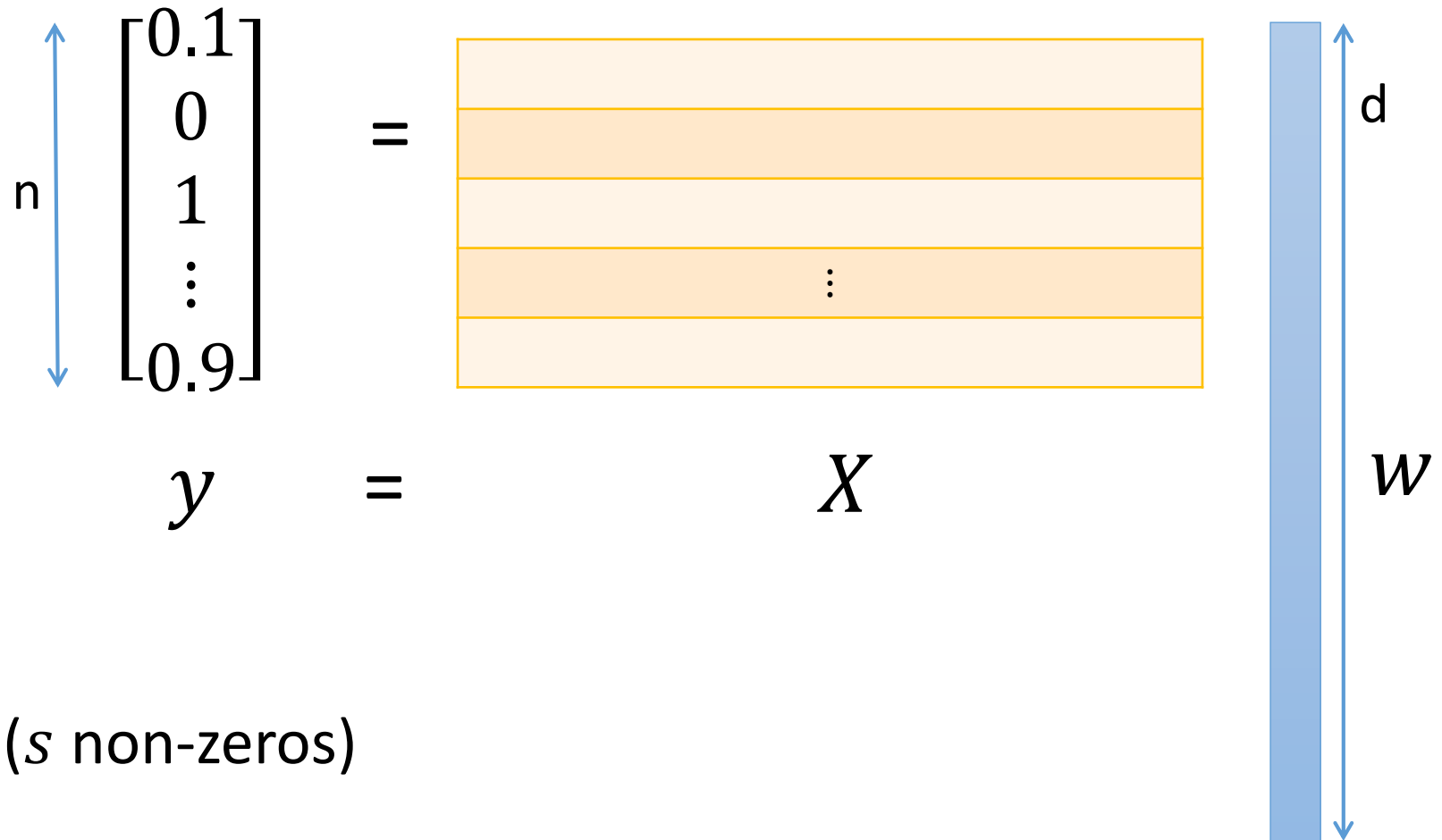
# Results

- Sparse Regression [[Garg & Khandekar. ICML 2008](#); [J., Kar, Tewari. NIPS14](#)]
  - $L_0$  – constraint
- Matrix Completion/Regression [[J., Netrapalli, Sanghavi. STOC 2013](#); [Hardt & Wooters. COLT 2014](#)]
  - Low-rank constraint
- Robust Regression [[Loh & Wainwright. NIPS 2013](#) ; [Bhatia, J., Kar. Submitted, 2015](#)]
- Tensor Factorization and Completion [[Anandkumar et al. Arxiv 2012](#); [J., Oh. NIPS14](#)]
  - Low-tensor rank constraint
- Dictionary Learning [[Agarwal et al. COLT 2014](#); [Arora et al. COLT 2014](#)]
  - Non-convex bilinear form + Sparsity constraint
- Phase Sensing [[Netrapalli, J., Sanghavi. NIPS13](#) ; [Candes et al. Arxiv'2014](#)]
  - System of quadratic equations
- Low-rank matrix approximation [[Bhojanapalli, J., Sanghavi. SODA15](#)]

# Outline

- Sparse Linear Regression
  - Lasso
- Iterative Hard Thresholding
  - Our Results
- Low-rank Matrix Regression
- Low-rank Matrix Completion
- Conclusions

# Sparse Linear Regression



The diagram illustrates the equation  $y = Xw$  for sparse linear regression. On the left, a vertical vector  $y$  is shown with a blue double-headed arrow indicating its length  $n$ . The vector contains the values  $0.1$ ,  $0$ ,  $1$ , a vertical ellipsis, and  $0.9$ . This is followed by an equals sign. To the right is a matrix  $X$ , represented by a grid of orange rectangles. The first row is light orange, the second is medium orange, the third is light orange, the fourth is medium orange and contains a vertical ellipsis, and the fifth is light orange. A blue vertical bar to the right of the matrix represents the weight vector  $w$ , with a blue double-headed arrow indicating its length  $d$ . Below the matrix, the label  $X$  is centered, and below the weight vector, the label  $w$  is centered. The overall equation is  $y = Xw$ .

