

Exploring Algorithmic Limits of Matrix Rank Minimization Under Affine Constraints

Bo Xin, *Member, IEEE*, David Wipf, *Member, IEEE*, Yizhou Wang, *Member, IEEE*, and Wen Gao, *Fellow, IEEE*

Abstract—Many applications require recovering a matrix of minimal rank within an affine constraint set, with matrix completion a notable special case. Because the problem is NP-hard in general, it is common to replace the matrix rank with the nuclear norm, which acts as a convenient convex surrogate. While elegant theoretical conditions elucidate when this replacement is likely to be successful, they are highly restrictive and convex algorithms fail when the ambient rank is too high or when the constraint set is poorly structured. Nonconvex alternatives fare somewhat better when carefully tuned; however, convergence to locally optimal solutions remains a continuing source of failure. Against this backdrop, we derive a deceptively simple and parameter-free probabilistic PCA-like algorithm that is capable, over a wide battery of empirical tests, of successful recovery even at the theoretical limit where the number of measurements equals the degrees of freedom in the unknown low-rank matrix. Somewhat surprisingly, this is possible even when the affine constraint set is highly ill-conditioned. While proving general recovery guarantees remains evasive for nonconvex algorithms, Bayesian-inspired or otherwise, we nonetheless show conditions whereby the underlying cost function has a unique stationary point located at the global optimum; no existing cost function we are aware of satisfies this property. The algorithm has also been successfully deployed on a computer vision application involving image rectification and a standard collaborative filtering benchmark.

Index Terms—Rank minimization, affine constraints, matrix completion, matrix recovery, empirical Bayes.

I. INTRODUCTION

RECENTLY there has been a surge of interest in finding minimum rank matrices subject to some problem-specific constraints often characterized as an affine set [1]–[7]. Mathematically this involves solving

$$\min_{\mathbf{X}} \text{rank}[\mathbf{X}] \quad \text{s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}), \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times m}$ is the unknown matrix, $\mathbf{b} \in \mathbb{R}^p$ represents a vector of observations and $\mathcal{A} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^p$ denotes a linear mapping. An important special case of (1) commonly applied

to collaborative filtering is the matrix completion problem

$$\min_{\mathbf{X}} \text{rank}[\mathbf{X}] \quad \text{s.t. } \mathbf{X}_{ij} = (\mathbf{X}_0)_{ij}, (i, j) \in \Omega, \quad (2)$$

where \mathbf{X}_0 is a low-rank matrix we would like to recover, but we are only able to observe elements from the set Ω [1], [2]. Unfortunately however, both this special case and the general problem (1) are well-known to be NP-hard, and the rank penalty itself is non-smooth. Consequently, a popular alternative is to instead compute

$$\min_{\mathbf{X}} \sum_i f(\sigma_i[\mathbf{X}]) \quad \text{s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}), \quad (3)$$

where $\sigma_i[\mathbf{X}]$ denotes the i -th singular value of \mathbf{X} and f is usually a concave, non-decreasing function (or nearly so). In the special case where $f(z) = I[z \neq 0]$ (i.e., an indicator function) we retrieve the matrix rank; however, smoother surrogates such as $f(z) = \log z$ or $f(z) = z^q$ with $q \leq 1$ are generally preferred for optimization purposes. When $f(z) = z$, (3) reduces to convex nuclear norm minimization. A variety of celebrated theoretical results have quantified specific conditions, heavily dependent on the singular values of matrices in the nullspace of \mathcal{A} , where the minimum nuclear norm solution is guaranteed to coincide with that of minimal rank [1], [3], [6]. However, these guarantees typically only apply to a highly restrictive set of rank minimization problems, and in a practical setting non-convex algorithms can succeed in a much broader range of conditions [2], [5], [6].

In Section II we will summarize state-of-the-art non-convex rank minimization algorithms that operate under affine constraints and point out some of their shortcomings. This will be followed in Section III by the derivation of an alternative approach using Bayesian modeling techniques adapted from probabilistic PCA [8]. Section IV will then describe connections with nuclear norm minimization, convergence issues, and properties of global and local solutions. The latter includes special cases whereby any stationary point of the intrinsic cost function is guaranteed to have optimal rank, illustrating an underlying smoothing mechanism which leads to success over competing methods. We next discuss algorithmic enhancements in Section V that further improve recovery performance in practice. Section VI contains a wide variety of numerical comparisons that highlight the efficacy of this algorithm, while Section VII presents a computer vision application involving image rectification and a standard collaborative filtering benchmark. Technical proofs and algorithm update rule details are contained in the Appendix. Portions of this work have previously appeared in conference proceedings [9].

Manuscript received October 22, 2014; revised March 9, 2015 and November 23, 2015; accepted February 26, 2016. Date of publication April 7, 2016; date of current version. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Tareq Al-Naffouri. The authors would like to thank the support from the following grants: 973-2015CB351800, NSFC-61231010, NSFC-61527804, NSFC-61210005 and the Microsoft Research Asia Collaborative Research funding.

B. Xin, Y. Wang, and W. Gao are with the Department of Electrical Engineering and Computer Science, Peking University, Beijing 100871, China (e-mail: jimxinbo@gmail.com; yizhou.wang@pku.edu.cn; wgao@pku.edu.cn).

D. Wipf is with the Visual Computing group, Microsoft Research, Beijing 100080, China (e-mail: davidwip@microsoft.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2016.2551697

Before proceeding, we highlight several main contributions as follows:

- 1) Bayesian inspiration can take uncountably many different forms and parameterizations, but the devil is in the details and existing methods offer little opportunity for both theoretical inquiry and substantial performance gains solving (1). In this regard, we apply carefully-tailored modifications to a veteran probabilistic PCA model leading to systematic theoretical and empirical insights and advantages. Model justification is ultimately based on such meticulous technical considerations rather than merely the presumed qualitative legitimacy of any underlying prior distributions.
- 2) Non-convex algorithms have demonstrated some improvement in estimation accuracy over the celebrated convex nuclear norm; however, this typically requires the inclusion of one or more additional tuning parameters to incrementally inject additional objective function curvature and avoid bad local solutions. In contrast, for solving (1) our non-convex Bayesian-inspired algorithm requires no such parameters at all, and noisy relaxations necessitate only a single, standard trade-off parameter balancing data-fit and minimal rank.¹
- 3) Over a wide battery of controlled experiments with ground-truth data, our approach outperforms all existing algorithms that we are aware of, Bayesian, non-convex, or otherwise. This includes direct head-to-head comparisons using the exact experimental designs and code prepared by original authors. In fact, even when \mathcal{A} is ill-conditioned we are consistently able to solve (1) right up to the theoretical limit of any possible algorithm, which has never been demonstrated previously.

II. RELATED WORK

Here we focus on a few of the latest and most effective rank minimization algorithms, all developed within the last few years and evaluated favorably against the state-of-the-art.

A. General Non-Convex Methods

In the non-convex regime, effective optimization strategies attempt to at least locally minimize (3), often exceeding the performance of the convex nuclear norm. For example, [6] derives a family of *iterative reweighted least squares* (IRLS) algorithms applied to $f(z) = (z^2 + \gamma)^{q/2}$ with $q, \gamma > 0$ as tuning parameters. A related penalty also considered, which coincides with the limit as $q \rightarrow 0$ (up to an inconsequential scaling and translation), is $f(z) = \log(z^2 + \gamma)$, which maintains an intimate connection with rank given that

$$\log z = \lim_{q \rightarrow 0} q^{-1} (z^q - 1) \quad \text{and} \quad \lim_{q \rightarrow 0} z^q = I [z \neq 0], \quad (4)$$

where I is a standard indicator function. Consequently, when γ is small, $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ behaves much like a scaled

and translated version of the rank, albeit with nonzero gradients away from zero.

The IRLS0 algorithm from [6] represents the best-performing special case of the above, where $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ is minimized using a homotopy continuation scheme merged with IRLS. Here a fixed γ is replaced with a decreasing sequence $\{\gamma^k\}$, the rationale being that when γ^k is large, the cost function is relatively smooth and devoid of local minima. As the iterations k progress, γ^k is reduced, and the cost behaves more like the matrix rank function. However, because now we are more likely to be within a reasonably good basin of attraction, spurious local minima are more easily avoided. The downside of this procedure is that it requires a pre-defined heuristic for reducing γ^k , and this schedule may be problem specific. Moreover, there is no guarantee that a global solution will ever be found.

In a related vein, [5] derives a family of *iterative reweighted nuclear norm* (IRNN) algorithms that can be applied to virtually any concave non-decreasing function f , even when f is non-smooth, unlike IRLS. For effective performance however the authors suggest a continuation strategy similar to IRLS0. Moreover, additional tuning parameters are required for different classes of functions f and it remains unclear which choices are optimal. While the reported results are substantially better than when using the convex nuclear norm, in our experiments IRLS0 seems to perform slightly better, possibly because the quadratic least squares inner loop is less aggressive in the initial stages of optimization than weighted nuclear norm minimization, leading to a better overall trajectory. Regardless, all of these affine rank minimization algorithms fail well before the theoretical recovery limit is reached, when the number of observations p equals the number of degrees of freedom in the low-rank matrix we wish to recover. Specifically, for an $n \times m$, rank r matrix, the number of degrees of freedom is given by $r(m+n) - r^2$, hence $p = r(m+n) - r^2$ is the best-case boundary. In practice if \mathcal{A} is ill-conditioned or degenerate the achievable limit may be more modest.

A third approach relies on replacing the convex nuclear norm with a truncated non-convex surrogate [2]. While some competitive results for image inpainting via matrix completion are shown, in practice the proposed algorithm has many parameters to be tuned via cross-validation. Moreover, recent comparisons contained in [5] show that default settings perform relatively poorly.

Finally, a somewhat different class of non-convex algorithms can be derived using a straightforward application of alternating minimization [10]. The basic idea is to assume $\mathbf{X} = \mathbf{UV}^T$ for some low-rank matrices \mathbf{U} and \mathbf{V} and then solve

$$\min_{\mathbf{U}, \mathbf{V}} \|b - \mathcal{A}(\mathbf{UV}^T)\|_{\mathcal{F}} \quad (5)$$

via coordinate descent. The downside of this approach is that it can be sensitive to data correlations and requires that \mathbf{U} and \mathbf{V} be parameterized with the correct rank. In contrast, our emphasis here is on algorithms that require no prior knowledge whatsoever regarding the true rank. This is especially important in application extensions that may manage multiple low-rank

¹While not our emphasis here, similar to other Bayesian frameworks, even this trade-off parameter can ultimately be learned from the data if a true, parameter-free implementation is desired across noise levels.

183 matrices such that prior knowledge of all individual ranks is not
184 feasible.

185 B. Bayesian Methods

186 From a probabilistic perspective, previous work has applied
187 Bayesian formalisms to rank minimization problems, although
188 not specifically within an affine constraint set. For example,
189 [11]–[13] derive robust PCA algorithms built upon the linear
190 summation of a rank penalty and an element-wise sparsity
191 penalty. In particular, [12] applies an MCMC sampling approach
192 for posterior inference, but the resulting iterations are not scal-
193 able, subjectable to detailed analysis, nor readily adaptable to
194 affine constraints. In contrast, [11] applies a similar probabilis-
195 tic model but performs inference using a variational mean-field
196 approximation. While the special case of matrix completion
197 is considered, from an empirical standpoint its estimation accu-
198 racy is not competitive with the state-of-the-art non-convex
199 algorithms mentioned above. Finally, without the element-wise
200 sparsity component intrinsic to robust PCA (which is not our
201 focus here), [13] simply collapses to a regular PCA model with
202 a closed-form solution, so the challenges faced in solving (1) do
203 not apply. Consequently, general affine constraints really are a
204 key differentiating factor.

205 From a motivational angle, the basic probabilistic model with
206 which we begin our development can be interpreted as a care-
207 fully re-parameterized generalization of the probabilistic PCA
208 model from [8]. This will ultimately lead to a non-convex algo-
209 rithm devoid of the heuristic tuning strategies mentioned above,
210 but nonetheless still uniformly superior in terms of estimation
211 accuracy. We emphasize that, although we employ a Bayesian
212 entry point for our algorithmic strategy, final justification of the
213 underlying model will be entirely based on properties of the
214 underlying cost function that emerges, rather than any putative
215 belief in the actual validity of the assumed prior distributions
216 or likelihood function. This is quite unlike the vast majority of
217 existing Bayesian approaches.

218 C. Analytical Considerations

219 Turning to analytical issues, a number of celebrated theoret-
220 ical results dictate conditions whereby substitution of the rank
221 function with the convex nuclear norm in (1) is nonetheless guar-
222 anteed to still produce the minimal rank solution. For example,
223 if \mathcal{A} is a Gaussian iid measurement ensemble and $\mathbf{X}_0 \in \mathbb{R}^{n \times n}$
224 represents the optimal solution to (1) with $\text{rank}[\mathbf{X}_0] = r$, then
225 with high probability as the problem dimensions grow large, the
226 minimum nuclear norm feasible solution will equal \mathbf{X}_0 if the
227 number of measurements p satisfies $p \geq 3r(2n - r)$ [14].

228 The limitation of this type of result is two-fold. First, in the
229 above situation the true minimum rank solution only actually re-
230 quires $p \geq r(2n - r)$ measurements to be recoverable via brute
231 force solution of (1), and the remaining difference of a factor
232 of three can certainly be considerable in many practical situa-
233 tions (e.g., requiring 300 measurements is far more laborious
234 than only needing 100 measurements). Secondly though, and
235 far more importantly, all existing provable recovery guarantees
236 place extremely strong restrictions on the structure of \mathcal{A} , e.g.,

strong restrictions on the singular value decay of matrices in
the nullspace of \mathcal{A} . Such conditions are unlikely to ever hold in
realistic application settings, including the image rectification
example we describe in Section VII.A (in fact, these conditions
are usually incapable of even being checked). In contrast, the
algorithm we propose is empirically observed to only require
the theoretically minimal number of measurements even when
such nullspace conditions are violated in many cases. While a
general theoretical guarantee of this sort is obviously not poss-
ible, we do nonetheless provide several supporting theoretical
results indicative of why such performance is at least empirically
obtainable.

III. ALTERNATIVE ALGORITHM DERIVATION

In this section we first detail our basic distributional assump-
tions followed by development of the associated update rules
for inference.

A. Basic Model

In contrast to the majority of existing algorithms organized
around practical solutions to (3), here we adopt an alternative,
probabilistic starting point. We first define the Gaussian likeli-
hood function

$$p(\mathbf{b}|\mathbf{X}; \mathcal{A}, \lambda) \propto \exp\left[-\frac{1}{2\lambda} \|\mathcal{A}(\mathbf{X}) - \mathbf{b}\|_2^2\right], \quad (6)$$

noting that in the limit as $\lambda \rightarrow 0$ this will enforce the same
constraint set as in (1). Next we define an independent, zero-
mean Gaussian prior distribution with covariance $\nu_i \Psi$ on each
column of \mathbf{X} , denoted $\mathbf{x}_{:i}$ for all $i = 1, \dots, m$. This produces
the aggregate prior on \mathbf{X} given by

$$p(\mathbf{X}; \Psi, \boldsymbol{\nu}) = \prod_i \mathcal{N}(\mathbf{x}_{:i}; \mathbf{0}, \nu_i \Psi) \propto \exp\left[\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x}\right], \quad (7)$$

where $\Psi \in \mathbb{R}^{n \times n}$ is a positive semi-definite symmetric matrix,²
 $\boldsymbol{\nu} = [\nu_1, \dots, \nu_m]^\top$ is a non-negative vector, $\mathbf{x} = \text{vec}[\mathbf{X}]$
(column-wise vectorization), and $\overline{\Psi} = \text{diag}[\boldsymbol{\nu}] \otimes \Psi$, with \otimes
denoting the Kronecker product. It is important to stress here
that we do not necessarily believe that the unknown \mathbf{X} actually
follows such a Gaussian distribution per se. Rather, we adopt
(7) primarily because it will lead to an objective function with
desirable properties related to solving (1).

Moving forward, given both likelihood and prior are Gaus-
sian, the posterior $p(\mathbf{X}|\mathbf{b}; \Psi, \boldsymbol{\nu}, \mathcal{A}, \lambda)$ is also Gaussian, with
mean given by an $\widehat{\mathbf{X}}$ such that

$$\widehat{\mathbf{x}} = \text{vec}\left[\widehat{\mathbf{X}}\right] = \overline{\Psi} \mathbf{A}^\top (\lambda \mathbf{I} + \mathbf{A} \overline{\Psi} \mathbf{A}^\top)^{-1} \mathbf{b}. \quad (8)$$

²Technically Ψ must be positive definite for the inverse in (7) to be de-
fined. However, we can accommodate the semi-definite case using the fol-
lowing convention. Without loss of generality assume that $\overline{\Psi} = \mathbf{R} \mathbf{R}^\top$
for some matrix \mathbf{R} . We then qualify that $p(\mathbf{X}; \Psi, \boldsymbol{\nu}) = 0$ if $\mathbf{x} \notin \text{span}[\mathbf{R}]$,
and $p(\mathbf{X}; \Psi, \boldsymbol{\nu}) \propto \exp[-\frac{1}{2} \mathbf{x}^\top (\mathbf{R}^\top)^\dagger \mathbf{R} \mathbf{x}]$ otherwise. Equivalently, through-
out the paper for convenience (and with slight abuse of notation) we define
 $\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x} = \infty$ when $\mathbf{x} \notin \text{span}[\mathbf{R}]$, and $\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x} = \mathbf{x}^\top (\mathbf{R}^\top)^\dagger \mathbf{R} \mathbf{x}$ other-
wise. This will come in handy, for example, when interpreting the bound in
(12) below. Note also that the final cost function (10) we will ultimately be
minimizing requires no such inverse anyway.

274 Here $\mathbf{A} \in \mathbb{R}^{p \times nm}$ is a matrix defining the linear operator \mathcal{A}
 275 such that $\mathbf{b} = \mathbf{A}\mathbf{x}$ reproduces the feasible region in (1). From
 276 this expression it is clear that, if Ψ represents a low-rank co-
 277 variance matrix, then each column of $\widehat{\mathbf{X}}$ will be constrained
 278 to a low-dimensional subspace resulting overall in a low-rank
 279 estimate as desired. Of course for this simple strategy to be suc-
 280 cessful we require some way of determining a viable Ψ and the
 281 scaling vector ν .

282 A common Bayesian strategy in this regard is to marginalize
 283 over \mathbf{X} and then maximize the resulting likelihood function
 284 with respect to Ψ and ν [15], [13], [16]. This involves solving

$$\max_{\Psi \in H^+, \nu \geq 0} \int p(\mathbf{b}|\mathbf{X}; \mathcal{A}, \lambda) p(\mathbf{X}; \Psi, \nu) d\mathbf{X}, \quad (9)$$

285 where H^+ denotes the set of positive semi-definite and symmet-
 286 ric $n \times n$ matrices. After a -2 log transformation and applica-
 287 tion of a standard convolution-of-Gaussians integration, solving
 288 (9) is equivalent to minimizing the cost function

$$\mathcal{L}(\Psi, \nu) = \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} + \log |\Sigma_b|, \quad (10)$$

289 where

$$\Sigma_b = \mathbf{A} \bar{\Psi} \mathbf{A}^\top + \lambda \mathbf{I} \text{ and } \bar{\Psi} = \text{diag}[\nu] \otimes \Psi. \quad (11)$$

290 Here Σ_b is the covariance of \mathbf{b} given Ψ and ν .

291 B. Update Rules

292 Minimizing (10) is a non-convex optimization problem, and
 293 we employ standard upper bounds for this purpose leading to an
 294 EM-like algorithm, somewhat related to [8]. In particular, we
 295 compute separate bounds, parameterized by auxiliary variables,
 296 for both the first and second terms of $\mathcal{L}(\Psi, \nu)$. While the gen-
 297 eral case can easily be handled and may be applicable for more
 298 challenging problems, here for simplicity and ease of presenta-
 299 tion we consider minimizing $\mathcal{L}(\Psi) \triangleq \mathcal{L}(\Psi, \nu = \mathbf{1})$, meaning
 300 all elements of ν are fixed at one (and such is the case for all
 301 experiments reported herein, although we are currently explor-
 302 ing situations where this added generality could be especially
 303 helpful).

304 Based on [16], for the first term in (10) we have

$$\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + \mathbf{x}^\top \bar{\Psi}^{-1} \mathbf{x} \quad (12)$$

305 with equality whenever \mathbf{x} satisfies (8). For the second term we
 306 use

$$\begin{aligned} \log |\Sigma_b| &\equiv m \log |\Psi| + \log |\lambda \mathbf{A}^\top \mathbf{A} + \bar{\Psi}^{-1}| \\ &\leq m \log |\Psi| + \text{tr}[\Psi^{-1} \nabla_{\Psi^{-1}}] + C, \end{aligned} \quad (13)$$

307 where because $\log |\lambda \mathbf{A}^\top \mathbf{A} + \bar{\Psi}^{-1}|$ is concave with respect to
 308 Ψ^{-1} , we can upper bound it using a first-order approximation
 309 with a bias term C that is independent of Ψ . Equality is obtained
 310 when the gradient satisfies

$$\nabla_{\Psi^{-1}} = \sum_{i=1}^m \Psi - \Psi \mathbf{A}_i^\top (\mathbf{A} \bar{\Psi} \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_i \Psi, \quad (14)$$

311 where $\mathbf{A}_i \in \mathbb{R}^{p \times n}$ is defined such that $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_m]$.
 312 Finally given the upper bounds from (12) and (13) with \mathbf{X}

and $\nabla_{\Psi^{-1}}$ fixed, we can compute the optimal Ψ in closed form
 by optimizing the relevant Ψ -dependent terms via

$$\begin{aligned} \Psi^{\text{opt}} &= \arg \min_{\Psi} \text{tr}[\Psi^{-1} (\mathbf{X} \mathbf{X}^\top + \nabla_{\Psi^{-1}})] + m \log |\Psi| \\ &= \frac{1}{m} [\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi^{-1}}]. \end{aligned} \quad (15)$$

315 By agnostically starting with $\Psi = \mathbf{I}$ and then iteratively com-
 316 puting (8), (14), and (15), we can then obtain an estimate for Ψ ,
 317 and more importantly, a corresponding estimate for \mathbf{X} given by
 318 (8) at convergence. We refer to this basic procedure as BARM
 319 for *Bayesian Affine Rank Minimization*. The next section will
 320 describe in detail why it is particularly well-suited for solving
 321 problems such as (1).

322 IV. PROPERTIES OF BARM

323 Here we first describe a close but perhaps not intuitively-
 324 obvious relationship between the BARM objective function and
 325 canonical nuclear norm minimization. We then discuss desirable
 326 properties of global and local minima before concluding with a
 327 brief examination of convergence issues.

328 A. Connections with Nuclear Norm Minimization

329 On the surface, it may appear that minimizing (10) is com-
 330 pletely unrelated to the convex problem

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \text{ s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}) \quad (16)$$

331 that is most commonly associated with practical rank mini-
 332 mization implementations. However, a close connection can be
 333 revealed by considering the modified objective function

$$\mathcal{L}'(\Psi) = \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} + \text{tr}[\bar{\Psi}], \quad (17)$$

334 which represents nothing more than (10), with $\nu = \mathbf{1}$ and with
 335 $\log |\Sigma_b|$ being replaced by $\text{tr}[\bar{\Psi}]$. Now suppose we minimize
 336 (17) with respect to $\Psi \in H^+$ obtaining some Ψ^* . We then go
 337 on to compute an estimate of \mathbf{X} using (8). Note that if we apply
 338 the bound from (12) to the first term in (17), then this estimate
 339 for \mathbf{X} equivalently solves

$$\min_{\Psi \in H^+, \mathbf{X}} \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + \mathbf{x}^\top \bar{\Psi}^{-1} \mathbf{x} + \text{tr}[\bar{\Psi}], \quad (18)$$

340 with $\mathbf{x} = \text{vec}[\mathbf{X}]$ as before. If we first optimize over Ψ , it is eas-
 341 ily demonstrated that the optimal value of Ψ equals $(\mathbf{X} \mathbf{X}^\top)^{1/2}$.
 342 Plugging this value into (18), simplifying, and then applying the
 343 definition of the nuclear norm, we arrive at

$$\min_{\mathbf{X}} \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + 2\|\mathbf{X}\|_*, \quad (19)$$

344 Furthermore, in the limit $\lambda \rightarrow 0$ (applied outside of the
 345 minimization), (19) becomes equivalent to (16). For more
 346 information regarding the duality relationship between vari-
 347 ance/covariance space and coefficient space, at least in the
 348 related context of compressive sensing models, please refer
 349 to [16].

350 Consequently, we may conclude that the central distinc-
 351 tion between the proposed BARM cost function and nuclear
 352 norm minimization is an intrinsic \mathcal{A} -dependent penalty function

353 $\log |\Sigma_b|$ which is applied in covariance space. In Section IV.B
 354 we will examine desirable properties of this non-convex sub-
 355 stitution, highlighting our desire to treat the underlying BARM
 356 probabilistic model as an independent cost function that may be
 357 subject to technical analysis independent of its Bayesian origins.

358 B. Global/Local Minima Analysis

359 As discussed in Section II one nice property of the
 360 $\sum_i \log(\sigma_i[\mathbf{X}])$ penalty employed (approximately) by IRLS0
 361 [6] is that it can be viewed as a smooth version of the matrix
 362 rank function while still possessing the same set of minimum,
 363 both global and local, over the affine constraint set, at least if we
 364 consider the limiting situation of $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ when γ
 365 becomes small so that we may avoid the distracting singularity
 366 of $\log 0$. Additionally, it possesses an attractive form of scale
 367 invariance, meaning that if \mathbf{X}^* is an optimal feasible solution,
 368 a block-diagonal rescaling of \mathbf{A} nevertheless leads to an equiv-
 369 alent rescaling of the optimum (without the need for solving
 370 an additional optimization problem using the new \mathbf{A}). This is
 371 very much unlike the nuclear norm or other non-convex surro-
 372 gates that penalize the singular values of \mathbf{X} in a scale-dependent
 373 manner.

374 In contrast, the proposed algorithm is based on a very differ-
 375 ent Gaussian statistical model with seemingly a more tenuous
 376 connection with rank minimization. Encouragingly however,
 377 the proposed cost function enjoys the same global/local minima
 378 properties as $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ with $\gamma \rightarrow 0$. Before present-
 379 ing these results, we define $\text{spark}[\mathbf{A}]$ as the smallest number
 380 of linearly dependent columns in matrix \mathbf{A} [17]. All proofs are
 381 deferred to the Appendix.

