

# Dropping convexity for faster semi-definite optimization

Srinadh Bhojanapalli, Anastasios Kyrillidis and Sujay Sanghavi



## [ Semi-definite optimization ]

We consider the following optimization problem:

$$(\text{SDP}) \quad \min_{X \succeq 0} f(X)$$

### Main assumptions:

- $X$  is a symmetric, positive semi-definite in  $\mathbb{R}^{n \times n}$ .
- $f : \mathbb{S}_+^n \rightarrow \mathbb{R}$  has Lipschitz continuous gradients:
$$\|\nabla f(X) - \nabla f(Y)\|_F \leq M \cdot \|X - Y\|_F, \quad X, Y \in \mathbb{S}_+^n$$
- $f$  might also be (restricted) strongly convex:
$$f(Y) \geq f(X) + \langle \nabla f(X), Y - X \rangle + \frac{m}{2} \|Y - X\|_F^2, \quad X, Y \in \mathbb{S}_+^n$$

## [ Common practice: low rankness ]

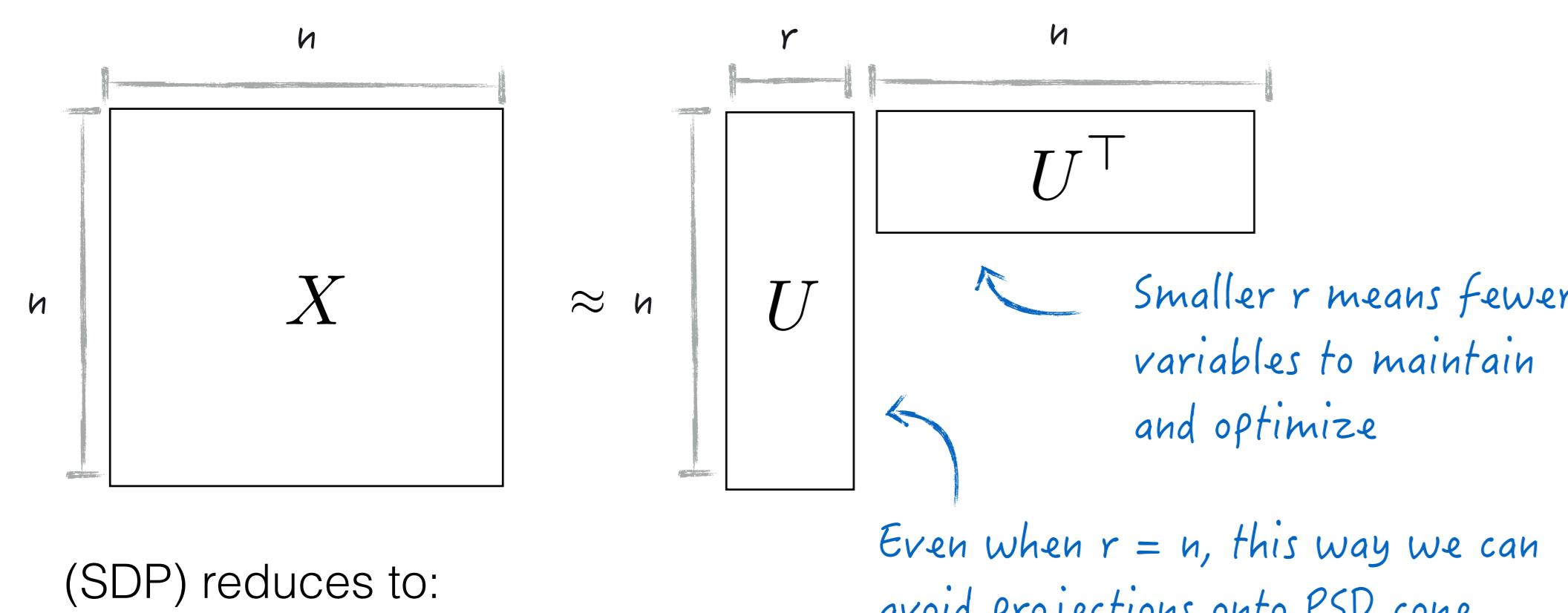
In modern machine learning tasks, one is interested in finding/learning a **low rank** matrix (Occam's razor):

- Low-rankness models better the underlying task.
- Low-rankness results into better computational complexity (fewer variables to optimize).
- Low-rankness might prevent over-fitting in machine learning problems.

