

Provable Low Rank Plus Sparse Matrix Separation Via Nonconvex Regularizers*

April Sagan[†] and John E. Mitchell[‡]

Abstract. This paper considers a large class of problems where we seek to recover a low rank matrix and/or sparse vector from some set of measurements. While methods based off of convex relaxations suffer from a (possibly large) estimator bias, and other nonconvex methods require the rank or sparsity to be known a priori, we utilize nonconvex regularizers to minimize the rank and l_0 norm without the estimator bias from the convex relaxation. We present a novel analysis of the alternating proximal gradient descent algorithm applied to such problems, and bound the error between the iterates and the ground truth sparse and low rank matrices. The algorithm and error bound can be applied to sparse optimization, matrix completion, and robust principal component analysis as special cases of our results.

Key words.

AMS subject classifications.

1. Introduction. In order to better understand large data-set and to make inferences about them, it is helpful to understand the underlying patterns in the dataset. Even when the underlying pattern is highly nonlinear, the data matrix can be approximated as being low rank, an observation that enables techniques to analyze the data in terms of a low dimensional latent space, such as principal Component Analysis (PCA), identifying outliers through Robust PCA (RPCA), and accurately inferring data points from very few observations of a data matrix through matrix completion.

Data analysis techniques based upon this low rank property have received much attention in the past decade, with impressive computational results on large matrices and theoretical results guaranteeing the success of RPCA and matrix completion. Many of these results are based off of minimizing the nuclear norm of a matrix (defined as the sum of the singular values) as a surrogate for the rank function, similar to minimizing the l_1 norm to promote sparsity in a vector.

While the convex relaxation is an incredibly useful technique in many applications, minimizing the nuclear norm of a matrix has been shown to introduce a (sometimes very large) estimator bias. Intuitively, we expect to see this bias because if we hope to recover a rank r matrix, we must impose enough weight on the nuclear norm term so that the $(r + 1)$ th singular value is zero. By the nature of the nuclear norm, this requires also putting weight on minimizing the first r singular values, resulting in a bias towards zero proportional to the spectral norm of the noise added to the true data matrix.

Fortunately, recent work has shown that the estimator bias from convex regularizers can be reduced (or even eliminated, for well conditioned matrices) by using nonconvex regularizers,

*Submitted to the editors DATE.

Funding: This work was funded by the National Science Foundation ...

[†]Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, NY (aprilsagan1729@gmail.com, <http://www.aprilsagan.net>).

[‡]Department of Mathematical Sciences, Rensselaer Polytechnic Institute, Troy, NY (mitchj@rpi.edu)

such as the Schatten-p norm or the minimax concave penalty (MCP). It has been shown that for sparse optimization, the nonconvexity introduced from these regularizers does not create a further burden in the optimization process – in the right circumstances, the nonconvex problem has just one minimizer. Similar results for rank minimization problems have been previously unavailable, a gap that we have aimed to fill in this paper.

1.1. Summary of Contributions. In this paper, we focus on the nonconvex, unconstrained optimization problem where we find a low-rank matrix $L \in \mathbb{R}^{d_1 \times d_2}$ and sparse vector $s \in \mathbb{R}^{d_2}$.

$$(1.1) \quad \min_{L,s} \frac{\lambda_L}{d_1 d_2} \Phi_{\gamma_L}(L) + \frac{\lambda_s}{d_s} \phi_{\gamma_s}(s) + \frac{1}{2n} \|\mathcal{A}_L(L) + A_S s - b\|_2^2$$

We denote $\phi : \mathbb{R} \rightarrow \mathbb{R}_+$ to be a concave function used to promote sparsity in both the singular values of L and individual entries in s . We overload the notation to allow for ϕ_{γ_s} to be a function of a vector $x \in \mathbb{R}^{d_s}$ whose range is \mathbb{R}_+ , and we denote $\Phi_{\gamma_L} : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}_+$ as a surrogate to the rank function:

$$\phi_{\gamma_s}(x) = \sum_{i=1}^{d_s} \phi_{\gamma_s}(x_i), \quad \Phi_{\gamma_L}(X) = \sum_{i=1}^{\min(d_1, d_2)} \phi_{\gamma_L}(\sigma_i(X)).$$

where $\sigma_i(X)$ denotes the i th largest singular value of X . We restrict our focus to nonconvex regularizers that are *amenable regularizers*, as described in [20], and defined below. Some key properties of amenable regularizers are stated in Appendix D.

Definition 1.1. A function $\phi_{\gamma}(t) : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is amenable if it satisfies the following criteria.

1. $\phi_{\gamma}(0) = 0$
2. ϕ_{γ} is non decreasing
3. For $t > 0$, the function $\frac{\phi_{\gamma}(t)}{t}$ is non increasing in t .
4. the function ϕ_{γ} is differentiable for all $t \neq 0$ and subdifferentiable at $t = 0$ with $\lim_{t \rightarrow 0^+} \phi'_{\gamma}(t) = 1$.
5. The function $\phi_{\gamma}(t)$ is ν weakly convex. That is, the function $\rho_{\nu} := \phi_{\gamma}(t) + \frac{\nu}{2} t^2$ is convex.

The linear mappings $\mathcal{A}_L : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^n$ and $A_S \in \mathbb{R}^{d_s \times n}$ serve as the *observation models* of the underlying low rank matrices and sparse vectors. Most commonly, we are interested in the observation model

$$[\mathcal{A}_{\Omega^{obs}}(X)]_k = X_{i_k, j_k}$$

for $(i_k, j_k) \in \Omega^{obs}$, where $\Omega^{obs} \subseteq \{1, \dots, d_1\} \times \{1, \dots, d_2\}$ is the set of indices where we have a measurement of the low rank matrix we hope to reconstruct.

We present a very simple alternating algorithm to find a stationary point of (1.1). Though similar algorithms have been previously studied and shown to converge, we present a novel analysis of this algorithm and show that not only does this algorithm linearly converge, but it converges to the exact low rank matrix L^* and sparse vector s^* when no Gaussian noise is present. Furthermore, when Gaussian noise is present, we obtain an error bound that matches the minimax optimal rate. Our bound greatly improves upon bounds obtained in previous results analysing the convex relaxation and quantify the observation based on computational results that nonconvex regularizers reduce the impact of noise on the quality of the estimator.

1.2. Related Works. The matrix completion problem is as follows: given the values of a low rank matrix for only a sparse set of indices, we seek to determine the rest of the values of the matrix. While the problem of finding the minimum rank matrix that fits the observations is NP hard in general, it has been shown that under some assumptions, the global minimizer to the convex problem

$$\min_{X \in \mathbb{R}^{d_1 \times d_2}} \|X\|_* \text{ s.t. } X_{ij} = M_{ij} \forall (i, j) \in \Omega^{obs}$$

is exactly M , where Ω^{obs} is the set of indices of M we have observed. If M is rank r and μ -incoherent (as defined in section 3.1), then with high probability for a set, Ω^{obs} of n indices chosen uniformly at random, M is the unique minimizer to the convex relaxation so long as $n > C\mu rd_1 \log^{1.2}(d_1)$ for some universal constant C [7, 4]. This condition was later improved to $n > C\mu rd_1 \log(d_1)$ [28].

Using a convex relaxation for Robust PCA has similar results. While PCA is a powerful technique, it has been shown to be less reliable when just a sparse set of data points are grossly corrupted, and so the goal of RPCA is to identify and remove such corruptions by separating the data matrix into the sum of a low rank and sparse matrix.

$$\min_{L, S} \|L\|_* + \lambda_0 \|S\|_1 \text{ st } L_{ij} + S_{ij} = M_{ij} \forall (i, j) \in \Omega^{obs}$$

The convex relaxation was shown to give the exact solution when every entry of M is observed in [10], and when only partially observed (under the same assumptions necessary in matrix completion) by [5].

In many cases of practical interest, our measurements may have some level of noise in addition to being only partially observed or having some corrupted entries. For matrix completion, we can relax the constraint using a penalty formulation as follows:

$$\min_{X \in \mathbb{R}^{d_1 \times d_2}} \lambda \|X\|_* + \sum_{(i,j) \in \Omega^{obs}} (X_{ij} - M_{ij})^2$$

Likewise, RPCA can be formulated as solving:

$$\min_{L, S \in \mathbb{R}^{d_1 \times d_2}} \lambda_L \|L\|_* + \lambda_s \|S\|_1 + \sum_{(i,j) \in \Omega^{obs}} (L_{ij} + S_{ij} - M_{ij})^2$$

Statistical guarantees on the performance of the first of these estimators are discussed in [6, 25, 24], and the later in [2]. Specific bounds and assumptions for these works are discussed in Section 4.

In order to reduce the estimator bias for l_1 minimization, [8] proposed an iteratively reweighted l_1 norm method to place more weight on minimizing smaller entries, and less on entries further from zero. This idea was generalized to minimizing any Amenable regularizes to promote sparsity. Theoretical results on the subject include algorithmic guarantees similar to the ones presented in this paper [33], and a proof that the nonconvex problem has no spurious local minimizers [20, 21].

Table 1

Table of common nonconvex regularizers used for l_0 and rank minimization, along with the associated proximal operator.

	$\phi_\gamma(x)$	$\text{prox}_{\phi_\gamma}^\tau(y)$
l_1	$ x $	$\text{sign}(y)(y - \frac{1}{\tau})_+$
MCP	$\begin{cases} x - \frac{x^2}{2\gamma} & x \leq \gamma \\ \frac{\gamma}{2} & x > \gamma \end{cases}$	$\text{sign}(y)\min\left(y , \frac{\gamma\tau}{\gamma\tau-1}(y - \frac{1}{\tau})_+\right)$

The same regularizers could be used as a surrogate to the rank function, as originally proposed by [23]. In [22, 9], the authors propose a generalization of the singular value thresholding algorithm proposed by [3], which was later applied to the problem of RPCA in [11, 18]. For the problem of matrix completion, the algorithms proposed by [35] and [29] achieve the fastest computational complexity in other state of the art methods.

