

Fun with Flags: Robust Principal Directions via Flag Manifolds

Nathan Mankovich
University of Valencia

Gustau Camps-Valls
University of Valencia

Tolga Birdal
Imperial College London

A. Theoretical Justifications & Discussions

On the unifying aspects of our framework. In our framework, the link between RPCA & Dual-PCA, established also in the discussed earlier works, emerges as a by-product of our unifying formulation. To elucidate, our flag-based framework allows for: (i) extending DPCP to manifold-valued data (fTDPCP), (ii) interpolating between L_1/L_2 -DPCP via the use of non-trivial flag types, and (iii) an efficient algorithms for computing flag-(tangent) DPCP for any flag type. To the best of our knowledge, Alg. 1 (main paper) is the only method for finding non-trivial flags of robust directions and when used for both fRPCA & fWPCA.

A.1. Proof of Prop. 3

Let us recall the proposition before delving into the proof.

Proposition 1 (Stiefel optimization of (weighted) fPCA). *Suppose we have weights $\{w_{ij}\}_{i=1, j=1}^{i=k, j=p}$ for a dataset $\{\mathbf{x}_j\}_{j=1}^p \subset \mathbb{R}^n$ along with a flag type $(n_1, n_2, \dots, n_k; n)$. We store the weights in the diagonal weight matrices $\{\mathbf{W}_i\}_{i=1}^k$ with diagonals $(\mathbf{W}_i)_{jj} = w_{ij}$. If*

$$\mathbf{U}^* = \arg \max_{\mathbf{U} \in \text{St}(n_k, n)} \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T \mathbf{U} \mathbf{I}_i) \quad (1)$$

where \mathbf{I}_i is determined as a function of the flag signature. For example, for $\mathcal{FL}(n+1)$:

$$(\mathbf{I}_i)_{l,s} = \begin{cases} 1, & l = s \in \{n_{i-1} + 1, n_{i-1} + 2, \dots, n_i\} \\ 0, & \text{otherwise} \end{cases}$$

Then $\llbracket \mathbf{U}^* \rrbracket = \llbracket \mathbf{U} \rrbracket^*$ is the weighted fPCA of the data with the given weights (e.g., solves ??) as long as we restrict ourselves to a region on $\mathcal{FL}(n+1)$ and $\text{St}(n_k, n)$ where weighted fPCA is convex.

Proof. First we will show that the flag and Stiefel objective functions are equivalent. Take

$$\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1) = \mathcal{FL}(n_1, n_2, \dots, n_k; n). \quad (2)$$

We decompose $\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_k]$ where $\mathbf{U}_i \in \mathbb{R}^{n \times m_i}$ and $\sum_{l=1}^i m_l = n_i$. Using \mathbf{I}_i (defined above) we have $\mathbf{U} \mathbf{I}_i \mathbf{U}^T = \mathbf{U}_i$.

Recall the objective function for both fRPCA and fD-PCP is

$$\mathbb{E}_j \left[\sum_{i=1}^k w_{ij} \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2^2 \right] = \sum_{j=1}^p \sum_{i=1}^k w_{ij} \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2^2, \quad (3)$$

$$= \sum_{j=1}^p \sum_{i=1}^k w_{ij} \|\mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2^2 \quad (4)$$

Using the definition of norms and $\mathbf{U}_i^T \mathbf{U}_i = \mathbf{I}$, Eq. (3) is equivalent to

$$\sum_{j=1}^p \sum_{i=1}^k w_{ij} \text{tr}(\mathbf{x}_j^T \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j) \quad (5)$$

Now, using properties of trace, matrix multiplication, and our handy $\{\mathbf{I}_i\}_{i=1}^k$ we reach our desired result

$$\sum_{j=1}^p \sum_{i=1}^k w_{ij} \text{tr}(\mathbf{U}_i^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i), \quad (6)$$

$$= \sum_{i=1}^k \text{tr} \left(\mathbf{U}_i^T \left(\sum_{j=1}^p w_{ij} \mathbf{x}_j \mathbf{x}_j^T \right) \mathbf{U}_i \right), \quad (7)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}_i^T (\mathbf{X} \mathbf{W}_i \mathbf{X}^T) \mathbf{U}_i), \quad (8)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}_i \mathbf{U}_i^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T), \quad (9)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T), \quad (10)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T \mathbf{U} \mathbf{I}_i). \quad (11)$$

So we have shown that the flag and Stiefel objective functions are equivalent.

Finally, we show $\llbracket \mathbf{U}^* \rrbracket = \llbracket \mathbf{U} \rrbracket^*$. Notice that the objective function for weighted flag PCA is invariant to different flag manifold representatives. First, let f denote the objective function in Eq. (11). Suppose \mathbf{U}^* solves

$\arg \max_{\mathbf{Y} \in St(n_k, n)} f(\mathbf{Y})$. Then take some other representative for $\llbracket \mathbf{U}^* \rrbracket$, namely $\mathbf{U}^* \mathbf{M}$ where

$$\mathbf{M} = \begin{bmatrix} \mathbf{M}_1 & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{M}_2 & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{M}_k \end{bmatrix} \text{ and } \mathbf{M}_1 \in O(m_i). \quad (12)$$

Then $f(\mathbf{U}^* \mathbf{M}) = f(\mathbf{U}^*)$ because

$$f(\mathbf{U}^* \mathbf{M}) = \sum_{i=1}^k \text{tr}((\mathbf{U}_i^* \mathbf{M})^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T (\mathbf{U}_i^* \mathbf{M})), \quad (13)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}_i^* \mathbf{M} \mathbf{M}^T \mathbf{U}_i^{*T} \mathbf{X} \mathbf{W}_i \mathbf{X}^T), \quad (14)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}_i^* \mathbf{U}_i^{*T} \mathbf{X} \mathbf{W}_i \mathbf{X}^T), \quad (15)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}_i^{*T} \mathbf{X} \mathbf{W}_i \mathbf{X}^T \mathbf{U}_i^*), \quad (16)$$

$$= f(\mathbf{U}^*). \quad (17)$$

So $f(\cdot)$ has the same value for any representative for $\llbracket \mathbf{U}^* \rrbracket$. Since $f(\mathbf{U}^*) \geq f(\mathbf{Y})$ for all $\mathbf{Y} \in St(n_k, n)$, then

$$f(\mathbf{U}^* \mathbf{M}) = f(\mathbf{U}^*) \geq f(\mathbf{Y}) = f(\mathbf{U}^* \mathbf{O}) \quad (18)$$

for all $\llbracket \mathbf{Y} \rrbracket \in \mathcal{FL}(n+1)$ where \mathbf{O} is of the same block structure as \mathbf{M} .