382 *Lemma 1:* Let $\mathbf{b} = \text{Avec}[\mathbf{X}]$, where $\mathbf{A} \in \mathbb{R}^{p \times nm}$ satisfies
 383 $\text{spark}[\mathbf{A}] = p + 1$. Also define r as the smallest rank of any fea-
 384 sible solution. Then if $r < p/m$, any global minimizer $\{\Psi^*, \nu^*\}$
 385 of (10) in the limit $\lambda \rightarrow 0$ is such that $\mathbf{x}^* = \bar{\Psi}^* \mathbf{A}^\top (\mathbf{A} \bar{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$
 386 is feasible and $\text{rank}[\mathbf{X}^*] = r$ with $\text{vec}[\mathbf{X}^*] = \mathbf{x}^*$.

387 *Lemma 2:* Additionally, let $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{D}$, where $\mathbf{D} = \text{diag}$
 388 $[\alpha_1 \mathbf{\Gamma}, \dots, \alpha_m \mathbf{\Gamma}]$ is a block-diagonal matrix with invertible
 389 blocks $\mathbf{\Gamma} \in \mathbb{R}^{n \times n}$ of unit norm scaled with coefficients $\alpha_i > 0$.
 390 Then iff $\{\Psi^*, \nu^*\}$ is a minimizer (global or local) to (10) in
 391 the limit $\lambda \rightarrow 0$, then $\{\mathbf{\Gamma}^{-1} \Psi^*, \text{diag}[\alpha]^{-1} \nu^*\}$ is a minimizer when
 392 $\tilde{\mathbf{A}}$ replaces \mathbf{A} . The corresponding estimates of \mathbf{X} are likewise
 393 in one-to-one correspondence.

394 *Remarks:* The assumption $r = \text{rank}[\mathbf{X}^*] < p/m$ in Lemma
 395 1 is completely unrestrictive, especially given that a unique,
 396 minimal-rank solution is only theoretically possible by any algo-
 397 rithm if $p \geq (n + m)r - r^2$, which is much more restrictive
 398 than $p > rm$. Hence the bound we require is well above that
 399 required for uniqueness anyway. Likewise the spark assumption
 400 will be satisfied for any \mathbf{A} with even an infinitesimal (con-
 401 tinuous) random component. Consequently, we are essentially
 402 always guaranteed that BARM possesses the same global opti-
 403 mum as the rank function. Regarding Lemma 2, no surrogate
 404 rank penalty of the form $\sum_i f(\sigma_i[\mathbf{X}])$ can achieve this result
 405 except for $f(z) = \log z$, or inconsequential limiting translations
 406 and rescalings of the log such as the indicator function $I[z \neq 0]$
 407 (which is related to the log via arguments in Section II).

While these results are certainly a useful starting point, the
 real advantage of adopting the BARM cost function is that lo-
 cally minimizing solutions are exceedingly rare, largely as a
 consequence of the marginalization process in (9), and in some
 cases provably so. A specialized example of this smoothing can
 be quantified in the following scenario.

Suppose \mathbf{A} is now block diagonal, with diagonal blocks \mathbf{A}_i
 such that $\mathbf{b}_i = \mathbf{A}_i \mathbf{x}_i$ producing the aggregate observation vec-
 tor $\mathbf{b} = [\mathbf{b}_1^\top, \dots, \mathbf{b}_m^\top]^\top$. While somewhat restricted, this situa-
 tion nonetheless includes many important special cases, includ-
 ing canonical matrix completion and generalized matrix com-
 pletion where elements of $\mathbf{Z} = \mathbf{W}\mathbf{X}_0$ are observed after some
 transformation \mathbf{W} , instead of \mathbf{X}_0 directly.

Theorem 1: Let $\mathbf{b} = \text{Avec}[\mathbf{X}]$, where \mathbf{A} is block diagonal,
 with blocks $\mathbf{A}_i \in \mathbb{R}^{p_i \times n}$. Moreover, assume $p_i > 1$ for all i
 and that $\cap_i \text{null}[\mathbf{A}_i] = \emptyset$. Then if $\min_{\mathbf{X}} \text{rank}[\mathbf{X}] = 1$ in the
 feasible region, any minimizer $\{\Psi^*, \nu^*\}$ of (10) (global or local)
 in the limit $\lambda \rightarrow 0$ is such that $\mathbf{x}^* = \bar{\Psi}^* \mathbf{A}^\top (\mathbf{A} \bar{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$
 is feasible and $\text{rank}[\mathbf{X}^*] = 1$ with $\text{vec}[\mathbf{X}^*] = \mathbf{x}^*$. Furthermore,
 no cost function in the form of (3) can satisfy the same result.
 In particular, there can always exist local and/or global minima
 with rank greater than one.

Remarks: This result implies that, under extremely mild con-
 ditions, which do not even depend on the concentration proper-
 ties of \mathbf{A} , the proposed cost function has no minima that are not
 global minima, at least in this rank-one case. (The minor techni-
 cal condition regarding nullspace intersections merely ensures
 that high-rank components cannot simultaneously “hide” in the
 nullspace of every measurement matrix \mathbf{A}_i ; the actual \mathbf{A} opera-
 tor may still be highly ill-conditioned.) Thus any algorithm with
 provable convergence to some local minimizer is guaranteed to
 obtain a globally optimal solution.³

Although a global optimal guarantee for finding a rank-one
 matrix sounds somewhat limited, such a guarantee is not possi-
 ble with any other penalty function of the standard form
 $\sum_i f(\sigma_i[\mathbf{X}])$, which is the typical recipe for rank minimization
 algorithms, convex or not. Moreover, finding rank one matrices
 subject to affine constraints represents a crucial component of
 applications such as phase retrieval [18], [19].

Additionally, if a unique rank-one solution exists to (1), then
 the unique minimizing solution to (10) will produce this \mathbf{X} via
 (8). Crucially, this will occur even when the minimal number
 of measurements $p = n + m - 1$ are available, unlike any other
 algorithm we are aware of that is blind to the true underlying
 rank.⁴ Moreover, as evident from the experiments, the proposed
 algorithm always successfully finds the global optimal in many
 situations where the underlying matrix has a rank much higher
 than one. Therefore, although we can only provide theoretical
 guarantee for the rank-one case, the underlying intuition that
 local minima are smoothed away arguably carries over to situa-
 tions where the rank is greater than one.

³Note also that with minimal additional effort, it can be shown that no sub-
 optimal stationary points of any kind, including saddle points, are possible.

⁴It is important to emphasize that the difficulty of estimating the optimal low-
 rank solution is based on the ratio of the d.o.f. in \mathbf{X} to the number of observations
 p . Consequently, estimating \mathbf{X} even with r small can be challenging when p is
 also small, meaning \mathbf{A} is highly overcomplete.

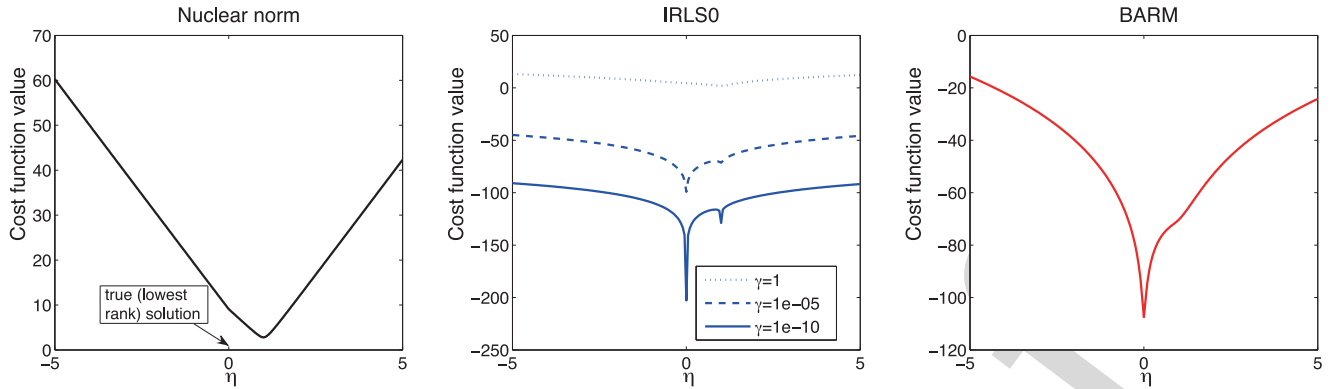


Fig. 1. Plots of different surrogates for matrix rank in a 1D feasible subspace. Here the convex nuclear norm does not retain the correct global minimum. In contrast, although the non-convex $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ penalty exhibits the correct minimum when γ is sufficiently small, it also contains spurious minima. Only BARM smooths away local minima while simultaneously retaining the correct global optima.

459 C. Visualization of BARM Local Minima Smoothing

460 To further explore the smoothing effect and complement The-
 461 orem 1, it helps to visualize rank penalty functions restricted to
 462 the feasible region. While the BARM algorithm involves mini-
 463 mizing (10), its implicit penalty function on \mathbf{X} can nonetheless
 464 be numerically obtained across the feasible region in a given
 465 subspace of interest; for other penalties such as the nuclear
 466 norm this is of course trivial. Practically it is convenient to ex-
 467 plore a 1D feasible subspace generated by $\mathbf{X}^* + \eta\mathbf{V}$, where
 468 \mathbf{X}^* is the true minimum rank solution, $\mathbf{V} \in \text{null}[\mathbf{A}]$, and η
 469 is a scalar. We may then plot various penalty function values
 470 as η is varied, tracing the corresponding 1D feasible subspace.
 471 We choose $\mathbf{V} = \mathbf{X}^1 - \mathbf{X}^*$, where \mathbf{X}^1 is a feasible solution
 472 with minimum nuclear norm; however, random selections from
 473 $\text{null}[\mathbf{A}]$ also show similar characteristics.

474 Fig. 1 provides a simple example of this process. \mathbf{A} is gener-
 475 ated randomly with all zeros and a single randomly placed
 476 ‘1’ in each row leading to a canonical matrix completion prob-
 477 lem. $\mathbf{X}^* \in \mathbb{R}^{5 \times 5}$ is randomly generated as $\mathbf{X}^* = \mathbf{u}\mathbf{v}^\top$, where
 478 \mathbf{u} and \mathbf{v} are iid $\mathcal{N}(0, 1)$ vectors, and so \mathbf{X}^* is rank one. Finally,
 479 $p = 10$ elements are observed, and therefore \mathbf{A} has 10 rows and
 480 $5 \times 5 = 25$ columns. η is varied from -5 to 5 and the values of
 481 the nuclear norm, $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$, and the implicit BARM
 482 cost function are displayed.

483 From the figure we observe that the minimum of the nuclear
 484 norm is not produced when the rank is smallest, which occurs
 485 when $\eta = 0$; hence the convex cost function fails for this prob-
 486 lem. Likewise, the $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ penalty used by IRLS0
 487 displays an incorrect global minimum when the tuning param-
 488 eter γ is large. In contrast, when γ is small, while the global
 489 minimum may now be correct, spurious local ditches have ap-
 490 peared in the cost function.⁵ Therefore, any success of the IRLS0
 491 algorithm depends heavily on a carefully balanced decaying se-
 492 quence of γ values, with the hope that initial iterations can steer
 493 the trajectory towards a desirable basin of attraction where local

⁵Technically speaking, these are not provably local minima since we are only considering a 1D subspace of the feasible region. However, it nonetheless illustrates the strong potential for troublesome local minima, especially in high dimensional practical problems.

minima are less problematic. One advantage of BARM then is
 that it is parameter free in this respect and yet still retains the
 correct global minimum, often without additional spurious local
 minima.

498 D. Convergence

499 Previous results of Section IV are limited to exploring aspects
 of the underlying BARM cost function. Regarding the BARM
 algorithm itself, by construction the updates generated by (8),
 (14), and (15) are guaranteed to reduce or leave unchanged
 $\mathcal{L}(\Psi)$ at each iteration. However, this is not technically suffi-
 cient to guarantee convergence to a stationary point of the cost
 function unless the additional conditions of Zangwill’s Global
 Convergence Theorem are satisfied [20]. However, provided we
 add a small regularization factor $\gamma \text{tr}[\Psi^{-1}]$, with $\gamma > 0$, then it
 can be shown that any cluster point of the resulting sequence of
 iterations $\{\Psi^k\}$ must be a stationary point. Moreover, because
 the sequence is bounded, there will always exist at least one
 cluster point, and therefore the algorithm is guaranteed to at
 least converge to a set of parameter values \mathcal{S} such that for any
 $\Psi^* \in \mathcal{S}$, $\mathcal{L}(\Psi^*) + \gamma \text{tr}[(\Psi^*)^{-1}]$ is a stationary point.

514 Finally, we should mention that this extra γ factor is akin to the
 homotopy continuation regularizer used by the IRLS0 algorithm
 [6] as discussed in Section II. However, whereas IRLS0 requires
 a carefully-chosen, decreasing sequence $\{\gamma^k\}$ with $\gamma^k > 0$ both
 to prove convergence and to avoid local minimum (and without
 this factor the algorithm performs very poorly in practice), for
 BARM a small, fixed factor only need be included as a technical
 necessity for proving formal convergence; in practice it can be
 fixed to exactly zero.

523 V. SYMMETRIZATION IMPROVEMENTS

524 Despite the promising theoretical attributes of BARM, there
 remains one important artifact of its probabilistic origins not
 found in more conventional existing rank minimization algo-
 rithms. In particular, other algorithms rely upon a symmetric
 penalty function that is independent of whether we are working
 with \mathbf{X} or \mathbf{X}^\top . All methods that reduce to (3) fall into this
 category, e.g., nuclear norm minimization, IRNN, or IRLS0. In

531 contrast, our method relies on defining a distribution with
 532 respect to the columns of \mathbf{X} . Consequently the underlying cost
 533 function is not identical when derived with respect to \mathbf{X} or
 534 \mathbf{X}^\top , a difference which will depend on \mathbf{A} . While globally opti-
 535 mal solutions should nonetheless be the same, the convergence
 536 trajectory could depend on this distinction leading to different
 537 local minima in certain circumstances. Although either con-
 538 struction leads to low-rank solutions, we may nonetheless expect
 539 improvement if we can somehow symmetrize the algorithm
 540 formulation.

541 To accomplish this, we consider a Gaussian prior on $\mathbf{x} =$
 542 $\text{vec}[\mathbf{X}]$ with a covariance formed using a block-wise averaging
 543 of covariances defined over rows and columns, denoted Ψ_r and
 544 Ψ_c respectively. The overall covariance is then given by the
 545 Kronecker sum

$$\overline{\Psi} = 1/2 (\Psi_r \otimes \mathbf{I} + \mathbf{I} \otimes \Psi_c). \quad (20)$$

546 The estimation process then proceeds in a similar fashion as
 547 before but with modifications and alternate upper-bounds that
 548 accommodate for this merger. For reported experimental results
 549 this symmetric version of BARM is used, with complete up-
 550 date rules listed in the Appendix and computational complexity
 551 evaluated in Section VI.E.

552 VI. EXPERIMENTAL VALIDATION

553 This section compares BARM with existing state-of-the-art
 554 affine rank minimization algorithms. For BARM, in all noise-
 555 less cases we simply used $\lambda = 10^{-10}$ (effectively zero), and
 556 hence no tuning parameters are required. Likewise, nuclear
 557 norm minimization [1], [4] requires no tuning parameters beyond
 558 implementation-dependent control parameters frequently
 559 used to enhance convergence speed (however the global mini-
 560 mum is unaltered given that the problem is convex). For the
 561 IRLS0 algorithm, we used our own implementation as the algo-
 562 rithm is straightforward and no code was available for the
 563 case of general \mathcal{A} ; we based the required decreasing γ_k se-
 564 quence on suggestions from [6]. IRLS0 code is available from
 565 the original authors for matrix completion; however, the results
 566 obtained with this code are not better than those obtained with
 567 our version. For the IRNN algorithm, we did not have access
 568 to code for general \mathcal{A} , nor specific details of how various pa-
 569 rameters should be set in the general case. Note also that IRNN
 570 has multiple parameters to tune even in noiseless problems un-
 571 like BARM. Therefore we report results directly from [5] where
 572 available. Note that both [5] and [6] show superior results to a
 573 number of other algorithms; we do not generally compare with
 574 these others given that they are likely no longer state-of-the-art
 575 and may clutter the presentation.

576 As stated previously, our focus here is on algorithms that do
 577 not require knowledge of the true rank of the optimal solution,
 578 and hence we do not include comparisons with [10] or the nor-
 579 malized hard thresholding algorithm from [21]. Regardless, we
 580 have nonetheless conducted numerous experiments with these
 581 algorithms, and even when the correct rank is provided, results
 582 are inferior to BARM, especially when correlated measurements
 583 are used. However, we do show limited empirical results with

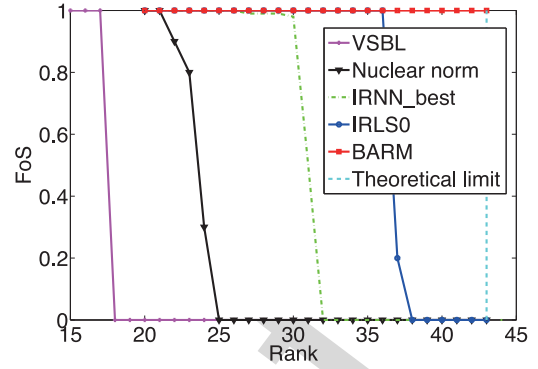


Fig. 2. Matrix completion comparisons (avg of 10 trials).

584 the variational sparse Bayesian algorithm (VSBL) from [11]
 585 because of its Bayesian origins, although the underlying param-
 586 eterization is decidedly different from BARM. But these results
 587 are limited to matrix completion as VSBL presently does not
 588 handle general affine constraints. Results from VSBL were ob-
 589 tained using publicly available code from the authors.

590 A. Matrix Completion

591 We begin with the matrix completion problem from (2), in
 592 part because this allows us to compare our results with the latest
 593 algorithms even when code is not available. For this purpose we
 594 reproduce the exact same experiment from [5], where a rank r
 595 matrix is generated as $\mathbf{X}_0 = \mathbf{T}\mathbf{M}_L\mathbf{T}\mathbf{M}_R$, with $\mathbf{T}\mathbf{M}_L \in \mathbb{R}^{n \times r}$
 596 and $\mathbf{T}\mathbf{M}_R \in \mathbb{R}^{r \times m}$ ($n = m = 150$) as iid $\mathcal{N}(0, 1)$ random ma-
 597 trices. 50% of all entries are then hidden uniformly at random.
 598 The *relative error* (REL) given by $\|\mathbf{X}_0 - \widehat{\mathbf{X}}\|_{\mathcal{F}} / \|\mathbf{X}_0\|_{\mathcal{F}}$
 599 is computed for each trial and averaged as r is varied. Likewise,
 600 we compute the *frequency of success* (FoS) score, which mea-
 601 sures the percentage of trials where the REL is below 10^{-3} .
 602 Results are shown in Fig. 2 where BARM is the only algorithm
 603 capable of reaching the theoretical recovery limit, beyond which
 604 $p = 0.5 \times 150^2 = 11250$ is surpassed by the number of degrees
 605 of freedom in \mathbf{X}_0 , in this case $2 \times 150 \times 44 - 44^2 = 11264$.
 606 Note that FoS values were reported in [5] over a wide range of
 607 non-convex IRNN algorithms. The green curve represents the
 608 best performing candidate from this pool as tuned by the original
 609 authors; REL values were unavailable. Interestingly, although
 610 VSBL is based on a somewhat related probabilistic model to
 611 BARM, the underlying parameterization, cost function, and up-
 612 date rules are entirely different and do not benefit from strong
 613 theoretical underpinnings. Hence performance does not always
 614 match recent state-of-the-art algorithms, although from a com-
 615 putational standpoint it is quite efficient.

616 Besides BARM, the IRLS0 algorithm also displayed better
 617 performance than the other methods. This motivated us to re-
 618 produce some of the matrix completion experiments from [6] so
 619 as to provide direct head-to-head comparisons with the authors'
 620 original implementation. For this purpose, \mathbf{X}_0 is conveniently
 621 generated in the same way as above; however, values of $n, m,$
 622 r , and the percentage of missing entries are varied while eval-
 623 uating reconstructions using FoS. While [6] tests a variety of

TABLE I
MATRIX COMPLETION RESULTS OF BARM WITH IRLS0 ON THE THREE
HARDEST PROBLEMS FROM [6]. PUBLISHED RESULTS IN [6] INCLUDED FOR
COMPARISON

Problem		IRLS0	IHT	FPCA	Opts	BARM
FR	n(=m)	r	FoS	FoS	FoS	FoS
0.78	500	20	0.9	0	0	1
0.8	40	9	1	0	0.5	1
0.87	100	14	0.5	0	0	1

624 combinations of these values to explore varying degrees of
625 problem difficulty, here we only reproduce the most challeng-
626 ing cases to see if BARM is still able to produce superior
627 reconstruction accuracy. In this respect problem difficulty is
628 measured by the *degrees of freedom ratio* (FR) given by FR
629 $= r(n + m - r)/p$ as defined in [6]. We also only include ex-
630 periments where algorithms are blind to the true rank of \mathbf{X}_0 .⁶
631 Results are shown in Table I, where we have also displayed
632 the published results of three additional algorithms that were
633 compared with IRLS0 in [6], namely, IHT [22], FPCA [23]
634 and Optspace [24]. From the table we observe that, in the most
635 difficult problem considered in [6], IRLS0 achieved only a 0.5
636 FoS score (meaning failure 50% of the time) while BARM still
637 achieves a perfect 1.0. Note that when FR is high, the problem
638 of recovering the underlying matrix is essentially much harder.
639 This happens in a manner that more local minima are induced
640 (due to increased rank) and/or much larger search space are
641 exposed (due to decreased number of observations/constraints).
642 In these cases, the equivalency of the global optimal with con-
643 vex relaxation usually does not hold, whereas for the existing
644 non-convex surrogates, there is no reason to assume any local
645 minima are not present. However, since BARM has an implicit
646 mechanism of smoothing local minima (though maybe not all
647 of them), it works more robustly in these situations.

648 B. General \mathbf{A}

649 Next we consider the more challenging problem involving
650 arbitrary affine constraints. The desired low-rank \mathbf{TX}_0 is gen-
651 erated in the same way as above. We then consider two types
652 of linear mappings where \mathbf{A} is generated as: (i) an iid $\mathcal{N}(0, 1)$,
653 $p \times n^2$ matrix, and (ii) $\sum_{i=1}^p i^{-1/2} \mathbf{u}_i \mathbf{v}_i^\top$, where $\mathbf{u}_i \in \mathbb{R}^p$ and
654 $\mathbf{v}_i \in \mathbb{R}^{n^2}$ are iid $\mathcal{N}(0, 1)$ vectors. The latter is meant to ex-
655 plore less-than-ideal conditions where the linear operator dis-
656 plays correlations and may be somewhat ill-conditioned. Fig. 3
657 displays aggregate results when \mathbf{X}_0 is 50×50 and 100×100 ,
658 including the underlying REL scores for additional comparison.
659 In both cases $p = 1000$ observations are used, and therefore the
660 corresponding measurement matrices \mathbf{A} are 1000×2500 and
661 1000×10000 respectively. We then vary r from 1 up to the
662 theoretical limit corresponding to problem size. Again we ob-
663 serve that BARM is consistently able to work up to the limit,
664 even when the \mathbf{A} operator is no longer an ideal Gaussian. In

⁶Note that IRLS0 can be modified to account for the true rank if such knowl-
edge were available.

665 general, we have explored a wide range of empirical conditions
666 too lengthy to report here, and it is only very rarely, and always
667 near the theoretical boundary, where BARM occasionally may
668 not succeed. We explore such failure cases in the next section.

669 C. Failure Case Analysis

670 Thus far we have not shown any cases where BARM actually
671 fails. Of course solving (1) for general \mathbf{A} is NP-hard so recovery
672 failures certainly must exist in some circumstances when using
673 a polynomial-time algorithm such as BARM. Although we cer-
674 tainly cannot explore every possible scenario, it behooves us
675 to probe more carefully for conditions under which such errors
676 may occur. One way to accomplish this is to push the problem
677 difficulty even further towards the theoretical limit by reducing
678 the number of measurements p as follows.

679 With the number of observations fixed at $p = 1000$ and a
680 general measurement matrix \mathbf{A} , the previous section examined
681 the recovery of 50×50 and 100×100 matrices as the rank was
682 varied from 1 to the recovery limit ($r = 11$ for the 50×50 case;
683 $r = 5$ for the 100×100 case). However, it is still possible to
684 make the problem even more challenging by fixing r at the limit
685 and then reducing p until it exactly equals the degrees of freedom
686 $2n^2 - r^2$. With $\{n = 50, r = 11\}$ this occurs at $p = 979$, for
687 $\{n = 100, r = 5\}$ this occurs at $p = 975$.

688 We examined the BARM algorithm under these conditions
689 with 10 additional trials using the uncorrelated \mathbf{A} for each prob-
690 lem size. Encouragingly, BARM was still 30% successful with
691 $\{n = 50, r = 11\}$, and 40% successful with $\{n = 100, r = 5\}$.
692 However, it is interesting to further examine the nature of these
693 failure cases. In Fig. 4 we have averaged the singular values of
694 $\widehat{\mathbf{X}}$ in all the failure cases. We notice that, although the recovery
695 was technically classified as a failure since the relative error
696 (REL) was above the stated threshold, the estimated matrices
697 are of almost exactly the correct minimal rank. Hence BARM
698 has essentially uncovered an alternative solution with minimal
699 rank that is nonetheless feasible by construction. We therefore
700 speculate that right at the theoretical limit, when \mathbf{A} is maxi-
701 mally overcomplete ($p \times n^2 = 979 \times 2500$ or 975×10000 for
702 the two problem sizes), there exists multiple feasible matri-
703 ces with singular value spectral cut-off points indistinguishable
704 from the optimal solution. Importantly, when the other algo-
705 rithms we tested failed, the failure is much more dramatic and
706 a clear spectral cut-off at the correct rank is not apparent.