Other approaches to low rank optimization rely on, instead of minimizing a surrogate to the rank function, constraining the matrix to be a given rank. This can be done using constrained optimization to optimize over the set of all rank r matrices [32], or by using the low-rank factorization of a matrix $X = UV^T$ for $U \in \mathbb{R}^{d_1 \times r}$ and $V \in \mathbb{R}^{d_2 \times r}$ [34, 30].

Previous work on theoretical results pertaining to rank and sparsity constrained methods have consisted of algorithmic guarantees ensuring that we can obtain a matrix sufficiently close to the ground truth low rank and/or sparse matrix for both matrix completion [17] and RPCA [26, 37, 36]. Additionally, both of these problems have been shown to have no spurious local minimizers, and so the ground truth matrices are the only minimizers under some assumptions [13] [12].

2. Alternating Proximal Gradient Descent Algorithm. Many different methods have been shown to be effective when minimizing nonconvex relaxations of the l_0 and rank functions, including iterative reweighted methods, and methods based off of low rank factorization. In this paper, we focus on the most commonly used technique: alternating proximal gradient descent.

Consider the objective function $F(x) = \phi(x) + g(x)$, where $g(x)$ is convex and differentiable, and $\phi(x)$ is weakly convex. Instead of minimizing $F(x)$ directly, the proximal gradient descent method approximates $g(x)$ by a quadratic, strongly convex function $\bar{g}_k(x)$ centered about the point x^k .

$$\bar{g}_k(x) = g(x^k) + \langle \nabla g(x^k), x - x^k \rangle + \frac{1}{2\tau} \|x - x^k\|^2$$

At each iteration, we now minimize the function $F_k(x) = \phi(x) + \bar{g}_k(x)$. For sufficiently small τ , this function is strongly convex, and the proximal gradient descent algorithm is guaranteed to converge.

The proximal gradient descent method applied to the function $F(x) = \phi(x) + g(x)$ iteratively solve the following problem:

$$x^{k+1} = \underset{x}{\operatorname{argmin}} \phi(x) + \frac{1}{2\tau} \|x - (x^k - \tau \nabla g(x^k))\|^2 := \text{prox}_\phi^\tau(x^k - \tau \nabla g(x^k))$$

where we defined the *proximal operator* of a function as the minimum of a combination of the function and the distance from a given point. For many functions ϕ that we are interested in, the proximal operator has a closed form solution, which we show in Table 1.

For each of the sparsity promoting regularizers in Table 1, the proximal operator is also dubbed as a *shrinkage operator* or a *thresholding operator* because when the input is less than τ , the output is 0. Otherwise, the input is moved towards zero or, for some nonconvex regularizers, is unchanged. So, we can view the proximal gradient algorithm as iteratively taking a step in the gradient direction of $g(x)$, and then applying the proximal operator to promote sparsity.

When applied to the optimization problem in Equation (1.1), we have

$$(2.1a) \quad \tilde{L}^{k+1} = L^k - \tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* (\mathcal{A}_L(L^k) + A_s s^k - b)$$

$$(2.1b) \quad L^{k+1} = \text{prox}_{\Phi_{\gamma_L}^{\tau_L \lambda_L}}(\tilde{L}^{k+1})$$

$$(2.1c) \quad \tilde{s}^{k+1} = s^k - \tau_S \frac{d_s}{n} A_s^T (\mathcal{A}_L(L^{k+1}) + A_s s^k - b)$$

$$(2.1d) \quad s^{k+1} = \text{prox}_{\phi_{\gamma_s}^{\tau_s \lambda_s}}(\tilde{s}^{k+1})$$

To further simplify the problem, the following proposition will allow L subproblem to be solved in each singular value separately.

Proposition 2.1. *Consider the optimization problem*

$$(2.2) \quad \min_X \sum_i \phi_\gamma(\sigma_i(X)) + \frac{1}{2\tau} \|X - Y\|_F^2$$

where ϕ_γ is a ν weakly convex function. If $\nu < \tau$, then equation (2.2) is strongly convex and the minimizer X^* has the same singular vectors as Y with singular values given by

$$\begin{aligned} \sigma_i(X^*) &= \underset{x}{\operatorname{argmin}} \phi_\gamma(x) + \frac{1}{2\tau} (x - \sigma_i(Y))^2 \\ &:= \text{prox}_{\phi_\gamma}^{\tau}(\sigma_i(Y)) \end{aligned}$$

where $\text{prox}_{\phi_\gamma}^{\tau}(\sigma_i(Y))$ is the proximal operator.

Proof. By convexity of ϕ ,

$$\phi(\lambda_i(X_2)) \geq \phi(\lambda_i(X_1)) + \phi'(\lambda_i(X_1))(\lambda_i(X_2) - \lambda_i(X_1))$$

Summing over all $i = 1, \dots, n$,

$$\Phi(X_2) \geq \Phi(X_1) + \langle \nabla \phi(\lambda(X_1)), \lambda(X_1) - \lambda(X_2) \rangle$$

By Corollary 1,

$$\Phi(X_2) \geq \Phi(X_1) + \langle U_1 \nabla \phi(\lambda(X_1)) U_1^T, X_1 - X_2 \rangle$$

And, as $\nabla \Phi(X) = U \nabla \phi(\lambda(X)) U^T$, Φ is convex.

Consider if the global optimizer X^* was not of the form $U\Sigma V^T$, for a diagonal matrix Σ . Let $\bar{X} = U\Sigma^*V^T$, where $\Sigma_{ii}^* = \sigma_i(X^*)$. By the Hoffman-Wielandt inequality (Corollary A.1),

$$\|X^* - Y\|_F^2 \leq \sum_i (\sigma_i(X^*) - \sigma_i(Y^*))^2 = \|X^* - Y\|^2.$$

150 And, because $\Phi(\bar{X}) = \Phi(X^*)$, \bar{X} must also be a global minimizer, contradicting the premise
151 that X^* is the sole global minimizer. ■

Proposition 2.1 tells us that L^k has the singular vectors of \tilde{L}^k and singular values given by

$$\sigma_i(L^k) = \text{prox}_{\phi_\gamma}^{\tau_L \lambda_L}(\sigma_i(\tilde{L}^k)).$$

Likewise, the subproblem in S can be solved in each entry of s individually.

Algorithm 2.1 Alternating Proximal Gradient Descent for Low-Rank Plus Sparse Optimization (APGD)

```

for  $k = 1, \dots$  do
   $g^{k+1} = -\tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* (\mathcal{A}_L(U^k \Sigma^k V^k) + A_s s^k - b)$ 
   $[U^{k+1}, \tilde{\Sigma}^{k+1}, V^{k+1}] = \text{LRSSVD}(U^k, \Sigma^k, V^k, g^{k+1})$ 
   $\Sigma^{k+1} = \text{prox}(\tilde{\Sigma}^{k+1})$ 
   $\tilde{s}^k = s^k - \tau_s \frac{d_s}{n} A_s^T (\mathcal{A}_L(U \Sigma V) + A_s s^k - b)$ 
   $s^{k+1} = \text{prox}^{\lambda_s}(\tilde{s}^{k+1})$ 
end for
```

152
153 The slowest operation in the alternating proximal gradient method is by far the singular
154 value decomposition. However, in practice we can reduce the number of operations by cal-
155 culating the truncated singular value decomposition only using the first r_0 singular values,
156 where r_0 is an upper bound on the rank, and enforce that the remaining singular values are
157 zero. Alternatively, we can calculate each singular value in descending order and stop when
158 a singular value falls below λ_L , as all remaining singular values will be set to zero by the
159 proximal operator. So, in the case of RPCA where each entry is observed, each iteration has
160 a computational complexity of $\mathcal{O}(d_1 d_2 r_0)$, which matches other state of the art methods.

161 In the case of matrix completion, however, only a sparse set of entries of the low rank
162 matrix are observed, which could be used to increase the efficiency by reducing the amount
163 of computation needed to find the singular value decomposition of L at each iterations. For
164 a low rank matrix with a low rank factorization $U\Sigma V$, we refer to the problem of finding the
165 SVD of the rank r approximation to the matrix $U\Sigma V + g$ for a sparse matrix g as Low Rank
166 plus Sparse SVD (LRSSVD), originally proposed by [16].

167 The LRSSVD task can be accomplished efficiently using the same methods as if we were
168 to find the SVD of any other matrix, such as the Power Iteration method. Recall that the
169 computational complexity of the Power Iteration is limited by the amount of operations needed
170 to multiply the matrix by a vector. Because the computational complexities of calculating
171 both $u(U\Sigma V^T + Y)$ for $u \in \mathbb{R}^{d_1}$ and $(U\Sigma V^T + Y)v$ for $v \in \mathbb{R}^{d_2}$ are $\mathcal{O}((d_1 + d_2)r_0 + n)$, we
172 can calculate the top r_0 singular values and vectors of $X + Y$ with only $\mathcal{O}((d_1 + d_2)r_0^2 + nr_0)$
173 operations. The other operations in Algorithm 2.1 take no more time than the LRPSSVD.

Algorithm 2.2 Singular Value Decomposition for a Low Rank Plus Sparse Matrix

Input: U, Σ, V, Y
Output: Singular value decomposition of $U\Sigma V^T + Y$
Initialization: $\tilde{V} = V$

- 1: **for** $k = 1, \dots$, **do**
 - 2: $\tilde{U}^k = (U\Sigma V^T + Y)\tilde{V}(\tilde{V}\tilde{V}^T)^\dagger$
 - 3: $\tilde{V}^k = (\tilde{U}\tilde{U}^T)Z^\dagger\tilde{U}^T(U\Sigma V^T + Y)$
 - 4: **end for**
 - 5: $[Q^U, R^U] = \text{QR}(\tilde{U}^k)$, $[Q^V, R^V] = \text{QR}(\tilde{V}^k)$
 - 6: $[U^R, \Sigma^R, V^R] = \text{SVD}(R^U \tilde{\Sigma} R^{V^T})$
 - 7: $\tilde{U} = Q^U U^R$, $\tilde{U} = Q^V V^R$ ($\tilde{U}, \tilde{\Sigma}, \tilde{V}$)
-

174 The gradient in the L direction, g^k , requires calculating $\Sigma_{ii}^k U_i^k (V_i^k)^T$ for each entry in the
 175 support of \mathcal{A}_L . In the case of matrix completion, this is nr operations, which matches the
 176 computational complexity per iteration for state of the art matrix completion algorithms.

