

A Simple Prior-free Method for Non-Rigid Structure-from-Motion Factorization

Yuchao Dai¹, Hongdong Li², Mingyi He¹

Northwestern Polytechnical University, China¹
Australian National University, Australia²

Abstract

This paper proposes a simple “prior-free” method for solving non-rigid structure-from-motion factorization problems. Other than using the basic low-rank condition, our method does not assume any extra prior knowledge about the nonrigid scene or about the camera motions. Yet, it runs reliably, produces optimal result, and does not suffer from the inherent basis-ambiguity issue which plagued many conventional nonrigid factorization techniques.

Our method is easy to implement, which involves solving no more than an SDP (semi-definite programming) of small and fixed size, a linear Least-Squares or trace-norm minimization. Extensive experiments have demonstrated that it outperforms most of the existing linear methods of nonrigid factorization. This paper offers not only new theoretical insight, but also a practical, everyday solution, to non-rigid structure-from-motion.

1. Introduction

This paper revisits the classical geometric computer vision problem of non-rigid structure-from-motion (NRSFM). We focus on the *factorization* framework for NRSFM, originally proposed by Bregler *et al.* in [7], as an important extension to the well-known Tomasi-Kanade factorization from rigid scene to nonrigid scene, assuming that the nonrigid shape deformation follows a low-order linear combination model. To date, a large body of researches has been devoted to this topic, and numerous different methods/algorithms have been proposed. However, despite all the efforts, this problem remains a difficult and still active research topic (see *e.g.*, [5, 11, 25, 21, 3, 14]).

One of the primary causes to such difficulty is due to the inherent *basis ambiguity* of the nonrigid problem ([26]). To overcome this, most existing work rely on introducing various *prior* knowledge to the problem at hand. They do so by assuming various constraints about the nonrigid scene, about the nonrigid shape bases, about the coefficients, about the deformation, about the shape itself, or about the camera motion *etc.* For instance, many methods require that the

camera moves smoothly, or the deformation trajectory is slow and smooth. However, these additional constraints not only limit the practical applicability of the methods, but also obscure a clear theoretical understanding to the problem. We would like to answer: in order to solve the nonrigid factorization effectively, are these additional priors *essential*?

In this paper, we propose a novel and simple solution to non-rigid factorization. Our method does not assume any extra prior knowledge about the problem other than the low-rank constraint, hence it is “prior-free”. Nevertheless, it does not suffer from the basis ambiguity difficulty, but is able to recover both camera motion and non-rigid shape accurately and reliably. Experiments on both synthetic and real benchmark datasets show that: the proposed method, being prior-free, outperforms most other (often prior-based) linear factorization methods.

To better present this paper, and also to put our contributions in context, we briefly review some recent progress in non-rigid factorization.

1.1. Related work

Ever since Bregler’s seminal work [7] in 2000, researchers have been actively applying the factorization framework to various nonrigid problems. However, they soon noticed that, different from its rigid counterpart, the non-rigid factorization appeared to be much more difficult. In a 2004 paper, Xiao *et al.* proved that the problem itself is indeed ill-posed or under-constrained, in the sense that, based on the orthonormality constraint alone, one cannot recover the non-rigid shape bases and the corresponding shape coefficients *uniquely* [26]. There is always a fundamental ambiguity between the shape bases and the shape coefficients.

To resolve this ambiguity, Xiao *et al.* suggested to add extraneous “basis constraints” so as to make the system well-constrained. In the same spirit of adding extra priors to regularize an otherwise under-constrained problem, Torresani *et al.* [25] introduced Gaussian prior on the shape coefficients. Del Bue introduced special shape priors [11]. Akhter *et al.* proposed to use a fixed set of DCT bases in

the trajectory dual space [3]. Temporally smooth deformation prior has also been used such as in [5][1]. Gatardo *et al.* assumed the time-trajectory of a single point is smooth [14]. Other priors imposed speciality on the model such as assuming a quadratic model [13], local-rigidity [24], or extend to non-linear models [22].

Akhter *et al.* made an important theoretical progress in [2] which reveals that: although the ambiguity in shape basis is inherent, the 3D shape itself can be recovered uniquely without ambiguity. In a slightly earlier paper, Hartley and Vidal proved a similar result but under perspective camera model [16]. Despite the significance of these theoretical results, neither paper has provided a practical algorithm to the problem, and our work aims to fill this gap.

2. Problem Statement

2.1. Formulation

The task of nonrigid factorization is to factorize an *image measurement matrix* W as the product of *camera motion* (projection) matrix M and a nonrigid shape matrix S , such that $W = MS$. We assume the measurement matrix is already centralized, therefore the camera matrix reduces to pure rotation [7].

