

# Recovering Low-Rank Matrices from Binary Measurements

---

Simon Foucart\* and Richard G. Lynch — Texas A&M University

---

## Abstract

This article studies the approximate recovery of low-rank matrices acquired through binary measurements. Two types of recovery algorithms are considered, one based on hard singular value thresholding and the other one based on semidefinite programming. In case no thresholds are introduced before binary quantization, it is first shown that the direction of the low-rank matrices can be well approximated. Then, in case nonadaptive thresholds are incorporated, it is shown that both the direction and the magnitude can be well approximated. Finally, by allowing the thresholds to be chosen adaptively, we exhibit a recovery procedure for which low-rank matrices are fully approximated with error decaying exponentially with the number of binary measurements. In all cases, the procedures are robust to prequantization error. The underlying arguments are essentially deterministic: they rely only on an unusual restricted isometry property of the measurement process, which is established once and for all for Gaussian measurement processes.

*Key words and phrases:* low-rank recovery, one-bit compressive sensing, quantization, hard singular value thresholding, semidefinite programming, adaptivity.

*AMS classification:* 94A12, 65F10, 90C22.

---

## 1 Introduction

This paper considers the problem of recovering low-rank matrices  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  from  $m$  binary measurements  $y_1, \dots, y_m$ , each of them given as the sign of a linear functional applied to  $\mathbf{X}$ , i.e., as

$$(1) \quad y_i = \text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F), \quad \langle \mathbf{A}_i, \mathbf{X} \rangle_F := \text{tr}(\mathbf{A}_i^* \mathbf{X}) = \sum_{k=1}^{n_1} \sum_{\ell=1}^{n_2} (A_i)_{k,\ell} X_{k,\ell},$$

for some  $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n_1 \times n_2}$ . In short, we write

$$(2) \quad \mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X}) \in \{\pm 1\}^m,$$

---

\*S. F. is partially supported by the NSF under the grant DMS-1622134

where  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  is the linear map defined by  $(\mathcal{A}\mathbf{Z})_i = \sum_{k=1}^{n_1} \sum_{\ell=1}^{n_2} (A_i)_{k,\ell} Z_{k,\ell}$ ,  $i \in \llbracket 1 : m \rrbracket$ . We shall also allow for thresholds  $\tau_1, \dots, \tau_m \in \mathbb{R}$  to be incorporated before binary quantization, leading to

$$(3) \quad \mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X} - \boldsymbol{\tau}) = [\text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F - \tau_i)]_{i=1}^m \in \{\pm 1\}^m.$$

This scenario of low-rank matrix recovery from binary measurements is a natural extension to the scenario of the sparse vector recovery from binary measurements (and of one-bit compressive sensing) and the techniques used here are adapted from the ones used there. In fact, general results on recovery of structured signals from nonlinear observations, such as the ones from [10], are already informative when particularized to our situation. But our contribution to the specific setting of low-rank matrix recovery from binary measurements goes further in several directions:

- the recovery procedure is not a generalized Lasso, instead it is either a hard singular value thresholding algorithm or a semidefinite program in the spirit of [7] (rather than of [8]);
- it can be enhanced to estimate not only the direction of low-rank matrices but also their magnitudes, thanks to the presence of thresholds;
- it can handle adversarial prequantization error;
- perhaps most importantly, probabilistic arguments and deterministic arguments are separated owing to an unusual restricted isometry property for the measurement map  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$  (this parallels simplified arguments put forward in [5, Section 8.4] for one-bit compressive sensing, which incidentally enables the removal of a logarithmic factor from the number of binary measurements obtained by adapting directly the techniques of [7]).

In Section 2, we will first establish the main theoretical tool at the basis of all the forthcoming arguments. In Section 3, we will prove that the direction of low-rank matrices  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  can be well estimated for measurements of the type (2) using reconstruction procedures based on hard singular value thresholding and on semidefinite programming. In Section 4, we will prove that it is also possible, if we allow thresholds in the binary measurements as in (3), to fully estimate the low-rank matrices  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$  using related reconstruction procedures. In Section 5, we will further prove that an adaptive choice of the thresholds can reduce the reconstruction error drastically. Finally, we report in Section 6 on some modest numerical experiments.

But before launching into the technicalities, we remark that we can, and from now on we will, assume that  $n_1 = n_2 =: n$ . Indeed, if  $n_1 > n_2$ , say, then instead of recovering the rectangular matrix  $\mathbf{X} \in \mathbb{R}^{n_1 \times n_2}$ , we would aim at recovering the square matrix  $\mathbf{X}' = \begin{bmatrix} \mathbf{X} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{n_1 \times n_1}$ , which has the same rank as  $\mathbf{X}$  and is in one-to-one correspondence with  $\mathbf{X}$ . Clearly, a suitable measurement map  $\mathcal{A}' : \mathbb{R}^{n_1 \times n_1} \rightarrow \mathbb{R}^m$  induces a suitable measurement map  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$ .

## 2 A Restricted Isometry Property

Our deterministic arguments rely on just one technical property of the map  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ , which is a version of the rank restricted isometry property modified to feature the  $\ell_1$ -norm as the inner norm. Precisely, we shall say that  $\mathcal{A}$  satisfies the modified rank restricted isometry property of order  $r$  with constant  $\delta \in (0, 1)$  — in short,  $\mathcal{A}$  satisfies  $\text{MRRIP}(r, \delta)$  — if

$$(4) \quad (1 - \delta)\|\mathbf{Z}\|_F \leq \|\mathcal{A}(\mathbf{Z})\|_1 \leq (1 + \delta)\|\mathbf{Z}\|_F \quad \text{whenever } \text{rank}(\mathbf{Z}) \leq r.$$

As a matter of fact, we will sometimes need this property to be valid not only for matrices  $\mathbf{Z} \in \mathbb{R}^{n \times n}$  of rank at most  $r$ , but also for matrices  $\mathbf{Z} \in \mathbb{R}^{n \times n}$  of effective rank  $\text{effrank}(\mathbf{Z}) := \|\mathbf{Z}\|_* / \|\mathbf{Z}\|_F$  at most  $r$ . Here,  $\|\mathbf{Z}\|_*$  represents the nuclear norm (aka trace norm or Schatten 1-norm) of  $\mathbf{Z}$ , i.e., the sum of its singular values. Thus, we shall write that  $\mathcal{A}$  satisfies  $\text{MRRIP}^{\text{eff}}(r, \delta)$  if

$$(5) \quad (1 - \delta)\|\mathbf{Z}\|_F \leq \|\mathcal{A}(\mathbf{Z})\|_1 \leq (1 + \delta)\|\mathbf{Z}\|_F \quad \text{whenever } \text{effrank}(\mathbf{Z}) \leq r.$$

The main result of this section states that properly normalized Gaussian random linear maps satisfy the modified restricted isometry property for genuinely and effectively low-rank matrices with high probability, as established below.

**Theorem 1.** There exist absolute constants  $c, C > 0$  such that, if  $m \geq C\delta^{-3}nr$ , resp.  $m \geq C\delta^{-5}nr$ , then, with failure probability at most  $2\exp(-c\delta^2m)$ , a random measurement map  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$  for which the  $(A_i)_{k,\ell}$  are independent Gaussian random variables with mean zero and standard deviation  $\sqrt{\pi/2}/m$  satisfies  $\text{MRRIP}(r, \delta)$ , resp.  $\text{MRRIP}^{\text{eff}}(r, \delta)$ .

**Remark.** We have opted to include a low-key proof that does not give optimal powers of  $\delta^{-1}$ . However, specifying results from [9] or [11], combined with mean width estimations from [8], would yield the same conclusions under the weaker assumption  $m \geq C\delta^{-2}nr$ .