- But:  $n \ll d$
- $w$ :  $s$  –sparse ( $s$  non-zeros)



# Sparse Linear Regression

$$\begin{aligned} \min_w & \|y - Xw\|^2 \\ \text{s.t.} & \|w\|_0 \leq s \end{aligned}$$

- $\|y - Xw\|^2 = \sum_i (y_i - \langle x_i, w \rangle)^2$
- $\|w\|_0$ : number of non-zeros
- NP-hard problem in general ☹
  - $L_0$ : non-convex function

# Convex Relaxation

$$\begin{aligned} \min_w & \|y - Xw\|^2 \\ \text{s.t.} & \|w\|_0 \leq s \end{aligned}$$

- Relaxed Problem:

$$\begin{aligned} \min_w & \|y - Xw\|^2 \\ \text{s.t.} & \|w\|_1 \leq \mu(s) \end{aligned}$$



$$\min_w \|y - Xw\|^2 + \lambda \|w\|_1$$

Lasso Problem

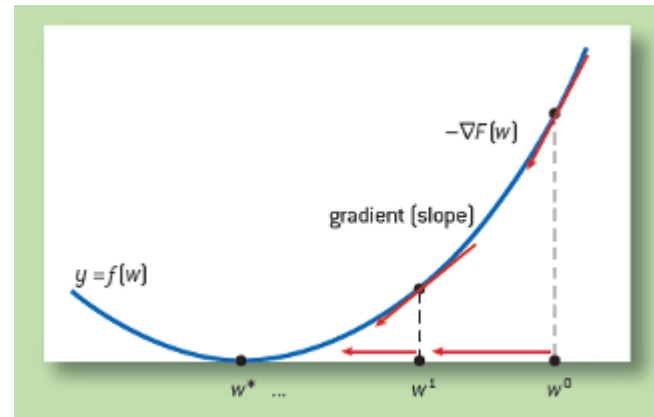
Non-differentiable

- $\|w\|_1 = \sum_i |w_i|$ 
  - Known to promote sparsity
- Pros: a) Principled approach, b) Solid theoretical guarantees
- Cons: Slow to optimize

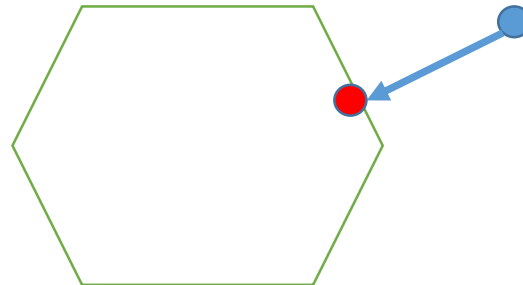
# Our Approach : Projected Gradient Descent

$$\min_w f(w) = \|y - Xw\|^2$$
$$s.t. \|w\|_0 \leq s$$

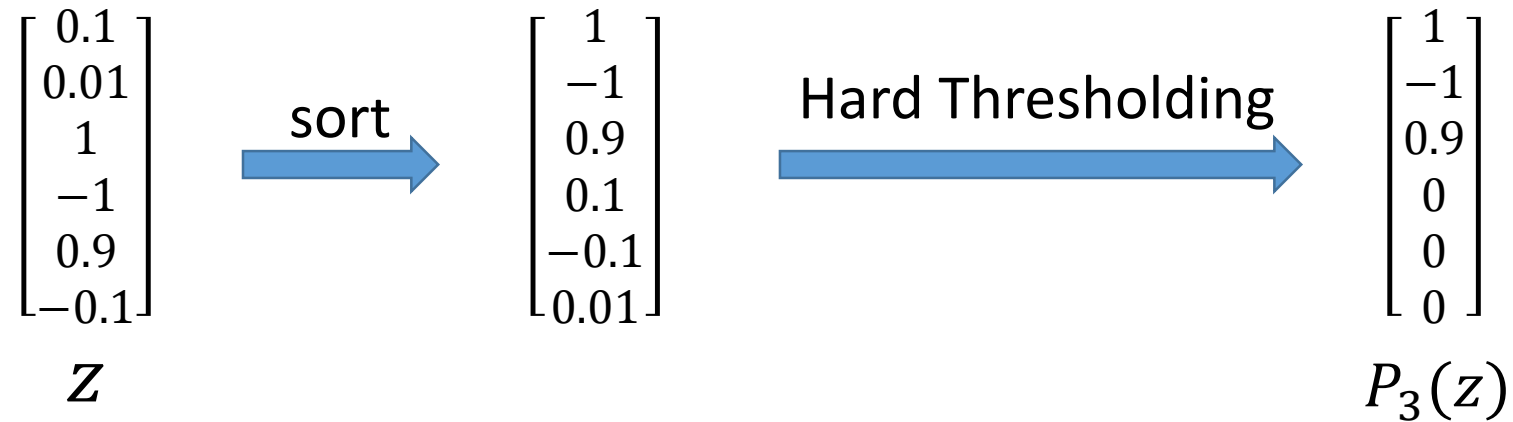
- $w_{t+1} = w_t - \partial_{w_t} f(w_t)$



