

Bilinear Parameterization For Differentiable Rank-Regularization

Marcus Valtonen Örnhag¹ Carl Olsson^{1,2} Anders Heyden¹

¹Centre for Mathematical Sciences
Lund University

²Department of Electrical Engineering
Chalmers University of Technology

{marcusv, calle, heyden}@math.lth.se

Abstract

Low rank approximation is a commonly occurring problem in many computer vision and machine learning applications. There are two common ways of optimizing the resulting models. Either the set of matrices with a given rank can be explicitly parametrized using a bilinear factorization, or low rank can be implicitly enforced using regularization terms penalizing non-zero singular values. While the former approach results in differentiable problems that can be efficiently optimized using local quadratic approximation, the latter is typically not differentiable (sometimes even discontinuous) and requires first order subgradient or splitting methods. It is well known that gradient based methods exhibit slow convergence for ill-conditioned problems.

In this paper we show how many non-differentiable regularization methods can be reformulated into smooth objectives using bilinear parameterization. This allows us to use standard second order methods, such as Levenberg–Marquardt (LM) and Variable Projection (VarPro), to achieve accurate solutions for ill-conditioned cases. We show on several real and synthetic experiments that our second order formulation converges to substantially more accurate solutions than competing state-of-the-art methods.

1. Introduction

Low rank models have been applied to numerous vision applications ranging from high level shape and deformation to pixel appearance models [49, 6, 53, 23, 2, 22, 51, 11]. When the sought rank is known, a commonly occurring formulation is the least squares minimization

$$\min_{\text{rank}(X) \leq r} \|\mathcal{A}X - b\|^2, \quad (1)$$

where $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^p$ is a linear operator, and $\|\cdot\|$ is the standard Euclidean vector norm. In general, this is a difficult non-convex problem and some versions are even known to be NP-hard [27]. In structure from motion, a pop-

ular approach [7] is to optimize over a bilinear factorization $X = BC^T$, where B is $m \times r$ and C is $n \times r$, and solve

$$\min_{B, C} \|\mathcal{A}BC^T - b\|^2. \quad (2)$$

Since the rank is bounded by the number of columns in B and C this approach explicitly parametrizes the set of matrices of rank r . While bilinear approaches often perform well [31, 16] they can have local minima [7]. Recent works [31, 32, 33, 35] have, however, shown that properly implemented, LM and VarPro approaches are remarkably robust to local minima, achieve quadratic convergence and give impressive reconstruction results. Recently [25, 3, 24] was able to give conditions which guarantee that there are no "spurious" local minimizers (meaning that all local minimizers are close to or identical to the global solution). They use the notion of restricted isometry property (RIP) [46] which assumes that the operator \mathcal{A} fulfills

$$(1 - \delta_r)\|X\|_F^2 \leq \|\mathcal{A}X\|^2 \leq (1 + \delta_r)\|X\|_F^2, \quad (3)$$

with $0 \leq \delta_r < 1$, if $\text{rank}(X) \leq r$. If the isometry constant δ_r is sufficiently small [25, 24, 3] prove that every local minimizer is optimal (or near optimal). Similarly, for the matrix completion problem [24] showed that there are no spurious local minima under uniformly distributed missing data. While the above theoretical assumptions generally do not hold for computer vision problems such as structure from motion, these results still give some intuition as to why bilinear parameterization often works well.

An alternative approach is to optimize directly over the entries of X and enforce low rank using regularization terms. Applying a robust function f to the singular values $\sigma_i(X) = 1, \dots, N = \min(m, n)$ results in a low-rank inducing objective

$$\min_X \mathcal{R}(X) + \|\mathcal{A}X - b\|^2, \quad (4)$$

where $\mathcal{R}(X) = \sum_{i=1}^N f(\sigma_i(X))$. Besides controlling the rank of the solution the generality of the function f offers

increased modeling capability compared to (1) and can for example be used to add priors on the size of the non-zero singular values.

The most popular regularization approach is undoubtedly the nuclear norm, $f(\sigma_i(X)) = \sigma_i(X)$, due to its convexity [18, 46, 45, 9, 10]. Under the RIP assumption exact or approximate recovery with the nuclear norm can then be guaranteed [46, 10]. On the other hand, since it penalizes large singular values, it suffers from a shrinking bias [8, 11, 36]. Ideally f should penalize small singular values (assumed to stem from measurement noise) harder than the large ones. Therefore non-increasing derivatives on $[0, \infty)$, or concavity, has been shown to give stronger relaxations [44, 39, 34, 41, 12, 48, 28]. These non-convex formulations usually only come with local convergence guarantees. Two exceptions are [36, 42] which gave optimality guarantees for (4) with $f = f_\mu$ as in (8).

The regularization term is generally not differentiable as a function of X . Thus, optimization methods based on local quadratic approximation become infeasible. Figure 1 gives a simple illustration on a 1-dimensional example of how non-differentiability occurs at the origin. In addition it is well known that the singular values become non-differentiable functions of the matrix elements when they are non distinct. To circumvent these issues subgradient and splitting methods are often employed [12, 48, 28, 40, 36]. It is well known from basic optimization theory (e.g. [5]) that gradient based methods exhibit slow convergence for ill-conditioned problems. It has also been observed (e.g. [4]) that splitting methods rapidly reduce the objective value the first couple of iterations, while convergence to the exact solution can be slow. In this paper we show that there are computer vision problems where these approaches make very little improvements at all, returning a solution that is far from optimal. In contrast, bilinear formulations with either LM or VarPro can be made to yield accurate results in few iterations [31].

An alternative approach that unifies bilinear parameterization with regularization approaches is based on the observation [46] that the nuclear norm $\|X\|_*$ of a matrix X can be expressed as $\|X\|_* = \min_{BC^T=X} \frac{\|B\|_F^2 + \|C\|_F^2}{2}$. Thus when $f(\sigma_i(X)) = \mu\sigma_i(X)$, where μ is a scalar controlling the strength of the regularization, optimization of (4) can be formulated as

$$\min_{B,C} \mu \frac{\|B\|_F^2 + \|C\|_F^2}{2} + \|ABC^T - b\|^2. \quad (5)$$

Optimizing directly over the factors has the advantages that the number of variables is much smaller and one may add constraints if a particular factorization is sought. Surprisingly, while (5) is non-convex, using the convexity of the underlying regularization problem (4) it can be shown that any local minimizer B, C with $\text{rank}(BC^T) < k$, where k

is the number of columns in B and C , is globally optimal [1, 29]. Additionally, the objective function is two times differentiable and second order methods can be employed.

In this paper we develop new regularizing terms that, similar to (5), work on the bilinear factors. However, in contrast to previous approaches we investigate formulations that exhibit less shrinking bias and go beyond convex penalties. Specifically, we prove that $\mathcal{R}(X) = \min_{X=BC^T} \tilde{\mathcal{R}}(B, C)$, where

$$\tilde{\mathcal{R}}(B, C) = \sum_{i=1}^k f\left(\frac{\|B_i\|^2 + \|C_i\|^2}{2}\right), \quad (6)$$

k is the number of columns, and B_i and C_i are the i :th columns of B and C , respectively. The result holds for a general class of concave penalty functions f , a few of which are illustrated in Figure 1. In view of the above result, we propose to minimize

$$\tilde{\mathcal{R}}(B, C) + \|ABC^T - b\|^2. \quad (7)$$

Rather than resorting to splitting or subgradient methods we present an algorithm that uses a quadratic approximation of the objective. Under the assumption that f is differentiable, we show that our quadratic approximation reduces to a weighted version of (5) to which we can apply VarPro. We show on several computer vision problems that our approach outperforms state-of-the-art methods such as [47, 12, 48, 28, 4].

While our problem is non-convex (both in the X parameterization (4) and in the B, C parameterization (7)) we show that in some cases it is still possible to give global optimality guarantees. Building on the results of [42] we characterize the local minima of the new formulation with the choice

$$f(x) = f_\mu(x) := \mu - \max(\sqrt{\mu} - x, 0)^2. \quad (8)$$

Specifically, for this choice, we give conditions that ensure that when a RIP constraint [46] holds a local minimizer of (7) is a global solution of both

$$\min_{\text{rank}(X) \leq r} \mathcal{R}(X) + \|AX - b\|^2, \quad (9)$$

where $\mathcal{R}(X) = \sum_i f_\mu(\sigma_i(X))$, and

$$\min_{\text{rank}(X) \leq r} \mu \text{rank}(X) + \|AX - b\|^2. \quad (10)$$

In summary our main contributions are:

- A new stronger non-convex regularization term for bilinear parameterizations with less/no shrinking bias.
- A new iteratively reweighed VarPro algorithm optimizing accurate quadratic approximations.

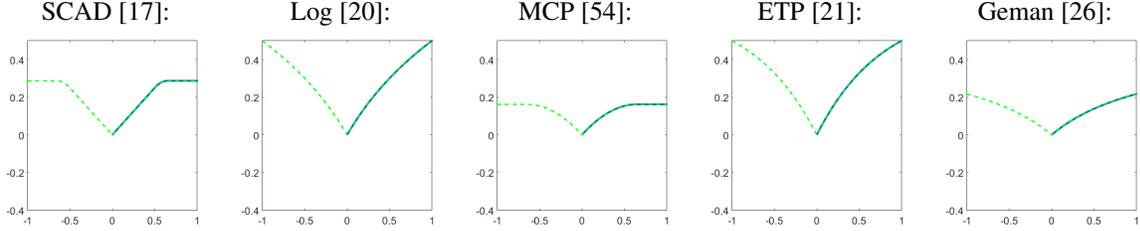


Figure 1. A few commonly occurring robust penalties of the form $f(\sigma)$, with $\sigma \in [0, \infty)$ and f differentiable everywhere (blue graph). The green dashed graph shows how non-differentiability occurs at the origin when applying the penalty to a 1×1 matrix $x \in \mathbb{R}$. In this case $\sigma(x) = |x|$ and therefore $f(\sigma(x)) = f(|x|)$. Note also that (8) is a special case of MCP.

- Theoretical conditions that guarantee optimal recovery under the RIP constraint.
- An experimental evaluation that shows that our methods outperforms state-of-the-art methods on several real computer vision problems.

1.1. Related Work

Our work is very much inspired by a recent series of papers by Hong *et al.* [31, 32, 33, 35] which show that bilinear formulations can be made remarkably robust to local minima, and achieve impressive reconstruction results for uncalibrated structure from motion problems, using the so called VarPro method. Our work represents an attempt to unify this line of work with regularization based alternatives, leveraging the benefits of them both.

An approach that is closely related to ours is that of [8] which uses (5) to unify the use of a regularized objective and factorization. They show that if the obtained solution has lower rank than its number of columns it is globally optimal. In practice [8] observes that the shrinking bias of the nuclear norm makes it too weak to enforce a low rank when the data is noisy. Therefore, a “continuation” approach where the size of the factorization is gradually reduced is proposed. While this yields solutions with lower rank, the optimality guarantees no longer apply.

Bach *et al.* [1] showed that

$$\|X\|_{s,t} := \min_{X=BC^T} \sum_{i=1}^k \frac{\|B_i\|_s^2 + \|C_i\|_t^2}{2}, \quad (11)$$

is convex for any choice of vector norms $\|\cdot\|_s$ and $\|\cdot\|_t$. In [29] it was shown that a more general class of 2-homogeneous factor penalties result in a convex regularization similar to (11). The property that a local minimizer B, C with $\text{rank}(BC^T) < k$, is also extended to this case. Still, because of convexity, it is clear that these formulations will suffer from a similar shrinking bias as the nuclear norm. Shang *et al.* [47] showed that penalization with the Schatten semi-norms $\|X\|_q = \sqrt[q]{\sum_{i=1}^N \sigma_i(X)^q}$, for $q = 1/2$ and $2/3$, can be achieved using a convex penalty on the factors

B and C . A generalization to general values of q is given in [52]. While this reduces shrinking bias to some extent, it results in a non-differentiable and non-convex formulation that is optimized with ADMM.