*Proof.* For a fixed  $\mathbf{Z} \in \mathbb{R}^{n \times n}$ , it is easy to notice that

$$(6) \quad \mathbb{E}(\|\mathcal{A}(\mathbf{Z})\|_1) = \|\mathbf{Z}\|_F.$$

Indeed, for all  $i \in \llbracket 1 : m \rrbracket$ , we see that  $\mathcal{A}(\mathbf{Z})_i = \sum_{k,\ell=1}^n (A_i)_{k,\ell} Z_{k,\ell} = (\sqrt{\pi/2}/m)\|\mathbf{Z}\|_F g_i$ , where the  $g_i$  are independent standard normal random variables, and as such have first absolute moment equal to  $\sqrt{2/\pi}$ . It is also easy to notice that the map  $f : (\mathbf{B}_1, \dots, \mathbf{B}_m) \mapsto \sum_{i=1}^m \left| \sum_{k,\ell=1}^n \frac{\sqrt{\pi/2}}{m} (B_i)_{k,\ell} Z_{k,\ell} \right|$  is  $L$ -Lipschitz with  $L := \frac{\sqrt{\pi/2}}{\sqrt{m}}\|\mathbf{Z}\|_F$ . By concentration of measure (see e.g. [6, Theorem 8.34 or Theorem 8.40]), it follows that

$$(7) \quad \mathbb{P}(|\|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F| > \varepsilon \|\mathbf{Z}\|_F) \leq 2 \exp\left(-\frac{(\varepsilon \|\mathbf{Z}\|_F)^2}{2L^2}\right) = 2 \exp\left(-\frac{m\varepsilon^2}{\pi}\right).$$

First, to deal with genuine low rank, we consider a  $\rho$ -net  $\{\mathbf{Z}_1, \dots, \mathbf{Z}_K\}$  for the set

$$(8) \quad \mathcal{S}_r := \{\mathbf{Z} \in \mathbb{R}^{n \times n} : \|\mathbf{Z}\|_F \leq 1, \text{rank}(\mathbf{Z}) \leq r\}.$$

According to [4, Lemma 3.1], we can take  $K \leq (1 + 6/\rho)^{(2n+1)r}$ , hence  $K \leq \exp(18nr/\rho)$ . We place ourselves in the situation where

$$(9) \quad \left| \|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F \right| \leq \varepsilon \|\mathbf{Z}_k\|_F \leq \varepsilon \quad \text{for all } k = 1, \dots, K,$$

which, by (7) and a union bound, occurs with failure probability at most

$$(10) \quad K \times 2 \exp\left(-\frac{m\varepsilon^2}{\pi}\right) \leq 2 \exp\left(\frac{18nr}{\rho} - \frac{m\varepsilon^2}{\pi}\right).$$

Let us introduce the smallest constant  $\delta' \geq 0$  such that

$$(11) \quad \left| \|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F \right| \leq \delta' \|\mathbf{Z}\|_F \quad \text{for all } \mathbf{Z} \in \mathbb{R}^{n \times n} \text{ with } \text{rank}(\mathbf{Z}) \leq r.$$

Given  $\mathbf{Z} \in \mathcal{S}_r$ , we select  $k \in [1 : K]$  with  $\|\mathbf{Z} - \mathbf{Z}_k\|_F \leq \rho$  and we observe that

$$(12) \quad \begin{aligned} \left| \left| \|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F \right| - \left| \|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F \right| \right| &\leq \left| (\|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F) - (\|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F) \right| \\ &= \left| (\|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathcal{A}(\mathbf{Z}_k)\|_1) - (\|\mathbf{Z}\|_F - \|\mathbf{Z}_k\|_F) \right| \\ &\leq \|\mathcal{A}(\mathbf{Z} - \mathbf{Z}_k)\|_1 + \|\mathbf{Z} - \mathbf{Z}_k\|_F. \end{aligned}$$

Since  $\mathbf{Z} - \mathbf{Z}_k$  has rank at most  $2r$ , it can be written as  $\mathbf{Z} - \mathbf{Z}_k = \mathbf{M} + \mathbf{N}$  where both  $\mathbf{M}, \mathbf{N} \in \mathbb{R}^{n \times n}$  have rank at most  $r$  and  $\|\mathbf{Z} - \mathbf{Z}_k\|_F^2 = \|\mathbf{M}\|_F^2 + \|\mathbf{N}\|_F^2$ . Then

$$(13) \quad \begin{aligned} \|\mathcal{A}(\mathbf{Z} - \mathbf{Z}_k)\|_1 &\leq \|\mathcal{A}(\mathbf{M})\|_1 + \|\mathcal{A}(\mathbf{N})\|_1 \leq \delta' \|\mathbf{M}\|_F + \delta' \|\mathbf{N}\|_F \leq \delta' \sqrt{2} \sqrt{\|\mathbf{M}\|_F^2 + \|\mathbf{N}\|_F^2} \\ &= \sqrt{2} \delta' \|\mathbf{Z} - \mathbf{Z}_k\|_F, \end{aligned}$$

and in turn

$$(14) \quad \left| \left| \|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F \right| - \left| \|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F \right| \right| \leq (\sqrt{2} \delta' + 1) \|\mathbf{Z} - \mathbf{Z}_k\|_F \leq (\sqrt{2} \delta' + 1) \rho.$$

In view of (9) and (14), we therefore have

$$(15) \quad \begin{aligned} \left| \|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F \right| &\leq \left| \|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F \right| + \left| \left| \|\mathcal{A}(\mathbf{Z})\|_1 - \|\mathbf{Z}\|_F \right| - \left| \|\mathcal{A}(\mathbf{Z}_k)\|_1 - \|\mathbf{Z}_k\|_F \right| \right| \\ &\leq \varepsilon + (\sqrt{2} \delta' + 1) \rho. \end{aligned}$$

Taking the supremum over all  $\mathbf{Z}$  yields

$$(16) \quad \delta' \leq \varepsilon + (\sqrt{2} \delta' + 1) \rho, \quad \text{or} \quad \delta' \leq \frac{1}{1 - \sqrt{2} \rho} (\varepsilon + \rho).$$

Choosing  $\varepsilon = \frac{\delta}{2}$  and  $\rho = \frac{\sqrt{2}-1}{2}\delta \leq \frac{\sqrt{2}-1}{2}$  at the beginning ensures that

$$(17) \quad \delta' \leq \frac{1}{1 - (1 - 1/\sqrt{2})} \left( \frac{1}{2} + \frac{\sqrt{2}-1}{2} \right) \delta = \delta.$$

This means that  $\mathcal{A}$  satisfies MRRIP( $r, \delta$ ), as desired. Overall, the failure probability is at most

$$(18) \quad 2 \exp \left( \frac{36nr}{(\sqrt{2}-1)\delta} - \frac{m\delta^2}{4\pi} \right) \leq \exp \left( -\frac{m\delta^2}{8\pi} \right),$$

provided  $\frac{36nr}{(\sqrt{2}-1)\delta} \leq \frac{m\delta^2}{8\pi}$ , i.e.,  $m \geq \frac{288\pi}{\sqrt{2}-1}\delta^{-3}nr =: C\delta^{-3}nr$ .

Second, to deal with effective low rank, we claim that we can find a  $\rho$ -net for the set

$$(19) \quad \mathcal{S}_r^{\text{eff}} := \{\mathbf{Z} \in \mathbb{R}^{n \times n} : \|\mathbf{Z}\|_F \leq 1, \text{effrank}(\mathbf{Z}) \leq r\}$$

with cardinality at most  $\exp(72nr/\rho^3)$ , and the rest of the argument follows the lines of the genuine low-rank case so closely that it is not reproduced. In order to justify our covering number estimate, we start by considering a  $(\rho/2)$ -net  $\{\mathbf{Z}_1, \dots, \mathbf{Z}_K\}$  for  $\mathcal{S}_t$  with  $t := \lceil \rho^{-2}r \rceil \in [\rho^{-2}r, 2\rho^{-2}r]$  and with  $K \leq (1 + 6/(\rho/2))^{(2n+1)t} \leq \exp(72nr/\rho^3)$ . Then, given  $\mathbf{Z} \in \mathcal{S}_r^{\text{eff}}$ , we denote by  $\mathbf{Z}'$  its best rank- $t$  approximant (with respect to the Frobenius norm), so that, in view of [6, Theorem 2.5],

$$(20) \quad \|\mathbf{Z} - \mathbf{Z}'\|_F = \left[ \sum_{j=t+1}^n \sigma_j(\mathbf{Z}) \right]^{1/2} \leq \frac{1}{2\sqrt{t}} \sum_{i=1}^n \sigma_i(\mathbf{Z}) = \frac{1}{2\sqrt{t}} \|\mathbf{Z}\|_*.$$

Next, we select  $k \in \llbracket 1 : K \rrbracket$  such that  $\|\mathbf{Z}' - \mathbf{Z}_k\|_F \leq \rho/2$ . We then observe, since  $\mathbf{Z}$  is of effective rank  $r$ , that

$$(21) \quad \|\mathbf{Z} - \mathbf{Z}_k\|_F \leq \|\mathbf{Z} - \mathbf{Z}'\|_F + \|\mathbf{Z}' - \mathbf{Z}_k\|_F \leq \frac{1}{2\sqrt{t}} \|\mathbf{Z}\|_* + \frac{\rho}{2} \leq \frac{\sqrt{r}}{2\sqrt{t}} + \frac{\rho}{2} \leq \frac{\rho}{2} + \frac{\rho}{2} = \rho.$$

This shows that  $\{\mathbf{Z}_1, \dots, \mathbf{Z}_K\}$  is a  $\rho$ -net for  $\mathcal{S}_r^{\text{eff}}$ , hence it concludes the proof.  $\square$

We emphasize that, just like the standard rank restricted isometry property, this modified rank restricted isometry property holds with a number of measurements only proportional to  $nr$ , the number of ‘degrees of freedom’ of  $n \times n$  matrices of rank  $r$ . An extra logarithmic factor, which is necessary in sparse vector recovery, does not appear in this situation.