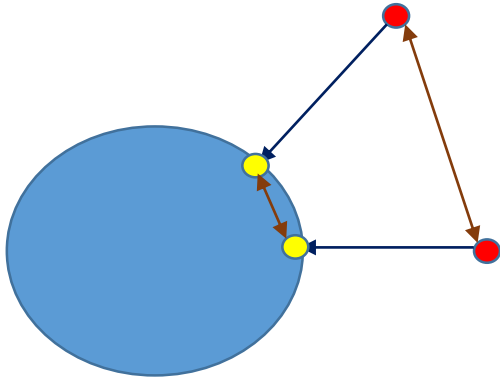
- $w_{t+1} = P_S(w_{t+1})$



# Projection onto $L_0$ ball?



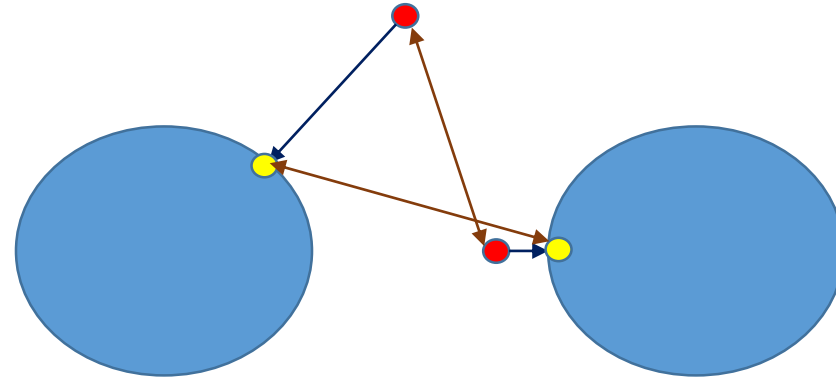
# Convex-projections vs Non-convex Projections



$$\|P_C(a) - P_C(b)\| \leq \|a - b\|$$

$C$ : convex set

1st order Optimality condition



$$\|P_C(a) - a\| \leq \|u - a\|, \quad \forall u \in C$$

$C$ : non-convex set

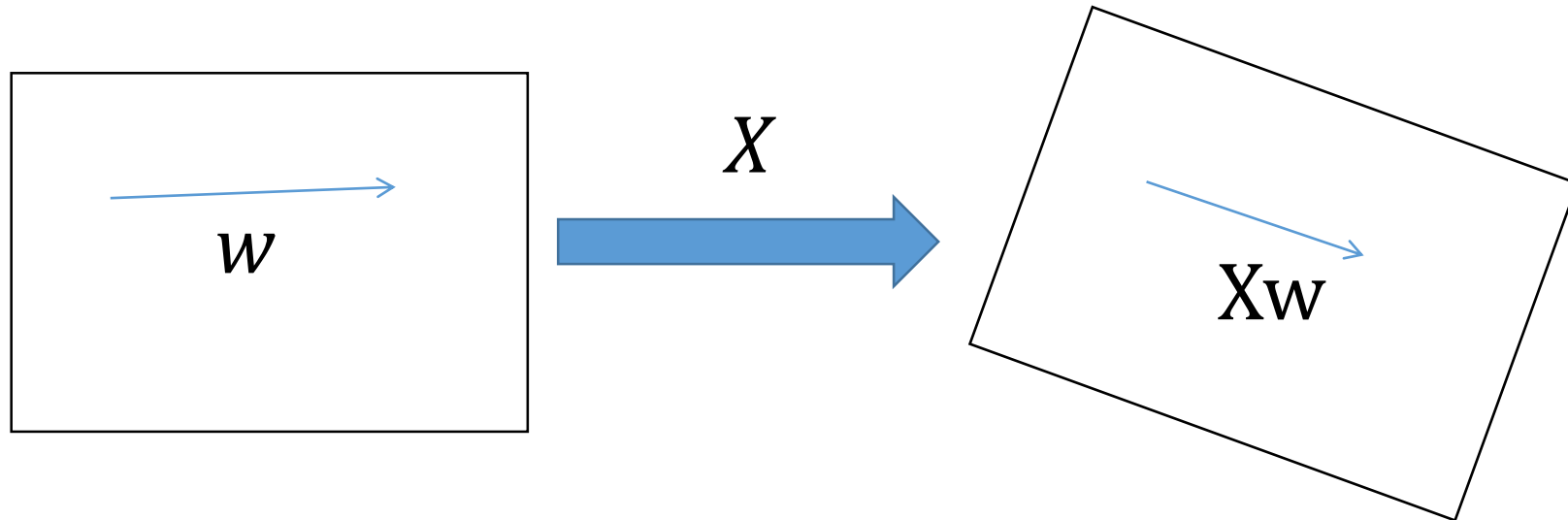
0-th order Optimality condition

- 0 order condition sufficient for convergence of Proj. Grad. Descent?
- In general, **NO** 😞
- But, for certain *specially structured* problems, **YES!!!**

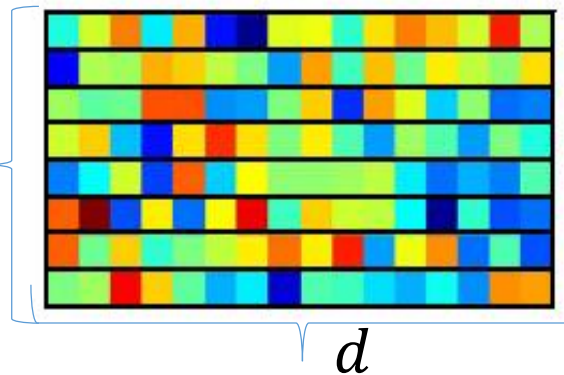
# Restricted Isometry Property (RIP)

- $X$  satisfies RIP if, for all **sparse** vectors  $X$  acts as an Isometry
- Formally: For all  $s$ -sparse  $\mathbf{w}$

$$(1 - \delta_s) \|\mathbf{w}\|^2 \leq \|X\mathbf{w}\|^2 \leq (1 + \delta_s) \|\mathbf{w}\|^2$$



# Popular RIP Ensembles

$$n = O\left(\frac{s}{\delta_s^2} \log \frac{d}{s}\right)$$


$X$

- Most popular examples:
  - $X_{ij} \sim N(0, 1/\sqrt{n})$
  - $X_{ij} = +\frac{1}{\sqrt{n}}$  (w.p.  $\frac{1}{2}$ ) and  $-\frac{1}{\sqrt{n}}$  (w.p.  $\frac{1}{2}$ )