**Examples:** matrix completion, phase retrieval, quantum state tomography, etc.

## [ Low rankness + PSD constraint ]

Any PSD matrix  $X \in \mathbb{R}^{n \times n}$  of rank  $r$  can be written as  $X = UU^\top$ , i.e.,



(SDP) reduces to:

$$(\text{Non-convex SDP}) \quad \min_{U \in \mathbb{R}^{n \times r}} g(U)$$

where  $r \leq n$

## [ Contributions ]

- Factored Gradient Descent** (FGD) algorithm:  $U^+ = U - \eta \cdot \nabla f(UU^\top) \cdot U$  Non-convex!!!
- Novel** step size selection  $\eta$ .
- (Local) convergence of FGD under common assumptions:
  - Linear convergence** for smooth and strongly convex  $f$ .
  - Sublinear** convergence for just smooth  $f$ .
- Global convergence after proper initialization.

## [ Factored gradient descent (FGD) ]

[ Algorithm: Factored gradient descent ]

**Input:** function  $f$ , target rank  $r$ , # of iterations  $K$ .

- Compute  $X^0$  (see below).
- Set  $U \in \mathbb{R}^{n \times r}$  such that  $X^0 = UU^\top$
- Set step size  $\eta$ .
- For  $k = 0$  to  $K - 1$

Compute

$$U^+ = U - \eta \cdot \nabla f(UU^\top) \cdot U$$

Set  $U = U^+$

**Output:**  $X = UU^\top$

## Step size selection

$$\eta = \frac{C}{M\|X^0\|_2 + \|\nabla f(X^0)\|_2}$$

for some constant  $C > 0$ .

## Intuition for $\eta$ selection.

- Let  $f$  be a separable function:

$$f(X) = \sum_{ij} f_{ij}(X_{ij})$$

- For  $r = 1$ , we have  $X = uu^\top$ .
- Define  $f(uu^\top) \equiv g(u)$ . Then:

$$\nabla g(u) = \nabla f(uu^\top) \cdot u$$

$$\nabla^2 g(u) = \text{mat}(\text{diag}(\nabla^2 f(uu^\top)) \cdot \text{vec}(uu^\top) + \nabla f(uu^\top))$$

- Convex optimization suggests:

$$\eta < \frac{1}{\|\nabla^2 g(\cdot)\|_2} \propto \frac{1}{M\|X\|_2 + \|\nabla f(X)\|_2}$$

## Initialization

- Option #1:** run standard convex algorithms for a few iterations.

Let  $X^+ = \mathcal{P}_+(X - \frac{1}{M}\nabla f(X))$  be the proj. gradient descent update. Then,  $\|X^+ - X\|_F \leq \frac{c}{\sqrt{r}\tau(X_r)} \cdot \sigma_r(X)$  implies:

$$\text{DIST}(U_r, U_r^*) \leq \frac{c'}{\tau(X_r^*)} \sigma_r(U_r^*), \quad \text{for constants } c, c' > 0.$$

- Option #2:** Set  $X^0 := \frac{1}{\|\nabla f(0) - \nabla f(e_1 e_1^\top)\|_F} \mathcal{P}_+(-\nabla f(0))$ . Then,

$$\text{DIST}(U_r^0, U_r^*) \leq 4\sqrt{2}r\tau(X_r^*) \cdot \sqrt{\kappa^2 - 2/\kappa + 1} \cdot \sigma_r(U_r^*)$$

## [ Theoretical guarantees ]

- Define the following *distance metric*:

$$\text{DIST}(U, V) := \min_{R: R \in \mathcal{O}} \|U - VR\|_F.$$

between matrices  $U, V \in \mathbb{R}^{n \times r}$ . Here,  $R$  is restricted such that  $R^\top R = I_{r \times r}$ .

## FGD Assumptions #1

- FGD is initialized with a “good” starting point  $X^0 = U^0(U^0)^\top$  such that:

$$\text{DIST}(U^0, U_r^*) \leq \rho \sigma_r(U_r^*) \quad \text{for } \rho := \frac{1}{100} \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \stackrel{\text{smooth}}{\text{smooth}}$$

$$\text{DIST}(U^0, U_r^*) \leq \rho' \sigma_r(U_r^*) \quad \text{for } \rho' := \frac{1}{100\kappa} \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \stackrel{\text{strongly convex}}{\text{strongly convex}}$$

## $\frac{1}{k}$ -convergence rate for **smooth** $f$

After  $k$  iterations, FGD algorithm finds a solution  $X^k$  such that:

$$f(X^k) - f(X_r^*) \leq \frac{\frac{6}{\eta^*} \cdot \text{DIST}(U^0, U_r^*)^2}{k + \frac{6}{\eta^*} \cdot \frac{\text{DIST}(U^0, U_r^*)^2}{f(X^0) - f(X_r^*)}}.$$