It is important to note that many of the above methods that are considered state-of-the-art have been developed for low level vision tasks such as image denoising, inpainting, alignment and background subtraction. The ground truth for these models are often of higher rank than models in *e.g.* structure from motion, making it possible to obtain good results with weaker regularization. Additionally, as we will see in the experiments, more difficult data terms prevent rapid convergence of the splitting methods they often employ.

2. Non-Convex Penalties and Shrinking Bias

In this section we will show how to formulate regularization terms of the type

$$\mathcal{R}(X) = \sum_{i=1}^N f(\sigma_i(X)), \quad (12)$$

by penalizing the factors of the factorization $X = BC^T$. We assume that B and C have k columns, making $\sigma_i(X) = 0$ if $i > k$ and $\text{rank}(X) \leq k$. Note, however, that we are aiming to achieve a lower rank using the regularization term. In many applications, the sought rank is unknown and should be determined by the regularization. We therefore set k large enough not to exclude the optimal solution. As we shall see in Section 3, this ability to over-parameterize can be used to ensure optimality.

Theorem 1. *If f is concave, non-decreasing on $[0, \infty)$ and $f(0) = 0$ then*

$$\mathcal{R}(X) = \min_{BC^T=X} \sum_{i=1}^k f(\|B_i\| \|C_i\|), \quad (13)$$

where B_i and C_i , $i = 1, \dots, k$ are the columns of B and C respectively.

Proof. The result is a consequence of the fact that \mathcal{R} will fulfill a triangle inequality $\mathcal{R}(X + Y) \leq \mathcal{R}(X) + \mathcal{R}(Y)$ under the assumptions on f . This is clear from Theorem 4.4 in [50] which shows that

$$\sum_{i=1}^N f(\sigma_i(X + Y)) \leq \sum_{i=1}^N (f(\sigma_i(X)) + f(\sigma_i(Y))). \quad (14)$$

Applying this to $X = BC^T = \sum_{i=1}^k B_i C_i^T$ we see that

$$\mathcal{R}(X) = \mathcal{R}\left(\sum_{i=1}^k B_i C_i^T\right) \leq \sum_{i=1}^k \mathcal{R}(B_i C_i^T). \quad (15)$$

Since $\text{rank}(B_i C_i^T) = 1$ we also have

$$\mathcal{R}(B_i C_i^T) = f(\sigma_1(B_i C_i^T)) = f(\|B_i C_i^T\|_F). \quad (16)$$

Lastly, since $\|B_i C_i^T\|_F = \|B_i\| \|C_i\|$ we get

$$\mathcal{R}(X) \leq \sum_{i=1}^k f(\|B_i\| \|C_i\|). \quad (17)$$

To see that equality can be achieved, let $B_i = \sqrt{\sigma_i(X)} U_i$ and $C_i = \sqrt{\sigma_i(X)} V_i$, where $X = \sum_{i=1}^k \sigma_i(X) U_i V_i^T$ is the SVD of X . Then, $BC^T = X$ and $f(\|B_i\| \|C_i\|) = f(\sigma_i(X))$. \square

While the above result allows optimization over the factors B and C we note that it yields an objective that is non-differentiable at $\|B_i\| \|C_i\| = 0$. Next we reformulate the objective to achieve a differentiable problem formulation.

Corollary 1. *Under the assumptions of Theorem 1, it follows that $\mathcal{R}(X) = \min_{X=BC^T} \tilde{\mathcal{R}}(B, C)$, where*

$$\tilde{\mathcal{R}}(B, C) = \sum_{i=1}^k f\left(\frac{\|B_i\|^2 + \|C_i\|^2}{2}\right). \quad (18)$$

If f is differentiable then $\tilde{\mathcal{R}}(B, C)$ is also differentiable.

Proof. By the rule of arithmetic and geometric means

$$\|B_i\| \|C_i\| \leq \frac{1}{2} (\|B_i\|^2 + \|C_i\|^2), \quad (19)$$

with equality if $\|B_i\| = \|C_i\|$ which is achieved when $B_i = \sqrt{\sigma_i(X)} U_i$ and $C_i = \sqrt{\sigma_i(X)} V_i$. Since f is assumed to be non-decreasing, it follows from (13), that $\mathcal{R}(X) = \min_{X=BC^T} \tilde{\mathcal{R}}(B, C)$. The differentiability of $\tilde{\mathcal{R}}(B, C)$ is now trivially checked using the chain rule. \square

We are particularly interested in the case (8) since, with this choice, it is known that the global minimizer of (4) is the same as that of $\mu \text{rank}(X) + \|\mathcal{A}X - b\|^2$ if $\|\mathcal{A}\| < 1$, see [13] for a proof. Note that f_μ is a special case of the MCP

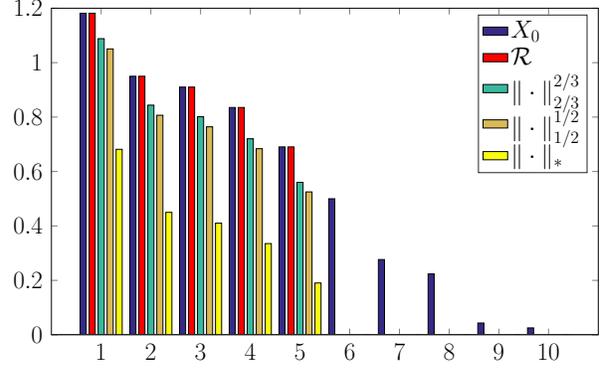


Figure 2. Singular values obtained when minimizing $\|X - X_0\|_F^2$ with the four regularizers $\mathcal{R}(X)$ with $f = f_\mu$, $\|X\|_{1/2}^{1/2}$, $\|X\|_{2/3}^{2/3}$ and $\|X\|_*$. Large singular values are left unchanged by \mathcal{R} .

class [54]. With this choice $\tilde{\mathcal{R}}(B, C)$ is differentiable and the second derivatives are also defined almost everywhere except in the transition $\frac{\|B_i\|^2 + \|C_i\|^2}{2} = \sqrt{\mu}$ where the function switches from quadratic to constant.

We conclude this section by comparing the shrinking bias of our approach and three others that can also be optimized over the factorization. Theorem 1 makes it possible to compute the global optimizer of $\tilde{\mathcal{R}}(B, C) + \|BC^T - X_0\|_F^2$ since the equivalent problem $\mathcal{R}(X) + \|X - X_0\|_F^2$ has closed form solution in the X -parameterization. It is shown in [36] that with $f = f_\mu$ the solution is obtained by thresholding the singular values at $\sqrt{\mu}$. Similarly, closed form solutions are also available when regularizing $\|X - X_0\|_F^2$ with $\|\cdot\|_{1/2}$, $\|\cdot\|_{2/3}$ and $\|\cdot\|_*$ [47]. In Figure 2 we show the singular values obtained when regularizing $\|X - X_0\|_F^2$ with these four options, and for comparison the singular values of X_0 . For all methods we have selected regularization weights as small as possible so that the five smallest singular values are completely suppressed, which minimizes the bias. While all choices, except \mathcal{R} , subtract a part from the singular values that should be retained, the Schatten norms reduce the bias significantly compared to the nuclear norm. For the Schatten norms the bias is larger for singular values that are close to the threshold since the derivative of σ^q , $0 < q < 1$, decreases with increasing σ . For problem instances where there is a clear separation in size between singular values that should be retained and those that should be suppressed, it is likely that this can be done with negligible bias. Since $f'_\mu(\sigma) = 0$ when $\sigma \geq \sqrt{\mu}$ this method does not affect the first five singular values.

3. Overparameterization and Optimality

The results of the previous section show that a global optimizer (B, C) of (7) gives a solution BC^T which is globally optimal in (4). On the other hand, optimizing (7) over B and C introduces additional stationary points, due

Table 1. Distance to ground truth (normalized) mean valued over 20 problem instances for different percentages of missing data, missing data patterns and noise levels σ . Best results are marked in bold.

Missing data (%)		PCP [9]	WNNM [28]	Unifying [8]	LpSq [40]	S12L12 [47]	S23L23 [47]	IRNN [12]	APGL [48]	$\ \cdot\ _*$ [4]	\mathcal{R} [36]	Our
Uniform ($\sigma = 0.0$)	0	0.0000	0.0000	0.0000	0.0000	0.0002	0.0002	0.0000	0.0000	0.1727	0.0000	0.0000
	10	0.0885	0.0028	0.0713	0.0213	0.0309	0.0071	0.0000	0.0000	0.1998	0.0000	0.0000
	20	0.2720	0.2220	0.1491	0.0170	0.0412	0.0209	0.0000	0.0000	0.2223	0.0128	0.0000
	30	0.7404	0.4787	0.7499	0.0003	0.0818	0.0895	0.0000	0.0014	0.2897	0.2346	0.0000
	40	1.0000	0.6097	0.9553	0.1083	0.1666	0.1360	0.0000	0.0017	0.3374	0.2198	0.0000
	50	1.0000	0.7170	1.0000	0.0315	0.1376	0.1001	0.0003	0.0301	0.4266	0.2930	0.0000
Tracking ($\sigma = 0.0$)	0	0.0000	0.0000	0.0000	0.0000	0.0002	0.0002	0.0000	0.0000	0.1810	0.0000	0.0000
	10	0.3160	0.2734	0.1534	0.0839	0.1296	0.1233	0.0772	0.0834	0.2193	0.0793	0.0658
	20	0.4877	0.4499	0.3017	0.1650	0.2389	0.2456	0.1010	0.1786	0.3436	0.2494	0.1018
	30	0.5821	0.5395	0.5486	0.2520	0.3289	0.3160	0.1189	0.2572	0.4299	0.3421	0.1189
	40	0.7072	0.6317	0.7376	0.2853	0.4084	0.4110	0.1417	0.2913	0.4825	0.5004	0.1385
	50	0.8125	0.7257	0.9521	0.4178	0.4267	0.4335	0.2466	0.4047	0.5754	0.6503	0.2214
Tracking ($\sigma = 0.1$)	0	0.0409	0.0207	0.0407	0.0450	0.0437	0.0435	0.0448	0.0191	0.1581	0.0166	0.0166
	10	0.3157	0.2734	0.1585	0.0848	0.0529	0.0518	0.0625	0.0696	0.2312	0.0488	0.0438
	20	0.4771	0.4338	0.3480	0.1394	0.0995	0.0982	0.1090	0.1188	0.3109	0.2071	0.0983
	30	0.5801	0.5225	0.4726	0.2026	0.2468	0.2592	0.1646	0.1993	0.3820	0.3465	0.1475
	40	0.7122	0.6148	0.8638	0.2225	0.3292	0.3252	0.1357	0.2110	0.4800	0.4599	0.1273
	50	0.7591	0.6819	0.9216	0.4105	0.4883	0.4811	0.3342	0.3639	0.5652	0.5930	0.3329

to the non-linear parameterization, that are not present in (4). One such point is $(B, C) = (0, 0)$ where the gradients of $\|ABC^T - b\|^2$ with respect to B and C vanish (in contrast to the gradient w.r.t. X). In this section we show that by overparametrizing, in the sense that we use B and C with more columns than the rank of the solution we seek, it is still possible to use properties of (4) to show optimality in (7). We will exclusively use f_μ from (8), assume that B and C have $2k$ columns and study locally optimal solutions with $\text{rank}(BC^T) < k$. The size of B and C makes it possible to parametrize line segments between such points and utilize convexity properties, see proof of Theorem 3. The following result (which is proven in Appendix A) gives conditions that ensure that local minimality in (7) implies that (4) grows in all “low rank” directions.