Recall $\llbracket \mathbf{U}^* \rrbracket \in \mathcal{FL}(n+1)$ maximizes f , so $f(\mathbf{U}) \geq f(\mathbf{Y})$ for all $\llbracket \mathbf{Y} \rrbracket \in \mathcal{FL}(n+1)$ and since $f(\cdot)$ has the same value for any representative of $\llbracket \mathbf{Y} \rrbracket$, we have $f(\mathbf{U}) \geq f(\mathbf{Y})$ for all $\mathbf{Y} \in St(n_k, n)$.

Recall, that $f(\mathbf{U}^*) \geq f(\mathbf{Y})$ for all $\mathbf{Y} \in St(n_k, n)$. So $f(\mathbf{U}^*) = f(\mathbf{U})$. Since f has a unique maximizer over $\mathcal{FL}(n+1)$, we have $\llbracket \mathbf{U}^* \rrbracket = \llbracket \mathbf{U} \rrbracket^* = \arg \max_{\llbracket \mathbf{Y} \rrbracket \in \mathcal{FL}(n+1)} f(\mathbf{Y})$. \square

A.2. Proof of Prop. 4

Let us recall the proposition before delving into the proof.

Proposition 2 (Stiefel optimization for flagged Robust (Dual-)PCAs). *We can formulate fRPCA, fWPCA, fDPCP, and fWDPCP as optimization problems over the Stiefel*

manifold using $\llbracket \mathbf{U} \rrbracket^ = \llbracket \mathbf{U}^* \rrbracket$ and the following:*

$$\mathbf{U}^* = \begin{cases} \arg \max_{\mathbf{U} \in St(n, n_k)} \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{P}_i^+ \mathbf{U} \mathbf{I}_i), & (\text{fRPCA}) \\ \arg \min_{\mathbf{U} \in St(n, n_k)} \sum_{i=1}^k \text{tr}(\mathbf{P}_i^- - \mathbf{U}^T \mathbf{P}_i^- \mathbf{U} \mathbf{I}_i), & (\text{fWPCA}) \end{cases} \quad (19)$$

$$\mathbf{U}^* = \begin{cases} \arg \min_{\mathbf{U} \in St(n, n_k)} \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{P}_i^+ \mathbf{U} \mathbf{I}_i), & (\text{fDPCP}) \\ \arg \max_{\mathbf{U} \in St(n, n_k)} \sum_{i=1}^k \text{tr}(\mathbf{P}_i^- - \mathbf{U}^T \mathbf{P}_i^- \mathbf{U} \mathbf{I}_i) & (\text{fWDPCP}) \end{cases} \quad (20)$$

where $\mathbf{P}^- = \mathbf{X} \mathbf{W}_i^- (\llbracket \mathbf{U} \rrbracket) \mathbf{X}^T$, $\mathbf{P}^+ = \mathbf{X} \mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket) \mathbf{X}^T$ and $\mathbf{W}_i^- (\llbracket \mathbf{U} \rrbracket)$, $\mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket)$ are defined in ?? as long as we restrict ourselves to a region on $\mathcal{FL}(n+1)$ and $St(n_k, n)$ where flag robust and dual PCAs are convex.

Proof. First, we write the objective functions for fRPCA and fDPCP over $St(n_k, n)$ using ?? to define each \mathbf{W}_i^+ as

$$f^+(\mathbf{U}) = \mathbb{E} \left[\sum_{i=1}^k \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2 \right], \quad (21)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2, \quad (22)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\text{tr}(\mathbf{x}_j^T \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j)}, \quad (23)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\text{tr}(\mathbf{U}_i^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i)}, \quad (24)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U} \mathbf{I}_i)}, \quad (25)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \frac{\text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U} \mathbf{I}_i)}{\sqrt{\text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U} \mathbf{I}_i)}}, \quad (26)$$

$$= \sum_{i=1}^k \text{tr} \left(\mathbf{U}^T \sum_{j=1}^p \frac{\mathbf{x}_j \mathbf{x}_j^T}{\|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2} \mathbf{U} \mathbf{I}_i \right), \quad (27)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{W}_i^+ \mathbf{X}^T \mathbf{U} \mathbf{I}_i), \quad (28)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{P}_i^+ \mathbf{U} \mathbf{I}_i). \quad (29)$$

Now we write the objective functions for fWPCA and

fWDPCP over $St(n_k, n)$ using ?? to define each \mathbf{W}_i^- as

$$f^-(\mathbf{U}) = \mathbb{E} \left[\sum_{i=1}^k \|\mathbf{x}_j - \pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2 \right], \quad (30)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \|\mathbf{x}_j - \pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2, \quad (31)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\text{tr}(\mathbf{x}_j^T \mathbf{x}_j - \mathbf{x}_j^T \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j)}, \quad (32)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\mathbf{x}_j^T \mathbf{x}_j - \text{tr}(\mathbf{U}_i^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i)}, \quad (33)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \sqrt{\mathbf{x}_j^T \mathbf{x}_j - \text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i)}, \quad (34)$$

$$= \sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{x}_j^T \mathbf{x}_j - \text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i)}{\sqrt{\mathbf{x}_j^T \mathbf{x}_j - \text{tr}(\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i)}}, \quad (35)$$

$$= \sum_{j=1}^p \frac{\mathbf{x}_j \mathbf{x}_j^T}{\|\mathbf{x}_j - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} \quad (36)$$

$$- \sum_{i=1}^k \text{tr} \left(\mathbf{U}^T \sum_{j=1}^p \frac{\mathbf{x}_j \mathbf{x}_j^T}{\|\mathbf{x}_j - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} \mathbf{U}_i \right) \quad (37)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{X} \mathbf{W}_i^- \mathbf{X}^T - \mathbf{U}^T \mathbf{X} \mathbf{W}_i^- \mathbf{X}^T \mathbf{U}_i), \quad (38)$$

$$= \sum_{i=1}^k \text{tr}(\mathbf{P}^- - \mathbf{U}^T \mathbf{P}^- \mathbf{U}_i). \quad (39)$$

