# Robust Recovery via Implicit Bias of Discrepant Learning Rates for Double Over-parameterization
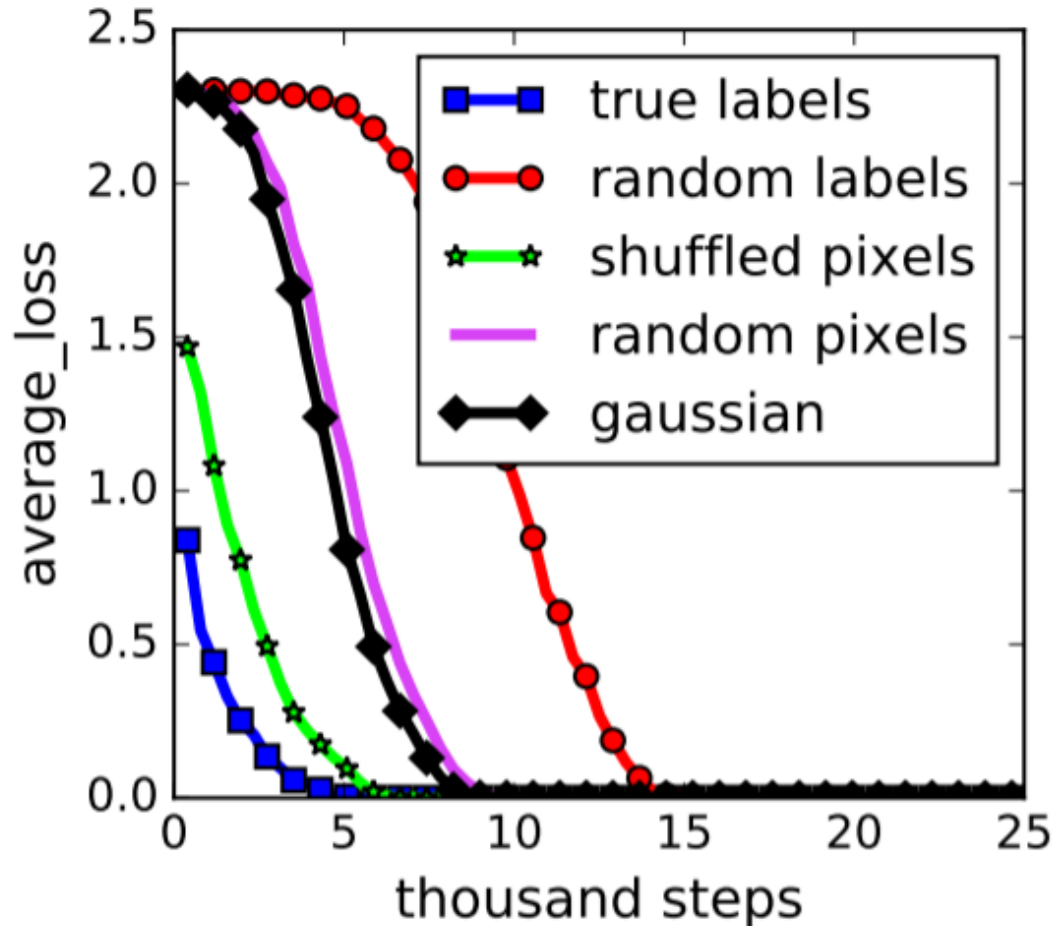
Chong You[†],     Zhihui Zhu[‡],     Qing Qu[#],     Yi Ma[†]