707 This motivates a looser success criteria than FoS to account
708 for the possibility of multiple (nearly) optimal solutions that
709 may not necessarily be close with respect to relative error. For
710 this purpose we define the *frequency of rank success* (FoRS) as
711 the percentage of trials whereby a feasible solution $\widehat{\mathbf{X}}$ is found
712 such that $\sigma_r[\widehat{\mathbf{X}}]/\sigma_{r+1}[\widehat{\mathbf{X}}] > 10^3$, where $\sigma_i[\cdot]$ denotes the i -th
713 singular value of a matrix and r is the rank of the true low-rank
714 \mathbf{X}_0 . In words, FoRS measures the percentage of trials such that
715 roughly a rank r solution is recovered, regardless of proximity
716 to \mathbf{X}_0 .

717 Under this new criteria, all of the failure cases with respect to
718 FoS described above, for both problem sizes, become successes;
719 however, none of the other algorithms show improvement under

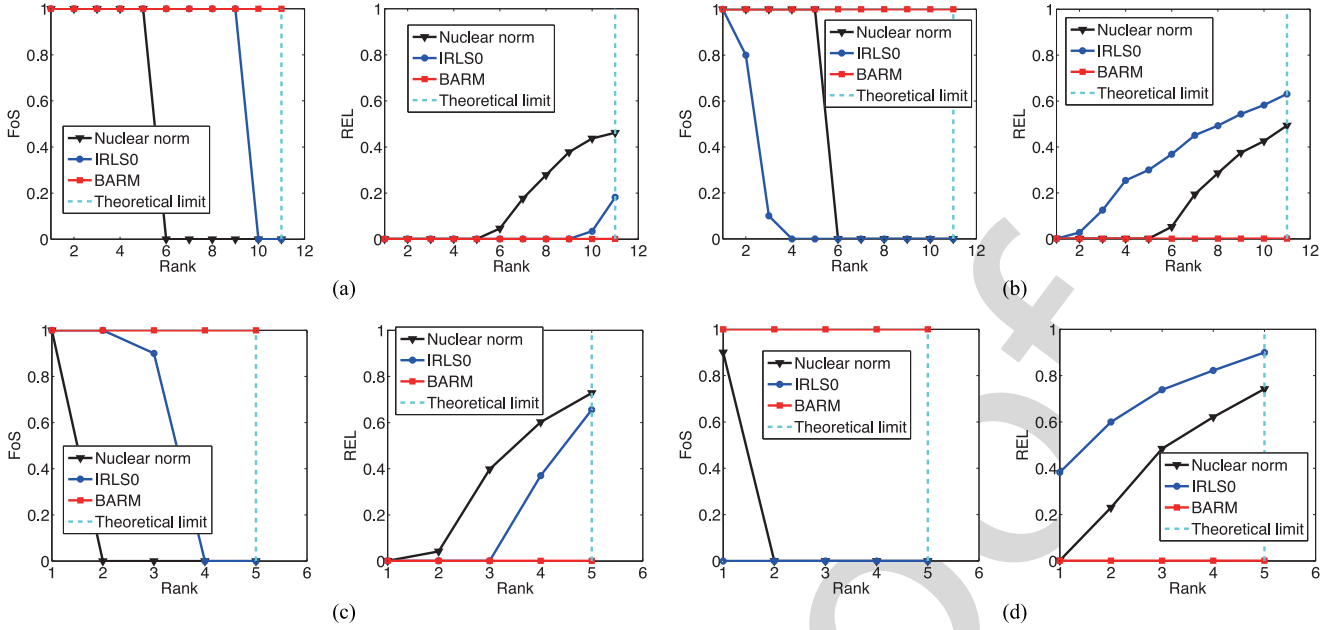


Fig. 3. Comparisons with general affine constraints (avg of 10 trials). (a) 50×50 , \mathbf{A} uncorrelated, (b) 50×50 , \mathbf{A} correlated, (c) 100×100 , \mathbf{A} uncorrelated, and (d) 100×100 , \mathbf{A} correlated.

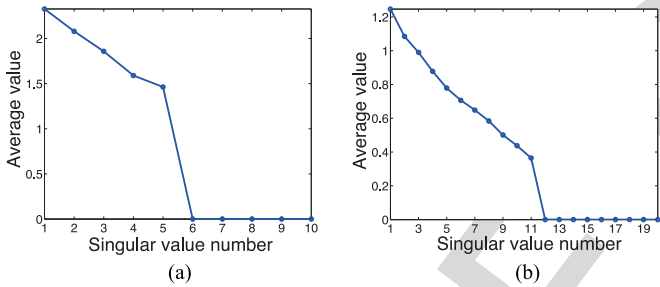


Fig. 4. Singular value averages of failure cases. In both cases solutions of minimal rank are obtained even though $\hat{\mathbf{X}} \neq \mathbf{X}_0$. (a) 50×50 and (b) 100×100 .

TABLE II
FURTHER MATRIX COMPLETION COMPARISONS OF BARM WITH IRLS0 BY REDUCING THE NUMBER OF MEASUREMENTS IN THE HARDEST PROBLEM FROM [6]. RESULTS WITH BOTH FoS AND FoRS METRICS ARE REPORTED (AVG OF 10 TRIALS)

Problem		IRLS0		BARM		
FR	$n(=m)$	r	FoS	FoRS	FoS	FoRS
0.9	100	14	0	0	1	1
0.95	100	14	0	0	0.8	1
0.99	100	14	0	0	0.7	1

720 this criteria, indicating that their original failures involved actual
 721 sub-optimal rank solutions. Something similar happens when we
 722 revisit the matrix completion experiments. For example, based
 723 on Table I the most difficult case involves $FR = 0.87$; however,
 724 by further reducing p , we can push FR towards 1.0 to further
 725 investigate the break-down point of BARM. Results are shown
 726 in Table II. While IRLS0 (which is the top performing algorithm

727 in [6] and in our experiments besides BARM) fails 100% of the
 728 time via both metrics, BARM can achieve an FoS of 0.7 even
 729 when $FR = 0.99$ and an FoRS of 1.0 in all cases.

730 We therefore adopt a more challenging measurement structure
 731 for \mathbf{A} to better evaluate the limits of BARM performance to
 732 reveal potential failures by both FoS and FoRS metrics. Specif-
 733 ically, we first applied 2-D *discrete cosine transform* (DCT) to
 734 \mathbf{X}_0 and then randomly sampled p of the resulting DCT coef-
 735 ficients. Because both the DCT and the sampling sub-process
 736 are linear operations on the entries of \mathbf{X}_0 , the whole process is
 737 representable via a matrix \mathbf{A} , which encodes highly structured
 738 information. Fig. 5 depicts the results using problem sizes con-
 739 sistent with Fig. 3; note that the FoRS metric has replaced the
 740 REL metric for comparison purposes.

741 Two things stand out from the analysis. First, while the other
 742 algorithms display almost identical behavior under either metric,
 743 BARM failures under the FoS criteria are mostly converted to
 744 successes by the FoRS metric by recovering a matrix of near-
 745 optimal rank. Secondly, even though certain unequivocal fail-
 746 ures emerge near the limits with this challenging DCT-based
 747 sampling matrix, BARM outperforms the other algorithms using
 748 either metric by a large margin.

749 To summarize, we have demonstrated that BARM is capa-
 750 ble of recovering a low-rank matrix right up to the theoretical
 751 limit in a variety of scenarios using different types of mea-
 752 surement processes. Moreover, even in cases where it fails, it
 753 often nonetheless still produces a feasible $\hat{\mathbf{X}}$ with rank nearly
 754 identical to the generative low-rank \mathbf{X}_0 , suggesting that multi-
 755 ple optimal solutions may be possible in challenging borderline
 756 cases. But when true unequivocal failures do occur, such fail-
 757 ures tend to be near the theoretical boundary, and with greater
 758 likelihood when the dictionary displays significant structure
 759 (or correlations). While certainly we envision that, out of the

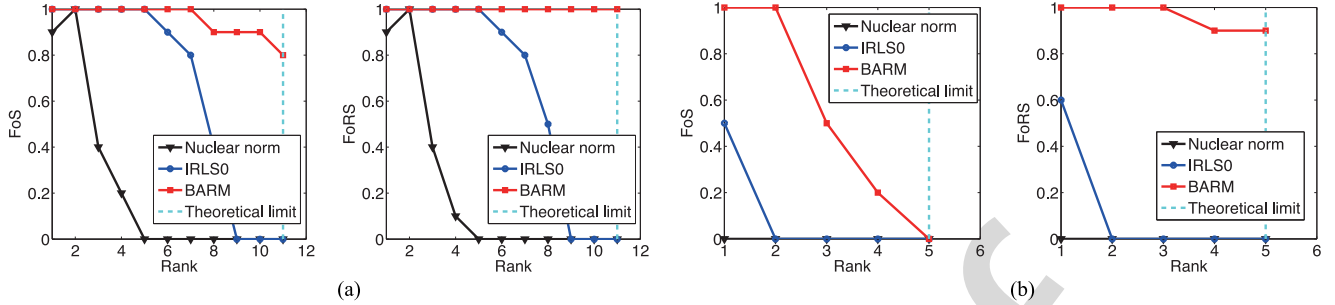


Fig. 5. Comparisons with structured affine constraints using both FoS and FoFS evaluation metrics (avg of 10 trials). (a) 50×50 , \mathbf{A} sub-sampled DCT, (b) 100×100 , \mathbf{A} sub-sampled DCT.

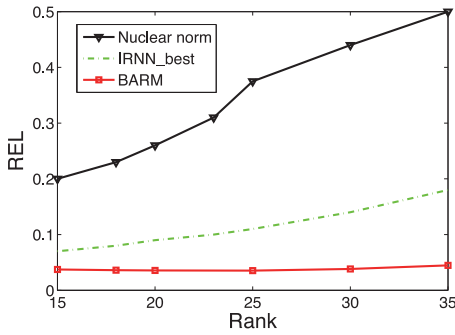


Fig. 6. Test with noisy data.

760 infinite multitude of testing situations further significant pock-
 761 ets of BARM failure can be revealed, we nonetheless feel that
 762 BARM is quite promising relative to existing algorithms.

763 D. Additional Noisy Tests

764 We also briefly present results that demonstrate the robustness
 765 of BARM to noise. For this purpose we reproduce the noisy
 766 experiment from [5] designed for validating IRNN algorithms.
 767 The simulated data are generated in the exact same way as was
 768 used to produce Fig. 2, only now instead of observing elements
 769 of \mathbf{X}_0 directly, we observe $\mathbf{X}_0 + 0.1 \times \mathbf{E}$, where elements
 770 of \mathbf{E} are iid $\mathcal{N}(0, 1)$. Although in [5] a heuristic strategy is
 771 introduced and tuned for adaptively setting all parameters (four
 772 in total), we simply applied BARM with $\lambda = 10^{-3}$ (so only a
 773 single parameter need be adjusted, and actually a wide range
 774 of λ values produces similar performance anyway). Results are
 775 shown in Fig. 6 where we compare BARM directly with the best
 776 result reported in [5] over the range $r = 15$ to $r = 35$. The
 777 nuclear norm solution is also included for reference. Overall, the
 778 BARM solution is stable and exhibits superior accuracy relative
 779 to the others.

780 E. Computational Complexity

781 Finally, regarding computational complexity, for general \mathbf{A}
 782 the BARM updates can be implemented to scale linearly in the
 783 elements of \mathbf{X} and quadratically in the number of observations
 784 p (the special case of matrix completion is decidedly much
 785 cheaper because of the special structure that can be exploited).
 786 In our experiments, for relatively easy problems on the order of

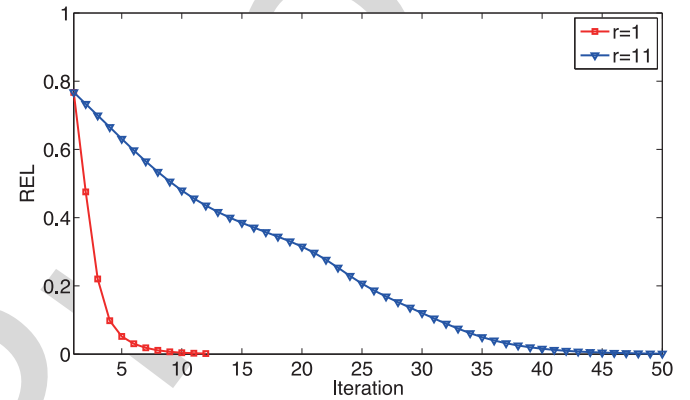


Fig. 7. Empirical convergence of BARM.

787 10 iterations are required, while for difficult recovery problems
 788 near the theoretical recovery boundary this may increase by a
 789 factor of 10 or so. This is somewhat expected though since as we
 790 near the theoretical limit, \mathbf{A} becomes highly overcomplete, and
 791 candidate solutions become much more difficult to differentiate.

792 To show this effect empirically, we compare two separate trials
 793 from Fig. 3(a), the first when $r = 1$ (relatively easy), the sec-
 794 ond when $r = 11$ (relatively hard).⁷ In Fig. 7 we plot the value
 795 of REL in both cases versus the iteration number of BARM.

796 VII. APPLICATION EXAMPLES

797 Many real-world problems from disparate fields can be for-
 798 mulated as the search for a low-rank matrix under affine con-
 799 straints [1], [3], [4], [25]. Here we briefly consider two such
 800 examples: low-rank image rectification and collaborative filter-
 801 ing for recommender systems. The former implicitly involves
 802 a general sampling operator \mathbf{A} , while the latter reduces to a
 803 standard matrix completion problem.

804 A. Low-Rank Image Rectification

805 In [4], the *transform invariant low-rank textures* (TILT) al-
 806 gorithm is derived for rectifying images containing low-rank

⁷Note that $r = 1$ is only relatively easy here because the number of obser-
 vations is sufficient for the larger $r = 11$ case; if only the minimal number
 of measurements are available then even $r = 1$ can be challenging for many
 algorithms.

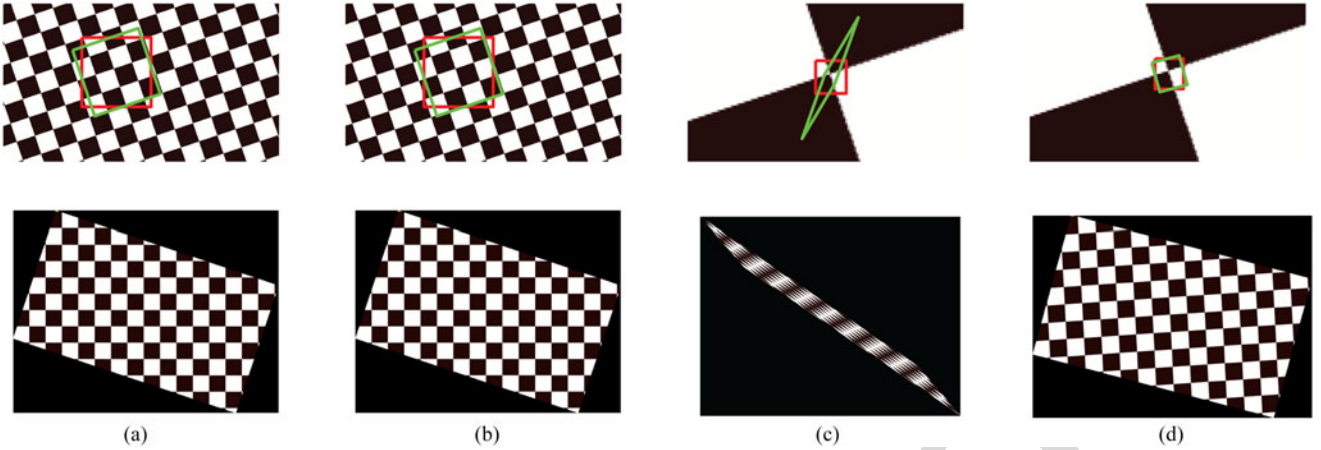


Fig. 8. Image rectification comparisons using a checkboard image. *Top*: Original image with observed region (red box) and estimated transformation (green box). *Bottom*: Rectified image estimates. (a) Nuclear norm (easy), (b) BARM (easy), (c) Nuclear norm (hard), (d) BARM (hard).

807 textures that have been transformed using an unknown operator
 808 τ from some group (e.g., a homography). For a given observed
 809 image \mathbf{Y} , the basic idea is to construct a first-order Taylor series
 810 approximation around the current rectified image estimate $\widehat{\mathbf{X}}$
 811 and solve

$$\min_{\mathbf{X}, \delta} \text{rank}[\mathbf{X}] \text{ s.t. } \mathbf{X} = \mathbf{Y} + \sum_i \mathbf{J}_i(\widehat{\mathbf{X}}) \delta_i, \quad (21)$$

812 where $\mathbf{J}_i(\widehat{\mathbf{X}})$ is the Jacobian matrix with respect to \mathbf{X} of
 813 the i -th parameter τ_i describing the transformation, with $\tau =$
 814 $[\tau_1, \tau_2, \dots]^\top$. Optimization over the vector of first-order differ-
 815 ences $\delta = [\delta_1, \delta_2, \dots]^\top$ can be accomplished in closed form by
 816 projecting both sides of the constraint to the orthogonal comple-
 817 ment of the span of all $\mathbf{J}_i(\widehat{\mathbf{X}})$. Let P_{J^c} represent this projection
 818 operator. The feasible region in (21) then becomes

$$P_{J^c}(\mathbf{X}) = P_{J^c}(\mathbf{Y}) + P_{J^c} \left(\sum_i \mathbf{J}_i(\widehat{\mathbf{X}}) \delta_i \right) = P_{J^c}(\mathbf{Y}) \quad (22)$$

819 The resulting problem then reduces exactly to (1) when we
 820 define $\mathcal{A} = P_{J^c}$ and $\mathbf{b} = \text{vec}[P_{J^c}(\mathbf{Y})]$. Once \mathbf{X} is computed in
 821 this way, we then update each $\mathbf{J}_i(\widehat{\mathbf{X}})$ and repeat until conver-
 822 gence.

823 While the original TILT algorithm substitutes the nuclear
 824 norm for $\text{rank}[\mathbf{X}]$, we embedded the BARM algorithm into
 825 the posted TILT source code [4] for comparison purposes (note
 826 that we disabled an additional sparse error term for both algo-
 827 rithms to simplify comparisons, and it is not necessary anyway
 828 in many regimes). Figs. 8 and 9 display results on both two
 829 easy examples, where the number of observations p is large,
 830 and two more difficult problems where the number observa-
 831 tions is small. While both algorithms succeed on the easy cases,
 832 when the observations are constrained by a small image window,
 833 only BARM is successful in accurately rectifying the images.
 834 This may be due, at least in part, to the fact that the implicit
 835 \mathcal{A} operator contains significant structure that is not consistent
 836 with the required nullspace properties required for nuclear norm
 837 minimization success.

B. Collaborative Filtering of MovieLens Data

838

839 Collaborative filtering, a technique used by many recom- 839
 840 mender systems, is a popular representative application of low- 840
 841 rank matrix completion. Typically the rows (or columns) of \mathbf{X}_0 841
 842 index users, the columns (or rows) denote items, and each entry 842
 843 $(\mathbf{X}_0)_{ij}$ is the rating/score of user i applied to item j . Given 843
 844 that we can observe some subset of elements of \mathbf{X}_0 , the task 844
 845 of collaborative filtering is to predict all or some of the miss- 845
 846 ing ratings. In general this would be impossible; however, if we 846
 847 have access to some prior knowledge, e.g., \mathbf{X}_0 is low-rank, then 847
 848 estimation may be feasible. 848

849 While our interest here is not in recommender systems or 849
 850 collaborative filtering per se, we nonetheless evaluate BARM 850
 851 using the 1M MovieLens dataset⁸ as this appears to represent 851
 852 one of the most common evaluation benchmarks. We emphasize 852
 853 at the outset that the strict validity of any low-rank assumptions 853
 854 underlying this data is debatable, and it remains entirely unclear 854
 855 whether the true globally optimal or lowest rank solution consis- 855
 856 tent with the observations, even if computable, would necessar- 856
 857 ily lead to the best prediction of the unknown ratings. In fact, the 857
 858 reported performance of various existing rank-minimization algo- 858
 859 rithms tends to cluster around almost the same value, implying 859
 860 that collaborative filtering may not provide the most discrimina- 860
 861 tive data type with which to compare. In most cases, it appears 861
 862 that tuning parameters and other heuristic modifications play 862
 863 a larger role than the underlying algorithmic distinctions funda- 863
 864 mental to finding optimal low-rank estimates. Nonetheless, 864
 865 we apply BARM for completeness and convention, adopting an 865
 866 additional simple mean-offset estimation term from [25] that is 866
 867 particularly suitable for this problem. 867

868 In [6], IRLS0 is compared with only two other algorithms on 868
 869 MovieLens data, but the performance is no better. Therefore, 869
 870 we choose to compare directly with [25], which both derives 870
 871 an IRLS-like algorithm and shows comparisons with a much 871
 872 wider variety of alternative algorithms using a strict evalua- 872
 873 tion protocol that is standard in the literature. Specifically, the 873

⁸<http://www.grouplens.org/>

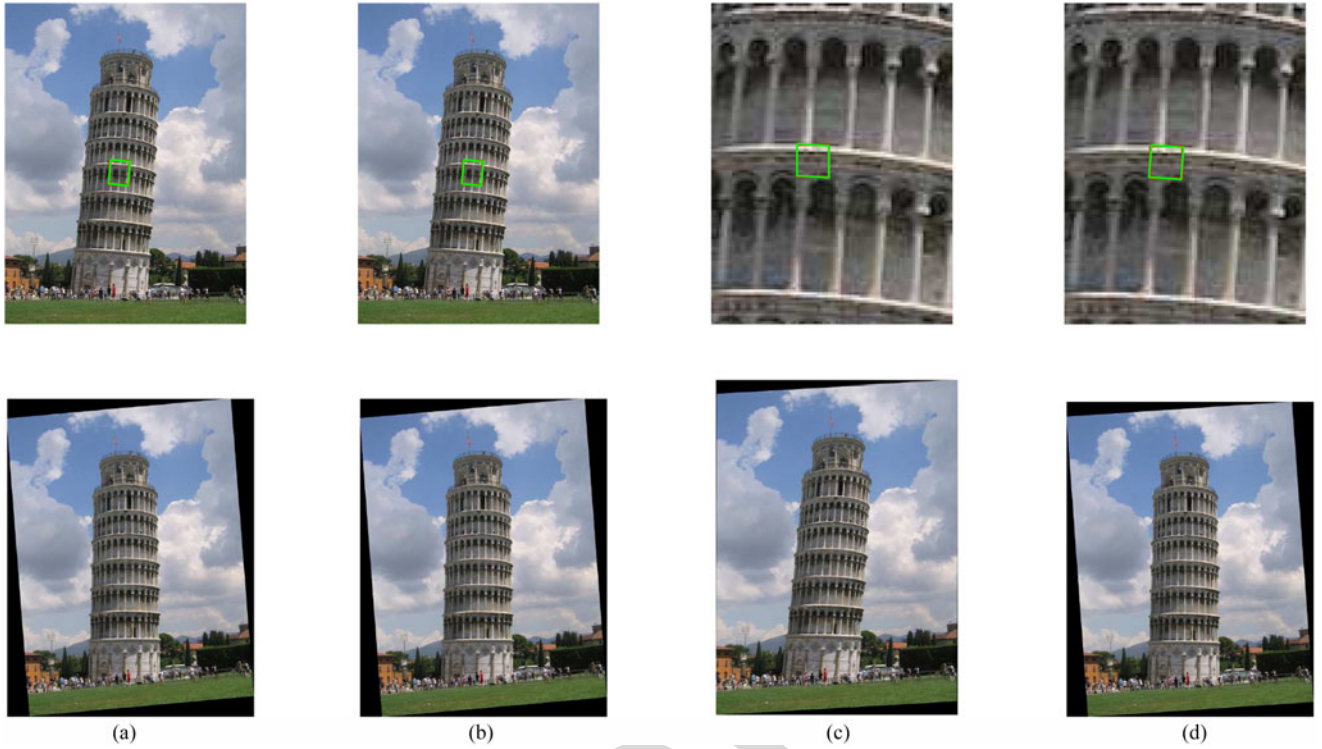


Fig. 9. Image rectification comparisons using a landmark photo. *Top*: Original image with observed region (red box) and estimated transformation (green box). *Bottom*: Rectified image estimates. (a) Nuclear norm (easy), (b) BARM (easy), (c) Nuclear norm (hard), (d) BARM (hard).

874 1M MovieLens dataset, which contains 1 million ratings in the
 875 range $\{1, \dots, 5\}$ for 3900 movies from 6040 unique users, is
 876 assessed under two test-protocols: *weak generalization*, which
 877 measures the ability to predict other items rated by the same
 878 user, and *strong generalization*, which measures the ability to
 879 predict items by novel users. 5 000 users are randomly selected
 880 for the weak generalization, and likewise 1 000 users are ex-
 881 tracted for the strong generalization. Each experiment is then
 882 run three times and the averaged results are reported. The per-
 883 formance metric is *normalized mean absolute error* (NMAE)
 884 given as

$$\text{NMAE} = \frac{\left(\sum_{i,j \in \text{supp}(\mathbf{X}_0)} \frac{|(\mathbf{X}_0)_{ij} - \hat{\mathbf{X}}_{ij}|}{|\text{supp}(\mathbf{X}_0)|} \right)}{(rt_{\max} - rt_{\min})},$$

885 where rt_{\max} and rt_{\min} are the maximum and minimum ratings
 886 possible.

887 We followed the same setup and reported results using BARM
 888 in Table III along with results from [25] for comparison. This
 889 includes the additional algorithms URP [26], Attitude [27],
 890 MMMF [28], IPCF [29], E-MMMF [30], GPLVM [31], NBMC
 891 [32], and IRLS/GM [25], [6]. From this table we observe that
 892 for the easier weak generalization problem BARM is a close
 893 second best, while for the more challenging strong generaliza-
 894 tion BARM is actually the best. Of course it is also immediately
 895 apparent that all algorithms fall within a relatively narrow per-
 896 formance range of approximately five percentage points. Con-
 897 sequently, we cannot unequivocally conclude that the attributes
 898 of BARM which make it suitable for optimally minimizing rank

TABLE III
 COLLABORATIVE FILTERING ON 1M MOVIELENS DATASET. RESULTS FROM
 [25] ARE IN ITALIC FOR COMPARISON PURPOSES

	Weak NMAE	Hard NMAE
<i>URP</i>	0.4341	0.4444
<i>Attitude</i>	0.4320	0.4375
<i>MMMF</i>	0.4156	0.4203
<i>IPCF</i>	0.4096	0.4113
<i>E-MMMF</i>	0.4029	0.4071
<i>GPLVM</i>	0.4026	0.3994
<i>NBMC</i>	0.3916	0.3992
<i>IRLS/GM</i>	0.3959	0.3928
BARM	0.3942	0.3898

899 necessarily translate into a truly significant practical advantage
 900 on this collaborative filtering task. But we would argue that the
 901 same holds for any matrix completion algorithm.