Now, we can write the Lagrangians for these problems with the symmetric matrix of Lagrange multipliers $\mathbf{\Lambda}$ as

$$\begin{aligned} \mathcal{L}^+(\mathbf{U}) &= f^+(\mathbf{U}) + \text{tr}(\mathbf{\Lambda}^+(\mathbf{I} - \mathbf{U}^T \mathbf{U})), \\ \mathcal{L}^-(\mathbf{U}) &= f^-(\mathbf{U}) + \text{tr}(\mathbf{\Lambda}^-(\mathbf{I} - \mathbf{U}^T \mathbf{U})). \end{aligned}$$

Then, we collect our gradients in the following equations

$$\nabla_{\mathbf{U}} \mathcal{L}^+ = \sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} - 2\mathbf{U} \mathbf{\Lambda}^+ \quad (40)$$

$$\nabla_{\mathbf{U}} \mathcal{L}^- = - \sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{x}_j \mathbf{x}_j^T - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} - 2\mathbf{U} \mathbf{\Lambda}^- \quad (41)$$

$$\nabla_{\mathbf{\Lambda}^+} \mathcal{L}^+ = \nabla_{\mathbf{\Lambda}^-} \mathcal{L}^- = \mathbf{I} - \mathbf{U}^T \mathbf{U}. \quad (42)$$

$$(43)$$

Then setting $\nabla_{\mathbf{U}} \mathcal{L}_1 = \mathbf{0}$, $\nabla_{\mathbf{\Lambda}_1} \mathcal{L}_1 = \mathbf{0}$, left multiplying

by \mathbf{U}^T , and playing with properties of trace results in

$$\sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T \mathbf{U}_i) = 2\text{tr}(\mathbf{\Lambda}^+), \quad (44)$$

$$\sum_{i=1}^k \text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{W}_i \mathbf{X}^T \mathbf{U}_i) = -2\text{tr}(\mathbf{\Lambda}^-). \quad (45)$$

Then we have the following cases: we choose

- (fRPCA) \mathbf{U}^* to maximize $\text{tr}(\mathbf{\Lambda}^+)$ so that we maximize f^+ ,
- (fDPCP) \mathbf{U}^* to minimize $\text{tr}(\mathbf{\Lambda}^+)$ so that we minimize f^+ ,
- (fWPCA) \mathbf{U}^* to minimize $-\text{tr}(\mathbf{\Lambda}^-)$ so that we minimize f^- ,
- (fWDPCP) \mathbf{U}^* to maximize $-\text{tr}(\mathbf{\Lambda}^-)$ so that we maximize f^- .

$$\sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{U}_i^T \mathbf{x}_j\|_2} = \mathbf{\Lambda}_1 \mathbf{U}, \quad (46)$$

$$\sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{U}_i^T \mathbf{x}_j\|_2} = \mathbf{\Lambda}_1, \quad (47)$$

$$\sum_{j=1}^p \sum_{i=1}^k (\mathbf{W}_i)_{jj} \mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i = \mathbf{\Lambda}_1, \quad (48)$$

$$\text{tr} \left(\sum_{i=1}^k \mathbf{U}^T \left(\sum_{j=1}^p (\mathbf{W}_i)_{jj} \mathbf{x}_j \mathbf{x}_j^T \right) \mathbf{U}_i \right) = \text{tr}(\mathbf{\Lambda}_1), \quad (49)$$

$$\sum_{i=1}^k \text{tr} \left(\mathbf{U}^T \left(\sum_{j=1}^p (\mathbf{W}_i)_{jj} \mathbf{x}_j \mathbf{x}_j^T \right) \mathbf{U}_i \right) = \text{tr}(\mathbf{\Lambda}_1), \quad (50)$$

$$h_{\llbracket \mathbf{U} \rrbracket}(\mathbf{U}) = \text{tr}(\mathbf{\Lambda}_1). \quad (51)$$

Similarly, setting $\nabla_{\mathbf{U}} \mathcal{L}_2 = \mathbf{0}$, $\nabla_{\mathbf{\Lambda}_2} \mathcal{L}_2 = \mathbf{0}$ and leveraging ?? to define $\{\mathbf{W}_i\}_i$ results in

$$\sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{x}_j - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} = \mathbf{\Lambda}_2 \mathbf{U}, \quad (52)$$

$$- \sum_{j=1}^p \sum_{i=1}^k \frac{\mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i}{\|\mathbf{x}_j - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2} = \mathbf{\Lambda}_2, \quad (53)$$

$$- \sum_{j=1}^p \sum_{i=1}^k (\mathbf{W}_i)_{jj} \mathbf{U}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{U}_i = \mathbf{\Lambda}_2, \quad (54)$$

$$- \sum_{i=1}^k \text{tr}(\mathbf{U}^T (\mathbf{X} \mathbf{W}_i \mathbf{X}^T) \mathbf{U}_i) = \text{tr}(\mathbf{\Lambda}_2), \quad (55)$$

$$-h_{\llbracket \mathbf{U} \rrbracket}(\mathbf{U}) = \text{tr}(\mathbf{\Lambda}_2). \quad (56)$$

Finally, using a similar argument to that for the proof of the Stiefel optimization of fPCA leveraging assumed convexity, we have that $\llbracket \mathbf{U}^* \rrbracket = \llbracket \mathbf{U} \rrbracket^*$. \square

A.3. Proof of Prop. 5

We now prove the convergence of our algorithm. Let us recall the proposition from the main paper before delving into the proof.

Proposition 3 (Convergence of ?? for fDPCP). *?? for fD-PCP converges as long as $\|\mathbf{U}\mathbf{I}_i\mathbf{U}^T\mathbf{x}_j\|_2 \geq \epsilon \forall i, j$ and we restrict ourselves to a region on $\mathcal{FL}(n+1)$ and $St(n_k, n)$ where fDPCP is convex.*

Proof. This proof follows closely to what was done in [1]. First let $f^+ : \mathcal{FL}(n+1) \times \mathcal{FL}(n+1) \rightarrow \mathbb{R}$ denote the fDPCP objective function and $T : \mathcal{FL}(n+1) \rightarrow \mathcal{FL}(n+1)$ denote an iteration of ?. Then, assuming that $\|\mathbf{U}\mathbf{I}_i\mathbf{U}^T\mathbf{x}_j\|_2 \geq \epsilon$ for $i = 1, 2, \dots, k$ and $j = 1, 2, \dots, p$, we define the function $h : \mathcal{FL}(d+1) \times \mathcal{FL}(d+1) \rightarrow \mathbb{R}$ as

$$h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket) = \sum_{i=1}^p \text{tr}(\mathbf{Z}^T \mathbf{X} \mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket) \mathbf{X}^T \mathbf{Z} \mathbf{I}_i), \quad (57)$$

using the definition in ?? for $\mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket)$. Some algebra reduces $h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket)$ to

$$h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket) = \sum_{i=1}^p \sum_{j=1}^k \frac{\|\mathbf{Z} \mathbf{I}_i \mathbf{Z}^T \mathbf{x}_j\|_2^2}{\|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2}. \quad (58)$$