[†]Department of EECS, University of California, Berkeley

[‡]Department of Electrical and Computer Engineering, University of Denver

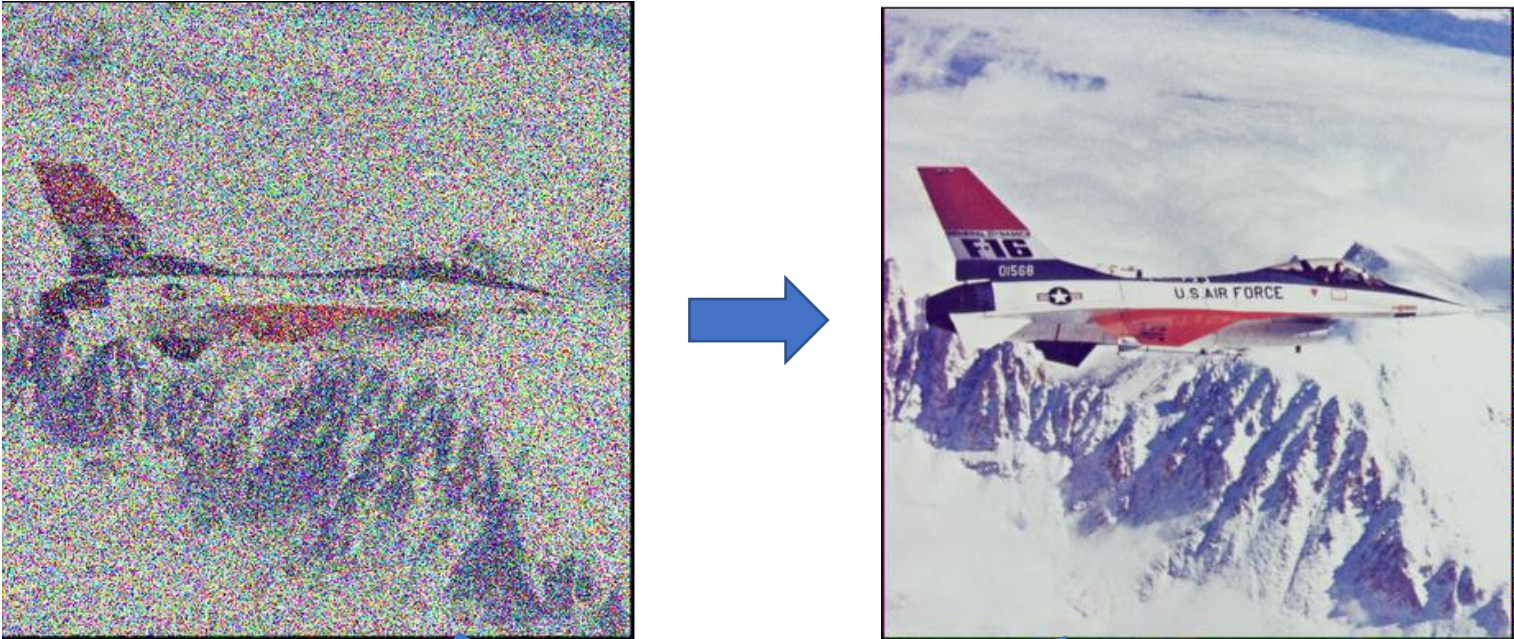[#]Center for Data Science, New York University

# Over-parameterization and Overfitting
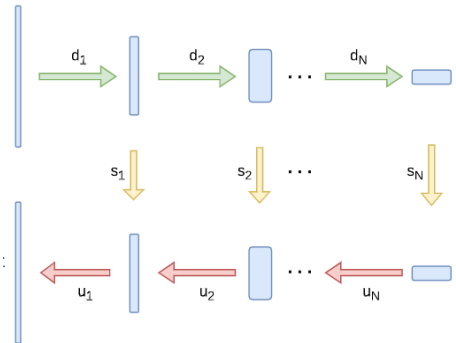


(Zhang et. al., ICLR'17)

- A classification network can overfit to label corruption

# Deep Image Prior (DIP) for Image Recovery



$$\min_{\boldsymbol{\theta}} \| \underbrace{\boldsymbol{y}}_{\text{corrupted input}} - \underbrace{f(\boldsymbol{\theta})}_{\text{recovered image}} \|_1$$