Here, we assume  $X^* = X_r^*$  and define:

$$\eta^* = \frac{1}{M\|X^*\|_2 + \|\nabla f(X^*)\|_2}$$

## Linear rate under **strong convexity**

Per iteration, FGD algorithm satisfies the following iteration invariant:

$$\text{DIST}(U^+, U^*)^2 \leq \alpha \cdot \text{DIST}(U, U^*)^2.$$

where

$$\alpha = 1 - \frac{m\sigma_r(X^*)}{64(M\|X^*\|_2 + \|\nabla f(X^*)\|_2)}$$

Furthermore:

$$\text{DIST}(U^+, U_r^*) \leq \rho' \sigma_r(U_r^*).$$

- For cases where  $\text{rank}(X^*) > r$

## FGD Assumptions #2

- Assuming

$$\|X^* - X_r^*\|_F \leq \frac{1}{100} \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \sigma_r(X^*) \stackrel{\text{smooth}}{\text{smooth}}$$

$$\|X^* - X_r^*\|_F \leq \frac{1}{200\kappa^{1.5}} \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \sigma_r(X^*) \stackrel{\text{strongly convex}}{\text{strongly convex}}$$

we further have

$$\text{DIST}(U^+, U^*)^2 \leq \alpha \cdot \text{DIST}(U, U^*)^2 + \beta \cdot \|X^* - X_r^*\|_F^2$$

for some (controllable)  $\beta$ .

## [ Related work ]

Reference	Conv. rate	Initialization	Output rank
[Haz08]	$1/\epsilon$ (Smooth $f$ )	$X^0 = 0$	$1/\epsilon$
[Lau12]	$1/\epsilon$ (Smooth $f$ )	$X^0 = 0$	$1/\epsilon$
[CW15]	$1/\epsilon$ (Local Asm.)	Application dependent	$r$
[CW15]	$\log(1/\epsilon)$ (Local Asm.)	Application dependent	$r$
This work	$1/\epsilon$ (Smooth $f$ )	SVD / top- $r$	$r$
This work	$\log(1/\epsilon)$ (Smooth, RSC $f$ )	SVD / top- $r$	$r$

...and for the affine rank minimization problem:

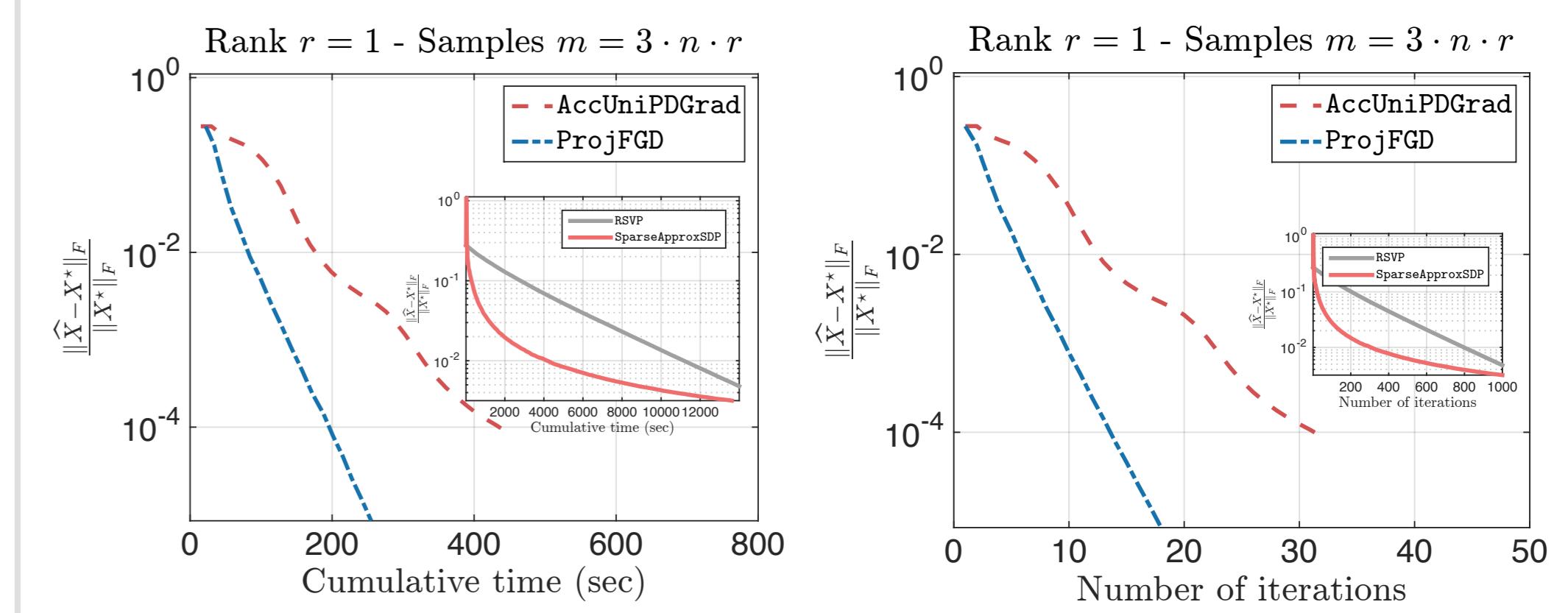
Reference	$\text{DIST}(U^+, U^*)^2 \leq \alpha \cdot \text{DIST}(U, U^*)^2$	$\text{DIST}(U^0, U^*) \leq \dots$
[JNS13]	$\alpha = \frac{1}{16}$	$\sqrt{6\delta} \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \sigma_r(U_r^*)$
[TBSR15]	$\alpha = 1 - \frac{c_1}{\tau(U_r^*)^4}$	$\frac{c_1}{4} \sigma_r(U_r^*)$
[ZL15]	$\alpha = 1 - \frac{c_2}{\tau(U_r^*)^4 \cdot r}$	$\sqrt{\frac{3}{16}} \sigma_r(U_r^*)$
[CW15]	$\alpha = 1 - \frac{c_3}{\tau(U_r^*)^{10}}$	$(1 - \tau) \cdot \sigma_r(U_r^*)$
This work	$\alpha = 1 - \frac{c_4}{\tau(U_r^*)^2}$	$\frac{1}{100} \cdot \frac{\sigma_r(X^*)}{\sigma_1(X^*)} \sigma_r(U_r^*)$