Theorem 2. *Assume that $(\bar{B}, \bar{C}) \in \mathbb{R}^{m \times 2k} \times \mathbb{R}^{n \times 2k}$, where $\bar{B} = U\sqrt{\Sigma}$ and $\bar{C} = V\sqrt{\Sigma}$, and $\bar{X} = U\Sigma V^T$, is a local minimizer of (7) with $\text{rank}(\bar{X}) < k$ and let $\mathcal{N}(X) = \mathcal{R}(X) + \|AX - b\|^2$. Then $\mathcal{R}(\bar{X}) = \tilde{\mathcal{R}}(\bar{B}, \bar{C})$ and the directional derivatives $\mathcal{N}'_{\Delta X}(\bar{X})$, where $\Delta X = \bar{X} - \bar{X}$ and $\text{rank}(\bar{X}) \leq k$, are non-negative.*

Note that there can be local minimizers for which $\tilde{\mathcal{R}}(\bar{B}, \bar{C}) > \mathcal{R}(\bar{B}\bar{C}^T)$ since $\tilde{\mathcal{R}}$ is non-convex. From an algorithmic point of view we can, however, escape such points by taking the current iterate and recompute the factorization of $\bar{B}\bar{C}^T$ using SVD. If the SVD of $\bar{B}\bar{C}^T = \sum_{i=1}^r \sigma_i U_i V_i^T$ we update \bar{B} and \bar{C} to $\bar{B}_i = \sqrt{\sigma_i} U_i$ and $\bar{C}_i = \sqrt{\sigma_i} V_i$, which we know reduces the energy and gives $\tilde{\mathcal{R}}(\bar{B}, \bar{C}) = \mathcal{R}(\bar{B}\bar{C}^T)$.

Theorem 2 allows us to derive optimality conditions using the properties of (4). As a simple example, consider the case where $\|AX\|^2 \geq \|X\|^2$, which makes (4) convex [13],

and let B and C have $2k$ columns. Suppose that we find a local minimizer (\bar{B}, \bar{C}) fulfilling the assumptions of Theorem 2. Then the derivative along a line segment towards any other low rank matrix is non-decreasing, and therefore $\bar{B}\bar{C}^T$ is the global optimum of (4) over the set of matrices with $\text{rank} \leq k$ by convexity.

Below we give a result that goes beyond convexity and applies to the important class [46] of problems that obey the RIP constraint (3). Let \mathcal{A}^* denote the adjoint operator of \mathcal{A} , then:

Theorem 3. *Assume that (\bar{B}, \bar{C}) is a local minimizer of (7), fulfilling the assumptions of Theorem 2. If the singular values of $Z = (I - \mathcal{A}^* \mathcal{A})\bar{B}\bar{C}^T + \mathcal{A}^* b$ fulfill $\sigma_i(Z) \notin [(1 - \delta_{2k})\sqrt{\mu}, \frac{\sqrt{\mu}}{(1 - \delta_{2k})}]$ then $\bar{B}\bar{C}^T$ is the solution of (9) and (10).*

The proof builds on the results of [42] and is given in Appendix A. The assumption that the singular values of Z are not too close to the threshold $\sqrt{\mu}$ is a natural restriction which is valid when the noise level is not too large. In case of exact data, i.e. $b = AX_0$, where $\text{rank}(X_0) = r$ it is trivially fulfilled for any choice of μ such that $\sqrt{\mu} < (1 - \delta_{2k})\sigma_r(X_0)$ since we have $Z = X_0$. For additional details on Z 's dependence on noise see [14].

The above result is similar in spirit to those of [46, 29], which show that, in the convex case, having $2k$ columns and $\text{rank } 2k - 1$ is enough to ensure that a local minimizer is global. For the proof in our non-convex case we need rank at most $k - 1$. Presently, it is not clear if our assumption can be relaxed to match that of the convex case or not.

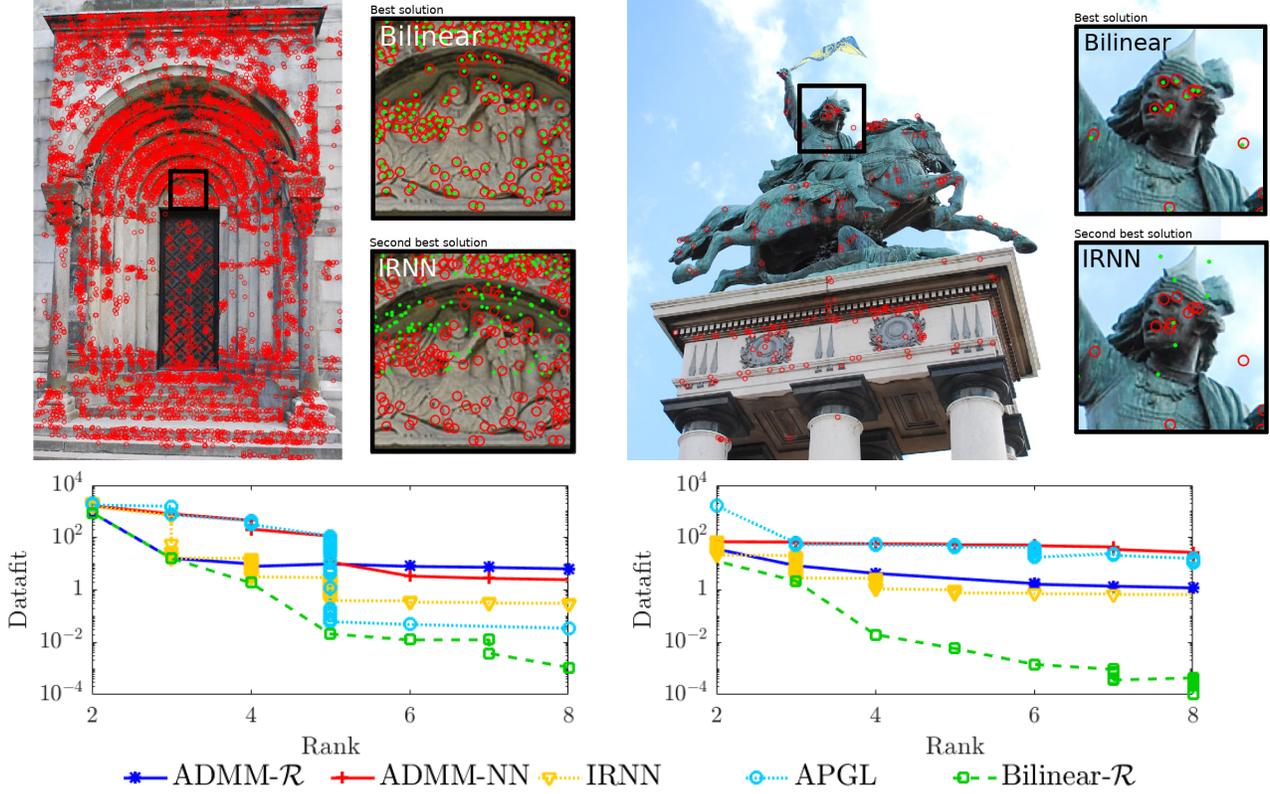


Figure 3. Comparison of reprojection error obtained using the bilinear formulation and ADMM, for datasets *Door* and *Vercingetorix* [43]. The red circles mark the feature points and the green dots the projected image points obtained from the different methods. The best rank 4 solution for the respective method was used. The control parameter $\eta = 0.5$ in both experiments.

4. An Iterative Reweighted VarPro Algorithm

In this section we give a brief overview of our algorithm for minimizing (7). A more detailed description is given in Appendix B.

Given a current iterate, $B^{(t)}$ and $C^{(t)}$, the first step of our algorithm is to replace the term $\tilde{\mathcal{R}}(B, C)$ with a quadratic function. To do this we note that by the Taylor expansion $f(x) \approx f(x_0) + f'(x_0)(x - x_0)$, minimizing $f(x)$ and $f'(x_0)x$ around x_0 is roughly the same (ignoring constants). Inserting $x_0 = \frac{\|B_i^{(t)}\|^2 + \|C_i^{(t)}\|^2}{2}$ and $x = \frac{\|B_i\|^2 + \|C_i\|^2}{2}$ now gives our approximation

$$\sum_{i=1}^k w_i^{(t)} (\|B_i\|^2 + \|C_i\|^2) + \|ABC^T - b\|^2, \quad (20)$$

where $w_i^{(t)} = \frac{1}{2} f' \left(\frac{\|B_i^{(t)}\|^2 + \|C_i^{(t)}\|^2}{2} \right)$. Here $B_i^{(t)}$ and $C_i^{(t)}$ are the i :th columns of $B^{(t)}$ and $C^{(t)}$, respectively. Minimizing (20) over C is now a least squares problem with closed form solution. Inserting this solution into the original problem gives a nonlinear problem in B alone, which is what VarPro solves. We use the so called Ruhe and Wedin (RW2) approximation with a dampening term

$\lambda \|B - B^{(t)}\|_F^2$, see [33] for details. In each step of the VarPro algorithm we update the weights $w_i^{(t)}$.

As previously mentioned, there can be stationary points for which $\tilde{\mathcal{R}}(B, C) > \mathcal{R}(BC^T)$. In each iteration we therefore take the current iterate and recompute the factorization of $B^{(t)}C^{(t)T}$ using SVD. If the SVD of $B^{(t)}C^{(t)T} = \sum_{i=1}^r \sigma_i U_i V_i^T$ we update $B^{(t)}$ and $C^{(t)}$ to $B_i^{(t)} = \sqrt{\sigma_i} U_i$ and $C_i^{(t)} = \sqrt{\sigma_i} V_i$ which we know reduces the energy and gives $\tilde{\mathcal{R}}(B^{(t)}, C^{(t)}) = \mathcal{R}(B^{(t)}C^{(t)T})$.

Our approach can be seen as iteratively reweighted nuclear norm minimization [12]; however, our bilinear formulation allows us to use quadratic approximation, thus benefiting from second order convergence in the neighborhood of a local minimum.

5. Experiments

In this section we will show the versatility and strength of the proposed method, focusing on computer vision problems. In Section 5.2 we show an example where state-of-the-art methods fail to achieve a value close to global optimality. We include two more examples of real problems, in Appendix C: background extraction and photometric stereo.

In both cases our method shows superior performance. In the main paper we focus on the trade-off between datafit and rank, but show, in the examples in the supplementary material, the added benefits of convergence speed using the proposed method. This is done by minimizing the same energy with ADMM and the proposed method, in which case the splitting schemes can be tediously slow. In all experiments our proposed method is initialized randomly, with zero mean and unit variance.

5.1. Synthetic Missing Data Problem

Let \odot denote the Hadamard product, and consider the missing data formulation

$$\min_X \mu \text{rank}(X) + \|W \odot (X - M)\|_F^2, \quad (21)$$

where M is a measurement matrix and W a missing data mask with entries $w_{ij} = 1$ if the entry is known, and zero otherwise.