### 3 Estimating the Direction Only

In this section, we suppose that genuinely or effectively low-rank matrices  $\mathbf{X} \in \mathbb{R}^{n \times n}$  are acquired via  $\mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X}) \in \{\pm 1\}^m$ . In this setting, one can only hope to recover the direction of  $\mathbf{X}$ , as

all  $c\mathbf{X}$ ,  $c > 0$ , produce the same sign measurements. We study two reconstruction procedures, one based on hard singular value thresholding and one based on semidefinite programming. As we shall see, they are both robust to a prequantization error  $\mathbf{e} \in \mathbb{R}^m$  that corrupts the sign measurement vector to  $\mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})$ .

### 3.1 Hard singular value thresholding

The output of the hard singular value thresholding procedure is merely

$$(22) \quad \mathbf{X}^{\text{ht}} =: \text{best rank-}r \text{ approximant to } \mathcal{A}^*(\mathbf{y}) = \sum_{i=1}^m y_i \mathbf{A}_i.$$

The main theorem of this subsection is stated below. Note that, according to Theorem 1, its assumption is met with high probability by (properly normalized) Gaussian measurement maps  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ , provided  $m \geq C\delta^{-3}nr$  (or  $m \geq C\delta^{-2}nr$  if the remark after Theorem 1 is taken into account).

**Theorem 2.** Under MRRIP( $2r, \delta$ ), if  $\mathbf{X} \in \mathbb{R}^{n \times n}$  satisfying  $\text{rank}(\mathbf{X}) \leq r$  and  $\|\mathbf{X}\|_F = 1$  is acquired via  $\mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) \in \{\pm 1\}^m$ , then

$$(23) \quad \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{ht}}}{\|\mathbf{X}^{\text{ht}}\|_F} \right\|_F \leq \sqrt{80\delta} + 8\sqrt{\|\mathbf{e}\|_1}.$$

**Remark.** The square-root scaling of  $\|\mathbf{e}\|_1$  in Theorem 2 is peculiar. It is probably an artifact of the oversimplicity of the argument in Lemma 3 below. We believe that the recovery error should depend linearly on  $\|\mathbf{e}\|_1$ , which is supported by the numerical experiment presented in Subsection 6.1.

The whole argument is based on the one simple lemma stated below, which already shows that normalized low-rank matrices sharing the same sign measurements are necessarily close to one another. In this lemma, the operator  $P_{\mathcal{S}}$  associated to a subspace of  $\mathcal{S}$  of  $\mathbb{R}^{n \times n}$  denotes the orthogonal projector onto the space  $\mathcal{S}$  with respect to the Frobenius inner product.

**Lemma 3.** Under MRRIP( $r, \delta$ ), if  $\mathbf{X} \in \mathbb{R}^{n \times n}$  belongs to a space  $\mathcal{S} = \text{span}\{\mathbf{u}_1 \mathbf{v}_1^*, \dots, \mathbf{u}_r \mathbf{v}_r^*\}$  and satisfies  $\|\mathbf{X}\|_F = 1$ , then

$$(24) \quad \|\mathbf{X} - P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F^2 \leq 5\delta + 4\|\mathbf{e}\|_1.$$

*Proof.* By expanding the square in the left-hand side of (24), we have

$$(25) \quad \begin{aligned} & \|\mathbf{X} - P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F^2 \\ &= \|\mathbf{X}\|_F^2 - 2\langle \mathbf{X}, P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}))) \rangle_F + \|P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F^2. \end{aligned}$$

Firstly, we take into account that  $\|\mathbf{X}\|_F^2 = 1$ . Secondly, thanks to  $\mathbf{X} \in \mathcal{S}$ , we observe that

$$(26) \quad \begin{aligned} \langle \mathbf{X}, P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}))) \rangle_F &= \langle \mathbf{X}, \mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})) \rangle_F = \langle \mathcal{A}\mathbf{X}, \text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) \rangle \\ &= \|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1 - \langle \mathbf{e}, \text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) \rangle \geq \|\mathcal{A}\mathbf{X}\|_1 - 2\|\mathbf{e}\|_1 \\ &\geq (1 - \delta) - 2\|\mathbf{e}\|_1, \end{aligned}$$

where  $\text{MRRIP}(r, \delta)$  was used in the latter inequality. Thirdly, relying again on  $\text{MRRIP}(r, \delta)$ , we notice that

$$(27) \quad \begin{aligned} \|P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F^2 &= \langle \mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})), P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}))) \rangle_F \\ &= \langle \text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}), \mathcal{A}(P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))) \rangle \\ &\leq \|\mathcal{A}(P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}))))\|_1 \\ &\leq (1 + \delta)\|P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F, \end{aligned}$$

so that a simplification gives

$$(28) \quad \|P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F \leq (1 + \delta).$$

Putting these three facts together in (25) yields

$$(29) \quad \begin{aligned} \|\mathbf{X} - P_{\mathcal{S}}(\mathcal{A}^*(\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e})))\|_F^2 &\leq 1 - 2((1 - \delta) - 2\|\mathbf{e}\|_1) + (1 + \delta)^2 = 4\delta + \delta^2 + 4\|\mathbf{e}\|_1 \\ &\leq 5\delta + 4\|\mathbf{e}\|_1, \end{aligned}$$

which is the desired result.  $\square$

Now that this key lemma has been established, the main result about direction recovery by hard singular value thresholding can be easily deduced.

*Proof of Theorem 2.* Since  $\mathbf{X}^{\text{ht}}$  is the best rank- $r$  approximant to  $\mathcal{A}^*(\mathbf{y})$ , it is a better approximant than  $\mathbf{X}$  itself, so we have

$$(30) \quad \|\mathcal{A}^*(\mathbf{y}) - \mathbf{X}^{\text{ht}}\|_F^2 \leq \|\mathcal{A}^*(\mathbf{y}) - \mathbf{X}\|_F^2.$$

Introducing  $\mathbf{X}$  in the left-hand side, expanding the square, and simplifying leads to

$$(31) \quad \|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F^2 \leq 2\langle \mathbf{X} - \mathbf{X}^{\text{ht}}, \mathbf{X} - \mathcal{A}^*(\mathbf{y}) \rangle_F.$$

Note that  $\mathbf{X}$  and  $\mathbf{X}^{\text{ht}}$  are of rank at most  $r$ , so that both  $\mathbf{X}$  and  $\mathbf{X} - \mathbf{X}^{\text{ht}}$  belong to a space  $\mathcal{S} = \text{span}\{\mathbf{u}_1 \mathbf{v}_1^*, \dots, \mathbf{u}_{2r} \mathbf{v}_{2r}^*\}$ . Hence, also using Lemma 3, we derive

$$(32) \quad \begin{aligned} \|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F^2 &\leq 2\langle \mathbf{X} - \mathbf{X}^{\text{ht}}, P_{\mathcal{S}}(\mathbf{X} - \mathcal{A}^*(\mathbf{y})) \rangle = 2\langle \mathbf{X} - \mathbf{X}^{\text{ht}}, \mathbf{X} - P_{\mathcal{S}}(\mathcal{A}^*(\mathbf{y})) \rangle \\ &\leq 2\|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F \|\mathbf{X} - P_{\mathcal{S}}(\mathcal{A}^*(\mathbf{y}))\|_F \leq 2\|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F \sqrt{5\delta + 4\|\mathbf{e}\|_1}. \end{aligned}$$

This immediately implies that

$$(33) \quad \|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F \leq 2\sqrt{5\delta + 4\|\mathbf{e}\|_1},$$

which already proves that  $\mathbf{X}$  is well approximated by the unnormalized vector  $\mathbf{X}^{\text{ht}}$ . To realize that it is also well approximated by the normalization  $\mathbf{X}^{\text{ht}}/\|\mathbf{X}^{\text{ht}}\|_F$  of this vector, we remark that the latter is the best normalized approximant to  $\mathbf{X}^{\text{ht}}$  and hence it is a better approximant than  $\mathbf{X}$ , so that  $\|\mathbf{X}^{\text{ht}} - \mathbf{X}^{\text{ht}}/\|\mathbf{X}^{\text{ht}}\|_F\|_F \leq \|\mathbf{X}^{\text{ht}} - \mathbf{X}\|_F$ , and in turn

$$(34) \quad \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{ht}}}{\|\mathbf{X}^{\text{ht}}\|_F} \right\|_F \leq \|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F + \left\| \mathbf{X}^{\text{ht}} - \frac{\mathbf{X}^{\text{ht}}}{\|\mathbf{X}^{\text{ht}}\|_F} \right\|_F \leq 2\|\mathbf{X} - \mathbf{X}^{\text{ht}}\|_F \leq 4\sqrt{5\delta + 4\|\mathbf{e}\|_1}.$$