177 **3. Analysis of APGD Algorithm.** In this section, we present the main result of the pa-
 178 per: a recursive bound on the difference of the iterates of the alternating proximal gradient
 179 algorithm and the ground truth low rank matrix and sparse vector. We present the bound for
 180 the most general case, and give results on specific problems in the following section.

181 **3.1. Restricted Isometry and Orthogonality Properties.** In order to bound the error
 182 in the output of our algorithm relative to the underlying ground truth low-rank and sparse
 183 matrices L^* and s^* , we must first make a number of assumptions about L^* , s^* , and the
 184 observation models \mathcal{A}_L and A_s .

185 First, we must assure that the low rank matrix L^* can be separated from a sparse matrix
 186 – that is, L^* is not sparse itself. Not only is this necessary for low-rank plus sparse decompo-
 187 sition, but for the problem of matrix completion, this assumption is necessary to assure that a
 188 sparse set of observations is a good representation of the entire matrix. For example, consider
 189 the matrix consisting of zeros in every entry besides one entry, where the value is 1. We must
 190 observe every entry in the matrix to assure that we can reconstruct the matrix exactly, due
 191 to the fact that we must observe the nonzero entry and every entry in its row and column.
 192 To exclude such ill-posed problems from our analysis, we will assume that L^* is *incoherent*,
 193 as defined as Chapter 1.

194 We define the projection of a matrix onto the sparse space Ω as

$$195 \quad (3.1) \quad \mathcal{P}_\Omega(X) = \begin{cases} X_{ij} & (i, j) \in \Omega \\ 0 & \text{else} \end{cases}$$

196 and onto the tangent space \mathcal{T} as

$$197 \quad (3.2) \quad \mathcal{P}_\mathcal{T}(X) = UU^T X + XVV^T - UU^T XVV^T$$

198 Next, we discuss the conditions that the observation models \mathcal{A}_L and A_s must satisfy
 199 in order to recover the ground truth low rank and sparse matrix, known as the restricted

isometry property. Loosely, the RIP states that for any two vectors in Ω (or matrices \mathcal{T}), we can obtain a sufficiently accurate estimate of the distance between the two through the observation model A_s (or \mathcal{A}_L).

Definition 3.1. *The linear mapping A_s satisfies the (α, κ) sparse Restricted Isometry Property if, for any x satisfying $\|x\|_0 \leq \alpha d_s$*

$$(1 - \kappa_S)\|x\|^2 \leq \tau_s \frac{n}{d_s} \|A_s x\|^2 \leq (1 + \kappa_S)\|A_s x\|^2$$

for some constant τ_s . Likewise, the linear mapping \mathcal{A} satisfies the (μ, r, κ) low rank Restricted Isometry Property if for any X in a μ -incoherent rank r tangent space \mathcal{T} ,

$$(1 - \kappa_L)\|X\|_F^2 \leq \tau_L \frac{n}{d_1 d_2} \|\mathcal{A}_L X\|_F^2 \leq (1 + \kappa_L)\|X\|_F^2$$

for some constant τ_L .

In some cases, it may be more useful to use the following characterization of the RIP, which bounds the difference between the operator $\tau \mathcal{A}^* \mathcal{A}$ and the identity operator when restricted to sparse vectors or low-rank and incoherent matrices.

Proposition 3.2. *For a matrix $A_s \in \mathbb{R}^{n \times d_s}$ satisfying the (α, κ_s) sparse RIP,*

$$\left\| \frac{\tau_s n}{d_s} \mathcal{P}_\Omega A_s^T A_s \mathcal{P}_\Omega - \mathcal{P}_\Omega \right\|^2 \leq \kappa_s.$$

Likewise, for any linear mapping $\mathcal{A}_L : \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}^n$ satisfying the (μ, r, κ_L) low rank RIP,

$$\left\| \frac{\tau_L n}{d_1 d_2} \mathcal{P}_\mathcal{T} \mathcal{A}_L^* \mathcal{A}_L \mathcal{P}_\mathcal{T} - \mathcal{P}_\mathcal{T} \right\|^2 \leq \kappa_L$$

Finally, we discuss the interplay between the set of sparse matrices and the low rank, incoherent tangent space, and their observation models. We hope to be able to separate the measurement vector b into two parts: one in the span of $\mathcal{A}_L \mathcal{P}_\mathcal{T}$, and one in the span of $A_s \mathcal{P}_\Omega$. In order for this to be able to be done quickly, we require there to be no non-trivial vectors in the intersection of the two sets, which is equivalent to saying that $\|\mathcal{P}_\Omega A_s^T \mathcal{A}_L \mathcal{P}_\mathcal{T}\| < 1$. Under some assumptions, this norm is actually close to zero, a concept we refer to as restricted orthogonality, which we define here and verify applies to the problems we are interested in in Section 4.

Definition 3.3. *Linear maps \mathcal{A}_L and A_s satisfy the restricted orthogonality property over the sets \mathcal{T} and Ω (respectively) when*

$$\frac{d_s}{n} \|\mathcal{P}_\Omega A_s^* \mathcal{A}_L \mathcal{P}_\mathcal{T}\|^2 \leq \kappa_{SL}, \quad \frac{d_1 d_2}{n} \|\mathcal{P}_\mathcal{T} \mathcal{A}_L^* A_s \mathcal{P}_\Omega\|^2 \leq \kappa_{LS}$$

3.2. Main Result. Define the difference between the iterates of the alternating proximal gradient descent algorithm and the ground truth low rank matrix and sparse vector at iteration k as $\Delta_L^k = L^* - L^k$ and $\Delta_s^k = s^* - s^k$. Our main result in the most general form gives a bound on the norm of Δ_L^k and Δ_s^k in terms of the differences at the previous iteration, Δ_L^{k-1} and Δ_s^{k-1} .

Theorem 3.4. Let L^k and s^k be the sequences generated by Algorithm 2.1. Assume that

$$b = \mathcal{A}_L(L^*) + A_s s^* + \mathcal{E} \in \mathbb{R}^n,$$

where $L^* \in \mathbb{R}^{d_1 \times d_2}$ is a rank r and μ -incoherent matrix, and $s^* \in \mathbb{R}^{d_s}$ is a sparse vector with $\text{supp}(s^*) = \Omega$, and the linear mappings \mathcal{A}_L and A_s satisfy the $(2r, 3\mu, \kappa_L)$ -low rank RIP and the (α, κ_s) -sparse RIP respectively, and together satisfy the ROP with constant κ_{L_s} . If $\lambda_L \geq \|\mathcal{A}_L^*(\mathcal{E})\|_2 + \|\mathcal{A}_L^* A_s \Delta_s^{k-1}\|_2$, then

$$\|\Delta_L^{k+1}\|_F^2 \leq \kappa_L \|\Delta_L^k\|_F^2 + \kappa_{L_s} \tau_L \|\Delta_s^k\|^2 + \frac{\tau_L d_1 d_2}{n} \|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 + \lambda_L r \phi'(\sigma_r(L^*))$$

Likewise, if $\lambda_S \geq \|A_s^T(\mathcal{E})\|_\infty + \|A_s^T \mathcal{A}_L \Delta_L^{k-1}\|_\infty$, then

$$\|\Delta_s^{k+1}\|^2 \leq \kappa_s \|\Delta_s^k\|^2 + \kappa_{sL} \tau_s \|\Delta_L^{k+1}\|_F^2 + \frac{\tau_s d_s}{n} \|\mathcal{P}_\Omega A_s^T \mathcal{E}\|^2 + \lambda_s \alpha d_s \phi'(s_{\min} - \lambda_s)$$

For each of these bounds, we can think of the first term as the estimation error introduced by the noise, and the second term as the approximation error, which accounts for the bias in the regularizer proportional to the derivative of the regularizer. Previous results for the nuclear norm and l_1 norm give similar bounds, but make the concession that the approximation error is the dominating term. Under some circumstances, that term is equal to zero in our bound.

3.3. Proof of Main Result. We start by presenting the the following two lemma regarding the proximal operator for the low rank regularizers, which we prove in the following section.

Lemma 3.5. Let L^* be a rank r , μ -incoherence matrix whose singular vectors form the tangent space T , and \bar{L} be defined as

$$(3.3) \quad \bar{L} = \underset{L \in \mathbb{R}^{d_1 \times d_2}}{\text{argmin}} \lambda \Phi(L) + \|L - L^* + \delta\|_F^2$$

Where Φ is an at most $\frac{1}{\lambda}$ weakly convex regularizer satisfying Assumption 1.1, and $\delta \in \mathbb{R}^{d_1 \times d_2}$ satisfies $\|\delta\|_2 \leq \lambda$. Define $\Delta_L = \bar{L} - L^*$. Then,

$$(3.4) \quad \|\Delta_L\|_F^2 \leq 2\|\mathcal{P}_T(\delta)\|_F^2 + \lambda r \phi'(\sigma_r(L^*))$$

In order to utilize the RIP and ROP conditions, we need to verify that L^{k+1} is low rank and incoherent, and that s^{k+1} is sparse, which we will do inductively. Assume that L^k is at most rank r , and that its tangent space is 2μ incoherent. Additionally, assume that $\text{supp}(S^k) \subseteq \text{supp}(S^*)$. Clearly, these conditions are met at the first iteration by initializing the algorithm with $L^0 = 0$ and $S^0 = 0$.

