

# Achieving Oracle Statistical Property for Low-Rank Matrix Recovery with Application to Image Denoising

Anonymous CVPR submission

Paper ID 10048

## Abstract

This paper studies the classical low-rank matrix recovery problem, a topic of significant interest in computer vision. The nuclear norm minimization (NNM) method is a commonly used convex approach to addressing this problem. However, equal regularization of singular values in NNM restricts its effectiveness in practical applications, such as image restoration, where singular values carry distinct physical meanings and should be treated differently. This limitation led to the development of the weighted nuclear norm minimization (WNNM) approach, which offers enhanced recovery capabilities. Despite the advantages of NNM and WNNM, their reliance on  $\ell_1$  homotopy penalties can introduce a non-negligible recovery bias. In response to this challenge, we introduce non-convex penalties for low-rank matrix recovery in our work. We propose an algorithm within the majorization-minimization framework, leading to the adaptive weighted nuclear norm minimization (AWNNM) method. Despite the non-convex nature of the problem, we rigorously analyze the statistical properties of the iterations produced by AWNNM and demonstrate its ability to achieve the optimal statistical rate, known as the oracle rate. Our theoretical results are validated through numerical experiments on synthetic data. Extensive experiments on real-world datasets for image denoising demonstrate the efficacy of AWNNM.

## 1. Introduction

The problem of low-rank matrix recovery is focused on estimating a provided matrix with a low-rank alternative. This process aids in simplifying a complex dataset by uncovering inherent patterns, thus maintaining crucial details [19, 29, 30, 32, 38]. By enhancing data analysis efficiency and accuracy, low-rank matrix recovery reduces computational complexity and minimizes storage requirements, making it particularly beneficial for handling large-scale datasets.

Low-rank matrix recovery is a fundamental concept with broad applications in various fields such as computer vision [7, 27, 30, 47, 50], natural language processing [21], pattern recognition [44], and collaborative filtering [35]. For example, low-rank matrix recovery is employed in image matting to separate video sequences into the background and move foreground components [26]. Similarly, in image compression, low-rank matrix recovery reduces data storage requirements without compromising image quality [45, 49]. Furthermore, research indicates that matrices generated by non-local similar blocks in natural images demonstrate low-rank structures, which are advantageous for enhancing image restoration performance [39].

Over the past few decades, various techniques for low-rank matrix recovery have been developed, typically falling into two primary categories: matrix factorization methods [1, 10, 18, 20, 32, 36, 38] and rank penalty methods [2, 4, 9, 12, 22–25, 34, 42]. Matrix factorization methods involve breaking down a matrix into the product of two lower-dimensional matrices. However, these methods face challenges due to their non-convex nature, leading to complex optimization landscapes with various local minima and saddle points, making it difficult to achieve global optimality. Furthermore, the non-uniqueness of the decomposition and the requirement to specify the matrix rank in advance create challenges, mainly when the actual rank is unknown or inaccurately estimated. On the other hand, rank penalty methods like nuclear norm minimization (NNM) provide a convex relaxation of the rank minimization problem, making them more appealing in theoretical contexts because of their advantageous optimization properties.

Among the various methods based on rank penalties, NNM is arguably the most widely used for low-rank matrix recovery, as it provides the tightest convex relaxation of the matrix rank [31]. While convex formulations guarantee global optimality, they treat all singular values uniformly, which may limit their effectiveness in specific applications [13, 14, 28]. This uniform treatment overlooks prior knowledge that larger singular values can be more critical in representing the latent data. To address this limita-

tion, [13] introduced the weighted nuclear norm minimization (WNNM) method to incorporate varying importance across singular values. However, both WNNM and NNM are based on  $\ell_1$  homotopy palatines, which introduces non-negligible estimation bias that compromises the estimation accuracy. In contrast, non-convex penalties, such as the smoothly clipped absolute deviation (SACD) penalty [11] and minimax concave penalty (MCP) [48], are preferred for their superior estimation accuracy and variable selection consistency [40]. Despite encouraging empirical findings in various studies [13–15, 33], the theoretical underpinnings of non-convex penalties for low-rank matrix recovery remain largely unexplored. The quest for a theoretical rationale supporting non-convex substitutes for matrix rank continues to be an unresolved issue.

In this paper, we provide a comprehensive investigation of the non-convex low-rank matrix recovery problem, aiming to bridge the gap between practical applications and theoretical foundations. We introduce a novel low-rank matrix recovery that incorporates non-convex penalties. We devise an algorithm based on the majorization-minimization (MM) framework to address the proposed problem, leading to the adaptive weighted nuclear norm minimization (AWNNM) method. Our investigation illustrates that our estimator effectively utilizes singular values of significant magnitudes, leading to enhanced statistical convergence rates. Furthermore, under a mild assumption on the magnitude of the singular values, we establish that the proposed estimator possesses the oracle property. This property enables accurate recovery of the actual rank of the underlying matrix and attaining improved convergence rates. The primary contributions of this paper are threefold:

- We introduce a novel approach for low-rank matrix recovery based on non-convex penalties called AWNNM. We develop an MM-based algorithm that iteratively solves a sequence of subproblems, each admitting a closed-form solution.
- We theoretically analyze the statistical properties of all iterates generated by the proposed MM-based algorithm, proving that our method achieves the oracle statistical rate in the Frobenius norm under weak assumptions. These advancements offer both computational efficiency and strong theoretical guarantees, enhancing the robustness and applicability of matrix estimation techniques.
- We validate our theoretical results through comprehensive numerical experiments on synthetic and real-world datasets, further demonstrating the effectiveness of our method in image restoration tasks.

## 2. Low-Rank Matrix Recovery

In this section, we first revisit the classical NNM and WNNM methods. We then introduce our proposed method, AWNNM, which is based on non-convex penalties. Fol-

lowing this, we provide an overview of the MM algorithm framework and present our MM-based multistage relaxation algorithm. We also outline the necessary technical assumptions for the theoretical analysis before establishing the statistical convergence rate of the proposed method.

### 2.1. Existing Methods

The rank minimization problem is typically given by

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 + \lambda \text{rank}(\mathbf{X}), \quad (1)$$

where  $\mathbf{Y} \in \mathbb{R}^{m \times n}$  is the observed noisy matrix,  $\mathbf{X}$  is the matrix to be recovered,  $\|\cdot\|_F$  denotes the Frobenius norm, and  $\lambda$  is a regularization parameter that balances the data fidelity and rank minimization. However, solving (1) is NP-hard. To enhance tractability, the NNM method utilizes the nuclear norm, defined as the sum of the singular values of the matrix  $\mathbf{X}$ , to approximate the rank  $\text{rank}(\mathbf{X})$ . Under certain incoherence assumptions, NNM has been shown to recover a near-optimal low-rank solution [3]. The NNM formulation is:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 + \lambda \sum_{i=1}^d \sigma_i(\mathbf{X}), \quad (146)$$

where  $\sigma_i(\mathbf{X})$  representing the  $i$ -th singular values of  $\mathbf{X}$  and  $d = \min\{m, n\}$ .

In many cases, the rows (or columns) of the matrix lie within a low-dimensional subspace, where the primary singular values capture the dominant projection directions. These prominent singular values are crucial for preserving the data structure and should, thus, be assigned smaller penalty coefficients. To address this, WNNM [13] introduces a weighting scheme for the singular values. This modification provides greater flexibility and can significantly improve matrix recovery, particularly in scenarios with structured noise or matrices with varying significance. The optimization problem of WNNM is defined as:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 + \lambda \sum_{i=1}^d w_i \sigma_i(\mathbf{X}), \quad (160)$$

where  $w_i$  denotes the non-negative weight for  $\sigma_i(\mathbf{X})$ . This weighted approach enables more precise rank recovery, enabling a more accurate recovery of the actual low-rank structure compared to the standard nuclear norm minimization method [17].

### 2.2. Proposed Method via Non-Convex Penalties

Both the NNM and WNNM methods are  $\ell_1$ -based method, which introduces non-negligible estimation bias, thus compromising solution accuracy [28]. In contrast, non-convex penalties are preferred due to their superior estimation accuracy and variable selection consistency, as they adaptively

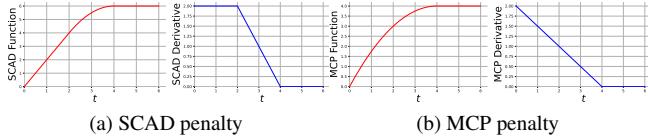


Figure 1. Illustration of the sparsity-inducing functions  $p_\lambda$  with their derivatives ( $\lambda = 2$  and  $a = 2$  for both SCAD and MCP)

adjust weights during the optimization process [40]. To mitigate the estimation bias inherent in  $\ell_1$  homotopy penalties, we propose a low-rank minimization approach using the non-convex penalty, also known as the adaptive weighted nuclear norm penalty. Specifically, we aim to solve the following non-convex optimization problem:

$$\min_{\mathbf{X}} \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \sum_{i=1}^d p_\lambda(\sigma_i(\mathbf{X})), \quad (2)$$

where  $p_\lambda$  is a univariate non-convex function with parameter  $\lambda$ . In this paper, we consider a class of functions  $p_\lambda$  satisfying the following assumption.

**Assumption 1.** The penalty function  $p_\lambda: \mathbb{R}^+ \rightarrow \mathbb{R}^+$  satisfies the following conditions:

- $p_\lambda(t)$  is symmetric and non-decreasing on  $[0, +\infty)$ , with  $p_\lambda(0) = 0$ , and is differentiable almost everywhere on  $(0, +\infty)$ .
- For all  $t_1 \geq t_2 \geq 0$ , it holds that  $0 \leq p'_\lambda(t_1) \leq p'_\lambda(t_2) \leq \lambda$ , and  $\lim_{t \rightarrow 0} p'_\lambda(t) = \lambda$ ;
- There exists an  $\alpha > 0$  such that  $p'_\lambda(t) = 0$  for  $t \geq \alpha\lambda$ .
- There exists some  $c \in (0, \alpha)$  such that  $p'_\lambda(c\lambda) \geq \frac{\lambda}{2}$ .

In Assumption 1, the first three conditions ensure sparsity and unbiasedness, while the last condition is introduced for the following theoretical analysis, which can always hold due to  $p'_\lambda(0) = \lambda$  and  $p'_\lambda(\alpha\lambda) = 0$ . Typical examples of the non-convex penalty function  $p_\lambda$  in Assumption 1 include:

- Smoothly clipped absolute deviation (SCAD) penalty: SCAD function is given by

$$p_\lambda(t) = \begin{cases} \lambda t, & 0 \leq t \leq \lambda \\ \frac{-t^2 + 2\lambda at - \lambda^2}{2(a-1)}, & \lambda < t \leq a\lambda \\ \frac{\lambda^2(a+1)}{2}, & t > a\lambda \end{cases}$$

for some  $a > 2$ . The choice  $a = 3.7$  is suggested in [11] based on a Bayesian argument.

- Minimax concave penalty (MCP): MCP function  $p_\lambda(\cdot)$  is defined as follows:

$$p_\lambda(t) = \begin{cases} \lambda t - \frac{t^2}{2a}, & 0 \leq t \leq a\lambda \\ \frac{1}{2}at^2, & t > a\lambda \end{cases}$$

for some  $a > 1$ .

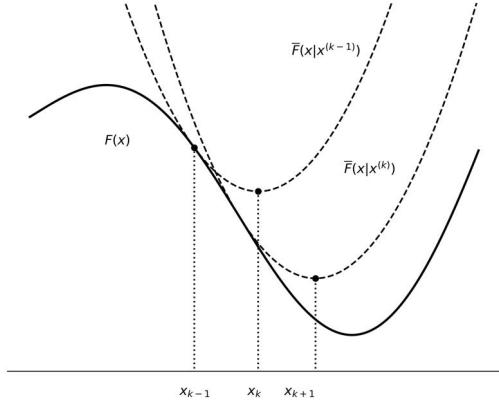


Figure 2. The MM procedure.

We give a pictorial illustration of these two functions together with their derivatives in Figure 1.

### 3. Optimization Algorithm

#### 3.1. Brief Review of the MM Algorithm Framework

The MM algorithm [16, 37] is an iterative optimization procedure involving two primary steps: majorization and minimization. Given a loss function  $F(\mathbf{x})$  to minimize, initialized at  $\mathbf{x}^{(0)}$ , the MM algorithm generates a sequence of feasible points  $\{\mathbf{x}^{(k)}\}_{k \geq 1}$ . At each iteration  $k$ , starting from  $\mathbf{x}^{(k-1)}$ , the majorization step involves constructing a surrogate function  $\bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)})$  that locally approximates the objective function  $F(\mathbf{x})$ , satisfying the following conditions:

$$\bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)}) \geq F(\mathbf{x}), \quad 219$$

$$\bar{F}(\mathbf{x}^{(k-1)} | \mathbf{x}^{(k-1)}) = F(\mathbf{x}^{(k-1)}). \quad 220$$

Next, in the minimization step, the algorithm update  $\mathbf{x}^{(k)}$  by minimizing the surrogate function:

$$\mathbf{x}^{(k)} \in \arg \min_{\mathbf{x}} \bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)}). \quad 223$$

Starting from an initial point  $\mathbf{x}^{(0)}$ , the alternating process of MM generates a sequence of iterates  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$  with non-increasing objective function values. The algorithm can be executed until a convergence criterion is met. See Figure 2 for an illustration.

#### 3.2. An MM-based Algorithm for Problem (2)

In this section, we derive an MM algorithm for problem (2). In the majorization step, we construct a surrogate function, which locally linearizes the non-convex penalty function  $p_\lambda$ . Specifically, we find an upper bound for  $p_\lambda(t)$  using its first-order Taylor expansion around the iterate  $x_0$ :

235  $p_\lambda(t) \leq p_\lambda(t_0) + p'_\lambda(t_0)(t - t_0).$

236 Based on the above result, we have the following inequality  
237 at iterate  $\mathbf{X}^{(k-1)}$ :

238 
$$\sum_{i=1}^d p_\lambda(\sigma_i(\mathbf{X})) \leq \sum_{i=1}^d p_\lambda(\sigma_i(\mathbf{X}^{(k-1)}))$$
  
239 
$$+ \sum_{i=1}^d w_i^{(k-1)} (\sigma_i(\mathbf{X}) - \sigma_i(\widehat{\mathbf{X}}^{(k-1)})),$$

240 where  $w_i^{(k-1)} = p'_\lambda(\sigma_i(\widehat{\mathbf{X}}^{(k-1)}))$ . This leads to the fol-  
241 lowing surrogate function:

242 
$$\bar{F}(\mathbf{X} | \widehat{\mathbf{X}}^{(k-1)}) = \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \sum_{i=1}^d w_i^{(k-1)} \sigma_i(\mathbf{X}).$$

243 In the minimization step, we compute the iterate that  
244 minimizes the surrogate function derived in the majoriza-  
245 tion step. This is achieved by solving the following opti-  
246 mization problem:

247 
$$\widehat{\mathbf{X}}^{(k)} \in \arg \min_{\mathbf{X}} \left\{ \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2 + \sum_{i=1}^d w_i^{(k-1)} \sigma_i(\mathbf{X}) \right\}. \quad (3)$$

248 The problem (3) is inherently non-convex. However, a  
249 closed-form solution can be obtained under Assumption 1,  
250 for the reason it states that  $p'(t)$  is non-increasing for  $t \geq 0$ .