# Proof under RIP

Assume:  $y = Xw^*$ ,  $\min_{w, \|w\|_0 \leq s} f(w) = \|y - Xw\|^2 = \|X(w - w^*)\|^2$

Recall:  $w_{t+1} = P_S(w_t - X^T X(w_t - w^*))$

Hard  
Thresholding

$$\|w_{t+1} - (w_t - X^T X(w_t - w^*))\| \leq \|w^* - (w_t - X^T X(w_t - w^*))\|$$

$I: \text{supp}(w_t) \cup \text{supp}(w_{t+1}) \cup \text{supp}(w^*)$ ,  $|I| \leq 3s$

$$\|w_{t+1} - (w_t - X_I^T X_I(w_t - w^*))\|^2 \leq \|w^* - (w_t - X_I^T X_I(w_t - w^*))\|^2$$

Triangle inequality

$$\begin{aligned} \|w_{t+1} - w^*\| &\leq 2\|(I - X_I^T X_I)(w_t - w^*)\| \\ &\leq 2\delta_{3s}\|w_t - w^*\| \end{aligned}$$

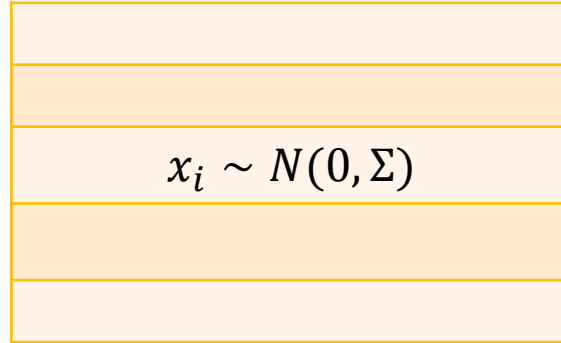
RIP



# What if RIP is not possible?

- $y_i = \langle x_i, w^* \rangle$

- $x_i \sim N(0, \Sigma)$



$$\Sigma = \begin{bmatrix} 1 & 1 - \epsilon & 0 \\ 1 - \epsilon & 1 & 0 \\ 0 & 0 & I_{d-2 \times d-2} \end{bmatrix}$$

- Eigenvalues of  $\Sigma = 2 - \epsilon, \epsilon$
- $\delta_s \geq \delta_2 = 1 - \epsilon$
- So,  $\delta_s < \frac{1}{2}$  doesn't hold even for infinite samples
  - Problem is solvable for  $O(d)$  samples using standard regression

# Iterative Hard Thresholding: Larger Sparsity

- $w_1 = 0$
- For  $t=1, 2, \dots$ 
  - $w_{t+1} = P_{s'}(w_t - \eta \nabla_w f(w_t))$
- $s' \geq s$

# Stronger Projection Guarantee

$$\|P_{s'}(z) - z\|_2^2 \leq \frac{d - s'}{d - s} \|P_s(z) - z\|_2^2$$

- $d$ : dim of  $z$
- $s \leq s'$

# Statistical Guarantees

$$\min_w f(w) = \|y - Xw\|^2$$

Statistically:  $n \geq \frac{\sigma^2 \cdot s \log d}{\epsilon^2}$

Known Computation Lower-bound:  $\frac{\kappa \cdot \sigma^2 \cdot s \log d}{\epsilon^2}$

Same as Lasso

- $w^*$ :  $s$  -sparse

$$\|\hat{w} - w^*\| \leq \epsilon,$$

$$n \geq \frac{\kappa^2 \cdot \sigma^2 \cdot s \log d}{\epsilon^2}$$

- $\kappa = \lambda_1(\Sigma) / \lambda_d(\Sigma)$
- Recall,  $w_{t+1} = P_{s'}(w_t - \eta \nabla_w f(w))$ 
  - $s' = \kappa^2 s$

# General Result for Any Function

- $f: R^d \rightarrow R$
- $f$ : satisfies RSC/RSS, i.e.,

$$\alpha_s \cdot I_{d \times d} \preceq H(w) \preceq L_s \cdot I_{d \times d}, \quad \text{if } \|w\|_0 \leq s$$