- **Idea**: CNN architecture encodes priors for clean images

Ulyanov, Vedaldi, Lempitsky, Deep Image Prior, *CVPR 2018*

# Over-Parameterization
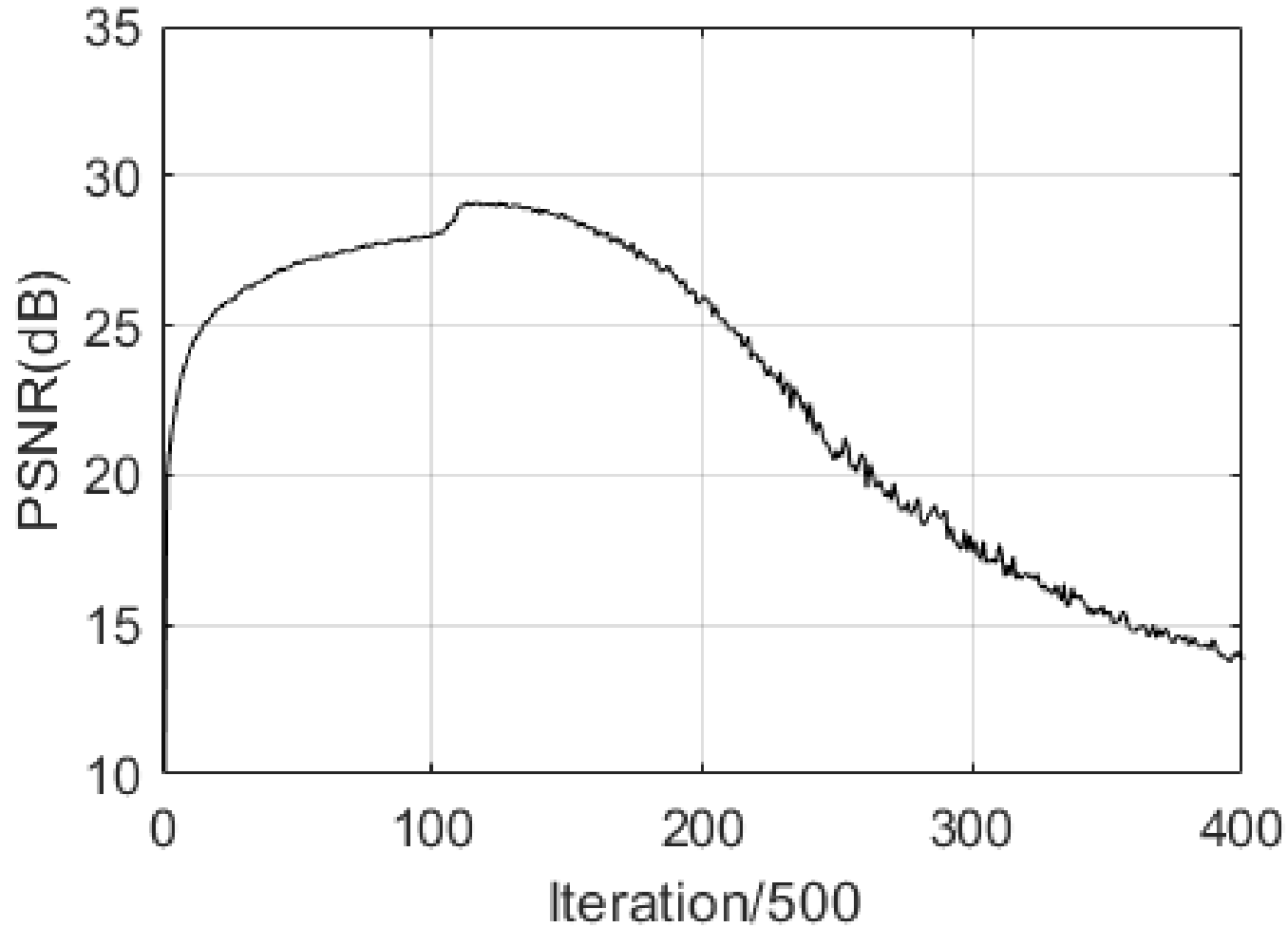
$$\sim 2 \text{ million} \quad \gg \quad \sim 0.1 \text{ million}$$
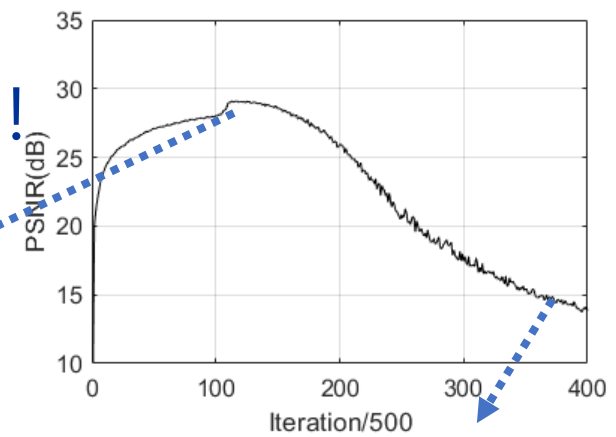
(#parameters in model) (DOF in an image)

**In principle, $f(\boldsymbol{\theta})$ can generate any image!**
(i.e., not only clean, but also corrupted images)

# Over-Parameterization -> Overfitting!

# Over-Parameterization -> Overfitting!



Early termination solution
(impractical!)

Global solution: $f(\boldsymbol{\theta}) \approx \boldsymbol{y}$
(overfitting!)