251 Let  $\mathbf{Y} \in \mathbb{R}^{m \times n}$  be a matrix with singular value  
252 decomposition (SVD) given by  $\mathbf{Y} = \mathbf{U}\Sigma\mathbf{V}^\top$ , where  
253  $\mathbf{U} \in \mathbb{R}^{m \times d}$  contains the left singular vectors,  $\mathbf{V} \in \mathbb{R}^{n \times d}$  contains the right singular vectors, and  $\Sigma = \text{diag}(\sigma_i(\mathbf{Y}), i = 1, \dots, d) \in \mathbb{R}^{d \times d}$  is a diagonal matrix of  
254 singular values. The singular values are ordered such that  
255  $\sigma_1(\mathbf{Y}) \geq \sigma_2(\mathbf{Y}) \geq \dots \geq \sigma_d(\mathbf{Y}) \geq 0$ . Due to this  
256 ordering, the weight sequence  $\{w_1, w_2, \dots, w_d\}$  satisfies  
257  $0 \leq w_1 \leq w_2 \leq \dots \leq w_d$ , following the anti-monotonic  
258 property. This ordering enables a closed-form solution us-  
259 ing a weighted singular value thresholding approach in the  
260 following lemma.

261 **Lemma 2.** [[6], Theorem 2.3] For  $\mathbf{X} \in \mathbb{R}^{m \times n}$  and  $0 \leq$   
262  $w_1 \leq w_2 \leq \dots \leq w_d$ , a global solution to the optimiza-  
263 tion problem in (3) is given by the weighted singular value  
264 thresholding

265 
$$\widehat{\mathbf{X}} := \mathcal{T}(\mathbf{Y}), \quad (4)$$

266 where  $\mathbf{U}\Sigma\mathbf{V}^\top$  is the SVD of  $\mathbf{Y}$ ,  $\mathcal{T}(\mathbf{Y}) = \mathbf{U}\mathcal{T}(\Sigma)\mathbf{V}^\top$ ,  
267  $\mathcal{T}(\Sigma) = \text{Diag}(\max(0, \sigma_i(\mathbf{Y}) - w_i), i = 1, \dots, d)$ .

268 The MM-based multistage convex relaxation algorithm  
269 is summarized in Algorithm 1. For simplicity, we start with  
270 a trivial initial value  $\widehat{\mathbf{X}}^{(0)} = \mathbf{0}$ .

---

**Algorithm 1:** Multistage Relaxation Algorithm for Solv-  
ing (2).

---

**Input:** Data matrix  $\mathbf{Y}$ , tuning parameter  $\lambda$ ;

**Initialize**  $\widehat{\mathbf{X}}^{(0)} = \mathbf{0}$ ,  $k = 1$ ;

**for**  $k = 1, 2, \dots, K$  **do**

update  $w_i^{(k-1)} = p'_\lambda(\sigma_i(\widehat{\mathbf{X}}^{(k-1)}))$

update  $\widehat{\mathbf{X}}^{(k)}$  by solving Problem (3) via Lemma 2;

$k = k + 1$ ;

**end for**

**Output:**  $\widehat{\mathbf{X}}^{(K)}$ .

---

## 4. Theoretical Analysis

We define the true low-rank matrix as  $\mathbf{X}^*$ . Let  $\mathcal{F}$  and its or-  
thogonal complement  $\mathcal{F}^\perp$  be the subspaces corresponding to the SVD of  $\mathbf{X}^*$ , defined as follows:

$$\mathcal{F}(\mathbf{X}^*) = \{\mathbf{W} \mid \mathcal{R}(\mathbf{W}) \subseteq \text{sp}(\mathbf{V}^*), \mathcal{C}(\mathbf{W}) \subseteq \text{sp}(\mathbf{U}^*)\},$$
  
$$\mathcal{F}^\perp(\mathbf{X}^*) = \{\mathbf{W} \mid \mathcal{R}(\mathbf{W}) \perp \text{sp}(\mathbf{V}^*), \mathcal{C}(\mathbf{W}) \perp \text{sp}(\mathbf{U}^*)\},$$

where  $\mathcal{R}(\mathbf{W})$  and  $\mathcal{C}(\mathbf{W})$  represent the row and column space of  $\mathbf{W}$ , respectively, and  $\text{sp}(\mathbf{V}^*)$ ,  $\text{sp}(\mathbf{U}^*)$  denote the subspace spanned by  $\mathbf{V}^*$  and  $\mathbf{U}^*$ , respectively. For brevity, we will use the shorthand notations of  $\mathcal{F}$  and  $\mathcal{F}^\perp$  when the dependence on  $\mathbf{X}^*$  is clear from the context. The projection operators onto the subspace  $\mathcal{F}$  and  $\mathcal{F}^\perp$  are denoted by  $\Pi_{\mathcal{F}}(\cdot)$  and  $\Pi_{\mathcal{F}^\perp}(\cdot)$ , respectively. Given matrix  $\mathbf{A}$ , the pro-  
jections are defined as:

$$\Pi_{\mathcal{F}}(\mathbf{A}) = \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{A} \mathbf{V}^* \mathbf{V}^{*\top},$$

$$\Pi_{\mathcal{F}^\perp}(\mathbf{A}) = (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{A} (\mathbf{I} - \mathbf{V}^* \mathbf{V}^{*\top}),$$

where  $\mathbf{I}$  is an identity matrix of the appropriate dimension.

Next, we introduce a useful assumption.

**Assumption 3** (Minimal Signal Strength). *The true matrix  $\mathbf{X}^*$  satisfies*

$$\min \{\sigma_i(\mathbf{X}^*) \mid i \in \mathcal{S}^*\} \geq (\alpha + c) \lambda \gtrsim \lambda,$$

where  $\mathcal{S}^* = \{i \mid \sigma_i(\Pi_{\mathcal{F}}(\mathbf{X}^*)) \neq 0\}$  is the index set associ-  
ated with the true rank subspace and  $\alpha$  and  $c$  are constants  
defined in Assumption 1.

Assumption 3, referred to as the minimum signal strength condition, is mild and commonly employed [14]. We now present the main theorem, which establishes the contraction property of the solution path  $\{\widehat{\mathbf{X}}^{(k)}\}_{k \geq 1}$ . Define  $f(\mathbf{X}) = \frac{1}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2$ .

**Theorem 4.** Suppose that Assumptions 1 and 3 hold. If  $\lambda \geq 4 \|\nabla f(\mathbf{X}^*)\|_2$ , then the optimal solution  $\widehat{\mathbf{X}}^{(k)}$  satisfies the following  $\delta$ -contraction property:

$$\|\widehat{\mathbf{X}}^{(k)} - \mathbf{X}^*\|_F \leq \underbrace{\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F}_{\text{oracle rate}} + \delta \underbrace{\|\widehat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F}_{\text{contraction}}, \quad (303)$$

304 for  $1 \leq k < K$ , where  $\mathcal{F}_{\mathcal{S}^*}$  is a subspace of  $\mathcal{F}$  associated  
 305 with  $\mathcal{S}^*$  and  $\delta \in (0, 1)$  is the contraction parameter.

306 **Remark 5.** The oracle rate refers to the statistical convergence  
 307 rate of the oracle estimator, which has prior knowledge of the true rank subspace  $\mathcal{F}(\mathbf{X}^*)$ . The oracle estimator  
 308  $\widehat{\mathbf{X}}^O$  is defined as  
 309

$$310 \quad \widehat{\mathbf{X}}^O = \arg \min_{\mathbf{X} \in \mathcal{F}(\mathbf{X}^*)} f(\mathbf{X}).$$

311 According to the definition, it is easy to obtain that  $\widehat{\mathbf{X}}^O$  sat-  
 312 isfies  $\|\widehat{\mathbf{X}}^O - \mathbf{X}^*\|_F \lesssim \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F$ .

313 Theorem 4 demonstrates that the estimation error be-  
 314 tween the optimal solution  $\widehat{\mathbf{X}}^{(k)}$  and the true value  $\mathbf{X}^*$  is  
 315 bounded by two terms, namely, the oracle rate and a con-  
 316 traction term. We now provide the explicit statistical rate of  
 317 convergence under the sub-Gaussian design.

318 **Corollary 6.** Let  $\mathbf{Y} = \mathbf{X}^* + \mathbf{E}$  and  $\mathbf{E}$  be a sub-Gaussian  
 319 random matrix with entries  $E_{ij}$  that are sub-Gaussian ran-  
 320 dom variables with zero mean and covariance  $\varepsilon$ . Suppose  
 321 that Assumptions 1 and 3 hold. If  $\lambda \asymp \sqrt{\log mn}$ , then the  
 322 optimal solution  $\widehat{\mathbf{X}}^{(1)}$  satisfies

$$323 \quad \|\widehat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F \lesssim \sqrt{r^* \log mn},$$

324 with high probability.

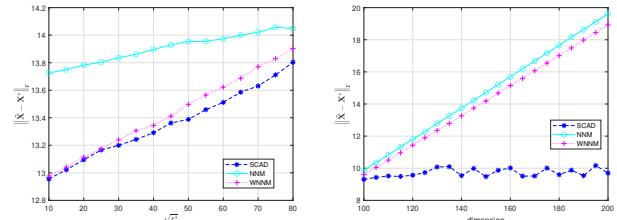
325 Corollary 6 follows directly from Theorem 4 for  $k = 1$ .  
 326 Notably, the first subproblem (i.e.,  $k = 1$ ) corresponds to  
 327 the nuclear norm penalized matrix recovery problem, which  
 328 results in a convergence rate of  $\sqrt{r^* \log mn}$ . In Theorem 4,  
 329 the contraction property is induced by the MM-based mul-  
 330 tistage convex relaxation algorithm. To achieve the oracle  
 331 rate, it is crucial to appropriately select a sufficiently large  
 332 value for  $K$ . The following result extends this analysis.

333 **Corollary 7.** Let  $\mathbf{Y} = \mathbf{X}^* + \mathbf{E}$  and  $\mathbf{E}$  be a sub-Gaussian  
 334 random matrix with entries  $E_{ij}$  are sub-Gaussian random  
 335 variables with zero mean and covariance  $\varepsilon$ . Suppose that  
 336 Assumptions 1 and 3 hold. If  $\lambda \asymp \sqrt{\log mn}$ , and  $K \gtrsim$   
 337  $\log \lambda \gtrsim \log \log mn$ , then the optimal solution  $\widehat{\mathbf{X}}^{(K)}$  satis-  
 338 fies

$$339 \quad \|\widehat{\mathbf{X}}^{(K)} - \mathbf{X}^*\|_F \lesssim \sqrt{r^*},$$

340 with high probability.

341 Corollary 7 follows directly from Theorem 4, which im-  
 342 plies that under weak assumptions, at most approximately  
 343  $\log \log mn$  convex problems need to be solved to achieve  
 344 the oracle rate  $\sqrt{r^*}$ . This result highlights that the pro-  
 345 posed method attains a faster statistical rate of convergence,  
 346 matching the oracle rate, compared to existing methods that  
 347 employ nuclear norm penalties in the Frobenius norm.



(a) Estimation error via varying  $r^*$  (b) Estimation error via varying dimension

Figure 3. The oracle rate for low-rank matrix recovery

## 5. Numerical Simulations

### 5.1. Simulations on Theoretical Properties

In this section, we evaluate the practical recovery performance of the proposed method and present numerical results for the low-rank matrix recovery problem. The oracle rate derived for the proposed method is provided in 3. Figure 3a illustrates the recovery performance for a matrix with dimensions  $100 \times 200$ , where each entry is sampled from a normal distribution. A scaled noise with a parameter of 0.1 is added to the matrix. We examine the effect of varying rank, and the results demonstrate that the proposed non-convex penalty outperforms both the weighted nuclear norm and the nuclear norm. This observation is consistent with our theoretical findings. Figure 3b presents results for a square matrix of different dimensions, with each entry generated from a normal distribution. A scaled noise with a parameter of 0.1 is also added, and the matrix rank is fixed at 50. The results confirm that the proposed non-convex penalty maintains superior performance compared to the weighted nuclear norm and the nuclear norm, highlighting that the recovery error remains stable as the matrix dimension increases. This contrasts with the weighted nuclear norm and nuclear norm, where the error increases with higher matrix dimensions.

### 5.2. Image Denoising Applications

To validate the proposed estimator, we apply it to a classical low-level vision task: image denoising. Image denoising is a critical pre-processing step in a variety of vision applications and serves as an effective framework for evaluating statistical image modeling techniques. This task aims to reconstruct the original image denoted as  $\mathbf{x}$ , from its noisy observation  $\mathbf{y}$ , where the noise component  $\mathbf{e}$  is assumed to be additive white Gaussian noise, such that  $\mathbf{y} = \mathbf{x} + \mathbf{e}$ .

Natural images often display significant redundancy, with recurring patterns and textures across different regions. This prior fact motivates the use of nonlocal self-similarity (NSS) techniques in image denoising. NSS relies on the observation that many patches within a natural image share similar structures, allowing denoising algorithms [ref] to re-

**Algorithm 2:** Image Denoising by AWNNM.

---

**Input:** Noisy image  $\mathbf{y}$ , tuning parameter  $\lambda$ ;  
**Initialize**  $\hat{\mathbf{x}}^{(0)} = \mathbf{y}$ ,  $\mathbf{y}^{(0)} = \mathbf{y}$ ;  
**for**  $k = 1, 2, \dots, K$  **do**  
    Residual correction  $\mathbf{y}^{(k)} = \hat{\mathbf{x}}^{(k-1)} + \delta (\mathbf{y} - \hat{\mathbf{y}}^{(k-1)})$ ;  
    **for** each patch  $\mathbf{y}_i$  in  $\mathbf{y}^{(k)}$  **do**  
        Search similar patches and group them into  $\mathbf{Y}_i$ ;  
        Apply Algorithm 1 to  $\mathbf{Y}_i$  to estimate  $\mathbf{X}_i$ ;  
    **end for**  
    Elementwisely average the estimates  $\mathbf{X}_i$  for each pixel to form the denoised image  $\hat{\mathbf{x}}^{(k)}$ ;  
**end for**  
**Output:**  $\hat{\mathbf{x}}^{(K)}$

---

387 construct each patch effectively by referencing information  
 388 from similar patches. By aggregating these reconstructed  
 389 patches, NSS techniques improve the overall quality of de-  
 390 noising.