- IHT guarantee:  $f(w_T) \leq f(w^*) + \epsilon$

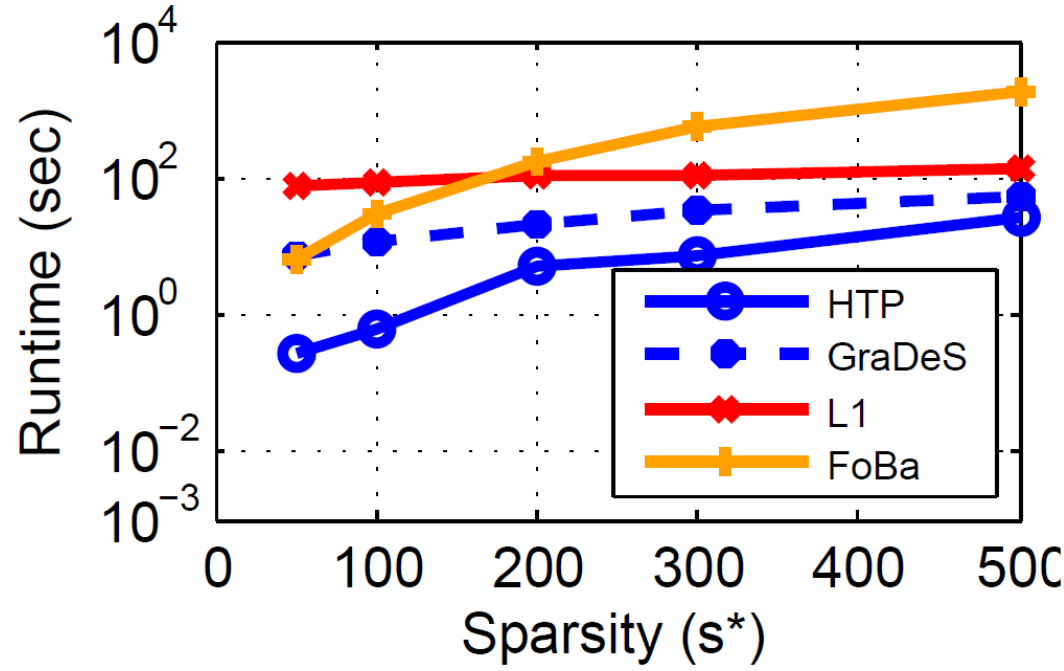
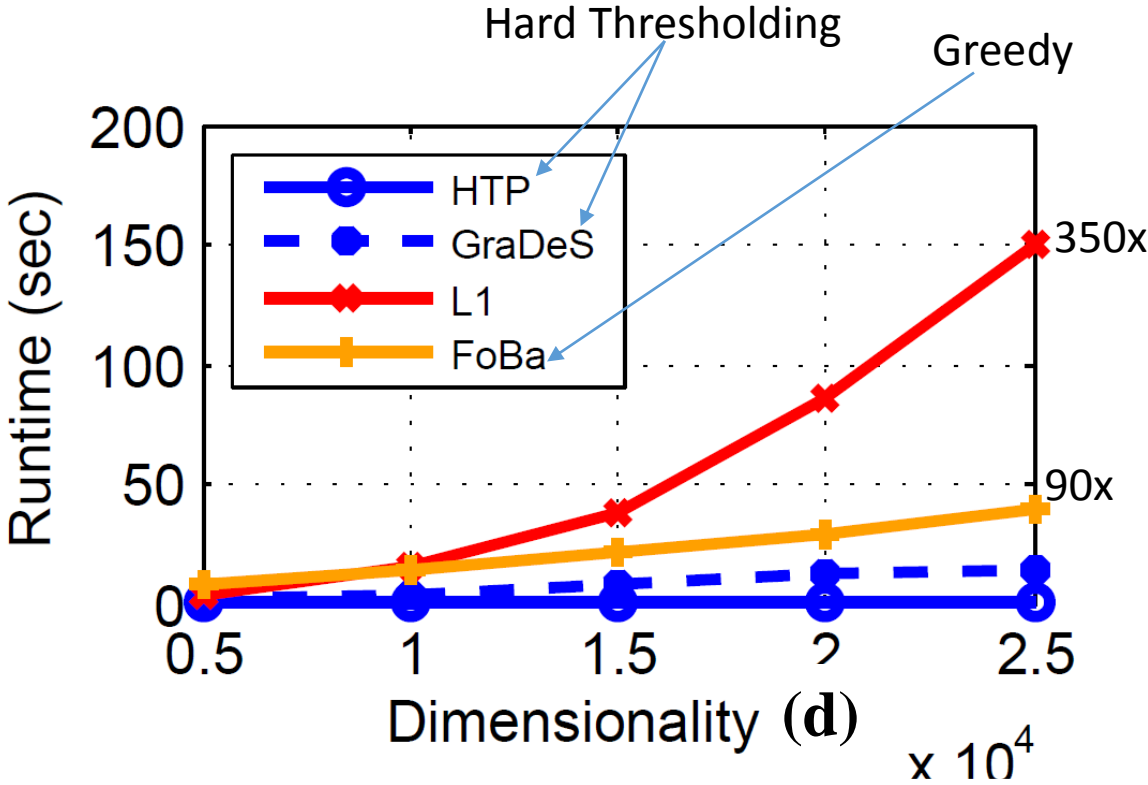
After  $T = O\left(\frac{\log\left(\frac{f(w^0)}{\epsilon}\right)}{\log\left(1 - \frac{L_{s'}}{\alpha_{s'}}\right)}\right)$  steps

- If  $\|w^*\|_0 \leq s$  and  $s' \geq 10 \frac{L_{s'}^2}{\alpha_{s'}} s$

# Extension to other Non-convex Procedures

- IHT-Fully Corrective
  - HTP [Foucart'12]
- CoSAMP [Tropp & Neadell'2008]
- Subspace Pursuit [Dai & Milenkovic'2008]
- OMPR [J., Tewari, Dhillon'2010]
- Partial hard thresholding and two-stage family [J., Tewari, Dhillon'2010]