# This Work: Over-Parameterization *Without* Overfitting!

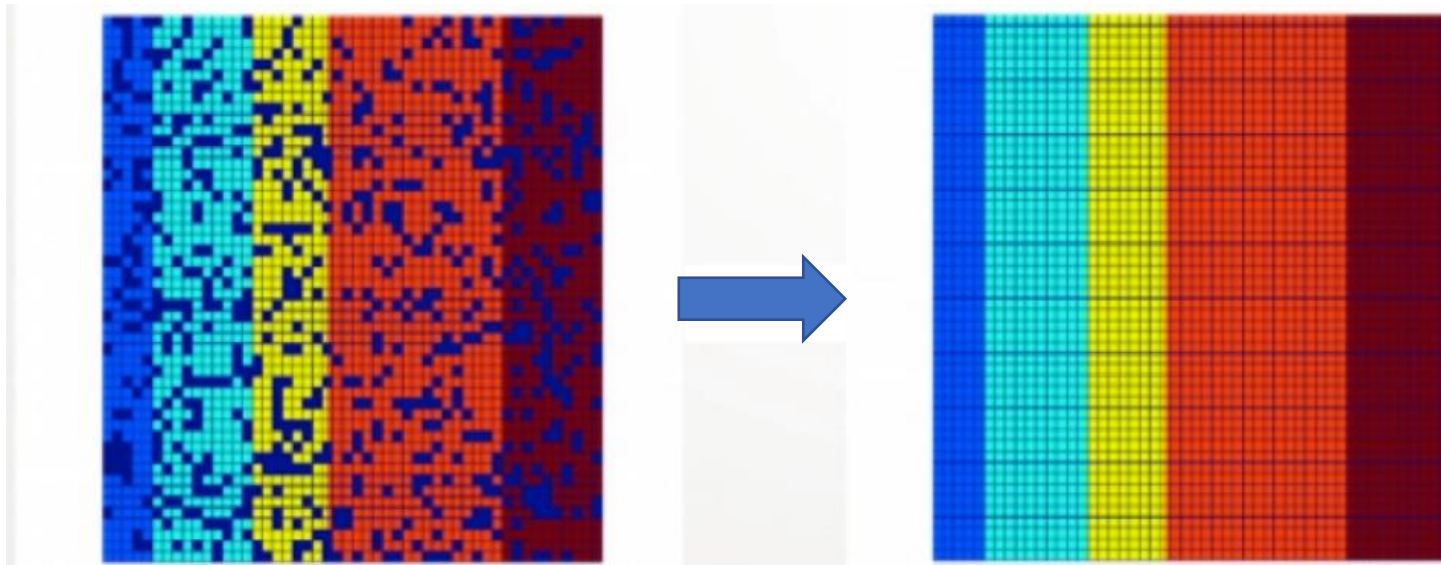Robust recovery of natural images

Simplification          Implication

Robust recovery of low-rank matrices

# Robust Recovery of Low-Rank Matrices

- **Goal**: Recover a rank-$r$ matrix $\mathbf{X}_\star \in \mathbb{R}^{n \times n}$ from (possibly corrupted) linear measurements $\mathbf{y} = \mathcal{A}(\mathbf{X}_\star) + \mathbf{s}_\star$

# Matrix Factorization – *Exact* Parameterization

| | (Classical) Exact-Parameterization |
|---|---|
| $\mathbf{s}_\star = 0$ (noiseless) | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2$ <br><br> − Gradient descent finds $\mathbf{X}_\star$ <br><br> (Classical) robust loss <br><br> $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_1$ |
| sparse $\mathbf{s}_\star$ | − Sub-grad method finds $\mathbf{X}_\star$ <br><br> Require knowing $r = \operatorname{rank}(\boldsymbol{X}_\star)$! |

- Li, Zhu, So, Vidal, Nonconvex Robust Low-rank Matrix Recovery, SIAM Journal on Optimization 2019
- Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro, Implicit Regularization in Matrix Factorization, NeurIPS 2017