In low-level vision applications such as denoising and image inpainting a uniformly random missing data pattern is often a reasonable approximation of the distribution; however, for structure from motion, the missing data pattern is often highly structured. To this end, we investigate two kinds of patterns: uniformly random and “tracking failure”. In order to construct realistic patterns of tracking failure, we use the method in [37]. This is done by randomly selecting if a track should have missing data (with uniform probability), then select (with uniform probability, starting after the first few frames) in which image tracking failure occurs. If a track is lost, it is not restarted.

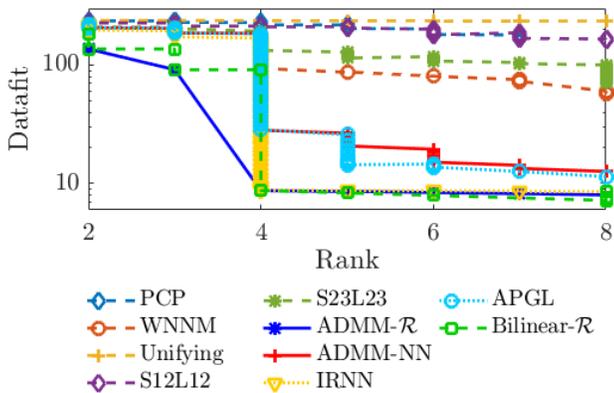


Figure 4. Rank vs datafit for the synthetic experiment in Section 5.1. No true low rank solution using LpSq [40] could be found, regardless of the choice of parameters.

We generate random ground truth matrices $M_0 \in \mathbb{R}^{32 \times 512}$ of rank 4, which can be expressed as $M_0 = UV^T$, where $U \in \mathbb{R}^{32 \times 4}$ and $V \in \mathbb{R}^{512 \times 4}$. The entries of U and V are normal distributed with zero mean

and unit variance. The measurement matrix $M = M_0 + N$, where N simulates noise and has normal distributed entries with zero mean and variance σ^2 .

Our proposed method is compared to a variety of different methods [8, 9, 28, 40, 47, 12, 48, 4, 36]. For the methods that need an initial estimate of the rank as input, the rank estimation heuristic by Shang *et al.* [47] is used. The regularization parameter is set to $\lambda = \sqrt{\max(m, n)}$, given a sought $m \times n$ matrix, as proposed by [9, 47]. In case other parameters should be provided, the one recommended from the respective authors have been used. The number of columns, for our proposed method, is set to $k = 8$, *i.e.* twice the rank of the original matrix M_0 . We exclusively use the f_μ regularization (8), and use $\sqrt{\mu} = \lambda$. Since f_μ is a special case of MCP, it is used for IRNN as well. Furthermore, we include the results for regularizing with nuclear norm [4] and f_μ (8) using ADMM, as proposed in [36]. Note that ADMM comes without optimality guarantees, however, it has been shown to work well for several computer vision problems in practice [36, 42]. Several of the compared methods solve the robust PCA problem, thus also include a sparse component, which is not taken into account.

The results are shown in Table 1. Note that most algorithms perform significantly better for the uniformly random missing data pattern, than compared to the structured missing data pattern. Our proposed method outperforms all other methods in this comparison.

Since the final rank of the estimated matrix is not necessarily the same as that of M_0 , we show the rank vs datafit obtained when varying the regularization parameter λ in Figure 4. It is evident from the results that the only candidates that yield an acceptable result for low rank solutions are ADMM with f_μ , IRNN with MCP and our proposed method.

5.2. pOSE: Pseudo Object Space Error

The Pseudo Object Space Error (pOSE) objective combines affine and projective camera models

$$\ell_{\text{OSE}} = \sum_{(i,j) \in \Omega} \|(P_{i,1:2} \tilde{\mathbf{x}}_j - (\mathbf{p}_{i,3}^T \tilde{\mathbf{x}}_j) \mathbf{m}_{i,j})\|^2, \quad (22)$$

$$\ell_{\text{Affine}} = \sum_{(i,j) \in \Omega} \|P_{i,1:2} \tilde{\mathbf{x}}_j - \mathbf{m}_{i,j}\|^2, \quad (23)$$

$$\ell_{\text{pOSE}} = (1 - \eta) \ell_{\text{OSE}} + \eta \ell_{\text{Affine}}, \quad (24)$$

where ℓ_{OSE} is the object space error and ℓ_{Affine} is the affine projection error. Here $P_{i,1:2}$ denotes the first two rows, $\mathbf{p}_{i,3}$ the third row of the i :th camera matrix, and $\tilde{\mathbf{x}}_j$ is the j :th 3D point in homogeneous coordinates. The control parameter $\eta \in [0, 1]$ determines the impact of the respective camera model. This objective was introduced in [35] to be used

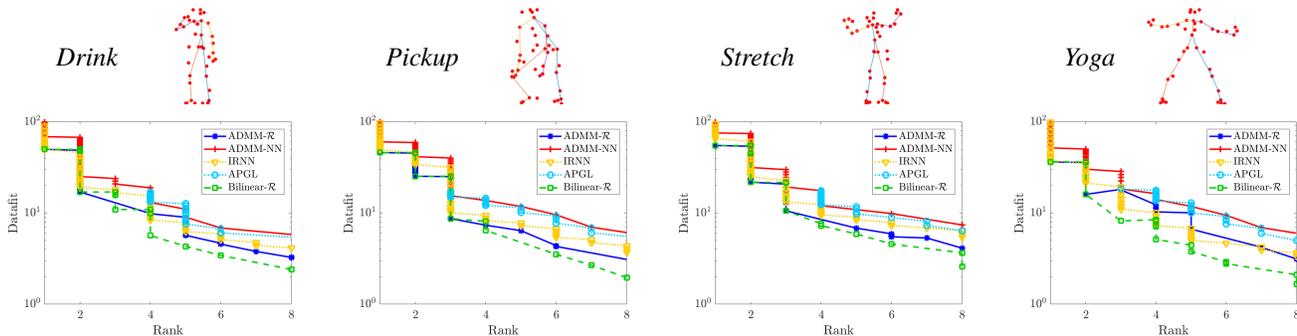


Figure 5. *Top row*: Example frames from the MOCAP dataset of the *drink*, *pickup*, *stretch* and *yoga* sequences. *Last row*: The bilinear method finds the same or a better datafit compared to the other methods for all ranks.

in a first stage of an initialization-free bundle adjustment pipeline, optimized using VarPro.

The ℓ_{pOSE} objective is linear, and acts on low-rank components P and X , which are constrained by $\text{rank}(PX^T) = 4$. Instead of enforcing the rank constraint, we replace it as before with a relaxation. By not enforcing the rank constraint we demonstrate the ability of the methods to make accurate trade-offs between minimizing the rank and fitting the data. Since the objective now becomes more complex, and is no longer compatible with the missing data formulations, only IRNN and APGL are directly applicable, as well as the ADMM approach using f_μ and nuclear norm. We use two real-life datasets with various amounts of camera locations and 3D points: *Door* with 12 images, resulting in seeking a matrix of size 36×8850 and *Vercingetorix* [43] with 69 images, resulting in seeking a matrix of size 207×1148 , both of which have rank 4.¹

As in the synthetic experiment from Section 5.1, the regularization parameter is varied and the resulting rank and datafit is stored and reported in Figure 3. To visualize the results, we considered the best rank 4 approximations, and show the reprojected points and the corresponding measured points obtained from the best method (ours in both cases) and the second best (IRNN in both cases), see Figure 3. As is readily seen by ocular inspection, the rank 4 solution obtained by our proposed method significantly outperforms those of other state-of-the-art methods.

5.3. Non-Rigid Structure From Motion

In this section we test our approach on non-rigid reconstruction (NRSfM) with the CMU Motion Capture (MOCAP) dataset. In NRSfM, the complexity of the deformations are controlled by some mild assumptions of the object shapes. Bregler *et al.* [6] suggested that the set of all possible configurations of the objects are spanned by a low dimensional linear basis of dimension K . In this setting, the non-rigid shapes $X_i \in \mathbb{R}^{3 \times n}$ can be represented

as $X_i = \sum_{k=1}^K c_{ik} B_k$, where $B_k \in \mathbb{R}^{3 \times n}$ are the basis shapes and $c_{ik} \in \mathbb{R}$ the shape coefficients. This way, the matrix X_i contains the world coordinates of point i , hence the observed image points are given by $x_i = R_i X_i$. We will assume orthographic cameras, *i.e.* $R_i \in \mathbb{R}^{2 \times 3}$ where $R_i R_i^T = I_2$. As proposed by Dai *et al.* [15], the problem can be turned into a low-rank factorization problem by reshaping and stacking the non-rigid shapes X_i . Let $X_i^\sharp \in \mathbb{R}^{1 \times 3n}$ denote the concatenation of the rows in X_i , and create $X^\sharp \in \mathbb{R}^{F \times 3n}$ by stacking X_i^\sharp . This allows us to decompose the matrix X^\sharp in the low-rank factors $X^\sharp = CB^\sharp$, where $C \in \mathbb{R}^{F \times K}$ contains the shape coefficients c_{ik} and $B^\sharp \in \mathbb{R}^{K \times 3n}$ is constructed as X^\sharp and contains the basis elements.

A suitable objective function is thus given by

$$\mu \text{rank}(X^\sharp) + \|RX - M\|_F^2, \quad (25)$$

where $R \in \mathbb{R}^{2F \times 3F}$ is a block-diagonal matrix with the camera matrices R_i on the main diagonal, $X \in \mathbb{R}^{3F \times n}$ is the concatenation of the 3D points X_i , and $M \in \mathbb{R}^{2F \times n}$ is the concatenated observed image points x_i . By replacing the rank penalty with a relaxation and minimize it using the proposed method and the methods used in the previous section. The regularization parameter is varied for the respective methods in order to obtain a rank 1–8 solution, and the respective datafit is reported in Figure 5, for four different sequences.

In all sequences, the best datafit for each rank level is obtained by our proposed method. IRNN and ADMM using f_μ is able to give the same, or very similar, datafit for lower ranks, but for solutions with rank larger than four our method consistently reports a lower value than the competing state-of-the-art methods.

6. Conclusions

In this paper we presented a unification of bilinear parameterization and rank regularization. Robust penalties for rank regularization has often been used together with splitting schemes, but it has been shown that such methods yield

¹ The datasets are available here: <http://www.maths.lth.se/matematiklth/personal/calle/dataset/dataset.html>.

unsatisfactory results for ill-posed problems in several computer vision applications. By using the bilinear formulation, the objective functions become differentiable, and convergence rates in the neighborhood of a local minimum are faster. Furthermore, we showed that theoretical optimality results known from the regularization formulations can be lifted to the bilinear formulation.

Lastly, the generality of the proposed framework allows for a wide range of problems, some of which, have not been amenable by state-of-the-art methods, but have been proven successful using our proposed method.