391 The structural redundancy among similar patches intro-  
 392 duces a low-rank tendency within each group of patches,  
 393 allowing specific patch matrices to be represented with  
 394 fewer components, retaining essential details while reduc-  
 395 ing noise. Specifically, for a given local patch  $\mathbf{y}_i$  in the im-  
 396 age, similar patches from other regions can be grouped into  
 397 a matrix  $\mathbf{Y}_i$ . Each group of similar patches, denoted as  $\mathbf{Y}_i$ ,  
 398 can be approximated by a low-rank component  $\mathbf{X}_i$  addi-  
 399 tion to Gaussian noise  $\mathbf{E}_i$ . We apply the proposed estimator  
 400 to  $\mathbf{Y}_i$  to estimate  $\mathbf{X}_i$ , thereby performing image denoising.  
 401 This process is formulated as follows:

$$402 \quad \hat{\mathbf{X}}_i \in \arg \min_{\mathbf{X}_i} \frac{1}{2} \|\mathbf{X}_i - \mathbf{Y}_i\|_{\text{F}}^2 + \sum_j p_{\lambda_i}(\sigma_j(\mathbf{X}_i)). \quad (5)$$

403 By applying this denoising procedure to each patch and  
 404 subsequently aggregating the denoised patches, we recon-  
 405 struct the image  $\mathbf{x}$ . In practice, multiple iterations of this  
 406 reconstruction process across all image patches can further  
 407 enhance the denoising results. The whole denoising algo-  
 408 rithm is summarized in Algorithm 2.

409 We compare the proposed AWNNM based image  
 410 restoration algorithm with several classical restoration  
 411 methods, including NNM, WNNM, BM3D [8], EPLL [51],  
 412 PCLR [5], PGPD [43], STROLLR [41], and SNSS [46].

- 413 • BM3D: Groups similar patches into 3D blocks for col-  
 laborative filtering, effectively preserving details, but may  
 struggle with complex noise patterns.
- 414 • NNM minimizes the nuclear norm to promote low-rank  
 structures. It works well in low-rank areas but introduces  
 bias in high-rank regions.
- 415 • EPLL maximizes the expected log-likelihood of patches  
 under a probabilistic prior, yielding good results but can  
 be computationally intensive.

- WNNM extends NNM by weighting singular values, pro-  
 viding flexibility for different noise levels but still limited  
 by non-uniform noise. 422
- PCLR clusters patches for localized low-rank regulariza-  
 tion, effective but sensitive to cluster accuracy in complex  
 regions. 425
- PGPD uses patch group priors to reduce noise, preserving  
 structure but reliant on effective priors. 428
- STROLLR maintains non-local similarities for structural  
 fidelity yet may struggle with unique features. 430
- SNSS preserves self-similarity, ideal for textures but less  
 effective in isolated features. 431
- AWNNM, the proposed algorithm, extends WNNM with  
 adaptive weights via non-convex penalties, which en-  
 hances flexibility and noise handling across various struc-  
 tures. 432

**5.2.1. Parameter Setting**

The proposed AWNNM algorithm involves several vital parameters, namely the regularization parameter  $\delta$ , the number of iterations  $K$ , and patch size, each tuned based on the noise level. The iterative regularization parameter  $\delta$  is fixed to 0.1 for all noise levels. Both  $K$  and patch size vary with noise level: larger patches and a higher number of iterations are required for higher noise levels to achieve optimal results. In practice, we set the patch size and iteration number as follows: for  $\sigma_n \leq 20$ ,  $20 < \sigma_n \leq 40$ ,  $40 < \sigma_n \leq 60$  and  $60 < \sigma_n$ , the patch sizes are  $6 \times 6$ ,  $7 \times 7$ ,  $8 \times 8$  and  $9 \times 9$ , respectively. The number of iterations  $K$  is set to 8, 10, 12, and 12, according to the noise levels, on the noise levels. For competing methods, we obtained source codes from the original authors and used default parameter settings to ensure a fair comparison across all methods. Our non-convex approach supports various penalty functions. In this experiment, we use the SCAD function with  $a = 3.7$  to generate adaptive weights. The tuning parameter  $\lambda$  is chosen by five-fold cross-validation.

**5.2.2. Results on Synthetic Images**

We benchmark against several classical denoising methods on 20 commonly used test images with synthetic noises. Additive White Gaussian Noise (AWGN) with zero mean and variance  $\sigma_n^2$  is applied to these test images to simulate various noise conditions. We test four different noise levels: low noise ( $\sigma_n = 10$ ), medium noise ( $\sigma_n = 30$  and  $\sigma_n = 50$ ), and high noise ( $\sigma_n = 100$ ).

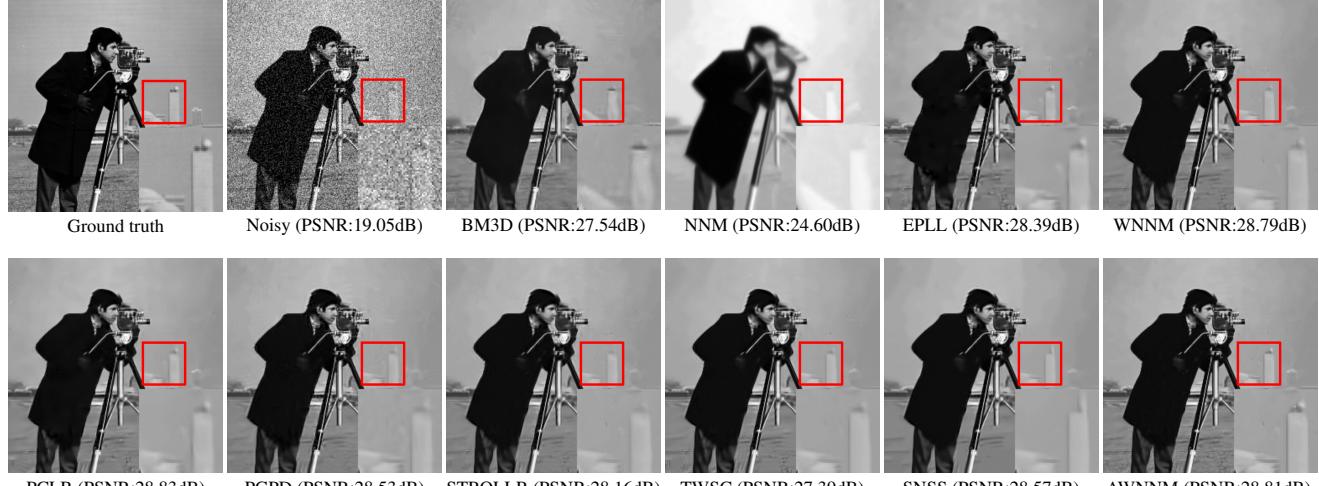
The Peak Signal-to-Noise Ratio (PSNR) is the primary metric to assess the quality of the denoised images quantitatively. PSNR is defined as follows:

$$\text{PSNR} = 10 \cdot \log \left( \frac{\text{MAX}^2}{\text{MSE}} \right),$$

where MAX represents the maximum possible pixel value of the image, and MSE (Mean Squared Error) is the average

Table 1. Denoising results (PSNR) by different methods.

	$\sigma_n = 10$										$\sigma_n = 30$										
	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	
Alley	32.29	26.47	32.45	32.76	32.67	32.45	32.22	32.78	32.41	<b>32.79</b>	26.47	24.70	27.11	27.46	27.36	27.19	26.88	27.29	27.27	<b>27.47</b>	
Baboon	30.43	23.99	30.32	30.50	30.47	30.23	30.29	30.42	30.07	<b>30.52</b>	24.56	22.61	24.04	24.64	24.56	24.23	24.42	24.42	24.30	<b>24.66</b>	
Barbara	30.92	28.34	32.63	34.54	34.23	33.60	33.89	33.79	34.49	<b>34.57</b>	25.87	25.69	26.91	28.60	28.09	27.79	28.37	27.94	28.54	<b>28.63</b>	
Boat	33.49	27.27	33.18	33.46	33.47	33.12	33.00	33.45	33.17	<b>33.58</b>	27.05	25.06	27.60	27.70	27.77	27.48	27.32	27.32	26.95	27.49	<b>27.80</b>
Book	33.69	30.48	36.28	36.78	36.71	36.22	35.97	36.39	36.23	<b>36.79</b>	28.25	28.19	31.31	31.58	31.83	31.27	30.98	31.35	31.70	<b>31.86</b>	
Cameraman	33.08	26.53	34.40	34.46	34.46	34.14	33.79	34.04	34.30	<b>34.50</b>	27.54	24.60	28.39	28.79	28.83	28.53	28.16	27.39	28.57	<b>28.81</b>	
Couple	34.14	27.12	33.11	33.44	33.44	33.24	32.75	32.62	33.21	<b>33.45</b>	28.53	24.81	27.37	27.48	27.48	27.32	27.17	26.87	27.28	<b>27.49</b>	
Dice	36.06	35.77	40.54	41.08	41.20	40.04	40.34	40.74	41.02	<b>41.24</b>	31.17	32.04	34.58	34.92	35.38	35.12	34.50	36.68	35.28	<b>35.39</b>	
F16	34.04	28.06	34.32	34.59	34.63	34.40	33.96	33.94	34.42	<b>34.64</b>	28.31	25.73	28.65	28.74	28.79	28.62	28.27	27.12	28.62	<b>28.83</b>	
Fingerprint	30.17	27.27	30.08	30.66	30.45	30.45	30.36	29.91	30.42	<b>30.69</b>	24.16	23.34	24.02	24.67	24.54	24.57	24.69	24.49	24.57	<b>24.80</b>	
Girl	33.08	36.18	39.35	39.59	39.66	38.75	39.44	39.44	<b>39.73</b>	39.67	27.27	32.06	34.12	34.03	34.53	34.52	34.51	31.80	34.34	<b>34.56</b>	
Hallway	39.05	34.21	38.85	39.58	39.45	38.89	38.97	39.79	39.43	<b>39.81</b>	34.49	30.22	33.23	33.62	33.83	33.59	33.08	32.07	<b>34.02</b>	33.96	
House	32.82	32.36	35.76	36.92	36.83	36.56	36.67	36.35	36.89	<b>36.94</b>	26.94	29.29	31.41	32.30	32.17	32.24	31.67	31.31	<b>32.68</b>	32.61	
Pentagon	33.35	26.68	30.77	31.69	31.40	30.85	30.85	31.18	31.19	<b>33.70</b>	26.38	24.40	25.50	26.29	25.84	25.63	25.69	24.67	25.94	<b>26.41</b>	
Peppers	33.45	28.16	34.58	34.98	34.96	34.65	34.32	34.94	34.86	<b>35.01</b>	28.35	25.81	29.27	29.43	29.56	28.60	28.90	28.10	29.39	<b>29.58</b>	
Plaza	34.19	30.62	34.52	34.68	34.64	34.27	34.20	34.37	34.33	<b>34.71</b>	28.55	28.23	29.95	29.98	30.03	29.87	29.52	28.03	29.78	<b>30.10</b>	
Statue	30.32	28.65	35.48	35.68	35.74	35.30	35.02	35.34	35.45	<b>35.69</b>	24.77	26.58	29.59	29.57	29.68	29.55	29.11	26.44	29.45	<b>29.72</b>	
Traffic	32.73	26.59	33.01	33.32	33.30	32.88	32.52	32.79	32.91	<b>33.33</b>	25.86	24.52	27.38	27.55	27.56	25.30	26.87	25.21	27.29	<b>27.57</b>	
Validemossa	32.79	26.05	33.49	33.79	33.82	33.49	33.15	33.81	33.54	<b>33.85</b>	27.84	24.18	27.51	27.69	27.72	27.50	27.15	27.13	27.50	<b>27.85</b>	
Yard	34.36	28.90	34.16	34.77	34.53	34.36	34.26	34.19	34.53	<b>34.79</b>	29.97	26.53	29.07	29.21	29.22	29.09	29.02	28.86	29.00	<b>30.01</b>	
AVE.	33.22	28.98	34.36	34.86	34.80	34.39	34.31	34.51	34.63	<b>35.01</b>	27.62	26.43	28.87	29.22	29.24	28.81	28.82	27.81	29.15	<b>29.40</b>	
	$\sigma_n = 50$										$\sigma_n = 100$										
	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	
Alley	23.84	23.25	25.05	25.36	25.27	25.17	24.91	25.35	25.23	<b>25.42</b>	20.67	21.10	22.63	22.94	22.84	22.83	22.49	22.36	23.00	<b>23.01</b>	
Baboon	22.81	21.60	22.47	22.70	22.68	22.51	22.64	22.60	22.51	<b>22.72</b>	20.56	19.99	20.85	20.92	20.87	20.97	20.94	20.86	20.85	<b>20.98</b>	
Barbara	24.02	23.59	25.11	26.13	25.61	22.72	25.85	25.97	26.07	<b>26.19</b>	21.94	20.54	22.52	22.82	22.78	22.72	22.57	22.07	22.85	<b>22.91</b>	
Boat	25.95	23.16	25.24	25.38	25.36	25.17	25.06	24.94	25.29	<b>25.96</b>	21.54	20.92	22.48	22.61	22.63	22.62	22.37	21.97	<b>22.68</b>	22.67	
Book	25.68	26.45	28.99	29.39	29.48	29.05	28.56	29.24	29.35	<b>29.51</b>	21.39	23.54	25.83	26.59	26.34	26.09	25.53	26.30	<b>26.75</b>	26.69	
Cameraman	25.91	23.12	26.09	26.45	26.56	26.46	25.83	26.33	26.38	<b>26.67</b>	21.71	21.05	22.94	23.37	23.48	23.23	22.71	23.24	23.39	<b>23.51</b>	
Couple	25.39	23.08	24.97	25.10	25.04	25.01	24.86	24.86	25.00	<b>25.66</b>	20.21	20.78	22.35	22.52	22.47	22.48	22.13	21.91	22.42	<b>22.56</b>	
Dice	31.69	29.49	31.65	32.76	32.52	32.52	31.44	32.64	33.04	<b>33.30</b>	23.52	26.06	27.99	29.17	28.54	28.73	27.53	28.80	<b>29.67</b>	29.61	
F16	25.04	23.96	26.10	26.34	26.21	26.12	25.71	26.28	26.24	<b>26.37</b>	20.46	21.52	23.06	23.39	23.24	23.30	23.09	22.51	23.27	<b>23.47</b>	
Fingerprint	26.03	20.62	21.52	22.54	22.33	22.30	22.51	22.58	22.38	<b>22.61</b>	19.79	18.11	17.59	20.23	19.53	19.57	20.25	19.86	19.88	<b>20.41</b>	
Girl	24.68	29.96	31.66	32.05	32.30	32.39	31.79	30.57	32.16	<b>32.32</b>	21.51	27.84	28.41	29.13	29.03	29.11	28.22	28.32	<b>29.54</b>	29.48	
Hallway	30.24	28.00	30.60	31.35	31.09	31.06	30.12	28.18	31.48	<b>31.38</b>	22.35	25.04	27.15	27.92	27.46	27.64	26.54	27.14	<b>28.31</b>	28.31	
House	24.26	26.59	29.02	30.30	29.78	29.93	29.06	30.14	<b>30.52</b>	30.43	24.09	23.34	25.50	26.71	25.96	26.16	25.35	26.10	<b>27.09</b>	26.98	
Pentagon	24.59	22.53	23.71	24.35	23.96	24.02	23.96	23.69	24.26	<b>24.37</b>	20.82	21.17	21.83	22.10	21.95	22.13	21.85	21.77	22.15	<b>22.16</b>	
Peppers	26.97	23.60	26.76	26.93	27.03	26.80	26.51	27.04	26.94	<b>27.11</b>	22.23	20.50	23.20	23.50	23.69	23.44	22.79	22.10	23.71	<b>23.72</b>	
Plaza	25.65	26.68	28.18	28.22	28.30	28.22	27.78	28.01	28.24	<b>28.33</b>	20.67	24.98	25.85	25.83	25.83	26.03	25.50	25.81	25.83	<b>25.86</b>	
Statue	26.57	25.02	27.25	27.41	27.35	27.31	27.12	27.30	27.31	<b>27.42</b>	19.26	22.72	24.28	24.66	24.43	24.58	24.09	24.13	24.64	<b>24.69</b>	
Traffic	28.95	22.84	25.25	25.45	25.42	25.30	24.84	25.27	25.30	<b>25.49</b>	22.96	20.51	22.47	22.70	22.68	22.54	21.98	22.27	22.63	<b>22.72</b>	
Validemossa	27.28	22.57	25.06	25.28	25.29	25.15	24.88	25.27	25.19	<b>25.35</b>	20.32	20.32	21.89	22.30	22.26	22.15	21.94	22.26	22.32	<b>22.36</b>	
Yard	29.51	24.96	27.02	27.13	27.11	27.08	26.80	27.09	26.97	<b>27.17</b>	24.02	22.66	24.12	24.43	24.41	24.40	23.85	23.06	24.38	<b>24.49</b>	
AVE.	26.25	24.55	26.59	27.04	26.93	26.71	26.51	26.67	26.99	<b>27.18</b>	21.60	22.13	23.64	24.19	24.02	24.04	23.59	23.64	24.25	<b>24.33</b>	