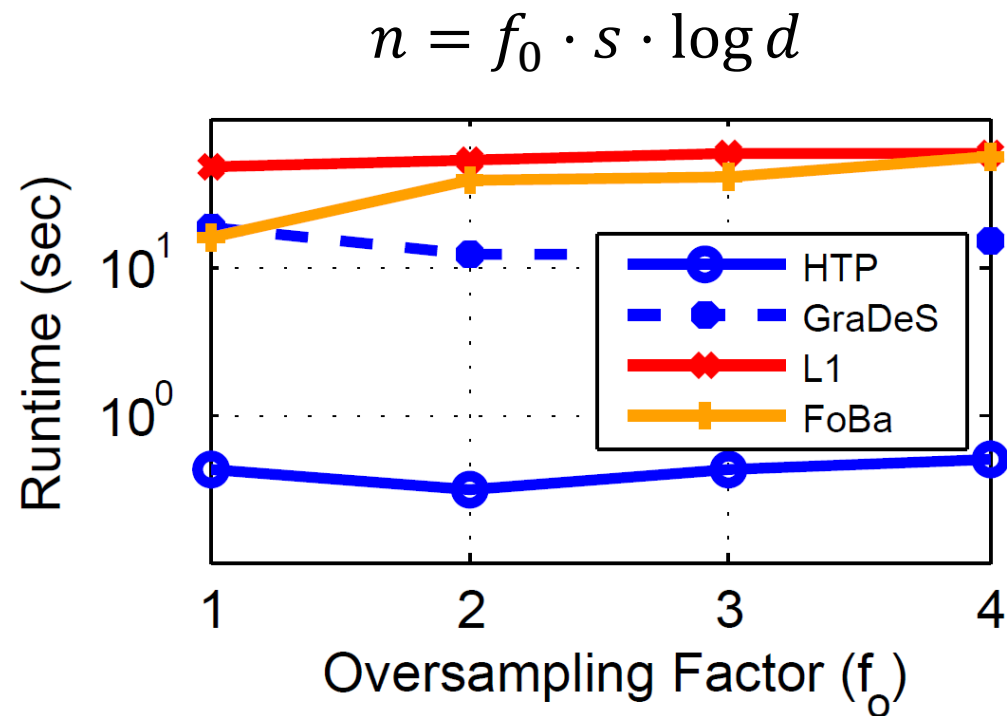
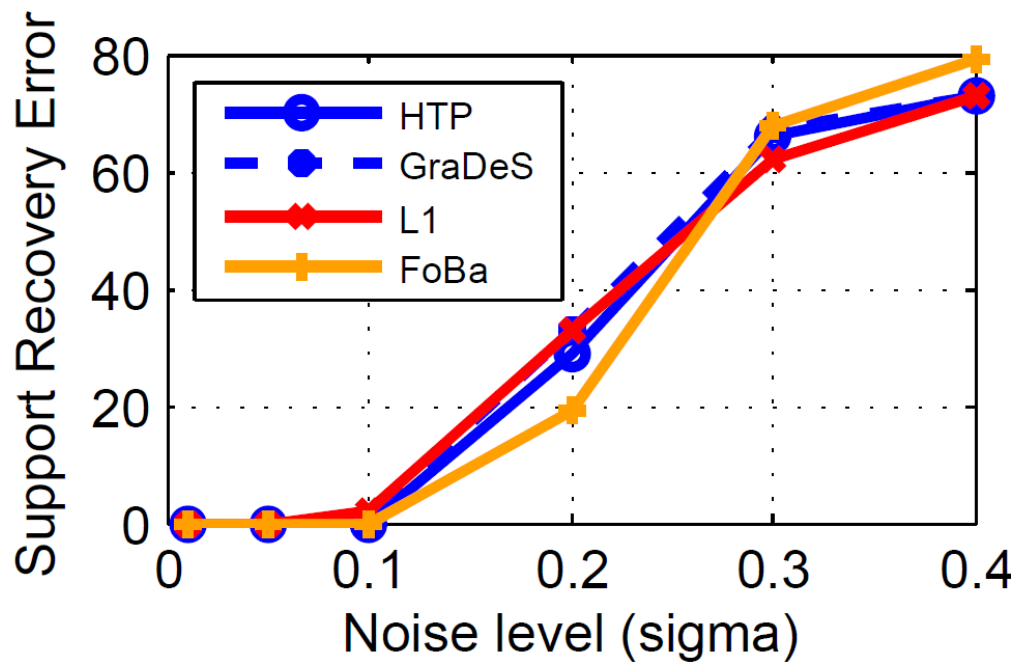
# Empirical Results



$n = 2 \cdot s \cdot \log(d), s = 300, \kappa = 1$

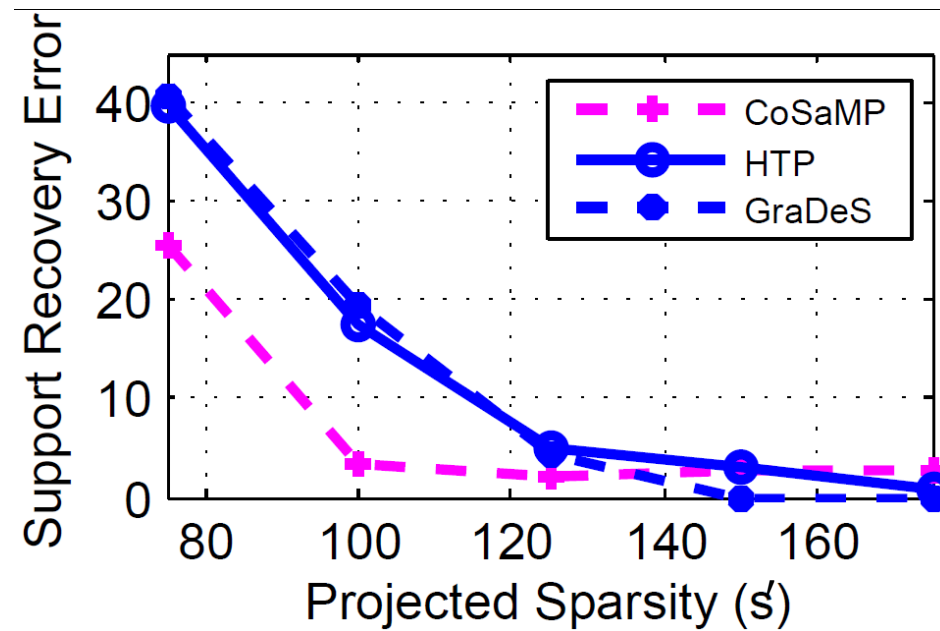
$n = 2 \cdot s \cdot \log(d), d = 20,000, \kappa = 1$