References

- [1] F. R. Bach. Convex relaxations of structured matrix factorizations. *CoRR*, abs/1309.3117, 2013.
- [2] R. Basri, D. Jacobs, and I. Kemelmacher. Photometric stereo with general, unknown lighting. *International Journal of Computer Vision*, 72(3):239–257, May 2007.
- [3] S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. In *Annual Conference in Neural Information Processing Systems (NIPS)*. 2016.
- [4] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.*, 3(1):1–122, 2011.
- [5] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [6] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3d shape from image streams. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2000.
- [7] A. M. Buchanan and A. W. Fitzgibbon. Damped newton algorithms for matrix factorization with missing data. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005.
- [8] R. Cabral, F. De la Torre, J. P. Costeira, and A. Bernardino. Unifying nuclear norm and bilinear factorization approaches for low-rank matrix decomposition. In *International Conference on Computer Vision (ICCV)*, 2013.
- [9] E. J. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *J. ACM*, 58(3):11:1–11:37, 2011.
- [10] E. J. Candès and B. Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [11] L. Canyi, J. Tang, S. Yan, and Z. Lin. Generalized nonconvex nonsmooth low-rank minimization. *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [12] L. Canyi, J. Tang, S. Yan, and Z. Lin. Nonconvex nonsmooth low-rank minimization via iteratively reweighted nuclear norm. *IEEE Transactions on Image Processing*, 25, 10 2015.
- [13] M. Carlsson. On convexification/optimization of functionals including an l2-misfit term. *arXiv preprint arXiv:1609.09378*, 2016.
- [14] M. Carlsson, D. Gerosa, and C. Olsson. An unbiased approach to compressed sensing. *arXiv preprint, arXiv:1806.05283*, 2018.
- [15] Y. Dai, H. Li, and M. He. A simple prior-free method for non-rigid structure-from-motion factorization. *International Journal of Computer Vision*, 107(2):101–122, 2014.
- [16] A. Eriksson and A. Hengel. Efficient computation of robust weighted low-rank matrix approximations using the L_1 norm. *IEEE Trans. Pattern Anal. Mach. Intell.*, 34(9):1681–1690, 2012.
- [17] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- [18] M. Fazel, H. Hindi, and S. P. Boyd. A rank minimization heuristic with application to minimum order system approximation. In *American Control Conference*, 2001.
- [19] R. T. Frankot and R. Chellappa. A method for enforcing integrability in shape from shading algorithms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10:439–451, 1988.
- [20] J. H. Friedman. Fast sparse regression and classification. *International Journal of Forecasting*, 28(3):722 – 738, 2012.
- [21] C. Gao, N. Wang, Q. R. Yu, and Z. Zhang. A feasible nonconvex relaxation approach to feature selection. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, 2011.
- [22] R. Garg, A. Roussos, and L. Agapito. A variational approach to video registration with subspace constraints. *International Journal of Computer Vision*, 104(3):286–314, 2013.
- [23] R. Garg, A. Roussos, and L. de Agapito. Dense variational reconstruction of non-rigid surfaces from monocular video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2013.
- [24] R. Ge, C. Jin, and Y. Zheng. No spurious local minima in nonconvex low rank problems: A unified geometric analysis. *arXiv preprint, arXiv:1704.00708*, 2017.
- [25] R. Ge, J. D. Lee, and T. Ma. Matrix completion has no spurious local minimum. In *Annual Conference on Neural Information Processing Systems (NIPS)*, 2016.
- [26] D. Geman and C. Yang. Nonlinear image recovery with half-quadratic regularization. *IEEE Transactions on Image Processing*, 4:932 – 946, 08 1995.
- [27] N. Gillis and F. Glinuer. Low-rank matrix approximation with weights or missing data is np-hard. *SIAM Journal on Matrix Analysis and Applications*, 32(4), 2011.
- [28] S. Gu, Q. Xie, D. Meng, W. Zuo, X. Feng, and L. Zhang. Weighted nuclear norm minimization and its applications to low level vision. *International Journal of Computer Vision*, 121, 07 2016.
- [29] B. D. Haeffele and R. Vidal. Structured low-rank matrix factorization: Global optimality, algorithms, and applications. *CoRR*, abs/1708.07850, 2017.
- [30] J. He, L. Balzano, and A. Szlam. Incremental gradient on the grassmannian for online foreground and background separation in subsampled video. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1568–1575, June 2012.

- [31] J. H. Hong and A. Fitzgibbon. Secrets of matrix factorization: Approximations, numerics, manifold optimization and random restarts. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [32] J. H. Hong, C. Zach, A. Fitzgibbon, and R. Cipolla. Projective bundle adjustment from arbitrary initialization using the variable projection method. In *European Conference on Computer Vision (ECCV)*, 2016.
- [33] J. H. Hong, C. Zach, A. Fitzgibbon, and R. Cipolla. Projective bundle adjustment from arbitrary initialization using the variable projection method. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [34] Y. Hu, D. Zhang, J. Ye, X. Li, and X. He. Fast and accurate matrix completion via truncated nuclear norm regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(9):2117–2130, 2013.
- [35] J. Hyeong Hong and C. Zach. pose: Pseudo object space error for initialization-free bundle adjustment. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [36] V. Larsson and C. Olsson. Convex low rank approximation. *International Journal of Computer Vision*, 120(2):194–214, 2016.
- [37] V. Larsson and C. Olsson. Compact matrix factorization with dependent subspaces. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4361–4370, 07 2017.
- [38] A. S. Lewis. The convex analysis of unitarily invariant matrix functions. *Journal of Convex Analysis*, 2(1):173–183, 1995.
- [39] K. Mohan and M. Fazel. Iterative reweighted least squares for matrix rank minimization. In *Annual Allerton Conference on Communication, Control, and Computing*, pages 653–661, 2010.
- [40] F. Nie, H. Wang, X. Cai, H. Huang, and C. H. Q. Ding. Robust matrix completion via joint Schatten p -norm and l_p -norm minimization. In *ICDM*, pages 566–574, 2012.
- [41] T. H. Oh, Y. W. Tai, J. C. Bazin, H. Kim, and I. S. Kweon. Partial sum minimization of singular values in robust PCA: Algorithm and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(4):744–758, 2016.
- [42] C. Olsson, M. Carlsson, F. Andersson, and V. Larsson. Non-convex rank/sparsity regularization and local minima. *Proceedings of the International Conference on Computer Vision*, 2017.
- [43] C. Olsson and O. Enqvist. Stable structure from motion for unordered image collections. In *Proceedings of the Scandinavian Conference on Image Analysis (SCIA)*, pages 524–535, 2011.
- [44] S. Oymak, A. Jalali, M. Fazel, Y. C. Eldar, and B. Hassibi. Simultaneously structured models with application to sparse and low-rank matrices. *IEEE Transactions on Information Theory*, 61(5):2886–2908, 2015.
- [45] S. Oymak, K. Mohan, M. Fazel, and B. Hassibi. A simplified approach to recovery conditions for low rank matrices. In *IEEE International Symposium on Information Theory Proceedings (ISIT)*, pages 2318–2322, 2011.
- [46] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.*, 52(3):471–501, Aug. 2010.
- [47] F. Shang, J. Cheng, Y. Liu, Z. Luo, and Z. Lin. Bilinear factor matrix norm minimization for robust PCA: Algorithms and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(9):2066–2080, Sep. 2018.
- [48] K.-C. Toh and S. Yun. An accelerated proximal gradient algorithm for nuclear norm regularized least squares problems. *Pacific Journal of Optimization*, 6, 09 2010.
- [49] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization method. *International Journal of Computer Vision*, 9(2):137–154, 1992.
- [50] M. Uchiyama. Subadditivity of eigenvalue sums. *Proceedings of The American Mathematical Society*, 134:1405–1412, 05 2005.
- [51] N. Wang, T. Yao, J. Wang, and D.-Y. Yeung. A probabilistic approach to robust matrix factorization. In *European Conference on Computer Vision (ECCV)*, 2012.
- [52] C. Xu, Z. Lin, and H. Zha. A unified convex surrogate for the Schatten- p norm. In *Proceedings of the Conference on Artificial Intelligence (AAAI)*, 2017.
- [53] J. Yan and M. Pollefeys. A factorization-based approach for articulated nonrigid shape, motion and kinematic chain recovery from video. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(5):865–877, 2008.
- [54] C.-H. Zhang et al. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894–942, 2010.

A. Proofs

In this section we present the proofs of Theorems 2 and 3. Our analysis will make use of the differentiable objective

$$\mathcal{D}(B, C) := \tilde{\mathcal{R}}(B, C) + \|ABC^T - b\|^2, \quad (26)$$

the non-convex function

$$\mathcal{N}(X) := \mathcal{R}(X) + \|AX - b\|^2, \quad (27)$$

and the convex function

$$\mathcal{C}(X) = \mathcal{R}(X) + \|X - Z\|_F^2. \quad (28)$$

We will also use the functions

$$\tilde{G}(B, C) = \tilde{\mathcal{R}}(B, C) + \|BC^T\|_F^2, \quad (29)$$

$$G(X) = \mathcal{R}(X) + \|X\|_F^2, \quad (30)$$

$$H(X) = \|AX - b\|^2 - \|X\|_F^2. \quad (31)$$

Note that $\mathcal{D}(B, C) = \tilde{G}(B, C) + H(BC^T)$ and $\mathcal{N}(X) = G(X) + H(X)$. Throughout the section we use $f = f_\mu$ with f_μ as in (8) (of the main paper) but for simplicity of notation we will suppress the subscript μ . Furthermore, the subdifferential $\partial G(X)$ of G will be of importance. Let $g(x) = f(|x|) + x^2$. The scalar function g has

$$\partial g(x) = \begin{cases} 2x & |x| \geq \sqrt{\mu} \\ 2\sqrt{\mu}\text{sign}(x) & 0 < |x| \leq \sqrt{\mu} \\ 2\sqrt{\mu}[-1, 1] & x = 0 \end{cases}. \quad (32)$$

The following lemma shows how to compute ∂G for the matrix case using ∂g .

Lemma 1. *The subdifferential of $G(X)$ is given by*

$$\partial G(X) = \{U\partial g(\Sigma)V^T + M : \sigma_1(M) \leq 2\sqrt{\mu}, \\ U^T M = 0 \text{ and } MV^T = 0\} \quad (33)$$

where $X = U\Sigma V^T$ is the SVD and $\partial g(\Sigma)$ is the matrix of same size as Σ with diagonal elements $\partial g(\sigma_i)$.

Next we give the stationary point conditions for \mathcal{D} that are needed for proving Theorem 2.

Lemma 2. *Let $B = U\sqrt{\Sigma}$, $C = V\sqrt{\Sigma}$ and $X = U\Sigma V^T$. If (B, C) is a stationary point of \mathcal{D} , then*

$$0 = B\partial G(\Sigma) + \nabla H(BC^T)C, \quad (34)$$

$$0 = \partial G(\Sigma)C^T + B^T\nabla H(BC^T). \quad (35)$$

We are now ready to prove Theorem 2.

Proof of Theorem 2. Let $\bar{X} = \bar{B}\bar{C}^T$, $\tilde{X} = \tilde{B}\tilde{C}^T$ and $\Delta X = \tilde{B}\tilde{C}^T - \bar{B}\bar{C}^T$. We first note that the limit

$$\mathcal{N}'_{\Delta X}(\bar{X}) = \lim_{t \searrow 0} \frac{\mathcal{N}(\bar{X} + t\Delta X) - \mathcal{N}(\bar{X})}{t}, \quad (36)$$

exists since \mathcal{N} is a sum of a finite convex function G and a differentiable function H . Our goal is now to show that the limit is non-negative. Suppose that we can find a factorization $B(t)C(t)^T = \bar{X} + t\Delta X$, such that $\mathcal{R}(\bar{X} + t\Delta X) = \tilde{\mathcal{R}}(B(t), C(t))$, $(B(t), C(t))$ is continuous and $(B(0), C(0)) = (\bar{B}, \bar{C})$. Then for small enough t we have

$$\mathcal{N}(\bar{X} + t\Delta X) - \mathcal{N}(\bar{X}) = \mathcal{D}(B(t), C(t)) - \mathcal{D}(\bar{B}, \bar{C}). \quad (37)$$

This quantity is clearly non-negative since (\bar{B}, \bar{C}) is a local minimizer of \mathcal{D} , which would prove that the limit (36) is non-negative. It is not difficult to see that this can be done when the two matrices \bar{X} and \tilde{X} have singular value decompositions with the same U and V . In what follows we will first show that all other cases can be reduced so that the matrices are of this form. When this is done we proceed to construct the factorization $B(t)C(t)^T$ which completes the proof.