Figure 4. Denoising results on image *Cameraman* by different methods (noise level  $\sigma_n = 30$ ).

468 squared difference between the pixels of the original and  
469 denoised images.  
470  
471  
472

Table 1 shows that the proposed AWNNM method almost achieves the highest PSNR, demonstrating a significant average improvement over the other methods. Fig. 4

and Fig. 5 illustrate a visual comparison, where AWNNM effectively recovers subtle features, such as the crack on the dice, which are missed by the other methods.

473  
474  
475

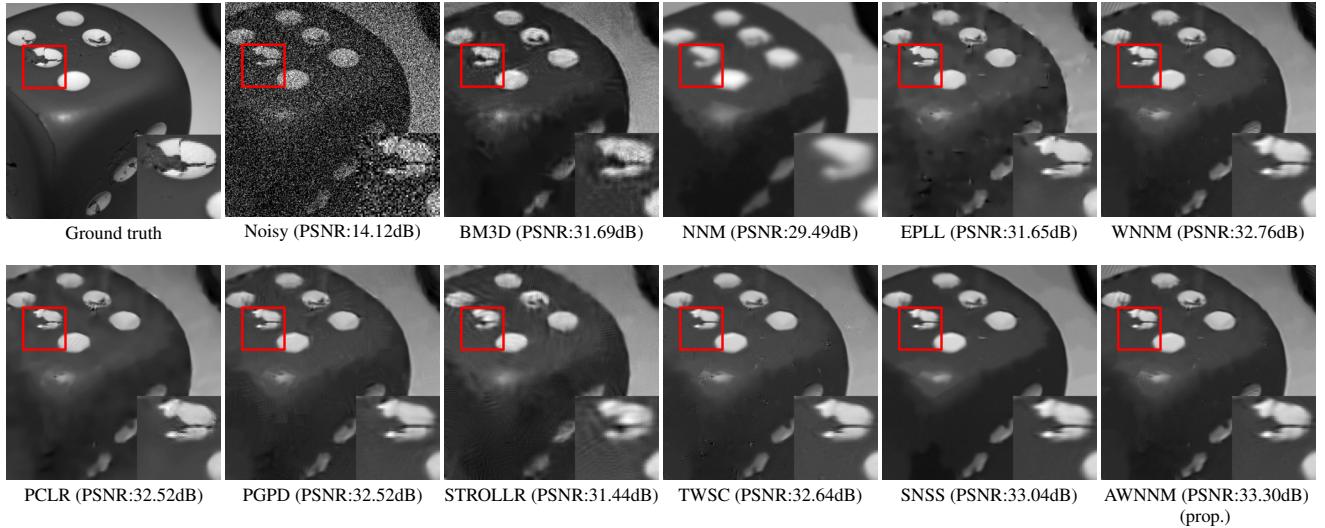
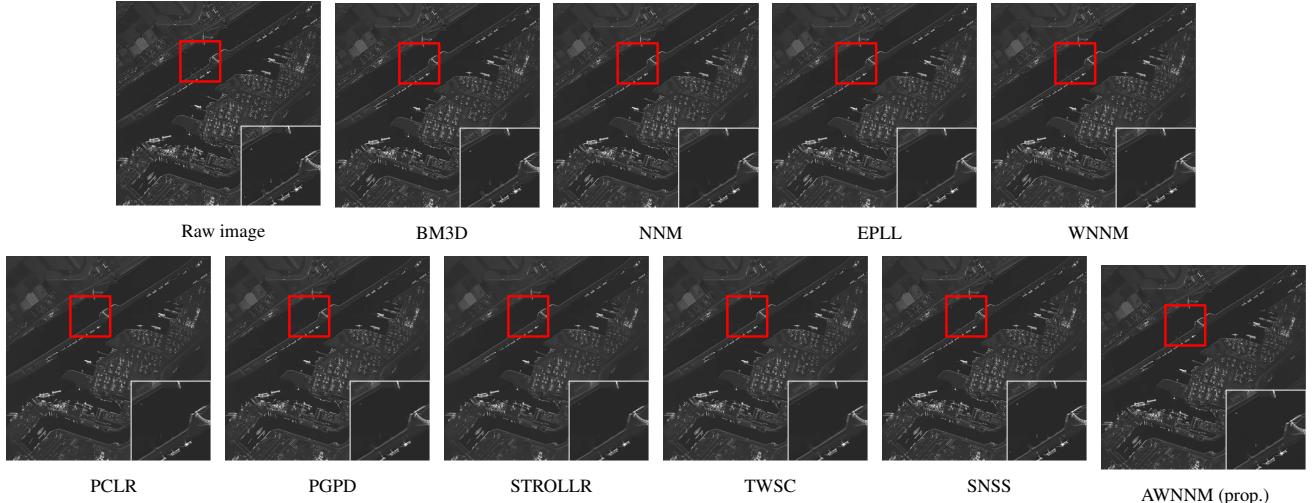
Figure 5. Denoising results on image *Dice* by different methods (noise level  $\sigma_n = 50$ ).

Figure 6. Denoising results on a real SAR image by all competing methods.

476

### 5.2.3. Results on Real-World Images

We also evaluate the denoising performance of the competing methods on real images<sup>1</sup>, including a grayscale SAR image. Fig. 6 shows a comparison of denoising results obtained by the competing methods. It is evident that the proposed AWNNM method preserves local structures with clearer textures in the image and produces the fewest visual artifacts compared to other methods.

## 6. Conclusion

In this paper, we have studied the low-rank matrix recovery problem by introducing a non-convex penalty function on the singular values to approximate the rank function. Our

proposed AWNNM method has been proven to exhibit a superior statistical rate of convergence compared to existing approaches. We then applied the proposed AWNNM method to the image-denoising task. Experimental results on widely used benchmark images demonstrate that AWNNM consistently outperforms state-of-the-art penalty-based algorithms. To the best of our knowledge, this is the first work to establish a statistical guarantee for image denoising based on low-rank matrix denoising.

## References

- [1] Aeron M Buchanan and Andrew W Fitzgibbon. Damped newton algorithms for matrix factorization with missing data. In *IEEE/CVF Computer Vision and Pattern Recognition Conference*, pages 316–322. IEEE, 2005.

<sup>1</sup>The SAT was downloaded at <https://www.iceye.com/resources/datasets>

488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501

- 502 [2] Jian-Feng Cai, Emmanuel J Candès, and Zuowei Shen. A  
503 singular value thresholding algorithm for matrix completion.  
504 *SIAM Journal on optimization*, 20(4):1956–1982, 2010.  
505 [3] Emmanuel J Candès and Terence Tao. The power of convex  
506 relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.  
507 [4] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright.  
508 Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):11, 2011.  
509 [5] Fei Chen, Lei Zhang, and Huimin Yu. External patch prior  
510 guided internal clustering for image denoising. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 603–611, 2015.  
511 [6] Kun Chen, Hongbo Dong, and Kung-Sik Chan. Reduced  
512 rank regression via adaptive nuclear norm penalization.  
513 *Biometrika*, 100(4):901–920, 2013.  
514 [7] Yongyong Chen, Xiaolin Xiao, and Yicong Zhou. Low-rank  
515 quaternion approximation for color image processing. *IEEE Transactions on Image Processing*, 29:1426–1439, 2020.  
516 [8] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and  
517 Karen Egiazarian. Image denoising by sparse 3-d transform-  
518 domain collaborative filtering. *IEEE Transactions on Image Processing*, 16(8):2080–2095, 2007.  
519 [9] David L Donoho, Matan Gavish, and Andrea Montanari.  
520 The phase transition of matrix recovery from gaussian mea-  
521 surements matches the minimax mse of matrix denoising.  
522 *Proceedings of the National Academy of Sciences*, 110(21):  
523 8405–8410, 2013.  
524 [10] Anders Eriksson and Anton van den Hengel. Efficient com-  
525 putation of robust low-rank matrix approximations in the  
526 presence of missing data using the  $L_1$  norm. In *IEEE/CVF Computer Vision and Pattern Recognition Conference*, pages  
527 771–778. IEEE, 2010.  
528 [11] Jianqing Fan and Runze Li. Variable selection via noncon-  
529 cave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360,  
530 2001.  
531 [12] Maryam Fazel, Haitham Hindi, and Stephen P Boyd. A rank  
532 minimization heuristic with application to minimum order  
533 system approximation. In *Proceedings of the 2001 American Control Conference*, pages 4734–4739. IEEE, 2001.  
534 [13] Shuhang Gu, Lei Zhang, Wangmeng Zuo, and Xiangchu  
535 Feng. Weighted nuclear norm minimization with applica-  
536 tion to image denoising. In *IEEE/CVF Computer Vision and Pattern Recognition Conference*, pages 2862–2869, 2014.  
537 [14] Huan Gui, Jiawei Han, and Quanquan Gu. Towards faster  
538 rates and oracle property for low-rank matrix estimation.  
539 In *International Conference on Machine Learning*, pages  
540 2300–2309. PMLR, 2016.  
541 [15] Yao Hu, Debping Zhang, Jieping Ye, Xuelong Li, and Xi-  
542 aofei He. Fast and accurate matrix completion via truncated  
543 nuclear norm regularization. *IEEE transactions on pattern analysis and machine intelligence*, 35(9):2117–2130, 2013.  
544 [16] David R Hunter and Kenneth Lange. A tutorial on mm al-  
545 gorithms. *The American Statistician*, 58(1):30–37, 2004.  
546 [17] Hui Ji, Chaoqiang Liu, Zuowei Shen, and Yuhong Xu. Ro-  
547 bust video denoising using low rank matrix completion. In  
548 *IEEE/CVF Computer Vision and Pattern Recognition Conference*, pages 1791–1798, 2010.  
549 [18] Qifa Ke and Takeo Kanade. Robust  $L_1$  norm factorization in  
550 the presence of outliers and missing data by alternative con-  
551 vex programming. In *IEEE/CVF Computer Vision and Pattern Recognition Conference*, pages 739–746. IEEE, 2005.  
552 [19] N Kishore Kumar and Jan Schneider. Literature survey on  
553 low rank approximation of matrices. *Linear and Multilinear Algebra*, 65(11):2212–2244, 2016.  
554 [20] Joonseok Lee, Seungyeon Kim, Guy Lebanon, Yoram  
555 Singer, and Samy Bengio. Llorma: Local low-rank matrix  
556 approximation. *Journal of Machine Learning Research*, 17  
557 (1):442–465, 2016.  
558 [21] Yixiao Li, Yifan Yu, Qingru Zhang, Chen Liang, Pengcheng  
559 He, Weizhu Chen, and Tuo Zhao. Losparse: Structured com-  
560 pression of large language models based on low-rank and  
561 sparse approximation. In *International Conference on Ma-  
562 chine Learning*, pages 20336–20350. PMLR, 2023.  
563 [22] Zhouchen Lin, Risheng Liu, and Zhixun Su. Linearized al-  
564 ternating direction method with adaptive penalty for low-  
565 rank representation. *Advances in Neural Information Pro-  
566 cessing Systems*, pages 612–620, 2011.  
567 [23] Guangcan Liu, Zhouchen Lin, and Yong Yu. Robust sub-  
568 space segmentation by low-rank representation. In *Proceed-  
569 ings of the 27th International Conference on Machine Learn-  
570 ing*, pages 663–670, 2010.  
571 [24] Risheng Liu, Zhouchen Lin, Fernando De la Torre, and  
572 Zhixun Su. Fixed-rank representation for unsupervised vi-  
573 sual learning. In *IEEE/CVF Computer Vision and Pattern  
574 Recognition Conference*, pages 598–605. IEEE, 2012.  
575 [25] Tianxiang Liu, Zhaosong Lu, Xiaojun Chen, and Yu-  
576 Hong Dai. An exact penalty method for semidefinite-box-  
577 constrained low-rank matrix optimization problems. *IMA Journal of Numerical Analysis*, 40(1):563–586, 2020.  
578 [26] Xin Liu, Guoying Zhao, Jiawen Yao, and Chun Qi. Back-  
579 ground subtraction based on low-rank and structured sparse  
580 decomposition. *IEEE Transactions on Image Processing*, 24  
581 (8):2502–2514, 2015.  
582 [27] Canyi Lu, Jinhui Tang, Shuicheng Yan, and Zhouchen Lin.  
583 Generalized nonconvex nonsmooth low-rank minimization.  
584 In *IEEE/CVF Computer Vision and Pattern Recognition  
585 Conference*, pages 4130–4137, 2014.  
586 [28] Canyi Lu, Jinhui Tang, Shuicheng Yan, and Zhouchen Lin.  
587 Nonconvex nonsmooth low rank minimization via iteratively  
588 reweighted nuclear norm. *IEEE Transactions on Image Pro-  
589 cessing*, 25(2):829–839, 2015.  
590 [29] Ivan Markovsky. Structured low-rank approximation and its  
591 applications. *Automatica*, 44(4):891–909, 2008.  
592 [30] Ivan Markovsky. *Low rank approximation: algorithms, im-  
593 plementation, applications*. Springer, 2011.  
594 [31] Bamdev Mishra, Gilles Meyer, Francis Bach, and Rodolphe  
595 Sepulchre. Low-rank optimization with trace norm penalty.  
596 *SIAM Journal on Optimization*, 23(4):2124–2149, 2013.  
597 [32] Yuji Nakatsukasa. Fast and stable randomized low-rank ma-  
598 trix approximation. *arXiv preprint arXiv:2009.11392*, 2020.  
599 [33] Feiping Nie, Hua Wang, Xiao Cai, Heng Huang, and Chris  
600 Ding. Robust matrix completion via joint schatten p-norm  
601