At iteration $k + 1$,

$$\begin{aligned} L^* - \bar{L}^{k+1} &= L^* - \left(L^k - \frac{\tau_L n}{d_1 d_2} \mathcal{A}_L^* (\mathcal{A}_L(L^k) + A_s(S^k) - b) \right) \\ &= L^* - L^k + \frac{\tau_L n}{d_1 d_2} \mathcal{A}_L^* \mathcal{A}_L L^k + \frac{\tau_L n}{d_1 d_2} \mathcal{A}_L^* A_s s^k - \frac{\tau_L n}{d_1 d_2} \mathcal{A}_L^* (\mathcal{A}_L(L^*) + A_s(s^*) + \mathcal{E}) \\ &= \left(\Delta_L^k - \frac{\tau_L n}{d_1 d_2} \mathcal{A}_L^* \mathcal{A}_L \Delta_L^k \right) - \frac{\tau_L n}{d_1 d_2} (\mathcal{A}_L^* A_s \Delta_s^k - \mathcal{A}_L^* \mathcal{E}) \end{aligned}$$

where the first equality comes from the definition of b , and the second inequality substitutes $\Delta_S^k = S^* - S^k$ and $\Delta_L^k = L^* - L^k$. By Lemma 3.5 (along with the triangle inequality), we have that

$$\begin{aligned} \|\Delta_L^{k+1}\|_F^2 &\leq \|\mathcal{P}_T(\Delta_L^k - \tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* \mathcal{A}_L \Delta_L^k)\|_F^2 + \tau_L \frac{d_1 d_2}{n} \|\mathcal{P}_T \mathcal{A}_L^* A_s \Delta_S^k\|_F^2 \\ &\quad + \tau_L \frac{d_1 d_2}{n} \|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 + \lambda_L r \phi'(\sigma_r(L^*) - 2\lambda_L \tau_L) \end{aligned}$$

Let \mathcal{T}^k denote the union of the tangent space of the rank r approximation L^k and \mathcal{T} so that $\Delta_L^k \in \mathcal{T}^k$. Then,

$$\begin{aligned} \|\mathcal{P}_T(\Delta_L^k - \tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* \mathcal{A}_L \Delta_L^k)\|_F^2 &\leq \|\mathcal{P}_{\mathcal{T}^k}(\Delta_L^k - \tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* \mathcal{A}_L \Delta_L^k)\|_F^2 \\ &\leq \|\mathcal{P}_{\mathcal{T}^k}(\tau_L \frac{d_1 d_2}{n} \mathcal{A}_L^* \mathcal{A}_L - \mathcal{I}) \mathcal{P}_{\mathcal{T}^k} \Delta_L^k\|_F^2 \leq \kappa_L \|\Delta_L^k\|_F^2 \end{aligned}$$

where the first inequality comes from the contractive property of \mathcal{P}_T , and the second inequality comes from the fact that $\mathcal{P}_{\mathcal{T}^k} \Delta_L^k = \Delta_L^k$. And, because \mathcal{T}^k has incoherence at most 3μ , we can apply the low-rank RIP to obtain the third inequality.

By the inductive hypothesis stating that $\text{supp}(S^k) \subseteq \text{supp}(S^*)$ (and thus, that S^k is α -sparse), we can use the ROP to claim that

$$\frac{d_1 d_2}{n} \|\mathcal{P}_T \mathcal{A}_L^* A_s \Delta_S^k\|_F^2 \leq \kappa_{SL} \|\Delta_S^k\|_F^2.$$

Combining these gives the desired bound on Δ_L^k .

Now, we must show that L^{k+1} is rank r and 2μ incoherent. By Weyl's inequality ([15], see Appendix A), we know that

$$\begin{aligned} \sigma_{r+1}(\tilde{L}^{k+1}) &\leq \sigma_{r+1}(L^*) + \|L^* - \tilde{L}^{k+1}\|_2 \\ &\leq \|\Delta_S^k\|_2 + \|\mathcal{A}_L^* \mathcal{E}\|_2 \end{aligned}$$

Because this is less than λ_L by our assumptions, L^{k+1} must be rank r .

In order to show that \tilde{L}^{k+1} is at most 2μ incoherent, we apply the following Lemma, which utilizes the Davis-Kahan inequality [?].

Lemma 3.6. *If $X \in \mathbb{R}^{d_1 \times d_2}$ (with $d_1 \geq d_2$) is a rank r , μ -incoherent matrix, and $\Delta \in \mathbb{R}^{d_1 \times d_2}$ satisfies $\|\Delta\|_2 \leq \frac{1}{2} \frac{\mu r}{d_1} \sigma_r(X)$, then the top r singular vectors of the matrix $X + \Delta$ form a 2μ -incoherent tangent space.*

Next, we bound Δ_s^k and showing $s^{k+1} \in \Omega$ in a similar manner. We present the following Lemma, which mirrors Lemma 3.5.

Lemma 3.7. *Let $s^* \in \mathbb{R}^{d_s}$ be an sparse vector with support Ω , and let \bar{S} be defined as*

$$(3.5) \quad \bar{s} = \underset{s \in \mathbb{R}^{d_s}}{\text{argmin}} \quad \lambda \phi(s) + \|s - s^* + \delta\|^2$$

Where ϕ is an at most $\frac{1}{\lambda}$ weakly convex amenable regularizer, and $\delta \in \mathbb{R}^{d_s}$ satisfies $\|\delta\|_\infty \leq \lambda$.

Then, $\text{supp}(\bar{s}) \subseteq \Omega$. Furthermore, we can bound the difference $\Delta_s = \bar{s} - s^*$ as

$$(3.6) \quad \|\Delta_s\|_F^2 \leq \|\mathcal{P}_\Omega(\delta)\|_F^2 + 2\lambda \sqrt{|\Omega|} \phi'(s_{\min} - \lambda).$$

By the update equation for \tilde{s}^{k+1} , we have

$$\begin{aligned}
 s^* - \tilde{s}^{k+1} &= s^* - (s^k - \mathcal{A}_s^*(\mathcal{A}_s(s^k) + \mathcal{A}_L(L^k) - b)) \\
 &= s^* - s^k + \frac{\tau_s}{d_s} \mathcal{A}_s^* \mathcal{A}_s s^k + \frac{\tau_s}{d_s} \mathcal{A}_L^* \mathcal{A}_L L^k - \frac{\tau_s}{d_s} \mathcal{A}_s^*(\mathcal{A}_s(s^*) + \mathcal{A}_L(L^*) + \mathcal{E}) \\
 &= (\Delta_s^k - \frac{\tau_s}{d_s} \mathcal{A}_s^* \mathcal{A}_s \Delta_s^k) - \frac{\tau_s}{d_s} \mathcal{A}_s^* \mathcal{A}_L \Delta_L^{k+1} - \frac{\tau_s}{d_s} \mathcal{A}_s^* \mathcal{E}
 \end{aligned}$$

By Lemma 3.7,

$$\begin{aligned}
 \|\Delta_s^{k+1}\|_F^2 &\leq \|\mathcal{P}_\Omega(\Delta_s^k - \mathcal{A}_s^* \mathcal{A}_s \Delta_s^k)\|_F^2 + \|\mathcal{P}_\Omega \mathcal{A}_s^* \mathcal{A}_L \Delta_L^{k+1}\|_F^2 \\
 &\quad + \|\mathcal{P}_\Omega \mathcal{A}_s^* \mathcal{E}\|_F^2 + \lambda_s |\Omega| \phi'_s(s_{\min} - \lambda_s)
 \end{aligned}$$

Applying the RIP to the first term gives:

$$\|\mathcal{P}_\Omega(\Delta_s^k - \frac{\tau_s}{d_s} \mathcal{A}_s^* \mathcal{A}_s \Delta_s^k)\|_F^2 \leq \kappa_s \|\Delta_s^k\|_F^2$$

And, applying the ROP to the second term gives us:

$$\|\frac{\tau_s}{d_s} \mathcal{P}_\Omega \mathcal{A}_s^* \mathcal{A}_L \Delta_L^{k+1}\|_F^2 = \|\frac{\tau_s}{d_s} \mathcal{P}_\Omega \mathcal{A}_s^* \mathcal{A}_L \mathcal{P}_{T^{k+1}} \Delta_L^{k+1}\|_F^2 \leq \frac{\tau_s}{d_s} \kappa_{sL} \|\Delta_L^{k+1}\|_F^2$$

Combining these terms gives us the desired result.

3.4. Proofs of Supporting Lemmas.

Proof of Lemma 3.7. First, we will show that if $s_i^* = 0$, then $\bar{s}_i = 0$. For i not in Ω ,

$$\bar{s}_i = \operatorname{argmin}_s \lambda \phi(s) + \frac{1}{2}(s - \delta_i)^2$$

The sub-gradient of the objective function evaluated at zero is $\{\delta_i + g \mid |g| \leq \lambda\}$. By our assumption that $\lambda \geq \delta_i$, 0 must be in this set, and so $\bar{s}_i = 0$ is a stationary point. Because the objective function is strongly convex by assumption, this is a global minimizer.

Next, consider if $i \in \Omega$. By first order necessary conditions for optimality,

$$\bar{s}_i = s_i^* + \delta_i - \phi'(\bar{s}_i) \leq s_i^* - \delta_i - \phi'(s_{ij}^* + \delta_i)$$

This gives the bound:

$$|\bar{s}_i - s_i^*| \leq |\delta_i| + \phi'(s_{\min} - \lambda)$$

Combining these facts,

$$\|\bar{s} - s^*\|^2 = \|\mathcal{P}_\Omega(\bar{s} - s^*)\|^2 \leq \|\mathcal{P}_\Omega(\delta)\|^2 + |\Omega| \phi'(s_{\min} - \lambda)^2 \quad \blacksquare$$

Proof of Lemma 3.5. Let

$$f(L) = \lambda \Phi(L) + \frac{1}{2} \|L - L^* + \delta\|_F^2$$

be the function that \bar{L} is the global minimizer of, that is, $0 \in \partial f(\bar{L})$. The subgradient of f at L^* is as follows:

$$\left\{ \lambda(U\Sigma^\phi V^T + W) + \delta \mid W \in T^\perp, \|W\|_2 \leq 1 \right\}$$

where $L^* = U\Sigma V^T$ and

$$\Sigma^\phi = \text{diag}(\phi'(\sigma_1(L^*)), \phi'(\sigma_2(L^*)), \dots, \phi'(\sigma_r(L^*))).$$

327 Because Φ is $\frac{1}{2\lambda}$ -weakly convex f is $\frac{1}{2}$ -strongly convex, which gives us:

328
$$\|L^* - \bar{L}\|_F \leq \inf 2\|\partial f(L^*)\|_F.$$

330 Consider the subgradient given by $W = -\frac{1}{\lambda}\mathcal{P}_{T^\perp}(\delta)$. Note the that $\|W\|_2 \leq 1$ as we assume
331 $\lambda \geq \|\delta\|_2$.