The directional derivatives can be computed using the sub-differential

$$\mathcal{N}'_{\Delta X} = \max_{Z \in \partial G(\bar{B}\bar{C}^T)} \langle 2Z, \Delta X \rangle + \langle \nabla H(\bar{B}\bar{C}^T), \Delta X \rangle. \quad (38)$$

By Lemma 1, the first term becomes

$$\langle U\partial G(\Sigma)V^T + M, \Delta X \rangle = \langle U\partial G(\Sigma)V^T, \tilde{B}\tilde{C}^T \rangle \\ + \langle M, \tilde{B}\tilde{C}^T \rangle \\ - \langle U\partial G(\Sigma)V^T, \bar{B}\bar{C}^T \rangle. \quad (39)$$

The columns of \tilde{B} can be written as a linear combination of the columns in \bar{B} and those of a matrix \tilde{B}_\perp with at most k columns that are perpendicular to \bar{B} . Similarly, the columns of \tilde{C} can be written as a linear combination of the columns in \bar{C} and those of a matrix \tilde{C}_\perp with at most k columns that are perpendicular to \bar{C} . Therefore, we may write

$$\tilde{B}\tilde{C}^T = [\bar{B} \quad \tilde{B}_\perp] \begin{bmatrix} K_{11} & K_{12} \\ K_{21} & K_{22} \end{bmatrix} \begin{bmatrix} \bar{C}^T \\ \tilde{C}_\perp^T \end{bmatrix} \\ = \bar{B}K_{11}C^T + \bar{B}K_{12}\tilde{C}_\perp^T \\ + \tilde{B}_\perp K_{21}\bar{C}^T + \tilde{B}_\perp K_{22}\tilde{C}_\perp^T, \quad (40)$$

where $\tilde{B}^T\tilde{B}_\perp = 0$ and $\tilde{C}^T\tilde{C}_\perp = 0$. Our goal is now to show that the terms K_{12} and K_{21} and the off diagonal elements of K_{11} vanish from (38) and can be assumed to be zero.

For the last term of (39) we have

$$\langle U\partial G(\Sigma)V^T, \tilde{B}\tilde{C}^T \rangle = \langle \partial G(\Sigma), U^T\tilde{B}\tilde{C}^TV \rangle \\ = \langle \partial G(\Sigma), \Sigma \rangle, \quad (41)$$

which is clearly independent of \tilde{B} and \tilde{C} . The first term of (39) reduces to

$$\begin{aligned} \langle U\partial G(\Sigma)V^T, \tilde{B}\tilde{C}^T \rangle &= \langle U\partial G(\Sigma)V^T, \bar{B}K_{11}\bar{C}^T \rangle \\ &= \langle \bar{B}^T U\partial G(\Sigma)V^T \bar{C}, K_{11} \rangle \quad (42) \\ &= \langle \Sigma\partial G(\Sigma), K_{11} \rangle. \end{aligned}$$

Note that the off diagonal elements of K_{11} vanish from this expression since $\Sigma\partial G(\Sigma)$ is diagonal. Similarly, the second term of (39) reduces to

$$\langle M, \tilde{B}\tilde{C}^T \rangle = \langle M, \bar{B}_\perp K_{22}\bar{C}_\perp^T \rangle. \quad (43)$$

We now consider the second term of (38)

$$\begin{aligned} \langle \nabla H(\bar{B}\bar{C}^T), \Delta X \rangle &= \\ \langle \nabla H(\bar{B}\bar{C}^T), \bar{B}K_{11}\bar{C}^T + \bar{B}K_{12}\bar{C}_\perp^T & \quad (44) \\ + \bar{B}_\perp K_{21}\bar{C}^T + \bar{B}_\perp K_{22}\bar{C}_\perp^T - \bar{B}\bar{C}^T \rangle. \end{aligned}$$

For the first term we have

$$\begin{aligned} \langle \nabla H(\bar{B}\bar{C}^T), \bar{B}K_{11}\bar{C}^T \rangle &= \langle \nabla H(\bar{B}\bar{C}^T)\bar{C}, \bar{B}K_{11} \rangle \\ &= -\langle \bar{B}\partial G(\Sigma), \bar{B}K_{11} \rangle \\ &= -\langle \bar{B}^T \bar{B}\partial G(\Sigma), K_{11} \rangle \quad (45) \\ &= -\langle \Sigma\partial G(\Sigma), K_{11} \rangle. \end{aligned}$$

Again the off diagonal elements of K_{11} vanish. For the second term of (44) we have

$$\begin{aligned} \langle \nabla H(\bar{B}\bar{C}^T), \bar{B}K_{12}\bar{C}_\perp^T \rangle &= \langle \bar{B}^T \nabla H(\bar{B}\bar{C}^T), K_{12}\bar{C}_\perp^T \rangle \\ &= -\langle \partial G(\Sigma)\bar{C}^T, K_{12}\bar{C}_\perp \rangle \\ &= -\langle \partial G(\Sigma)\bar{C}^T \bar{C}_\perp, K_{12} \rangle = 0. \end{aligned} \quad (46)$$

Similarly, the third term is $\langle \nabla H(\bar{B}\bar{C}^T), \bar{B}_\perp K_{21}\bar{C}^T \rangle = 0$. Thus

$$\begin{aligned} \langle \nabla H(\bar{B}\bar{C}^T), \Delta X \rangle &= \langle \nabla H(\bar{B}\bar{C}^T), \bar{B}_\perp^T K_{22}\bar{C}_\perp^T \rangle \\ &\quad - \langle \Sigma\partial G(\Sigma), K_{11} \rangle \quad (47) \\ &\quad - \langle \nabla H(\bar{B}\bar{C}^T), \bar{B}\bar{C}^T \rangle. \end{aligned}$$

Summarizing we see that we have now proven that all the terms in (39) are independent of K_{12} , K_{21} as well as the off diagonal terms of K_{11} . They therefore do not affect the value of $\mathcal{N}'_{\Delta X}$ and can be assumed to be zero. We can now write ΔX as

$$\Delta X = [U \quad U_\perp] \begin{bmatrix} (D - I)\Sigma & 0 \\ 0 & \tilde{\Sigma} \end{bmatrix} \begin{bmatrix} V^T \\ V_\perp^T \end{bmatrix}, \quad (48)$$

where D are the diagonal elements of K_{11} and $U_\perp \tilde{\Sigma} V_\perp^T$ is the SVD of $\bar{B}_\perp K_{22}\bar{C}_\perp^T$. Note that $U_\perp^T U = 0$ since U and U_\perp span orthogonal subspaces. Similarly $V_\perp^T V = 0$.

We now consider the directional derivative (36) with $\bar{B} = U\sqrt{\Sigma}$, $\bar{C} = V\sqrt{\Sigma}$. It is clear that for small t the matrix $\bar{X} + t\Delta X$ has the singular value decomposition

$$[U \quad U_\perp] \begin{bmatrix} ((1-t)I + tD)\Sigma & 0 \\ 0 & t\tilde{\Sigma} \end{bmatrix} \begin{bmatrix} V^T \\ V_\perp^T \end{bmatrix}. \quad (49)$$

We now let

$$B(t) = [U \quad U_\perp] \sqrt{\begin{bmatrix} ((1-t)I + tD)\Sigma & 0 \\ 0 & t\tilde{\Sigma} \end{bmatrix}}, \quad (50)$$

$$C(t) = [V \quad V_\perp] \sqrt{\begin{bmatrix} ((1-t)I + tD)\Sigma & 0 \\ 0 & t\tilde{\Sigma} \end{bmatrix}}. \quad (51)$$

Then, we clearly have $\tilde{\mathcal{R}}(B(t), C(t)) = \mathcal{R}(X + t\Delta X)$ for small enough t , which completes the proof. \square

Next we will prove Theorem 3. Our results build on those of [42] and we remind the reader that we exclusively use $f_\mu(\sigma) = \mu - \max(\sqrt{\mu} - \sigma, 0)^2$ throughout this section, but suppress the subscript μ . We will use the fact that the directional derivatives in a local minimum are non-negative for all low rank directions to show that (\bar{B}, \bar{C}) minimizes the non-convex \mathcal{N} over matrices of rank $< k$ in Theorem 3. For this we will need the following result:

Lemma 3. *If \bar{X} is a solution to $\min_{\text{rank}(X) \leq k} \mathcal{C}(X)$ with $\text{rank}(\bar{X}) < k$ and the singular values of Z fulfill $\sigma_i(Z) \notin [(1 - \delta_{2k})\sqrt{\mu}, \frac{\sqrt{\mu}}{(1 - \delta_{2k})}]$ then \bar{X} also solves $\min_X \mathcal{C}(X)$.*

Proof of Lemma 3. By von Neumann's trace theorem it is easy to see that the problem $\min_{\text{rank}(X) \leq k} \mathcal{C}(X)$ reduces to a minimization over the singular values of X . We should thus find $\sigma_i(X)$ such that

$$\sum_{i=1}^n \underbrace{-\max(\sqrt{\mu} - \sigma_i(X), 0)^2 + (\sigma_i(X) - \sigma_i(Z))^2}_{:=g_i(\sigma_i(X))} \quad (52)$$

is minimized and at most k singular values are non-zero. The unconstrained minimizers of g_i can be written down in closed form: If $0 \leq \sqrt{\mu} < \sigma_i(Z)$ then $\sigma_i(X) = \sigma_i(Z)$ is optimal giving $g_i(\sigma_i(X)) = 0$. If $0 \leq \sigma_i(Z) < \sqrt{\mu}$ then $\sigma_i(X) = 0$ is optimal giving $g_i(\sigma_i(X)) = -\mu + \sigma_i(Z)^2$. Hence for any solution of $\min_{\text{rank}(X) \leq k} \mathcal{C}(X)$ we have $\sigma_i(X) = 0$ if $0 \leq \sigma_i(Z) \leq \sqrt{\mu}$. There are now two cases:

1. If $\sigma_{k+1}(Z) < \sqrt{\mu}$ then the sequence of unconstrained minimizers has at most k non-zero values. Thus, in this case the resulting X solves both $\min_X \mathcal{C}(X)$ and $\min_{\text{rank}(X) \leq k} \mathcal{C}(X)$.
2. If $\sigma_{k+1} > \sqrt{\mu}$ we will not be able to select $\sigma_i(X) = \sigma_i(Z)$ for all i where $0 \leq \sqrt{\mu} < \sigma_i(Z)$. Choosing $\sigma_i(X) = 0$ gives $g_i(0) = -\mu + \sigma_i(Z)^2 < 0$. Since

$\sigma_i(Z)$ is decreasing with i it is clear that the smallest value is obtained when selecting $\sigma_i(X) = \sigma_i(Z)$ for $i = 1, \dots, k$.

We now conclude that if $\text{rank}(\bar{X}) < k$ then we are in case 1 and therefore \bar{X} solves the unconstrained problem. \square

We are now ready to give the proof of Theorem 3.