- 616 and  $l_p$ -norm minimization. In *2012 IEEE 12th International*
- 617 *Conference on Data Mining*, pages 566–574. IEEE, 2012.
- 618 [34] Ankit Parekh and Ivan W Selesnick. Enhanced low-rank ma-
- 619 trix approximation. *IEEE Signal Processing Letters*, 23(4):
- 620 493–497, 2016.
- 621 [35] Hong Peng, Shuyi Hong, Linkai Luo, Qifeng Zhou, and Xi-
- 622 aqin Huang. A new approach of matrix factorization and
- 623 its application in recommender systems. In *2016 15th IEEE*
- 624 *International Conference on Machine Learning and Appli-*
- 625 *cations (ICMLA)*, pages 682–686, 2016.
- 626 [36] Nathan Srebro and Tommi Jaakkola. Weighted low-rank ap-
- 627 proximations. In *Proceedings of the 20th International Con-*
- 628 *ference on Machine Learning*, pages 720–727, 2003.
- 629 [37] Ying Sun, Prabhu Babu, and Daniel P Palomar.
- 630 Majorization-minimization algorithms in signal pro-
- 631 cessing, communications, and machine learning. *IEEE*
- 632 *Transactions on Signal Processing*, 65(3):794–816, 2017.
- 633 [38] Joel A Tropp and Robert J Webber. Randomized algorithms
- 634 for low-rank matrix approximation: Design, analysis, and
- 635 applications. *arXiv preprint arXiv:2306.12418*, 2023.
- 636 [39] Shenlong Wang, Lei Zhang, and Yan Liang. Nonlocal spec-
- 637 tral prior model for low-level vision. In *Computer Vision–*
- 638 *ACCV 2012: 11th Asian Conference on Computer Vision,*
- 639 *Daejeon, Korea, November 5–9, 2012, Revised Selected Pa-*
- 640 *pers, Part III 11*, pages 231–244. Springer, 2012.
- 641 [40] Zhaoran Wang, Han Liu, and Tong Zhang. Optimal compu-
- 642 tational and statistical rates of convergence for sparse non-
- 643 convex learning problems. *Annals of Statistics*, 42(6):2164,
- 644 2014.
- 645 [41] Bihani Wen, Yanjun Li, and Yoram Bresler. When spar-
- 646 sity meets low-rankness: Transform learning with non-local
- 647 low-rank constraint for image restoration. In *Proceedings of*
- 648 *the IEEE International Conference on Acoustics, Speech and*
- 649 *Signal Processing*, pages 2297–2301. IEEE, 2017.
- 650 [42] John Wright, Arvind Ganesh, Shankar Rao, Yigang Peng,
- 651 and Yi Ma. Robust principal component analysis: Exact re-
- 652 covery of corrupted low-rank matrices via convex optimiza-
- 653 tion. *Advances in neural information processing systems*, 22,
- 654 2009.
- 655 [43] Jun Xu, Lei Zhang, Wangmeng Zuo, David Zhang, and Xi-
- 656 angchu Feng. Patch group based nonlocal self-similarity
- 657 prior learning for image denoising. In *Proceedings of the*
- 658 *IEEE International Conference on Computer Vision*, pages
- 659 244–252, 2015.
- 660 [44] Ming Yang, Qilun Luo, Wen Li, and Mingqing Xiao. Non-
- 661 convex 3d array image data recovery and pattern recognition
- 662 under tensor framework. *Pattern recognition*, 122:108311,
- 663 2022.
- 664 [45] Zhiyuan Zha, Xin Yuan, Bihani Wen, Jiantao Zhou, Jiachao
- 665 Zhang, and Ce Zhu. From rank estimation to rank approxi-
- 666 mation: Rank residual constraint for image restoration. *IEEE*
- 667 *Transactions on Image Processing*, 29:3254–3269, 2020.
- 668 [46] Zhiyuan Zha, Xin Yuan, Jiantao Zhou, Ce Zhu, and Bi-
- 669 han Wen. Image restoration via simultaneous nonlocal self-
- 670 similarity priors. *IEEE Transactions on Image Processing*,
- 671 29:8561–8576, 2020.
- 672 [47] Lei Zhang and Wangmeng Zuo. Image restoration: From
- 673 sparse and low-rank priors to deep priors [lecture notes].
- 674 *IEEE Signal Processing Magazine*, 34(5):172–179, 2017.
- 675 [48] Tong Zhang. Analysis of multi-stage convex relaxation for
- 676 sparse regularization. *Journal of Machine Learning Re-*
- 677 *search*, 11(3), 2010.
- 678 [49] Xinfeng Zhang, Weisi Lin, Ruiqin Xiong, Xianming Liu,
- 679 Siwei Ma, and Wen Gao. Low-rank decomposition-based
- 680 restoration of compressed images via adaptive noise estima-
- 681 tion. *IEEE Transactions on Image Processing*, 25(9):4158–
- 682 4171, 2016.
- 683 [50] Xiaowei Zhou, Can Yang, Hongyu Zhao, and Weichuan Yu.
- 684 Low-rank modeling and its applications in image analysis.
- 685 *ACM Computing Surveys*, 47(2):36, 2014.
- 686 [51] Daniel Zoran and Yair Weiss. From learning models of natu-
- 687 ral image patches to whole image restoration. In *Proceedings*
- 688 *of the IEEE International Conference on Computer Vision*,
- 689 pages 479–486. IEEE, 2011.

# Achieving Oracle Statistical Property for Low-Rank Matrix Recovery with Application to Image Denoising

## Supplementary Material

### 690 7. Backgrounds

691 Here, we first introduce essential notations. We denote  
 692 the all-zero matrix and the identity matrix by  $\mathbf{0}$  and  $\mathbf{I}$ , re-  
 693 spectively, with dimensions inferred from context. Con-  
 694 sider a matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ . The smallest of  $\mathbf{X}$  are rep-  
 695 resented as  $\lambda_{\min}(\mathbf{X})$ . The  $p$ -norm of  $\mathbf{X}$  is defined as  
 696  $\|\mathbf{X}\|_p = \left( \sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^p \right)^{\frac{1}{p}}$  for a real  $p > 0$ . Notably,  
 697 the Frobenius norm and spectral norm are denoted by  $\|\mathbf{X}\|_F$   
 698 and  $\|\mathbf{X}\|_2$ , respectively, while  $\|\mathbf{X}\|_1$  represents the sum of  
 699 the absolute values of all entries of  $\mathbf{X}$ . The maximum-  
 700 absolute-value norm of  $\mathbf{X}$  is expressed as  $\|\mathbf{X}\|_{\max}$ , and  
 701 the minimum-absolute-value norm as  $\|\mathbf{X}\|_{\min}$ . The nu-  
 702 clear norm of  $\mathbf{X}$ , defined as the sum of the singular val-  
 703 ues of  $\mathbf{X}$ , is expressed as  $\|\mathbf{X}\|_* = \sum_{i=1}^d \sigma_i(\mathbf{X})$ . The  
 704 vectorization of  $\mathbf{X}$ , achieved by stacking its columns, is  
 705 denoted as  $\text{vec}(\mathbf{X})$ . Furthermore, the Kronecker product  
 706 and Hadamard product (also referred to as the entry-wise  
 707 product) of matrices  $\mathbf{X}$  and  $\mathbf{Y}$  are denoted by  $\mathbf{X} \otimes \mathbf{Y}$  and  
 708  $\mathbf{X} \odot \mathbf{Y}$ , respectively. The Euclidean inner product is de-  
 709 fined as  $\langle \mathbf{X}, \mathbf{Y} \rangle := \text{tr}(\mathbf{XY}^\top)$ , where  $\text{tr}(\cdot)$  denotes the  
 710 trace of a matrix. For a real-valued function  $f(\mathbf{X})$ , the gra-  
 711 dient  $\nabla f(\mathbf{X})$  is a  $d \times d$  matrix with the  $(i, j)$ -th element  
 712 given by  $\frac{\partial}{\partial X_{ij}} f(\mathbf{X})$ , denoted by  $\nabla_{ij} f(\mathbf{X})$ , while  $\nabla^2 f(\mathbf{X})$   
 713 represents the  $d^2 \times d^2$  Hessian matrix. For an index set  $\mathcal{E}$ ,  
 714 its cardinality is represented as  $|\mathcal{E}|$ , while its complement is  
 715 denoted by  $\bar{\mathcal{E}}$ . Given a vector  $\mathbf{w}$ ,  $\mathbf{w}_{\mathcal{E}}$  is defined such that  
 716 its  $i$ -th entry equals to  $w_i$  if  $i \in \mathcal{E}$  and is 0 otherwise. For  
 717 the functionals  $f(n)$  and  $g(n)$ , we express  $f(n) \gtrsim g(n)$   
 718 if  $f(n) \geq cg(n)$ ,  $f(n) \lesssim g(n)$  if  $f(n) \leq Cg(n)$ , and  
 719  $f(n) \asymp g(n)$  if  $cg(n) \leq f(n) \leq Cg(n)$  for some pos-  
 720 itive constants  $c$  and  $C$ . Additionally,  $\mathcal{O}_p(\cdot)$  is utilized to  
 721 indicate boundedness in probability.

722 We recall the problem presented in (3):

$$723 \hat{\mathbf{X}}^{(k)} \in \arg \min_{\mathbf{X}} \left\{ f(\mathbf{X}) + \sum_{i=1}^d w_i^{(k-1)} \sigma_i(\mathbf{X}) \right\}.$$

724 This weighted nuclear norm penalized problem can be gen-  
 725 erally expressed as:

$$726 \hat{\mathbf{X}}^{(k)} \in \arg \min_{\mathbf{X}} \{f(\mathbf{X}) + \|\mathbf{w} \odot \boldsymbol{\sigma}(\mathbf{X})\|_1\}, \quad (6)$$

727 where  $\mathbf{w} \in \mathbb{R}^d$  is a vector of regularization parameters with  
 728  $w_i \in [0, \lambda]$  for  $i = 1, \dots, d$ .

### 729 8. Proofs of Statistical Theory

#### 730 8.1. Technical Lemmata

731 **Lemma 8.** Consider the general problem in (6). Assume  
 732 that there exists a set  $\mathcal{E}$  such that

$$733 \mathcal{S}^* \subseteq \mathcal{E}, |\mathcal{E}| \leq 2r^*, \text{ and } \|\mathbf{w}_{\bar{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2}. \quad 733$$

734 If  $\lambda \geq 2 \|\nabla f(\mathbf{X}^*)\|_2$ , then the optimal solution  $\hat{\mathbf{X}}$  satisfies

$$735 \|\hat{\mathbf{X}} - \mathbf{X}^*\|_F \leq \|\mathbf{w}_{\mathcal{S}^*}\|_2 + \|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F \leq \frac{3}{2}\lambda\sqrt{r^*}. \quad 735$$

736 *Proof.* Given that  $\nabla^2 f(\mathbf{X}) = \mathbf{I} \otimes \mathbf{I}$ , the mean value theo-  
 737 rem guarantees the existence of  $\rho \in [0, 1]$  such that

$$738 \langle \nabla f(\mathbf{X}) - \nabla f(\mathbf{X}^*), \Delta \rangle \\ 739 = \text{vec}^\top(\Delta) \nabla^2 f(\mathbf{X}^* + \rho\Delta) \text{vec}(\Delta) = \|\Delta\|_F^2,$$

740 where  $\Delta = \mathbf{X} - \mathbf{X}^*$ .