VIII. CONCLUSION

902 This paper explores a conceptually-simple, parameter-free
 903 algorithm called BARM for matrix rank minimization under
 904 affine constraints that is capable of successful recovery empir-
 905 ically observed to approach the theoretical limit over a broad
 906 class of experimental settings (including many not shown here)
 907 unlike any existing algorithms, and long after any convex guar-
 908 antees break down. Our strategy in this effort has been to
 909 adopt Bayesian machinery for inspiring a principled cost func-
 910 tion; however, ultimate model justification is placed entirely in
 911

912 theoretical evaluation of desirable global and local minima prop-
 913 erties, and in the empirical recovery performance that inevitably
 914 results from these properties. Although in general non-convex
 915 algorithms are exponentially more challenging to analyze, in
 916 this regard we have at least attempted to contextualize BARM
 917 in the same manner as convex optimization-based approaches
 918 such as nuclear-norm minimization.

919 APPENDIX A

920 Here we provide brief proofs of Lemmas 1 and 2 as well as
 921 Theorem 1. We also address the augmented update rules that
 922 account for the revised, symmetrized cost function discussed in
 923 Section V.

924 A. Proof of Lemmas 1 and 2

925 Regarding Lemma 1, this result mirrors related ideas from
 926 [16] in the context of Bayesian compressive sensing. Hence,
 927 while a more rigorous presentation is possible, here we de-
 928 scribe the basic aspects of the adaptation. At any candidate
 929 minimizer of (10) in the limit $\lambda \rightarrow 0$, define \mathbf{W} such that
 930 $\mathbf{A}\bar{\Psi}\mathbf{A}^\top = \mathbf{W}\mathbf{W}^\top$. To be a minimizer, global or local, it must
 931 be that $\mathbf{b} \in \text{span}[\mathbf{W}]$. If this were not the case, then $\mathcal{L}(\Psi, \nu)$
 932 would diverge to infinity as $\lambda \rightarrow 0$ because $\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b}$ progresses
 933 to infinity at a faster rate than $\log |\Sigma_b|$ can compensate by ap-
 934 proaching minus infinity. Intuitively, in much the same way
 935 $\text{argmin}_z \frac{1}{z} + \log z = 1$, meaning the optimal z must lie in the
 936 ‘span’ of 1 else the overall objective will be driven to infinity.

937 Consequently, the only way to minimize the cost in the limit
 938 as $\lambda \rightarrow 0$ is to consider low-rank solutions within the constraint
 939 set that $\mathbf{b} \in \text{span}[\mathbf{W}]$, and it is equivalent to requiring that
 940 $\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq C$ for some constant C independent of λ (which
 941 ultimately corresponds with maintaining $\mathcal{A}(\mathbf{X}) = \mathbf{b}$ in the limit
 942 as well).

943 In this setting, while $0 \leq \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq C$ is bounded, the sec-
 944 ond term in $\mathcal{L}(\Psi, \nu)$ can be unbounded from below when
 945 $\text{rank}[\Psi]$ is sufficiently small. To see this note that

$$\log |\Sigma_b| = \sum_{i=1}^p \log (\sigma_i [\mathbf{A}\bar{\Psi}\mathbf{A}^\top] + \lambda), \quad (23)$$

946 where $\sigma_i [\cdot]$ denotes the i -th singular value of a matrix. While
 947 the maximum rank of $\mathbf{A}\bar{\Psi}\mathbf{A}^\top$ is obviously p , if $r \triangleq \text{rank}[\Psi] <$
 948 p/m and $\text{spark}[\mathbf{A}] = p + 1$ (maximal spark) as stipulated in the
 949 lemma statement, then $\text{rank}[\mathbf{A}\bar{\Psi}\mathbf{A}^\top] = mr$ and (23) becomes

$$\log |\Sigma_b| = \sum_{i=1}^{mr} \log (\sigma_i [\mathbf{A}\bar{\Psi}\mathbf{A}^\top] + \lambda) + (p - mr) \log \lambda. \quad (24)$$

950 Note that the spark assumption accomplishes two objectives
 951 in this context. First, it guarantees that a high rank Ψ cannot
 952 masquerade as a low rank Ψ behind the nullspace of some col-
 953 lection of columns \mathbf{A}_i . Secondly, it ensures that after assuming
 954 $r < p/m$, then $\text{rank}[\mathbf{A}\bar{\Psi}\mathbf{A}^\top] = mr$.

955 Consequently, in the limit where $\lambda \rightarrow 0$ (with the limit being
 956 taken outside of the minimization), (23) effectively scales as
 957 $(p - mr) \log \lambda$, and hence the overall cost is minimized when

Ψ has minimal rank. This in turn ensures that the corresponding
 \mathbf{X} will also have minimal rank, completing the proof sketch for
 Lemma 1.

Finally, Lemma 2 follows directly from the structure of the
 $\mathcal{L}(\Psi, \nu)$ cost function via simple reparameterizations. ■

963 B. Proof of Theorem 1

964 To begin we assume that $\mathbf{b}_i \neq 0, \forall i$, where \mathbf{b}_i denotes the
 965 sub-vector of \mathbf{b} such that $\mathbf{b}_i = \mathbf{A}_i \mathbf{x}_i$. If this were not the case
 966 we can always collapse \mathbf{X} by the corresponding column (which
 967 is indistinguishable from zero) and achieve an equivalent result.
 968 Given the assumptions of Theorem 1, the BARM cost function
 969 becomes

$$\mathcal{L}(\Psi, \nu) = \sum_{i=1}^m \mathbf{b}_i^\top (\nu_i \mathbf{A}_i \Psi \mathbf{A}_i^\top)^{-1} \mathbf{b}_i + \log |\nu_i \mathbf{A}_i \Psi \mathbf{A}_i^\top|. \quad (25)$$

970 If there exists a feasible rank one solution to $\mathbf{b} = \text{Avec}$
 $[\mathbf{X}]$, then there also exists a set of $\Psi'_i = \nu_i \Psi$ such that $\mathbf{b}_i \mathbf{b}_i^\top =$
 $\mathbf{A}_i \Psi'_i \mathbf{A}_i^\top$ for all i . To see this, note that $\mathbf{b}_i \mathbf{b}_i^\top = \mathbf{A}_i \mathbf{x}_i \mathbf{x}_i^\top$
 \mathbf{A}_i^\top . Because $\text{rank}[\mathbf{X}] = 1$, it also follows that $\mathbf{b}_i \mathbf{b}_i^\top = \alpha_i \mathbf{A}_i \mathbf{X}$
 $\mathbf{X}^\top \mathbf{A}_i^\top$, where $\alpha_i = \|\mathbf{x}_i \mathbf{x}_i^\top\| / \|\mathbf{X} \mathbf{X}^\top\|$. Therefore $\Psi'_i =$
 $\nu_i \mathbf{X} \mathbf{X}^\top$ achieves the desired result with $\nu_i = \alpha_i$.

976 Now suppose we have converged to any solution $\{\hat{\Psi}, \hat{\nu}\}$ with
 $\text{rank}[\hat{\Psi}] > 1$ and associated $\hat{\Sigma} = \mathbf{I} \otimes \hat{\Psi}$. Note that since $\mathbf{b}_i \neq$
 $0, \nu_i > 0$ for all i , otherwise a local minimum is not possible
 (the cost function would be driven to positive infinity).

980 Define $\hat{\Sigma}_{b_i} = \hat{\nu}_i \mathbf{A}_i \hat{\Psi} \mathbf{A}_i^\top$. Additionally we can assume that
 $\mathbf{b}_i^\top \hat{\Sigma}_{b_i}^{-1} \mathbf{b}_i$ is finite, meaning that \mathbf{b}_i lies in the span of the singular
 982 vectors of $\hat{\Sigma}_{b_i}$. (If this were not the case, the cost would be
 983 driven to infinity and we could not be at a minimizing solution
 984 anyway.) If $\{\hat{\Psi}, \hat{\nu}\}$ is a local minimum, then $\{\lambda_1 = 1, \lambda_2 = 0\}$
 985 must be a local minimum of the revised cost function

$$\mathcal{L}(\lambda_1, \lambda_2) = \sum_{i=1}^m \mathbf{b}_i^\top (\lambda_1 \hat{\Sigma}_{b_i} + \lambda_2 \mathbf{b}_i \mathbf{b}_i^\top)^{-1} \mathbf{b}_i + \log |\lambda_1 \hat{\Sigma}_{b_i} + \lambda_2 \mathbf{b}_i \mathbf{b}_i^\top|. \quad (26)$$

986 This is because $\mathbf{b}_i \mathbf{b}_i^\top$ represents a valid set of basis vectors for
 987 updating the covariance per the construction above involving
 Ψ'_i . First consider optimization over λ_1 . If $\lambda_1 = 1$ is a local
 988 minimum, then by taking gradients and equating to zero, we
 989 require that
 990

$$\sum_{i=1}^m \mathbf{b}_i^\top \hat{\Sigma}_{b_i}^{-1} \mathbf{b}_i = \sum_{i=1}^m \text{rank}[\hat{\Sigma}_{b_i}]. \quad (27)$$

991 Likewise, taking the gradient with respect to λ_2 we obtain

$$\frac{\partial \mathcal{L}(\lambda_1, \lambda_2)}{\partial \lambda_2} \Big|_{\lambda_1=1, \lambda_2=0} = \sum_{i=1}^m \mathbf{b}_i^\top \hat{\Sigma}_{b_i}^{-1} \mathbf{b}_i - \sum_{i=1}^m (\mathbf{b}_i^\top \hat{\Sigma}_{b_i}^{-1} \mathbf{b}_i)^2. \quad (28)$$

992 The nullspace condition (a very mild assumption) ensures
 993 that $\sum_{i=1}^m \text{rank}[\hat{\Sigma}_{b_i}] = k$ for some $k > m$ when $\text{rank}[\Psi] > 1$.
 994 To see this, observe that to achieve $\sum_{i=1}^m \text{rank}[\hat{\Sigma}_{b_i}] = m$ when
 995 $\text{rank}[\Psi] > 1$ requires that $\Psi = \mathbf{u} \mathbf{u}^\top + \mathbf{W} \mathbf{W}^\top$ where \mathbf{u} is a

996 vector and \mathbf{W} is a matrix (or vector) with columns in $\text{null}[\mathbf{A}_i]$,
 997 $\forall i$. If any such \mathbf{W} is not in this nullspace for some i , then given
 998 that $p_i > 1$, the associated $\mathbf{A}_i \Psi \mathbf{A}_i^\top$ will have rank greater than
 999 one, and the overall rank sum will exceed m .

1000 Consequently, (28) will always be negative. This is because
 1001 if $\sum_{i=1}^m z_i = k$ for any set of non-negative variables $\{z_i\}$, the
 1002 minimal value of $\sum_{i=1}^m z_i^2$ occurs when $z_i = k/m, \forall i$. In our
 1003 case, this implies that

$$\sum_{i=1}^m (\mathbf{b}_i^\top \widehat{\Sigma}_{\mathbf{b}_i}^{-1} \mathbf{b}_i)^2 \geq \sum_{i=1}^m (k/m)^2 > k > m. \quad (29)$$

1004 Therefore we can add a small contribution of $\mathbf{b}_i \mathbf{b}_i^\top$ to each
 1005 $\widehat{\Sigma}_{\mathbf{b}_i}$ and reduce the underlying cost function. Hence we cannot
 1006 have a local minimum, except when Ψ is equal to some Ψ^*
 1007 with $\text{rank}[\Psi^*] = 1$. Moreover, we may directly conclude that
 1008 $\mathbf{x}^* = \overline{\Psi}^* \mathbf{A}^\top (\mathbf{A} \overline{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$ is feasible and $\text{rank}[\mathbf{X}^*] = 1$ with
 1009 $\mathbf{x}^* = \text{vec}[\mathbf{X}^*]$.

1010 Regarding the last part of the theorem, we consider only
 1011 f that are concave non-decreasing functions (this is the only
 1012 reasonable choice for shrinking singular values to zero, and
 1013 the more general case naturally follows anyway with additional
 1014 effort, but minimal enlightenment). Without loss of generality
 1015 we may also assume that $f(0) = 0$ and $f(1) = 1$; we can always
 1016 apply an inconsequential translation and scaling such that these
 1017 conditions hold.⁹ Simple counter examples then demonstrate
 1018 that $f(\epsilon)$ must be greater than some constant C independent of
 1019 ϵ for all ϵ sufficiently small. To see this, note that we can always
 1020 rescale elements of \mathbf{A} such that a solution with rank greater
 1021 than one is preferred unless this condition holds. However, such
 1022 an f , which effectively must display infinite gradient at $f(0)$ to
 1023 guarantee a global solution is always rank one, will then always
 1024 display local minima for certain \mathbf{A} . This can easily be revealed
 1025 through simple counter-examples. ■

1026 C. Symmetrization Update Rules

1027 These iterative update rules follow from alternative upper
 1028 bounds tailored to the symmetric version of BARM. When both
 1029 Ψ_r and Ψ_c are fixed, \mathbf{x} is updated via the posterior mean cal-
 1030 culation

$$\begin{aligned} \widehat{\mathbf{x}} &= \text{vec}[\widehat{\mathbf{X}}] = \frac{1}{2} (\overline{\Psi}_r + \overline{\Psi}_c) \mathbf{A}^\top \\ &\quad \times \left[\lambda \mathbf{I} + \mathbf{A} \frac{1}{2} (\overline{\Psi}_r + \overline{\Psi}_c) \mathbf{A}^\top \right]^{-1} \mathbf{b}. \end{aligned} \quad (30)$$

1031 where $\overline{\Psi}_r = \Psi_r \otimes \mathbf{I}$ and $\overline{\Psi}_c = \mathbf{I} \otimes \Psi_c$. Likewise we update
 1032 $\nabla_{\Psi_r^{-1}}$ and $\nabla_{\Psi_c^{-1}}$ using

$$\nabla_{\Psi_r^{-1}} = \sum_{i=1}^m \Psi_r - \Psi_r \mathbf{A}_{r_i}^\top (\mathbf{A} \overline{\Psi}_r \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_{r_i} \Psi_r, \quad (31)$$

$$\nabla_{\Psi_c^{-1}} = \sum_{i=1}^n \Psi_c - \Psi_c \mathbf{A}_{c_i}^\top (\mathbf{A} \overline{\Psi}_c \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_{c_i} \Psi_c, \quad (32)$$

⁹The log function is a limiting case, but what follows holds nonetheless.

where $\mathbf{A}_{r_i} \in \mathbb{R}^{p \times m}$ is defined such that $\mathbf{A} = [\mathbf{A}_{r_1}^\top, \dots, \mathbf{A}_{r_m}^\top]^\top$ 1033
 and $\mathbf{A}_{c_i} \in \mathbb{R}^{p \times m}$ is defined such that $\mathbf{A} = [\mathbf{A}_{c_1}, \dots, \mathbf{A}_{c_n}]$. Fi- 1034
 nally given these values, with \mathbf{X} , $\nabla_{\Psi_r^{-1}}$ and $\nabla_{\Psi_c^{-1}}$ fixed, we can 1035
 compute the optimal Ψ_r and Ψ_c in closed form by optimizing 1036
 the relevant Ψ_r - and Ψ_c -dependent terms via 1037

$$\Psi_r^{\text{opt}} = \frac{1}{n} \left[\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi_r^{-1}} \right], \quad (33)$$

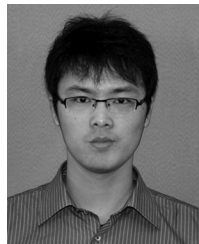
$$\Psi_c^{\text{opt}} = \frac{1}{m} \left[\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi_c^{-1}} \right]. \quad (34)$$

In practice the simple initialization $\Psi_r = \mathbf{I}$ and $\Psi_c = \mathbf{I}$ is 1038
 sufficient for obtaining good performance. 1039

REFERENCES 1040

- [1] E. J. Candès and B. Recht, "Exact matrix completion via convex optimiza- 1041
tion," *Found. Comput. Math.*, vol. 9, no. 6, pp. 717–772, 2009. 1042
- [2] Y. Hu, D. Zhang, J. Ye, X. Li, and X. He, "Fast and accurate matrix com- 1043
pletion via truncated nuclear norm regularization," *IEEE Trans. Pattern 1044
Anal. Mach. Intell. (PAMI)*, vol. 35, no. 9, pp. 2117–2130, 2013. 1045
- [3] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of 1046
subspace structures by low-rank representation," *IEEE Trans. Pattern 1047
Anal. Mach. Intell. (PAMI)*, vol. 35, no. 1, pp. 171–184, 2013. 1048
- [4] Z. Zhang, A. Ganesh, X. Liang, and Y. Ma, "Tilt: Transform invariant 1049
low-rank textures," *Int. J. Comput. Vis. (IJCV)*, vol. 99, no. 1, pp. 1–24, 1050
2012. 1051
- [5] C. Lu, J. Tang, S. Yan, and Z. Lin, "Generalized nonconvex nonsmooth 1052
low-rank minimization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog- 1053
nit. (CVPR)*, 2014. 1054
- [6] K. Mohan and M. Fazel, "Iterative reweighted algorithms for matrix rank 1055
minimization," *J. Mach. Learn. Res. (JMLR)*, vol. 13, no. 1, pp. 3441– 1056
3473, 2012. 1057
- [7] Z. Li, J. Liu, Y. Jiang, J. Tang, and H. Lu, "Low rank metric learning for 1058
social image retrieval," in *Pro. 20th ACM Int. Conf. Multimedia*, 2012, 1059
pp. 853–856. 1060
- [8] M. Tipping and C. Bishop, "Probabilistic principal component analysis," 1061
J. Roy. Statist. Soc. B, vol. 61, no. 3, pp. 611–622, 1999. 1062
- [9] B. Xin and D. Wipf, "Pushing the limits of affine rank minimization by 1063
adapting probabilistic pca," in *Proc. 32nd Int. Conf. Mach. Learn. (ICML)*, 1064
2015, pp. 419–427. 1065
- [10] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion 1066
using alternating minimization," in *Proc. 45th Annu. ACM Symp. Theory 1067
Comput.*, 2013, pp. 665–674. 1068
- [11] S. D. Babacan, M. Luessi, R. Molina, and A. K. Katsaggelos, "Sparse 1069
bayesian methods for low-rank matrix estimation," *IEEE Trans. Signal 1070
Process.*, vol. 60, no. 8, pp. 3964–3977, 2012. 1071
- [12] X. Ding, L. He, and L. Carin, "Bayesian robust principal component 1072
analysis," *IEEE Trans. Image Process.*, vol. 20, no. 12, pp. 3419–3430, 1073
2011. 1074
- [13] D. Wipf, "Non-convex rank minimization via an empirical Bayesian ap- 1075
proach," in *Proc. 28th Conf. Uncertainty Artif. Intell. (UAI)*, 2012. 1076
- [14] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex 1077
geometry of linear inverse problems," *Found. Comput. Math.*, vol. 12, 1078
no. 6, pp. 805–849, 2012. 1079
- [15] M. E. Tipping, "Sparse Bayesian learning and the relevance vector ma- 1080
chine," *J. Mach. Learn. Res. (JMLR)*, vol. 1, pp. 211–244, 2001. 1081
- [16] D. P. Wipf, B. D. Rao, and S. Nagarajan, "Latent variable Bayesian models 1082
for promoting sparsity," *IEEE Trans. Inf. Theory*, vol. 57, no. 9, pp. 6236– 1083
6255, 2011. 1084
- [17] D. L. Donoho and M. Elad, "Optimally sparse representation in general 1085
(nonorthogonal) dictionaries via l_1 minimization," *Proc. Nat. Acad. Sci.*, 1086
vol. 100, no. 5, pp. 2197–2202, 2003. 1087
- [18] E. J. Candès and X. Li, "Solving quadratic equations via phaselift when 1088
there are about as many equations as unknowns," *Found. Comput. Math.*, 1089
vol. 14, no. 5, pp. 1017–1026, 2014. 1090
- [19] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski, "Phase retrieval 1091
via matrix completion," *SIAM Rev.*, vol. 57, no. 2, pp. 225–251, 2015. 1092
- [20] W. I. Zangwill, *Nonlinear Programming: A Unified Approach*, Englewood 1093
Cliffs, NJ, USA: Prentice-Hall, 1969. 1094

- 1095 [21] J. Tanner and K. Wei, "Normalized iterative hard thresholding for matrix
1096 completion," *SIAM J. Scientif. Comput.*, vol. 35, no. 5, pp. S104–S125,
1097 2013.
- 1098 [22] P. Jain, R. Meka, and I. S. Dhillon, "Guaranteed rank minimization via
1099 singular value projection," in *Proc. Adv. Neural Inf. Process. Syst.*, 2010,
1100 pp. 937–945.
- 1101 [23] D. Goldfarb and S. Ma, "Convergence of fixed-point continuation algo-
1102 rithms for matrix rank minimization," *Found. Comput. Math.*, vol. 11, no.
1103 2, pp. 183–210, 2011.
- 1104 [24] R. H. Keshavan and S. Oh, "A gradient descent algorithm on the
1105 Grassman manifold for matrix completion," 2009 DOI: arXiv Preprint
1106 arXiv:0910.5260.
- 1107 [25] F. Léger, G. Yu, and G. Sapiro, "Efficient matrix completion with Gaussian
1108 models," 2010 DOI: arXiv Preprint arXiv:1010.4050.
- 1109 [26] B. Marlin, *Collaborative filtering: A machine learning perspective*, Ph.D.
1110 dissertation, Univ. of Toronto, Toronto, Canada ON, 2004.
- 1111 [27] B. M. Marlin, "Modeling user rating profiles for collaborative filtering,"
1112 in *Proc. Adv. Neural Inf. Process. Syst.*, 2003.
- 1113 [28] J. D. Rennie and N. Srebro, "Fast maximum margin matrix factorization
1114 for collaborative prediction," in *Proc. 22nd ACM Int. Conf. Mach. Learn.*,
1115 2005, pp. 713–719.
- 1116 [29] S.-T. Park and D. M. Pennock, "Applying collaborative filtering techniques
1117 to movie search for better ranking and browsing," in *Proc. 13th ACM
1118 SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2007, pp. 550–559.
- 1119 [30] D. DeCoste, "Collaborative prediction using ensembles of maximum margin
1120 matrix factorizations," in *Proc. 23rd ACM Int. Conf. Mach. Learn.*,
1121 2006, pp. 249–256.
- 1122 [31] N. D. Lawrence and R. Urtasun, "Non-linear matrix factorization with
1123 gaussian processes," in *Proc. 26th Annu. ACM Int. Conf. Mach. Learn.*,
1124 2009, pp. 601–608.
- 1125 [32] M. Zhou, C. Wang, M. Chen, J. Paisley, D. Dunson, and L. Carin, "Non-
1126 parametric Bayesian matrix completion," *Proc. IEEE SAM*, 2010.



1127
1128
1129
1130
1131
1132
1133

Bo Xin (M'XX) received the B.S. degree in electronic engineering from Dalian University of Technology, China, in 2011. He is currently working toward the Ph.D. degree in computer science at Peking University, China. His research interests include optimization, machine learning and computer vision.



and the 2006 NIPS Outstanding Paper Award.

David Wipf (M'XX) received the B.S. degree with highest honors from the University of Virginia, and the Ph.D. degree from UC San Diego, where he was an NSF IGERT Fellow. Later he was an NIH Post-doctoral Fellow at UC San Francisco. Since 2011 he has been with Microsoft Research in Beijing. His research interests include Bayesian learning techniques applied to signal/image processing and computer vision. He is the recipient of several awards including the 2012 Signal Processing Society Best Paper Award, the Biomag 2008 Young Investigator Award, and the 2006 NIPS Outstanding Paper Award.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146

Q2



Yizhou Wang (M'XX) received his Ph.D. in computer science from University of California at Los Angeles (UCLA) in 2005. He was a Research Staff of the Palo Alto Research Center (Xerox-PARC) from 2005 to 2008. He is currently a Professor of the Computer Science Department at Peking University (PKU), China. His research interests include computer vision, statistical modeling and learning.

1147
1148
1149
1150
1151
1152
1153
1154
1155



Wen Gao (F'XX) received M.S. degree in computer science from Harbin Institute of Technology in 1985, and Ph.D. degree in electronics engineering from the University of Tokyo in 1991. He was a Professor in computer science at Harbin Institute of Technology from 1991 to 1995 and a Professor in computer science at Institute of Computing Technology of Chinese Academy of Sciences from 1996 to 2005. He is currently a Professor at the School of Electronics Engineering and Computer Science, Peking University, China. He has been leading research efforts to

1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166

develop systems and technologies for video coding, face recognition, sign language recognition and synthesis, and multimedia retrieval. He earned many awards, which include five national awards for his research achievements and activities. He did many services to academic society, such as general co-chair of IEEE ICME07, and the head of Chinese delegation to the Moving Picture Expert Group (MPEG) of International Standard Organization (ISO). Since 1997, he is also the chairman of the working group responsible for setting a national Audio Video coding Standard (AVS) for China. He published four books and over 500 technical articles in refereed journals and proceedings in the areas of signal processing, image and video communication, computer vision, multimodal interface, pattern recognition, and bioinformatics.

1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178

Q3

- 1180 Q1. Author: For all conference paper references, provide page numbers if printed in proceeding or location of where conference
1181 was presented if not printed.
- 1182 Q2. Author: Please provide a forward facing headshot for Dr. Wipf. Otherwise the bio will have to run without photo.
- 1183 Q3. Author: Please provide initial year(s) of IEEE membership grade(s) for all authors.