Proof of Theorem 3. Since \mathcal{C} and \mathcal{N} has the same subdifferential (see [36]) at $\bar{X} = \bar{B}\bar{C}^T$ it is clear that the directional derivatives $\mathcal{C}'_{\Delta X}(\bar{X}) = \mathcal{N}'_{\Delta X}(\bar{X}) \geq 0$, where $\Delta X = \bar{X} - \bar{B}\bar{C}^T$ and $\text{rank}(\bar{X}) \leq k$. By convexity of \mathcal{C} it is then also clear that

$$\bar{B}\bar{C}^T \in \underset{\text{rank}(X) \leq k}{\text{arg min}} \mathcal{C}(X). \quad (53)$$

Since $\text{rank}(\bar{B}\bar{C}^T) < k$, $\bar{B}\bar{C}^T$ is also the unrestricted global minimizer of $\mathcal{C}(X)$ according to Lemma 3. By Lemma 3.1 of [42] it is then a stationary point of $\mathcal{N}(X)$.

What remains now is to prove that $\bar{X} = \bar{B}\bar{C}^T$ is a global minimizer of \mathcal{N} over all line segments $\bar{X} + t\Delta X$. This can be done by estimating the growth of the directional derivatives along such lines. For this purpose we consider the functions G and H defined as in (30) and (31). Note that \bar{X} is a stationary point of $\mathcal{N}(X) = G(X) + H(X)$ if and only if $-\nabla H(\bar{X}) = 2Z \in \partial G(\bar{X})$.

Since $\nabla H(\bar{X} + t\Delta X) - \nabla H(\bar{X}) = t\nabla H(\Delta X) = 2t(\mathcal{A}^* \mathcal{A} \Delta X - \Delta X)$ we have

$$\langle \nabla H(\bar{X} + t\Delta X) - \nabla H(\bar{X}), t\Delta X \rangle = 2t^2(\|\mathcal{A} \Delta X\|_F^2 - \|\Delta X\|_F^2), \quad (54)$$

and due to RIP $\|\mathcal{A} \Delta X\|_F^2 - \|\Delta X\|_F^2 \geq -\delta_{2r} \|\Delta X\|_F^2$. From Corollary 4.2 of [42] we see that for any $2Z' \in \partial G(\bar{X} + t\Delta X)$ we have

$$\langle Z' - Z, t\Delta X \rangle > t^2 \delta_{2r} \|\Delta X\|_F^2, \quad (55)$$

as long as $t \neq 0$. Since $G'_{\Delta X}(X) = \max_{2Z \in \partial G(X)} \langle 2Z, \Delta X \rangle$, $H'_{\Delta X}(X) = \langle \nabla H(X), \Delta X \rangle$ and $2Z + \nabla H'(\bar{X}) = 0$ we get

$$\mathcal{N}'_{\Delta X}(\bar{X} + t\Delta X) \geq \langle 2Z' + \nabla H(\bar{X} + t\Delta X), \Delta X \rangle > 0 \quad (56)$$

This shows that \bar{X} solves (9). That \bar{X} also solves (10) is now a consequence of the fact that $\mathcal{R}(X) \leq \mu \text{rank}(X)$ with equality if X have no singular values in the interval $(0, \sqrt{\mu}]$. Note that \bar{X} is the unrestricted minimizer of $\mathcal{C}(X)$, where the singular values of Z fulfill $\sigma_i(Z) \notin \left[(1 - \delta_{2k})\sqrt{\mu}, \frac{\sqrt{\mu}}{1 - \delta_{2k}} \right]$. Since the solution to this problem is hard thresholding \bar{X} has no singular values in $\left(0, \frac{\sqrt{\mu}}{1 - \delta_{2k}}\right] \supset (0, \sqrt{\mu}]$. \square

For completeness we give the proofs that were previously omitted.

Proof of Lemma 1. With some abuse of notation we define the function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ by $g(\mathbf{x}) = \sum_{i=1}^n g(x_i)$, where x_i , $i = 1, \dots, n$ are the elements of \mathbf{x} and $g(x) = f(|x|) + x^2$. The function g is an absolutely symmetric convex function and G can be written $G(X) = g \circ \sigma(X)$, where $\sigma(X)$ is the vector of singular values of X . Then according to [38] the matrix $Y \in \partial G(X)$ if and only if $Y = U' \text{diag}(\partial g \circ \sigma(X)) V'^T$ when $X = U' \text{diag}(\sigma(X)) V'^T$. (Here we use the full SVD with square orthogonal matrices U' and V' .) Now given a thin SVD $X = U \Sigma V^T$ all possible full SVD's of X can be written

$$X = [U \quad U_{\perp}] \begin{bmatrix} \Sigma & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V^T \\ V_{\perp}^T \end{bmatrix}, \quad (57)$$

where U_{\perp} and V_{\perp} are singular vectors corresponding to singular values that are zero. Note that U_{\perp} and V_{\perp} are not uniquely defined since their corresponding singular values are all zero. Therefore we get

$$Y = [U' \quad U'_{\perp}] \begin{bmatrix} \partial g(\Sigma) & 0 \\ 0 & D \end{bmatrix} \begin{bmatrix} V'^T \\ V'_{\perp}{}^T \end{bmatrix} = U' \partial g(\Sigma) V'^T + U_{\perp} D V_{\perp}^T, \quad (58)$$

where D is a diagonal matrix with elements in $2\sqrt{\mu}[-1, 1]$. It is clear that $\sigma_1(U_{\perp} D V_{\perp}^T) = \sigma_1(D) \leq 2\sqrt{\mu}$. Furthermore, since U_{\perp} and V_{\perp} can be any orthogonal bases of the spaces perpendicular to the column and row spaces of X , it is clear that any matrix M fulfilling $U^T M = 0$, $M V = 0$ and $\sigma_1(M) \leq 2\sqrt{\mu}$ can be written $M = U_{\perp} D V_{\perp}^T$, hence

$$\partial G(X) = \{U \partial g(\Sigma) V^T + M : \sigma_1(M) \leq 2\sqrt{\mu}, U^T M = 0, M V = 0\}. \quad (59)$$

\square

Proof of Lemma 2. The gradients of \tilde{G} are given by

$$\nabla_B \tilde{G}(B, C) = \nabla_B(\tilde{\mathcal{R}}(B, C)) + \nabla_B(\|BC^T\|_F^2). \quad (60)$$

For the first term we get

$$\nabla_{B_i} \tilde{\mathcal{R}}(B, C) = f' \left(\frac{\|B_i\|^2 + \|C_i\|^2}{2} \right) B_i. \quad (61)$$

With $B = U\sqrt{\Sigma}$ and $C = V\sqrt{\Sigma}$ we get

$$\nabla_B \tilde{\mathcal{R}}(B, C) = B \begin{bmatrix} f'(\sigma_1) & 0 & \dots \\ 0 & f'(\sigma_2) & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} = B f'(\Sigma), \quad (62)$$

which gives

$$\nabla_B \tilde{G}(B, C) = Bf'(\Sigma) + 2BC^T C = B(f'(\Sigma) + 2\Sigma). \quad (63)$$

For a non-zero σ we have $\partial g(\sigma) = \{f'(\sigma) + 2\sigma\}$ and therefore

$$\nabla_B \tilde{G}(B, C) = B(\partial G(\Sigma)), \quad (64)$$

where $g(X) = \mathcal{R}_\mu(X) + \|X\|_F^2$. Similarly we get

$$\nabla_C \tilde{G}(B, C) = C(\partial G(\Sigma)). \quad (65)$$

If (B, C) is a stationary point then

$$0 = B\partial G(\Sigma) + \nabla H(BC^T)C, \quad (66)$$

$$0 = C\partial G(\Sigma) + (\nabla H(BC^T))^T B. \quad (67)$$

The second equation can be re-written to the form stated in the lemma. \square

B. Implementation Details

In this section we present some more details on our Iteratively Reweighted VarPro approach. Recall that our approach consists of three main steps. In the first step we make a quadratic approximation (20) of the regularization term by replacing $\tilde{\mathcal{R}}(B, C)$ with $\sum_{i=1}^k w_i^{(t)} (\|B_i\|^2 + \|C_i\|^2)$ as described in Section 4.

In the second step we apply one step of VarPro with the Ruhe Wedin approximation, see [33] for details on the implementation. VarPro uses Jacobians with respect to both the B and C parameters. In our case we have two terms that needs to be linearized. The regularization term can be written

$$\|\text{diag}(w^{(t)})B\|_F^2 + \|\text{diag}(w^{(t)})C\|_F^2, \quad (68)$$

where $\text{diag}(w^{(t)})$ is a diagonal matrix with the weights $w_i^{(t)}$ in the diagonal. The residuals $\text{diag}(w^{(t)})B$ are already linear and by column stacking the variables we can write them as $J_B^{\text{reg}} \mathbf{b}$, where \mathbf{b} is a column stacked version of B . If B has k columns the matrix J_B^{reg} will consist of k copies of the matrix $\text{diag}(w^{(t)})$. Additionally, each row of J_B^{reg} has only one non-zero element making the matrix extremely sparse. Similarly, we obtain the contribution due to the second bilinear factor C , which can be written as $J_C^{\text{reg}} \mathbf{c}$. Here we use $\mathbf{c} = \text{vec}(C^T)$, as it alleviates the computations of the data terms, hence J_C^{reg} consists of a k copies of $\text{diag}(w^{(t)})$ permuted to match this design choice. Given a current iterate $(\mathbf{b}^{(t)}, \mathbf{c}^{(t)})$ we write the regularization term as $\|J_B^{\text{reg}} \delta \mathbf{b} + \mathbf{r}_B\|^2 + \|J_C^{\text{reg}} \delta \mathbf{c} + \mathbf{r}_C\|^2$, where $\mathbf{r}_B = J_B^{\text{reg}} \mathbf{b}^{(t)}$, $\mathbf{r}_C = J_C^{\text{reg}} \mathbf{c}^{(t)}$, $\mathbf{b} = \mathbf{b}^{(t)} + \delta \mathbf{b}$ and $\mathbf{c} = \mathbf{c}^{(t)} + \delta \mathbf{c}$.

Linearizing the residuals $ABC^T - b$ around $(\mathbf{b}^{(t)}, \mathbf{c}^{(t)})$ gives an expression of the form

$$J_B^{\text{data}} \delta \mathbf{b} + J_C^{\text{data}} \delta \mathbf{c} + \mathbf{r}^{\text{data}}. \quad (69)$$

The particular shape of the Jacobians in this expression depends on the application; however, in all of our applications they are sparse. For example, in the missing data problem each residual corresponds to an element of the matrix X which in turn only depends on k elements of B and C . Locally we may now write the objective function as

$$\|J_B \delta \mathbf{b} + J_C \delta \mathbf{c} + \mathbf{r}\|^2, \quad (70)$$

where

$$J_B = \begin{bmatrix} J_B^{\text{reg}} \\ 0 \\ J_B^{\text{data}} \end{bmatrix}, \quad J_C = \begin{bmatrix} 0 \\ J_C^{\text{reg}} \\ J_C^{\text{data}} \end{bmatrix}, \quad \mathbf{r} = \begin{bmatrix} \mathbf{r}_B \\ \mathbf{r}_C \\ \mathbf{r}^{\text{data}} \end{bmatrix}. \quad (71)$$

It was shown in [32] that each step of VarPro is equivalent to first minimizing (70) with the additional dampening term $\lambda \|\delta \mathbf{b}\|^2$ and then performing an exact optimization of (20) over the C -variables (when fixing the B -variables to their new values). Since we also have a reweighing we only do one iteration with VarPro before updating the weights $w^{(t)}$.