# More Empirical Results





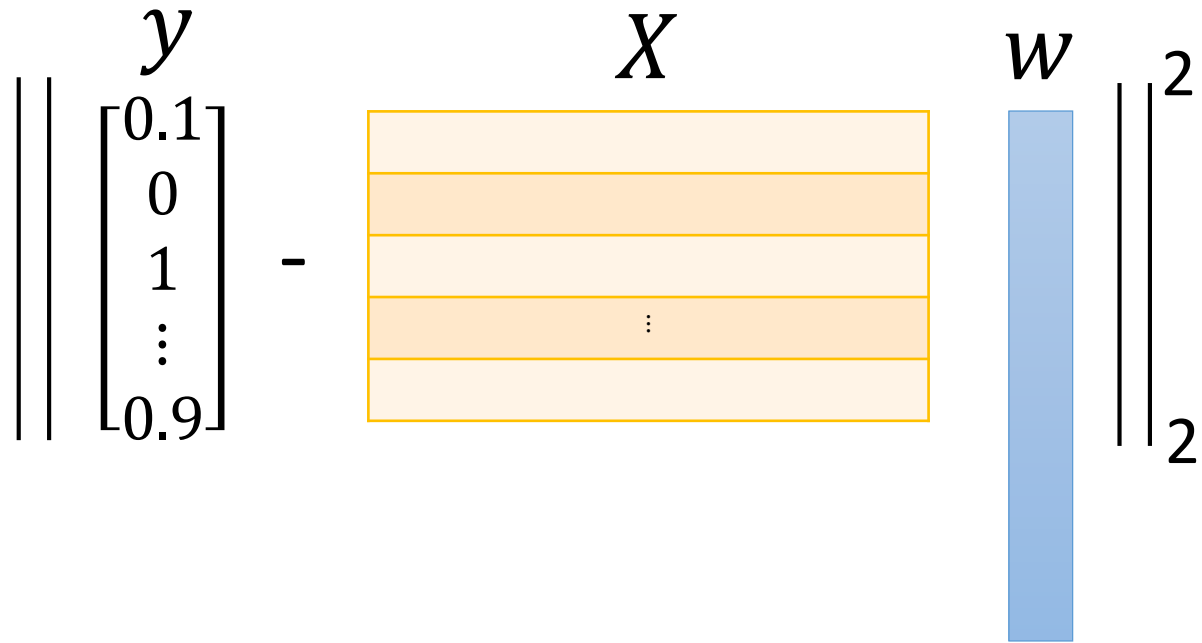
# Empirical Results: poor condition number



$$n = 2 \cdot s \cdot \log(d), s = 50, d = 20,000$$

$$\kappa = 50$$

# Low-rank Matrix Regression



- $W$ :  $d_1 \times d_2$  matrix
- $\text{rank}(W) = r \ll \min(d_1, d_2)$

# Low-rank Matrix Regression

$$\begin{aligned} \min_W f(W) &= \|y - X \cdot W\|^2 \\ \text{s.t. } \text{rank}(W) &\leq r \end{aligned}$$

- Convex relaxation:  $\text{rank}(W) \Rightarrow \|W\|_*$ 
  - $\|W\|_*$  = sum of singular values of  $W$
  - Several interesting results: [\[Recht et al.'2007, Negahban et al'2009\]](#)...
- Projected Gradient Descent:
  - $W_{t+1} = P_k(W_t - \eta \nabla_W f(W_t)), \forall t$
  - $k \geq r$
- $P_k(Z) = U_k \Sigma_k V_k^T$  where  $Z = U \Sigma V^T$

# Statistical Guarantees

$$y_i = \langle x_i, W^* \rangle + \eta_i$$



- $x_i \sim N(0, \Sigma) \in R^d$
- $\eta_i \sim N(0, \sigma^2)$
- $W^* \in R^{d_1 \times d_2}$ ,  $\text{rank}(W^*) = r$

$$\|\widehat{W} - W^*\|_2 \leq \frac{\sigma \cdot \kappa \cdot \sqrt{r(d_1 + d_2) \log(d_1 + d_2)}}{\sqrt{n}}$$

- $\kappa = \frac{\lambda_1(\Sigma)}{\lambda_d(\Sigma)}$ ,  $k = \kappa^2 r$

# Low-rank Matrix Completion

|        |   | users |   |   |   |   |   |   |   |   |    |    |    |
|--------|---|-------|---|---|---|---|---|---|---|---|----|----|----|
|        |   | 1     | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| movies | 1 | 1     |   | 3 |   |   | 5 |   |   | 5 |    | 4  |    |
|        | 2 |       |   | 5 | 4 |   |   | 4 |   |   | 2  | 1  | 3  |
|        | 3 | 2     | 4 |   | 1 | 2 |   | 3 |   | 4 | 3  | 5  |    |
|        | 4 |       | 2 | 4 |   | 5 |   |   | 4 |   |    | 2  |    |
|        | 5 |       |   | 4 | 3 | 4 | 2 |   |   |   |    | 2  | 5  |
|        | 6 | 1     |   | 3 |   | 3 |   |   | 2 |   |    | 4  |    |

 - unknown rating     - rating between 1 to 5

$$\min_W \sum_{(i,j) \in \Omega} (W_{ij} - M_{ij})^2$$

s. t    **rank**( $W$ )  $\leq r$

$\Omega$ : set of known entries

- Special case of low-rank matrix regression
- However, assumptions required by the regression analysis not satisfied

# Guarantees

- Projected Gradient Descent:
  - $W_{t+1} = P_r(W_t - \eta \nabla_W f(W_t)), \quad \forall t$
- Show  $\epsilon$ -approximate recovery in  $\log \frac{1}{\epsilon}$  iterations
- Assuming:
  - $M$ : incoherent
  - $\Omega$ : uniformly sampled
  - $|\Omega| \geq n \cdot r^5 \cdot \log^3 n$
- First near linear time algorithm for **exact** Matrix Completion with finite samples