# Matrix Factorization – *Over* Parameterization

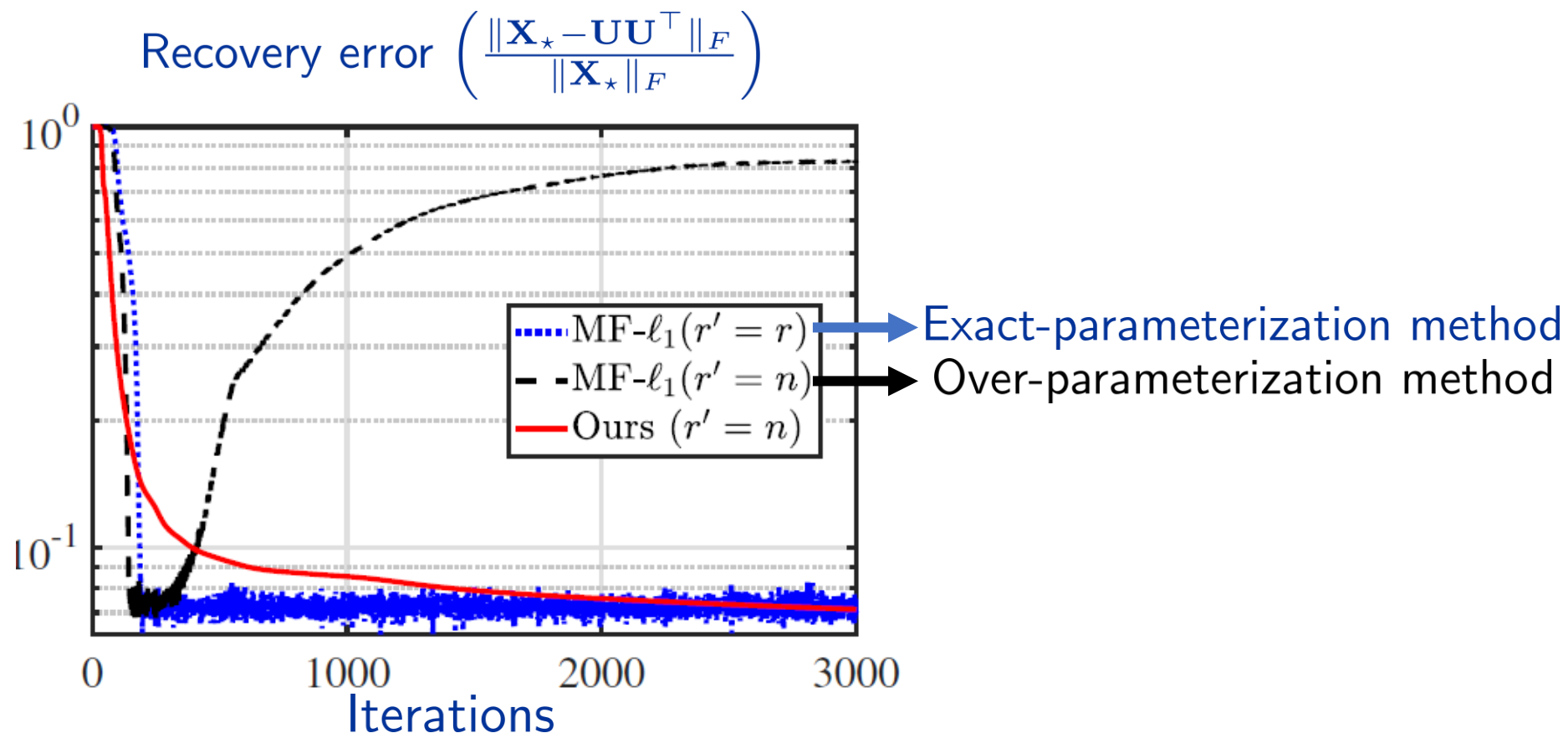|  | (Classical)<br>Exact-Parameterization | (Modern)<br>Over-Parameterization |
|---|---|---|
| $\mathbf{s}_\star = 0$<br>(noiseless) | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2$<br><br>– Gradient descent finds $\mathbf{X}_\star$ | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times n}} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2$<br><br>– Gradient descent finds $\mathbf{X}_\star$ |
| sparse $\mathbf{s}_\star$ | (Classical)<br>robust loss<br><br>$\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_1$<br><br>– Sub-grad method finds $\mathbf{X}_\star$<br>Require knowing $r = \mathrm{rank}(\boldsymbol{X}_\star)$! | ? |

- Li, Zhu, So, Vidal, Nonconvex Robust Low-rank Matrix Recovery, SIAM Journal on Optimization 2019
- Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro, Implicit Regularization in Matrix Factorization, NeurIPS 2017

# Matrix Factorization – *Over* Parameterization

|  | (Classical)<br>Exact-Parameterization | (Modern)<br>Over-Parameterization |
|---|---|---|
| $\mathbf{s}_\star = 0$<br>(noiseless) | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}}\|\mathbf{y}-\mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2$<br><br>– Gradient descent finds $\mathbf{X}_\star$ | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times n}}\|\mathbf{y}-\mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_2^2$<br><br>– Gradient descent finds $\mathbf{X}_\star$ |
| sparse $\mathbf{s}_\star$ | $\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times r}}\|\mathbf{y}-\mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_1$<br><br>– Sub-grad method finds $\mathbf{X}_\star$ | $\boxed{\displaystyle\min_{\mathbf{U}\in\mathbb{R}^{n\times n}}\|\mathbf{y}-\mathcal{A}(\mathbf{U}\mathbf{U}^\top)\|_1}$<br><br>– **Does this work?** |

(Classical) robust loss

Use classical approach?

Require knowing $r = \mathrm{rank}(\boldsymbol{X}_\star)$!

- Li, Zhu, So, Vidal, Nonconvex Robust Low-rank Matrix Recovery, SIAM Journal on Optimization 2019
- Gunasekar, Woodworth, Bhojanapalli, Neyshabur, Srebro, Implicit Regularization in Matrix Factorization, NeurIPS 2017

# *Failure* of Classical Approach for Over-Param. Models



Recovery error $\left( \frac{\|\mathbf{X}_\star - \mathbf{U}\mathbf{U}^\top\|_F}{\|\mathbf{X}_\star\|_F} \right)$