The required estimate (23) is now a consequence of  $\sqrt{5\delta + 4\|\mathbf{e}\|_1} \leq \sqrt{5\delta} + 2\sqrt{\|\mathbf{e}\|_1}$ .  $\square$

### 3.2 Semidefinite programming

The output of the semidefinite procedure is, in the absence of prequantization error,

$$(35) \quad \mathbf{X}^{\text{sdp}} = \operatorname{argmin} \|\mathbf{Z}\|_* \quad \text{subject to } \operatorname{sgn}(\mathcal{A}\mathbf{Z}) = \mathbf{y} \quad \text{and} \quad \|\mathcal{A}\mathbf{Z}\|_1 = 1.$$

This optimization program does not immediately appear as a semidefinite program, but it can be recast as one by adding slack variables  $c \in \mathbb{R}$  and  $\mathbf{M}, \mathbf{N} \in \mathbb{R}^{n \times n}$ , because the condition  $\|\mathbf{Z}\|_* \leq c$  can be rephrased as the condition  $(\operatorname{tr}(\mathbf{M}) + \operatorname{tr}(\mathbf{N}))/2 \leq c$  subject to  $\begin{bmatrix} \mathbf{M} & \mathbf{Z} \\ \mathbf{Z}^* & \mathbf{N} \end{bmatrix} \succeq \mathbf{0}$ . As for the constraints  $\operatorname{sgn}(\mathcal{A}\mathbf{Z}) = \mathbf{y}$  and  $\|\mathcal{A}\mathbf{Z}\|_1 = 1$ , they reduce to the linear constraints  $y_i \langle \mathbf{A}_i, \mathbf{Z} \rangle_F \geq 0$ ,  $i \in [1 : m]$ , and  $\sum_{i=1}^m y_i \langle \mathbf{A}_i, \mathbf{Z} \rangle_F = 1$ . In the presence of prequantization error  $\mathbf{e} \in \mathbb{R}^m$  for which a bound  $\|\mathbf{e}\|_1 \leq \eta$  is available, we refine the optimization procedure slightly and consider

$$(36) \quad (\mathbf{X}^{\text{sdp}}, \mathbf{e}^{\text{sdp}}) = \operatorname{argmin}_{(\mathbf{Z}, \mathbf{w}) \in \mathbb{R}^{n \times n} \times \mathbb{R}^m} \|\mathbf{Z}\|_* \quad \text{subject to} \quad \begin{cases} \operatorname{sgn}(\mathcal{A}\mathbf{Z} + \mathbf{w}) = \mathbf{y}, \\ \|\mathcal{A}\mathbf{Z} + \mathbf{w}\|_1 = 1, \\ \|\mathbf{w}\|_1 \leq (10/7)\eta. \end{cases}$$

The value  $10/7$  is somewhat arbitrary — it has been chosen to make the subsequent arguments look nice. Note that the constraint  $\|\mathbf{w}\|_1 \leq (10/7)\eta$  can still be recast as linear constraints, e.g. by introducing slack variables  $\mathbf{w}^+, \mathbf{w}^- \in \mathbb{R}_+^m$  with  $\mathbf{w}^+ - \mathbf{w}^- = \mathbf{w}$  and  $\sum_{i=1}^m (w_i^+ + w_i^-) \leq (10/7)\eta$ . The main theorem of this subsection reads as follows. Once again, we emphasize that its assumption is met with high probability by (properly normalized) Gaussian measurement maps  $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ , provided  $m \geq C\delta^{-5}nr$  (or  $m \geq C\delta^{-2}nr$  if the remark after Theorem 1 is taken into account).

**Theorem 4.** Under  $\text{MRRIP}^{\text{eff}}(16r, \delta)$ , if  $\mathbf{X} \in \mathbb{R}^{n \times n}$  satisfying  $\text{effrank}(\mathbf{X}) \leq r$  and  $\|\mathbf{X}\|_F = 1$  is acquired via  $\mathbf{y} = \operatorname{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) \in \{\pm 1\}^m$  with  $\|\mathbf{e}\|_1 \leq \eta$ , then

$$(37) \quad \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{sdp}}}{\|\mathbf{X}^{\text{sdp}}\|_F} \right\|_F \leq C\sqrt{\delta} + D\sqrt{\eta}$$

for some absolute constants  $C, D > 0$ .



The argument relies on a lemma about the effective rank of the output of the semidefinite program.

**Lemma 5.** Suppose that  $\eta \leq 1/10$  and that  $\mathcal{A}$  satisfies  $\text{MRRIP}(16r, 1/5)$ . If  $\mathbf{X} \in \mathbb{R}^{n \times n}$  satisfies  $\text{rank}(\mathbf{X}) \leq r$  and  $\|\mathbf{X}\|_F = 1$ , then any convex combination of  $\mathbf{X}$  and  $\mathbf{X}^{\text{sdp}}$  is of effective rank at most  $16r$ . The same conclusion holds if  $\text{effrank}(\mathbf{X}) \leq r$  and  $\|\mathbf{X}\|_F = 1$ , provided  $\mathcal{A}$  satisfies  $\text{MRRIP}^{\text{eff}}(16r, 1/5)$ .

*Proof.* We set  $t = 16r$  and  $\delta = 1/5$ . Taking note of

$$(38) \quad \|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1 \geq \|\mathcal{A}\mathbf{X}\|_1 - \|\mathbf{e}\|_1 \geq (1 - \delta) - \eta \geq 1 - \frac{1}{5} - \frac{1}{10} = \frac{7}{10},$$

we first point out that the couple  $(\mathbf{X}, \mathbf{e})/\|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1$  is feasible for the optimization program in (36). It follows that

$$(39) \quad \|\mathbf{X}^{\text{sdp}}\|_* \leq \left\| \frac{\mathbf{X}}{\|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1} \right\|_* \leq \frac{\sqrt{r}\|\mathbf{X}\|_F}{7/10} = \frac{10}{7}\sqrt{r}.$$

Let  $\widehat{\mathbf{X}} := (1 - \lambda)\mathbf{X} + \lambda\mathbf{X}^{\text{sdp}}$ ,  $\lambda \in [0, 1]$ , be a convex combination of  $\mathbf{X}$  and  $\mathbf{X}^{\text{sdp}}$ . We also introduce  $\widehat{\mathbf{e}} := (1 - \lambda)\mathbf{e} + \lambda\mathbf{e}^{\text{sdp}}$ . We notice that

$$(40) \quad \|\widehat{\mathbf{X}}\|_* \leq (1 - \lambda)\|\mathbf{X}\|_* + \lambda\|\mathbf{X}^{\text{sdp}}\|_* \leq (1 - \lambda)\sqrt{r} + \lambda\frac{10}{7}\sqrt{r} = \left(1 + \frac{3\lambda}{7}\right)\sqrt{r}.$$

Moreover, thanks to  $\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) = \text{sgn}(\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}) = \text{sgn}(\mathcal{A}\widehat{\mathbf{X}} + \widehat{\mathbf{e}}) = \mathbf{y}$ , we have

$$(41) \quad \|\mathcal{A}\widehat{\mathbf{X}} + \widehat{\mathbf{e}}\|_1 = (1 - \lambda)\|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1 + \lambda\|\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}\|_1 \geq (1 - \lambda)\frac{7}{10} + \lambda = \frac{7}{10} \left(1 + \frac{3\lambda}{7}\right),$$

while  $\|\widehat{\mathbf{e}}\|_1$  can be bounded from above as

$$(42) \quad \|\widehat{\mathbf{e}}\|_1 \leq (1 - \lambda)\|\mathbf{e}\|_1 + \lambda\|\mathbf{e}^{\text{sdp}}\|_1 \leq (1 - \lambda)\eta + \lambda\frac{10}{7}\eta = \left(1 + \frac{3\lambda}{7}\right)\eta \leq \frac{1}{10} \left(1 + \frac{3\lambda}{7}\right),$$

so that  $\|\mathcal{A}\widehat{\mathbf{X}}\|_1$  can be bounded from below as

$$(43) \quad \|\mathcal{A}\widehat{\mathbf{X}}\|_1 \geq \|\mathcal{A}\widehat{\mathbf{X}} + \widehat{\mathbf{e}}\|_1 - \|\widehat{\mathbf{e}}\|_1 \geq \left(\frac{7}{10} - \frac{1}{10}\right) \left(1 + \frac{3\lambda}{7}\right) = \frac{3}{5} \left(1 + \frac{3\lambda}{7}\right).$$

Besides, by applying the sort-and-split technique, we write  $\widehat{\mathbf{X}} = \widehat{\mathbf{X}}_0 + \widehat{\mathbf{X}}_1 + \widehat{\mathbf{X}}_2 + \dots$ , where the  $\widehat{\mathbf{X}}_k$ 's are rank- $r$  matrices defined from the singular value decomposition  $\widehat{\mathbf{X}} = \sum_{i=1}^n \sigma_i(\widehat{\mathbf{X}}) \mathbf{u}_i \mathbf{v}_i^*$  of  $\widehat{\mathbf{X}}$  by  $\widehat{\mathbf{X}}_k = \sum_{i=kt+1}^{(k+1)t} \sigma_i(\widehat{\mathbf{X}}) \mathbf{u}_i \mathbf{v}_i^*$ , and we observe that

$$(44) \quad \begin{aligned} \|\mathcal{A}\widehat{\mathbf{X}}\|_1 &\leq \|\mathcal{A}\widehat{\mathbf{X}}_0\|_1 + \|\mathcal{A}\widehat{\mathbf{X}}_1\|_1 + \|\mathcal{A}\widehat{\mathbf{X}}_2\|_1 + \dots \leq (1 + \delta) \left( \|\widehat{\mathbf{X}}_0\|_F + \sum_{k \geq 1} \|\widehat{\mathbf{X}}_k\|_F \right) \\ &\leq (1 + \delta) \left( \|\widehat{\mathbf{X}}\|_F + \sum_{k \geq 1} \frac{\|\widehat{\mathbf{X}}_{k-1}\|_*}{\sqrt{t}} \right) \leq (1 + \delta) \left( \|\widehat{\mathbf{X}}\|_F + \frac{1}{\sqrt{t}} \|\widehat{\mathbf{X}}\|_* \right) \\ &\leq (1 + \delta) \left( \|\widehat{\mathbf{X}}\|_F + \frac{\sqrt{r}}{\sqrt{t}} \left(1 + \frac{3\lambda}{7}\right) \right) = \frac{6}{5} \left( \|\widehat{\mathbf{X}}\|_F + \frac{1}{4} \left(1 + \frac{3\lambda}{7}\right) \right). \end{aligned}$$

Combining the lower bound (43) and upper bound (44) on  $\|\mathcal{A}\hat{\mathbf{X}}\|_1$  gives  $\|\hat{\mathbf{X}}\|_F \geq \left(1 + \frac{3\lambda}{7}\right)/4$ . We finally arrive at

$$(45) \quad \frac{\|\hat{\mathbf{X}}\|_*}{\|\hat{\mathbf{X}}\|_F} \leq \frac{(1 + 3\lambda/7)\sqrt{r}}{(1 + 3\lambda/7)/4} = \sqrt{16r}.$$

This means that  $\hat{\mathbf{X}}$  is effectively of rank at most  $16r$ , as announced.  $\square$

With this crucial lemma now established, the main result about direction recovery by semidefinite programming follows from a curvature argument already found in [9].

*Proof of Theorem 4.* Suppose first that  $\delta \leq 1/5$  and  $\eta \leq 1/10$ . The parallelogram identity gives

$$(46) \quad \left\| \frac{\mathbf{X} - \mathbf{X}^{\text{sdp}}}{2} \right\|_F^2 + \left\| \frac{\mathbf{X} + \mathbf{X}^{\text{sdp}}}{2} \right\|_F^2 = \frac{\|\mathbf{X}\|_F^2 + \|\mathbf{X}^{\text{sdp}}\|_F^2}{2} = \frac{1 + \|\mathbf{X}^{\text{sdp}}\|_F^2}{2}.$$

According to Lemma 5, both  $\mathbf{X}^{\text{sdp}}$  and  $(\mathbf{X} + \mathbf{X}^{\text{sdp}})/2$  are effectively of rank at most  $16r$ . Thus, on the one hand,

$$(47) \quad \|\mathbf{X}^{\text{sdp}}\|_F \leq \frac{1}{1-\delta} \|\mathcal{A}\mathbf{X}^{\text{sdp}}\|_F \leq \frac{1}{1-\delta} \left( \|\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}\|_1 + \|\mathbf{e}^{\text{sdp}}\|_1 \right) \leq \frac{1 + (10/7)\eta}{1-\delta},$$

while on the other hand, using  $\text{sgn}(\mathcal{A}\mathbf{X} + \mathbf{e}) = \text{sgn}(\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}) = \mathbf{y}$ ,

$$(48) \quad \begin{aligned} \left\| \frac{\mathbf{X} + \mathbf{X}^{\text{sdp}}}{2} \right\|_F &\geq \frac{1}{1+\delta} \left\| \mathcal{A} \left( \frac{\mathbf{X} + \mathbf{X}^{\text{sdp}}}{2} \right) \right\|_1 \\ &\geq \frac{1}{2(1+\delta)} \left( \|\mathcal{A}\mathbf{X} + \mathbf{e} + \mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}\|_1 - \|\mathbf{e}\|_1 - \|\mathbf{e}^{\text{sdp}}\|_1 \right) \\ &= \frac{1}{2(1+\delta)} \left( \|\mathcal{A}\mathbf{X} + \mathbf{e}\|_1 + \|\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}\|_1 - \|\mathbf{e}\|_1 - \|\mathbf{e}^{\text{sdp}}\|_1 \right) \\ &\geq \frac{1}{2(1+\delta)} \left( \|\mathcal{A}\mathbf{X}\|_1 + \|\mathcal{A}\mathbf{X}^{\text{sdp}} + \mathbf{e}^{\text{sdp}}\|_1 - 2\|\mathbf{e}\|_1 - \|\mathbf{e}^{\text{sdp}}\|_1 \right) \\ &\geq \frac{1}{2(1+\delta)} \left( (1-\delta) + 1 - 2\eta - \frac{10}{7}\eta \right) = \frac{1 - \delta/2 - (12/7)\eta}{1+\delta}. \end{aligned}$$

Substituting (47) and (48) into (46) leads to

$$(49) \quad \left\| \frac{\mathbf{X} - \mathbf{X}^{\text{sdp}}}{2} \right\|_F^2 \leq \frac{1 + \left( \frac{1 + (10/7)\eta}{1-\delta} \right)^2}{2} - \left( \frac{1 - \delta/2 - (12/7)\eta}{1+\delta} \right)^2 \leq C'\delta + D'\eta$$

for some absolute constants  $C', D' > 0$  that could be determined explicitly if needed. In turn, since  $\mathbf{X}^{\text{sdp}}/\|\mathbf{X}^{\text{sdp}}\|_F$  is the best normalized approximant to  $\mathbf{X}^{\text{sdp}}$ , we obtain

$$(50) \quad \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{sdp}}}{\|\mathbf{X}^{\text{sdp}}\|_F} \right\|_F \leq \left\| \mathbf{X} - \mathbf{X}^{\text{sdp}} \right\|_F + \left\| \mathbf{X}^{\text{sdp}} - \frac{\mathbf{X}^{\text{sdp}}}{\|\mathbf{X}^{\text{sdp}}\|_F} \right\|_F \leq 2 \left\| \mathbf{X} - \mathbf{X}^{\text{sdp}} \right\|_F \leq C\sqrt{\delta} + D\sqrt{\eta}$$

for some constants  $C, D > 0$  that could again be determined explicitly. Up to possibly changing these constants, the same estimate remains true if  $\eta \geq 1/10$ , since then

$$(51) \quad \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{sdp}}}{\|\mathbf{X}^{\text{sdp}}\|_F} \right\|_F \leq \|\mathbf{X}\|_F + \left\| \frac{\mathbf{X}^{\text{sdp}}}{\|\mathbf{X}^{\text{sdp}}\|_F} \right\|_F = 2 \leq 2\sqrt{10}\sqrt{\eta} \leq C\sqrt{\delta} + D\sqrt{\eta}.$$

A similar reasoning would take care of the case  $\delta \geq 1/5$ . The proof is now complete.  $\square$

## 4 Estimating the Magnitude as well as the Direction

In this section, we assume that a bound  $\|\mathbf{X}\|_F \leq \gamma$  is available for the Frobenius norm of low-rank matrices  $\mathbf{X} \in \mathbb{R}^{n \times n}$ . We show that, in this case, we can estimate not only the direction  $\mathbf{X}/\|\mathbf{X}\|_F$ , but also the magnitude  $\|\mathbf{X}\|_F$ , i.e., it is possible to fully estimate  $\mathbf{X}$ , provided the binary measurements feature well-chosen thresholds  $\tau_1, \dots, \tau_m \in \mathbb{R}$ . Namely, these measurements take the form  $y_i = \text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F - \tau_i)$ , or more realistically, in the presence of prequantization error  $\mathbf{e} \in \mathbb{R}^m$ ,

$$(52) \quad y_i = \text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F + e_i - \tau_i), \quad i \in \llbracket 1 : m \rrbracket.$$

For the recovery algorithm, one can choose between a hard singular value thresholding procedure or a semidefinite procedure. The argument involves a combination of the results from Section 3 and an augmentation trick already used in [3]. Precisely, matrices  $\mathbf{X} \in \mathbb{R}^{n \times n}$  of (effective) rank at most  $r$  are augmented to matrices  $\tilde{\mathbf{X}} \in \mathbb{R}^{(n+1) \times (n+1)}$  of (effective) rank at most  $r + 1$  via

$$(53) \quad \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} & 0 \\ 0 & \gamma \end{bmatrix}.$$

One then works in the augmented space  $\mathbb{R}^{(n+1) \times (n+1)}$  to approximate the direction of  $\tilde{\mathbf{X}}$ , which provides a full approximation of the original  $\mathbf{X}$ . This can be done because the thresholded binary measurements in the original space can be interpreted in the augmented space as the unthresholded binary measurements

$$(54) \quad y_i = \text{sgn}(\langle \tilde{\mathbf{A}}_i, \tilde{\mathbf{X}} \rangle_F + e_i) = \text{sgn}((\tilde{\mathcal{A}}\tilde{\mathbf{X}})_i + e_i), \quad i \in \llbracket 1 : m \rrbracket,$$

where the augmented matrices  $\tilde{\mathbf{A}}_1, \dots, \tilde{\mathbf{A}}_m \in \mathbb{R}^{(n+1) \times (n+1)}$  are given by

$$(55) \quad \tilde{\mathbf{A}}_i = \begin{bmatrix} \mathbf{A}_i & \mathbf{u}_i \\ \mathbf{v}_i^* & -\tau_i/\gamma \end{bmatrix}$$

for arbitrary vectors  $\mathbf{u}_1, \mathbf{v}_1, \dots, \mathbf{u}_m, \mathbf{v}_m \in \mathbb{R}^n$ .

We now formalize the main result of this section. It is worth pointing out, once again, that the necessary number of measurements does not comprise any spurious logarithmic factor.

**Theorem 6.** There exist absolute constants  $c, C > 0$  such that, if  $m \geq C\delta^{-3}nr$ , if independent random matrices  $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n \times n}$  are populated by independent Gaussian variables with mean zero and standard deviation  $\sqrt{\pi/2}/m$ , and if random thresholds  $\tau_1, \dots, \tau_m \in \mathbb{R}$  are Gaussian variables with mean zero and standard deviation  $\gamma\sqrt{\pi/2}/m$ , independent from one another and from the  $\mathbf{A}_i$ 's, then the following holds with failure probability at most  $2\exp(-c\delta^2m)$ :

every matrix  $\mathbf{X} \in \mathbb{R}^{n \times n}$  which has rank at most  $r$ , satisfies  $\|\mathbf{X}\|_F \leq \gamma$ , and is acquired via  $y_i = \text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F + e_i - \tau_i)$ ,  $i \in \llbracket 1 : m \rrbracket$ , is approximated with error

$$(56) \quad \left\| \mathbf{X} - \frac{\gamma}{x^{\text{ht}}} \mathbf{X}^{\text{ht}} \right\|_F \leq C\sqrt{\delta}\gamma + D\sqrt{\|\mathbf{e}\|_1}\sqrt{\gamma},$$

$$(57) \quad \left\| \mathbf{X} - \frac{\gamma}{x^{\text{sdp}}} \mathbf{X}^{\text{sdp}} \right\|_F \leq C\sqrt{\delta}\gamma + D\sqrt{\eta}\sqrt{\gamma}, \quad \eta \text{ being an a priori bound on } \|\mathbf{e}\|_1.$$

The matrices  $\mathbf{X}^{\text{ht}}, \mathbf{X}^{\text{sdp}} \in \mathbb{R}^{n \times n}$  and the scalars  $x^{\text{ht}}, x^{\text{sdp}} \in \mathbb{R}$  are produced by<sup>1</sup>

$$(58) \quad \begin{bmatrix} \mathbf{X}^{\text{ht}} & * \\ * & x^{\text{ht}} \end{bmatrix} = \text{best rank-}(r+1) \text{ approximant to } \tilde{\mathcal{A}}^*(\mathbf{y}) = \sum_{i=1}^m y_i \tilde{\mathbf{A}}_i,$$

$$(59) \quad \left( \begin{bmatrix} \mathbf{X}^{\text{sdp}} & * \\ * & x^{\text{sdp}} \end{bmatrix}, \mathbf{e}^{\text{sdp}} \right) = \underset{(\tilde{\mathbf{Z}}, \mathbf{w}) \in \mathbb{R}^{(n+1) \times (n+1)} \times \mathbb{R}^m}{\text{argmin}} \|\tilde{\mathbf{Z}}\|_* \quad \text{subject to} \quad \begin{cases} \text{sgn}(\tilde{\mathcal{A}}\tilde{\mathbf{Z}} + \mathbf{w}) = \mathbf{y}, \\ \|\tilde{\mathcal{A}}\tilde{\mathbf{Z}} + \mathbf{w}\|_1 = 1, \\ \|\mathbf{w}\|_1 \leq (10/7)\eta, \end{cases}$$

with extra randomness injected through the presence in (55) of independent random vectors  $\mathbf{u}_1, \mathbf{v}_1, \dots, \mathbf{u}_m, \mathbf{v}_m \in \mathbb{R}^n$  populated by independent Gaussian variables with mean zero and standard deviation  $\sqrt{\pi/2}/m$ .

The following observation is central to the argument.

**Lemma 7.** If  $\tilde{\mathbf{X}}, \tilde{\mathbf{Z}} \in \mathbb{R}^{(n+1) \times (n+1)}$  take the form

$$(60) \quad \tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{X} & * \\ * & x \end{bmatrix}, \quad \tilde{\mathbf{Z}} = \begin{bmatrix} \mathbf{Z} & * \\ * & z \end{bmatrix},$$

then

$$(61) \quad \left\| \frac{\mathbf{X}}{x} - \frac{\mathbf{Z}}{z} \right\|_F \leq \frac{\|\tilde{\mathbf{X}}\|_F \|\tilde{\mathbf{Z}}\|_F}{|x||z|} \left\| \frac{\tilde{\mathbf{X}}}{\|\tilde{\mathbf{X}}\|_F} - \frac{\tilde{\mathbf{Z}}}{\|\tilde{\mathbf{Z}}\|_F} \right\|_F.$$

*Proof.* A vector analog of this statement appeared in [3, Lemma 8]. One can follow the steps presented there and adapt them to our situation. Alternatively, one can vectorize matrices in  $\mathbb{R}^{(n+1) \times (n+1)}$  to vectors in  $\mathbb{R}^{n^2+2n+1}$ , apply [3, Lemma 8] to bound the  $\ell_2$ -norm of the difference of two vectors in  $\mathbb{R}^{n^2+2n}$  by the right-hand side of (61), and observe that the left-hand side of (61) is itself bounded by the  $\ell_2$ -norm of this difference.  $\square$

<sup>1</sup>The facts that  $x^{\text{ht}} \neq 0$  and  $x^{\text{sdp}} \neq 0$  are established in the proof.

Let us now exploit the observation made in Lemma 7 to prove the result stated above.

*Proof of Theorem 6.* The assumptions guarantee that the map  $\tilde{\mathcal{A}} : \mathbb{R}^{(n+1) \times (n+1)} \rightarrow \mathbb{R}^m$  satisfies MRRIP(16(r+1),  $\delta$ ). Therefore, Theorems 2 and 4 imply that

$$(62) \quad \left\| \frac{\tilde{\mathbf{X}}}{\|\tilde{\mathbf{X}}\|_F} - \frac{\tilde{\mathbf{X}}^\circ}{\|\tilde{\mathbf{X}}^\circ\|_F} \right\|_F \leq C\sqrt{\delta} + D\sqrt{\frac{\eta^\circ}{\|\tilde{\mathbf{X}}\|_F}} := \varepsilon^\circ,$$

where  $\tilde{\mathbf{X}}^\circ = \begin{bmatrix} \mathbf{X}^\circ & * \\ * & x^\circ \end{bmatrix}$  is either given by (58), in which case  $\eta^\circ := \|\mathbf{e}\|_1$ , or by (59), in which case  $\eta^\circ := \eta$ . Looking at the bottom right entry, we derive that

$$(63) \quad \left| \frac{\gamma}{\|\tilde{\mathbf{X}}\|_F} - \frac{x^\circ}{\|\tilde{\mathbf{X}}^\circ\|_F} \right| \leq \varepsilon^\circ,$$

so that, in view of  $\|\tilde{\mathbf{X}}\|_F = \sqrt{\|\mathbf{X}\|_F^2 + \gamma^2} \leq \sqrt{\gamma^2 + \gamma^2} = \sqrt{2}\gamma$ ,

$$(64) \quad \frac{|x^\circ|}{\|\tilde{\mathbf{X}}^\circ\|_F} \geq \frac{\gamma}{\|\tilde{\mathbf{X}}\|_F} - \varepsilon^\circ \geq \frac{1}{\sqrt{2}} - \varepsilon^\circ \geq \frac{1}{2},$$

provided  $\varepsilon^\circ$  is small enough. Lemma 7, together with (62), now gives

$$(65) \quad \left\| \frac{\mathbf{X}}{\gamma} - \frac{\mathbf{X}^\circ}{x^\circ} \right\|_F \leq \frac{\|\tilde{\mathbf{X}}\|_F}{\gamma} \frac{\|\tilde{\mathbf{X}}^\circ\|_F}{|x^\circ|} \varepsilon^\circ \leq \sqrt{2} \times 2 \times \varepsilon^\circ.$$

Multiplying throughout by  $\gamma$ , we obtain

$$(66) \quad \left\| \mathbf{X} - \frac{\gamma}{x^\circ} \mathbf{X}^\circ \right\|_F \leq 2\sqrt{2}\varepsilon^\circ\gamma = 2\sqrt{2} \left( C\sqrt{\delta}\gamma + D\sqrt{\frac{\eta^\circ}{\|\tilde{\mathbf{X}}\|_F}}\gamma \right),$$

which, in view of  $\|\tilde{\mathbf{X}}\|_F \geq \gamma$ , yields the estimates required in (56) and (57), up to an adjustment of the constants  $C, D > 0$ .  $\square$

## 5 Exponential Decay of the Recovery Error

In this section, we still work under the assumption that an a priori bound  $\|\mathbf{X}\|_F \leq \gamma$  is available for the rank- $r$  matrices  $\mathbf{X} \in \mathbb{R}^{n \times n}$  of interest. In the absence of prequantization error, Section 4 essentially showed that random thresholds  $\boldsymbol{\tau} \in \mathbb{R}^m$  in  $\mathbf{y} = \text{sgn}(\mathcal{A}\mathbf{X} - \boldsymbol{\tau})$  allow for the computation of approximants satisfying

$$(67) \quad \|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq C\gamma\lambda^{-1/6},$$

where  $\lambda := m/(nr)$  represents an oversampling factor — note that no logarithmic factors are featured here. We now aim at proving that one can find a recovery procedure making the error  $\|\mathbf{X} - \widehat{\mathbf{X}}\|_F$  decay exponentially fast with the oversampling factor  $\lambda$ , provided the thresholds  $\boldsymbol{\tau}$  are chosen adaptively. The vector analog of this result was established in [2] and we follow the same guiding idea here.

The procedure intertwines measurement and reconstruction processes by dividing the  $m = qT$  binary measurements into  $T$  batches of  $q \asymp nr$  binary measurements, each of them having the form

$$(68) \quad \mathbf{y}^{(t)} = \text{sgn}(\mathcal{A}^{(t)}\mathbf{X} + \mathbf{e}^{(t)} - \boldsymbol{\tau}^{(t)}) \in \{\pm 1\}^q, \quad t \in \llbracket 1 : T \rrbracket.$$

Precisely, the procedure iteratively constructs, starting from  $\mathbf{X}^{(0)} = \mathbf{0}$ , a sequence  $(\boldsymbol{\tau}^{(t)})_{t \in \llbracket 0 : T-1 \rrbracket}$  of thresholds and a sequence  $(\mathbf{X}^{(t)})_{t \in \llbracket 0 : T \rrbracket}$  of matrices according to

$$(69) \quad \boldsymbol{\tau}^{(t)} := \mathcal{A}^{(t)}\mathbf{X}^{(t)} + \frac{\gamma}{2^t} \frac{\sqrt{\pi/2}}{q} \mathbf{g}^{(t)}, \quad \mathbf{g}^{(t)} \in \mathbb{R}^q \text{ populated by independent } \mathcal{N}(0, 1) \text{ entries,}$$

$$(70) \quad \widehat{\mathbf{X} - \mathbf{X}^{(t)}} := \Delta(\mathbf{y}^{(t)}), \quad \Delta \text{ one of the reconstruction maps from Theorem 6,}$$

$$(71) \quad \mathbf{X}^{(t+1)} := \left[ \mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} \right]_{(r)}, \quad \text{i.e., the best rank-}r \text{ approximant to } \mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}}.$$

**Theorem 8.** For  $t \in \llbracket 1 : T \rrbracket$ , let  $\mathcal{A}^{(t)} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^q$  be maps induced by independent random matrices  $\mathbf{A}_1^{(t)}, \dots, \mathbf{A}_q^{(t)}$  populated by independent Gaussian variables with mean zero and standard deviation  $\sqrt{\pi/2}/q$  and let  $\mathbf{g}^{(t)} \in \mathbb{R}^q$  be random vectors populated by standard Gaussian variables independent from one another and from the  $\mathbf{A}_i^{(t)}$ 's. There exist absolute constants  $c, c', c'' > 0$  such that, if  $q \approx c'nr$ , then the following holds with failure probability at most  $2T \exp(-c''q)$ :

every matrix  $\mathbf{X} \in \mathbb{R}^{n \times n}$  which has rank at most  $r$ , satisfies  $\|\mathbf{X}\|_F \leq \gamma$ , and is acquired via the  $m = qT$  binary measurements  $y_i^{(t)} = \text{sgn}(\langle \mathbf{A}_i, \mathbf{X} \rangle_F + e_i - \tau_i)$ ,  $i \in \llbracket 1 : q \rrbracket$ ,  $t \in \llbracket 1 : T \rrbracket$ , where  $\|\mathbf{e}^{(t)}\|_1 \leq \nu(\gamma/2^t)$  for an appropriately small constant  $\nu > 0$ , is approximated by the output  $\mathbf{X}^{(T)}$  of the procedure described in (69)-(70)-(71) with error

$$(72) \quad \|\mathbf{X} - \mathbf{X}^{(T)}\|_F \leq \gamma \exp(-c\lambda), \quad \lambda := \frac{m}{nr}.$$

*Proof.* We shall prove by induction on  $t \in \llbracket 0 : T \rrbracket$  that

$$(73) \quad \|\mathbf{X} - \mathbf{X}^{(t)}\|_F \leq \frac{\gamma}{2^t}.$$

Indeed, the resut for  $t = T$  implies the desired bound (72), since

$$(74) \quad \|\mathbf{X} - \mathbf{X}^{(T)}\|_F \leq \frac{\gamma}{2^T} = \gamma \exp(-\ln(2)T) = \gamma \exp\left(-\ln(2)\frac{m}{q}\right) = \gamma \exp\left(-c\frac{m}{nr}\right).$$

The inductive hypothesis clearly holds for  $t = 0$ . Let us now suppose that it holds for  $t \in \llbracket 0 : T-1 \rrbracket$ . We observe that

$$(75) \quad \widehat{\mathbf{X} - \mathbf{X}^{(t)}} = \Delta\left(\text{sgn}\left(\mathcal{A}^{(t)}(\mathbf{X} - \mathbf{X}^{(t)}) + \mathbf{e}^{(t)} - \frac{\gamma}{2^t} \frac{\sqrt{\pi/2}}{q} \mathbf{g}^{(t)}\right)\right).$$

Since  $\mathbf{X} - \mathbf{X}^{(t)}$  has rank at most  $2r$  and satisfies  $\|\mathbf{X} - \mathbf{X}^{(t)}\|_F \leq \gamma/2^t$ , Theorem 6 applied with a small enough  $\delta \in (0, 1)$  yields

$$(76) \quad \|\mathbf{X} - \mathbf{X}^{(t)} - \widehat{\mathbf{X} - \mathbf{X}^{(t)}}\|_F \leq C\sqrt{\delta}\frac{\gamma}{2^t} + D\sqrt{\frac{\nu\gamma}{2^t}}\sqrt{\frac{\gamma}{2^t}} \leq \frac{1}{4}\frac{\gamma}{2^t}.$$

But since  $\mathbf{X}^{(t+1)}$  is the best rank- $r$  approximant to  $\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}}$ , we have

$$(77) \quad \|\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} - \mathbf{X}^{(t+1)}\|_F \leq \|\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} - \mathbf{X}\|_F.$$

In turn, we deduce from (77) and (76) that

$$(78) \quad \begin{aligned} \|\mathbf{X} - \mathbf{X}^{(t+1)}\|_F &\leq \|\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} - \mathbf{X}^{(t+1)}\|_F + \|\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} - \mathbf{X}\|_F \\ &\leq 2\|\mathbf{X}^{(t)} + \widehat{\mathbf{X} - \mathbf{X}^{(t)}} - \mathbf{X}\|_F = 2\|\mathbf{X} - \mathbf{X}^{(t)} - \widehat{\mathbf{X} - \mathbf{X}^{(t)}}\|_F \\ &\leq 2\frac{1}{4}\frac{\gamma}{2^t} = \frac{\gamma}{2^{t+1}}. \end{aligned}$$

This shows that the induction hypothesis holds for  $t + 1$ . The proof is therefore complete.  $\square$

## 6 Numerical Experiments

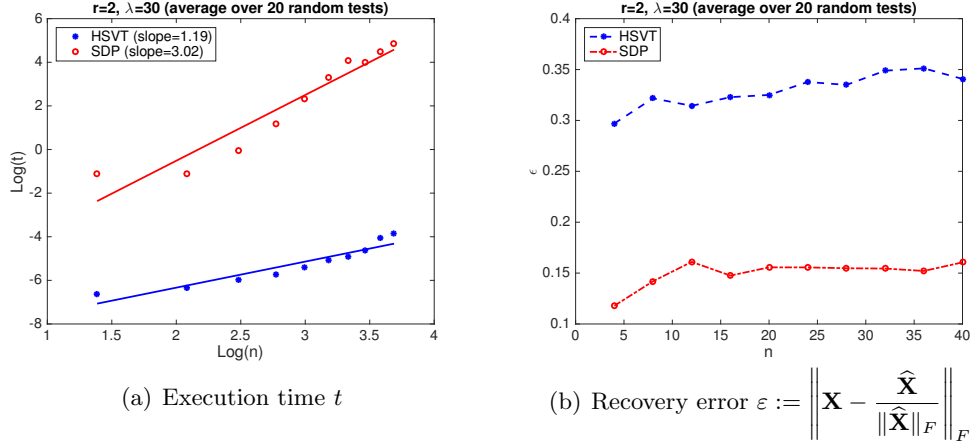
The goal of this final section is merely illustrative. The investigations showcased here are far from exhaustive, since they mostly serve as a confirmation that the algorithms considered in the paper are indeed implementable. Our own implementations (by no means optimized) can be downloaded from the first author's webpage as part of the MATLAB reproducible containing the code to generate the experiments below. In order to execute the semidefinite optimization procedures, CVX [1], a package for specifying and solving convex programs, is required.

### 6.1 Direction estimation

In this subsection, we look at the procedures (22) and (35)-(36) for estimating the direction of a low-rank matrix acquired via binary measurements without thresholds.

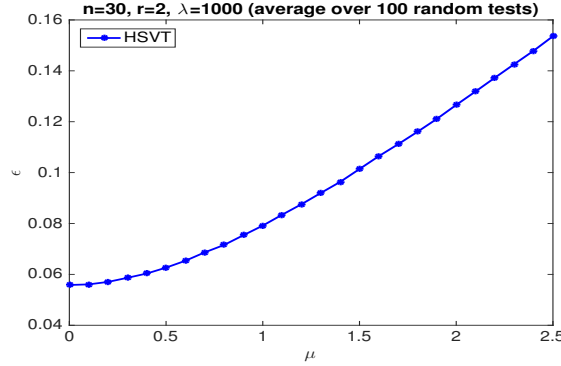
#### Hard singular value thresholding vs semidefinite programming

In a first experiment, we consider the situation where there is no prequantization error ( $\mathbf{e} = \mathbf{0}$ ). For fixed values of the rank  $r$  and of the oversampling factor  $\lambda$ , we vary  $n$  and record both the execution time and the recovery error (averaged over a reasonable number of tests). Figure 1(a) supports the intuitive prediction that the hard singular value thresholding procedure is much faster


 Figure 1: Comparison of the procedures (22) and (35) when the dimension  $n$  varies

than the semidefinite procedure (and the slopes even suggest a cost  $\Theta(n^c)$  with  $c$  close to 1 vs close to 3), while Figure 1(b) supports the other intuitive prediction that the semidefinite procedure is more accurate than the hard singular value thresholding procedure (note also that the oversampling factor  $\lambda$  is fixed, so the recovery error is more or less constant as  $n$  varies).

### Influence of the measurement error


 Figure 2: Recovery error  $\varepsilon := \left\| \mathbf{X} - \frac{\mathbf{X}^{\text{ht}}}{\|\mathbf{X}^{\text{ht}}\|_F} \right\|_F$  vs magnitude  $\mu := \|\mathbf{e}\|_1$  of prequantization error

In a second experiment, we assess the hard singular value thresholding procedure (22) in the presence of prequantization error  $\mathbf{e} \in \mathbb{R}^m$  to determine if the term  $\sqrt{\|\mathbf{e}\|_1}$  found in (23) reflects reality or if they are artifacts of our proofs. Figure 2 seems to suggest that the recovery error scales in reality like  $\|\mathbf{e}\|_1$ .



## 6.2 Full estimation

In this subsection, we verify empirically that incorporating thresholds before quantization allows one to estimate not only the direction but also the magnitude of low-rank matrices. However, since many binary measurements need to be made on low-rank matrices in order to observe any persuasive phenomenon, we exclude the demanding semidefinite procedures and focus on the swift hard singular values thresholding procedures. We work in the absence of prequantization error.

### Nonadaptive thresholds

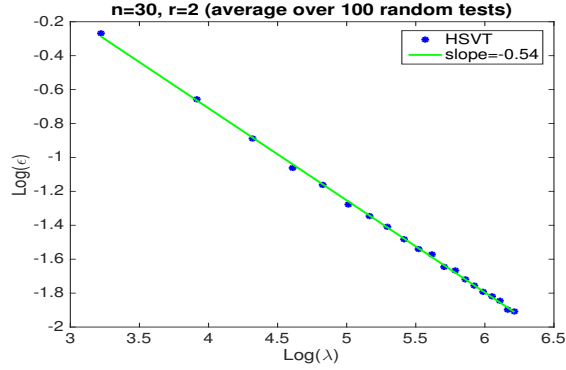


Figure 3: Recovery error  $\varepsilon := \|\mathbf{X} - \hat{\mathbf{X}}\|_F$  vs oversampling factor  $\lambda$  for the procedure (58)

We first validate that the procedure (58) results in a recovery error decaying like the inverse of a power of the oversampling factor  $\lambda$ , as indicated in (67). However, Figure 3 suggests that the power  $1/6$  (or even  $1/4$  if the remark after Theorem 1 is taken into account) is too pessimistic and could be replaced by  $1/2$ .

### Adaptive thresholds

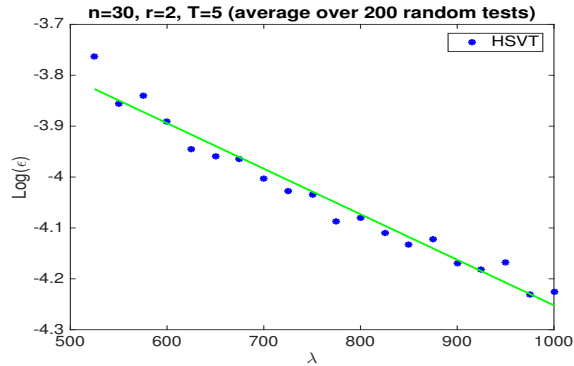


Figure 4: Recovery error  $\varepsilon := \|\mathbf{X} - \hat{\mathbf{X}}\|_F$  vs oversampling factor  $\lambda$  for the procedure (69)-(70)-(71)

Our final experiment confirms that the procedure based on binary measurements incorporating adaptive thresholds and described in (69)-(70)-(71) results in a reconstruction error that decays exponentially with the oversampling factor  $\lambda$ , see Figure 4.

## References

- [1] CVX Research, Inc. CVX: MATLAB software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx>, 2014.
- [2] R. Baraniuk, S. Foucart, D. Needell, Y. Plan, and M. Wotter. *Exponential decay of reconstruction error from binary measurements of sparse signals*. IEEE Transactions on Information Theory 63.6 (2017): 3368–3385.
- [3] R. Baraniuk, S. Foucart, D. Needell, Y. Plan, and M. Wotter. *One-bit compressive sensing of dictionary-sparse signals*. Information and Inference (to appear).
- [4] E. J. Candès and Y. Plan. *Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements*. IEEE Transactions on Information Theory 57.4 (2011): 2342–2359.
- [5] S. Foucart. *Flavors of Compressive Sensing*. In: Approximation Theory XV: San Antonio 2016, Springer Proceedings in Mathematics & Statistics (2017): 61–104.
- [6] S. Foucart and H. Rauhut. A Mathematical Introduction to Compressive Sensing. Birkhäuser (2013).
- [7] Y. Plan and R. Vershynin. *One-bit compressed sensing by linear programming*. Communications on Pure and Applied Mathematics 66.8 (2013): 1275–1297.
- [8] Y. Plan and R. Vershynin. *Robust 1-bit compressed sensing and sparse logistic regression: a convex programming approach*. IEEE Transactions on Information Theory 59.1 (2013): 482–494.
- [9] Y. Plan and R. Vershynin. *Dimension reduction by random hyperplane tessellations*. Discrete & Computational Geometry 51.2 (2014): 438–461.
- [10] Y. Plan and R. Vershynin. *The generalized Lasso with non-linear observations*. IEEE Transactions on Information Theory 62.3 (2016): 1528–1537.
- [11] G. Schechtman. *Two observations regarding embedding subsets of Euclidean spaces in normed spaces*. Advances in Mathematics 200.1 (2006): 125–135.