# Tale of two Lemmas

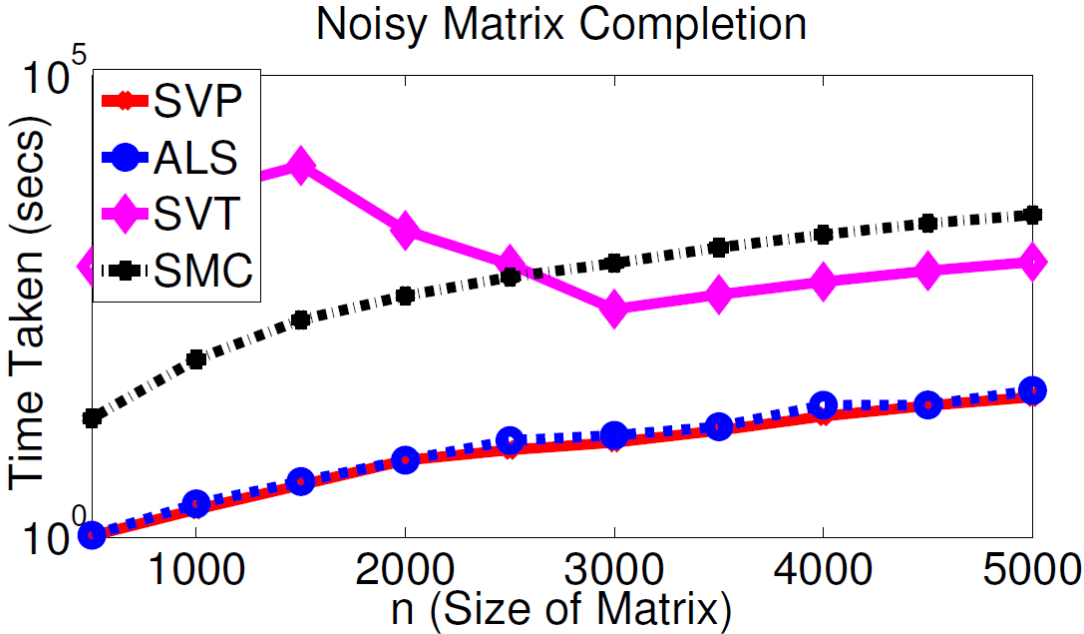
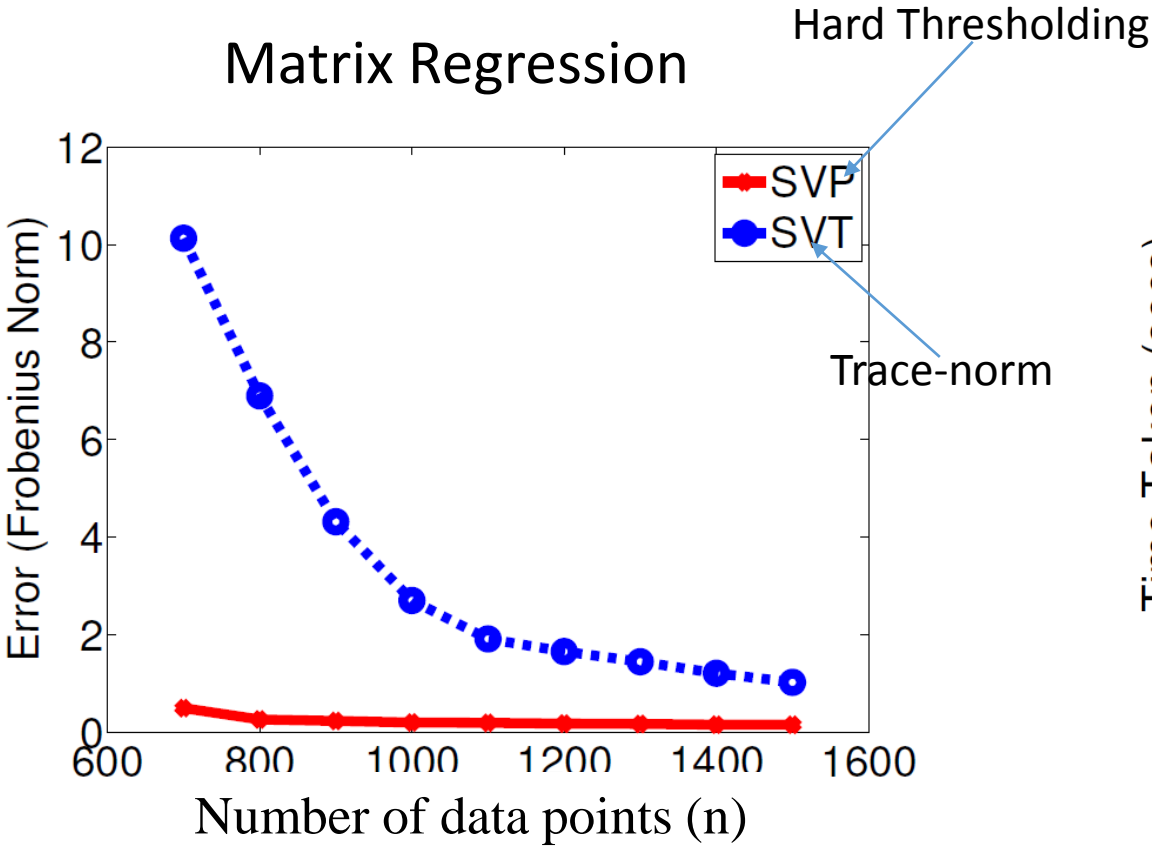
- Lemma 1: Perturbation bound with  $L_\infty$  bounds

$$\|P_r(M + E_t) - M\|_\infty \leq .5 \|E_t\|_\infty$$

- Standard bounds only give:  $\|P_r(M + E_t) - M\|_2 \leq 2\|E_t\|_2$
  - $M$ : incoherent
  - $E_t$ : zero-mean with small variance
- Lemma 2: Davis-Kahn style result for matrix perturbation
    - If  $\sigma_{k+1}(M) < .25 \sigma_k(M)$  and  $\|E_t\|_F \leq .25 \sigma_k(M)$

$$\|P_k(M + E_t) - P_k(M)\|_F \leq c(\sqrt{k} \|E_t\|_2 + \|E\|_F)$$

# Empirical Results



$$r = 100, |\Omega| = 5 r n \log n$$



# Summary

- High-dimensional problems
  - $n \ll d$
- Need to impose structure on  $w$
- Typical structures
  - Sparsity
  - Low-rank
  - Low-rank+Sparsity
- Non-convex sets but easy projection
- Proof of convergence (linear rate)
  - Under suitable generative model assumptions

# Future Work

- Generalized theory for such provable non-convex optimization
- Performance analysis on different models
- Empirical comparisons on “real-world” datasets

Questions?