Legend:
- MF-$\ell_1(r' = r)$
- MF-$\ell_1(r' = n)$
- Ours $(r' = n)$

Exact-parameterization method
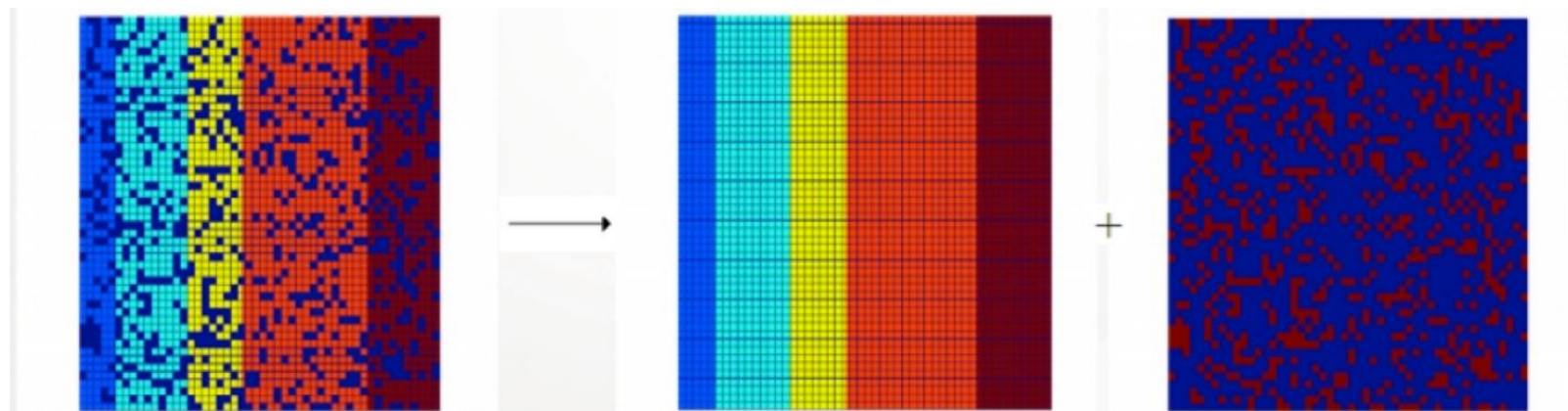Over-parameterization method

Iterations

Key Challenge: How to robustify over-parameterized models?

# A Double Over-Parameterization Method

- **Our Strategy:** Over-parameterize the noise $s$

Hadamard product

$$\min_{U \in \mathbb{R}^{n \times n}, g, h} \ell(U, g, h) := \| \underbrace{y} - ( \mathcal{A} \underbrace{(UU^\top)} + \underbrace{g \odot g - h \odot h} )\|_2^2$$
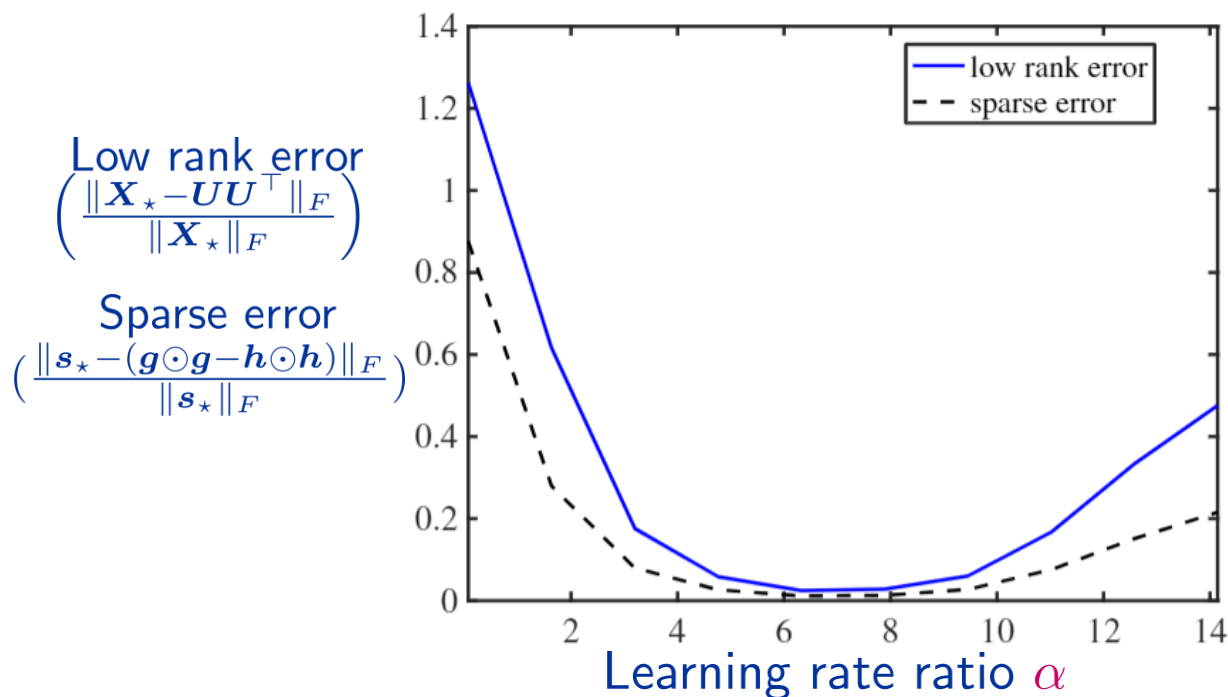
over-parameterizes $X$ \quad over-parameterizes $s$



- (Seemingly) Problematic: Even more parameters, infinitely many global solutions