332
$$\begin{aligned} \|L^* - \bar{L}\|_F &\leq 2\|\lambda U\Sigma^\phi V^T + \delta - \mathcal{P}_{T^\perp}(\delta)\|_F \\ &\leq 2\lambda\|U\Sigma^\phi V^T\|_F + 2\|\delta - \mathcal{P}_{T^\perp}(\delta)\|_F \\ &\leq 2\lambda\sqrt{r}\phi'(\sigma_r(L^*)) + 2\|\mathcal{P}_T(\delta)\|_F \end{aligned}$$

333
334
335 ■

336 The second inequality is the triangle inequality, and the third uses the fact that $\|\Sigma^\phi\|_2 =$
337 $\phi'(\sigma_r(L^*))$

Proof of Lemma 3.6. Let $U \in \mathbb{R}^{d_1 \times r}$ and $\tilde{U} \in \mathbb{R}^{d_1 \times r}$ denote the (top r) right singular vectors of X and $X + \Delta$ respectively. By the Davis Kahan theorem,

$$\text{dist}(U, \tilde{U}) \leq \frac{\|\Delta\|_2}{\sigma_r(X) + \|\Delta\|_2} \leq \frac{1}{2} \frac{\mu r}{d_1}$$

where the second inequality uses the fact that $\|\Delta\|_2 \leq \frac{\mu r}{2d_1} \sigma_r(X)$. Let u_i and \tilde{u}_i be the i^{th} eigenvector of X and $X + \Delta$ respectively, and let θ_i be the angle between the vectors.

$$\max(\tilde{u}_i) \leq \max(u_i) + 2\sin(\theta_i) \leq \frac{\mu r}{d_1} + \frac{\mu r}{d_1} = \frac{2\mu r}{d_1}$$

where we use the fact that

$$\text{dist}(U, \tilde{U}) = \max_i \sin(\theta_i)$$

338 So, the rank r approximation of $X + \Delta$ is 2μ incoherent. ■

339 **4. Results for Specific Models.** In this section, we utilize Theorem 3.4 to analyze an
340 application of the alternating proximal gradient descent algorithm to the problems of matrix
341 completion and RPCA.

342 **4.1. Matrix Completion.** We start by considering the problem of matrix completion.
343 Here, we have a sparse set of observed entries of L , $\mathcal{A}_L = \mathcal{A}_\Omega$, and we do not consider a sparse
344 vector s (i.e. $A_s = 0$).

345 We present a version of the RIP for the sampling operator from [28].

Lemma 4.1 ([28]). *Let Ω be a set of n entries of $\{1, \dots, d_1\} \times \{1, \dots, d_2\}$ drawn independently at random with uniform probability, with $n > 64\mu r(d_1 + d_2)\log(d_2)$. Then, with probability at least $1 - 2d_2^{-2}$,*

$$\frac{5}{6}\|X\|_F^2 \leq \frac{d_1 d_2}{n}\|\mathcal{A}_\Omega(X)\|^2 \leq \frac{7}{6}\|X\|_F^2$$

for any rank r , μ -incoherent matrix X .

Additionally, we will assume that the additive noise \mathcal{E} has entries that are mean zero i.i.d. variables. The effect the noise has on the estimator is reduced due to the fact that very little of $\mathcal{A}_\Omega^*(\mathcal{E})$ will lie in the tangent space \mathcal{T} . To formalize this intuition, we cite the following lemma from [37].

Lemma 4.2. *Assume that the entries of Ω are chosen uniformly at random from $\{1, \dots, d_1\} \times \{1, \dots, d_2\}$, and the entries of \mathcal{E} are mean zero i.i.d. random variables with variance ν^2 . For some universal constants C_1 and C_2 , the following hold:*

$$\|\mathcal{A}_\Omega^*(\mathcal{E})\|_2 \leq C_1 \nu \sqrt{pd \log(d)} \quad \text{and} \quad \|\mathcal{A}_\Omega^*(\mathcal{E})\|_\infty \leq C_2 \nu \sqrt{p \log(d)}$$

We now present an error bound of for a stationary point of Algorithm 1 applied to matrix completion.

Theorem 4.3. *Let $b = \mathcal{A}_\Omega X^* + \mathcal{E}$ for a rank r , μ -incoherent matrix $X^* \in \mathbb{R}^{d_1 \times d_2}$, and the entries of \mathcal{E} are mean zero i.i.d. random variables with variance ν^2 . There exists universal constants C_1 and C_2 such that under the same assumptions as Lemma 4.1, if $\lambda > C_1 \nu \sqrt{pd \log(d)}$, then the iterates of Algorithm 1 linearly converge to a point \bar{X} such that $\|X^* - \bar{X}\|_F^2$ is less than:*

$$C_2 \frac{d_1 d_2}{n} \left(\underbrace{r \nu^2 d \log(d)}_{\text{optimal error rate}} + \underbrace{\lambda r \phi'(\sigma_r(X^*))}_{\text{bias term}} \right)$$

with convergence rate $\frac{1}{6} \frac{d_1 d_2}{n}$.

The two terms of the error bound account for the optimal error rate and a bias term. The optimal error rate is the error bound if we know the tangent space of X^* a priori, that is, the difference between X^* and the solution to the optimization problem

$$\min_{X \in \mathcal{T}} \|\mathcal{A}_\Omega(X) - b\|_F^2.$$

The oracle rate is further discussed in [6] and [25].

Proof. We can apply Theorem 1 with $s^k = s^* = 0$, $\tau_L = 1$, and $\kappa_L = \frac{1}{6}$ (by Lemma 3) to get the bound

$$\|\Delta^{k+1}\|_F^2 \leq \frac{1}{6}\|\Delta^k\|_F^2 + \left(\frac{d_1 d_2}{n}\|\mathcal{P}_\mathcal{T} \mathcal{A}_\Omega^* \mathcal{E}\|_F^2 + \lambda r \phi'(\sigma_r(X^*) - 2\lambda)\right)$$

Initializing with $X^0 = 0$, we have the error at each iteration as follows:

$$\|\Delta^k\|^2 \leq \left(\frac{1}{6}\right)^k \|X^*\|_F^2 + \frac{5}{6} \left(1 - \left(\frac{1}{6}\right)^k\right) \left(\frac{d_1 d_2}{n}\|\mathcal{P}_\mathcal{T} \mathcal{A}_\Omega^* \mathcal{E}\|_F^2 + \lambda r \phi'(\sigma_r(X^*) - 2\lambda)\right)$$

355 Taking the limit as $k \rightarrow \infty$, and applying the bounds from Lemma 4.2 gives the desired
 356 result. ■

357 Perhaps counter-intuitively, a choice of step size, τ , that minimizes the loss function (i.e.
 358 the steepest descent step size) is not always the step size that leads to the fastest convergence
 359 rate. To see this, we compare the error bound at iteration k in both cases. The steepest
 360 descent step size ($\tau = n$) would give

$$361 \quad L^* - \tilde{L}^k = L^* - L^{k-1} - \mathcal{P}_\Omega(L^* - L^{k-1}) = P_{\Omega^c}(\Delta_L)$$

363 and the stepsize informed by the RIP ($\tau = \frac{d_1 d_2}{n}$) gives

$$364 \quad L^* - \tilde{L}^k = L^* - L^{k-1} - \frac{d_1 d_2}{n} \mathcal{P}_\Omega(L^* - L^{k-1}) = \Delta_L - \frac{d_1 d_2}{n} P_\Omega(\Delta_L).$$

While the $\tau = 1$ gives a significantly smaller error when simply comparing $L^* - \tilde{L}^k$, the
 error bound for $L^* - L^k$ comes from projecting $L^* - \tilde{L}^k$ onto \mathcal{T} . Without further information,
 the best bound we can get when using $\tau = n$ would be

$$\|L^* - L^k\|_F^2 \leq 1 - \frac{(1 - \kappa)n}{d_1 d_2} \|L^* - L^{k-1}\|_F^2.$$

366 This convergence rate approaches 1 asymptotically when we consider the information theo-
 367 retic minimum number of measurement for large d_1 and d_2 . However, when $\tau = \frac{d_1 d_2}{n}$, the
 368 convergence rate remains constant:

$$369 \quad \|L^* - L^k\|_F^2 \leq \|\Delta_L - \frac{d_1 d_2}{n} \mathcal{P}_\mathcal{T} \mathcal{P}_\Omega \mathcal{P}_\mathcal{T}(\Delta_L)\|_F^2 \leq \kappa \|\Delta_L\|_F^2$$

371 The first inequality uses the fact that $\Delta_L \in \mathcal{T}$ and the second uses the RIP.

372 **4.2. Robust PCA.** Next, we will use Theorem 3.4 to analyze the APGD algorithm applied
 373 to the problem of RPCA. Specifically, we are interested in the special case of (1.1) where the
 374 $A_s = I_n$, and $A_L = \mathcal{A}_{\Omega^{obs}}$.

375 In order for RPCA to be possible, we need the nonzero entries of $\mathcal{A}_{\Omega^{obs}}^* s^*$ to be sufficiently
 376 well-distributed throughout the rows and columns – if the sparse corruptions affected the same
 377 row or column of L^* , then this would also be a low-rank perturbation and thus be impossible
 378 to separate from L^* without further information. So, we will assume that $\mathcal{A}_L^* s^*$ is α -sparse,
 379 defined as follows.