From Eq. (57), we see that $h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket)$ is the weighted flag PCA objective function of $\{\mathbf{x}_j\}_{j=1}^p$ with weights on the diagonals of $\mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket)$. The weighted flagged orthogonal PCA (f \perp PCA) optimization problem with weights on the diagonals $\mathbf{W}_i^+ (\llbracket \mathbf{U} \rrbracket)$ can be solved using a similar algorithm to Alg. 1 by just minimizing instead of maximizing (see Alg. 2). Thus minimizing $h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket)$ over $\llbracket \mathbf{Z} \rrbracket$ is an iteration of Alg. 2 for fDPCP which means

$$T(\llbracket \mathbf{U} \rrbracket) = \underset{\llbracket \mathbf{Z} \rrbracket \in \mathcal{FL}(d+1)}{\text{arg min}} h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket). \quad (59)$$

Using this, we have

$$h(T(\llbracket \mathbf{U} \rrbracket), \llbracket \mathbf{U} \rrbracket) \leq h(\llbracket \mathbf{U} \rrbracket, \llbracket \mathbf{U} \rrbracket). \quad (60)$$

By the definition of h

$$h(\llbracket \mathbf{U} \rrbracket, \llbracket \mathbf{U} \rrbracket) = \sum_{i=1}^p \sum_{j=1}^k \frac{\|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2^2}{\|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2}, \quad (61)$$

$$= \sum_{i=1}^p \sum_{j=1}^k \|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2, \quad (62)$$

$$= f(\llbracket \mathbf{U} \rrbracket). \quad (63)$$

This means, we have

$$h(T(\llbracket \mathbf{U} \rrbracket), \llbracket \mathbf{U} \rrbracket) \leq f(\llbracket \mathbf{U} \rrbracket). \quad (64)$$

Now we use the identity from algebra: $\frac{a^2}{b} \geq 2a - b$ for any $a, b \in \mathbb{R}$ and $b > 0$. Let

$$a = \|\mathbf{Z} \mathbf{I}_i \mathbf{Z}^T \mathbf{x}_j\|_2 \text{ and } b = \|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2. \quad (65)$$

Then

$$h(\llbracket \mathbf{Z} \rrbracket, \llbracket \mathbf{U} \rrbracket) \geq 2 \sum_{j=1}^p \sum_{i=1}^k \|\mathbf{Z} \mathbf{I}_i \mathbf{Z}^T \mathbf{x}_j\|_2 \quad (66)$$

$$- \sum_{j=1}^p \sum_{i=1}^k \|\mathbf{U} \mathbf{I}_i \mathbf{U}^T \mathbf{x}_j\|_2, \quad (67)$$

$$= 2f(\llbracket \mathbf{Z} \rrbracket) - f(\llbracket \mathbf{U} \rrbracket). \quad (68)$$

Now, take $\llbracket \mathbf{Z} \rrbracket = T(\llbracket \mathbf{U} \rrbracket)$. This gives us

$$h(T(\llbracket \mathbf{U} \rrbracket), \llbracket \mathbf{U} \rrbracket) \geq 2f(T(\llbracket \mathbf{U} \rrbracket)) - f(\llbracket \mathbf{U} \rrbracket). \quad (69)$$

Then, combining Eq. 69 with Eq. 64, we have

$$2f(T(\llbracket \mathbf{U} \rrbracket)) - f(\llbracket \mathbf{U} \rrbracket) \leq f(\llbracket \mathbf{U} \rrbracket), \quad (70)$$

$$f(T(\llbracket \mathbf{U} \rrbracket)) \leq f(\llbracket \mathbf{U} \rrbracket). \quad (71)$$

Finally, notice that the real sequence with terms $f^+(T(\llbracket \mathbf{U}^{(m-1)} \rrbracket)) = f^+(\llbracket \mathbf{U}^{(m)} \rrbracket) \in \mathbb{R}$ is bounded below by 0 and is decreasing. So it converges as $m \rightarrow \infty$. \square

B. Further Notes on Flagified PCA

We now generalize PCA and its variants using flags by grouping eigenvectors using the flag type. The PCA optimization problem is naturally an optimization problem on the Stiefel manifold, $St(k, n) := \{\mathbf{U} \in \mathbb{R}^{k \times n} : \mathbf{U}^T \mathbf{U} = \mathbf{I}\}$. Suppose $\mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \in St(k, n)$ are the $k < n$ principal components of a data matrix \mathbf{X} . These are naturally ordered according to their decreasing associated objective function values¹. This results in the nested subspace structure

$$\llbracket \mathbf{U} \rrbracket = [\mathbf{u}_1] \subset [\mathbf{u}_1, \mathbf{u}_2] \subset \dots \subset [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \subset \mathbb{R}^n. \quad (72)$$

So one can think of $\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(1, 2, \dots, k; n)$, and consequently, reformulate PCA as an optimization problem over $\mathcal{FL}(1, 2, \dots, k; n)$. Thinking of \mathbf{U} as $\llbracket \mathbf{U} \rrbracket$ emphasizes the nested subspace structure of the principal components according to their associated objective function values.

What if we have multiple principal components with the same objective function value? In other words, suppose we have at least one eigenvalue of $\mathbf{X}\mathbf{X}^T$ with a geometric multiplicity greater than 1? For example, assume

¹The objective function values are also referred to as explained variances, eigenvalues and squared singular values

our dataset has a large variance on some 2-plane, and all other directions orthogonal to that plane have smaller, unequal variance. Then, the first two principal components, \mathbf{u}_1 and \mathbf{u}_2 , will have the same objective function value in $\mathcal{F}\mathcal{L}(2, 3, \dots, k; n)$. Additionally, any rotation of the two vectors within the plane $\text{span}(\mathbf{u}_1, \mathbf{u}_2)$ will still produce the same objective function values. So, $\mathcal{F}\mathcal{L}(2, 3, \dots, k; n)$ is no longer a convex optimization problem over $St(k, n)$ because the first two principal components are not unique. However, if we remove $[\mathbf{u}_1] \subset [\mathbf{u}_1, \mathbf{u}_2]$ from the nested subspace structure and consider $[\mathbf{U}] \in \mathcal{F}\mathcal{L}(2, 3, \dots, k; n)$ as

$$[\mathbf{U}] = [\mathbf{u}_1, \mathbf{u}_2] \subset [\mathbf{u}_1, \mathbf{u}_2, \mathbf{u}_3] \subset \dots \subset [\mathbf{u}_1, \dots, \mathbf{u}_k] \subset \mathbb{R}^n. \quad (73)$$

Then we have a unique solution to $\mathcal{F}\mathcal{L}(2, 3, \dots, k; n)$ in place of $St(k, n)$. In practice, it is unlikely that we will have two eigenvectors with the same eigenvalue. However, we can consider two eigenvalues the same as long as $|\lambda_i - \lambda_j| < \epsilon$ for some $\epsilon > 0$.