741 Hence, we obtain

$$742 \|\Delta\|_F^2 \leq \langle \nabla f(\mathbf{X}) - \nabla f(\mathbf{X}^*), \Delta \rangle. \quad 742$$

743 Let  $\hat{\Delta} = \hat{\mathbf{X}} - \mathbf{X}^*$ . Applying the inequality (7) with  
 744  $\mathbf{X} = \hat{\mathbf{X}}$  and  $\Delta = \hat{\Delta}$  yields

$$745 \|\hat{\Delta}\|_F^2 \leq \langle \nabla f(\hat{\mathbf{X}}) - \nabla f(\mathbf{X}^*), \hat{\Delta} \rangle. \quad 745$$

746 Since  $\hat{\mathbf{X}}$  is the minimizer of the loss  $f(\mathbf{X})$  and the  
 747 penalty  $\sum_{i=1}^m w_i^{(k-1)} \sigma_i(\mathbf{X}) = \langle \text{Diag}(\mathbf{w}), \text{Diag}(\boldsymbol{\sigma}(\mathbf{X})) \rangle$ ,  
 748 we have the optimality condition  $\nabla f(\hat{\mathbf{X}}) + \text{Diag}(\mathbf{w}) \cdot \boldsymbol{\Xi} = \mathbf{0}$  for the subgradient  $\boldsymbol{\Xi} \in \partial \|\text{Diag}(\boldsymbol{\sigma}(\mathbf{X}))\|_1$  at  $\mathbf{X} = \hat{\mathbf{X}}$ .  
 749 Therefore, we have

$$750 \|\hat{\Delta}\|_F^2 \leq \langle -\nabla f(\mathbf{X}^*) - \text{Diag}(\mathbf{w}) \cdot \boldsymbol{\Xi}, \hat{\Delta} \rangle \quad 751$$

$$752 = -\underbrace{\langle \nabla f(\mathbf{X}^*), \hat{\Delta} \rangle}_{\text{I}} - \underbrace{\langle \text{Diag}(\mathbf{w}) \cdot \boldsymbol{\Xi}, \hat{\Delta} \rangle}_{\text{II}}. \quad 752$$

753 It remains to bound terms I and II, respectively. For term  
 754 I, projecting the support of  $\nabla f(\mathbf{X}^*)$  and  $\hat{\Delta}$  to  $\mathcal{E}$  and  $\bar{\mathcal{E}}$ , and

755 then using the matrix Hölder inequality, we obtain

$$\begin{aligned}
 756 & \left\langle \nabla f(\mathbf{X}^*), \widehat{\Delta} \right\rangle \\
 757 &= \left\langle \Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*)), \Pi_{\mathcal{F}_{\mathcal{E}}}(\widehat{\Delta}) \right\rangle \\
 758 &\quad + \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\nabla f(\mathbf{X}^*)), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 759 &\geq - \|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F \left\| \Pi_{\mathcal{F}_{\mathcal{E}}}(\widehat{\Delta}) \right\|_F \\
 760 &\quad - \|\Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\nabla f(\mathbf{X}^*))\|_2 \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_*.
 \end{aligned}$$

761 For term II, projecting the support of  $\text{Diag}(\mathbf{w}) \cdot \Xi$  and  
762  $\widehat{\Delta}$  to  $\mathcal{S}^*$ ,  $\mathcal{E} \setminus \mathcal{S}^*$ , and  $\bar{\mathcal{E}}$ , we obtain

$$\begin{aligned}
 763 & \left\langle \text{Diag}(\mathbf{w}) \cdot \Xi, \widehat{\Delta} \right\rangle \\
 764 &= \left\langle \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\Delta}) \right\rangle \\
 765 &\quad + \left\langle \Pi_{\mathcal{F}_{\mathcal{E} \setminus \mathcal{S}^*}}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\mathcal{E} \setminus \mathcal{S}^*}}(\widehat{\Delta}) \right\rangle \\
 766 &\quad + \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 767 &= \left\langle \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\Delta}) \right\rangle \\
 768 &\quad + \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 769 &\geq - \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\text{Diag}(\mathbf{w}) \cdot \Xi)\|_F \left\| \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\Delta}) \right\|_F \\
 770 &\quad + \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 771 &\geq - \|\mathbf{w}_{\mathcal{S}^*}\|_2 \left\| \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\Delta}) \right\|_F + \|\mathbf{w}_{\bar{\mathcal{E}}}^c\|_{\min} \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_*, 
 \end{aligned}$$

772 where the second equality is due to the fact that

$$773 \left\langle \Pi_{\mathcal{F}_{\mathcal{E} \setminus \mathcal{S}^*}}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\mathcal{E} \setminus \mathcal{S}^*}}(\widehat{\Delta}) \right\rangle = 0,$$

774 and the last inequality is due to

$$\begin{aligned}
 775 & \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 776 &= \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Sigma}) \right\rangle \\
 777 &= \left\langle \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\text{Diag}(\mathbf{w}) \cdot \Xi), \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\rangle \\
 778 &\geq \|\mathbf{w}_{\bar{\mathcal{E}}}^c\|_{\min} \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_*. 
 \end{aligned}$$

Substituting the above results into (9) yields

$$\begin{aligned}
 780 & \left\| \widehat{\Delta} \right\|_F^2 \leq \left\| \Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*)) \right\|_F \left\| \Pi_{\mathcal{F}_{\mathcal{E}}}(\widehat{\Delta}) \right\|_F \\
 781 &\quad + \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\nabla f(\mathbf{X}^*)) \right\|_2 \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_* \\
 782 &\quad - \|\mathbf{w}_{\bar{\mathcal{E}}}^c\|_{\min} \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_* \\
 783 &\quad + \|\mathbf{w}_{\mathcal{S}^*}\|_2 \left\| \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\Delta}) \right\|_F \\
 784 &\leq (\|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F + \|\mathbf{w}_{\mathcal{S}^*}\|_2) \left\| \Pi_{\mathcal{F}_{\mathcal{E}}}(\widehat{\Delta}) \right\|_F \\
 785 &\quad + \left\| \Pi_{\mathcal{F}_{\bar{\mathcal{E}}}^c}(\widehat{\Delta}) \right\|_* (\|\nabla f(\mathbf{X}^*)\|_2 - \|\mathbf{w}_{\bar{\mathcal{E}}}^c\|_{\min}) \\
 786 &\leq (\|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F + \|\mathbf{w}_{\mathcal{S}^*}\|_2) \left\| \Pi_{\mathcal{F}_{\mathcal{E}}}(\widehat{\Delta}) \right\|_F \\
 787 &\leq (\|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F + \|\mathbf{w}_{\mathcal{S}^*}\|_2) \left\| \widehat{\Delta} \right\|_F, \quad (10)
 \end{aligned}$$

where the third inequality is due to  $\|\mathbf{w}_{\bar{\mathcal{E}}}^c\|_{\min} \geq \frac{\lambda}{2} \geq \|\nabla f(\mathbf{X}^*)\|_2$ . Dividing by  $\left\| \widehat{\Delta} \right\|_F$  on both sides of the inequality (10), we have

$$\begin{aligned}
 791 & \left\| \widehat{\Delta} \right\|_F \leq \|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_F + \|\mathbf{w}_{\mathcal{S}^*}\|_2 \\
 792 &\leq \sqrt{r^*} \|\Pi_{\mathcal{F}_{\mathcal{E}}}(\nabla f(\mathbf{X}^*))\|_2 + \sqrt{r^*} \|\mathbf{w}_{\mathcal{S}^*}\|_{\max} \\
 793 &\leq \frac{3}{2} \lambda \sqrt{r^*}.
 \end{aligned}$$

□ 794

**Lemma 9.** Suppose that Assumption 1 holds. Consider the problem (3). Define the set  $\mathcal{E}^{(k)}$  by

$$795 \mathcal{E}^{(k)} = \mathcal{S}^* \cup \mathcal{S}^{(k)}, \text{ with } \mathcal{S}^{(k)} = \left\{ i \mid w_i^{(k-1)} \leq p_{\lambda}'(u) \right\}, \quad 797$$

where  $u = c\lambda$  and  $c = \frac{3}{2}$  is the same to that given in Assumption 1. If  $\lambda \geq 4 \|\nabla f(\mathbf{X}^*)\|_2$ , then for  $k \geq 1$ , we have  $|\mathcal{E}^{(k)}| \leq 2r^*$  and the optimal solution  $\widehat{\Sigma}^{(k)}$  satisfies

$$\begin{aligned}
 801 & \left\| \widehat{\mathbf{X}}^{(k)} - \mathbf{X}^* \right\|_F \leq \left\| \mathbf{w}_{\mathcal{S}^*}^{(k-1)} \right\|_2 + \left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}} \nabla f(\mathbf{X}^*) \right\|_F \\
 802 &\leq \frac{3}{2} \lambda \sqrt{r^*}.
 \end{aligned}$$

*Proof.* We first prove  $|\mathcal{E}^{(k)}| \leq 2r^*$  holds by induction. For  $k = 1$ , we have  $w_i^{(0)} = \lambda \geq p_{\lambda}'(u)$  and thus  $\mathcal{S}^{(0)} = \emptyset$  and  $\mathcal{E}^{(1)} = \mathcal{S}^*$ , which implies  $|\mathcal{E}^{(1)}| \leq 2r^*$  holds. Assume  $|\mathcal{E}^{(k)}| \leq 2r^*$  holds at  $k - 1$ , i.e.,  $|\mathcal{E}^{(k-1)}| \leq 2r^*$  holds for some  $k \geq 2$ . Next, we will prove  $|\mathcal{E}^{(k)}| \leq 2r^*$  holds at  $k$ . For any  $i \in \mathcal{S}^{(k)}$ , we obtain  $\sigma_i(\widehat{\mathbf{X}}^{(k-1)}) \geq u$  and further

809 have

$$\begin{aligned} \sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} &\leq \frac{\left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)} \setminus \mathcal{S}}}(\widehat{\mathbf{X}}^{(k-1)}) \right\|_F}{u} \\ &\leq \frac{\left\| \widehat{\mathbf{X}}^{(k-1)} - \mathbf{X}^* \right\|_F}{u}. \end{aligned} \quad (11)$$

812 For any  $i \in \overline{\mathcal{S}^{(k-1)}}$ , we have

813  $\mathbf{w}_i^{(k-2)} = p'_\lambda \left( \sigma_i(\widehat{\mathbf{X}}^{(k-2)}) \right) \geq p'_\lambda(u) \geq \frac{\lambda}{2},$

814 which implies

815  $\left\| \mathbf{w}_{\overline{\mathcal{E}^{(k-1)}}}^{(k-2)} \right\|_{\min} \geq \left\| \mathbf{w}_{\overline{\mathcal{S}^{(k-1)}}}^{(k-2)} \right\|_{\min} \geq p'_\lambda(u) \geq \lambda/2.$

816 One also has  $|\mathcal{E}^{(k-1)}| \leq 2r^*$  and  $\mathcal{S}^* \subseteq \mathcal{E}^{k-1}$ . Applying 817 Lemma 8 with  $\widehat{\mathbf{X}} = \widehat{\mathbf{X}}^{(k-1)}$ ,  $\mathcal{E} = \mathcal{E}^{(k-1)}$ , and  $\mathbf{w}_{\mathcal{S}^*} =$  818  $\mathbf{w}_{\mathcal{S}^*}^{(k-2)}$  yields

819  $\left\| \widehat{\mathbf{X}}^{(k-1)} - \mathbf{X}^* \right\|_F \leq \frac{3}{2} \lambda \sqrt{r^*}.$

820 Substituting the above result into the inequality (11) yields

821  $\sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} \leq \frac{3}{2u} \lambda \sqrt{r^*} = \sqrt{r^*}.$

822 Thus, we have

823  $|\mathcal{E}^{(k)}| = |\mathcal{S}^* \cup (\mathcal{E}^{(k)} \setminus \mathcal{S}^*)| = |\mathcal{S}^*| + |\mathcal{E}^{(k)} \setminus \mathcal{S}^*| \leq 2r^*,$

824 completing the induction. Then, by the definition of  $\mathcal{E}^{(k)}$ , 825 we have

826  $\left\| \mathbf{w}_{\overline{\mathcal{E}^{(k)}}} \right\|_{\min} \geq p'_\lambda(u) \geq \lambda/2.$

827 Applying Lemma 8 with  $\widehat{\mathbf{X}} = \widehat{\mathbf{X}}^{(k)}$ ,  $\mathcal{E} = \mathcal{E}^{(k)}$ , and  $\mathbf{w}_{\mathcal{S}^*} =$  828  $\mathbf{w}_{\mathcal{S}^*}^{(k-1)}$ , the optimal solution  $\widehat{\mathbf{X}}^{(k)}$  to (3) satisfies

$$\begin{aligned} \left\| \widehat{\mathbf{X}}^{(k)} - \mathbf{X}^* \right\|_F &\leq \left\| \mathbf{w}_{\mathcal{S}^*}^{(k-1)} \right\|_2 + \left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*)) \right\|_F \\ &\leq \frac{3}{2} \lambda \sqrt{r^*}. \end{aligned}$$

831  $\square$ 832 **Lemma 10.** For any norm  $\|\cdot\|$ , the following inequality 833 holds:

834  $\|\mathbf{w}_{\mathcal{S}^*}\| \leq \left\| \mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \right\| + \lambda u^{-1} \left\| \mathbf{X}_{\mathcal{S}^*} - \mathbf{X}_{\mathcal{S}}^* \right\|,$

835 where  $\mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) =$   
836  $[p'_\lambda(\sigma_1(\mathbf{X}^*) - u), \dots, p'_\lambda(\sigma_d(\mathbf{X}^*) - u)]$

Proof. We consider two cases for each index  $i \in \mathcal{S}^*$ .**Case I:** If  $|\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{X}^*)| \geq u$ , under Assumption 1, 838 we have 839

840  $p'_\lambda(\sigma_i(\mathbf{X})) \leq \lambda \leq \lambda u^{-1} |\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{X}^*)|.$

**Case II:** If  $|\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{X}^*)| < u$ , then due to the non-increasing nature of  $p'_\lambda(\cdot)$ , it follows that 841 842

843  $p'_\lambda(\sigma_i(\mathbf{X})) \leq p'_\lambda(\sigma_i(\mathbf{X}^*) - u).$

By combining both cases, we establish that for all  $i \in \mathcal{S}^*$ , the following inequality holds: 844 845

846  $p'_\lambda(\sigma_i(\mathbf{X})) \leq p'_\lambda(\sigma_i(\mathbf{X}^*) - u) + \lambda u^{-1} |\sigma_i(\mathbf{X}) - \sigma_i(\mathbf{X}^*)|.$

Expressing this inequality in vector form for all indices in  $\mathcal{S}^*$ , applying the any norm  $\|\cdot\|$  and utilizing the triangle inequality, we obtain 847 848 849

$$\begin{aligned} &\left\| \mathbf{w}_{\mathcal{S}^*} \right\| \\ &\leq \left\| \mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \right\| \\ &\quad + \lambda u^{-1} \left\| \sigma_{\mathcal{S}^*}(\mathbf{X}) - \sigma_{\mathcal{S}^*}(\mathbf{X}^*) \right\|. \end{aligned} \quad (12) \quad 850 \quad 851 \quad 852$$

Recognizing that the singular values are Lipschitz continuous functions of the matrix entries, we have 853 854

855  $\left\| \sigma_{\mathcal{S}^*}(\mathbf{X}) - \sigma_{\mathcal{S}^*}(\mathbf{X}^*) \right\| \leq \left\| \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}) - \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^*) \right\|.$

Substituting this back into the inequality (12), we obtain 856

$$\begin{aligned} \left\| \mathbf{w}_{\mathcal{S}^*} \right\| &\leq \left\| \mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \right\| \\ &\quad + \lambda u^{-1} \left\| \mathbf{X}_{\mathcal{S}^*} - \mathbf{X}_{\mathcal{S}}^* \right\|. \end{aligned} \quad 857 \quad 858$$

 $\square$  859

## 8.2. Proof of Theorem 4

Proof. By Lemma 9, we have

860  $\left\| \widehat{\mathbf{X}}^{(k)} - \mathbf{X}^* \right\|_F \leq \left\| \mathbf{w}_{\mathcal{S}^*}^{(k-1)} \right\|_2 + \left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*)) \right\|_F. \quad 861 \quad 862$

Using Lemma 10 yields that

863 
$$\begin{aligned} \left\| \mathbf{w}_{\mathcal{S}^*} \right\| &\leq \left\| \mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \right\| \\ &\quad + \lambda u^{-1} \left\| \mathbf{X}_{\mathcal{S}^*} - \mathbf{X}_{\mathcal{S}}^* \right\|. \end{aligned} \quad 864 \quad 865$$

Plugging two inequalities above obtains that 866

$$\begin{aligned} \left\| \widehat{\mathbf{X}}^{(k)} - \mathbf{X}^* \right\|_F &\leq \left\| \mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \right\|_2 \\ &\quad + \lambda u^{-1} \left\| \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\widehat{\mathbf{X}}^{(k-1)}) - \Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^*) \right\|_F \\ &\quad + \left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*)) \right\|_F \\ &\leq \lambda u^{-1} \left\| \widehat{\mathbf{X}}^{(k-1)} - \mathbf{X}^* \right\|_F \\ &\quad + \left\| \Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*)) \right\|_F, \end{aligned} \quad 867 \quad 868 \quad 869 \quad 870 \quad 871$$

872 where  $\mathbf{w}(\sigma_{\mathcal{S}^*}(\mathbf{X}^*) - u\mathbf{1}_{\mathcal{S}^*}) \leq \mathbf{w}(\gamma\lambda \cdot \mathbf{1}_{\mathcal{S}}) = \mathbf{0}_{\mathcal{S}^*}$  due to  
 873  $\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^*)\|_{\min} \geq (\gamma + c)\lambda \gtrsim \lambda$  in Assumption 3. In the  
 874 next, we bound term  $\|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*))\|_F$ . Separating the  
 875 support of  $\Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*))$  to  $\mathcal{S}^*$  and  $\mathcal{E}^{(k)} \setminus \mathcal{S}^*$  and then  
 876 using triangle inequality, we obtain

$$\begin{aligned} 877 & \|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*))\|_F \\ 878 & \leq \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F + \|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)} \setminus \mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F \\ 879 & \leq \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F + \sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} \|\nabla f(\mathbf{X}^*)\|_2. \end{aligned}$$

880 From the assumption, we know  $\lambda \geq 4 \|\nabla f(\mathbf{X}^*)\|_2$ . Following the proof of Lemma 9,  
 881  $\sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|}$  can be bounded by

$$\begin{aligned} 883 & \sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} \leq \frac{\|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)} \setminus \mathcal{S}^*}}(\hat{\mathbf{X}}^{(k-1)})\|_F}{u} \\ 884 & \leq \frac{\|\hat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F}{u}. \end{aligned}$$

885 Therefore,  $\|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*))\|_F$  can be simplified to

$$\begin{aligned} 886 & \|\Pi_{\mathcal{F}_{\mathcal{E}^{(k)}}}(\nabla f(\mathbf{X}^*))\|_F \leq \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F \\ 887 & \quad + \frac{\lambda}{4u} \|\hat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F, \end{aligned}$$

888 which yields the contraction property

$$\begin{aligned} 889 & \|\hat{\mathbf{X}}^{(k)} - \mathbf{X}^*\|_F \\ 890 & \leq \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F + 1.25\lambda u^{-1} \|\hat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F \\ 891 & = \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F + \delta \|\hat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F, \end{aligned}$$

892 where  $\delta = 1.25\lambda u^{-1}$ . Consequently, we obtain

$$\begin{aligned} 893 & \|\hat{\mathbf{X}}^{(k)} - \mathbf{X}^*\|_F \leq \frac{1}{1-\delta} \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F \\ 894 & \quad + \delta^{k-1} \|\hat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F. \end{aligned}$$

895  $\square$

### 896 8.3. Concentration Inequalities

897 **Lemma 11.** Let  $\mathbf{E}$  be a sub-Gaussian random matrix with  
 898  $E_{ij}$  are sub-Gaussian random variables with zero mean and  
 899 covariance  $\varepsilon$ . Consider that  $\mathbf{Y} = \mathbf{X}^* + \mathbf{E}$ . There exists  
 900 some constant  $c_1, c_2$ , the following tail bound holds

$$901 \mathbb{P}(|X_{ij}^* - Y_{ij}| > t) \leq c_1 \exp(-c_2 t^2).$$

**Lemma 12.** Under the same conditions in Lemma 11, if taking  $\lambda = \sqrt{\frac{3 \log mn}{c_2}} \asymp \sqrt{\log mn}$ , then the following result holds

$$\mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_F \leq \lambda) \geq 1 - \frac{c_1}{mn}.$$

*Proof.* Applying Lemma 11 and union bound, for any  $\lambda$  such that  $0 < \lambda < t_0$ , we obtain

$$\begin{aligned} 908 & \mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_F > \lambda) \\ 909 & \leq c_1 mn \exp(-c_2 \lambda^2) \\ 910 & = c_1 \exp(-c_2 \lambda^2 + \log mn). \end{aligned}$$

By taking  $\lambda = \sqrt{\frac{2 \log mn}{c_2}} \asymp \sqrt{\log mn}$ , we obtain

$$\begin{aligned} 912 & \mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_F \leq \lambda) \\ 913 & \geq 1 - c_1 \exp(-c_2 \lambda^2 + \log mn) \\ 914 & = 1 - \frac{c_1}{mn}. \end{aligned}$$

$\square$

**Lemma 13.** Under the same conditions in Lemma 11, the following result holds

$$\mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_2 \leq \lambda) \geq 1 - \frac{c_1}{mn}.$$

*Proof.* Recall that the spectral norm of a matrix is bounded above by its Frobenius norm

$$\|\mathbf{X}^* - \mathbf{Y}\|_2 \leq \|\mathbf{X}^* - \mathbf{Y}\|_F.$$

Therefore, we have

$$\mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_2 > \lambda) \leq \mathbb{P}(\|\mathbf{X}^* - \mathbf{Y}\|_F > \lambda).$$

Based on Lemma 12, we complete the proof.

$\square$

**Lemma 14.** Under the same conditions in Lemma 11, the following result holds

$$\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_F = \mathcal{O}_P(\sqrt{r^*}).$$

*Proof.* Applying Lemma 11 and union bound, for any  $M$  such that  $0 < M \sqrt{\frac{1}{n}} < t_0$ , we obtain

$$\begin{aligned} 931 & \mathbb{P}(\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_2 > M) \\ 932 & \leq \mathbb{P}(\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_F > M) \\ 933 & \leq c_1 mn \exp(-c_2 M^2) \\ 934 & = c_1 \exp(-c_2 M^2 + \log mn). \end{aligned}$$

By taking  $M$  such that  $\sqrt{\frac{\log mn}{c_2}} < M$  and  $M \rightarrow \infty$  in above inequality obtains

$$\limsup_{M \rightarrow \infty} \frac{1}{n} \mathbb{P}(\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_2 > M) = 0. \quad (13)$$

This proof is completed by applying  $\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_F \leq \sqrt{r^*} \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_2$ .

$\square$

940 **8.4. Proof of Corollary 6**

941 *Proof.* If  $\lambda$  satisfies  $\lambda \asymp \sqrt{\log nm}$ , then by Lemma 13,  
 942  $\lambda \geq 4 \|\nabla f(\mathbf{X}^*)\|_2$  holds with high probability. Applying  
 943 Lemma 9 with  $k = 1$ , we obtain

$$944 \quad \|\hat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F \leq \frac{3}{2} \lambda \sqrt{r^*}.$$

945 If  $\lambda \asymp \sqrt{\log mn}$ , then  $\|\hat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F \lesssim \sqrt{r^* \log mn}$ .  $\square$

946 **8.5. Proof of Corollary 7**

947 *Proof.* If  $\lambda$  satisfies  $\lambda \asymp \sqrt{\log mn}$ , then by Lemma 13,  
 948  $\lambda \geq 4 \|\nabla f(\Sigma^*)\|_2$  holds with high probability. Applying  
 949 Theorem 4, we obtain

$$\begin{aligned} 950 \quad & \|\hat{\mathbf{X}}^{(k)} - \mathbf{X}^*\|_F \\ 951 \quad & \leq \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\Sigma^*))\|_F + \delta \|\hat{\mathbf{X}}^{(k-1)} - \mathbf{X}^*\|_F \\ 952 \quad & \leq \frac{1}{1-\delta} \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\Sigma^*))\|_F + \delta^{k-1} \|\hat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F \\ 953 \quad & \leq \frac{1}{1-\delta} \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\Sigma^*))\|_F + \delta^{k-1} \frac{3}{2} \lambda \sqrt{r^*}, \end{aligned}$$

954 where the last inequality is due to  $\|\hat{\mathbf{X}}^{(1)} - \mathbf{X}^*\|_F \leq \frac{3}{2} \lambda \sqrt{r^*}$ ,  
 955 which follow from Lemma 9 with  $k = 1$ . One has

$$956 \quad \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\nabla f(\mathbf{X}^*))\|_F = \|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_F.$$

957 By Lemma 14,  $\|\Pi_{\mathcal{F}_{\mathcal{S}^*}}(\mathbf{X}^* - \mathbf{Y})\|_F = \mathcal{O}_P(\sqrt{r^*})$ . If  $K \geq$   
 958  $1 + \frac{\log \lambda}{\log \delta^{-1}} \gtrsim \log \lambda \gtrsim \log \log mn$ , then we have

$$959 \quad \delta^{k-1} \lambda \sqrt{r^*} \leq \frac{1}{\lambda} \lambda \sqrt{r^*} \leq \sqrt{r^*}.$$

960 Combining the above results yields that  $\|\hat{\mathbf{X}}^{(K)} - \mathbf{X}^*\|_F =$   
 961  $\mathcal{O}_P(\sqrt{r^*})$ .  $\square$

963 **9. More Experimental Setup**

964 We also check the performance of our proposed estimator  
 965 on the BSD100 datasets. The BSD100 dataset contains a di-  
 966 verse set of natural images with complex textures and real-  
 967 world scenes, posing significant challenges for denoising  
 968 algorithms. The following experiments demonstrate the su-  
 969 perior performance of our algorithm in handling real-world  
 970 image noise and enhancing the visual quality of denoised  
 971 images.

972 • We primarily evaluate the denoising performance by com-  
 973 paring the visual quality of the output images. As shown  
 974 in Fig 7 - Fig 10, the AWNNM algorithm outperforms  
 975 competing methods in preserving fine details and effec-  
 976 tively suppressing noise. This is particularly evident in

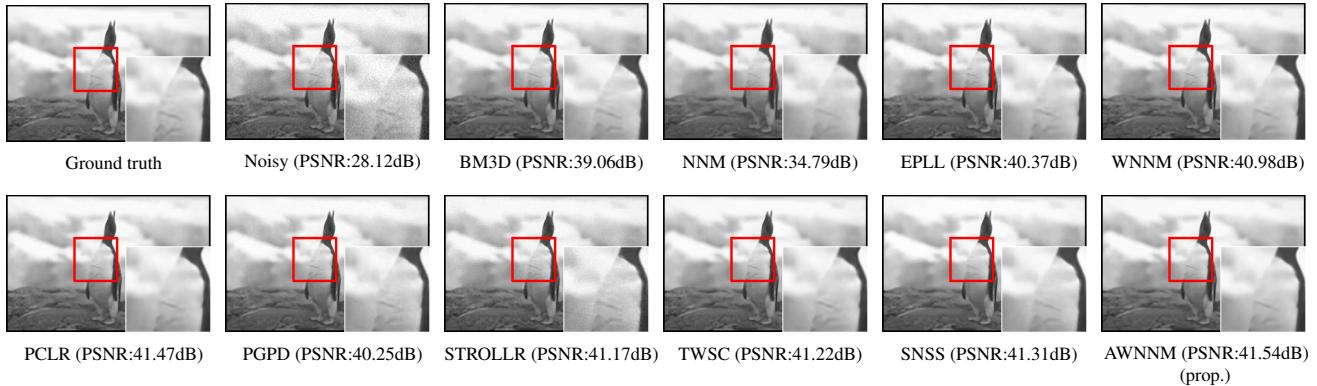
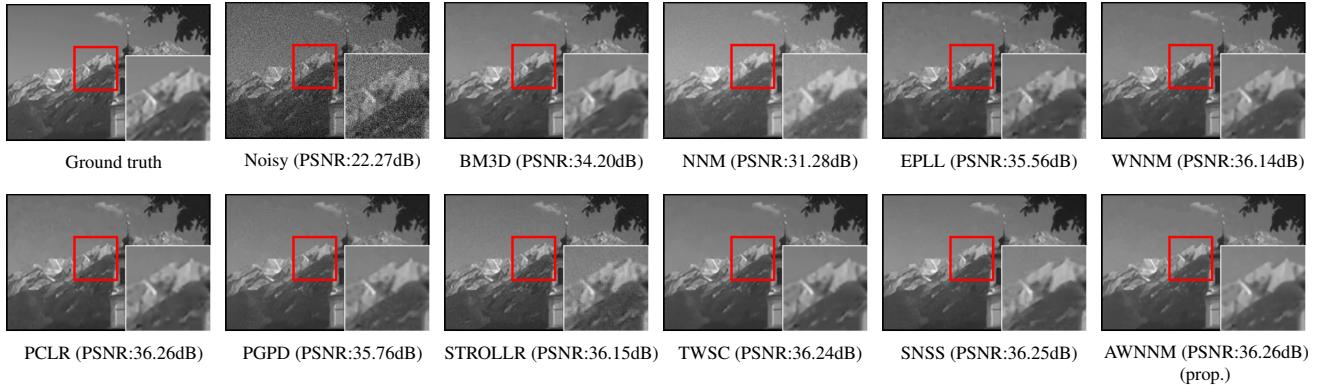
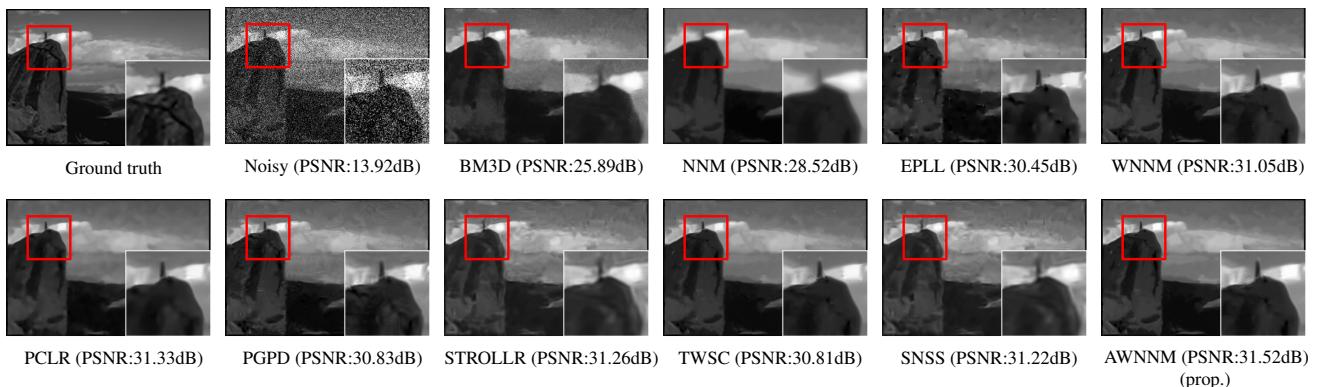
977 high-texture regions, where AWNNM demonstrates ex-  
 978 ceptional robustness and fidelity. These results highlight  
 979 the effectiveness of AWNNM for real-world denoising  
 980 tasks.