380 **Definition 4.4.** *The matrix S is α -sparse for $0 < \alpha < 1$ if the proportion of nonzero entries*
 381 *in any row or column is less than α . That is,*

$$382 \quad (4.1) \quad \|S_{i:}\|_0 \leq \alpha d_1, \|S_{:j}\|_0 \leq \alpha d_2 \quad \forall i, j$$

383 In order to verify that the ROP property holds, we present the following Lemma:

384 **Lemma 4.5.** *Let \mathcal{T} be a rank r , μ -incoherent tangent space, and let Ω be an α -sparse*
 385 *subspace. Then,*

$$386 \quad (4.2) \quad \|\mathcal{P}_\mathcal{T} \mathcal{P}_\Omega\| \leq 2\alpha\mu r, \quad \|\mathcal{P}_\Omega \mathcal{P}_\mathcal{T}\| \leq 2\alpha\mu r$$

Proof of Lemma 4. By the triangle inequality, for a matrix $S \in \Omega$,

$$\|P_T(S)\|_F^2 \leq \|UU^T S\|_F^2 + \|SVV^T\|_F^2 + \|UU^T SVV^T\|_F^2$$

Because U is an orthonormal matrix,

$$\begin{aligned} \|UU^T S\|_F^2 &= \text{trace}(UU^T SS^T UU^T) = \text{trace}(U^T UU^T SS^T U) \\ &= \text{trace}(U^T SS^T U) = \|U^T S\|_F^2 \end{aligned}$$

We now expand this norm and use the incoherence property to obtain the desired result:

$$\begin{aligned} \|U^T S\|_F^2 &= \sum_{k=1}^r \sum_{i=1}^{d_1} \langle U_k, S_i \rangle^2 \leq \sum_{k=1}^r \sum_{i=1}^{d_1} \left(\sum_{j \in \Omega_i} U_{kj}^2 \right) \|S_i\|^2 \\ &= \sum_{i=1}^{d_1} \|S_i\|^2 \sum_{j \in \Omega_i} \|U_j\|^2 \leq \sum_{i=1}^{d_1} \|S_i\|^2 \alpha \mu r \leq \alpha \mu r \|S\|_F^2 \end{aligned}$$

We can now give a bound for the stationary point of the APGD algorithm applied to RPCA.

Theorem 4.6. Let $b = \mathcal{A}_{\Omega^{obs}}(L^*) + s^* + \mathcal{E}$ for a rank r , μ -incoherent matrix $L^* \in \mathbb{R}^{d_1 \times d_2}$, a sparse vector $s^* \in \mathbb{R}^n$ with $\|s^*\|_\infty \leq 2\|L^*\|_\infty$, and a vector $\mathcal{E} \in \mathbb{R}^n$ whose entries are independent Gaussian variables with mean 0 and standard deviation σ . Under the same assumptions on Ω^{obs} as Lemma 4.1, and assuming that $\mathcal{A}_\Omega^* s^*$ is α -sparse and $\alpha \mu r \leq \frac{1}{64}$, if $\lambda_L \geq \frac{1}{6} + \|\mathcal{A}_\Omega^* \mathcal{E}\|_2$ and $\lambda_s \geq \frac{\mu r}{d_1} + \|\mathcal{E}\|_\infty$, then the iterates of Algorithm 1 linearly converge to a point \bar{L}, \bar{s} satisfying

$$\|\Delta_L\|_F^2 \leq C_1 \frac{d_1 d_2}{n} r \nu^2 d \log(d), \quad \|\Delta_S\|_F^2 \leq C_2 \frac{d_1 d_2}{n} r \nu^2 d \log(d)$$

with convergence rate $\frac{1}{6}$.

In order to apply Theorem 1, we first show the following properties about the parameters λ_L and λ_s .

$$\begin{aligned} \lambda_L &\geq \|\mathcal{A}_L^* \mathcal{E}\|_2 + \|\mathcal{A}_L^* s^*\|_2 \\ &\geq \|\mathcal{A}_L^* \mathcal{E}\|_2 + \|\mathcal{A}_L^* \delta_s^k\|_2 \quad \forall k \end{aligned}$$

where the first inequality comes from the assumption on λ_L and the second from the fact that the two norm of s^k is decreasing.

$$\begin{aligned} \lambda_s &\geq \|\mathcal{E}\|_\infty + \frac{\mu r}{d_1 d_2} \\ &\geq \|\mathcal{E}\|_\infty + \|L^*\|_\infty \\ &\geq \|\mathcal{E}\|_\infty + \|\Delta_L^k\|_\infty \quad \forall k \end{aligned}$$

where the first inequality is our assumption on λ_s , the second follows from incoherence of L^* , and the last from the fact that

$$\|\Delta_L^k\|_\infty \leq \|\Delta_L^0\|_\infty = \|L^*\|_\infty.$$

In order to show the bias term is 0, we use our assumption that $\lambda_L \leq \frac{1}{4}\sigma_r(L^*)$ and the fact that the MCP regularizer can satisfy $\phi'_L(3\lambda_L) = 0$ while still maintaining $\frac{1}{2\lambda_L}$ -weak convexity. Likewise, when we assume $\lambda_s \leq \frac{1}{4}s_{min}$, we can choose ϕ_s such that $\phi'_s(s_{min} - \lambda_s) \leq \phi'_s(3\lambda_s) = 0$ while still maintaining $\frac{1}{2\lambda_s}$ weak convexity.

We can now apply the first part of Theorem 1 to obtain the bound:

$$\begin{aligned} \|\Delta_L^{k+1}\|_F^2 &\leq \frac{1}{6}\|\Delta_L^k\|_F^2 + \frac{1}{64}\|\Delta_s^k\|_F^2 + \frac{d_1 d_2}{n}\|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 \\ \|\Delta_S^{k+1}\|_F^2 &\leq \frac{1}{64}\|\Delta_L^{k+1}\|_F^2 + \|\mathcal{P}_\Omega \mathcal{E}\|_F^2 \\ &\leq \frac{1}{384}\|\Delta_L^k\|_F^2 + \frac{1}{4096}\|\Delta_s^k\|_F^2 + \frac{d_1 d_2}{64n}\|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 + \|\mathcal{P}_\Omega \mathcal{E}\|_F^2 \end{aligned}$$

The limit point of this sequence gives:

$$\begin{aligned} \|\Delta_L\|_F^2 &\leq \frac{6}{5} \frac{d_1 d_2}{n} \|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 \\ \|\Delta_S\|_F^2 &\leq \frac{d_1 d_2}{32n} \|\mathcal{P}_T \mathcal{A}_L^* \mathcal{E}\|_F^2 + \frac{6}{5} \|\mathcal{P}_\Omega \mathcal{E}\|_F^2 \end{aligned}$$

From here, we can apply the bounds from Lemma 4.2 to get the desired results.

5. Numerical Results. We implemented Algorithm 1 in Matlab R2020a, for which the code is available at GitHub. All results in this section are obtained with the Matlab version in order to accurately compare to other algorithms which are only available in Matlab, and are run on a Windows 10 desktop with an AMD Phenom 3.40 GHz processor and 8 Gb of RAM.

5.1. Matrix Completion. We compare our method for matrix completion to another method utilizing nonconvex regularizer from [35] (FaNCL), along with a method to minimize the nuclear norm IALM, from [19], and a rank constrained method, LMaFit [34]. The results are shown in Table 2.

We start by comparing the performance of the methods on randomly generated low rank matrices of varying size, rank, percentage of observed entries, and standard deviation of the noise in the measurements (shown relative to the mean absolute value of the low rank matrix). In each of the cases, our method performs exactly as well as FaNCL and LMaFit when the correct rank is given. IALM performs equally well in the first case, and slightly worse than the remaining three cases due to the fact that the nuclear norm biases the result towards zero.

Next, we show the results on common data sets for recommendation systems, the Jester data set [14] and MovieLens 1M [1]. For each of these two data sets, we partition the observations into five folds, fit a low rank model to four of the folds and calculate the accuracy on the remaining fold. We do this for each of the five folds and present the average normalized mean absolute error. In both cases, our method outperforms the other three algorithms we compare to.

Table 2

Comparison of four different matrix completion algorithms on randomly generated low rank matrices and common recommendation data sets. The algorithm LMaFit reconstructs a matrix of a given rank k . The table shows the results when the algorithm is given the exact rank ($k = r$) and an incorrect rank ($k = 2r$). Time is given in seconds

Jester Dataset									ML 1M			
d_1	1000		1000		5000		5000		24983		6040	
d_2	500		500		1000		1000		100		3952	
r	5		5		10		10		-			
n/d_1d_2	0.3		0.1		0.2		0.05		0.58		0.034	
std(\mathcal{E})	0.1		0.02		0.1		0.02		-			
	RFNE	T	RFNE	T	RFNE	T	RFNE	T	NMAE	T	NMAE	T
APGD	3.28e-4	0.7	2.90e-4	1.1	1.69e-4	10	1.96e-4	29	0.159	21	0.172	172
FaNCL	3.28e-4	1.8	2.90e-4	4.3	1.69e-4	31	2.49E-4	57	0.183	42	0.200	42
IALM	3.28e-4	2.6	2.92e-4	2.8	1.72e-4	32	1.99e-4	27	0.163	77	0.183	216
LMaFit ($2r$)	4.68e2	21	1.20e3	6.2	4.99e2	197	1.60e3	40	-			
LMaFit(r)	3.28e-4	0.4	2.90e-4	0.3	1.69e-4	3.9	1.96e-4	3.9	0.168	7.4	9.174	44

Table 3

Accuracy and run-time for RPCA on four videos of fish.

	Marine Snow		Small Aquaculture		Caustics		Two Fish		Fish Swarm	
	Acc	Time (s)	Acc	Time (s)	Acc	Time (s)	Acc	Time (s)	Acc	Time (s)
APGD	0.92	87	0.92	592	0.88	459	0.90	467	0.75	473
LMaFit	0.59	258	0.65	316	0.7	273	0.62	291	0.55	310
AltProj	0.84	143	0.62	109	0.86	83	0.76	109	0.62	130
RPCA-GD	0.92	94	0.72	77	0.87	76	0.96	96	0.66	93
IALM	0.57	379	0.56	286	0.56	288	0.56	342	0.56	305

5.2. Robust PCA. We compare our method to several other prominent RPCA methods, including LMaFit [30], AltProj [26], RPCA-GD [36], and IALM [19]. We compared with many other methods included in the LRSLibrary [31], however we only include results from the aforementioned algorithms as they gave the most accurate results for matrices a significant amount with Gaussian noise, a test case we emphasise in this section.

It is worth noting that, while our algorithm does not require an estimate of the rank of L^* or the sparsity of S^* a priori, BLANK requires both, and LMaFit and AltProj require an estimate of the rank. However, we found that LMaFit and AltProj still perform very well when this estimate is unreliable, as LMaFit includes a rank-estimation scheme and AltProj starts by performing a rank 1 projection, and increases the rank up until the estimate given. We provide all methods with an upper bound on the rank equal to twice the rank of L^* and an upper bound on the number of corrupted entries equal to twice that of S^* .

The most commonly utilized test case in the literature for RPCA is the task of separating the background and foreground of a video. In this scenario, each frame of the video is represented as a column vector in M . If the background is static (or, at least, in some way

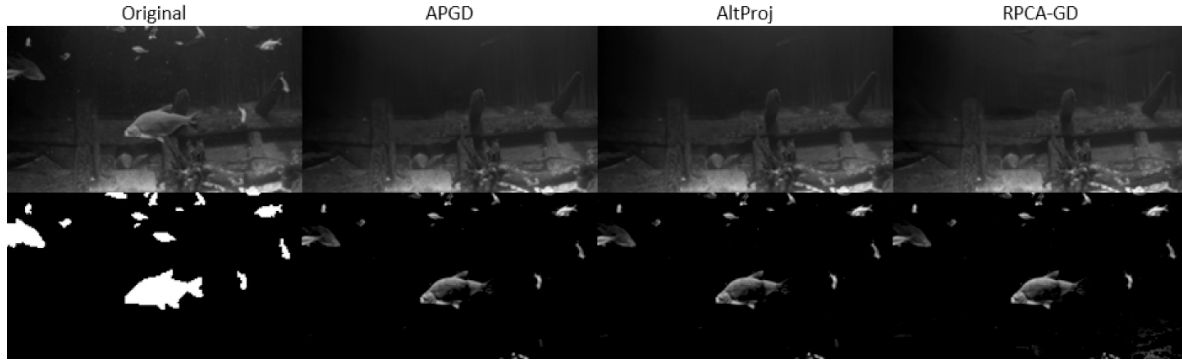


Figure 1. Frames from the Marine Snow video. The first row of the first column shows the original frame, with the image segmentation for that frame below it. The remaining three columns show the background and foreground obtained by three different methods.

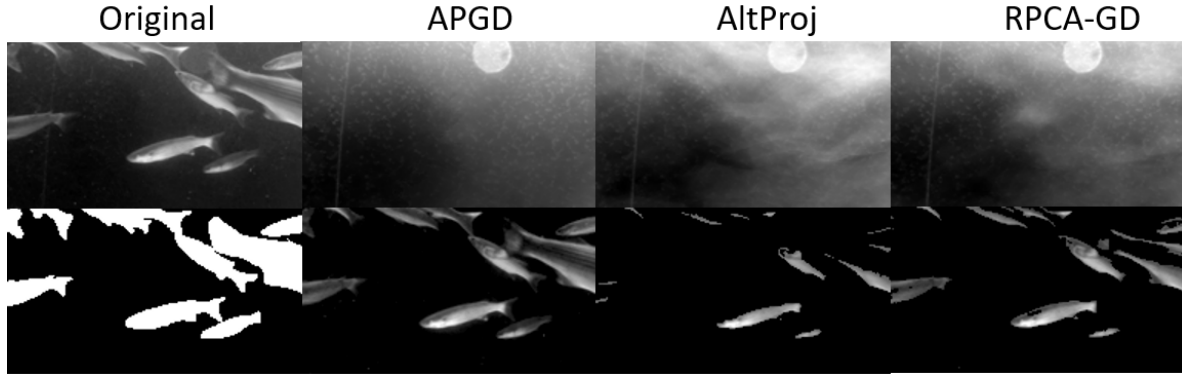


Figure 2. Frames from the Small Aquaculture video.

repetitive), then we can expect the background of each frame to be represented as a low rank matrix, and the foreground as a sparse matrix. See [5] for further details.

We evaluated our algorithm on the Underwater Change Detection dataset [27]. Out of the 1100 frames in each video, the ground truth image segmentation is included for last 100 frames. This allows use to present an objective and accurate metric of how well each method is able to identify the foreground of the image.

In Table 3, we present the runtime and the accuracy of determining which pixels contain a fish for our method compared to the four previously mentioned approaches. To calculate the accuracy, we average of the true positive ratio and true negative ratio. In three of the five videos, our method achieves the highest accuracy, where as in the other two the RPCA-GD algorithm preforms slightly better. The recovered background and foreground for our method, RPCA-GD and AltProj are shown in Figures 1, 2, 3, 4, 5, along with the original frame and segmented image.

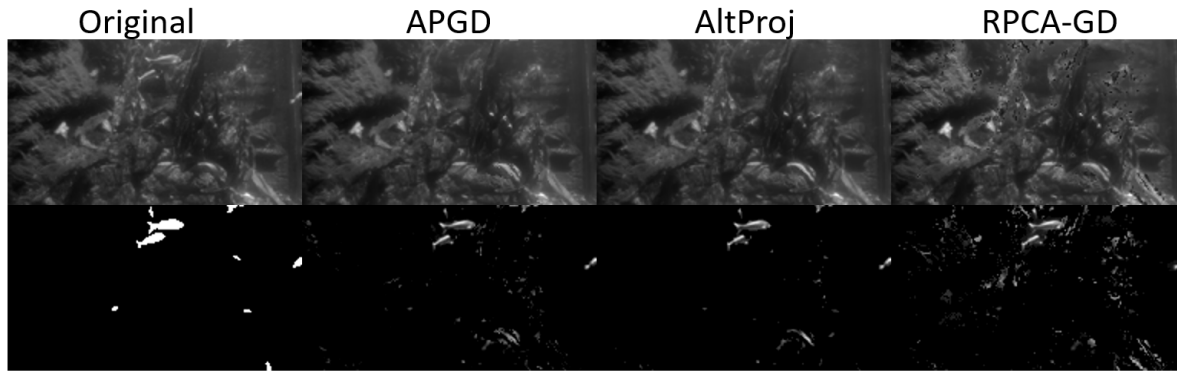


Figure 3. Frames from the Small Aquaculture video.

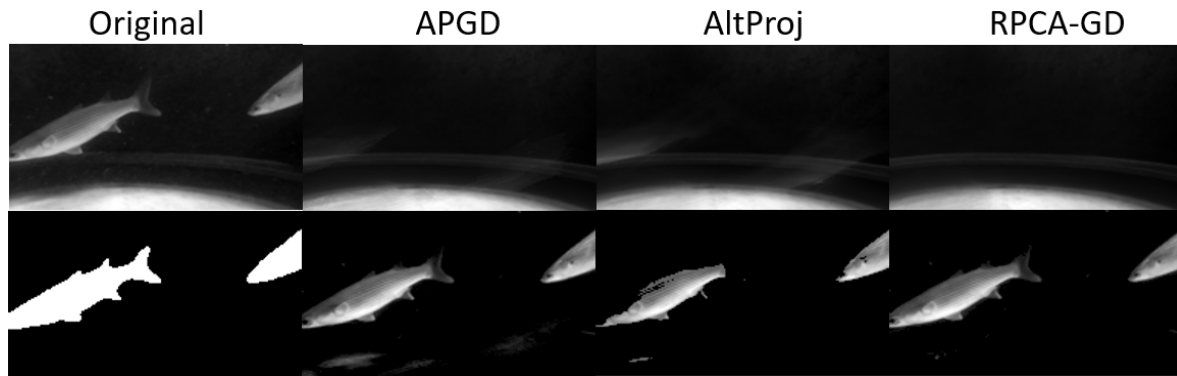


Figure 4. Frames from the Small Aquaculture video.

6. Conclusions. We have shown a novel convergence analysis of the alternating proximal gradient descent algorithm applied to the problems of matrix completion and RPCA with nonconvex regularizers, and bound the difference from the ground truth low rank matrix and sparse vector. Future work on the topic could include extending our analysis to data that lies on more complicated, nonlinear manifolds.

REFERENCES

- [1] *MovieLens*. <https://grouplens.org/datasets/movielens/>. Accessed: 2019-11-21.
- [2] A. AGARWAL, S. NEGAHBAN, AND M. J. WAINWRIGHT, *Noisy matrix decomposition via convex relaxation: Optimal rates in high dimensions*, Ann. Statist., 40 (2012), pp. 1171–1197, <https://doi.org/10.1214/12-AOS1000>, <https://doi.org/10.1214/12-AOS1000>.
- [3] J.-F. CAI, E. J. CANDÈS, AND Z. SHEN, *A singular value thresholding algorithm for matrix completion*, 2008, <https://arxiv.org/abs/0810.3286>.
- [4] E. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Commun. ACM, 55 (2012), p. 111–119, <https://doi.org/10.1145/2184319.2184343>, <https://doi.org/10.1145/2184319.2184343>.
- [5] E. J. CANDÈS, X. LI, Y. MA, AND J. WRIGHT, *Robust principal component analysis?*, J. ACM, 58 (2011),

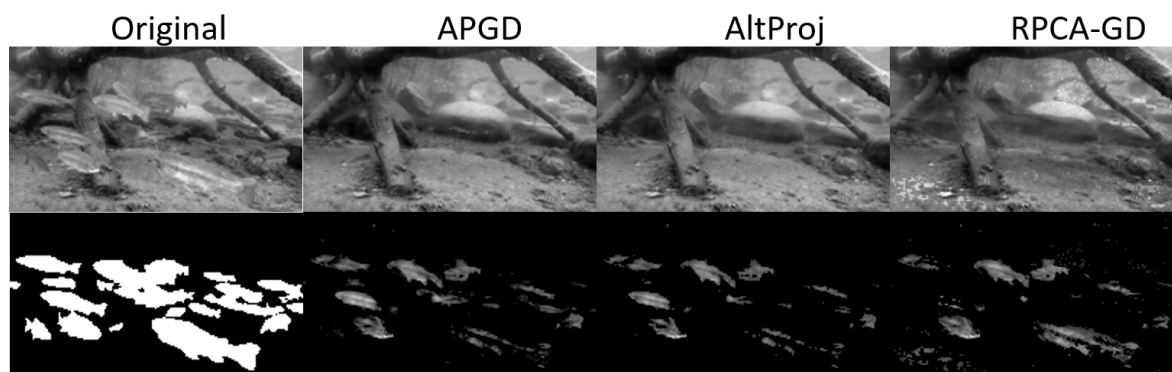


Figure 5. Frames from the Small Aquaculture video.