Motivated by this example, we state a generalization of PCA, which optimizes over flags of a given type.

Definition 1 (Flagified PCA (fPCA) [4]). *A flag of principal components is the solution to:*

$$\arg \max_{[\mathbf{U}] \in \mathcal{F}\mathcal{L}(n+1)} \mathbb{E} \left[\sum_{i=1}^k \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2^2 \right] \quad (74)$$

Ye *et al.* find a solution Eq. (74) using Newton’s method on the flag manifold [6] and Nguyen offers a method for solving such a problem using RTR on flag manifolds [3]. These algorithms produce the same basis vectors for flags regardless of flag type. These basis vectors are different than those found using standard PCA. But, for $[\mathbf{U}] \in \mathcal{F}\mathcal{L}(n_1, n_2, \dots, n_k, n)$ that solves Eq. (74) using either Newton’s method or RTR, the column space of $\mathbf{U}_{:,n_k}$ is the same as the span of the first n_k principal components. This is because the objective function in Eq. (74) is invariant to ordering the columns of \mathbf{U} .

Variants on flagified PCA that maximize $\text{tr}(\mathbf{U}^T \mathbf{X} \mathbf{X}^T \mathbf{U})^q$ over $\mathcal{F}\mathcal{L}(n+1)$ are coined “nonlinear eigenflags” and are difficult to solve for $q = 2$ [6]. Yet, methods from Mankovich *et al.* can be adapted to solve such problems, especially for $q = 1/2$. Another variant of fPCA is weighted fPCA where we assume a weight for each subspace dimension in the flag i and each data point j as $w_{ij} \in \mathbb{R}$. We propose this formulation in the manuscript.

C. DPCP-IRLS and the Grassmannian

This concept was first unearthed in [2]. Expanding the matrix norm we have

$$\|\mathbf{X}^T \mathbf{B}\|_{1,2} = \sum_{j=1}^p \|\mathbf{B}^T \mathbf{x}_j\|_2, \quad (75)$$

$$= \sum_{j=1}^p \sqrt{\sum_{i=1}^k |\mathbf{b}_i^T \mathbf{x}_j|^2}, \quad (76)$$

$$= \sum_{j=1}^p \sqrt{\mathbf{x}_j^T \mathbf{B} \mathbf{B}^T \mathbf{x}_j}, \quad (77)$$

$$= \sum_{j=1}^p \sqrt{\text{tr}(\mathbf{B}^T \mathbf{x}_j \mathbf{x}_j^T \mathbf{B})}. \quad (78)$$

This can be phrased using principal angles as

$$\arg \min_{\mathbf{B}^T \mathbf{B} = \mathbf{I}} \sum_{j=1}^p \cos \theta([\mathbf{x}_j], [\mathbf{B}]). \quad (79)$$

Suppose $\{[\mathbf{X}_j]\}_{j=1}^p \subset \text{Gr}(k, n)$. Namely, $\mathbf{X}_j \in \mathbb{R}^{n \times k}$ where $\mathbf{X}_j^T \mathbf{X}_j = \mathbf{I}$ for each j . A natural generalization of DPCP-IRLS is the optimization problem on the Grassmannian,

$$\arg \min_{[\mathbf{B}] \in \text{Gr}(k, n)} \sum_{j=1}^p \|\cos \theta([\mathbf{X}_j], [\mathbf{B}])\|_2. \quad (80)$$

This can also be solved by an IRLS scheme.

The “flagified” version of Eq. (80) is

$$\arg \min_{[\mathbf{B}] \in \mathcal{F}\mathcal{L}(n+1)} \sum_{j=1}^p \|\cos \theta([\mathbf{X}_j], [\mathbf{B}])\|_2. \quad (81)$$

D. Novel Flagified Robust and Dual PCA and TPCA Variants

We present the intuition behind the geometry of Robust and Dual PCA versus TPCA in Fig. 1. Then we provide a visual comparison between Euclidean and manifold variants of RPCA and DPCP in Fig. 2.

Tab. 1 summarizes our novel flagified robust and dual PCA variants and emphasizes that flag types other than $(1, 2, \dots, k; n)$ and $(k; n)$ produce novel principal directions that are “in between” L_1 and L_2 formulations.

Finally Tab. 2 summarizes the naming schemes of all of the algorithms introduced in this paper

E. Rest of the Proposed Algorithms

In the paper, we proposed three new algorithms. We now present these algorithms as well as the objective functions

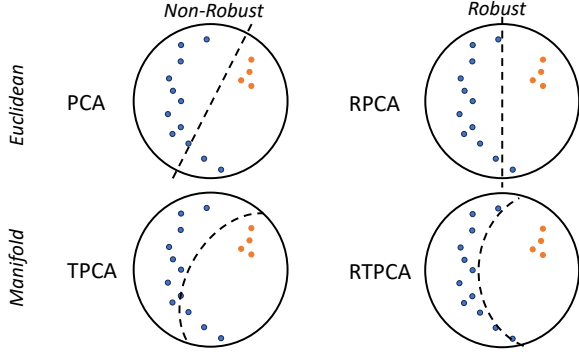


Figure 1. Inliers (blue) and outliers (orange) on the 2-sphere. The first row are Euclidean algorithms and the second row are manifold (tangent space) algorithms. The dashed lines are the first principal subspace (first row) and geodesic (second row) spanned by the first principal direction. Note: first principal subspaces pass through the center of the sphere and first principal geodesics are great circles on the sphere.