- 981 In addition to visual comparisons, we also compute com-  
 982 monly used metric PSNR to further validate the perfor-  
 983 mance of the AWNNM algorithm. As evidenced by the  
 984 results presented in the figures and tables, our algorithm  
 985 achieves higher scores on PNSR compared to other meth-  
 986 ods. This indicates that AWNNM not only provides high-  
 987 quality denoising results but also preserves structural in-  
 988 formation and texture details effectively, even under chal-  
 989 lenging noise conditions.

990 The experimental results on the BSD100 dataset demon-  
 991 strate the outstanding performance of the proposed  
 992 AWNNM algorithm for real-world image denoising. Its  
 993 ability to balance texture preservation, noise suppression,  
 994 and detail restoration surpasses existing methods, under-  
 995 scoring its potential for widespread applications in image  
 996 processing.

Table 2. Denoising results (PSNR) by different methods on BSD100.

	$\sigma_n = 10$										$\sigma_n = 20$									
	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM
001	35.97	28.42	36.98	37.63	37.76	37.26	37.46	37.67	37.71	<b>37.78</b>	31.97	28.00	33.03	33.39	33.58	33.16	33.27	33.49	33.51	<b>33.60</b>
002	37.01	31.23	39.51	40.41	40.60	39.73	39.75	40.59	40.55	<b>40.65</b>	32.26	30.77	35.41	36.14	36.42	35.89	36.17	36.34	36.40	<b>36.46</b>
003	38.18	31.41	39.34	40.05	40.31	39.28	40.27	40.15	40.20	<b>40.32</b>	33.92	30.83	35.42	35.86	36.31	35.52	36.27	36.06	36.11	<b>36.39</b>
004	37.27	31.50	38.50	39.03	39.19	38.53	38.56	39.12	39.16	<b>39.24</b>	33.26	30.77	34.43	34.65	34.48	34.43	34.48	34.80	34.64	<b>34.66</b>
005	36.13	29.13	37.59	38.37	<b>38.50</b>	37.90	38.24	38.47	38.36	<b>38.50</b>	31.73	28.73	33.51	33.99	34.23	33.76	34.16	34.13	34.11	<b>34.26</b>
006	39.06	34.79	40.37	40.98	41.47	40.25	41.17	41.22	41.31	<b>41.54</b>	34.36	33.73	36.26	36.83	37.35	36.57	37.86	37.11	37.52	<b>37.87</b>
007	36.05	28.22	37.09	37.76	37.88	37.32	37.28	37.88	37.79	<b>37.89</b>	31.98	27.78	33.03	33.42	33.57	33.14	33.41	33.57	33.55	<b>33.62</b>
008	35.82	27.75	36.85	37.45	37.63	37.04	37.58	37.61	37.63	<b>37.65</b>	31.77	27.39	32.76	33.01	33.23	32.79	33.29	33.17	33.22	<b>33.31</b>
009	36.77	29.14	37.91	38.66	38.78	38.09	38.36	38.77	38.76	<b>38.79</b>	32.82	28.67	33.82	34.38	34.46	33.99	34.23	34.52	34.49	<b>34.48</b>
010	37.27	30.88	38.44	38.94	39.14	38.43	38.80	39.05	39.09	<b>39.16</b>	33.19	30.25	34.41	34.59	34.88	34.36	34.80	34.78	34.81	<b>34.89</b>
011	36.56	28.14	37.28	38.28	38.30	37.76	38.26	38.34	38.27	<b>38.32</b>	32.61	27.76	33.30	34.19	34.23	33.78	34.21	34.32	34.29	<b>34.25</b>
012	36.19	28.64	37.11	37.76	37.92	37.28	37.71	37.85	37.88	<b>37.97</b>	32.16	28.16	33.05	33.37	<b>33.57</b>	33.09	33.26	33.49	33.52	<b>33.57</b>
013	36.74	30.93	37.76	38.34	38.51	37.94	38.36	38.45	38.51	<b>38.56</b>	32.72	30.25	33.72	33.99	34.22	33.81	33.99	34.11	34.16	<b>34.23</b>
014	38.39	31.84	39.60	40.35	40.50	39.53	40.55	40.46	40.52	<b>40.58</b>	34.20	31.28	35.56	36.14	<b>36.26</b>	35.76	36.15	36.24	36.25	<b>36.26</b>
015	36.45	28.53	37.44	38.04	38.23	37.59	38.26	38.12	38.20	<b>38.26</b>	32.37	28.16	33.44	33.77	33.97	33.51	33.68	33.86	33.89	<b>33.98</b>
016	34.95	25.25	35.85	36.83	36.95	36.35	36.61	36.93	36.95	<b>36.96</b>	30.96	25.04	31.85	32.55	32.71	32.30	32.65	32.71	32.76	
017	39.63	35.47	40.68	41.24	41.30	40.33	41.25	41.36	41.32	<b>41.37</b>	35.25	34.17	36.63	36.96	37.10	36.64	37.15	37.13	37.15	<b>37.18</b>
018	39.68	35.17	40.62	41.29	41.65	40.34	41.37	41.39	41.57	<b>41.73</b>	35.44	34.03	36.69	37.11	37.60	36.78	37.64	37.34	37.54	<b>37.64</b>
019	37.49	30.08	38.17	39.04	39.13	38.50	38.67	39.12	39.14	<b>39.16</b>	33.42	29.59	34.30	34.84	<b>35.05</b>	34.66	34.96	34.94	35.02	<b>35.05</b>
020	37.66	29.37	38.65	39.52	39.71	38.78	39.56	39.63	39.68	<b>39.75</b>	33.69	29.00	32.76	35.29	35.54	34.90	35.55	35.46	35.55	<b>35.61</b>
AVE.	37.16	30.29	38.29	39.00	39.17	38.41	38.90	39.11	39.13	<b>39.21</b>	33.00	29.27	34.17	34.72	34.49	34.44	34.84	34.88	34.92	<b>35.00</b>
	$\sigma_n = 40$										$\sigma_n = 60$									
	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM
001	27.62	26.14	29.15	29.45	<b>29.68</b>	29.48	29.64	29.38	29.42	<b>29.68</b>	25.01	23.01	26.99	27.28	27.53	27.46	26.99	27.14	<b>27.62</b>	
002	26.97	28.59	31.28	32.20	32.40	32.24	32.37	32.22	32.31	<b>32.49</b>	23.97	24.93	28.94	30.05	30.23	29.98	29.95	30.28		<b>30.87</b>
003	29.57	28.68	31.34	31.98	32.34	31.81	32.26	31.95	32.23	<b>32.39</b>	26.94	25.55	28.92	29.83	29.85	29.36	29.84	29.49	24.76	<b>29.94</b>
004	29.04	28.57	30.50	30.73	31.01	30.76	31.01	30.72	30.87	<b>31.04</b>	26.74	25.62	28.33	28.63	<b>28.91</b>	28.51	28.86	28.50	28.67	<b>28.91</b>
005	26.83	26.77	29.47	30.00	30.28	30.08	30.16	29.94	30.17	<b>30.30</b>	23.83	23.46	27.24	27.89	28.16	27.97	28.01	27.73	28.11	<b>28.32</b>
006	28.82	30.72	32.09	32.91	33.55	33.08	33.41	32.98	33.22	<b>33.62</b>	25.63	26.63	29.77	30.78	31.52	30.74	30.36	30.62	30.98	<b>31.53</b>
007	28.03	25.94	29.08	29.41	29.53	29.30	29.28	29.39	29.44	<b>29.54</b>	25.80	23.32	26.84	27.23	27.14	27.17	27.01	27.12	27.24	
008	27.87	25.82	28.85	29.11	29.22	29.04	29.19	29.02	29.22	<b>29.25</b>	25.49	23.55	26.68	27.02	26.96	26.85	26.73	27.01	<b>27.07</b>	
009	28.74	26.75	29.77	30.27	30.28	30.11	30.21	30.25	30.28	<b>30.31</b>	25.58	24.08	27.52	27.96	27.89	27.88	27.72	27.81	<b>28.01</b>	
010	28.83	28.26	30.43	30.60	30.88	30.66	30.89	30.58	30.67	<b>30.93</b>	26.67	25.45	28.17	28.45	28.69	28.62	28.23	28.35		<b>28.73</b>
011	28.08	25.75	29.25	30.20	30.21	29.86	30.16	30.15	30.19	<b>30.22</b>	24.81	22.39	26.91	<b>27.89</b>	27.79	27.43	27.69	27.77		<b>27.89</b>
012	28.04	26.33	29.04	29.27	29.49	29.23	29.35	29.19	29.21	<b>29.51</b>	25.78	23.81	26.82	27.09	27.17	27.10	26.98		<b>27.19</b>	
013	28.62	28.14	29.81	30.09	<b>30.34</b>	30.20	30.28	30.02	30.23	<b>30.34</b>	26.39	25.29	27.71	28.03	28.29	27.96	27.88	28.13	<b>28.35</b>	
014	29.81	29.21	31.37	32.12	32.11	31.90	32.14	32.11	32.14	<b>32.15</b>	27.12	26.28	29.02	29.83	29.67	29.52	29.54	29.67		<b>29.85</b>
015	28.13	26.44	29.52	29.78	29.96	29.77	29.74	29.72	29.87	<b>29.99</b>	25.57	23.83	27.35	27.62	27.80	27.67	27.44	27.68		<b>27.89</b>
016	26.96	23.43	27.88	28.48	28.55	28.20	28.54	28.43	28.46	<b>28.59</b>	24.66	20.87	25.59	26.20	26.16	25.87	25.99	26.15		<b>26.27</b>
017	29.47	31.49	32.63	33.10	33.35	33.12	33.30	33.13	33.24	<b>33.42</b>	25.89	25.82	30.45	31.05	31.33	30.83	31.02	30.75	31.22	<b>31.42</b>
018	30.63	31.32	32.60	33.25	33.74	33.30	33.74	33.34	33.46	<b>33.77</b>	27.80	27.56	30.20	31.10	31.49	30.83	31.26	30.81	31.23	<b>31.52</b>
019	28.99	27.54	30.61	31.20	31.46	31.21	31.45	31.18	31.27	<b>31.51</b>	26.02	24.48	28.47	29.25	29.42	29.06	29.09	29.35		<b>29.44</b>
020	29.33	27.16	30.83	31.32	31.58	31.17	31.57	31.26	31.46	<b>31.63</b>	26.08	24.41	28.57	29.08	29.30	28.90	28.97	28.31	28.46	<b>29.39</b>
AVE.	28.52	27.65	30.28	30.77	31.00	30.73	30.93	30.75	30.87	<b>31.03</b>	25.79	24.65	28.02	28.61	28.77	28.48	28.58	28.36	28.34	<b>28.87</b>
	$\sigma_n = 80$										$\sigma_n = 100$									
	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM	BM3D	NNM	EPLL	WNNM	PCLR	PGPD	STROLLR	TWSC	SNSS	AWNNM
001	22.81	24.14	25.53	25.80	26.07	25.86	25.74	25.44	25.65	<b>26.10</b>	20.80	23.64	24.41	24.68	25.00	24.67	<b>25.09</b>	24.26	24.38	24.73
002	21.83	24.57	27.27	28.41	28.65	28.26	28.26	28.30	28.44	<b>28.66</b>	20.04	23.63	25.94	27.17	<b>27.40</b>	27.03	27.33	26.90	27.12	27.29
003	24.76	24.70	27.20	28.15	<b>28.26</b>	27.60	27.83	27.66	28.13	<b>28.26</b>	22.9									

Figure 7. Denoising results on image 006 by different methods (noise level  $\sigma_n = 10$ ).Figure 8. Denoising results on image 014 by different methods (noise level  $\sigma_n = 20$ ).Figure 9. Denoising results on image 017 by different methods (noise level  $\sigma_n = 60$ ).

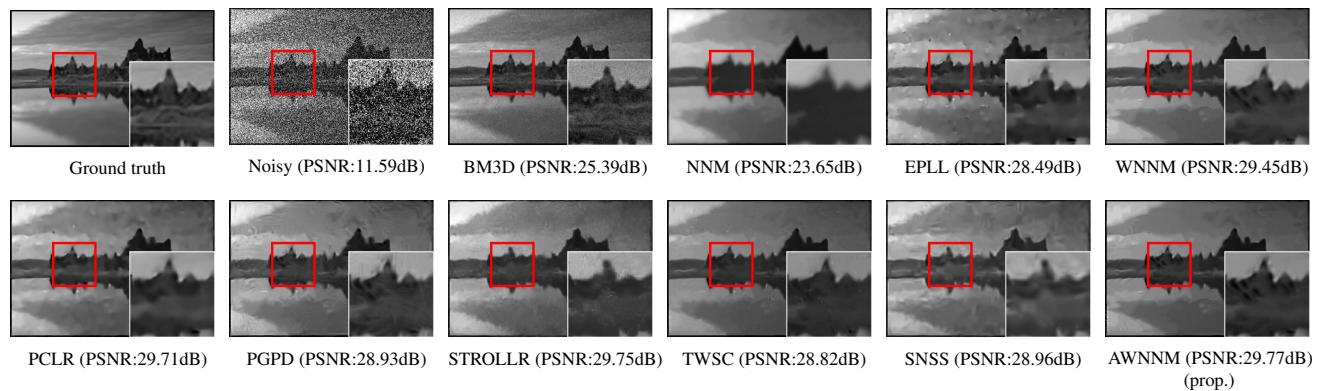


Figure 10. Denoising results on image 018 by different methods (noise level  $\sigma_n = 80$  ).