The above procedure can return stationary points for which $\tilde{\mathcal{R}}(B, C) > \mathcal{R}(BC^T)$. Our last step is designed to escape such points by taking the current iterate and recompute the factorization of $\tilde{B}\tilde{C}^T$ using SVD. If the SVD of $\tilde{B}\tilde{C}^T = \sum_{i=1}^r \sigma_i U_i V_i^T$ we update \tilde{B} and \tilde{C} to $\tilde{B}_i = \sqrt{\sigma_i} U_i$ and $\tilde{C}_i = \sqrt{\sigma_i} V_i$ which we know reduces the energy and gives $\tilde{\mathcal{R}}(\tilde{B}, \tilde{C}) = \mathcal{R}(\tilde{B}\tilde{C}^T)$. Therefore we proceed by refactorizing the current iterate using SVD in each iteration. The detailed steps of the bilinear method are summarized in Algorithm 1.

C. Additional Experiments on Real Data

C.1. pOSE: Pseudo Object Space Error

In this section we compare the energies over time for ADMM optimizing the same energy [36], *i.e.* with the regularizer \mathcal{R} , and $f = f_\mu$ as in (8) (of the main paper), and our proposed method. We let the bilinear method run until convergence, and let ADMM execute the same time in seconds. As a comparison we use the nuclear norm relaxation and the discontinuous rank regularization. The results of the experiment are shown in Figure 6.

Again, note that the bilinear method optimizes the same energy as ADMM- \mathcal{R}_μ , and that, despite the initial fast lowering of the objective value, the ADMM approach fails to reach the global optimum, within the allotted 150 seconds. This holds true for all methods employing ADMM. In all experiments, the control parameter $\eta = 0.5$, and the μ parameter was chosen to be smaller than all non-zero singular values of the best known optimum (obtained using VarPro). For a fair comparison, the μ -value for the nuclear norm relaxation, was modified due to the shrinking bias, and was chosen to be the smallest value of μ for which a solution with accurate rank was obtained. Due to this modification,

Input: Robust penalty function f , linear operator \mathcal{A} and regularization parameter μ , damping parameter λ .
Initialize B and C with random entries
while not converged **do**
 Compute weights $w^{(t)}$ from current iterate (B, C)
 Compute the vectorizations $\mathbf{b} = \text{vec}(B)$, $\mathbf{c} = \text{vec}(C^T)$
 Compute residuals \mathbf{r}_B , \mathbf{r}_C , and Jacobians J_B^{data} and J_C^{data} depending on \mathcal{A}
 Compute residual \mathbf{r}^{reg} , and Jacobians J_B^{reg} and J_C^{reg}
 Create full residual \mathbf{r} and Jacobians J_B and J_C
 Compute $\tilde{J}^T \tilde{J} + \lambda I = J_B^T (I - J_C J_C^+) J_B + \lambda I$
 Compute $\mathbf{b}' = \mathbf{b} - (\tilde{J}^T \tilde{J} + \lambda I)^{-1} J_B^T \mathbf{r}$ and reshape into matrix B'
 Compute C' by minimizing (20) with fixed B'
 if $\mathcal{R}(B' C'^T) + \|\mathcal{A}(B' C'^T) - \mathbf{b}\|^2 < \mathcal{R}(B C^T) + \|\mathcal{A}(B C^T) - \mathbf{b}\|^2$ **then**
 $[U, \Sigma, V] = \text{svd}(B' C'^T)$
 Update $B = U \sqrt{\Sigma}$ and $C = V \sqrt{\Sigma}$
 Decrease λ
 else
 Increase λ
 end
end

Algorithm 1: Outline of the bilinear method.

the energy it minimizes is not directly correlated to the others, but is shown for completeness. Furthermore, the iteration speed of ADMM is significantly faster than for VarPro, and therefore we show the elapsed time (in seconds) for all methods. The reported values are averaged over 50 instances with random initialization.

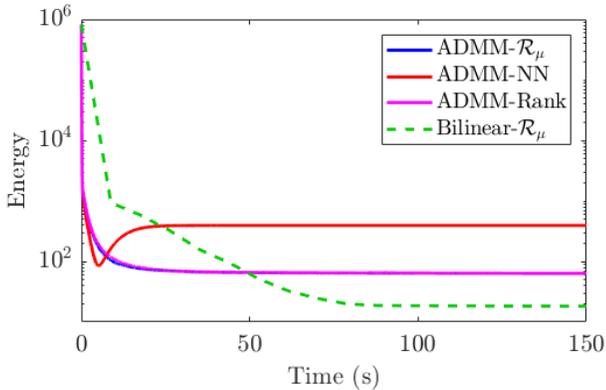


Figure 6. The average energy for the pOSE problem over 50 instances with random initializations, for test sequence *Door*. (Note that the energy for ADMM-Rank and ADMM- \mathcal{R}_μ are very similar).

C.2. Background Extraction

The missing data problem formulation can also be used in *e.g.* background extraction, where the goal is to separate

the foreground from the background in a video sequence. For this experiment, security footage of an airport is used. The frame size is 144×176 pixels, and we use the first 200 frames, as in [30]. The camera does not move, hence the background is static.

By concatenating the vectorization of the frames into a matrix we expect it to be additively decomposable in terms of a low rank matrix (background) and a sparse matrix (foreground). We follow the setup used in [8], and crop the width to half of the height, and shift it 20 pixels to the right after 100 frames to simulate a virtual pan of the camera. This increases the complexity of the background, as it is no longer static. Lastly, we randomly drop 70 % of the entries. To allow for smaller singular values, we use Geman, as it is a robust penalty with shrinking bias. The results are shown in Figure 8.

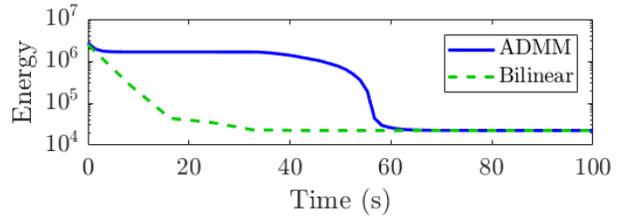


Figure 7. Energy minimization comparison for the background extraction experiment.

Initially ADMM struggles to find the correct balance between lowering the rank and fitting the data, which is seen in Figure 7, where the objective is almost unaffected the first forty seconds. At this point, the bilinear method has already converged.

C.3. Photometric Stereo

Photometric stereo can be used for estimating depth and surface orientation from images of the same object and view with varying lighting directions. Assuming M lighting directions and N pixels define $I \in \mathbb{R}^{M \times N}$, where I_{ij} is the light intensity for lighting direction i and pixel j . Assuming Lambertian reflectance, uniform albedo and a distant light source, $I = LN$, where $L \in \mathbb{R}^{M \times 3}$ contain the lighting directions and $N \in \mathbb{R}^{3 \times N}$ the unknown surface normals. Thus, the resulting problem is to find a rank 3 approximation of the intensity matrix I .

We use the Harvard Photometric Stereo testset [19], which contains images of various objects from varying lighting direction. The images are scaled to 160×125 pixels, and only the foreground pixels are used in the optimization. Similar to [8], we introduce missing data by thresholding dark pixels with pixel value less than 40 and bright pixels with pixel value more than 205. The measurement matrix is reconstructed using the bilinear method and the ADMM equivalent with the \mathcal{R}_μ regularization. The result



Figure 8. Background extraction using Geman. Samples from frame no. 40, 70, 100, 130, 170 and 200. *Top row*: Original images. *Middle row*: Training data with 70 % missing data. *Bottom row*: Reconstruction of background (bilinear method).

is shown in Figure 9. We let the bilinear method run until convergence and let the ADMM equivalent run for the same time in seconds, at which point the objective value is still decreasing when ADMM is interrupted; however, the reduction is almost negligible. In all cases ADMM fails to converge to a low rank solution in the same time as the bilinear method, which yields a consistent result.

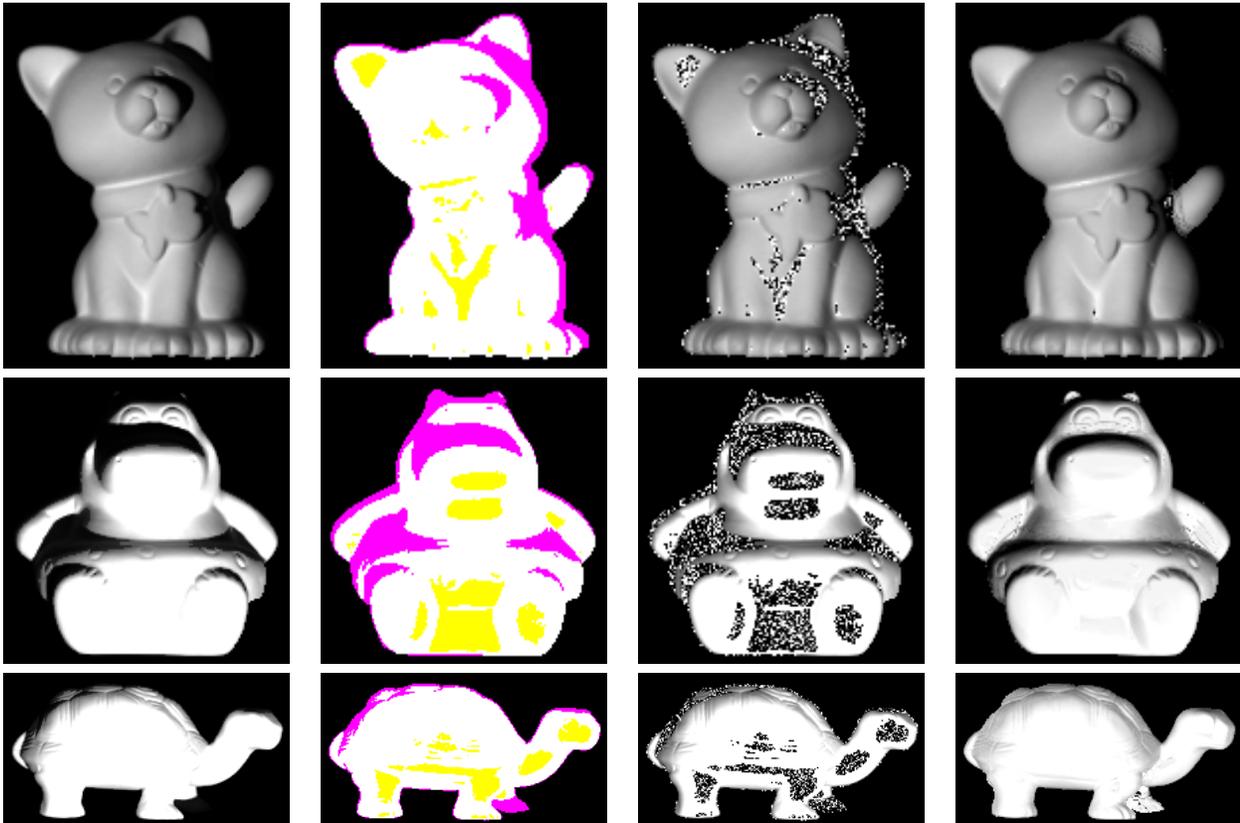


Figure 9. Images from the photometric stereo experiment. *From left to right:* (a) Ground truth image, (b) missing data mask with static background (black), dark pixels (purple), bright pixels (yellow), (c) reconstruction using ADMM, and (d) reconstruction using the Bilinear formulation.