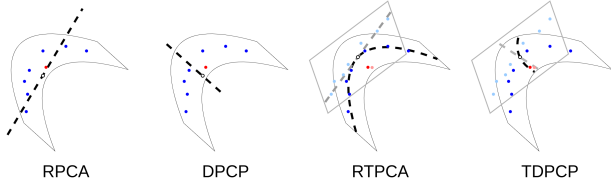


Figure 2. Given manifold valued data with inliers (blue) and outliers (red). The dashed black lines are the 1st principal component for RPCA and DPCP, for RTPCA and TDPCP this is the 1st principal geodesic. For RPCA and RTPCA this line / geodesic should contain the inliers. Due to the reversal in the objective, for DPCP and TDPCP this geodesic should contain the outliers.

they minimize. First, Alg. 1 finds a solution to weighted flagged PCA

$$\llbracket \mathbf{U} \rrbracket^* = \arg \max_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k w_{ij} \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2^2 \right]. \quad (82)$$

Second, Alg. 2 finds a solution to weighted flagged orthogonal PCA ($f\perp$ PCA)

$$\llbracket \mathbf{U} \rrbracket^* = \arg \min_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k w_{ij} \|\pi_{\mathbf{U}_i}(\mathbf{x}_j)\|_2^2 \right]. \quad (83)$$

Flagified (Dual-)PCA	Robust PCA Variant
$f\text{RPCA}(1, \dots, k)$	L_1 -RPCA
$f\text{RPCA}(\cdot)$	–
$f\text{RPCA}(k)$	L_2 -RPCA
$f\text{WPCA}(1, \dots, k)$	L_1 -WPCA
$f\text{WPCA}(\cdot)$	–
$f\text{WPCA}(k)$	L_2 -WPCA
$f\text{DPCP}(1, \dots, k)$	L_1 -DPCP
$f\text{DPCP}(\cdot)$	–
$f\text{DPCP}(k)$	L_2 -DPCP
$f\text{RTPCA}(1, \dots, k)$	L_1 -RTPCA
$f\text{RTPCA}(\cdot)$	–
$f\text{RTPCA}(k)$	L_2 -RTPCA
$f\text{WTPCA}(1, \dots, k)$	L_1 -WTPCA
$f\text{WTPCA}(\cdot)$	–
$f\text{WTPCA}(k)$	L_2 -WTPCA
$f\text{TDPCP}(1, \dots, k)$	L_1 -TDPCP
$f\text{TDPCP}(\cdot)$	–
$f\text{TDPCP}(k)$	L_2 -TDPCP

Table 1. Flag types for Euclidean optimization (first half) and manifold optimization (second half). Flag optimization in these algorithms provides a new objective functions which live in between L_1 and L_2 robust PCA formulations. Note: we remove the number of the ambient dimension in the flag signature for less redundant notation and we assume we are computing the first k principal components.

Lastly, Alg. 3 approximates solutions to

$$\llbracket \mathbf{U} \rrbracket^* \approx \begin{cases} \arg \max_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k w_{ij} d(\boldsymbol{\mu}, \pi_{\mathbf{U}_i}(\mathbf{x}_j))^2 \right], & (f\text{TPCA}) \\ \arg \min_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k w_{ij} d(\boldsymbol{\mu}, \pi_{\mathbf{U}_i}(\mathbf{x}_j))^2 \right], & (f\perp\text{TPCA}) \\ \arg \max_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k d(\boldsymbol{\mu}, \pi_{\mathbf{U}_i}(\mathbf{x}_j)) \right], & (f\text{RTPCA}) \\ \arg \min_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k d(\mathbf{x}_j, \pi_{\mathbf{U}_i}(\mathbf{x}_j)) \right], & (f\text{WTPCA}) \\ \arg \min_{\llbracket \mathbf{U} \rrbracket \in \mathcal{FL}(n+1)} \mathbb{E}_j \left[\sum_{i=1}^k d(\boldsymbol{\mu}, \pi_{\mathbf{U}_i}(\mathbf{x}_j)) \right], & (f\text{TDPCP}) \end{cases} \quad (84)$$

F. Extra Experiments

Impact of flag-type on cluster detection. To assess the impact of flag-type, we generate a dataset $\{\mathbf{x}_j\}_{j=1}^{300} \subset \mathbb{R}^{10}$ with 3 clusters (C_1, C_2, C_3) in which we curate the flag type corresponding to the data structure: $\mathcal{FL}(2, 5, 7; 10)$. To do this we sample $\{\mathbf{x}_j\}_{j=1}^{300} \subset \mathbb{R}^{10}$ with 3 clusters. The l th

Abbreviation	Name
PCA	Principal Component Analysis
RPCA	Robust PCA
WPCA	Weiszfeld PCA
DPCP	Dual Principal Component Pursuit
WDPCP	Weiszfeld DPCP
fPCA	Flagified PCA
fRPCA	Flagified RPCA
fWPCA	Flagified WPCA
fDPCP	Flagified DPCP
fWDPCP	Flagified WDPCP
\mathcal{T} PCA	Tangent PCA
R \mathcal{T} PCA	Robust \mathcal{T} PCA
W \mathcal{T} PCA	Weiszfeld \mathcal{T} PCA
\mathcal{T} DPCP	Tangent DPCP
W \mathcal{T} DPCP	Tangent WDPCP
f \mathcal{T} PCA	Flagified \mathcal{T} PCA
fR \mathcal{T} PCA	Flagified R \mathcal{T} PCA
fW \mathcal{T} PCA	Flagified W \mathcal{T} PCA
f \mathcal{T} DPCP	Flagified \mathcal{T} DPCP
fW \mathcal{T} DPCP	Flagified W \mathcal{T} DPCP

Table 2. The names of the major algorithms covered in this work.

Algorithm 1: Weighted fPCA

Inputs: Dataset $\{\mathbf{x}_j \in \mathbb{R}^n\}_{j=1}^p$,

weights $\{w_{ij}\}_{i,j=1}^{i=k,j=p} \subset \mathbb{R}$,

flag type $(n+1)$

Output: Weighted flagified principal directions

$[\mathbf{U}]^* \in \mathcal{FL}(n+1)$

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

$(\mathbf{W}_i)_{jl} \leftarrow \begin{cases} w_{ij}, & j = l \\ 0, & \text{elsewhere} \end{cases}$

$\mathbf{U}^* \leftarrow$ Solve ?? with $\{\mathbf{W}_i\}_{i=1}^k$ via Stiefel-CGD.