# Implicit Algorithmic Bias

- A gradient descent with discrepant learning rates algorithm:

$$\boldsymbol{U}_{k+1} = \boldsymbol{U}_k - \tau \cdot \nabla\ell(\boldsymbol{U}_k, \boldsymbol{g}_k, \boldsymbol{h}_k)$$

$$\begin{bmatrix} \boldsymbol{g}_{k+1} \\ \boldsymbol{h}_{k+1} \end{bmatrix} = \begin{bmatrix} \boldsymbol{g}_k \\ \boldsymbol{h}_k \end{bmatrix} - \alpha \cdot \tau \cdot \begin{bmatrix} \nabla_{\boldsymbol{g}}\ell(\boldsymbol{U}_k, \boldsymbol{g}_k, \boldsymbol{h}_k) \\ \nabla_{\boldsymbol{h}}\ell(\boldsymbol{U}_k, \boldsymbol{g}_k, \boldsymbol{h}_k) \end{bmatrix}$$

- Implicit bias of learning rate ratio $\alpha$

Low rank error
$\left( \frac{\|\boldsymbol{X}_\star - \boldsymbol{U}\boldsymbol{U}^\top\|_F}{\|\boldsymbol{X}_\star\|_F} \right)$

Sparse error
$\left( \frac{\|\boldsymbol{s}_\star - (\boldsymbol{g}\odot\boldsymbol{g} - \boldsymbol{h}\odot\boldsymbol{h})\|_F}{\|\boldsymbol{s}_\star\|_F} \right)$



Learning rate ratio $\alpha$

# Implicit Algorithmic Bias: Main Theory

- **Q1:** How does $\alpha$ (learning rate ratio) control the solution?

**Theorem: [You*, Zhu*, Qu, Ma'2020]**

- $U_0, g_0, h_0$ are infinitesimally small, $\tau$ is infinitesimally small
- $\mathcal{A}$ is symmetric and commutable (i.e., $A_i A_j = A_j A_i$ for each $i \neq j$)
- $X_\infty = \lim_{k \to \infty} U_k U_k^\top$ and $s_\infty = \lim_{k \to \infty} (g_k \odot g_k - h_k \odot h_k)$ exist and produce a global solution (i.e., $y = \mathcal{A}(X_\infty) + s_\infty$),

then $(X_\infty, s_\infty)$ is a solution to the following problem

$$\min_{X, s} \|X\|_* + \frac{1}{\alpha} \|s\|_1 \quad \text{s.t.} \quad y = \mathcal{A}(X) + s$$

- **Q2:** What value of $\alpha$ should I use?

  **Theorem [Candes et al. '09]** Using $\alpha = \sqrt{n}$, the solution to the convex optimization above is $(X_\star, s_\star)$ under certain conditions

  – No parameter tuning is required
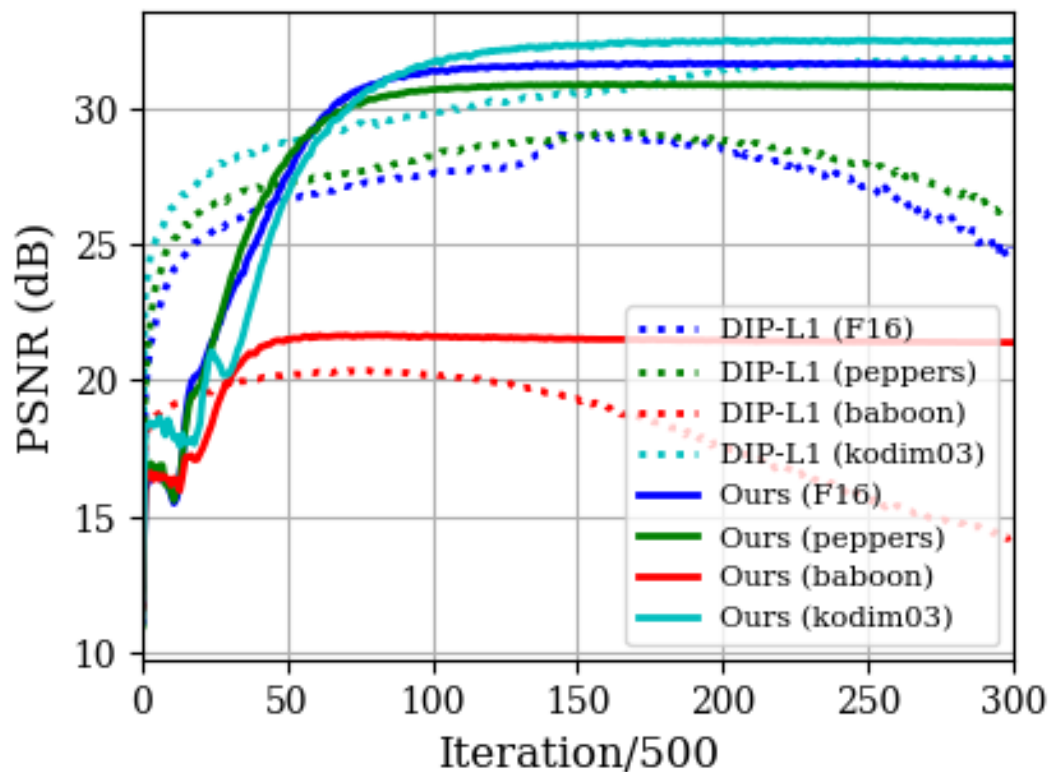
# Extension to Image Recovery