IEEE Proof

Exploring Algorithmic Limits of Matrix Rank Minimization Under Affine Constraints

Bo Xin, *Member, IEEE*, David Wipf, *Member, IEEE*, Yizhou Wang, *Member, IEEE*, and Wen Gao, *Fellow, IEEE*

Abstract—Many applications require recovering a matrix of minimal rank within an affine constraint set, with matrix completion a notable special case. Because the problem is NP-hard in general, it is common to replace the matrix rank with the nuclear norm, which acts as a convenient convex surrogate. While elegant theoretical conditions elucidate when this replacement is likely to be successful, they are highly restrictive and convex algorithms fail when the ambient rank is too high or when the constraint set is poorly structured. Nonconvex alternatives fare somewhat better when carefully tuned; however, convergence to locally optimal solutions remains a continuing source of failure. Against this backdrop, we derive a deceptively simple and parameter-free probabilistic PCA-like algorithm that is capable, over a wide battery of empirical tests, of successful recovery even at the theoretical limit where the number of measurements equals the degrees of freedom in the unknown low-rank matrix. Somewhat surprisingly, this is possible even when the affine constraint set is highly ill-conditioned. While proving general recovery guarantees remains evasive for nonconvex algorithms, Bayesian-inspired or otherwise, we nonetheless show conditions whereby the underlying cost function has a unique stationary point located at the global optimum; no existing cost function we are aware of satisfies this property. The algorithm has also been successfully deployed on a computer vision application involving image rectification and a standard collaborative filtering benchmark.

Index Terms—Rank minimization, affine constraints, matrix completion, matrix recovery, empirical Bayes.

I. INTRODUCTION

RECENTLY there has been a surge of interest in finding minimum rank matrices subject to some problem-specific constraints often characterized as an affine set [1]–[7]. Mathematically this involves solving

$$\min_{\mathbf{X}} \text{rank}[\mathbf{X}] \quad \text{s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}), \quad (1)$$

where $\mathbf{X} \in \mathbb{R}^{n \times m}$ is the unknown matrix, $\mathbf{b} \in \mathbb{R}^p$ represents a vector of observations and $\mathcal{A} : \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^p$ denotes a linear mapping. An important special case of (1) commonly applied

to collaborative filtering is the matrix completion problem

$$\min_{\mathbf{X}} \text{rank}[\mathbf{X}] \quad \text{s.t. } \mathbf{X}_{ij} = (\mathbf{X}_0)_{ij}, (i, j) \in \Omega, \quad (2)$$

where \mathbf{X}_0 is a low-rank matrix we would like to recover, but we are only able to observe elements from the set Ω [1], [2]. Unfortunately however, both this special case and the general problem (1) are well-known to be NP-hard, and the rank penalty itself is non-smooth. Consequently, a popular alternative is to instead compute

$$\min_{\mathbf{X}} \sum_i f(\sigma_i[\mathbf{X}]) \quad \text{s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}), \quad (3)$$

where $\sigma_i[\mathbf{X}]$ denotes the i -th singular value of \mathbf{X} and f is usually a concave, non-decreasing function (or nearly so). In the special case where $f(z) = I[z \neq 0]$ (i.e., an indicator function) we retrieve the matrix rank; however, smoother surrogates such as $f(z) = \log z$ or $f(z) = z^q$ with $q \leq 1$ are generally preferred for optimization purposes. When $f(z) = z$, (3) reduces to convex nuclear norm minimization. A variety of celebrated theoretical results have quantified specific conditions, heavily dependent on the singular values of matrices in the nullspace of \mathcal{A} , where the minimum nuclear norm solution is guaranteed to coincide with that of minimal rank [1], [3], [6]. However, these guarantees typically only apply to a highly restrictive set of rank minimization problems, and in a practical setting non-convex algorithms can succeed in a much broader range of conditions [2], [5], [6].

In Section II we will summarize state-of-the-art non-convex rank minimization algorithms that operate under affine constraints and point out some of their shortcomings. This will be followed in Section III by the derivation of an alternative approach using Bayesian modeling techniques adapted from probabilistic PCA [8]. Section IV will then describe connections with nuclear norm minimization, convergence issues, and properties of global and local solutions. The latter includes special cases whereby any stationary point of the intrinsic cost function is guaranteed to have optimal rank, illustrating an underlying smoothing mechanism which leads to success over competing methods. We next discuss algorithmic enhancements in Section V that further improve recovery performance in practice. Section VI contains a wide variety of numerical comparisons that highlight the efficacy of this algorithm, while Section VII presents a computer vision application involving image rectification and a standard collaborative filtering benchmark. Technical proofs and algorithm update rule details are contained in the Appendix. Portions of this work have previously appeared in conference proceedings [9].

Manuscript received October 22, 2014; revised March 9, 2015 and November 23, 2015; accepted February 26, 2016. Date of publication April 7, 2016; date of current version. The associate editor coordinating the review of this manuscript and approving it for publication was Dr. Tareq Al-Naffouri. The authors would like to thank the support from the following grants: 973-2015CB351800, NSFC-61231010, NSFC-61527804, NSFC-61210005 and the Microsoft Research Asia Collaborative Research funding.

B. Xin, Y. Wang, and W. Gao are with the Department of Electrical Engineering and Computer Science, Peking University, Beijing 100871, China (e-mail: jimxinbo@gmail.com; yizhou.wang@pku.edu.cn; wgao@pku.edu.cn).

D. Wipf is with the Visual Computing group, Microsoft Research, Beijing 100080, China (e-mail: davidwip@microsoft.com).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSP.2016.2551697

Before proceeding, we highlight several main contributions as follows:

- 1) Bayesian inspiration can take uncountably many different forms and parameterizations, but the devil is in the details and existing methods offer little opportunity for both theoretical inquiry and substantial performance gains solving (1). In this regard, we apply carefully-tailored modifications to a veteran probabilistic PCA model leading to systematic theoretical and empirical insights and advantages. Model justification is ultimately based on such meticulous technical considerations rather than merely the presumed qualitative legitimacy of any underlying prior distributions.
- 2) Non-convex algorithms have demonstrated some improvement in estimation accuracy over the celebrated convex nuclear norm; however, this typically requires the inclusion of one or more additional tuning parameters to incrementally inject additional objective function curvature and avoid bad local solutions. In contrast, for solving (1) our non-convex Bayesian-inspired algorithm requires no such parameters at all, and noisy relaxations necessitate only a single, standard trade-off parameter balancing data-fit and minimal rank.¹
- 3) Over a wide battery of controlled experiments with ground-truth data, our approach outperforms all existing algorithms that we are aware of, Bayesian, non-convex, or otherwise. This includes direct head-to-head comparisons using the exact experimental designs and code prepared by original authors. In fact, even when \mathcal{A} is ill-conditioned we are consistently able to solve (1) right up to the theoretical limit of any possible algorithm, which has never been demonstrated previously.

II. RELATED WORK

Here we focus on a few of the latest and most effective rank minimization algorithms, all developed within the last few years and evaluated favorably against the state-of-the-art.

A. General Non-Convex Methods

In the non-convex regime, effective optimization strategies attempt to at least locally minimize (3), often exceeding the performance of the convex nuclear norm. For example, [6] derives a family of *iterative reweighted least squares* (IRLS) algorithms applied to $f(z) = (z^2 + \gamma)^{q/2}$ with $q, \gamma > 0$ as tuning parameters. A related penalty also considered, which coincides with the limit as $q \rightarrow 0$ (up to an inconsequential scaling and translation), is $f(z) = \log(z^2 + \gamma)$, which maintains an intimate connection with rank given that

$$\log z = \lim_{q \rightarrow 0} q^{-1} (z^q - 1) \quad \text{and} \quad \lim_{q \rightarrow 0} z^q = I [z \neq 0], \quad (4)$$

where I is a standard indicator function. Consequently, when γ is small, $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ behaves much like a scaled

¹While not our emphasis here, similar to other Bayesian frameworks, even this trade-off parameter can ultimately be learned from the data if a true, parameter-free implementation is desired across noise levels.

and translated version of the rank, albeit with nonzero gradients away from zero.

The IRLS0 algorithm from [6] represents the best-performing special case of the above, where $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ is minimized using a homotopy continuation scheme merged with IRLS. Here a fixed γ is replaced with a decreasing sequence $\{\gamma^k\}$, the rationale being that when γ^k is large, the cost function is relatively smooth and devoid of local minima. As the iterations k progress, γ^k is reduced, and the cost behaves more like the matrix rank function. However, because now we are more likely to be within a reasonably good basin of attraction, spurious local minima are more easily avoided. The downside of this procedure is that it requires a pre-defined heuristic for reducing γ^k , and this schedule may be problem specific. Moreover, there is no guarantee that a global solution will ever be found.

In a related vein, [5] derives a family of *iterative reweighted nuclear norm* (IRNN) algorithms that can be applied to virtually any concave non-decreasing function f , even when f is non-smooth, unlike IRLS. For effective performance however the authors suggest a continuation strategy similar to IRLS0. Moreover, additional tuning parameters are required for different classes of functions f and it remains unclear which choices are optimal. While the reported results are substantially better than when using the convex nuclear norm, in our experiments IRLS0 seems to perform slightly better, possibly because the quadratic least squares inner loop is less aggressive in the initial stages of optimization than weighted nuclear norm minimization, leading to a better overall trajectory. Regardless, all of these affine rank minimization algorithms fail well before the theoretical recovery limit is reached, when the number of observations p equals the number of degrees of freedom in the low-rank matrix we wish to recover. Specifically, for an $n \times m$, rank r matrix, the number of degrees of freedom is given by $r(m+n) - r^2$, hence $p = r(m+n) - r^2$ is the best-case boundary. In practice if \mathcal{A} is ill-conditioned or degenerate the achievable limit may be more modest.

A third approach relies on replacing the convex nuclear norm with a truncated non-convex surrogate [2]. While some competitive results for image inpainting via matrix completion are shown, in practice the proposed algorithm has many parameters to be tuned via cross-validation. Moreover, recent comparisons contained in [5] show that default settings perform relatively poorly.

Finally, a somewhat different class of non-convex algorithms can be derived using a straightforward application of alternating minimization [10]. The basic idea is to assume $\mathbf{X} = \mathbf{UV}^T$ for some low-rank matrices \mathbf{U} and \mathbf{V} and then solve

$$\min_{\mathbf{U}, \mathbf{V}} \|b - \mathcal{A}(\mathbf{UV}^T)\|_{\mathcal{F}} \quad (5)$$

via coordinate descent. The downside of this approach is that it can be sensitive to data correlations and requires that \mathbf{U} and \mathbf{V} be parameterized with the correct rank. In contrast, our emphasis here is on algorithms that require no prior knowledge whatsoever regarding the true rank. This is especially important in application extensions that may manage multiple low-rank

183 matrices such that prior knowledge of all individual ranks is not
184 feasible.

185 B. Bayesian Methods

186 From a probabilistic perspective, previous work has applied
187 Bayesian formalisms to rank minimization problems, although
188 not specifically within an affine constraint set. For example,
189 [11]–[13] derive robust PCA algorithms built upon the linear
190 summation of a rank penalty and an element-wise sparsity
191 penalty. In particular, [12] applies an MCMC sampling approach
192 for posterior inference, but the resulting iterations are not scal-
193 able, subjectable to detailed analysis, nor readily adaptable to
194 affine constraints. In contrast, [11] applies a similar probabilis-
195 tic model but performs inference using a variational mean-field
196 approximation. While the special case of matrix completion
197 is considered, from an empirical standpoint its estimation accu-
198 racy is not competitive with the state-of-the-art non-convex
199 algorithms mentioned above. Finally, without the element-wise
200 sparsity component intrinsic to robust PCA (which is not our
201 focus here), [13] simply collapses to a regular PCA model with
202 a closed-form solution, so the challenges faced in solving (1) do
203 not apply. Consequently, general affine constraints really are a
204 key differentiating factor.

205 From a motivational angle, the basic probabilistic model with
206 which we begin our development can be interpreted as a care-
207 fully re-parameterized generalization of the probabilistic PCA
208 model from [8]. This will ultimately lead to a non-convex algo-
209 rithm devoid of the heuristic tuning strategies mentioned above,
210 but nonetheless still uniformly superior in terms of estimation
211 accuracy. We emphasize that, although we employ a Bayesian
212 entry point for our algorithmic strategy, final justification of the
213 underlying model will be entirely based on properties of the
214 underlying cost function that emerges, rather than any putative
215 belief in the actual validity of the assumed prior distributions
216 or likelihood function. This is quite unlike the vast majority of
217 existing Bayesian approaches.

218 C. Analytical Considerations

219 Turning to analytical issues, a number of celebrated theoret-
220 ical results dictate conditions whereby substitution of the rank
221 function with the convex nuclear norm in (1) is nonetheless guar-
222 anteed to still produce the minimal rank solution. For example,
223 if \mathcal{A} is a Gaussian iid measurement ensemble and $\mathbf{X}_0 \in \mathbb{R}^{n \times n}$
224 represents the optimal solution to (1) with $\text{rank}[\mathbf{X}_0] = r$, then
225 with high probability as the problem dimensions grow large, the
226 minimum nuclear norm feasible solution will equal \mathbf{X}_0 if the
227 number of measurements p satisfies $p \geq 3r(2n - r)$ [14].

228 The limitation of this type of result is two-fold. First, in the
229 above situation the true minimum rank solution only actually re-
230 quires $p \geq r(2n - r)$ measurements to be recoverable via brute
231 force solution of (1), and the remaining difference of a factor
232 of three can certainly be considerable in many practical situa-
233 tions (e.g., requiring 300 measurements is far more laborious
234 than only needing 100 measurements). Secondly though, and
235 far more importantly, all existing provable recovery guarantees
236 place extremely strong restrictions on the structure of \mathcal{A} , e.g.,

strong restrictions on the singular value decay of matrices in
the nullspace of \mathcal{A} . Such conditions are unlikely to ever hold in
realistic application settings, including the image rectification
example we describe in Section VII.A (in fact, these conditions
are usually incapable of even being checked). In contrast, the
algorithm we propose is empirically observed to only require
the theoretically minimal number of measurements even when
such nullspace conditions are violated in many cases. While a
general theoretical guarantee of this sort is obviously not possi-
ble, we do nonetheless provide several supporting theoretical
results indicative of why such performance is at least empirically
obtainable.

249 III. ALTERNATIVE ALGORITHM DERIVATION

250 In this section we first detail our basic distributional assump-
251 tions followed by development of the associated update rules
252 for inference.

253 A. Basic Model

254 In contrast to the majority of existing algorithms organized
255 around practical solutions to (3), here we adopt an alternative,
256 probabilistic starting point. We first define the Gaussian likeli-
257 hood function

$$p(\mathbf{b}|\mathbf{X}; \mathcal{A}, \lambda) \propto \exp\left[-\frac{1}{2\lambda} \|\mathcal{A}(\mathbf{X}) - \mathbf{b}\|_2^2\right], \quad (6)$$

258 noting that in the limit as $\lambda \rightarrow 0$ this will enforce the same
259 constraint set as in (1). Next we define an independent, zero-
260 mean Gaussian prior distribution with covariance $\nu_i \Psi$ on each
261 column of \mathbf{X} , denoted $\mathbf{x}_{:i}$ for all $i = 1, \dots, m$. This produces
262 the aggregate prior on \mathbf{X} given by

$$p(\mathbf{X}; \Psi, \nu) = \prod_i \mathcal{N}(\mathbf{x}_{:i}; \mathbf{0}, \nu_i \Psi) \propto \exp\left[\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x}\right], \quad (7)$$

263 where $\Psi \in \mathbb{R}^{n \times n}$ is a positive semi-definite symmetric matrix,²
264 $\nu = [\nu_1, \dots, \nu_m]^\top$ is a non-negative vector, $\mathbf{x} = \text{vec}[\mathbf{X}]$
265 (column-wise vectorization), and $\overline{\Psi} = \text{diag}[\nu] \otimes \Psi$, with \otimes
266 denoting the Kronecker product. It is important to stress here
267 that we do not necessarily believe that the unknown \mathbf{X} actually
268 follows such a Gaussian distribution per se. Rather, we adopt
269 (7) primarily because it will lead to an objective function with
270 desirable properties related to solving (1).

271 Moving forward, given both likelihood and prior are Gaus-
272 sian, the posterior $p(\mathbf{X}|\mathbf{b}; \Psi, \nu, \mathcal{A}, \lambda)$ is also Gaussian, with
273 mean given by an $\widehat{\mathbf{X}}$ such that

$$\widehat{\mathbf{x}} = \text{vec}\left[\widehat{\mathbf{X}}\right] = \overline{\Psi} \mathbf{A}^\top (\lambda \mathbf{I} + \mathbf{A} \overline{\Psi} \mathbf{A}^\top)^{-1} \mathbf{b}. \quad (8)$$

²Technically Ψ must be positive definite for the inverse in (7) to be defined. However, we can accommodate the semi-definite case using the following convention. Without loss of generality assume that $\overline{\Psi} = \mathbf{R} \mathbf{R}^\top$ for some matrix \mathbf{R} . We then qualify that $p(\mathbf{X}; \Psi, \nu) = 0$ if $\mathbf{x} \notin \text{span}[\mathbf{R}]$, and $p(\mathbf{X}; \Psi, \nu) \propto \exp[-\frac{1}{2} \mathbf{x}^\top (\mathbf{R}^\top)^\dagger \mathbf{R} \mathbf{x}]$ otherwise. Equivalently, throughout the paper for convenience (and with slight abuse of notation) we define $\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x} = \infty$ when $\mathbf{x} \notin \text{span}[\mathbf{R}]$, and $\mathbf{x}^\top \overline{\Psi}^{-1} \mathbf{x} = \mathbf{x}^\top (\mathbf{R}^\top)^\dagger \mathbf{R} \mathbf{x}$ otherwise. This will come in handy, for example, when interpreting the bound in (12) below. Note also that the final cost function (10) we will ultimately be minimizing requires no such inverse anyway.

Here $\mathbf{A} \in \mathbb{R}^{p \times nm}$ is a matrix defining the linear operator \mathcal{A} such that $\mathbf{b} = \mathbf{A}\mathbf{x}$ reproduces the feasible region in (1). From this expression it is clear that, if Ψ represents a low-rank covariance matrix, then each column of $\widehat{\mathbf{X}}$ will be constrained to a low-dimensional subspace resulting overall in a low-rank estimate as desired. Of course for this simple strategy to be successful we require some way of determining a viable Ψ and the scaling vector ν .

A common Bayesian strategy in this regard is to marginalize over \mathbf{X} and then maximize the resulting likelihood function with respect to Ψ and ν [15], [13], [16]. This involves solving

$$\max_{\Psi \in H^+, \nu \geq 0} \int p(\mathbf{b}|\mathbf{X}; \mathcal{A}, \lambda) p(\mathbf{X}; \Psi, \nu) d\mathbf{X}, \quad (9)$$

where H^+ denotes the set of positive semi-definite and symmetric $n \times n$ matrices. After a -2 log transformation and application of a standard convolution-of-Gaussians integration, solving (9) is equivalent to minimizing the cost function

$$\mathcal{L}(\Psi, \nu) = \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} + \log |\Sigma_b|, \quad (10)$$

where

$$\Sigma_b = \mathbf{A} \bar{\Psi} \mathbf{A}^\top + \lambda \mathbf{I} \text{ and } \bar{\Psi} = \text{diag}[\nu] \otimes \Psi. \quad (11)$$

Here Σ_b is the covariance of \mathbf{b} given Ψ and ν .

B. Update Rules

Minimizing (10) is a non-convex optimization problem, and we employ standard upper bounds for this purpose leading to an EM-like algorithm, somewhat related to [8]. In particular, we compute separate bounds, parameterized by auxiliary variables, for both the first and second terms of $\mathcal{L}(\Psi, \nu)$. While the general case can easily be handled and may be applicable for more challenging problems, here for simplicity and ease of presentation we consider minimizing $\mathcal{L}(\Psi) \triangleq \mathcal{L}(\Psi, \nu = \mathbf{1})$, meaning all elements of ν are fixed at one (and such is the case for all experiments reported herein, although we are currently exploring situations where this added generality could be especially helpful).

Based on [16], for the first term in (10) we have

$$\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + \mathbf{x}^\top \bar{\Psi}^{-1} \mathbf{x} \quad (12)$$

with equality whenever \mathbf{x} satisfies (8). For the second term we use

$$\begin{aligned} \log |\Sigma_b| &\equiv m \log |\Psi| + \log |\lambda \mathbf{A}^\top \mathbf{A} + \bar{\Psi}^{-1}| \\ &\leq m \log |\Psi| + \text{tr}[\Psi^{-1} \nabla_{\Psi^{-1}}] + C, \end{aligned} \quad (13)$$

where because $\log |\lambda \mathbf{A}^\top \mathbf{A} + \bar{\Psi}^{-1}|$ is concave with respect to Ψ^{-1} , we can upper bound it using a first-order approximation with a bias term C that is independent of Ψ . Equality is obtained when the gradient satisfies

$$\nabla_{\Psi^{-1}} = \sum_{i=1}^m \Psi - \Psi \mathbf{A}_i^\top (\mathbf{A} \bar{\Psi} \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_i \Psi, \quad (14)$$

where $\mathbf{A}_i \in \mathbb{R}^{p \times n}$ is defined such that $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_m]$. Finally given the upper bounds from (12) and (13) with \mathbf{X}

and $\nabla_{\Psi^{-1}}$ fixed, we can compute the optimal Ψ in closed form by optimizing the relevant Ψ -dependent terms via

$$\begin{aligned} \Psi^{\text{opt}} &= \arg \min_{\Psi} \text{tr}[\Psi^{-1} (\mathbf{X} \mathbf{X}^\top + \nabla_{\Psi^{-1}})] + m \log |\Psi| \\ &= \frac{1}{m} [\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi^{-1}}]. \end{aligned} \quad (15)$$

By agnostically starting with $\Psi = \mathbf{I}$ and then iteratively computing (8), (14), and (15), we can then obtain an estimate for Ψ , and more importantly, a corresponding estimate for \mathbf{X} given by (8) at convergence. We refer to this basic procedure as BARM for *Bayesian Affine Rank Minimization*. The next section will describe in detail why it is particularly well-suited for solving problems such as (1).

IV. PROPERTIES OF BARM

Here we first describe a close but perhaps not intuitively-obvious relationship between the BARM objective function and canonical nuclear norm minimization. We then discuss desirable properties of global and local minima before concluding with a brief examination of convergence issues.

A. Connections with Nuclear Norm Minimization

On the surface, it may appear that minimizing (10) is completely unrelated to the convex problem

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* \text{ s.t. } \mathbf{b} = \mathcal{A}(\mathbf{X}) \quad (16)$$

that is most commonly associated with practical rank minimization implementations. However, a close connection can be revealed by considering the modified objective function

$$\mathcal{L}'(\Psi) = \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} + \text{tr}[\bar{\Psi}], \quad (17)$$

which represents nothing more than (10), with $\nu = \mathbf{1}$ and with $\log |\Sigma_b|$ being replaced by $\text{tr}[\bar{\Psi}]$. Now suppose we minimize (17) with respect to $\Psi \in H^+$ obtaining some Ψ^* . We then go on to compute an estimate of \mathbf{X} using (8). Note that if we apply the bound from (12) to the first term in (17), then this estimate for \mathbf{X} equivalently solves

$$\min_{\Psi \in H^+, \mathbf{x}} \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + \mathbf{x}^\top \bar{\Psi}^{-1} \mathbf{x} + \text{tr}[\bar{\Psi}], \quad (18)$$

with $\mathbf{x} = \text{vec}[\mathbf{X}]$ as before. If we first optimize over Ψ , it is easily demonstrated that the optimal value of Ψ equals $(\mathbf{X} \mathbf{X}^\top)^{1/2}$. Plugging this value into (18), simplifying, and then applying the definition of the nuclear norm, we arrive at

$$\min_{\mathbf{X}} \frac{1}{\lambda} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + 2\|\mathbf{X}\|_*, \quad (19)$$

Furthermore, in the limit $\lambda \rightarrow 0$ (applied outside of the minimization), (19) becomes equivalent to (16). For more information regarding the duality relationship between variance/covariance space and coefficient space, at least in the related context of compressive sensing models, please refer to [16].

Consequently, we may conclude that the central distinction between the proposed BARM cost function and nuclear norm minimization is an intrinsic \mathcal{A} -dependent penalty function

353 $\log |\Sigma_b|$ which is applied in covariance space. In Section IV.B
 354 we will examine desirable properties of this non-convex sub-
 355 stitution, highlighting our desire to treat the underlying BARM
 356 probabilistic model as an independent cost function that may be
 357 subject to technical analysis independent of its Bayesian origins.

358 B. Global/Local Minima Analysis

359 As discussed in Section II one nice property of the
 360 $\sum_i \log(\sigma_i[\mathbf{X}])$ penalty employed (approximately) by IRLS0
 361 [6] is that it can be viewed as a smooth version of the matrix
 362 rank function while still possessing the same set of minimum,
 363 both global and local, over the affine constraint set, at least if we
 364 consider the limiting situation of $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ when γ
 365 becomes small so that we may avoid the distracting singularity
 366 of $\log 0$. Additionally, it possesses an attractive form of scale
 367 invariance, meaning that if \mathbf{X}^* is an optimal feasible solution,
 368 a block-diagonal rescaling of \mathbf{A} nevertheless leads to an equiv-
 369 alent rescaling of the optimum (without the need for solving
 370 an additional optimization problem using the new \mathbf{A}). This is
 371 very much unlike the nuclear norm or other non-convex surro-
 372 gates that penalize the singular values of \mathbf{X} in a scale-dependent
 373 manner.

374 In contrast, the proposed algorithm is based on a very differ-
 375 ent Gaussian statistical model with seemingly a more tenuous
 376 connection with rank minimization. Encouragingly however,
 377 the proposed cost function enjoys the same global/local minima
 378 properties as $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ with $\gamma \rightarrow 0$. Before present-
 379 ing these results, we define $\text{spark}[\mathbf{A}]$ as the smallest number
 380 of linearly dependent columns in matrix \mathbf{A} [17]. All proofs are
 381 deferred to the Appendix.

382 *Lemma 1:* Let $\mathbf{b} = \text{Avec}[\mathbf{X}]$, where $\mathbf{A} \in \mathbb{R}^{p \times nm}$ satisfies
 383 $\text{spark}[\mathbf{A}] = p + 1$. Also define r as the smallest rank of any fea-
 384 sible solution. Then if $r < p/m$, any global minimizer $\{\Psi^*, \nu^*\}$
 385 of (10) in the limit $\lambda \rightarrow 0$ is such that $\mathbf{x}^* = \bar{\Psi}^* \mathbf{A}^\top (\mathbf{A} \bar{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$
 386 is feasible and $\text{rank}[\mathbf{X}^*] = r$ with $\text{vec}[\mathbf{X}^*] = \mathbf{x}^*$.

387 *Lemma 2:* Additionally, let $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{D}$, where $\mathbf{D} = \text{diag}$
 388 $[\alpha_1 \mathbf{\Gamma}, \dots, \alpha_m \mathbf{\Gamma}]$ is a block-diagonal matrix with invertible
 389 blocks $\mathbf{\Gamma} \in \mathbb{R}^{n \times n}$ of unit norm scaled with coefficients $\alpha_i > 0$.
 390 Then iff $\{\Psi^*, \nu^*\}$ is a minimizer (global or local) to (10) in
 391 the limit $\lambda \rightarrow 0$, then $\{\mathbf{\Gamma}^{-1} \Psi^*, \text{diag}[\alpha]^{-1} \nu^*\}$ is a minimizer when
 392 $\tilde{\mathbf{A}}$ replaces \mathbf{A} . The corresponding estimates of \mathbf{X} are likewise
 393 in one-to-one correspondence.

394 *Remarks:* The assumption $r = \text{rank}[\mathbf{X}^*] < p/m$ in Lemma
 395 1 is completely unrestrictive, especially given that a unique,
 396 minimal-rank solution is only theoretically possible by any algo-
 397 rithm if $p \geq (n + m)r - r^2$, which is much more restrictive
 398 than $p > rm$. Hence the bound we require is well above that
 399 required for uniqueness anyway. Likewise the spark assumption
 400 will be satisfied for any \mathbf{A} with even an infinitesimal (con-
 401 tinuous) random component. Consequently, we are essentially
 402 always guaranteed that BARM possesses the same global opti-
 403 mum as the rank function. Regarding Lemma 2, no surrogate
 404 rank penalty of the form $\sum_i f(\sigma_i[\mathbf{X}])$ can achieve this result
 405 except for $f(z) = \log z$, or inconsequential limiting translations
 406 and rescalings of the log such as the indicator function $I[z \neq 0]$
 407 (which is related to the log via arguments in Section II).

While these results are certainly a useful starting point, the
 real advantage of adopting the BARM cost function is that lo-
 cally minimizing solutions are exceedingly rare, largely as a
 consequence of the marginalization process in (9), and in some
 cases provably so. A specialized example of this smoothing can
 be quantified in the following scenario.

Suppose \mathbf{A} is now block diagonal, with diagonal blocks \mathbf{A}_i
 such that $\mathbf{b}_i = \mathbf{A}_i \mathbf{x}_i$, producing the aggregate observation vec-
 tor $\mathbf{b} = [\mathbf{b}_1^\top, \dots, \mathbf{b}_m^\top]^\top$. While somewhat restricted, this situa-
 tion nonetheless includes many important special cases, includ-
 ing canonical matrix completion and generalized matrix com-
 pletion where elements of $\mathbf{Z} = \mathbf{W}\mathbf{X}_0$ are observed after some
 transformation \mathbf{W} , instead of \mathbf{X}_0 directly.

Theorem 1: Let $\mathbf{b} = \text{Avec}[\mathbf{X}]$, where \mathbf{A} is block diagonal,
 with blocks $\mathbf{A}_i \in \mathbb{R}^{p_i \times n}$. Moreover, assume $p_i > 1$ for all i
 and that $\cap_i \text{null}[\mathbf{A}_i] = \emptyset$. Then if $\min_{\mathbf{X}} \text{rank}[\mathbf{X}] = 1$ in the
 feasible region, any minimizer $\{\Psi^*, \nu^*\}$ of (10) (global or local)
 in the limit $\lambda \rightarrow 0$ is such that $\mathbf{x}^* = \bar{\Psi}^* \mathbf{A}^\top (\mathbf{A} \bar{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$
 is feasible and $\text{rank}[\mathbf{X}^*] = 1$ with $\text{vec}[\mathbf{X}^*] = \mathbf{x}^*$. Furthermore,
 no cost function in the form of (3) can satisfy the same result.
 In particular, there can always exist local and/or global minima
 with rank greater than one.

Remarks: This result implies that, under extremely mild con-
 ditions, which do not even depend on the concentration proper-
 ties of \mathbf{A} , the proposed cost function has no minima that are not
 global minima, at least in this rank-one case. (The minor techni-
 cal condition regarding nullspace intersections merely ensures
 that high-rank components cannot simultaneously “hide” in the
 nullspace of every measurement matrix \mathbf{A}_i ; the actual \mathbf{A} opera-
 tor may still be highly ill-conditioned.) Thus any algorithm with
 provable convergence to some local minimizer is guaranteed to
 obtain a globally optimal solution.³

Although a global optimal guarantee for finding a rank-one
 matrix sounds somewhat limited, such a guarantee is not possi-
 ble with any other penalty function of the standard form
 $\sum_i f(\sigma_i[\mathbf{X}])$, which is the typical recipe for rank minimization
 algorithms, convex or not. Moreover, finding rank one matrices
 subject to affine constraints represents a crucial component of
 applications such as phase retrieval [18], [19].

Additionally, if a unique rank-one solution exists to (1), then
 the unique minimizing solution to (10) will produce this \mathbf{X} via
 (8). Crucially, this will occur even when the minimal number
 of measurements $p = n + m - 1$ are available, unlike any other
 algorithm we are aware of that is blind to the true underlying
 rank.⁴ Moreover, as evident from the experiments, the proposed
 algorithm always successfully finds the global optimal in many
 situations where the underlying matrix has a rank much higher
 than one. Therefore, although we can only provide theoretical
 guarantee for the rank-one case, the underlying intuition that
 local minima are smoothed away arguably carries over to situa-
 tions where the rank is greater than one.

³Note also that with minimal additional effort, it can be shown that no sub-
 optimal stationary points of any kind, including saddle points, are possible.

⁴It is important to emphasize that the difficulty of estimating the optimal low-
 rank solution is based on the ratio of the d.o.f. in \mathbf{X} to the number of observations
 p . Consequently, estimating \mathbf{X} even with r small can be challenging when p is
 also small, meaning \mathbf{A} is highly overcomplete.

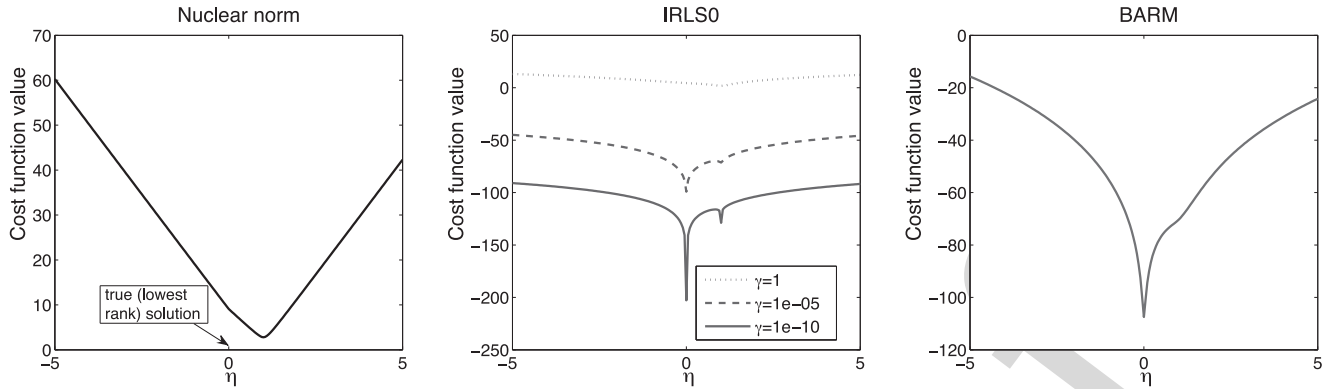


Fig. 1. Plots of different surrogates for matrix rank in a 1D feasible subspace. Here the convex nuclear norm does not retain the correct global minimum. In contrast, although the non-convex $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ penalty exhibits the correct minimum when γ is sufficiently small, it also contains spurious minima. Only BARM smooths away local minimum while simultaneously retaining the correct global optima.

459 C. Visualization of BARM Local Minima Smoothing

460 To further explore the smoothing effect and complement The-
 461 orem 1, it helps to visualize rank penalty functions restricted to
 462 the feasible region. While the BARM algorithm involves mini-
 463 mizing (10), its implicit penalty function on \mathbf{X} can nonetheless
 464 be numerically obtained across the feasible region in a given
 465 subspace of interest; for other penalties such as the nuclear
 466 norm this is of course trivial. Practically it is convenient to ex-
 467 plore a 1D feasible subspace generated by $\mathbf{X}^* + \eta\mathbf{V}$, where
 468 \mathbf{X}^* is the true minimum rank solution, $\mathbf{V} \in \text{null}[\mathbf{A}]$, and η
 469 is a scalar. We may then plot various penalty function values
 470 as η is varied, tracing the corresponding 1D feasible subspace.
 471 We choose $\mathbf{V} = \mathbf{X}^1 - \mathbf{X}^*$, where \mathbf{X}^1 is a feasible solution
 472 with minimum nuclear norm; however, random selections from
 473 $\text{null}[\mathbf{A}]$ also show similar characteristics.

474 Fig. 1 provides a simple example of this process. \mathbf{A} is gen-
 475 erated randomly with all zeros and a single randomly placed
 476 ‘1’ in each row leading to a canonical matrix completion prob-
 477 lem. $\mathbf{X}^* \in \mathbb{R}^{5 \times 5}$ is randomly generated as $\mathbf{X}^* = \mathbf{u}\mathbf{v}^\top$, where
 478 \mathbf{u} and \mathbf{v} are iid $\mathcal{N}(0, 1)$ vectors, and so \mathbf{X}^* is rank one. Finally,
 479 $p = 10$ elements are observed, and therefore \mathbf{A} has 10 rows and
 480 $5 \times 5 = 25$ columns. η is varied from -5 to 5 and the values of
 481 the nuclear norm, $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$, and the implicit BARM
 482 cost function are displayed.

483 From the figure we observe that the minimum of the nuclear
 484 norm is not produced when the rank is smallest, which occurs
 485 when $\eta = 0$; hence the convex cost function fails for this prob-
 486 lem. Likewise, the $\sum_i \log(\sigma_i[\mathbf{X}]^2 + \gamma)$ penalty used by IRLS0
 487 displays an incorrect global minimum when the tuning param-
 488 eter γ is large. In contrast, when γ is small, while the global
 489 minimum may now be correct, spurious local ditches have ap-
 490 peared in the cost function.⁵ Therefore, any success of the IRLS0
 491 algorithm depends heavily on a carefully balanced decaying se-
 492 quence of γ values, with the hope that initial iterations can steer
 493 the trajectory towards a desirable basin of attraction where local

⁵Technically speaking, these are not provably local minima since we are only considering a 1D subspace of the feasible region. However, it nonetheless illustrates the strong potential for troublesome local minima, especially in high dimensional practical problems.

minima are less problematic. One advantage of BARM then is
 that it is parameter free in this respect and yet still retains the
 correct global minimum, often without additional spurious local
 minima.

498 D. Convergence

499 Previous results of Section IV are limited to exploring aspects
 of the underlying BARM cost function. Regarding the BARM
 algorithm itself, by construction the updates generated by (8),
 (14), and (15) are guaranteed to reduce or leave unchanged
 $\mathcal{L}(\Psi)$ at each iteration. However, this is not technically suffi-
 cient to guarantee convergence to a stationary point of the cost
 function unless the additional conditions of Zangwill’s Global
 Convergence Theorem are satisfied [20]. However, provided we
 add a small regularization factor $\gamma \text{tr}[\Psi^{-1}]$, with $\gamma > 0$, then it
 can be shown that any cluster point of the resulting sequence of
 iterations $\{\Psi^k\}$ must be a stationary point. Moreover, because
 the sequence is bounded, there will always exist at least one
 cluster point, and therefore the algorithm is guaranteed to at
 least converge to a set of parameter values \mathcal{S} such that for any
 $\Psi^* \in \mathcal{S}$, $\mathcal{L}(\Psi^*) + \gamma \text{tr}[(\Psi^*)^{-1}]$ is a stationary point.

514 Finally, we should mention that this extra γ factor is akin to the
 homotopy continuation regularizer used by the IRLS0 algorithm
 [6] as discussed in Section II. However, whereas IRLS0 requires
 a carefully-chosen, decreasing sequence $\{\gamma^k\}$ with $\gamma^k > 0$ both
 to prove convergence and to avoid local minimum (and without
 this factor the algorithm performs very poorly in practice), for
 BARM a small, fixed factor only need be included as a technical
 necessity for proving formal convergence; in practice it can be
 fixed to exactly zero.

523 V. SYMMETRIZATION IMPROVEMENTS

524 Despite the promising theoretical attributes of BARM, there
 remains one important artifact of its probabilistic origins not
 found in more conventional existing rank minimization algo-
 rithms. In particular, other algorithms rely upon a symmetric
 penalty function that is independent of whether we are working
 with \mathbf{X} or \mathbf{X}^\top . All methods that reduce to (3) fall into this
 category, e.g., nuclear norm minimization, IRNN, or IRLS0. In

531 contrast, our method relies on defining a distribution with
 532 respect to the columns of \mathbf{X} . Consequently the underlying cost
 533 function is not identical when derived with respect to \mathbf{X} or
 534 \mathbf{X}^\top , a difference which will depend on \mathbf{A} . While globally opti-
 535 mal solutions should nonetheless be the same, the convergence
 536 trajectory could depend on this distinction leading to different
 537 local minima in certain circumstances. Although either con-
 538 struction leads to low-rank solutions, we may nonetheless expect
 539 improvement if we can somehow symmetrize the algorithm
 540 formulation.

541 To accomplish this, we consider a Gaussian prior on $\mathbf{x} =$
 542 $\text{vec}[\mathbf{X}]$ with a covariance formed using a block-wise averaging
 543 of covariances defined over rows and columns, denoted Ψ_r and
 544 Ψ_c respectively. The overall covariance is then given by the
 545 Kronecker sum

$$\overline{\Psi} = 1/2 (\Psi_r \otimes \mathbf{I} + \mathbf{I} \otimes \Psi_c). \quad (20)$$

546 The estimation process then proceeds in a similar fashion as
 547 before but with modifications and alternate upper-bounds that
 548 accommodate for this merger. For reported experimental results
 549 this symmetric version of BARM is used, with complete up-
 550 date rules listed in the Appendix and computational complexity
 551 evaluated in Section VI.E.

552 VI. EXPERIMENTAL VALIDATION

553 This section compares BARM with existing state-of-the-art
 554 affine rank minimization algorithms. For BARM, in all noise-
 555 less cases we simply used $\lambda = 10^{-10}$ (effectively zero), and
 556 hence no tuning parameters are required. Likewise, nuclear
 557 norm minimization [1], [4] requires no tuning parameters bey-
 558 ond implementation-dependent control parameters frequently
 559 used to enhance convergence speed (however the global mini-
 560 mum is unaltered given that the problem is convex). For the
 561 IRLS0 algorithm, we used our own implementation as the al-
 562 gorithm is straightforward and no code was available for the
 563 case of general \mathcal{A} ; we based the required decreasing γ_k se-
 564 quence on suggestions from [6]. IRLS0 code is available from
 565 the original authors for matrix completion; however, the results
 566 obtained with this code are not better than those obtained with
 567 our version. For the IRNN algorithm, we did not have access
 568 to code for general \mathcal{A} , nor specific details of how various pa-
 569 rameters should be set in the general case. Note also that IRNN
 570 has multiple parameters to tune even in noiseless problems un-
 571 like BARM. Therefore we report results directly from [5] where
 572 available. Note that both [5] and [6] show superior results to a
 573 number of other algorithms; we do not generally compare with
 574 these others given that they are likely no longer state-of-the-art
 575 and may clutter the presentation.

576 As stated previously, our focus here is on algorithms that do
 577 not require knowledge of the true rank of the optimal solution,
 578 and hence we do not include comparisons with [10] or the nor-
 579 malized hard thresholding algorithm from [21]. Regardless, we
 580 have nonetheless conducted numerous experiments with these
 581 algorithms, and even when the correct rank is provided, results
 582 are inferior to BARM, especially when correlated measurements
 583 are used. However, we do show limited empirical results with

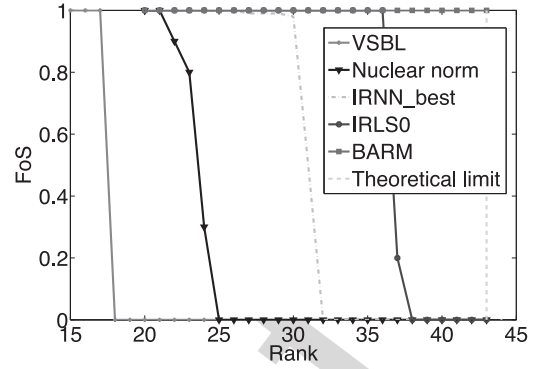


Fig. 2. Matrix completion comparisons (avg of 10 trials).

584 the variational sparse Bayesian algorithm (VSBL) from [11]
 585 because of its Bayesian origins, although the underlying param-
 586 eterization is decidedly different from BARM. But these results
 587 are limited to matrix completion as VSBL presently does not
 588 handle general affine constraints. Results from VSBL were ob-
 589 tained using publicly available code from the authors.

590 A. Matrix Completion

591 We begin with the matrix completion problem from (2), in
 592 part because this allows us to compare our results with the latest
 593 algorithms even when code is not available. For this purpose we
 594 reproduce the exact same experiment from [5], where a rank r
 595 matrix is generated as $\mathbf{X}_0 = \mathbf{TM}_L \mathbf{TM}_R$, with $\mathbf{TM}_L \in \mathbb{R}^{n \times r}$
 596 and $\mathbf{TM}_R \in \mathbb{R}^{r \times m}$ ($n = m = 150$) as iid $\mathcal{N}(0, 1)$ random ma-
 597 trices. 50% of all entries are then hidden uniformly at random.
 598 The *relative error* (REL) given by $\|\mathbf{X}_0 - \widehat{\mathbf{X}}\|_{\mathcal{F}} / \|\mathbf{X}_0\|_{\mathcal{F}}$ is
 599 computed for each trial and averaged as r is varied. Likewise,
 600 we compute the *frequency of success* (FoS) score, which mea-
 601 sures the percentage of trials where the REL is below 10^{-3} .
 602 Results are shown in Fig. 2 where BARM is the only algorithm
 603 capable of reaching the theoretical recovery limit, beyond which
 604 $p = 0.5 \times 150^2 = 11250$ is surpassed by the number of degrees
 605 of freedom in \mathbf{X}_0 , in this case $2 \times 150 \times 44 - 44^2 = 11264$.
 606 Note that FoS values were reported in [5] over a wide range of
 607 non-convex IRNN algorithms. The green curve represents the
 608 best performing candidate from this pool as tuned by the original
 609 authors; REL values were unavailable. Interestingly, although
 610 VSBL is based on a somewhat related probabilistic model to
 611 BARM, the underlying parameterization, cost function, and up-
 612 date rules are entirely different and do not benefit from strong
 613 theoretical underpinnings. Hence performance does not always
 614 match recent state-of-the-art algorithms, although from a com-
 615 putational standpoint it is quite efficient.

616 Besides BARM, the IRLS0 algorithm also displayed better
 617 performance than the other methods. This motivated us to re-
 618 produce some of the matrix completion experiments from [6] so
 619 as to provide direct head-to-head comparisons with the authors'
 620 original implementation. For this purpose, \mathbf{X}_0 is conveniently
 621 generated in the same way as above; however, values of $n, m,$
 622 r , and the percentage of missing entries are varied while eval-
 623 uating reconstructions using FoS. While [6] tests a variety of

TABLE I
MATRIX COMPLETION RESULTS OF BARM WITH IRLS0 ON THE THREE
HARDEST PROBLEMS FROM [6]. PUBLISHED RESULTS IN [6] INCLUDED FOR
COMPARISON

Problem		IRLS0	IHT	FPCA	Opts	BARM
FR	n(=m)	r	FoS	FoS	FoS	FoS
0.78	500	20	0.9	0	0	1
0.8	40	9	1	0	0.5	1
0.87	100	14	0.5	0	0	1

624 combinations of these values to explore varying degrees of
625 problem difficulty, here we only reproduce the most challeng-
626 ing cases to see if BARM is still able to produce superior
627 reconstruction accuracy. In this respect problem difficulty is
628 measured by the *degrees of freedom ratio* (FR) given by FR
629 $= r(n + m - r)/p$ as defined in [6]. We also only include ex-
630 periments where algorithms are blind to the true rank of \mathbf{X}_0 .⁶
631 Results are shown in Table I, where we have also displayed
632 the published results of three additional algorithms that were
633 compared with IRLS0 in [6], namely, IHT [22], FPCA [23]
634 and Optspace [24]. From the table we observe that, in the most
635 difficult problem considered in [6], IRLS0 achieved only a 0.5
636 FoS score (meaning failure 50% of the time) while BARM still
637 achieves a perfect 1.0. Note that when FR is high, the problem
638 of recovering the underlying matrix is essentially much harder.
639 This happens in a manner that more local minima are induced
640 (due to increased rank) and/or much larger search space are
641 exposed (due to decreased number of observations/constraints).
642 In these cases, the equivalency of the global optimal with con-
643 vex relaxation usually does not hold, whereas for the existing
644 non-convex surrogates, there is no reason to assume any local
645 minima are not present. However, since BARM has an implicit
646 mechanism of smoothing local minima (though maybe not all
647 of them), it works more robustly in these situations.

648 B. General \mathbf{A}

649 Next we consider the more challenging problem involving
650 arbitrary affine constraints. The desired low-rank \mathbf{TX}_0 is gen-
651 erated in the same way as above. We then consider two types
652 of linear mappings where \mathbf{A} is generated as: (i) an iid $\mathcal{N}(0, 1)$,
653 $p \times n^2$ matrix, and (ii) $\sum_{i=1}^p i^{-1/2} \mathbf{u}_i \mathbf{v}_i^\top$, where $\mathbf{u}_i \in \mathbb{R}^p$ and
654 $\mathbf{v}_i \in \mathbb{R}^{n^2}$ are iid $\mathcal{N}(0, 1)$ vectors. The latter is meant to ex-
655 plore less-than-ideal conditions where the linear operator dis-
656 plays correlations and may be somewhat ill-conditioned. Fig. 3
657 displays aggregate results when \mathbf{X}_0 is 50×50 and 100×100 ,
658 including the underlying REL scores for additional comparison.
659 In both cases $p = 1000$ observations are used, and therefore the
660 corresponding measurement matrices \mathbf{A} are 1000×2500 and
661 1000×10000 respectively. We then vary r from 1 up to the
662 theoretical limit corresponding to problem size. Again we ob-
663 serve that BARM is consistently able to work up to the limit,
664 even when the \mathbf{A} operator is no longer an ideal Gaussian. In

⁶Note that IRLS0 can be modified to account for the true rank if such knowl-
edge were available.

665 general, we have explored a wide range of empirical conditions
666 too lengthy to report here, and it is only very rarely, and always
667 near the theoretical boundary, where BARM occasionally may
668 not succeed. We explore such failure cases in the next section.

669 C. Failure Case Analysis

670 Thus far we have not shown any cases where BARM actually
671 fails. Of course solving (1) for general \mathbf{A} is NP-hard so recovery
672 failures certainly must exist in some circumstances when using
673 a polynomial-time algorithm such as BARM. Although we cer-
674 tainly cannot explore every possible scenario, it behooves us
675 to probe more carefully for conditions under which such errors
676 may occur. One way to accomplish this is to push the problem
677 difficulty even further towards the theoretical limit by reducing
678 the number of measurements p as follows.

679 With the number of observations fixed at $p = 1000$ and a
680 general measurement matrix \mathbf{A} , the previous section examined
681 the recovery of 50×50 and 100×100 matrices as the rank was
682 varied from 1 to the recovery limit ($r = 11$ for the 50×50 case;
683 $r = 5$ for the 100×100 case). However, it is still possible to
684 make the problem even more challenging by fixing r at the limit
685 and then reducing p until it exactly equals the degrees of freedom
686 $2n^2 - r^2$. With $\{n = 50, r = 11\}$ this occurs at $p = 979$, for
687 $\{n = 100, r = 5\}$ this occurs at $p = 975$.

688 We examined the BARM algorithm under these conditions
689 with 10 additional trials using the uncorrelated \mathbf{A} for each prob-
690 lem size. Encouragingly, BARM was still 30% successful with
691 $\{n = 50, r = 11\}$, and 40% successful with $\{n = 100, r = 5\}$.
692 However, it is interesting to further examine the nature of these
693 failure cases. In Fig. 4 we have averaged the singular values of
694 $\widehat{\mathbf{X}}$ in all the failure cases. We notice that, although the recovery
695 was technically classified as a failure since the relative error
696 (REL) was above the stated threshold, the estimated matrices
697 are of almost exactly the correct minimal rank. Hence BARM
698 has essentially uncovered an alternative solution with minimal
699 rank that is nonetheless feasible by construction. We therefore
700 speculate that right at the theoretical limit, when \mathbf{A} is maxi-
701 mally overcomplete ($p \times n^2 = 979 \times 2500$ or 975×10000 for
702 the two problem sizes), there exists multiple feasible matri-
703 ces with singular value spectral cut-off points indistinguishable
704 from the optimal solution. Importantly, when the other algo-
705 rithms we tested failed, the failure is much more dramatic and
706 a clear spectral cut-off at the correct rank is not apparent.

707 This motivates a looser success criteria than FoS to account
708 for the possibility of multiple (nearly) optimal solutions that
709 may not necessarily be close with respect to relative error. For
710 this purpose we define the *frequency of rank success* (FoRS) as
711 the percentage of trials whereby a feasible solution $\widehat{\mathbf{X}}$ is found
712 such that $\sigma_r[\widehat{\mathbf{X}}]/\sigma_{r+1}[\widehat{\mathbf{X}}] > 10^3$, where $\sigma_i[\cdot]$ denotes the i -th
713 singular value of a matrix and r is the rank of the true low-rank
714 \mathbf{X}_0 . In words, FoRS measures the percentage of trials such that
715 roughly a rank r solution is recovered, regardless of proximity
716 to \mathbf{X}_0 .