- 498 <https://doi.org/10.1145/1970392.1970395>, <https://doi.org/10.1145/1970392.1970395>.
- 499 [6] E. J. CANDLES AND Y. PLAN, *Matrix completion with noise*, Proceedings of the IEEE, 98 (2010), pp. 925–
- 500 936, <https://doi.org/10.1109/JPROC.2009.2035722>.
- 501 [7] E. J. CANDLES AND T. TAO, *The power of convex relaxation: Near-optimal matrix completion*, IEEE
- 502 Transactions on Information Theory, 56 (2010), pp. 2053–2080, [https://doi.org/10.1109/TIT.2010.](https://doi.org/10.1109/TIT.2010.2044061)
- 503 2044061.
- 504 [8] E. J. CANDLES, M. B. WAKIN, AND S. P. BOYD, *Enhancing sparsity by reweighted l_1 minimization*,
- 505 Journal of Fourier analysis and applications, 14 (2008), pp. 877–905.
- 506 [9] C. X. S. Y. Z. L. CANYI LU, CHANGBO ZHU, *Generalized singular value thresholding*, Proceedings of the
- 507 Twenty-Ninth AAAI Conference on Artificial Intelligence, (2015), pp. 1805–1811.
- 508 [10] V. CHANDRASEKARAN, S. SANGHAVI, P. A. PARRILO, AND A. S. WILLSKY, *Rank-sparsity incoher-*
- 509 *ence for matrix decomposition*, SIAM Journal on Optimization, 21 (2011), pp. 572–596, [https://doi.](https://doi.org/10.1137/090761793)
- 510 [org/10.1137/090761793](https://doi.org/10.1137/090761793), <https://doi.org/10.1137/090761793>, [https://arxiv.org/abs/https://doi.org/](https://arxiv.org/abs/https://doi.org/10.1137/090761793)
- 511 [10.1137/090761793](https://doi.org/10.1137/090761793).
- 512 [11] R. CHARTRAND, *Nonconvex splitting for regularized low-rank + sparse decomposition*, IEEE Transactions
- 513 on Signal Processing, 60 (2012), pp. 5810–5819, <https://doi.org/10.1109/TSP.2012.2208955>.
- 514 [12] R. GE, C. JIN, AND Y. ZHENG, *No spurious local minima in nonconvex low rank problems: A unified*
- 515 *geometric analysis*, vol. 70 of Proceedings of Machine Learning Research, International Convention
- 516 Centre, Sydney, Australia, 06–11 Aug 2017, PMLR, pp. 1233–1242, [http://proceedings.mlr.press/](http://proceedings.mlr.press/v70/ge17a.html)
- 517 [v70/ge17a.html](http://proceedings.mlr.press/v70/ge17a.html).
- 518 [13] R. GE, J. D. LEE, AND T. MA, *Matrix completion has no spurious local minimum*, in Advances in Neural
- 519 Information Processing Systems, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, eds.,
- 520 vol. 29, Curran Associates, Inc., 2016, pp. 2973–2981, [https://proceedings.neurips.cc/paper/2016/](https://proceedings.neurips.cc/paper/2016/file/7fb8ceb3bd59c7956b1df66729296a4c-Paper.pdf)
- 521 [file/7fb8ceb3bd59c7956b1df66729296a4c-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/7fb8ceb3bd59c7956b1df66729296a4c-Paper.pdf).
- 522 [14] K. GOLDBERG, T. ROEDER, D. GUPTA, AND C. PERKINS, *Eigentaste: A constant time collaborative*
- 523 *filtering algorithm*, Inf. Retr., 4 (2001), p. 133–151, <https://doi.org/10.1023/A:1011419012209>, <https://doi.org/10.1023/A:1011419012209>.
- 524 <https://doi.org/10.1023/A:1011419012209>.
- 525 [15] T. A. HORN AND C. R. JOHNSON, *Topics in Matrix Analysis*, Cambridge, 1991.
- 526 [16] P. JAIN, R. MEKA, AND I. DHILLON, *Guaranteed rank minimization via singular value projection*, in Ad-
- 527 *vances in Neural Information Processing Systems*, J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel,
- 528 and A. Culotta, eds., vol. 23, Curran Associates, Inc., 2010, pp. 937–945, [https://proceedings.neurips.](https://proceedings.neurips.cc/paper/2010/file/08d98638c6fcd194a4b1e6992063e944-Paper.pdf)
- 529 [cc/paper/2010/file/08d98638c6fcd194a4b1e6992063e944-Paper.pdf](https://proceedings.neurips.cc/paper/2010/file/08d98638c6fcd194a4b1e6992063e944-Paper.pdf).
- 530 [17] P. JAIN, P. NETRAPALLI, AND S. SANGHAVI, *Low-rank matrix completion using alternating minimization*,
- 531 in Proceedings of the Forty-Fifth Annual ACM Symposium on Theory of Computing, STOC '13, New
- 532 York, NY, USA, 2013, Association for Computing Machinery, p. 665–674, [https://doi.org/10.1145/](https://doi.org/10.1145/2488608.2488693)
- 533 [2488608.2488693](https://doi.org/10.1145/2488608.2488693), <https://doi.org/10.1145/2488608.2488693>.

- [18] Z. KANG, C. PENG, AND Q. CHENG, *Robust pca via nonconvex rank approximation*, in 2015 IEEE International Conference on Data Mining, 2015, pp. 211–220, <https://doi.org/10.1109/ICDM.2015.15>.
- [19] Z. LIN, M. CHEN, AND Y. MA, *The augmented lagrange multiplier method for exact recovery of corrupted low-rank matrices*, arXiv preprint arXiv:1009.5055, (2010).
- [20] P.-L. LOH AND M. J. WAINWRIGHT, *Regularized m -estimators with nonconvexity: Statistical and algorithmic theory for local optima*, Journal of Machine Learning Research, 16 (2015), pp. 559–616, <http://jmlr.org/papers/v16/loh15a.html>.
- [21] P.-L. LOH AND M. J. WAINWRIGHT, *Support recovery without incoherence: A case for nonconvex regularization*, Ann. Statist., 45 (2017), pp. 2455–2482, <https://doi.org/10.1214/16-AOS1530>.
- [22] C. LU, J. TANG, S. YAN, AND Z. LIN, *Nonconvex nonsmooth low rank minimization via iteratively reweighted nuclear norm*, IEEE Transactions on Image Processing, 25 (2016), pp. 829–839, <https://doi.org/10.1109/TIP.2015.2511584>.
- [23] K. MOHAN AND M. FAZEL, *Iterative reweighted algorithms for matrix rank minimization*, J. Mach. Learn. Res., 13 (2012), p. 3441–3473.
- [24] S. NEGAHBAN AND M. J. WAINWRIGHT, *Estimation of (near) low-rank matrices with noise and high-dimensional scaling*, The Annals of Statistics, (2011), pp. 1069–1097.
- [25] S. NEGAHBAN AND M. J. WAINWRIGHT, *Restricted strong convexity and weighted matrix completion: Optimal bounds with noise*, The Journal of Machine Learning Research, 13 (2012), pp. 1665–1697.
- [26] P. NETRAPALLI, U. NIRANJAN, S. SANGHAVI, A. ANANDKUMAR, AND P. JAIN, *Non-convex robust pca*, in Advances in Neural Information Processing Systems, 2014, pp. 1107–1115.
- [27] M. RADOLKO, F. FARHADIFARD, AND U. F. VON LUKAS, *Dataset on underwater change detection*, in OCEANS 2016 MTS/IEEE Monterey, 2016, pp. 1–8, <https://doi.org/10.1109/OCEANS.2016.7761129>.
- [28] B. RECHT, *A simpler approach to matrix completion*, J. Mach. Learn. Res., 12 (2011), p. 3413–3430.
- [29] A. SAGAN AND J. E. MITCHELL, *Low-rank factorization for rank minimization with nonconvex regularizers*, 2020, <https://arxiv.org/abs/2006.07702>.
- [30] Y. SHEN, Z. WEN, AND Y. ZHANG, *Augmented lagrangian alternating direction method for matrix separation based on low-rank factorization*, Optimization Methods and Software, 29 (2014), pp. 239–263.
- [31] A. SOBRAL, T. BOUWMANS, AND E.-H. ZAHZAH, *Lrslibrary: Low-rank and sparse tools for background modeling and subtraction in videos*, in Robust Low-Rank and Sparse Matrix Decomposition: Applications in Image and Video Processing, CRC Press, Taylor and Francis Group.
- [32] B. VANDEREYCKEN, *Low-rank matrix completion by riemannian optimization*, SIAM Journal on Optimization, 23 (2013), pp. 1214–1236.
- [33] Z. WANG, H. LIU, AND T. ZHANG, *Optimal computational and statistical rates of convergence for sparse nonconvex learning problems*, Ann. Statist., 42 (2014), pp. 2164–2201, <https://doi.org/10.1214/14-AOS1238>, <https://doi.org/10.1214/14-AOS1238>.
- [34] Z. WEN, W. YIN, AND Y. ZHANG, *Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm*, Mathematical Programming Computation, 4 (2012), pp. 333–361.
- [35] Q. YAO, J. T.-Y. KWOK, AND B. HAN, *Efficient nonconvex regularized tensor completion with structure-aware proximal iterations*, vol. 97 of Proceedings of Machine Learning Research, Long Beach, California, USA, 09–15 Jun 2019, PMLR, pp. 7035–7044, <http://proceedings.mlr.press/v97/yao19a.html>.
- [36] X. YI, D. PARK, Y. CHEN, AND C. CARAMANIS, *Fast algorithms for robust pca via gradient descent*, in Advances in Neural Information Processing Systems, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, eds., vol. 29, Curran Associates, Inc., 2016, pp. 4152–4160, <https://proceedings.neurips.cc/paper/2016/file/b5f1e8fb36cd7fbeb798e8639ac79e9-Paper.pdf>.
- [37] X. ZHANG, L. WANG, AND Q. GU, *A unified framework for nonconvex low-rank plus sparse matrix recovery*, vol. 84 of Proceedings of Machine Learning Research, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018, PMLR, pp. 1097–1107, <http://proceedings.mlr.press/v84/zhang18c.html>.