- **Goal:** Recover an image $\mathbf{X}_\star \in \mathbb{R}^{H \times W \times 3}$ from $\mathbf{y} = \mathbf{X}_\star + \mathbf{s}_\star$, where $\mathbf{s}_\star$ is salt-and-pepper (sparse) noise
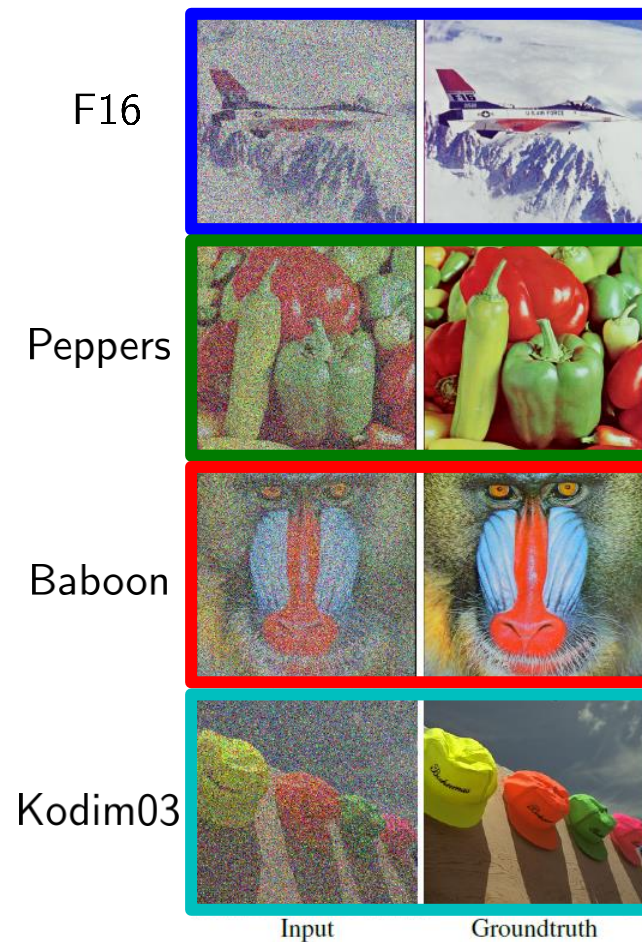
- **Method:**

$$\min_{\boldsymbol{\theta}, \boldsymbol{g}, \boldsymbol{h}} \| \underbrace{\mathbf{y}} - ( \underbrace{f(\boldsymbol{\theta})} + \underbrace{\boldsymbol{g} \odot \boldsymbol{g} - \boldsymbol{h} \odot \boldsymbol{h}} )\|_2^2$$



  – Solve this problem using gradient descent with learning rate ratio $\alpha$

- No parameter tuning is required:
  – No theory for optimal $\alpha$, but $\alpha = 500$ works well for all cases (i.e., with different images/corruption levels/network widths)

# Results for Varying Images



No early stop, no parameter tuning!

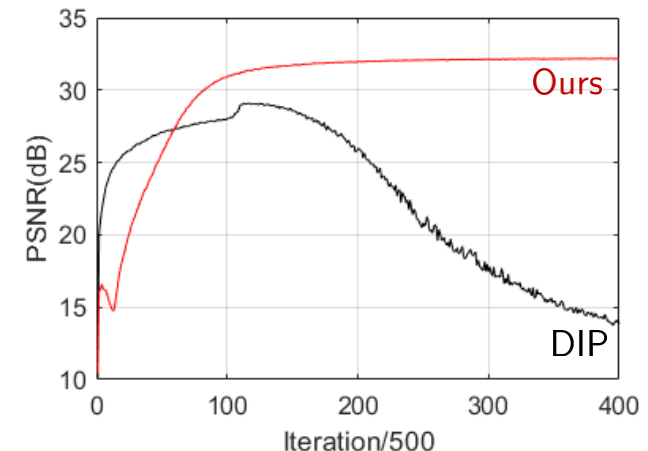# Take-Home: Over-parameterization without Overfitting!

- Classical robust loss function approaches will overfit for over-parameterized models

- This work proposes a

  *double over-parameterization model*

  which does not overfit by exploiting

  *implicit bias of discrepant learning rates*

# Thank you for your attention!

Acknowledgement: