

IMPLICIT BIAS OF PROJECTED SUBGRADIENT METHOD GIVES PROVABLE ROBUST RECOVERY OF SUBSPACES OF UNKNOWN CODIMENSION

Paris Giampouras, Benjamin D. Haeffele and René Vidal

Mathematical Institute for Data Science

Johns Hopkins University

Baltimore, MD, USA

{parisg, bhaeffele, rvidal}@jhu.edu

ABSTRACT

Robust subspace recovery (RSR) is the problem of learning a subspace from sample data points corrupted by outliers. Dual Principal Component Pursuit (DPCP) is a robust subspace recovery method that aims to find a basis for the orthogonal complement of the subspace by minimizing the sum of the distances of the points to the subspaces subject to orthogonality constraints on the basis. Prior work has shown that DPCP can provably recover the correct subspace in the presence of outliers as long as the true dimension of the subspace is known. In this paper, we show that if the orthogonality constraints –adopted in previous DPCP formulations– are relaxed and random initialization is used instead of spectral one, DPCP can provably recover a subspace of *unknown dimension*. Specifically, we propose a very simple algorithm based on running multiple instances of a projected sub-gradient descent method (PSGM), with each problem instance seeking to find one vector in the null space of the subspace. We theoretically prove that under mild conditions this approach succeeds with high probability. In particular, we show that 1) all of the problem instances will converge to a vector in the nullspace of the subspace and 2) the ensemble of problem instance solutions will be sufficiently diverse to fully span the nullspace of the subspace thus also revealing its true unknown codimension. We provide empirical results that corroborate our theoretical results and showcase the remarkable implicit rank regularization behavior of the PSGM algorithm that allows us to perform RSR without knowing the subspace dimension.

1 INTRODUCTION

Robust subspace recovery (RSR) refers to the problem of identifying an underlying linear subspace (with dimension less than the ambient data dimension) from sample data points that are potentially corrupted with outliers (i.e., points that do not lie in the linear subspace). Many methods for RSR have been proposed in the literature over the past several years (Xu et al., 2012; You et al., 2017a; Lerman & Maunu, 2018). Formulations based on convex relaxations and decompositions of the data matrix into a low-rank matrix plus a matrix of sparse corruptions – either entrywise-sparse corruptions as in Candès et al. (2011) or columnwise-sparse corruptions as in Xu et al. (2012); McCoy & Tropp (2011) – can, in certain situations, be shown to provably recover the true subspace when the dimension is unknown. However, these theoretical guarantees often require the dimension of the subspace, d , to be significantly less than the ambient dimension of the data, D , and these methods are not suitable for the more challenging regime of subspaces of high relative dimension (i.e., when $\frac{d}{D} \approx 1$).

Dual Principal Component Pursuit. Recently, progress has been made towards solving the RSR problem in the high relative dimension regime by a formulation termed Dual Principal Component Pursuit (DPCP). As implied by its name, DPCP follows a *dual* perspective of RSR by aiming to recover a basis for the orthogonal complement of the inliers’ subspace. As shown in (Tsakiris & Vidal, 2018), DPCP is provably robust in recovering subspaces of high relative dimension. However, a key limitation of DPCP is that it requires *a priori* knowledge of the true subspace dimension.

DPCP for $c = 1$. Let $\tilde{\mathbf{X}} \in \mathbb{R}^{D \times (N+M)}$ denote the data matrix defined as $\tilde{\mathbf{X}} = [\mathbf{X} \ \mathbf{O}]\mathbf{\Gamma}$, where $\mathbf{X} \in \mathbb{R}^{D \times N}$ is a matrix containing N inliers as its columns, $\mathbf{O} \in \mathbb{R}^{D \times M}$ is a matrix containing M outliers, and $\mathbf{\Gamma}$ is an unknown permutation matrix. DPCP was first formulated by Tsakiris & Vidal (2018) for handling subspaces of codimension $c = D - d$ equal to 1 (i.e., the subspace is a hyperplane with dimension $d = D - 1$). In this case, DPCP is formulated as the optimization problem

$$\min_{\mathbf{b} \in \mathbb{R}^D} \|\tilde{\mathbf{X}} \succ \mathbf{b}\|_1 \quad \text{s.t.} \quad \|\mathbf{b}\|_2 = 1, \quad (1)$$

which is nonconvex due to the spherical constraint imposed on the normal vector $\mathbf{b} \in S^{D-1}$ of the $D - 1$ dimensional hyperplane. Tsakiris & Vidal (2018) showed that the global minimizer of (1) is a normal vector of the underlying true hyperplane when both inliers and outliers are well-distributed or the ratio between the number of inliers and number of outliers is sufficiently small. Following a probabilistic point of view, (Zhu et al., 2018) presented an improved theoretical analysis of DPCP giving further insights on the remarkable robustness of DPCP in recovering the true underlying subspaces even in datasets heavily corrupted by outliers. Moreover, the authors introduced a projected subgradient method which converges to a normal vector of the true subspace at a linear rate.

Recursive DPCP for known $c > 1$. Zhu et al. (2018) also proposed an extension to DPCP to subspaces with codimension $c > 1$ via a projected subgradient algorithm that tries to learn c normal vectors to the subspace in a recursive manner. Specifically, after convergence to a normal vector, the projected subgradient algorithm is initialized with a vector *orthogonal* to the previously estimated normal vector. However, for that approach to be successful, knowledge of the true subspace codimension c becomes critical. Specifically, if an underestimate of the true codimension c is assumed the recovered basis for the null space, $\hat{\mathbf{B}}$, will fail to span the whole null space, \mathcal{S}_γ . On the other hand, an overestimate of c will lead to columns of $\hat{\mathbf{B}}$ corresponding to vectors that lie in \mathcal{S} .

Orthogonal DPCP for known c . Zhu et al. (2019) proposed an alternative to (1) which attempts to solve for c normal vectors to the subspace at once by minimizing the sum of the distances from the points to the subspace:

$$\min_{\mathbf{B} \in \mathbb{R}^{D \times c}} \|\tilde{\mathbf{X}} \succ \mathbf{B}\|_{1,2} \quad \text{s.t.} \quad \mathbf{B}^\top \mathbf{B} = \mathbf{I}. \quad (2)$$

The authors also proposed an optimization algorithm based on the projected Riemannian subgradient method (RSGM), which builds on similar ideas as the projected subgradient method of Zhu et al. (2018) and enjoys a linear converge rate when the step size is selected based on a geometrically diminishing rule. Ding et al. (2021) provided a geometric analysis of (2) which shows the merits of DPCP in handling a) datasets highly contaminated by outliers (in the order of $M = \mathcal{O}(N^2)$) and b) subspaces of high relative dimension. However, a key shortcoming of this approach is that, because all minimizers of (2) are orthogonal matrices, a prerequisite for recovering the correct orthogonal complement of the inliers subspace is the *a priori* knowledge of the true codimension c (see Fig. 1).

Contributions. In this work, we address this key limitation by proposing a framework that allows us to perform robust subspace recovery in the high relative subspace dimension regime **without** requiring *a priori* knowledge of the true subspace dimension. In particular, our proposed approach is based on the simple idea of solving multiple, parallel instances of the DPCP formulation for subspaces of codimension one,

$$\min_{\mathbf{B} \in \mathbb{R}^{D \times c^\ell}} \sum_{i=1}^{c^\ell} \|\tilde{\mathbf{X}} \succ \mathbf{b}_i\|_1 \quad \text{s.t.} \quad \|\mathbf{b}_i\|_2 = 1, \quad i = 1, 2, \dots, c^\ell \quad (3)$$

where c^ℓ is assumed to be an upper bound of c , i.e., $c^\ell \geq c$. Contrary to (2), the objective function in (3) decouples the columns \mathbf{b}_i of matrix $\mathbf{B} = [\mathbf{b}_1 \ \mathbf{b}_2 \ \dots \ \mathbf{b}_{c^\ell}]$ and thus can be solved in a parallel manner by independently applying a projected subgradient algorithm (referred to as PSGM) from c^ℓ different *random initializations*. Moreover, we observe that with random initialization we can get vectors sufficiently spread on the sphere that lead PSGM (initialized with those vectors) to return normal vectors of \mathcal{S} . These are all *linearly independent* when $c^\ell \leq c$ and thus can span \mathcal{S}_γ when $c^\ell = c$. If $c^\ell > c$ then PSGM will return $c^\ell - c$ redundant vectors that will still lie in \mathcal{S}_γ yet they will be linearly dependent (see Figure 1). That being said, we show that this simple strategy permits us to robustly recover the true subspace even without knowledge of the true codimension c .

As is detailed in Sections 3 and 4, this remarkable behavior of PSGM originates from the *implicit bias* that is induced in the optimization problem due to a) the relaxation of *orthogonality constraints* in (3) and b) the *random initialization* scheme that is adopted. Our specific contributions are as follows:

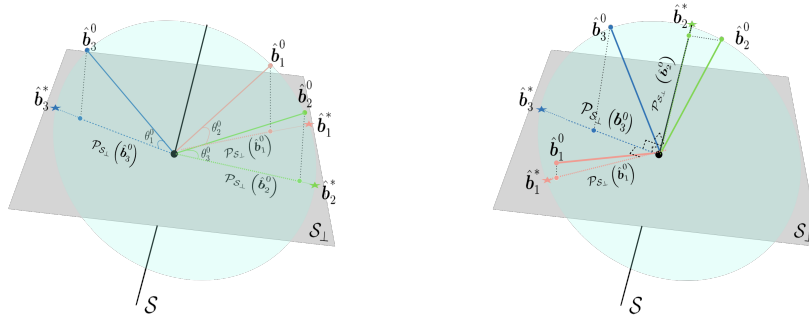


Figure 1: Graphical illustration of the recovered normal vectors of S by (left) the proposed DPCP-PSGM approach and (right) methods that use spectral initialization and impose orthogonality constraints. Initial vectors $\mathbf{b}_1^0, \mathbf{b}_2^0, \mathbf{b}_3^0$ are randomly initialized and are non-orthogonal in (left) and spectrally initialized (hence orthogonal) in (right). Note that in (left) $\text{rank}(\hat{\mathbf{B}})$ (where $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1, \hat{\mathbf{b}}_2, \hat{\mathbf{b}}_3]$) equals to the true codimension $c = 2$ of S and $\text{span}(\hat{\mathbf{B}}) \equiv S_\perp$ while in (right) $\hat{\mathbf{B}}$ is orthogonal hence $\text{rank}(\hat{\mathbf{B}}) = 3$ with $\hat{\mathbf{b}}_2 \in S$.

1. First, we study a continuous version of (3) where the inliers and outliers are drawn from continuous measures. We show that this induces a benign landscape on the DPCP objective, which can be analyzed more easily. Specifically, we prove that the DPCP problem in (3) can be solved via a projected subgradient algorithm that implicitly biases solutions towards low-rank matrices $\hat{\mathbf{B}} \in \mathbb{R}^{D \times c^\ell}$ whose columns are the projections of the randomly initialized columns of \mathbf{B}^0 onto S_\perp . As a result, $\hat{\mathbf{B}}$ almost surely spans S_\perp as long as it is randomly initialized with $c^\ell \geq c$.
2. Second, we analyze the discrete version which is more challenging, yet of more practical interest, showing that iterates of DPCP-PSGM converge to a *scaled* and *perturbed* version of the initial matrix \mathbf{B}^0 . This compelling feature of DPCP-PSGM allows to derive a sufficient condition and a probabilistic bound guaranteeing when the matrix $\hat{\mathbf{B}} \in \mathbb{R}^{D \times c^\ell}$ spans S_\perp .
3. We provide empirical results both on simulated and a real datasets, corroborating our theory and showing the robustness of our approach even without knowledge of the true subspace codimension.

2 RELATED WORK

Subspace Recovery. Learning underlying low-dimensional subspace representations of data has been a central topic of interest in machine learning research. Principal Component Analysis (PCA) has been the most celebrated method of this kind and is based on the minimization of the perpendicular distances of the data points from the estimates linear subspace, Jolliffe & Cadima (2016). Albeit, it is originally formulated as nonconvex optimization problem, PCA can be easily solved in closed form using a singular value decomposition (SVD) operation (see e.g. Vidal et al. (2016)). Despite its great success, PCA is prone to failure when handling datasets that contain *outliers* i.e., data points whose deviation from the inliers’ subspace is “large” in the ℓ_2 norm sense.

Robust Subspace Recovery (RSR). To remedy this weakness of PCA *robust subspace recovery (RSR)* methods attempt to identify the outliers in the dataset and recover the true underlying low-dimensional subspace of the inliers Lerman & Maunu (2018); Maunu et al. (2019). A classical approach to this problem is RANSAC Fischler & Bolles (1981), which is given a time budget and randomly chooses per iteration d points and then fits a d -dimensional subspace to those points and checks how the proposed subspace fits the remaining data points. RANSAC then outputs the subspace that agrees with the largest number of points. However, RANSAC’s reliance on randomly sampling points to propose subspaces can be highly inefficient when the number of outliers is high (as well as the fact that RANSAC also needs knowledge of the true subspace dimension, d). The need to tackle inherent shortcomings of RANSAC pertaining to computational complexity issues inspired alternative convex formulations of RSR, Xu et al. (2012); You et al. (2017b); Rahmani & Atia (2017); Zhang & Lerman (2014). In Xu et al. (2012) the authors decompose the data matrix as a sum of a low-rank and a column-sparse component. However, theoretical guarantees obtained for convex

formulations only hold for subspaces of relatively low-dimensional subspaces i.e., for $d \ll D$ where d and D denote the subspace and the ambient dimension, respectively. To the best of our knowledge, existing RSR algorithms rely heavily on one of two key assumptions. 1) The subspace is very low-dimensional relative to the ambient dimension ($d \ll D$) or 2) *The subspace dimension is a priori known*. Undoubtedly, the second hypothesis is rather strong in real-world applications, and many applications also do not satisfy the first assumption. Moreover, heuristic strategies for selecting the dimension of the subspace are hard to be applied in the RSR setting since they incur computationally prohibitive procedures, Lerman & Maunu (2018).

Relation to Orthogonal Dictionary Learning (ODL). Note that objective functions in the form of (3) show up beyond RSR problems i.e., in orthogonal dictionary learning (ODL), sparse blind deconvolution, etc., Qu et al. (2020). Specifically, based on a similar formulation the authors in Bai et al. (2019) proved that $c^d = \mathcal{O}(c \log c)$ independent random initial vectors suffice in order to recover with high probability a dictionary of size $D \times c$ with high accuracy. In this paper we aim to recover a basis of the orthogonal complement of a subspace of unknown dimension instead of accurately estimating a dictionary hence our goal differs from that in Bai et al. (2019).

Implicit bias in Robust Recovery Problems. The notions of implicit bias and implicit regularization have been used interchangeably in the nonconvex optimization literature for describing the tendency of optimization algorithms to converge to global minima of minimal complexity with favorable generalization properties in overparameterized models, Gunasekar et al. (2018). In the context of robust recovery, the authors in You et al. (2020) showed that Robust PCA can be suitably reparametrized in such a way to favor low-rank and sparse solutions without using any explicit regularization. In this work, we use the term implicit bias for describing the convergence of DPCP-PSGM to low-rank solutions, which are not necessarily global minimizers, that span the orthogonal complement of the subspace when a) orthogonality constraints in DPCP formulation are relaxed b) DPCP is overparameterized i.e., $c^d \geq c$ and c) PSGM randomly initialized.

3 DUAL PRINCIPAL COMPONENT PURSUIT AND THE PROJECTED SUBGRADIENT METHOD

We re-write the DPCP formulation given in (1) as

$$\min_{\mathbf{B} \in \mathbb{R}^{D \times d}} \min_{c^d} \|\tilde{\mathbf{X}} \succ \mathbf{b}\|_1 = \|\mathbf{X} \succ \mathbf{b}\|_1 + \|\mathbf{O} \succ \mathbf{b}\|_1 \quad \text{s.t.} \quad \|\mathbf{b}\|_2 = 1 \quad (4)$$

In Zhu et al. (2018), the authors proposed a projected subgradient descent algorithm for addressing (4) that consists of a subgradient step followed by a projection onto the sphere i.e.,

$$\mathbf{b}^{k+1} = \hat{\mathbf{b}}^k - \mu^k (\mathbf{X} \text{Sgn}(\mathbf{X} \succ \mathbf{b}^k) + \mathbf{O} \text{Sgn}(\mathbf{O} \succ \mathbf{b}^k)) \quad \text{and} \quad \hat{\mathbf{b}}^{k+1} = \mathcal{P}_{\mathbb{S}^{D-1}}(\mathbf{b}^{k+1}), \quad (5)$$

where μ^k is the -adaptively updated per iteration- step size and $\hat{\mathbf{b}}^k$ is the unit ℓ_2 norm vector corresponding to the k th iteration.

The convergence properties of the projected subgradient algorithm described above depend on specific quantities denoted as $c_{\mathbf{X}, \min}$ and $c_{\mathbf{X}, \max}$ that reflect the geometry of the problem and are defined as $c_{\mathbf{X}, \min} = \frac{1}{N} \min_{\mathbf{b} \in \mathcal{S}^{d-1} \setminus \mathcal{S}} \|\mathbf{X} \succ \mathbf{b}\|_1$ and $c_{\mathbf{X}, \max} = \frac{1}{N} \max_{\mathbf{b} \in \mathcal{S}^{d-1} \setminus \mathcal{S}} \|\mathbf{X} \succ \mathbf{b}\|_1$. Note that the more well distributed the inliers are in the subspace \mathcal{S} the higher the value of the quantity $c_{\mathbf{X}, \min}$ (called as permeance statistic which first appeared in Lerman et al. (2015)) as it becomes harder to find a vector \mathbf{b} in the subspace \mathcal{S} that is orthogonal to the inliers. Moreover, $c_{\mathbf{X}, \min}$ and $c_{\mathbf{X}, \max}$ converge to the same value as $N \rightarrow \infty$ provided the inliers are uniformly distributed in the subspace i.e., $c_{\mathbf{X}, \min} \rightarrow c_d$, $c_{\mathbf{X}, \max} \rightarrow c_d$, where c_d is given as the average height of the unit hemisphere on \mathbb{R}^d ,

$$c_d := \frac{(d-2)!!}{(d-1)!!} \begin{cases} \frac{2}{\pi}, & \text{if } d \text{ is even,} \\ 1, & \text{if } d \text{ is odd} \end{cases} \quad \text{where } k!! = \begin{cases} k(k-2)(k-4) \cdots 4 \cdot 2, & k \text{ is even,} \\ k(k-2)(k-4) \cdots 3 \cdot 1, & k \text{ is odd} \end{cases} \quad (6)$$

Similarly to $c_{\mathbf{X}, \min}$, $c_{\mathbf{X}, \max}$, we will also be interested in quantities $c_{\mathbf{O}, \min}$, $c_{\mathbf{O}, \max}$ which indicate how well-distributed the outliers are in the ambient space. These quantities are defined as $c_{\mathbf{O}, \min} = \min_{\mathbf{b} \in \mathcal{S}^{D-1} \setminus \mathcal{S}} \frac{1}{M} \|\mathbf{O} \succ \mathbf{b}\|_1$ and $c_{\mathbf{O}, \max} = \max_{\mathbf{b} \in \mathcal{S}^{D-1} \setminus \mathcal{S}} \frac{1}{M} \|\mathbf{O} \succ \mathbf{b}\|_1$. $c_{\mathbf{O}, \max}$ can be viewed as the *dual* permeance statistic and is bounded away from small values while its difference from $c_{\mathbf{O}, \min}$ tends to

zero as $M \rightarrow \infty$. Further, if the outliers are uniformly distributed on the sphere, then $c_{\mathcal{O},\max} \rightarrow c_D$ and $c_{\mathcal{O},\min} \rightarrow c_D$ where c_D is defined as in (6), (Zhu et al., 2018).

Finally, we also define the quantities $\eta_{\mathcal{O}} = \frac{1}{M} \max_{\mathbf{b} \in \mathcal{S}^{D-1}} \|(\mathbf{I} - \mathbf{b}\mathbf{b}^\top)\mathcal{O}\text{Sgn}(\mathcal{O}^\top \mathbf{b})\|_2$ and $\eta_{\mathcal{X}} = \frac{1}{M} \max_{\mathbf{b} \in \mathcal{S}^{D-1}} \|(\mathcal{P}_{\mathcal{S}} - \mathbf{b}\mathbf{b}^\top)\mathbf{X}\text{Sgn}(\mathbf{X}^\top \mathbf{b})\|_2$. As $M \rightarrow \infty$ and assuming outliers in \mathcal{O} are well-distributed we get $\mathcal{O}\text{Sgn}(\mathcal{O}^\top \mathbf{b}) \rightarrow c_D \mathbf{b}$ thus $\eta_{\mathcal{O}} \rightarrow 0$ (Tsakiris & Vidal, 2018). Likewise, $\eta_{\mathcal{X}} \rightarrow 0$ as $N \rightarrow \infty$ provided that inliers are uniformly distributed in the d -dimensional subspace. The following theorem (see full version in Appendix) provides convergence guarantees of the projected subgradient method that was proposed in Zhu et al. (2018) for addressing problem (1).

Theorem 1 (Informal Theorem 3 of Zhu et al. (2018)) Let $\{\hat{\mathbf{b}}_k\}$ the sequence generated by the projected subgradient algorithm in Zhu et al. (2018), with initialization $\hat{\mathbf{b}}_0$ such that

$$\theta_0 < \arctan\left(\frac{Nc_{\mathcal{X},\min}}{N\eta_{\mathcal{X}} + M\eta_{\mathcal{O}}}\right) \text{ and } Nc_{\mathcal{X},\min} \geq N\eta_{\mathcal{X}} + M\eta_{\mathcal{O}} \quad (7)$$

where θ_0 denotes the principal angle of $\hat{\mathbf{b}}^0$ from \mathcal{S}^\perp . If the step size μ^k is updated according to a piecewise geometrically diminishing rule given as

$$\mu^k = \begin{cases} \mu^0, & k < K_0 \\ \mu^0 \beta^{\lfloor (k - K_0)/K \rfloor + 1}, & k \geq K_0 \end{cases} \quad (8)$$

where $\beta < 1$, $\lfloor \cdot \rfloor$ is the floor function, then the iterates $\hat{\mathbf{b}}^k$ converge to a normal vector of \mathcal{S} .

4 DUAL PRINCIPAL COMPONENT PURSUIT IN SUBSPACES OF UNKNOWN CODIMENSION

Current theoretical results provide guarantees for recovering the true inlier subspace, when the proposed algorithms know *a priori* of the subspace codimension c , which is a rather strong requirement and is far from being true in real word applications. Here we describe our proposed approach, which consists of removing the orthogonality constraint on \mathbf{B} , along with a theoretical analysis that gives guarantees of recovering the true underlying subspace even when the true codimension c is unknown. First we analyze a continuous version of DPCP, which arises when the number of inliers and outliers are distributed according to continuous measures and their number tends to ∞ . The continuous DPCP incurs an optimization problem with a benign landscape that allows us to better illustrate the favorable properties of DPCP-PSGM when it comes to the convergence of its iterates. Then we extend the results to the discrete case that deals with a finite number of inliers and outliers yielding a more challenging optimization landscape.

4.1 PSGM’S ITERATES CONVERGENCE IN THE CONTINUOUS VERSION OF DPCP

The following lemma provides the continuous version of the discrete objective function given in (3).

Lemma 2 In the continuous case, the discrete DPCP problem given in (3) is reformulated as,

$$\begin{aligned} \min_{\mathbf{B} \in \mathbb{R}^{D \times c^0}} \sum_{i=1}^{c^0} (p \mathbb{E}_{\mu_{\mathcal{S}^D-1}}[f_{\mathbf{b}_i}] + (1-p) \mathbb{E}_{\mu_{\mathcal{S}^D-1} \setminus \mathcal{S}}[f_{\mathbf{b}_i}]) &= \sum_{i=1}^{c^0} \|\mathbf{b}_i\|_2 (pc_D + (1-p)c_d \cos(\phi_i)) \\ \text{s.t. } \|\mathbf{b}_i\|_2 &= 1, \quad i = 1, 2, \dots, c^0 \end{aligned} \quad (9)$$

where $f_{\mathbf{b}} : \mathcal{S}^{D-1} \rightarrow \mathbb{R}$, $f_{\mathbf{b}}(\mathbf{z}) = |\mathbf{z}^\top \mathbf{b}|$, ϕ_i is the principal angle of \mathbf{b}_i from the inliers subspace \mathcal{S} and p is the probability of occurrence of an outlier.

Note that $\mu_{\mathcal{S}^D-1}$, $\mu_{\mathcal{S}^D-1 \setminus \mathcal{S}}$ are the continuous measures associated with the outliers and inliers, respectively. Evidently, (9) attains its global minimum for vectors \mathbf{b}_i s that are orthogonal to the inliers’ subspace. Based on (3) and due to Lemma 2, we can now minimize the objective function of the “continuous version” of DPCP by employing a projected subgradient methods (PSGM) that

performs the following steps per iteration ¹

$$\mathbf{b}_i^{k+1} = \hat{\mathbf{b}}_i^k - \mu_i^k (pc_D \hat{\mathbf{b}}_i^k + (1-p)c_d \hat{\mathbf{s}}_i^k) \text{ and } \hat{\mathbf{b}}_i^{k+1} = \mathcal{P}_{S_\gamma}(\mathbf{b}_i^{k+1}), i = 1, 2, \dots, c^\ell \quad (10)$$

Lemma 3 *A projected subgradient algorithm consisting of the steps described in (10) using a piecewise geometrically diminishing step size rule (see (8) in Theorem 1) will almost surely asymptotically converge to a matrix $\hat{\mathbf{B}} \in \mathbb{R}^{D \times c^\ell}$ whose columns $\hat{\mathbf{b}}_i$, $i = 1, 2, \dots, c^\ell$ will be normal vectors of the inliers' subspace when randomly initialized with vectors $\hat{\mathbf{b}}_i^0 \in S^{D-1}$, $i = 1, 2, \dots, c^\ell$ uniformly distributed over the sphere S^{D-1} .*

Lemma 3 allows us to claim that we can always recover $c^\ell \geq c$ normal vectors to the inliers' subspace using a PSGM algorithm consisting of steps given in (10). However, this does not tell the whole story yet, since our ultimate objective is to recover a matrix $\hat{\mathbf{B}}$ that spans S_γ . Thus, it remains to show that the rank of $\hat{\mathbf{B}}$ is equal to the true and *unknown* codimension of the inliers' subspace c . Next we prove that by initializing with a $\hat{\mathbf{B}}_0$ such that $\text{rank}(\hat{\mathbf{B}}_0) = c^\ell$ (i.e., $\hat{\mathbf{B}}_0$ is initialized to be full-rank), we can guarantee that we can solve the continuous version of DPCP using PSGM and converge to a $\hat{\mathbf{B}}$ such that $\text{rank}(\hat{\mathbf{B}}) = c$ thus getting $\text{span}(\hat{\mathbf{B}}) \equiv S_\gamma$ (along with recovering the true subspace dimension). By projecting the PSGM iterates given in (10) onto S_γ we have,

$$\mathcal{P}_{S_\gamma}(\mathbf{b}_i^{k+1}) = (1 - \mu_i^k pc_D) \mathcal{P}_{S_\gamma}(\hat{\mathbf{b}}_i^k) \text{ and } \mathcal{P}_{S_\gamma}(\hat{\mathbf{b}}_i^{k+1}) = \mathcal{P}_{S_\gamma}(\mathcal{P}_{S^{D-1}}(\mathbf{b}_i^{k+1})) \quad (11)$$

We hence observe that the projections of successive iterates of PSGM are scaled versions of the corresponding projections of the previous iterates. We can now state Lemma 4.

Lemma 4 *The PSGM iterates $\hat{\mathbf{b}}_i^k$, $i = 1, 2, \dots, c^\ell$, $k = 1, 2, \dots$ given in (10), when randomly initialized with $\hat{\mathbf{b}}_i^0$ s, $i = 1, 2, \dots, c^\ell$ that are independently drawn from a spherical distribution with unit ℓ_2 norm converge almost surely to c^ℓ normal vectors of the inliers subspace S denoted as $\hat{\mathbf{b}}_i$, $i = 1, 2, \dots, c^\ell$ that are given by $\hat{\mathbf{b}}_i = \frac{\mathcal{P}_{S_\gamma}(\hat{\mathbf{b}}_i^0)}{k \mathcal{P}_{S_\gamma}(\hat{\mathbf{b}}_i^0)_{k_2}}$, $i = 1, 2, \dots, c^\ell$.*

Lemma 4 shows that the initialization of PSGM plays a pivotal role since it determines the direction of the recovered normal vectors $\{\hat{\mathbf{b}}_i\}_{i=1}^{c^\ell}$. Lemmas 3 and 4 pave the way for Theorem 5.

Theorem 5 *Let $\hat{\mathbf{B}}^0 \in \mathbb{R}^{D \times c^\ell}$ where $c^\ell \geq c$ with c denoting the true codimension of the inliers subspace S , consisting of unit ℓ_2 norm column vectors $\hat{\mathbf{b}}_i^0 \in S^{D-1}$, $i = 1, 2, \dots, c^\ell$ that are independently drawn from uniform distribution over the sphere S^{D-1} . A PSGM algorithm initialized with $\hat{\mathbf{B}}^0$ will almost surely converge to a matrix $\hat{\mathbf{B}}$ such that $\text{span}(\hat{\mathbf{B}}) \equiv S_\gamma$.*

From Theorem 5 we observe that in the benign scenario where inliers and outliers are distributed under continuous measures, we can recover the correct orthogonal complement of the inlier's subspace even when we are oblivious to its true codimension. Remarkably, this is achieved by exploiting the *implicit bias* induced by multiple random initializations of the PSGM algorithm for solving the DPCP formulation given in (3), which is free of orthogonality constraints.

4.2 PSGM'S ITERATES CONVERGENCE IN THE DISCRETE VERSION OF DPCP

From this analysis of the continuous version of DPCP we now extend to the discrete version, which is of more practical relevance for finite data, yet also presents more challenges. To begin, we reformulate the DPCP objective as follows

$$\sum_{i=1}^{c^\ell} \|\tilde{\mathbf{X}}^> \mathbf{b}_i\|_1 = \sum_{i=1}^{c^\ell} \|\mathbf{X}^> \mathbf{b}_i\|_1 + \|\mathbf{O}^> \mathbf{b}_i\|_1 = M \sum_{i=1}^{c^\ell} \mathbf{b}_i^> \mathbf{o}_{\mathbf{b}_i} + N \sum_{i=1}^{c^\ell} \mathbf{b}_i^> \mathbf{x}_{\mathbf{b}_i} \quad (12)$$

where $\mathbf{x}_{\mathbf{b}_i}$ and $\mathbf{o}_{\mathbf{b}_i}$ are called as *average inliers* and *average outliers* terms, defined as $\mathbf{x}_{\mathbf{b}_i} = \frac{1}{N} \sum_{j=1}^N \text{Sgn}(\mathbf{b}_i^> \mathbf{x}_j) \mathbf{x}_j$ and $\mathbf{o}_{\mathbf{b}_i} = \frac{1}{M} \sum_{j=1}^M \text{Sgn}(\mathbf{b}_i^> \mathbf{o}_j) \mathbf{o}_j$.

In Algorithm 1, we give the projected subgradient method (DPCP-PSGM) applied on the DPCP problem given in (3).

¹Note that $\partial \|\mathbf{b}\|_2 = \frac{\mathbf{b}}{\|\mathbf{b}\|_2}$ for $\mathbf{b} \neq \mathbf{0}$ and $\|\mathbf{b}_i\|_2 \cos(\phi_i) = \mathbf{b}_i^\top \hat{\mathbf{s}}_i$ where $\hat{\mathbf{s}}_i = \frac{\mathcal{P}_S(\mathbf{b}_i)}{\|\mathcal{P}_S(\mathbf{b}_i)\|_2}$.

Algorithm 1: DPCP-PSGM algorithm for solving (3)

Result: $\hat{\mathbf{B}} = [\hat{\mathbf{b}}_1^k, \hat{\mathbf{b}}_2^k, \dots, \hat{\mathbf{b}}_{c^\theta}^k]$
Initialize: Randomly sample $\hat{\mathbf{b}}_0^1, \hat{\mathbf{b}}_0^2, \dots, \hat{\mathbf{b}}_0^{c^\theta}$ from a uniform distribution on S^{D-1} ;
for $k = 1, 2, \dots$ **do**
 for $i = 1, 2, \dots, c^\theta$ **do**
 Update the step-size according to a specific rule;
 $\mathbf{b}_i^{k+1} = \hat{\mathbf{b}}_i^k - \mu_i^k (M \mathbf{o}_{\hat{\mathbf{b}}_i^k}^k + N \mathbf{x}_{\hat{\mathbf{b}}_i^k}^k)$;
 $\hat{\mathbf{b}}_i^{k+1} = \mathcal{P}_{S^{D-1}}(\mathbf{b}_i^{k+1})$;
 end
end

Note that the average outliers and inliers terms are discrete versions of the corresponding continuous average terms $c_D \mathbf{b}_i$ and $c_d \hat{\mathbf{s}}_i$ where $\hat{\mathbf{s}}_i = \mathcal{P}_S(\mathbf{b}_i)$, respectively, Tsakiris & Vidal (2018). We now express the sub-gradient step of Algorithm 1 as

$$\mathbf{b}_i^{k+1} = \hat{\mathbf{b}}_i^k - \mu_i^k \left(M(c_D \hat{\mathbf{b}}_i^k + \mathbf{e}_{\mathbf{O}}^{i,k}) + N(c_D \hat{\mathbf{s}}_i^k + \mathbf{e}_{\mathbf{X}}^{i,k}) \right) \quad (13)$$

where the quantities $\mathbf{e}_{\mathbf{O}}^{i,k} = \mathbf{o}_{\hat{\mathbf{b}}_i^k} - c_D \hat{\mathbf{b}}_i^k$ and $\mathbf{e}_{\mathbf{X}}^{i,k} = \mathbf{x}_{\hat{\mathbf{b}}_i^k} - c_D \hat{\mathbf{s}}_i^k$ account for the error between the continuous and discrete versions of the average outliers and the average inliers terms, respectively. Following a similar path as in the continuous case we next project the iterates of (13) onto \mathcal{S}_γ ,

$$\begin{aligned} \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^{k+1}) &= \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^k) - \mu_i^k \left(M \mathcal{P}_{\mathcal{S}_\gamma}(c_D \hat{\mathbf{b}}_i^k + \mathbf{e}_{\mathbf{O}}^{i,k}) + N \mathcal{P}_{\mathcal{S}_\gamma}(c_D \hat{\mathbf{s}}_i^k + \mathbf{e}_{\mathbf{X}}^{i,k}) \right) \\ &= (1 - \mu_i^k M c_D) \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^k) - \mu_i^k M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,k}) \end{aligned} \quad (14)$$

Remark. Eq. (14) reveals that DPCP-PSGM, applied on the discrete problem, gives rise to updates whose projections to \mathcal{S}_γ are **scaled** and **perturbed** versions of the previous estimates. The magnitude of perturbation depends on the discrepancy between the continuous and the discrete problem.

Recall that $\mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^k) = \frac{\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^k)}{k \mathbf{b}_i^k k_2}$, and we can rewrite the update of the 2nd iteration of DPCP-PSGM,

$$\begin{aligned} \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^2) &= \frac{(1 - \mu_i^1 M c_D)}{\|\mathbf{b}_i^1\|_2} \left((1 - \mu_i^0 M c_D) \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0) - \mu_i^0 M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,0}) \right) - \mu_i^1 M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,1}) \\ &= \frac{(1 - \mu_i^1 M c_D)(1 - \mu_i^0 M c_D)}{\|\mathbf{b}_i^1\|_2 \|\mathbf{b}_i^0\|_2} \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0) - \frac{(1 - \mu_i^1 M c_D)}{\|\mathbf{b}_i^1\|_2} \mu_i^0 M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,0}) - \mu_i^1 M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,1}) \end{aligned} \quad (15)$$

where we have assumed that $\|\mathbf{b}_i^0\|_2 = 1$. By repeatedly applying the same steps, we can reach to the following recursive expression for $\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^K)$,

$$\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^K) = \left(\prod_{k=0}^{K-1} \frac{(1 - \mu_i^k M c_D)}{\|\mathbf{b}_i^k\|_2} \right) \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^0) - \sum_{k=0}^{K-1} \left(\prod_{j=k+1}^{K-1} \frac{(1 - \mu_i^j M c_D)}{\|\mathbf{b}_i^j\|_2} \right) \mu_i^k M \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{e}_{\mathbf{O}}^{i,k}) \quad (16)$$

where for $j > K-1$ we set $\prod_{j=k+1}^{K-1} \frac{(1 - \mu_i^j M c_D)}{k \mathbf{b}_i^j k_2} = 1$.

By dividing (16) with $\prod_{k=0}^{K-1} \frac{(1 - \mu_i^k M c_D)}{k \mathbf{b}_i^k k_2}$ and by projecting onto the sphere S^{D-1} we get

$$\mathcal{P}_{S^{D-1}}(\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^K)) = \mathcal{P}_{S^{D-1}}(\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^0) - \mathcal{P}_{\mathcal{S}_\gamma}(\boldsymbol{\delta}_i^K)) \quad (17)$$

where $\boldsymbol{\delta}_i$ is defined as $\boldsymbol{\delta}_i^K = \sum_{k=0}^{K-1} \left(\prod_{j=0}^k \frac{k \mathbf{b}_i^j k_2}{(1 - \mu_i^j M c_D)} \right) \mu_i^k M \mathbf{e}_{\mathbf{O}}^{i,k}$.

Assumption 1. We assume that the principal angles θ_0^i for all \mathbf{b}_i^0 s satisfy the inequality $\theta_0^i < \arctan\left(\frac{N c_{\mathbf{X}, \min}}{N \eta_{\mathbf{X}} + M \eta_{\mathbf{O}}}\right) \forall i, i = 1, 2, \dots, c^\theta$.

Assumption 1 essentially assumes that the sufficient condition given in eq. (7) required by PSGM algorithm for converging to a normal vector is satisfied which is the same condition for success in Zhu et al. (2018). Under Assumption 1 we can invoke the convergence properties of PSGM given in Theorem 1 and get as $K \rightarrow \infty, \mathcal{P}_{S^D-1}(\mathcal{P}_{S^D}(\mathbf{b}^K)) \rightarrow \hat{\mathbf{b}} \in \mathcal{S}^D \cap S^{D-1}$. That being said, we denote $\hat{\mathbf{b}}_i = \mathcal{P}_{S^D-1}(\mathcal{P}_{S^D}(\mathbf{b}_i^0) - \mathcal{P}_{S^D}(\hat{\delta}_i))$, where $\hat{\delta}_i = \lim_{K \rightarrow \infty} \delta_i^{K^2}$. Following the same steps as in Section 4.1 and by defining matrices $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{c^0}]$, $\mathbf{B}^0 = [\mathbf{b}_1^0, \mathbf{b}_2^0, \dots, \mathbf{b}_{c^0}^0]$, $\hat{\Delta} = [\hat{\delta}_1, \hat{\delta}_2, \dots, \hat{\delta}_{c^0}]$ we can express the matrix \mathbf{B} as $\hat{\mathbf{B}} = \mathcal{P}_{S^D-1}(\mathcal{P}_{S^D}(\mathbf{B}^0 - \hat{\Delta}))$ where \mathbf{B} will now consist of normal vectors of the inliers' subspace. In order to guarantee that $\text{span}(\mathbf{B}) \equiv \mathcal{S}^D$ it thus suffices to ensure that $\text{rank}(\mathbf{B}) = c$. Here we show that a sufficient condition for this to hold is that the matrix $\mathbf{A} = \mathbf{B}^0 - \hat{\Delta}$ is full-rank.

Lemma 6 *If $\sigma_{c^0}(\mathbf{B}^0) > \|\hat{\Delta}\|_2$ then matrix $\mathbf{A} = \mathbf{B}^0 - \hat{\Delta}$ is full-rank.*

From Lemma 6 we can see that the success of DPCP-PSGM hinges on how well-conditioned the matrix \mathbf{B}^0 is. Specifically, it says that if a lower-bound on the smallest singular is satisfied then DPCP-PSGM is guaranteed to converge to the correct complement of the inlier without knowledge of the correct codimension c . From this, we can prove the following Theorem.

Theorem 7 *Let $\mathbf{B}^0 \in \mathbb{R}^{D \times c^0}$ with columns randomly sampled from a unit ℓ_2 norm spherical distribution where $c^0 \geq c$ with c denoting the true codimension of the inliers subspace \mathcal{S} that satisfies Assumption 1. If*

$$1 - C_1 \sqrt{\frac{c^0}{D}} - \frac{\epsilon}{\sqrt{D}} > \sqrt{c^0} \kappa (\eta_{\mathcal{O}} + c_{\mathcal{O}, \max} - c_d) \quad (18)$$

where $\kappa = \max_i \frac{M \mu_i^0}{\beta^{K_0/K} (1 - r_i)}$ and $r_i = \frac{(1 + \mu_i^0 (N(\eta_{\mathcal{X}} + c_{\mathcal{X}, \max}) + M(\eta_{\mathcal{O}} + c_{\mathcal{O}, \max})))}{1 - \mu_i^0 M c_D} \beta^{1/K}$ then with probability at least $1 - 2 \exp(-\epsilon^2 C_2)$ (where C_1, C_2 are absolute constants), Algorithm 1 with a piecewise geometrically diminishing step size rule will converge to a matrix $\hat{\mathbf{B}}$ such that $\text{span}(\hat{\mathbf{B}}) \equiv \mathcal{S}^D$.

Note that quantities β, K, K_0 are used in the step-size update rule that is used as defined in (8) (See also full version of Theorem 1 in Appendix). Theorem 7 shows that we can randomly initialize DPCP-PSGM, with a matrix $\hat{\mathbf{B}}^0$ whose number of columns c^0 is an overestimate of the true codimension c of the inliers' subspace and with columns sampled *independently* by a uniform distribution over the unit sphere and recover a matrix that will span the orthogonal complement of \mathcal{S} . The probability of success depends on the geometry of the problem since condition (18) is trivially satisfied (RHS of (18) tends to 0 since $\eta_{\mathcal{O}} \rightarrow 0$ and $c_{\mathcal{O}, \max} \rightarrow c_d$) in the continuous case which incurs a benign geometry. Moreover, the a less benign geometry would increase the value of $(\eta_{\mathcal{O}} + c_{\mathcal{O}, \max} - c_d)$ thus requiring a smaller initial codimension c^0 that would lead to larger values the LHS of (18).

5 NUMERICAL SIMULATIONS

In this section we demonstrate the effectiveness of the proposed DPCP formulation and the derived DPCP-PSGM algorithm in recovering orthogonal complements of subspace of unknown codimension. We compare the proposed algorithm with previously developed methods i.e., DPCP-IRLS Tsakiris & Vidal (2018) and the Riemannian Subgradient Method (RSGM) Zhu et al. (2019). Recall that both DPCP-IRLS and RSGM address DPCP problem by enforcing orthogonality constraints, and thus both algorithms are quite sensitive if the true codimension of \mathcal{S} is not known *a priori*. Further, they are both initialized using of spectral initialization i.e., $\hat{\mathbf{B}}^0 \in \mathbb{R}^{D \times c^0}$ which contains the first c^0 eigenvectors of matrix $\tilde{\mathcal{X}} \tilde{\mathcal{X}}^T$ as its columns, as proposed in Tsakiris & Vidal (2018).

Robustness to outliers in the unknown codimension regime. In this experiment we set the dimension of the ambient space to $D = 200$. We randomly generate N inliers uniformly distributed with unit ℓ_2 norm in a $d = 195$ dimensional subspace (hence for its codimension we have $c = D - d = 5$). Following a similar process we generate M outliers that live in the ambient space and are sampled

²For the sake of brevity we assume that the step size has been selected such that existence of the limit is guaranteed. We refer the reader to the proof of Lemma 8 for further details.

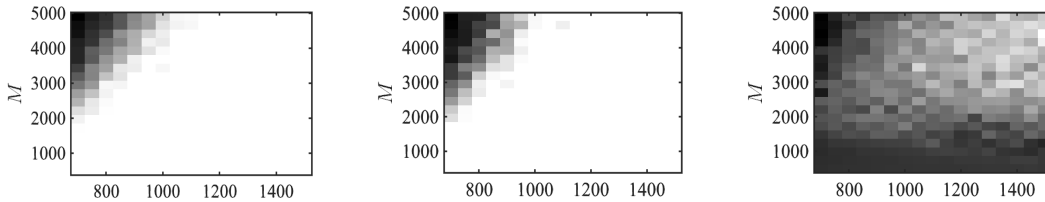


Figure 2: Distances of the recovered \hat{B} from the true orthogonal complements $\mathcal{S}_\mathcal{I}$ as recovered by the proposed DPCP-PSGM algorithm provided an overestimated of the true c i.e., $c^l = 10$ (left), RSGM provided the true c (middle) and RSGM provided $c^l = 10$ (right). Darker colors reflect higher values of distances while lighter colors indicate successful recoveries of B .

from a uniform distribution over the unit sphere. Fig. 2, illustrates the distances (see Appendix) of the recovered matrix \hat{B} as obtained by the proposed DPCP-PSGM algorithm initialized with an overestimate $c^l = 10$ of the true codimension c and two versions of RSGM i.e., RSGM when it is given as input true $c = 5$ and RSGM when being incognizant of c and hence it initialized with a $c^l = 10$ of c As is shown in Fig. 2 (right), RSGM fails to recover the correct orthogonal complement of \mathcal{S} when it is provided with an overestimate of the true c which is attributed to spectral initialization and the imposed orthogonality constraints. On the contrary, DPCP-PSGM displays a remarkably robust behavior (Fig. 2(middle)) even without knowing the true value of c , performing similarly to RSGM when the latter knows beforehand the correct codimension (Fig. 2(left)).

Recovery of the true codimension. Here we test DPCP-PSGM on the recovery of the true codimension c of the inliers’ subspace \mathcal{S} . Again, we set $D = 200$ and generate $N = 1500$ inliers as before. We vary the true codimension of \mathcal{S} from $c = 10$ to 20 and consider two different outliers’ ratios r , defined as $r = \frac{M}{M+N}$, namely $r = 0.6$ and $r = 0.7$. In both cases, DPCP-PSGM is initialized with the same overestimate of c i.e., $c^l = 30$. In Fig. 3 we report the estimated codimensions obtained by DPCP-PSGM for 10 independent trials of the experiments. It can be observed that DPCP-PSGM achieves 100% for all different codimensions for $r = 0.6$. Moreover, it shows a remarkable performance in estimating the correct c ’s even in the more challenging case corresponding to outliers’ ratios equal to 0.7. with the estimated codimensions being close to the true values even in the cases that it fails to exactly compute c . The results corroborate the theory showing that the DPCP-PSGM with random initialization biases the solutions of \hat{B} towards matrices with rank c .

6 CONCLUSIONS

We proposed a simple framework which allows us to perform robust subspace recovery without requiring a priori knowledge of the subspace codimension. This is based on Dual Principal Component Pursuit (DPCP) and thus is amenable to handling subspaces of high relative dimensions. We observed that a projected subgradient method (PSGM) induces implicit bias and converges to a matrix that spans a basis of the orthogonal complement of the inliers subspace even as long as a) we overestimate it codimension, b) lift orthogonality constraints enforced in previous DPCP formulations and c) use random initialization. Empirical results that corroborate the developed theory and showcase the merits of our approach.

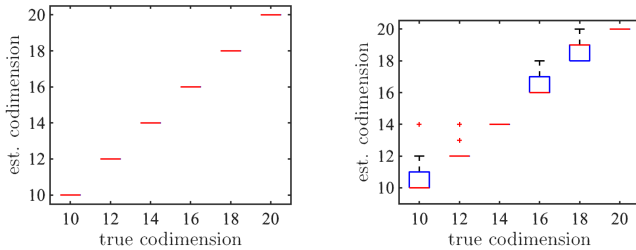


Figure 3: Estimated by DPCP-PSGM codimensions for two different outliers’ ratios $r = \frac{M}{M+N}$ (a) $r = 0.6$ (left) and (b) $r = 0.7$ (right)

Empirical results that corroborate the developed theory and showcase the merits of our approach.

Ethics Statement This work focuses on theoretical aspects of robust subspace recovery problem which is a well-established topic in machine learning research. The research conducted in the framework of this work raises no ethical issues or any violations vis-a-vis the ICLR Code of Ethics.

ACKNOWLEDGMENTS

We would like to thank Christian Kümmerle for helpful discussions on the probabilistic theorem that is used in Theorem 7. This work is partially supported by the European Union under the Horizon 2020 Marie-Skłodowska-Curie Global Fellowship program: HyPPOCRATES— H2020-MSCA-IF-2018, Grant Agreement Number: 844290, and the NSF Grants 1704458, 2031985 and 1934979.

REFERENCES

- Yu Bai, Qijia Jiang, and Ju Sun. Subgradient descent learns orthogonal dictionaries. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=Hk1Sf3CqKm>.
- Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM (JACM)*, 58(3):1–37, 2011.
- Tianyu Ding, Zhihui Zhu, Rene Vidal, and Daniel P Robinson. Dual principal component pursuit for robust subspace learning: Theory and algorithms for a holistic approach. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 2739–2748. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/ding21b.html>.
- Martin A Fischler and Robert C Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981.
- Paris V. Giampouras, Athanasios A. Rontogiannis, and Konstantinos D. Koutroumbas. Alternating iteratively reweighted least squares minimization for low-rank matrix factorization. *IEEE Transactions on Signal Processing*, 67(2):490–503, 2019. doi: 10.1109/TSP.2018.2883921.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 1832–1841. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/gunasekar18a.html>.
- Ian T Jolliffe and Jorge Cadima. Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 374(2065):20150202, 2016.
- Gilad Lerman and Tyler Maunu. An overview of robust subspace recovery. *Proceedings of the IEEE*, 106(8):1380–1410, 2018.
- Gilad Lerman, Michael B McCoy, Joel A Tropp, and Teng Zhang. Robust computation of linear models by convex relaxation. *Foundations of Computational Mathematics*, 15(2):363–410, 2015.
- Tyler Maunu, Teng Zhang, and Gilad Lerman. A well-tempered landscape for non-convex robust subspace recovery. *Journal of Machine Learning Research*, 20(37), 2019.
- Michael McCoy and Joel A Tropp. Two proposals for robust PCA using semidefinite programming. *Electronic Journal of Statistics*, 5:1123–1160, 2011.
- Qing Qu, Zhihui Zhu, Xiao Li, Manolis C Tsakiris, John Wright, and René Vidal. Finding the sparsest vectors in a subspace: Theory, algorithms, and applications. *arXiv preprint arXiv:2001.06970*, 2020.

- Mostafa Rahmani and George Atia. Coherence pursuit: Fast, simple, and robust subspace recovery. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pp. 2864–2873. PMLR, 06–11 Aug 2017. URL <https://proceedings.mlr.press/v70/rahmani17a.html>.
- Manolis C. Tsakiris and René Vidal. Dual principal component pursuit. *Journal of Machine Learning Research*, 19(18):1–50, 2018.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- René Vidal, Yi Ma, and S Shankar Sastry. *Generalized principal component analysis*, volume 5. Springer, 2016.
- Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust PCA via outlier pursuit. *IEEE Transactions on Information Theory*, 58(5):3047–3064, 2012.
- C. You, D. Robinson, and R. Vidal. Provable self-representation based outlier detection in a union of subspaces. In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4323–4332, 2017a.
- Chong You, Daniel P Robinson, and René Vidal. Provable self-representation based outlier detection in a union of subspaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3395–3404, 2017b.
- Chong You, Zhihui Zhu, Qing Qu, and Yi Ma. Robust recovery via implicit bias of discrepant learning rates for double over-parameterization. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 17733–17744. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/cd42c963390a9cd025d007dacfa99351-Paper.pdf>.
- Teng Zhang and Gilad Lerman. A novel m-estimator for robust PCA. *The Journal of Machine Learning Research*, 15(1):749–808, 2014.
- Z. Zhu, T. Ding, M. C. Tsakiris, D. P. Robinson, and R. Vidal. A linearly convergent method for non-smooth non-convex optimization on the grassmannian with applications to robust subspace and dictionary learning. In *Neural Information Processing Systems (NIPS)*, 2019.
- Zhihui Zhu, Yifan Wang, Daniel Robinson, Daniel Naiman, René Vidal, and Manolis Tsakiris. Dual principal component pursuit: Improved analysis and efficient algorithms. In *Advances in Neural Information Processing Systems 2018*, volume 31. Curran Associates, Inc., 2018. URL <https://proceedings.neurips.cc/paper/2018/file/af21d0c97db2e27e13572cbf59eb343d-Paper.pdf>.

A APPENDIX

Theorem 1 (Theorem 3 of Zhu et al. (2018)) Let $\{\hat{\mathbf{b}}_k\}$ the sequence generated by the projected subgradient method in Zhu et al. (2018), with initialization $\hat{\mathbf{b}}_0$ such that

$$\theta_0 < \arctan\left(\frac{Nc_{\mathbf{X},\min}}{N\eta_{\mathbf{X}} + M\eta_{\mathbf{O}}}\right) \quad (19)$$

where θ_0 denotes the principal angle of \mathbf{b}^0 from $\mathcal{S}_{\mathcal{?}}$, and

$$Nc_{\mathbf{X},\min} \geq N\eta_{\mathbf{X}} + M\eta_{\mathbf{O}} \quad (20)$$

Let $\mu^0 := \frac{1}{4 \max\{Nc_{\mathbf{X},\min}, M\eta_{\mathbf{O},\max}\}}$. If $\mu^0 \leq \mu^0$ and the step size μ^k is updated according to a piece-wise geometrically diminishing rule given as

$$\mu^k = \begin{cases} \mu^0, & k < K_0 \\ \mu^0 \beta^{\lfloor (k - K_0)/K \rfloor}, & k \geq K_0 \end{cases} \quad (21)$$

where $\beta < 1$, $\lfloor \cdot \rfloor$ is the floor function, and $K_0, K \in \mathbb{N}$ are chosen such that

$$\begin{aligned} K_0 &\geq K^\lambda(\mu^0), \\ K &\geq \left(\sqrt{2}\beta\mu^0(Nc_{\mathbf{X}} - (N\eta_{\mathbf{X}} + M\eta_{\mathbf{O}}))\right)^{-1} \end{aligned}$$

where,

$$K^\lambda(\mu) := \frac{\tan(\theta_0)}{\mu(Nc_{\mathbf{X},\min} - \max\{1, \tan(\theta_0)\}(N\eta_{\mathbf{X}} + M\eta_{\mathbf{O}}))} \quad (22)$$

then for the angle θ_k between $\hat{\mathbf{b}}^k$ and $\mathcal{S}_{\mathcal{?}}$ it holds

$$\tan(\theta_k) \leq \begin{cases} \max\{\tan(\theta_0), \frac{\mu^0}{2\mu^k}\}, & k < K_0 \\ \frac{\mu^0}{2\mu^k} \beta^{\lfloor (k - K_0)/K \rfloor}, & k \geq K_0 \end{cases} \quad (23)$$

A.1 PROOF OF LEMMA 2

Lemma 2 In the continuous case, the discrete DPCP problem given in (3) is reformulated as,

$$\begin{aligned} \min_{\mathbf{B} \in \mathbb{R}^{D \times c^0}} \sum_{i=1}^{c^0} (p\mathbb{E}_{\mu_{\mathcal{S}^D-1}}[f_{\mathbf{b}_i}] + (1-p)\mathbb{E}_{\mu_{\mathcal{S}^D-1 \setminus \mathcal{S}}}[f_{\mathbf{b}_i}]) &= \sum_{i=1}^{c^0} \|\mathbf{b}_i\|_2 (pc_D + (1-p)c_d \cos(\phi_i)) \\ \text{s.t. } \|\mathbf{b}_i\|_2 &= 1, \quad i = 1, 2, \dots, c^0 \end{aligned} \quad (24)$$

where $f_{\mathbf{b}} : \mathcal{S}^{D-1} \rightarrow \mathbb{R}$, $f_{\mathbf{b}}(\mathbf{z}) = |\mathbf{z}^T \mathbf{b}|$, ϕ_i is the principal angle of \mathbf{b}_i from the inliers subspace \mathcal{S} and p is the probability of occurrence of an outlier.

Proof We define the discrete measures $\mu_{\mathbf{X}}, \mu_{\mathbf{O}}$ associated with the inliers and outliers, respectively as,

$$\mu_{\mathbf{X}}(\mathbf{z}) = \frac{1}{N} \sum_{j=1}^N \delta(\mathbf{z} - \mathbf{o}_j), \quad \mu_{\mathbf{O}}(\mathbf{z}) = \frac{1}{M} \sum_j^N \delta(\mathbf{z} - \mathbf{x}_j) \quad (25)$$

where $\delta(\cdot)$ is the Dirac function. Recall that,

$$\int_{\mathbf{z} \in \mathcal{S}^{D-1}} g(\mathbf{z}) \delta(\mathbf{z} - \mathbf{z}_0) d\mu_{\mathcal{S}^D-1} = g(\mathbf{z}_0) \quad (26)$$

where $g : \mathcal{S}^{D-1} \rightarrow \mathbb{R}$ and $\mu_{\mathcal{S}^D-1}$ is the uniform measure on \mathcal{S}^{D-1} .

The DPCP objective for the discrete version of the problem divided by $M + N$ can be written as,

$$\begin{aligned}
\frac{1}{M+N} \sum_{i=1}^{c^\circ} \|\tilde{\mathbf{x}}^{\triangleright} \mathbf{b}_i\|_1 &= \frac{1}{M+N} \sum_{i=1}^{c^\circ} \left(\|\mathbf{x}^{\triangleright} \mathbf{b}_i\|_1 + \|\mathbf{o}^{\triangleright} \mathbf{b}_i\|_1 \right) = \frac{1}{M+N} \sum_i^{c^\circ} \left(\sum_{j=1}^N |\mathbf{x}_j^{\triangleright} \mathbf{b}_i| + \sum_{j=1}^M |\mathbf{o}_j^{\triangleright} \mathbf{b}_i| \right) \\
&= \frac{1}{M+N} \sum_{i=1}^{c^\circ} \left(\sum_{j=1}^N \int_{\mathbf{z} \in \mathcal{S}^{D-1}} |\mathbf{z}^{\triangleright} \mathbf{b}_i| \delta(\mathbf{z} - \mathbf{x}_j) d\mu_{\mathcal{S}^{D-1}} + \sum_j^M \int_{\mathbf{z} \in \mathcal{S}^{D-1}} |\mathbf{z}^{\triangleright} \mathbf{b}_i| \delta(\mathbf{z} - \mathbf{o}_j) d\mu_{\mathcal{S}^{D-1}} \right) \\
&= \frac{1}{M+N} \sum_{i=1}^{c^\circ} \left(\int_{\mathbf{z} \in \mathcal{S}^{D-1}} |\mathbf{z}^{\triangleright} \mathbf{b}_i| \sum_{j=1}^N \delta(\mathbf{z} - \mathbf{x}_j) d\mu_{\mathcal{S}^{D-1}} + \int_{\mathbf{z} \in \mathcal{S}^{D-1}} |\mathbf{z}^{\triangleright} \mathbf{b}_i| \sum_j^M \delta(\mathbf{z} - \mathbf{o}_j) d\mu_{\mathcal{S}^{D-1}} \right) \\
&= \sum_{i=1}^{c^\circ} (p E_{\mu_{\mathcal{X}}} [f_{\mathbf{b}_i}] + (1-p) E_{\mu_{\mathcal{O}}} [f_{\mathbf{b}_i}])
\end{aligned}$$

Note that $\mu_{\mathcal{X}}, \mu_{\mathcal{O}}$ arise by discretizing the continuous uniform measures $\mu_{\mathcal{S}^{D-1}}$ and $\mu_{\mathcal{S}^{D-1} \setminus \mathcal{S}}$ respectively ($\mu_{\mathcal{S}^{D-1} \setminus \mathcal{S}}$ denotes the uniform measure on $\mathcal{S}^{D-1} \cap \mathcal{S}$) and p is the probability of occurrence of an outlier i.e., $\frac{M}{M+N} \rightarrow p$ as $M, N \rightarrow \infty$ ($1-p$ corresponds to the probability of occurrence of an inlier). That being said, the continuous version of DPCP can be simply stated by replacing $\mu_{\mathcal{X}}, \mu_{\mathcal{O}}$ with $\mu_{\mathcal{S}^{D-1}}$ and $\mu_{\mathcal{S}^{D-1} \setminus \mathcal{S}}$ in equation 27 as follows,

$$\min_{\mathbf{B}} \sum_{i=1}^{c^\circ} (p E_{\mu_{\mathcal{S}^{D-1}}} [f_{\mathbf{b}_i}] + (1-p) E_{\mu_{\mathcal{S}^{D-1} \setminus \mathcal{S}}} [f_{\mathbf{b}_i}]) \quad (27)$$

The RHS of (9) immediately shows up by invoking Proposition 4 in Tsakiris & Vidal (2018). \blacksquare

A.2 PROOF OF LEMMA 3

Lemma 3: A projected subgradient algorithm consisting of the steps described in (10) using a piecewise geometrically diminishing step size rule (see (8) in Theorem 1) will almost surely asymptotically converge to a matrix $\hat{\mathbf{B}} \in \mathbb{R}^{D \times c^\circ}$ whose columns $\hat{\mathbf{b}}_i$, $i = 1, 2, \dots, c^\circ$ will be normal vectors of the inliers' subspace when randomly initialized with vectors $\hat{\mathbf{b}}_i^0 \in \mathcal{S}^{D-1}$, $i = 1, 2, \dots, c^\circ$ uniformly distributed over the sphere \mathcal{S}^{D-1} .

Proof:

The proof can be trivially obtained by noticing a) that the condition for convergence i.e., inequality (7) of the projected subgradient algorithm given in Theorem 1 becomes $\theta_i^0 < \frac{\pi}{2}$ in the continuous case (since $\eta_{\mathcal{X}} \rightarrow 0$, $\eta_{\mathcal{O}} \rightarrow 0$, $c_{\mathcal{X}, \min} \rightarrow c_d > 0$) and b) the set of unit ℓ_2 -norm vectors $\hat{\mathbf{b}}_i^0$ s, $i = 1, 2, \dots, c^\circ$ sampled independently by a uniform distribution over the sphere and whose principal angle θ_i^0 is $\frac{\pi}{2}$ form the inliers' subspace has measure 0. \blacksquare

A.3 PROOF OF THEOREM 5

Lemma 4 The PSGM iterates $\hat{\mathbf{b}}_i^k$, $i = 1, 2, \dots, c^\circ$, $k = 1, 2, \dots$ given in (10), when randomly initialized with $\hat{\mathbf{b}}_i^0$ s, $i = 1, 2, \dots, c^\circ$ that are independently drawn from a spherical distribution with unit ℓ_2 norm converge almost surely to c° normal vectors of the inliers subspace \mathcal{S} denoted as $\hat{\mathbf{b}}_i$, $i = 1, 2, \dots, c^\circ$ that are given by

$$\hat{\mathbf{b}}_i = \frac{\mathcal{P}_{\mathcal{S}^\circ}(\hat{\mathbf{b}}_i^0)}{\|\mathcal{P}_{\mathcal{S}^\circ}(\hat{\mathbf{b}}_i^0)\|_2}, \quad i = 1, 2, \dots, c^\circ \quad (28)$$

Proof Let us assume $\mathbf{b}_i^0 = \hat{\mathbf{b}}_i^0$. The iterates of subgradients steps of PSGM can be written in the following form,

$$\begin{aligned}\mathbf{b}_i^1 &= (1 - \mu_i^0 p c_D) \hat{\mathbf{b}}_i^0 - \mu_i^0 (1 - p) c_d \hat{\mathbf{s}} \\ \mathbf{b}_i^2 &= (1 - \mu_i^1 p c_D) \hat{\mathbf{b}}_i^1 - \mu_i^1 (1 - p) c_d \hat{\mathbf{s}} \\ &\vdots \\ \mathbf{b}_i^K &= (1 - \mu_i^{K-1} p c_D) \hat{\mathbf{b}}_i^{K-1} - \mu_i^{K-1} (1 - p) c_d \hat{\mathbf{s}}\end{aligned}\quad (29)$$

By projecting each update of PSGM onto \mathcal{S}_γ and since $\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^k) = \frac{\|\mathbf{b}_i^k\|}{\|\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^k)\|} \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^k)$ we have,

$$\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^{k+1}) = \frac{(1 - \mu_i^k p c_D)}{\|\mathbf{b}_i^k\|} \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^k) \quad (30)$$

We can thus easily derive the following form for $\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^K)$,

$$\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^K) = \left(\prod_{k=1}^{K-1} \frac{(1 - \mu_i^k p c_D)}{\|\mathbf{b}_i^k\|_2} \right) \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{b}_i^0) \quad (31)$$

We know from Theorem 1 and Lemma 3 when DPCP-PSGM is initialized with $\hat{\mathbf{b}}_i^0$, $i = 1, 2, \dots, c^\ell$'s randomly drawn according to a spherical distribution then it will almost surely converge as $K \rightarrow \infty$ to vectors $\hat{\mathbf{b}}_i$, $i = 1, 2, \dots, c^\ell$ i.e., $\hat{\mathbf{b}}_i^K \rightarrow \hat{\mathbf{b}}_i$ where $\hat{\mathbf{b}}_i \in \mathcal{S}_\gamma$. Hence $\mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^K) \rightarrow \hat{\mathbf{b}}_i$ as $K \rightarrow \infty$. Note that from Theorem 1 we have that $\mu_i^k \neq \frac{1}{p c_D} \quad \forall k = \{1, 2, \dots, K\}$ hence $\prod_{k=1}^{K-1} \frac{(1 - \mu_i^k p c_D)}{\|\mathbf{b}_i^k\|_2} \neq 0$. From 31 and after projecting on the unit sphere and we thus have $\hat{\mathbf{b}}_i = \frac{\mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0)}{k \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0)_{k_2}}$. ■

Theorem 5 Let $\hat{\mathbf{B}}^0 \in \mathbb{R}^{D \times c^\ell}$ where $c^\ell \geq c$ with c denoting the true codimension of the inliers subspace \mathcal{S} , consisting of unit ℓ_2 norm column vectors $\hat{\mathbf{b}}_i^0 \in S^{D-1}$, $i = 1, 2, \dots, c^\ell$ that are independently drawn from uniform distribution over the sphere S^{D-1} . A PSGM algorithm initialized with $\hat{\mathbf{B}}^0$ will almost surely converge to a matrix $\hat{\mathbf{B}}$ such that $\text{span}(\hat{\mathbf{B}}) \equiv \mathcal{S}_\gamma$.

Proof From Lemma 4 we have that for each initial unit norm vector \mathbf{b}_i^0 which corresponds to the i th column of \mathbf{B}^0 will almost surely converge to $\hat{\mathbf{b}}_i = \frac{\mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0)}{k \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_i^0)_{k_2}}$. We can thus write $\mathbf{B} = \mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{B}^0) \mathbf{\Gamma}$ where $\mathbf{\Gamma}$ is a full-rank diagonal matrix given by $\mathbf{\Gamma} = \text{diag}(\frac{1}{k \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_1^0)_{k_2}}, \frac{2}{k \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_2^0)_{k_2}}, \dots, \frac{1}{k \mathcal{P}_{\mathcal{S}_\gamma}(\hat{\mathbf{b}}_{c^\ell}^0)_{k_2}})$. Note that $\mathcal{P}_{\mathcal{S}_\gamma}$ is a linear projection and thus we can write $\mathcal{P}_{\mathcal{S}_\gamma}(\mathbf{B}^0) = \mathbf{B}_{\mathcal{S}_\gamma} \mathbf{B}_{\mathcal{S}_\gamma}^\top$ where $\mathbf{B}_{\mathcal{S}_\gamma} \in \mathbb{R}^{D \times c}$ is an orthonormal matrix which spans \mathcal{S}_γ . Note that the probability of sampling a low-rank matrix $\mathbf{B}_0 = [\mathbf{b}_1^0, \mathbf{b}_2^0, \dots, \mathbf{b}_{c^\ell}^0]$ when columns \mathbf{b}_i^0 's are randomly and independently drawn from a spherical distribution is zero. We thus have $\mathbf{B} = \mathbf{B}_{\mathcal{S}_\gamma} \mathbf{B}_{\mathcal{S}_\gamma}^\top \mathbf{B}^0 \mathbf{\Gamma}$ with $\text{rank}(\mathbf{B}) = c$. ■

Lemma 6 If $\sigma_{c^\ell}(\mathbf{B}^0) > \|\hat{\mathbf{\Delta}}\|_2$ then the rank of matrix $\mathbf{B}^0 - \hat{\mathbf{\Delta}}$ equals c^ℓ .

Proof Let $\mathbf{A} = \mathbf{B}^0 - \hat{\mathbf{\Delta}}$. From singular value perturbation inequalities we have $|\sigma_i(\mathbf{A}) - \sigma_i(\mathbf{B}^0)| \leq \|\hat{\mathbf{\Delta}}\|_2$, for $i = 1, 2, \dots, c^\ell$. Hence it holds,

$$-\|\hat{\mathbf{\Delta}}\|_2 \leq \sigma_i(\mathbf{A}) - \sigma_i(\mathbf{B}^0) \quad (32)$$

If $\sigma_{c^\ell}(\mathbf{B}^0) > \|\hat{\mathbf{\Delta}}\|_2$ then from equation 32 we get

$$\sigma_{c^\ell}(\mathbf{A}) > 0 \quad (33)$$

hence the matrix $\mathbf{B}^0 - \hat{\mathbf{\Delta}}$ will be full-rank. ■

A.4 PROOF OF THEOREM 7

We first give the following Lemmas:

Lemma 7 For the ℓ_2 norm of $\mathbf{e}_{\mathcal{O}}^{i,k}$ for any $k = 1, 2, \dots, K$ and $i = 1, 2, \dots, c^0$ it holds,

$$\|\mathbf{e}_{\mathcal{O}}^{i,k}\|_2 \leq \eta_{\mathcal{O}} + c_{\mathcal{O},\max} - c_d \quad (34)$$

Proof

$$\begin{aligned} \|\mathbf{e}_{\mathcal{O}}^{i,k}\|_2 &= \|\mathbf{o}_b - c_d \mathbf{b}\|_2 = \left\| \frac{1}{M} \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \mathbf{b}) - c_d \mathbf{b} \right\|_2 \\ &= \left\| \frac{1}{M} (\mathbf{I} - \mathbf{b} \mathbf{b}^{\triangleright}) \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \mathbf{b}) + \frac{1}{M} \mathbf{b} \mathbf{b}^{\triangleright} \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \mathbf{b}) - c_d \mathbf{b} \right\|_2 \\ &\leq \left\| \frac{1}{M} (\mathbf{I} - \mathbf{b} \mathbf{b}^{\triangleright}) \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \mathbf{b}) \right\|_2 + \left(\frac{1}{M} \|\mathcal{O}^{\triangleright} \mathbf{b}\|_1 - c_d \right) \|\mathbf{b}\|_2 \\ &\leq \eta_{\mathcal{O}} + c_{\mathcal{O},\max} - c_d \end{aligned} \quad (35)$$

where we have used the fact that $\|\mathbf{b}\|_2 = 1$. ■

Lemma 8 Let the step size of Algorithm 1 (DPCP-PSGM) μ_i^k being updated following the piecewise geometrically diminishing step size rule with

$$\beta < \left(\frac{1 - \mu_i^0 M c_D}{1 + \mu_i^0 (N(\eta_{\mathcal{X}} + c_{\mathcal{X},\max}) + M(\eta_{\mathcal{O}} + c_{\mathcal{O},\max}))} \right)^K.$$

For the spectral norm of $\hat{\Delta}$ it holds $\|\hat{\Delta}\|_2 \leq \sqrt{c^0} \kappa (\eta_{\mathcal{O}} + c_{\mathcal{O},\max} - c_d)$ where $\kappa = \max_i \frac{M \mu_i^0}{\beta^{K_0/K} (1 - r_i)}$ and $r_i = \frac{(1 + \mu_i^0 (N(\eta_{\mathcal{X}} + c_{\mathcal{X},\max}) + M(\eta_{\mathcal{O}} + c_{\mathcal{O},\max})))}{1 - \mu_i^0 M c_D} \beta^{1/K}$.

Proof We first bound the ℓ_2 norm of vectors \mathbf{b}_i^j 's. We have that $\forall i = 1, 2, \dots, c^0$ and $j = 1, 2, \dots, K$ it holds

$$\mathbf{b}_i^{j+1} = \hat{\mathbf{b}}_i^j - \mu_i^j \left(\mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j) + \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \hat{\mathbf{b}}_i^j) \right) \quad (36)$$

We define the quantities

$$\eta_{\mathcal{X}} := \max_{\mathbf{b} \in \mathcal{S}^D} \frac{1}{N} \|(\mathcal{P}_S - \hat{\mathbf{b}}_i^j \hat{\mathbf{b}}_i^{j,\triangleright}) \mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 \quad (37)$$

$$c_{\mathcal{X},\max} := \max_{\mathbf{b} \in \mathcal{S}^D} \frac{1}{N} \|\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j\|_1 \quad (38)$$

$\|\hat{\mathbf{b}}_i^j\|_2 = 1$ hence

$$\begin{aligned} \|\mathbf{b}_i^{j+1}\|_2 &\leq 1 + \mu_i^j \|\mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j) + \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 \\ &\leq 1 + \mu_i^j \left(\|\mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 + \|\mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 \right) \\ &\leq 1 + \mu_i^j \left(\|(\mathcal{P}_S - \hat{\mathbf{b}}_i^j \hat{\mathbf{b}}_i^{j,\triangleright}) \mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 + \|(\hat{\mathbf{b}}_i^j \hat{\mathbf{b}}_i^{j,\triangleright}) \mathcal{X} \text{Sgn}(\mathcal{X}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 \right) \\ &\quad + \left\| \left(1 - \hat{\mathbf{b}}_i^j \hat{\mathbf{b}}_i^{j,\triangleright} \right) \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \hat{\mathbf{b}}_i^j) \right\|_2 + \|\hat{\mathbf{b}}_i^j \hat{\mathbf{b}}_i^{j,\triangleright} \mathcal{O} \text{Sgn}(\mathcal{O}^{\triangleright} \hat{\mathbf{b}}_i^j)\|_2 \\ &\leq 1 + \mu_i^j (N(\eta_{\mathcal{X}} + c_{\mathcal{X},\max}) + M(\eta_{\mathcal{O}} + c_{\mathcal{O},\max})) \end{aligned}$$

Due to equation 39 and since μ_i^j follows a non-increasing path as $j \rightarrow k$, the scalar term $\prod_{j=0}^k \frac{k \mu_i^j k_2}{(1 - \mu_i^j M c_D)}$ is bounded above as follows,

$$\prod_{j=0}^k \frac{\|\mathbf{b}_i^j\|_2}{(1 - \mu_i^j M c_D)} \leq \left(\frac{(1 + \mu_i^0 (N(\eta_{\mathcal{X}} + c_{\mathcal{X},\max}) + M(\eta_{\mathcal{O}} + c_{\mathcal{O},\max})))}{1 - \mu_i^0 M c_D} \right)^k \quad (39)$$

We now focus on the geometrically diminishing step size rule given in equation 8. We have $\mu_i^k = \mu_i^0 \beta^{bk - K_0/K} c^{+1} < \mu_i^0 \beta^{(k - K_0)/K}$ for $k \geq K_0$ and $\mu_i^0 \beta^{(k - K_0)/K} > \mu_i^0$ for $k < K_0$. Hence we can get the following upper bound

$$\begin{aligned} & \lim_{K \uparrow} \sum_{k=0}^{K-1} \prod_{j=0}^k \frac{\|\mathbf{b}_i^j\|_2}{(1 - \mu_i^j M c_D)} \mu_i^k M < \\ & M \lim_{K \uparrow} \sum_{k=0}^{K-1} \left(\frac{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))}{1 - \mu_i^0 M c_D} \right)^k \mu_i^0 \beta^{(k - K_0)/K} \\ & \equiv \lim_{K \uparrow} \sum_{k=0}^{K-1} M \frac{1}{\beta^{K_0/K}} \mu_i^0 \left(\frac{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))}{1 - \mu_i^0 M c_D} \beta^{1/K} \right)^k \end{aligned}$$

The series $S = \sum_{k=0}^{K-1} \left(\frac{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))}{1 - \mu_i^0 M c_D} \beta^{1/K} \right)^k$ is geometric and if

$$\beta < \left(\frac{1 - \mu_i^0 M c_D}{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))} \right)^K \quad (40)$$

it converges as $K \rightarrow \infty$ to $\frac{1}{1 - r_i}$ where $r_i = \frac{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))}{1 - \mu_i^0 M c_D} \beta^{1/K}$.

Let us now bound the ℓ_2 norms of the columns of $\hat{\Delta}$. From Lemma 7 we have $\|e_{\boldsymbol{o}}^{i,k}\|_2 \leq \boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max} - c_d$. We can easily thus derive that $\|\hat{\boldsymbol{\delta}}_i\|_2 \leq \kappa(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max} - c_d)$. For the spectral norm of $\hat{\Delta}$ we thus have

$$\begin{aligned} \|\hat{\Delta}\|_2 &= \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\hat{\Delta}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\sum_{i=1}^{c^0} \hat{\boldsymbol{\delta}}_i x_i\|_2}{\|\mathbf{x}\|_2} \leq \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\sum_{i=1}^{c^0} \|\hat{\boldsymbol{\delta}}_i\|_2 |x_i|}{\|\mathbf{x}\|_2} \\ &\leq \max_i \|\hat{\boldsymbol{\delta}}_i\|_2 \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{x}\|_1}{\|\mathbf{x}\|_2} \leq \kappa(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max} - c_d) \sqrt{c^0} \end{aligned} \quad (41)$$

Where the last inequality arises since $\|\mathbf{x}\|_1 \leq \sqrt{c^0} \|\mathbf{x}\|_2$. ■

We then give the Theorem.

Theorem 9 (Theorem 5.58 of Vershynin (2010)) Let \mathbf{B} be a $D \times d$ matrix ($D \geq d$) whose columns \mathbf{b}_i are independent sub-gaussian isotropic random vectors in \mathbb{R}^D with $\|\mathbf{b}_i\|_2 = \sqrt{D}$ almost surely. Then for every $t \geq 0$ the inequality

$$\sqrt{D} - C\sqrt{d} - t \leq \sigma_{\min}(\mathbf{B}) \leq \sigma_{\max}(\mathbf{B}) \leq \sqrt{D} + C\sqrt{d} + t \quad (42)$$

with probability at least $1 - 2 \exp(-ct^2)$, where $C = C_K^0$, $c = c_K^0 > 0$ depend only on the subgaussian norm $K = \max_j \|\mathbf{b}_i\|_{\psi_2}$ of the columns.

The proof of Theorem 7 follows next.

Theorem 7 Let $\mathbf{B}^0 \in \mathbb{R}^{D \times c^0}$ with columns randomly sampled from a unit ℓ_2 norm spherical distribution where $c^0 \geq c$ with c denoting the true codimension of the inliers subspace \mathcal{S} that satisfies Assumption 1. If

$$1 - C_1 \sqrt{\frac{c^0}{D}} - \frac{\epsilon}{\sqrt{D}} > \sqrt{c^0} \kappa(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max} - c_d) \quad (43)$$

where $\kappa = \max_i \frac{M \mu_i^0}{\beta^{K_0/K} (1 - r_i)}$ and $r_i = \frac{(1 + \mu_i^0 (N(\boldsymbol{\eta}\boldsymbol{x} + c\boldsymbol{x}_{,\max}) + M(\boldsymbol{\eta}\boldsymbol{o} + c_{\boldsymbol{o},\max})))}{1 - \mu_i^0 M c_D} \beta^{1/K}$ then with probability at least $1 - 2 \exp(-\epsilon^2 C_2)$ (where C_1, C_2 are absolute constants), Algorithm 1 with a geometrically diminishing step size rule will converge to a matrix $\hat{\mathbf{B}}$ such that $\text{span}(\hat{\mathbf{B}}) \equiv \mathcal{S}$.

Table 1: Results on Washington DC AVIRIS hyperspectral image

Methods	F1-scores	
	r = 80%	r = 90%
DPCP-PSGM (unknown c)	0.994	0.993
RSGM (unknown c)	0	0
DPCP-IRLS (unknown c)	0	0
RSGM ($c = 5$)	0.999	0.993
DPCP-IRLS $c = 5$	1	0.995

Proof By Assumption 1 we have that all columns of \mathbf{B}_0 will satisfy the sufficient condition for convergence of DPCP-PSGM (Algorithm 1) to a normal vector of \mathcal{S} . From Lemma 6 and we use the inequality $\sigma_{c^0}(\mathbf{B}_0) > \|\hat{\Delta}\|_2$ which ensures full-rankness of $\hat{\mathbf{B}}$, which is the key ingredient in order to prove that $\text{span}(\hat{\mathbf{B}}) = \mathcal{S}$. We can then Use Theorem 9 for matrix \mathbf{B}_0 . Note that columns of \mathbf{B}_0 are drawn independently and are uniformly distributed on the unit sphere. Hence, columns of \mathbf{B}_0 are sampled by subgaussian distribution and the LHS of the inequality of the theorem appears if we scale with $\frac{1}{\sqrt{D}}$ so that to create unit-norm columns and use LHS of the inequality of Theorem 7. The RHS of the inequality is due to the upper bound of $\|\hat{\Delta}\|_2$ as stated in Lemma 8. The absolute constants C_1, C_2 depend only the subgaussian norm of the uniform distribution (they is no dependency on the dimensions of the problem). ■

B EXPERIMENTAL DETAILS AND ADDITIONAL MATERIAL

All experiments were conducted on a MacBook Pro 2.6Ghz 6-Core Intel Core i7, memory 16GB 2667 Mhz DDR using Matlab2019B. For computational purposes and in order to avoid fine-tuning of the piecewise geometrically diminishing (PGD) step size, the modified backtracking line-search (MBLS) step-size rule was adopted for DPCP-PSGM as proposed in Zhu et al. (2018). We define the distance between two subspaces spanned by matrices \mathbf{B} and \mathbf{A} as $\text{dist}(\mathbf{B}, \mathbf{A}) = \min_{\mathbf{Q} \in \mathcal{O}(D, c)} \|\mathbf{B} - \mathbf{A}\mathbf{Q}\|_F$ where $\mathcal{O}(D, c)$ denotes the Stiefel manifold of orthogonal matrices of rank c . Note that $\text{dist}(\mathbf{B}, \mathbf{A}) = 0 \iff \text{span}(\mathbf{B}) \equiv \text{span}(\mathbf{A})$ (see Zhu et al. (2019)).

B.1 OUTLIERS PURSUIT IN WASHINGTON DC MALL AVIRIS HSI

Hyperspectral images (HSIs) provide rich spectral information as compared to RGB images capturing a wide range of the electromagnetic spectrum. Washington DC Mall AVIRIS HSI contains contiguous spectral bands captured at 0.4 to 2.4 μm region of visible and infrared spectrum, Giampouras et al. (2019). In this experiments we randomly choose 10 out of its 210 spectral bands. Due to high coherence in the both the spectral and the spatial domain, pixels of HSIs admit representations in low-dimensional subspaces. Here, we use a 100×100 segment of the hyperspectral image selecting randomly 10 out of its $D = 210$ spectral bands. We form a matrix $\tilde{\mathcal{X}}$ of size 10×10000 whose columns correspond to different points in the 10-dimensional ambient space. Then we corrupt columns of $\tilde{\mathcal{X}}$ by replacing them with outliers that are generated uniformly at random with unit ℓ_2 norm for two different outliers' ratios i.e., $r = 0.8$ and $r = 0.9$. In the corrupted $\tilde{\mathcal{X}}$, the remaining clear pixels are considered as the inliers. Table 1 displays the F1 scores obtained by DPCP-PSGM, RSGM and DPCP-IRLS algorithm. The latter two algorithms are evaluated in two scenarios: a) codimension is initialized $c^0 = 5$ and b) $c^0 = 10$. Given the singular value distribution of the initial image, we infer that the dimension d of the inliers' subspace is less or equal than 5.

Hence, $c^0 = 5$ (recall $c = D - d$) is close to the true codimension value while $c^0 = 10$ is an overestimate thereof. From Table 1, we can see that the proposed DPCP-PSGM succeeds in both outliers' ratios regardless its unawareness of the true codimension value. On the other hand, DPCP-IRLS and RSGM fail when initialized with $c = 10$ and this is attributed to the restrictions induced

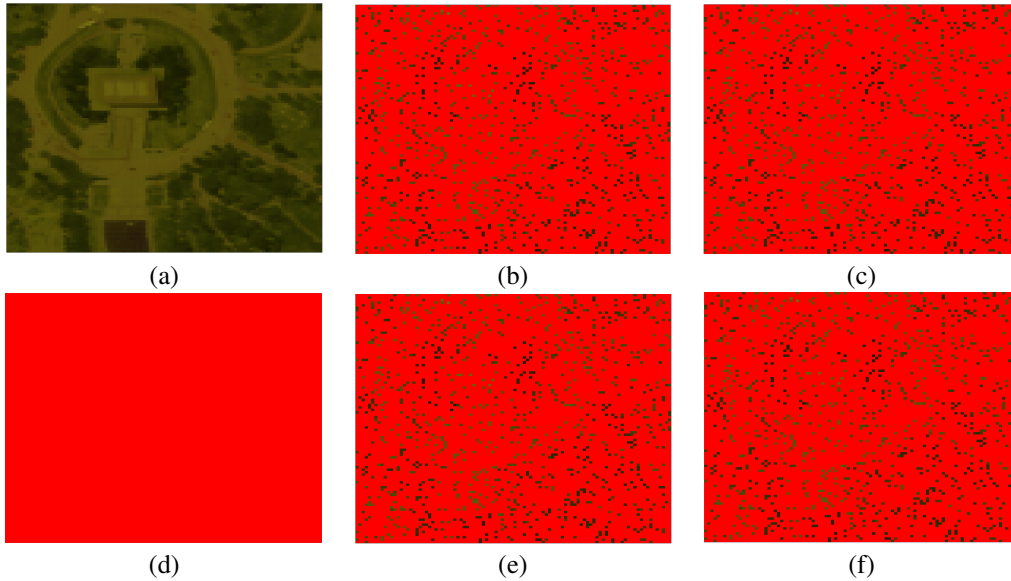


Figure 4: (a) False RGB color image of the clean version of Washington Mall AVIRIS HSI, (b) corrupted by outliers depicted with red and inliers corresponding the non-red pixels (c) annotated outliers as recovered by the proposed DPCP-PSGM method initialized with $c^l = 10$ (d) RSGM with $c^l = 10$, (e) RSGM with $c^l = 5$ and (f) DPCP-IRLS with $c^l = 5$.

due to the orthogonality constraints they both impose. In Fig. 4 we provide annotated versions of the clean HSI, its corrupted by outliers version for outliers' ratio $r = 90\%$, and the annotated outliers as recovered by the proposed DPCP-PSGM, RSGM, RSGM with $c^l = 5$ and DPCP-IRLS with $c^l = 5$.