## [ In practice... ]

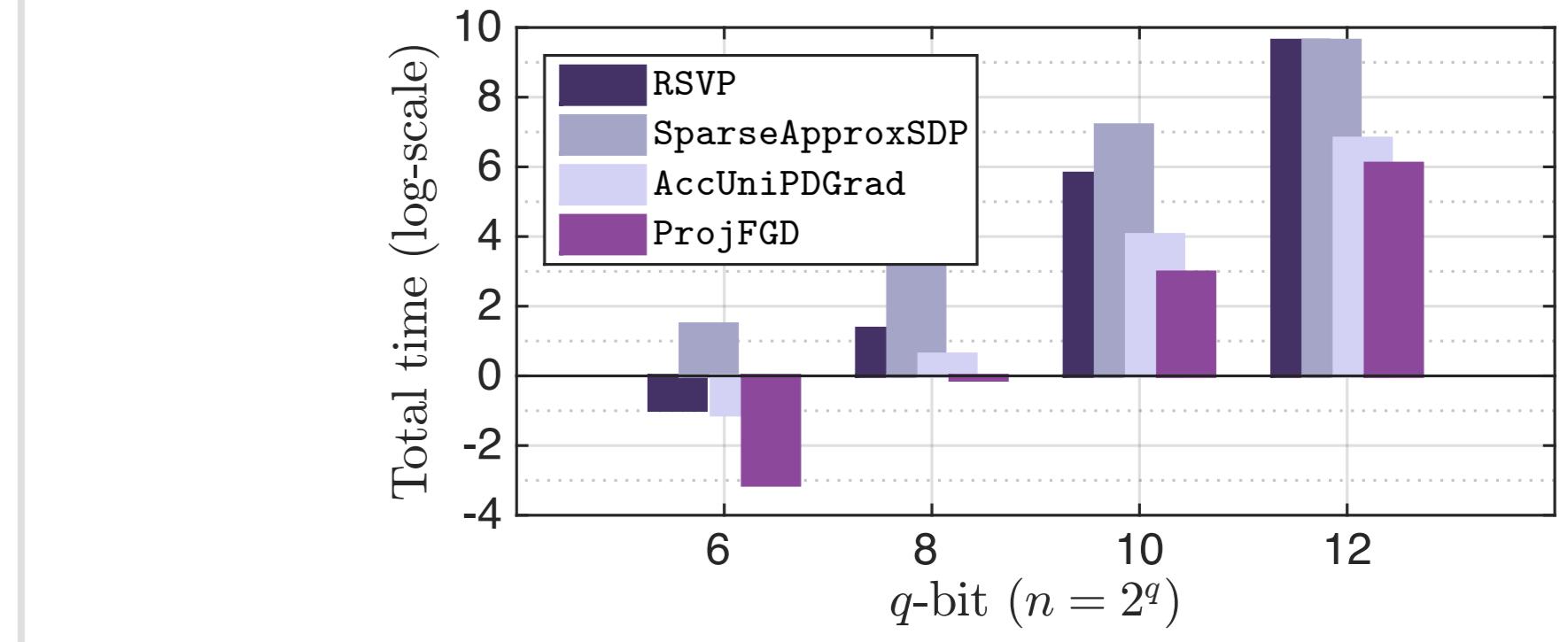
### Quantum state tomography:

$$\begin{aligned} \text{minimize}_{X \in \mathbb{R}^{n \times n}} \quad & f(X) \\ \text{subject to} \quad & \text{rank}(X) = r \\ & X \succeq 0, \text{TR}(X) \leq 1. \end{aligned} \quad \xrightarrow{\hspace{1cm}} \quad \begin{aligned} \text{minimize}_{U \in \mathbb{R}^{n \times r}} \quad & f(UU^\top) \\ \text{subject to} \quad & \|U\|_F^2 \leq 1. \end{aligned}$$

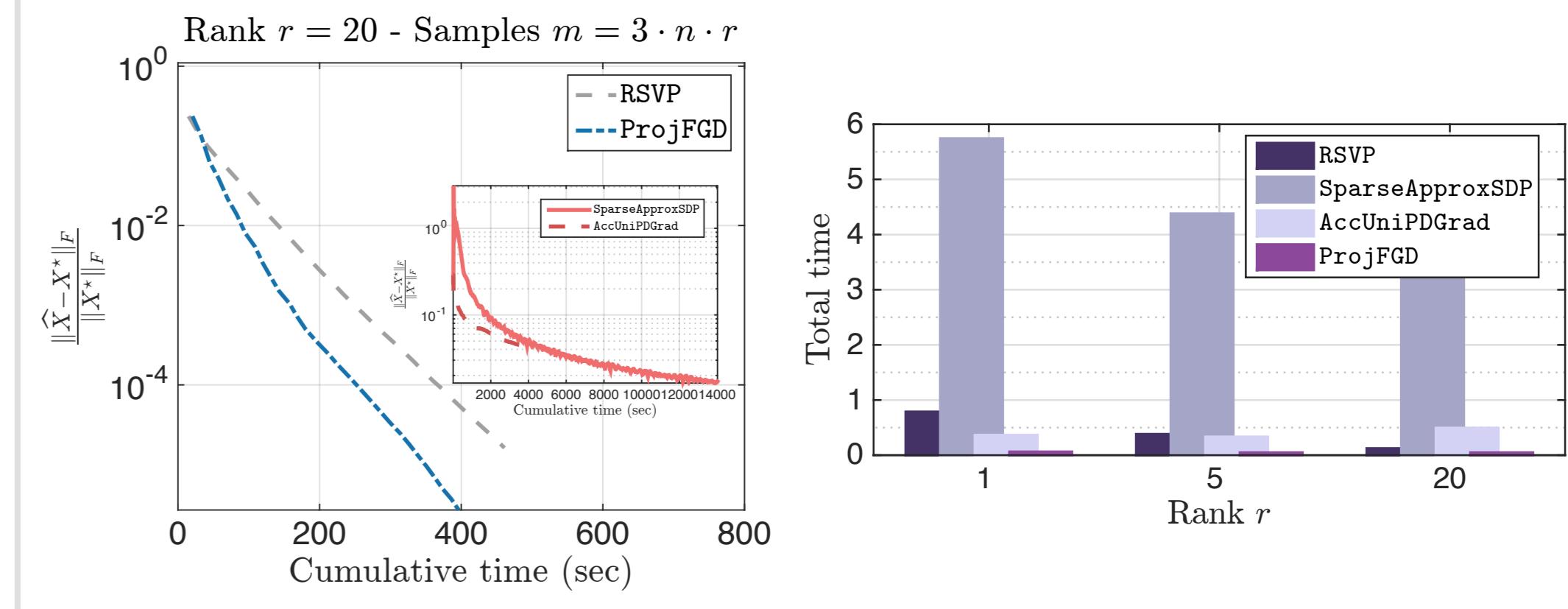
### Pure state density matrix recovery ( $r = 1$ )



### Faster convergence to a vicinity around $X^*$ :



### Almost pure state density matrix recovery ( $r > 1$ ; here, $r = \{5, 20\}$ )



## [ References ]

- Srinadh Bhojanapalli, Anastasios Kyrillidis and Sujay Sanghavi, “Dropping convexity for faster semi-definite optimization”, Arxiv preprint, submitted to STOC.