$[\mathbf{U}]^* \leftarrow [\mathbf{U}^*]$

entry of \mathbf{x} , $(\mathbf{x})_l \in \mathbb{R}$, is sampled from

$$C(1) : (\mathbf{x})_l \sim \begin{cases} \mathcal{U}[0, 1), & l \leq 2 \\ \mathcal{U}[0, 0.1), & l \geq 3 \end{cases}, \quad (85)$$

$$C(2) : (\mathbf{x})_l \sim \begin{cases} \mathcal{U}[0, 1), & 3 \leq i \leq 5 \\ \mathcal{U}[0, 0.1), & i \leq 2 \text{ or } i \geq 6 \end{cases}, \quad (86)$$

$$C(3) : (\mathbf{x})_l \sim \begin{cases} \mathcal{U}[0, 1), & i = 6, 7 \\ \mathcal{U}[0, 0.1), & i \leq 5 \text{ or } i \geq 8 \end{cases}. \quad (87)$$

We then compute 2 sets of $k = 7$ principal directions by running fWPCA with flag type $(2, 5, 7; 10)$

Algorithm 2: Weighted flag \perp PCA (f \perp PCA)

Inputs: Dataset $\{\mathbf{x}_j \in \mathbb{R}^n\}_{j=1}^p$,

weights $\{w_{ij}\}_{i,j=1}^{i=k,j=p} \subset \mathbb{R}$,

flag type $(n+1)$

Output: Weighted flagified principal directions

$[\mathbf{U}]^* \in \mathcal{FL}(n+1)$

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

$(\mathbf{W}_i)_{jl} \leftarrow \begin{cases} w_{ij}, & j = l \\ 0, & \text{elsewhere} \end{cases}$

$\mathbf{U}^* \leftarrow$ Minimize the objective in ?? with $\{\mathbf{W}_i\}_{i=1}^k$ via Stiefel-CGD.

$[\mathbf{U}]^* \leftarrow [\mathbf{U}^*]$

Algorithm 3: f \mathcal{T} PCA/fR \mathcal{T} PCA/fW \mathcal{T} PCA/f \mathcal{T} DPCP

Input: Dataset: $\{\mathbf{x}_j\}_{j=1}^p \subset \mathcal{M}$, flag type $(n+1)$,

fPCA Variant: $\Phi : \mathcal{W} \rightarrow \mathcal{FL}(n+1)$

Output: Flagified principal tangent directions $[\mathbf{U}]^*$

if robust then

$\mu \leftarrow$ KarcherMedian $(\{\mathbf{x}_j\}_{j=1}^p)$

else

$\mu \leftarrow$ KarcherMean $(\{\mathbf{x}_j\}_{j=1}^p)$

$\{\mathbf{v}_j\}_j \leftarrow \{\text{Exp}_\mu(\mathbf{x}_j)\}_j$

$\mathcal{W} \leftarrow \{\text{vec}(\mathbf{v}_j)\}_j$

$[\mathbf{U}]^* \leftarrow \Phi(\mathcal{W}, n+1)$

	Cluster 1	Cluster 2	Cluster 3
fWPCA(\cdot)	(7) (2, 5, 7)	(7) (2, 5, 7)	(7) (2, 5, 7)
AUC \uparrow	0.72 0.73	0.48 1.00	0.43 0.49

Table 3. AUC for cluster classification using fWPCA. We see higher AUCs when we match the flag type for fWPCA with the cluster dimensions (e.g., $(2, 5, 7)$).

(fWPCA(2, 5, 7)) and fWPCA(k) using ?? with 200 max. iters. Both of these methods result in a flag representative $\mathbf{U} = [\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3] \in \mathbb{R}^{10 \times 7}$ where $\mathbf{U}_1 \in \mathbb{R}^{10 \times 2}$, $\mathbf{U}_2 \in \mathbb{R}^{10 \times 3}$, and $\mathbf{U}_3 \in \mathbb{R}^{10 \times 2}$. We compute the reconstruction error for point j against each \mathbf{U}_i as $\sum_{j=1}^p \|\mathbf{x}_j - \mathbf{U}_i \mathbf{U}_i^T \mathbf{x}_j\|_2$. These errors are used for 3 classification tasks, predicting C_i using \mathbf{U}_i for $i = 1, 2, 3$. The corresponding AUC values are in Tab. 3. fWPCA(2, 5, 7) produces higher AUCs because it is optimized over a more optimal flag type, respecting the subspace structure of the data.

Data generation for ‘‘Convergence on 4-sphere’’. We first sample a random center $\mathbf{x} \in \mathbb{S}^4$, and then sample 100 inlier tangent vectors from $\mathcal{U}[0, .01)$. Another 20 outlier tangent vectors \mathbf{v} , have entries $v_1, v_2 \sim \mathcal{U}[0, .01)$ and $v_3, v_4, v_5 \sim \mathcal{U}[0, .1)$. We wrap these vectors to have our

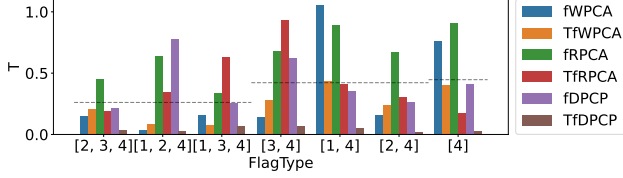


Figure 3. Smaller T corresponds to principal directions which are more similar to those computed with flag type $(1, 2, 3, 4; 5)$. The mean T values for each class of flag type are the horizontal dashed lines. Notice that, these mean values increase as we increase the distance between flag types. We truncate flag types by removing the ambient dimension (5) .

dataset, $\{\text{Exp}_{\mathbf{x}}(\mathbf{v})\}$.

Impact of flag type on principal directions. We run flagified robust PCA and \mathcal{T} PCA variants using ?? (with 200 max. iters.) with different flag types on data on $Gr(2, 4)$ with 100 inliers and 20 outliers sampled as described in the “Outlier detection on $Gr(2, 4)$ ” section of the manuscript. We call $(1, 2, 3, 4; 5)$ the “base” flag type. We use T to measure the different between principal directions from the base flag type $\{\mathbf{u}_1 \dots, \mathbf{u}_4\}$ and other principal directions $\{\mathbf{v}_1 \dots, \mathbf{v}_4\}$ as

$$T = \frac{1}{4} \sum_{i=1}^4 \theta(\mathbf{u}_i, \mathbf{v}_i)^2. \quad (88)$$

We plot T values for different flagified robust PCA and \mathcal{T} PCA variants in Fig. 3. We separate flag types into classes based on the number of nested subspaces. Flag types with the same number of nested subspaces are considered “closer” flags. We find that closer flag types have smaller T values. This experiment verifies that running flagified robust PCA variants with different flag types recover different principal directions and these differences are directly proportional to the “distance” between flag types. This also emphasizes that flag types other than $(1, \dots, k; n)$ and $(k; n)$ indeed recover novel principal directions. The direct utility of these gap-filling methods to real-world datasets is future work.

Outlier detection on $Gr(2, 4)$. We present the result of using PCA, fWPCA $(1, \dots, k)$, fWPCA (k) , fRPCA $(1, \dots, k)$, fRPCA (k) , fDPCP $(1, \dots, k)$, and fDPCP (k) on $Gr(2, 4)$ data for outlier detection in Fig. 4. This is the same data as the data used for ??; but, in this case, we run our algorithms on the vectorized matrix representatives for points on $Gr(2, 4)$ and do outlier detection using Euclidean distance and variances.

Hand reconstructions. We use the 2D Hands dataset and add “hairball” outliers by sampling from a normal distribution with mean 0 and standard deviation 10 ($\mathcal{N}(0, 10)$), then we divide by the Frobenius norm and mean center to obtain

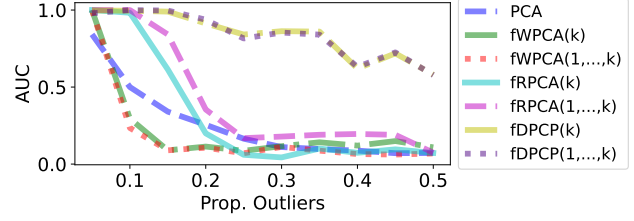


Figure 4. AUC of different algorithms for outlier detection using the first $k = 2$ principal directions of outlier-contaminated data on $Gr(2, 4)$. All algorithms other than PCA are optimized with ?? with 100 max. iters.

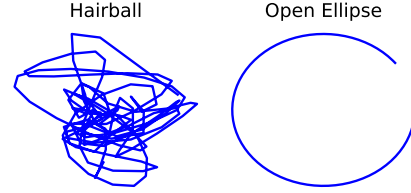


Figure 5. Examples of outliers used for contamination of the hands dataset. Hairballs are used in hand reconstruction and open ellipses are used in outlier detection.

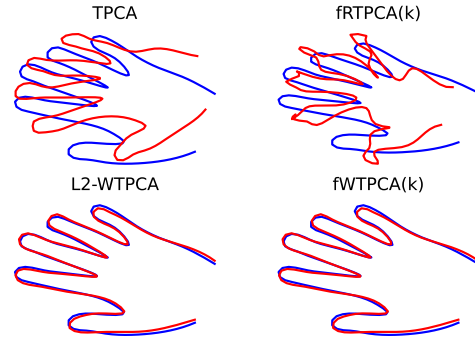


Figure 6. Reconstruction of hand 6 using the first principal direction computed on a dataset with 40 hands and 5 outliers. The cumulative reconstruction errors for the 40 inlier hands from L to R, Top to Bottom, are: 8.19, 6.20, 5.35, and 5.35.

a point on Σ_2^{56} . A figure with an example of an outlier ellipse and a hairball outlier is in Fig. 5.

We run fWTPCA $(1, \dots, k)$, L_1 -WTPCA using Alg. 1 from [5] run on the tangent space, fRTPCA $(1, \dots, k)$, and \mathcal{T} PCA to find different versions of the first $k = 1$ principal direction on a dataset with all 40 hands and 5 outliers. We compute reconstruction error for each method using the framework described in the $Gr(2, 4)$ experiments. Our cumulative reconstruction errors for the 40 inlier hands and a visualization of a hand reconstruction is in Fig. 6. L_1 -WTPCA and fWTPCA $(1, \dots, k)$ produce the lowest reconstruction errors on the hands and have the most sensible reconstructions. Additionally, Alg. 3 preforms just as well as Alg. 1 from [5] run on the tangent space.

We move on to computing cumulative inlier reconstruction errors as we gradually add outliers and report results

in Fig. 7. $fWTPCA$ have the most stable reconstruction errors followed by $fRTPCA$, then $TPCA$.

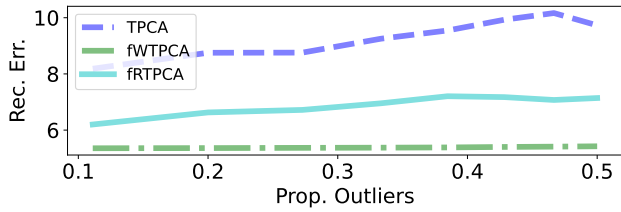


Figure 7. The cumulative reconstruction error of the 40 inlier hands using the first $k = 1$ principal direction where we gradually add hairball outliers.

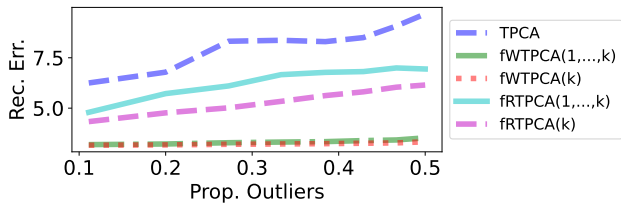


Figure 8. The cumulative reconstruction error of the 40 inlier hands using the first $k = 2$ principal directions where we gradually add hairball outliers.

References

- [1] Nathan Mankovich and Tolga Birdal. Chordal averaging on flag manifolds and its applications. In *ICCV*, pages 3881–3890, 2023. 4
- [2] Nathan J Mankovich. *Subspace and Network Averaging for Computer Vision and Bioinformatics*. PhD thesis, Colorado State University, 2023. 5
- [3] Du Nguyen. Closed-form geodesics and optimization for Riemannian logarithms of Stiefel and flag manifolds. *Journal of Optimization Theory and Applications*, 194(1), 2022. 5
- [4] Xavier Pennec. Barycentric subspace analysis on manifolds. *Annals of Statistics*, 46(6A), 2018. 5
- [5] Qianqian Wang, Quanxue Gao, Xinbo Gao, and Feiping Nie. $\ell_{2,p}$ -norm based PCA for image recognition. *IEEE Transactions on Image Processing*, 27(3):1336–1346, 2017. 8
- [6] Ke Ye, Ken Sze-Wai Wong, and Lek-Heng Lim. Optimization on flag manifolds. *Mathematical Programming*, 194(1):621–660, 2022. 5