717 Under this new criteria, all of the failure cases with respect to
718 FoS described above, for both problem sizes, become successes;
719 however, none of the other algorithms show improvement under

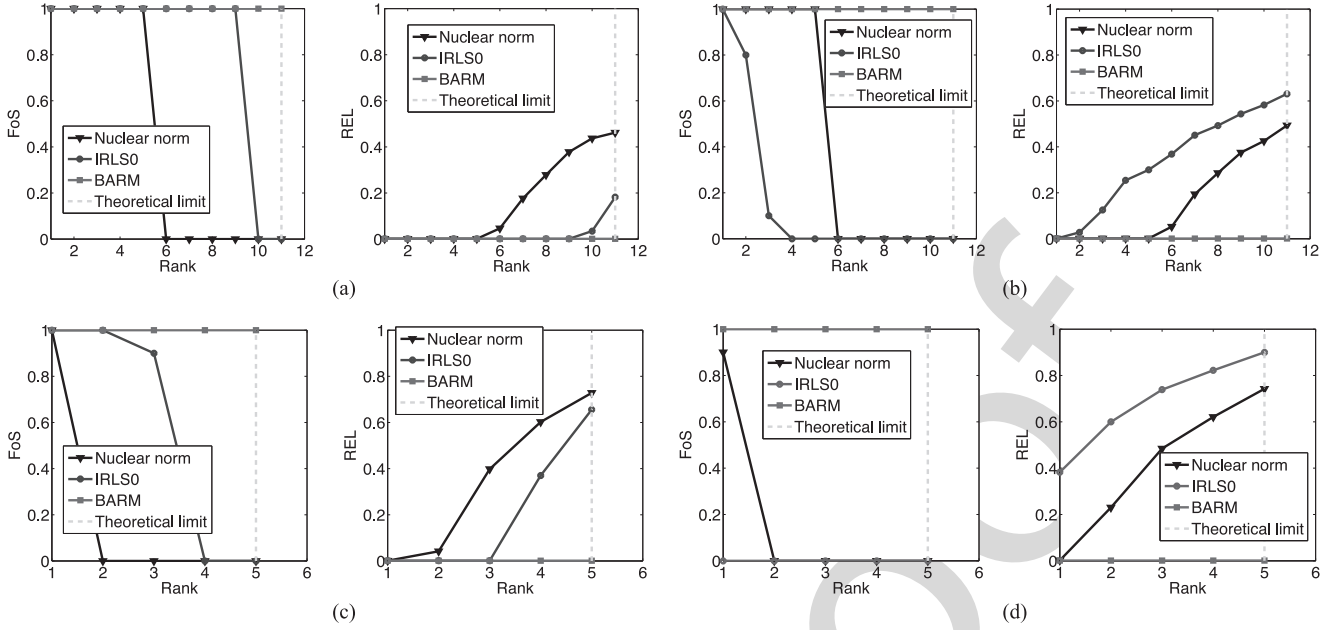


Fig. 3. Comparisons with general affine constraints (avg of 10 trials). (a) 50×50 , \mathbf{A} uncorrelated, (b) 50×50 , \mathbf{A} correlated, (c) 100×100 , \mathbf{A} uncorrelated, and (d) 100×100 , \mathbf{A} correlated.

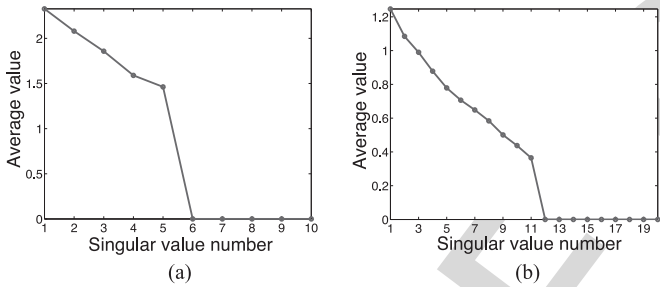


Fig. 4. Singular value averages of failure cases. In both cases solutions of minimal rank are obtained even though $\hat{\mathbf{X}} \neq \mathbf{X}_0$. (a) 50×50 and (b) 100×100 .

TABLE II
FURTHER MATRIX COMPLETION COMPARISONS OF BARM WITH IRLS0 BY REDUCING THE NUMBER OF MEASUREMENTS IN THE HARDEST PROBLEM FROM [6]. RESULTS WITH BOTH FoS AND FoRS METRICS ARE REPORTED (AVG OF 10 TRIALS)

Problem		IRLS0		BARM		
FR	$n(=m)$	r	FoS	FoRS	FoS	FoRS
0.9	100	14	0	0	1	1
0.95	100	14	0	0	0.8	1
0.99	100	14	0	0	0.7	1

720 this criteria, indicating that their original failures involved actual
 721 sub-optimal rank solutions. Something similar happens when we
 722 revisit the matrix completion experiments. For example, based
 723 on Table I the most difficult case involves $FR = 0.87$; however,
 724 by further reducing p , we can push FR towards 1.0 to further
 725 investigate the break-down point of BARM. Results are shown
 726 in Table II. While IRLS0 (which is the top performing algorithm

727 in [6] and in our experiments besides BARM) fails 100% of the
 728 time via both metrics, BARM can achieve an FoS of 0.7 even
 729 when $FR = 0.99$ and an FoRS of 1.0 in all cases.

730 We therefore adopt a more challenging measurement structure
 731 for \mathbf{A} to better evaluate the limits of BARM performance to
 732 reveal potential failures by both FoS and FoRS metrics. Specifi-
 733 cally, we first applied 2-D *discrete cosine transform* (DCT) to
 734 \mathbf{X}_0 and then randomly sampled p of the resulting DCT coef-
 735 ficients. Because both the DCT and the sampling sub-process
 736 are linear operations on the entries of \mathbf{X}_0 , the whole process is
 737 representable via a matrix \mathbf{A} , which encodes highly structured
 738 information. Fig. 5 depicts the results using problem sizes con-
 739 sistent with Fig. 3; note that the FoRS metric has replaced the
 740 REL metric for comparison purposes.

741 Two things stand out from the analysis. First, while the other
 742 algorithms display almost identical behavior under either metric,
 743 BARM failures under the FoS criteria are mostly converted to
 744 successes by the FoRS metric by recovering a matrix of near-
 745 optimal rank. Secondly, even though certain unequivocal fail-
 746 ures emerge near the limits with this challenging DCT-based
 747 sampling matrix, BARM outperforms the other algorithms using
 748 either metric by a large margin.

749 To summarize, we have demonstrated that BARM is capa-
 750 ble of recovering a low-rank matrix right up to the theoretical
 751 limit in a variety of scenarios using different types of mea-
 752 surement processes. Moreover, even in cases where it fails, it
 753 often nonetheless still produces a feasible $\hat{\mathbf{X}}$ with rank nearly
 754 identical to the generative low-rank \mathbf{X}_0 , suggesting that multi-
 755 ple optimal solutions may be possible in challenging borderline
 756 cases. But when true unequivocal failures do occur, such fail-
 757 ures tend to be near the theoretical boundary, and with greater
 758 likelihood when the dictionary displays significant structure
 759 (or correlations). While certainly we envision that, out of the

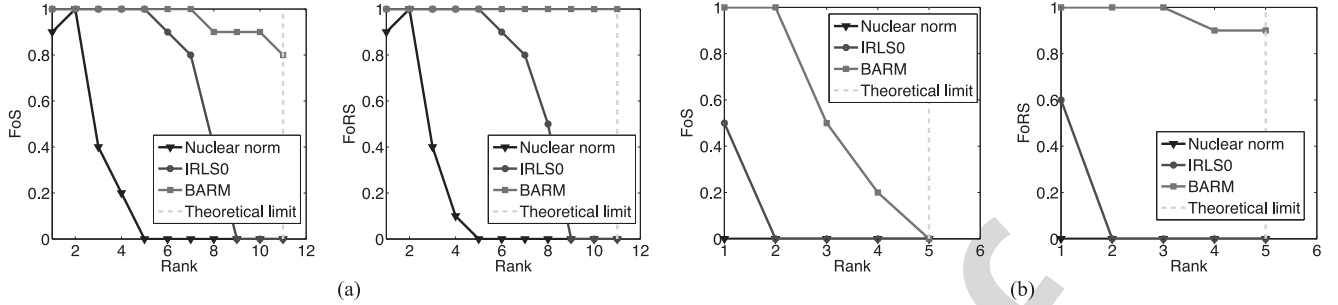


Fig. 5. Comparisons with structured affine constraints using both FoS and FoFS evaluation metrics (avg of 10 trials). (a) 50×50 , \mathbf{A} sub-sampled DCT, (b) 100×100 , \mathbf{A} sub-sampled DCT.

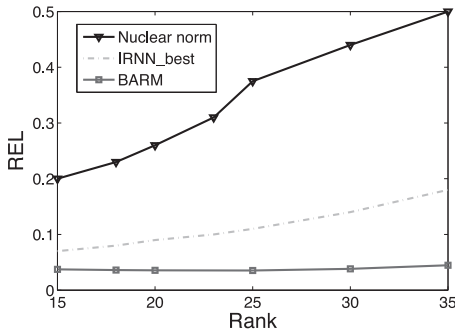


Fig. 6. Test with noisy data.

760 infinite multitude of testing situations further significant pock-
 761 ets of BARM failure can be revealed, we nonetheless feel that
 762 BARM is quite promising relative to existing algorithms.

763 D. Additional Noisy Tests

764 We also briefly present results that demonstrate the robustness
 765 of BARM to noise. For this purpose we reproduce the noisy
 766 experiment from [5] designed for validating IRNN algorithms.
 767 The simulated data are generated in the exact same way as was
 768 used to produce Fig. 2, only now instead of observing elements
 769 of \mathbf{X}_0 directly, we observe $\mathbf{X}_0 + 0.1 \times \mathbf{E}$, where elements
 770 of \mathbf{E} are iid $\mathcal{N}(0, 1)$. Although in [5] a heuristic strategy is
 771 introduced and tuned for adaptively setting all parameters (four
 772 in total), we simply applied BARM with $\lambda = 10^{-3}$ (so only a
 773 single parameter need be adjusted, and actually a wide range
 774 of λ values produces similar performance anyway). Results are
 775 shown in Fig. 6 where we compare BARM directly with the best
 776 result reported in [5] over the range $r = 15$ to $r = 35$. The
 777 nuclear norm solution is also included for reference. Overall, the
 778 BARM solution is stable and exhibits superior accuracy relative
 779 to the others.

780 E. Computational Complexity

781 Finally, regarding computational complexity, for general \mathbf{A}
 782 the BARM updates can be implemented to scale linearly in the
 783 elements of \mathbf{X} and quadratically in the number of observations
 784 p (the special case of matrix completion is decidedly much
 785 cheaper because of the special structure that can be exploited).
 786 In our experiments, for relatively easy problems on the order of

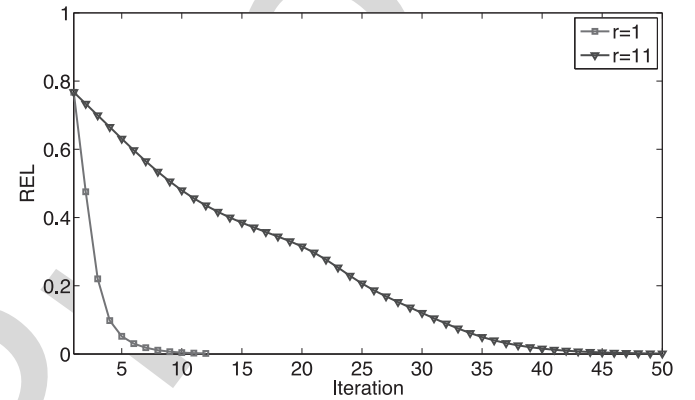


Fig. 7. Empirical convergence of BARM.

10 iterations are required, while for difficult recovery problems
 near the theoretical recovery boundary this may increase by a
 factor of 10 or so. This is somewhat expected though since as we
 near the theoretical limit, \mathbf{A} becomes highly overcomplete, and
 candidate solutions become much more difficult to differentiate.

To show this effect empirically, we compare two separate tri-
 als from Fig. 3(a), the first when $r = 1$ (relatively easy), the sec-
 ond when $r = 11$ (relatively hard).⁷ In Fig. 7 we plot the value
 of REL in both cases versus the iteration number of BARM.

796 VII. APPLICATION EXAMPLES

797 Many real-world problems from disparate fields can be for-
 798 mulated as the search for a low-rank matrix under affine con-
 799 straints [1], [3], [4], [25]. Here we briefly consider two such
 800 examples: low-rank image rectification and collaborative filter-
 801 ing for recommender systems. The former implicitly involves
 802 a general sampling operator \mathbf{A} , while the latter reduces to a
 803 standard matrix completion problem.

804 A. Low-Rank Image Rectification

805 In [4], the *transform invariant low-rank textures* (TILT) al-
 806 gorithm is derived for rectifying images containing low-rank

⁷Note that $r = 1$ is only relatively easy here because the number of obser-
 vations is sufficient for the larger $r = 11$ case; if only the minimal number
 of measurements are available then even $r = 1$ can be challenging for many
 algorithms.

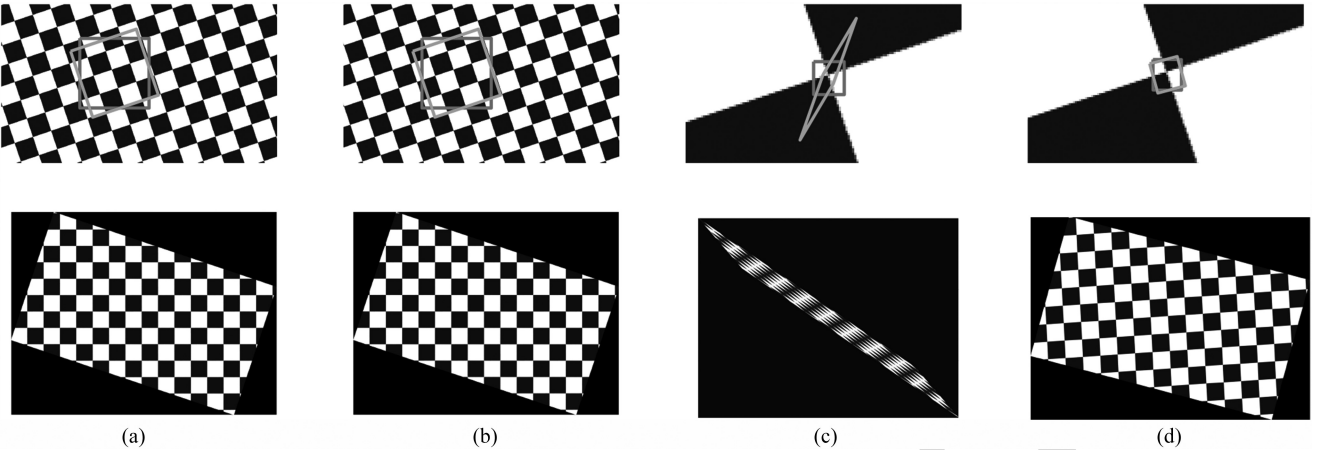


Fig. 8. Image rectification comparisons using a checkboard image. *Top*: Original image with observed region (red box) and estimated transformation (green box). *Bottom*: Rectified image estimates. (a) Nuclear norm (easy), (b) BARM (easy), (c) Nuclear norm (hard), (d) BARM (hard).

807 textures that have been transformed using an unknown operator
 808 τ from some group (e.g., a homography). For a given observed
 809 image \mathbf{Y} , the basic idea is to construct a first-order Taylor series
 810 approximation around the current rectified image estimate $\widehat{\mathbf{X}}$
 811 and solve

$$\min_{\mathbf{X}, \delta} \text{rank}[\mathbf{X}] \text{ s.t. } \mathbf{X} = \mathbf{Y} + \sum_i \mathbf{J}_i(\widehat{\mathbf{X}}) \delta_i, \quad (21)$$

812 where $\mathbf{J}_i(\widehat{\mathbf{X}})$ is the Jacobian matrix with respect to \mathbf{X} of
 813 the i -th parameter τ_i describing the transformation, with $\tau =$
 814 $[\tau_1, \tau_2, \dots]^\top$. Optimization over the vector of first-order differ-
 815 ences $\delta = [\delta_1, \delta_2, \dots]^\top$ can be accomplished in closed form by
 816 projecting both sides of the constraint to the orthogonal comple-
 817 ment of the span of all $\mathbf{J}_i(\widehat{\mathbf{X}})$. Let P_{J^c} represent this projection
 818 operator. The feasible region in (21) then becomes

$$P_{J^c}(\mathbf{X}) = P_{J^c}(\mathbf{Y}) + P_{J^c} \left(\sum_i \mathbf{J}_i(\widehat{\mathbf{X}}) \delta_i \right) = P_{J^c}(\mathbf{Y}) \quad (22)$$

819 The resulting problem then reduces exactly to (1) when we
 820 define $\mathcal{A} = P_{J^c}$ and $\mathbf{b} = \text{vec}[P_{J^c}(\mathbf{Y})]$. Once \mathbf{X} is computed in
 821 this way, we then update each $\mathbf{J}_i(\widehat{\mathbf{X}})$ and repeat until conver-
 822 gence.

823 While the original TILT algorithm substitutes the nuclear
 824 norm for $\text{rank}[\mathbf{X}]$, we embedded the BARM algorithm into
 825 the posted TILT source code [4] for comparison purposes (note
 826 that we disabled an additional sparse error term for both algo-
 827 rithms to simplify comparisons, and it is not necessary anyway
 828 in many regimes). Figs. 8 and 9 display results on both two
 829 easy examples, where the number of observations p is large,
 830 and two more difficult problems where the number observa-
 831 tions is small. While both algorithms succeed on the easy cases,
 832 when the observations are constrained by a small image window,
 833 only BARM is successful in accurately rectifying the images.
 834 This may be due, at least in part, to the fact that the implicit
 835 \mathcal{A} operator contains significant structure that is not consistent
 836 with the required nullspace properties required for nuclear norm
 837 minimization success.

B. Collaborative Filtering of MovieLens Data

838

839 Collaborative filtering, a technique used by many recom- 839
 840 mender systems, is a popular representative application of low- 840
 841 rank matrix completion. Typically the rows (or columns) of \mathbf{X}_0 841
 842 index users, the columns (or rows) denote items, and each entry 842
 843 $(\mathbf{X}_0)_{ij}$ is the rating/score of user i applied to item j . Given 843
 844 that we can observe some subset of elements of \mathbf{X}_0 , the task 844
 845 of collaborative filtering is to predict all or some of the miss- 845
 846 ing ratings. In general this would be impossible; however, if we 846
 847 have access to some prior knowledge, e.g., \mathbf{X}_0 is low-rank, then 847
 848 estimation may be feasible. 848

849 While our interest here is not in recommender systems or 849
 850 collaborative filtering per se, we nonetheless evaluate BARM 850
 851 using the 1M MovieLens dataset⁸ as this appears to represent 851
 852 one of the most common evaluation benchmarks. We emphasize 852
 853 at the outset that the strict validity of any low-rank assumptions 853
 854 underlying this data is debatable, and it remains entirely unclear 854
 855 whether the true globally optimal or lowest rank solution consis- 855
 856 tent with the observations, even if computable, would necessar- 856
 857 ily lead to the best prediction of the unknown ratings. In fact, the 857
 858 reported performance of various existing rank-minimization algo- 858
 859 rithms tends to cluster around almost the same value, implying 859
 860 that collaborative filtering may not provide the most discrimina- 860
 861 tive data type with which to compare. In most cases, it appears 861
 862 that tuning parameters and other heuristic modifications play 862
 863 a larger role than the underlying algorithmic distinctions fun- 863
 864 damental to finding optimal low-rank estimates. Nonetheless, 864
 865 we apply BARM for completeness and convention, adopting an 865
 866 additional simple mean-offset estimation term from [25] that is 866
 867 particularly suitable for this problem. 867

868 In [6], IRLS0 is compared with only two other algorithms on 868
 869 MovieLens data, but the performance is no better. Therefore, 869
 870 we choose to compare directly with [25], which both derives 870
 871 an IRLS-like algorithm and shows comparisons with a much 871
 872 wider variety of alternative algorithms using a strict evalua- 872
 873 tion protocol that is standard in the literature. Specifically, the 873

⁸<http://www.grouplens.org/>

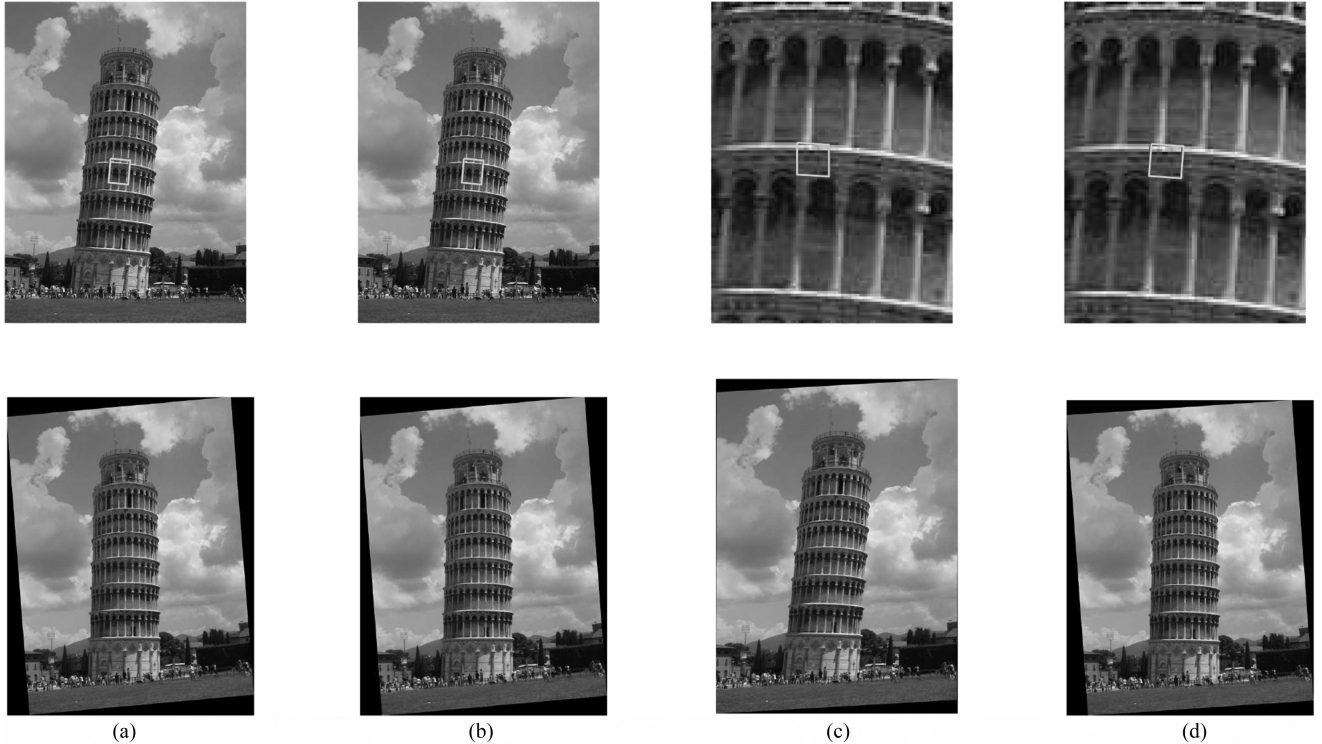


Fig. 9. Image rectification comparisons using a landmark photo. *Top*: Original image with observed region (red box) and estimated transformation (green box). *Bottom*: Rectified image estimates. (a) Nuclear norm (easy), (b) BARM (easy), (c) Nuclear norm (hard), (d) BARM (hard).

874 1M MovieLens dataset, which contains 1 million ratings in the
 875 range $\{1, \dots, 5\}$ for 3900 movies from 6040 unique users, is
 876 assessed under two test-protocols: *weak generalization*, which
 877 measures the ability to predict other items rated by the same
 878 user, and *strong generalization*, which measures the ability to
 879 predict items by novel users. 5 000 users are randomly selected
 880 for the weak generalization, and likewise 1 000 users are ex-
 881 tracted for the strong generalization. Each experiment is then
 882 run three times and the averaged results are reported. The per-
 883 formance metric is *normalized mean absolute error* (NMAE)
 884 given as

$$\text{NMAE} = \frac{\left(\sum_{i,j \in \text{supp}(\mathbf{X}_0)} \frac{|(\mathbf{X}_0)_{ij} - \hat{\mathbf{X}}_{ij}|}{|\text{supp}(\mathbf{X}_0)|} \right)}{(rt_{\max} - rt_{\min})},$$

885 where rt_{\max} and rt_{\min} are the maximum and minimum ratings
 886 possible.

887 We followed the same setup and reported results using BARM
 888 in Table III along with results from [25] for comparison. This
 889 includes the additional algorithms URP [26], Attitude [27],
 890 MMMF [28], IPCF [29], E-MMMF [30], GPLVM [31], NBMC
 891 [32], and IRLS/GM [25], [6]. From this table we observe that
 892 for the easier weak generalization problem BARM is a close
 893 second best, while for the more challenging strong generaliza-
 894 tion BARM is actually the best. Of course it is also immediately
 895 apparent that all algorithms fall within a relatively narrow per-
 896 formance range of approximately five percentage points. Con-
 897 sequently, we cannot unequivocally conclude that the attributes
 898 of BARM which make it suitable for optimally minimizing rank

TABLE III
 COLLABORATIVE FILTERING ON 1M MOVIELENS DATASET. RESULTS FROM
 [25] ARE IN ITALIC FOR COMPARISON PURPOSES

	Weak NMAE	Hard NMAE
<i>URP</i>	0.4341	0.4444
<i>Attitude</i>	0.4320	0.4375
<i>MMMF</i>	0.4156	0.4203
<i>IPCF</i>	0.4096	0.4113
<i>E-MMMF</i>	0.4029	0.4071
<i>GPLVM</i>	0.4026	0.3994
<i>NBMC</i>	0.3916	0.3992
<i>IRLS/GM</i>	0.3959	0.3928
BARM	0.3942	0.3898

899 necessarily translate into a truly significant practical advantage
 900 on this collaborative filtering task. But we would argue that the
 901 same holds for any matrix completion algorithm.

VIII. CONCLUSION

902 This paper explores a conceptually-simple, parameter-free
 903 algorithm called BARM for matrix rank minimization under
 904 affine constraints that is capable of successful recovery empir-
 905 ically observed to approach the theoretical limit over a broad
 906 class of experimental settings (including many not shown here)
 907 unlike any existing algorithms, and long after any convex guar-
 908 antees break down. Our strategy in this effort has been to
 909 adopt Bayesian machinery for inspiring a principled cost func-
 910 tion; however, ultimate model justification is placed entirely in
 911

912 theoretical evaluation of desirable global and local minima prop-
 913 erties, and in the empirical recovery performance that inevitably
 914 results from these properties. Although in general non-convex
 915 algorithms are exponentially more challenging to analyze, in
 916 this regard we have at least attempted to contextualize BARM
 917 in the same manner as convex optimization-based approaches
 918 such as nuclear-norm minimization.

919 APPENDIX A

920 Here we provide brief proofs of Lemmas 1 and 2 as well as
 921 Theorem 1. We also address the augmented update rules that
 922 account for the revised, symmetrized cost function discussed in
 923 Section V.

924 A. Proof of Lemmas 1 and 2

925 Regarding Lemma 1, this result mirrors related ideas from
 926 [16] in the context of Bayesian compressive sensing. Hence,
 927 while a more rigorous presentation is possible, here we de-
 928 scribe the basic aspects of the adaptation. At any candidate
 929 minimizer of (10) in the limit $\lambda \rightarrow 0$, define \mathbf{W} such that
 930 $\mathbf{A}\bar{\Psi}\mathbf{A}^\top = \mathbf{W}\mathbf{W}^\top$. To be a minimizer, global or local, it must
 931 be that $\mathbf{b} \in \text{span}[\mathbf{W}]$. If this were not the case, then $\mathcal{L}(\Psi, \nu)$
 932 would diverge to infinity as $\lambda \rightarrow 0$ because $\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b}$ progresses
 933 to infinity at a faster rate than $\log |\Sigma_b|$ can compensate by ap-
 934 proaching minus infinity. Intuitively, in much the same way
 935 $\text{argmin}_z \frac{1}{z} + \log z = 1$, meaning the optimal z must lie in the
 936 ‘span’ of 1 else the overall objective will be driven to infinity.

937 Consequently, the only way to minimize the cost in the limit
 938 as $\lambda \rightarrow 0$ is to consider low-rank solutions within the constraint
 939 set that $\mathbf{b} \in \text{span}[\mathbf{W}]$, and it is equivalent to requiring that
 940 $\mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq C$ for some constant C independent of λ (which
 941 ultimately corresponds with maintaining $\mathcal{A}(\mathbf{X}) = \mathbf{b}$ in the limit
 942 as well).

943 In this setting, while $0 \leq \mathbf{b}^\top \Sigma_b^{-1} \mathbf{b} \leq C$ is bounded, the sec-
 944 ond term in $\mathcal{L}(\Psi, \nu)$ can be unbounded from below when
 945 $\text{rank}[\Psi]$ is sufficiently small. To see this note that

$$\log |\Sigma_b| = \sum_{i=1}^p \log (\sigma_i [\mathbf{A}\bar{\Psi}\mathbf{A}^\top] + \lambda), \quad (23)$$

946 where $\sigma_i [\cdot]$ denotes the i -th singular value of a matrix. While
 947 the maximum rank of $\mathbf{A}\bar{\Psi}\mathbf{A}^\top$ is obviously p , if $r \triangleq \text{rank}[\Psi] <$
 948 p/m and $\text{spark}[\mathbf{A}] = p + 1$ (maximal spark) as stipulated in the
 949 lemma statement, then $\text{rank}[\mathbf{A}\bar{\Psi}\mathbf{A}^\top] = mr$ and (23) becomes

$$\log |\Sigma_b| = \sum_{i=1}^{mr} \log (\sigma_i [\mathbf{A}\bar{\Psi}\mathbf{A}^\top] + \lambda) + (p - mr) \log \lambda. \quad (24)$$

950 Note that the spark assumption accomplishes two objectives
 951 in this context. First, it guarantees that a high rank Ψ cannot
 952 masquerade as a low rank Ψ behind the nullspace of some col-
 953 lection of columns \mathbf{A}_i . Secondly, it ensures that after assuming
 954 $r < p/m$, then $\text{rank}[\mathbf{A}\bar{\Psi}\mathbf{A}^\top] = mr$.

955 Consequently, in the limit where $\lambda \rightarrow 0$ (with the limit being
 956 taken outside of the minimization), (23) effectively scales as
 957 $(p - mr) \log \lambda$, and hence the overall cost is minimized when

Ψ has minimal rank. This in turn ensures that the corresponding
 \mathbf{X} will also have minimal rank, completing the proof sketch for
 Lemma 1.

Finally, Lemma 2 follows directly from the structure of the
 $\mathcal{L}(\Psi, \nu)$ cost function via simple reparameterizations. ■

963 B. Proof of Theorem 1

964 To begin we assume that $\mathbf{b}_i \neq 0, \forall i$, where \mathbf{b}_i denotes the
 965 sub-vector of \mathbf{b} such that $\mathbf{b}_i = \mathbf{A}_i \mathbf{x}_i$. If this were not the case
 966 we can always collapse \mathbf{X} by the corresponding column (which
 967 is indistinguishable from zero) and achieve an equivalent result.
 968 Given the assumptions of Theorem 1, the BARM cost function
 969 becomes

$$\mathcal{L}(\Psi, \nu) = \sum_{i=1}^m \mathbf{b}_i^\top (\nu_i \mathbf{A}_i \Psi \mathbf{A}_i^\top)^{-1} \mathbf{b}_i + \log |\nu_i \mathbf{A}_i \Psi \mathbf{A}_i^\top|. \quad (25)$$

970 If there exists a feasible rank one solution to $\mathbf{b} = \text{Avec}$
 971 $[\mathbf{X}]$, then there also exists a set of $\Psi'_i = \nu_i \Psi$ such that $\mathbf{b}_i \mathbf{b}_i^\top =$
 972 $\mathbf{A}_i \Psi'_i \mathbf{A}_i^\top$ for all i . To see this, note that $\mathbf{b}_i \mathbf{b}_i^\top = \mathbf{A}_i \mathbf{x}_i \mathbf{x}_i^\top$
 973 \mathbf{A}_i^\top . Because $\text{rank}[\mathbf{X}] = 1$, it also follows that $\mathbf{b}_i \mathbf{b}_i^\top = \alpha_i \mathbf{A}_i \mathbf{X}$
 974 $\mathbf{X}^\top \mathbf{A}_i^\top$, where $\alpha_i = \|\mathbf{x}_i \mathbf{x}_i^\top\| / \|\mathbf{X} \mathbf{X}^\top\|$. Therefore $\Psi'_i =$
 975 $\nu_i \mathbf{X} \mathbf{X}^\top$ achieves the desired result with $\nu_i = \alpha_i$.

976 Now suppose we have converged to any solution $\{\widehat{\Psi}, \widehat{\nu}\}$ with
 977 $\text{rank}[\widehat{\Psi}] > 1$ and associated $\widehat{\Sigma} = \mathbf{I} \otimes \widehat{\Psi}$. Note that since $\mathbf{b}_i \neq$
 978 $0, \nu_i > 0$ for all i , otherwise a local minimum is not possible
 979 (the cost function would be driven to positive infinity).

980 Define $\widehat{\Sigma}_{b_i} = \widehat{\nu}_i \mathbf{A}_i \widehat{\Psi} \mathbf{A}_i^\top$. Additionally we can assume that
 981 $\mathbf{b}_i^\top \widehat{\Sigma}_{b_i}^{-1} \mathbf{b}_i$ is finite, meaning that \mathbf{b}_i lies in the span of the singular
 982 vectors of $\widehat{\Sigma}_{b_i}$. (If this were not the case, the cost would be
 983 driven to infinity and we could not be at a minimizing solution
 984 anyway.) If $\{\widehat{\Psi}, \widehat{\nu}\}$ is a local minimum, then $\{\lambda_1 = 1, \lambda_2 = 0\}$
 985 must be a local minimum of the revised cost function

$$\mathcal{L}(\lambda_1, \lambda_2) = \sum_{i=1}^m \mathbf{b}_i^\top (\lambda_1 \widehat{\Sigma}_{b_i} + \lambda_2 \mathbf{b}_i \mathbf{b}_i^\top)^{-1} \mathbf{b}_i + \log |\lambda_1 \widehat{\Sigma}_{b_i} + \lambda_2 \mathbf{b}_i \mathbf{b}_i^\top|. \quad (26)$$

986 This is because $\mathbf{b}_i \mathbf{b}_i^\top$ represents a valid set of basis vectors for
 987 updating the covariance per the construction above involving
 988 Ψ'_i . First consider optimization over λ_1 . If $\lambda_1 = 1$ is a local
 989 minimum, then by taking gradients and equating to zero, we
 990 require that

$$\sum_{i=1}^m \mathbf{b}_i^\top \widehat{\Sigma}_{b_i}^{-1} \mathbf{b}_i = \sum_{i=1}^m \text{rank}[\widehat{\Sigma}_{b_i}]. \quad (27)$$

991 Likewise, taking the gradient with respect to λ_2 we obtain

$$\left. \frac{\partial \mathcal{L}(\lambda_1, \lambda_2)}{\partial \lambda_2} \right|_{\lambda_1=1, \lambda_2=0} = \sum_{i=1}^m \mathbf{b}_i^\top \widehat{\Sigma}_{b_i}^{-1} \mathbf{b}_i - \sum_{i=1}^m (\mathbf{b}_i^\top \widehat{\Sigma}_{b_i}^{-1} \mathbf{b}_i)^2. \quad (28)$$

992 The nullspace condition (a very mild assumption) ensures
 993 that $\sum_{i=1}^m \text{rank}[\widehat{\Sigma}_{b_i}] = k$ for some $k > m$ when $\text{rank}[\Psi] > 1$.
 994 To see this, observe that to achieve $\sum_{i=1}^m \text{rank}[\widehat{\Sigma}_{b_i}] = m$ when
 995 $\text{rank}[\Psi] > 1$ requires that $\Psi = \mathbf{u} \mathbf{u}^\top + \mathbf{W} \mathbf{W}^\top$ where \mathbf{u} is a

996 vector and \mathbf{W} is a matrix (or vector) with columns in $\text{null}[\mathbf{A}_i]$,
 997 $\forall i$. If any such \mathbf{W} is not in this nullspace for some i , then given
 998 that $p_i > 1$, the associated $\mathbf{A}_i \Psi \mathbf{A}_i^\top$ will have rank greater than
 999 one, and the overall rank sum will exceed m .

1000 Consequently, (28) will always be negative. This is because
 1001 if $\sum_{i=1}^m z_i = k$ for any set of non-negative variables $\{z_i\}$, the
 1002 minimal value of $\sum_{i=1}^m z_i^2$ occurs when $z_i = k/m, \forall i$. In our
 1003 case, this implies that

$$\sum_{i=1}^m \left(\mathbf{b}_i^\top \widehat{\Sigma}_{\mathbf{b}_i}^{-1} \mathbf{b}_i \right)^2 \geq \sum_{i=1}^m (k/m)^2 > k > m. \quad (29)$$

1004 Therefore we can add a small contribution of $\mathbf{b}_i \mathbf{b}_i^\top$ to each
 1005 $\widehat{\Sigma}_{\mathbf{b}_i}$ and reduce the underlying cost function. Hence we cannot
 1006 have a local minimum, except when Ψ is equal to some Ψ^*
 1007 with $\text{rank}[\Psi^*] = 1$. Moreover, we may directly conclude that
 1008 $\mathbf{x}^* = \overline{\Psi}^* \mathbf{A}^\top (\mathbf{A} \overline{\Psi}^* \mathbf{A}^\top)^\dagger \mathbf{b}$ is feasible and $\text{rank}[\mathbf{X}^*] = 1$ with
 1009 $\mathbf{x}^* = \text{vec}[\mathbf{X}^*]$.

1010 Regarding the last part of the theorem, we consider only
 1011 f that are concave non-decreasing functions (this is the only
 1012 reasonable choice for shrinking singular values to zero, and
 1013 the more general case naturally follows anyway with additional
 1014 effort, but minimal enlightenment). Without loss of generality
 1015 we may also assume that $f(0) = 0$ and $f(1) = 1$; we can always
 1016 apply an inconsequential translation and scaling such that these
 1017 conditions hold.⁹ Simple counter examples then demonstrate
 1018 that $f(\epsilon)$ must be greater than some constant C independent of
 1019 ϵ for all ϵ sufficiently small. To see this, note that we can always
 1020 rescale elements of \mathbf{A} such that a solution with rank greater
 1021 than one is preferred unless this condition holds. However, such
 1022 an f , which effectively must display infinite gradient at $f(0)$ to
 1023 guarantee a global solution is always rank one, will then always
 1024 display local minima for certain \mathbf{A} . This can easily be revealed
 1025 through simple counter-examples. ■

1026 C. Symmetrization Update Rules

1027 These iterative update rules follow from alternative upper
 1028 bounds tailored to the symmetric version of BARM. When both
 1029 Ψ_r and Ψ_c are fixed, \mathbf{x} is updated via the posterior mean cal-
 1030 culation

$$\begin{aligned} \widehat{\mathbf{x}} &= \text{vec} \left[\widehat{\mathbf{X}} \right] = \frac{1}{2} (\overline{\Psi}_r + \overline{\Psi}_c) \mathbf{A}^\top \\ &\quad \times \left[\lambda \mathbf{I} + \mathbf{A} \frac{1}{2} (\overline{\Psi}_r + \overline{\Psi}_c) \mathbf{A}^\top \right]^{-1} \mathbf{b}. \end{aligned} \quad (30)$$

1031 where $\overline{\Psi}_r = \Psi_r \otimes \mathbf{I}$ and $\overline{\Psi}_c = \mathbf{I} \otimes \Psi_c$. Likewise we update
 1032 $\nabla_{\Psi_r^{-1}}$ and $\nabla_{\Psi_c^{-1}}$ using

$$\nabla_{\Psi_r^{-1}} = \sum_{i=1}^m \Psi_r - \Psi_r \mathbf{A}_{r_i}^\top (\mathbf{A} \overline{\Psi}_r \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_{r_i} \Psi_r, \quad (31)$$

$$\nabla_{\Psi_c^{-1}} = \sum_{i=1}^n \Psi_c - \Psi_c \mathbf{A}_{c_i}^\top (\mathbf{A} \overline{\Psi}_c \mathbf{A}^\top + \lambda \mathbf{I})^{-1} \mathbf{A}_{c_i} \Psi_c, \quad (32)$$

⁹The log function is a limiting case, but what follows holds nonetheless.

where $\mathbf{A}_{r_i} \in \mathbb{R}^{p \times m}$ is defined such that $\mathbf{A} = [\mathbf{A}_{r_1}^\top, \dots, \mathbf{A}_{r_m}^\top]^\top$ 1033
 and $\mathbf{A}_{c_i} \in \mathbb{R}^{p \times m}$ is defined such that $\mathbf{A} = [\mathbf{A}_{c_1}, \dots, \mathbf{A}_{c_n}]$. Fi- 1034
 nally given these values, with \mathbf{X} , $\nabla_{\Psi_r^{-1}}$ and $\nabla_{\Psi_c^{-1}}$ fixed, we can 1035
 compute the optimal Ψ_r and Ψ_c in closed form by optimizing 1036
 the relevant Ψ_r - and Ψ_c -dependent terms via 1037

$$\Psi_r^{\text{opt}} = \frac{1}{n} \left[\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi_r^{-1}} \right], \quad (33)$$

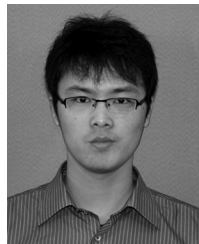
$$\Psi_c^{\text{opt}} = \frac{1}{m} \left[\widehat{\mathbf{X}} \widehat{\mathbf{X}}^\top + \nabla_{\Psi_c^{-1}} \right]. \quad (34)$$

In practice the simple initialization $\Psi_r = \mathbf{I}$ and $\Psi_c = \mathbf{I}$ is 1038
 sufficient for obtaining good performance. 1039

1040 REFERENCES

- [1] E. J. Candès and B. Recht, "Exact matrix completion via convex optimiza- 1041
tion," *Found. Comput. Math.*, vol. 9, no. 6, pp. 717–772, 2009. 1042
- [2] Y. Hu, D. Zhang, J. Ye, X. Li, and X. He, "Fast and accurate matrix com- 1043
pletion via truncated nuclear norm regularization," *IEEE Trans. Pattern 1044
Anal. Mach. Intell. (PAMI)*, vol. 35, no. 9, pp. 2117–2130, 2013. 1045
- [3] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma, "Robust recovery of 1046
subspace structures by low-rank representation," *IEEE Trans. Pattern 1047
Anal. Mach. Intell. (PAMI)*, vol. 35, no. 1, pp. 171–184, 2013. 1048
- [4] Z. Zhang, A. Ganesh, X. Liang, and Y. Ma, "Tilt: Transform invariant 1049
low-rank textures," *Int. J. Comput. Vis. (IJCV)*, vol. 99, no. 1, pp. 1–24, 1050
2012. 1051
- [5] C. Lu, J. Tang, S. Yan, and Z. Lin, "Generalized nonconvex nonsmooth 1052
low-rank minimization," in *Proc. IEEE Conf. Comput. Vis. Pattern Recog- 1053
nit. (CVPR)*, 2014. 1054
- [6] K. Mohan and M. Fazel, "Iterative reweighted algorithms for matrix rank 1055
minimization," *J. Mach. Learn. Res. (JMLR)*, vol. 13, no. 1, pp. 3441– 1056
3473, 2012. 1057
- [7] Z. Li, J. Liu, Y. Jiang, J. Tang, and H. Lu, "Low rank metric learning for 1058
social image retrieval," in *Pro. 20th ACM Int. Conf. Multimedia*, 2012, 1059
pp. 853–856. 1060
- [8] M. Tipping and C. Bishop, "Probabilistic principal component analysis," 1061
J. Roy. Statist. Soc. B, vol. 61, no. 3, pp. 611–622, 1999. 1062
- [9] B. Xin and D. Wipf, "Pushing the limits of affine rank minimization by 1063
adapting probabilistic pca," in *Proc. 32nd Int. Conf. Mach. Learn. (ICML)*, 1064
2015, pp. 419–427. 1065
- [10] P. Jain, P. Netrapalli, and S. Sanghavi, "Low-rank matrix completion 1066
using alternating minimization," in *Proc. 45th Annu. ACM Symp. Theory 1067
Comput.*, 2013, pp. 665–674. 1068
- [11] S. D. Babacan, M. Luessi, R. Molina, and A. K. Katsaggelos, "Sparse 1069
Bayesian methods for low-rank matrix estimation," *IEEE Trans. Signal 1070
Process.*, vol. 60, no. 8, pp. 3964–3977, 2012. 1071
- [12] X. Ding, L. He, and L. Carin, "Bayesian robust principal component 1072
analysis," *IEEE Trans. Image Process.*, vol. 20, no. 12, pp. 3419–3430, 1073
2011. 1074
- [13] D. Wipf, "Non-convex rank minimization via an empirical Bayesian ap- 1075
proach," in *Proc. 28th Conf. Uncertainty Artif. Intell. (UAI)*, 2012. 1076
- [14] V. Chandrasekaran, B. Recht, P. A. Parrilo, and A. S. Willsky, "The convex 1077
geometry of linear inverse problems," *Found. Comput. Math.*, vol. 12, 1078
no. 6, pp. 805–849, 2012. 1079
- [15] M. E. Tipping, "Sparse Bayesian learning and the relevance vector ma- 1080
chine," *J. Mach. Learn. Res. (JMLR)*, vol. 1, pp. 211–244, 2001. 1081
- [16] D. P. Wipf, B. D. Rao, and S. Nagarajan, "Latent variable Bayesian models 1082
for promoting sparsity," *IEEE Trans. Inf. Theory*, vol. 57, no. 9, pp. 6236– 1083
6255, 2011. 1084
- [17] D. L. Donoho and M. Elad, "Optimally sparse representation in general 1085
(nonorthogonal) dictionaries via l_1 minimization," *Proc. Nat. Acad. Sci.*, 1086
vol. 100, no. 5, pp. 2197–2202, 2003. 1087
- [18] E. J. Candès and X. Li, "Solving quadratic equations via phaselift when 1088
there are about as many equations as unknowns," *Found. Comput. Math.*, 1089
vol. 14, no. 5, pp. 1017–1026, 2014. 1090
- [19] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski, "Phase retrieval 1091
via matrix completion," *SIAM Rev.*, vol. 57, no. 2, pp. 225–251, 2015. 1092
- [20] W. I. Zangwill, *Nonlinear Programming: A Unified Approach*, Englewood 1093
Cliffs, NJ, USA: Prentice-Hall, 1969. 1094

- 1095 [21] J. Tanner and K. Wei, "Normalized iterative hard thresholding for matrix
1096 completion," *SIAM J. Scientif. Comput.*, vol. 35, no. 5, pp. S104–S125,
1097 2013.
- 1098 [22] P. Jain, R. Meka, and I. S. Dhillon, "Guaranteed rank minimization via
1099 singular value projection," in *Proc. Adv. Neural Inf. Process. Syst.*, 2010,
1100 pp. 937–945.
- 1101 [23] D. Goldfarb and S. Ma, "Convergence of fixed-point continuation algo-
1102 rithms for matrix rank minimization," *Found. Comput. Math.*, vol. 11, no.
1103 2, pp. 183–210, 2011.
- 1104 [24] R. H. Keshavan and S. Oh, "A gradient descent algorithm on the
1105 Grassman manifold for matrix completion," 2009 DOI: arXiv Preprint
1106 arXiv:0910.5260.
- 1107 [25] F. Léger, G. Yu, and G. Sapiro, "Efficient matrix completion with Gaussian
1108 models," 2010 DOI: arXiv Preprint arXiv:1010.4050.
- 1109 [26] B. Marlin, *Collaborative filtering: A machine learning perspective*, Ph.D.
1110 dissertation, Univ. of Toronto, Toronto, Canada ON, 2004.
- 1111 [27] B. M. Marlin, "Modeling user rating profiles for collaborative filtering,"
1112 in *Proc. Adv. Neural Inf. Process. Syst.*, 2003.
- 1113 [28] J. D. Rennie and N. Srebro, "Fast maximum margin matrix factorization
1114 for collaborative prediction," in *Proc. 22nd ACM Int. Conf. Mach. Learn.*,
1115 2005, pp. 713–719.
- 1116 [29] S.-T. Park and D. M. Pennock, "Applying collaborative filtering techniques
1117 to movie search for better ranking and browsing," in *Proc. 13th ACM
1118 SIGKDD Int. Conf. Knowl. Discov. Data Min.*, 2007, pp. 550–559.
- 1119 [30] D. DeCoste, "Collaborative prediction using ensembles of maximum margin
1120 matrix factorizations," in *Proc. 23rd ACM Int. Conf. Mach. Learn.*,
1121 2006, pp. 249–256.
- 1122 [31] N. D. Lawrence and R. Urtasun, "Non-linear matrix factorization with
1123 gaussian processes," in *Proc. 26th Annu. ACM Int. Conf. Mach. Learn.*,
1124 2009, pp. 601–608.
- 1125 [32] M. Zhou, C. Wang, M. Chen, J. Paisley, D. Dunson, and L. Carin, "Non-
1126 parametric Bayesian matrix completion," *Proc. IEEE SAM*, 2010.



1127
1128
1129
1130
1131
1132
1133

Bo Xin (M'XX) received the B.S. degree in electronic engineering from Dalian University of Technology, China, in 2011. He is currently working toward the Ph.D. degree in computer science at Peking University, China. His research interests include optimization, machine learning and computer vision.



and the 2006 NIPS Outstanding Paper Award.

David Wipf (M'XX) received the B.S. degree with highest honors from the University of Virginia, and the Ph.D. degree from UC San Diego, where he was an NSF IGERT Fellow. Later he was an NIH Post-doctoral Fellow at UC San Francisco. Since 2011 he has been with Microsoft Research in Beijing. His research interests include Bayesian learning techniques applied to signal/image processing and computer vision. He is the recipient of several awards including the 2012 Signal Processing Society Best Paper Award, the Biomag 2008 Young Investigator Award, and the 2006 NIPS Outstanding Paper Award.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146

Q2



Yizhou Wang (M'XX) received his Ph.D. in computer science from University of California at Los Angeles (UCLA) in 2005. He was a Research Staff of the Palo Alto Research Center (Xerox-PARC) from 2005 to 2008. He is currently a Professor of the Computer Science Department at Peking University (PKU), China. His research interests include computer vision, statistical modeling and learning.

1147
1148
1149
1150
1151
1152
1153
1154
1155



Wen Gao (F'XX) received M.S. degree in computer science from Harbin Institute of Technology in 1985, and Ph.D. degree in electronics engineering from the University of Tokyo in 1991. He was a Professor in computer science at Harbin Institute of Technology from 1991 to 1995 and a Professor in computer science at Institute of Computing Technology of Chinese Academy of Sciences from 1996 to 2005. He is currently a Professor at the School of Electronics Engineering and Computer Science, Peking University, China. He has been leading research efforts to develop systems and technologies for video coding, face recognition, sign language recognition and synthesis, and multimedia retrieval. He earned many awards, which include five national awards for his research achievements and activities. He did many services to academic society, such as general co-chair of IEEE ICME07, and the head of Chinese delegation to the Moving Picture Expert Group (MPEG) of International Standard Organization (ISO). Since 1997, he is also the chairman of the working group responsible for setting a national Audio Video coding Standard (AVS) for China. He published four books and over 500 technical articles in refereed journals and proceedings in the areas of signal processing, image and video communication, computer vision, multimodal interface, pattern recognition, and bioinformatics.

1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178

Q3

- 1180 Q1. Author: For all conference paper references, provide page numbers if printed in proceeding or location of where conference
1181 was presented if not printed.
- 1182 Q2. Author: Please provide a forward facing headshot for Dr. Wipf. Otherwise the bio will have to run without photo.
- 1183 Q3. Author: Please provide initial year(s) of IEEE membership grade(s) for all authors.

IEEE Proof