Based on the linear combination model, the non-rigid shape $S_i \in \mathbb{R}^{3 \times P}$ can be represented as a linear combination of K *shape bases* $B_k \in \mathbb{R}^{3 \times P}$ with *shape coefficients* c_{ik} as: $S_i = \sum_{k=1}^K c_{ik} B_k$. Under orthographic camera model, the coordinates of the 2D image points observed at frame i are given by: $W_i = R_i S_i$, where $R_i \in \mathbb{R}^{2 \times 3}$ is the first two rows of the i -th camera rotation, hence $R_i R_i^T = I_2$. Use this representation, and stack all the F frames of measurements and all the P points in a matrix form, we obtain:

$$\begin{aligned} W &= \begin{bmatrix} \mathbf{x}_{11} & \cdots & \mathbf{x}_{1F} \\ \vdots & & \vdots \\ \mathbf{x}_{F1} & \cdots & \mathbf{x}_{FP} \end{bmatrix} = \begin{bmatrix} R_1 S_1 \\ \vdots \\ R_F S_F \end{bmatrix} \\ &= \begin{bmatrix} c_{11} R_1 & \cdots & c_{1K} R_1 \\ \vdots & \ddots & \vdots \\ c_{F1} R_F & \cdots & c_{FK} R_F \end{bmatrix} \begin{bmatrix} B_1 \\ \vdots \\ B_K \end{bmatrix} \\ &= R(C \otimes I_3)B \doteq \Pi B. \end{aligned} \quad (1)$$

In this formula, we call $R = \text{blkdiag}(R_1, \dots, R_F) \in \mathbb{R}^{2F \times 3F}$ the camera motion (rotation) matrix. Since $\Pi \in \mathbb{R}^{2F \times 3K}$ and $B \in \mathbb{R}^{3K \times P}$, it is easy to see: $\text{rank}(W) \leq \min(\text{rank}(\Pi), \text{rank}(B)) \leq 3K$. Additionally, the shape matrix $S = (C \otimes I_3)B$ is low rank too, as $\text{rank}(S) \leq \min(\text{rank}(C \otimes I_3), \text{rank}(B)) \leq 3K$.

2.2. Orthonormality Constraint

From a measurement matrix W one can compute its rank- $3K$ decomposition $W = \hat{\Pi} \hat{B}$ via SVD. However, this decomposition is not unique as any nonsingular matrix $G \in$

$\mathbb{R}^{3K \times 3K}$ can be inserted between $\hat{\Pi}$ and \hat{B} to obtain a new valid factorization as $W = \hat{\Pi} \hat{B} = \hat{\Pi} G G^{-1} \hat{B} = \Pi B$.

A particular matrix G that rectifies $\hat{\Pi}$ to be a canonical Euclidean form is called the *Euclidean corrective matrix*, because once such a G is determined, one obtains $R(C \otimes I_3) = \hat{\Pi} G$ and the true shape bases $B = G^{-1} \hat{B}$.

Denote the i -th double rows of $\hat{\Pi}$ as $\hat{\Pi}_{2i-1:2i} \in \mathbb{R}^{2 \times 3K}$, and the k -th column-triplet of G as $G_k \in \mathbb{R}^{3K \times 3}$, we have:

$$\hat{\Pi}_{2i-1:2i} G_k = c_{ik} R_i, \quad i = 1, \dots, F, k = 1, \dots, K. \quad (2)$$

Orthonormality constraints (*i.e.* rotation constraints) in the Π matrix can be imposed to recover a Gram matrix $Q_k \in \mathbb{R}^{3K \times 3K}$ formed by $Q_k = G_k G_k^T$ as $\hat{\Pi}_{2i-1:2i} Q_k \hat{\Pi}_{2i-1:2i}^T = c_{ik}^2 I_2$. Since c_{ik} is not known, one can only establish two linear equations over Q_k as:

$$\hat{\Pi}_{2i-1} Q_k \hat{\Pi}_{2i-1}^T = \hat{\Pi}_{2i} Q_k \hat{\Pi}_{2i}^T, \quad \hat{\Pi}_{2i-1} Q_k \hat{\Pi}_{2i}^T = 0. \quad (3)$$

2.3. Inherent Ambiguity

In doing the above non-rigid factorization, Xiao *et al.* [26] discovered that, the solutions are however fundamentally ambiguous, in the sense that one cannot expect to find the shape bases and shape coefficients uniquely. Such an inherent ambiguity largely explains why nonrigid factorization is fundamentally more difficult than its rigid counterpart. Later, Akhter *et al.* [2] showed that, quite surprisingly, the fundamental ambiguity does not necessarily lead to an ambiguous shape. In addition, they further proved that using the orthonormality constraints alone is in fact sufficient to recover a unique (unambiguous) non-rigid shape (provided that a previously-overlooked rank-3 constraint on Q_k (Eq.-(3)) is accounted for). However, apart from its evident theoretical value, their paper did not propose any optimization algorithm (other than a local search method due to [6]) to efficiently find the correct G_k . Instead, the authors argued that “the real difficulty of in achieving good 3D reconstructions for nonrigid structures...is not the ambiguity of the [basis] constraints, but the complexity of the underlying non-linear optimization”. In this paper we will challenge this argument, by providing a simple yet efficient (optimization) solution to Nonrigid SFM Factorization.

3. Main Theory

From now on, let us assume the measurement matrix W is already truncated to rank $3K$ (by *e.g.* SVD), the number of shape bases K has been estimated, and all the shape bases are non-degenerate. We start with a known result of NRSFM factorization.

Theorem 3.1. *All the solutions of Q_k to linear system Eq.-(3) form a linear subspace of dimensionality $(2K^2 - K)$.*

This result is in fact a direct consequence of Xiao *et al.*'s Theorem in [26]. It shows that the above linear system is inherently under-determined (as by $(2K^2 - K)$ rank deficient), no matter how many image frames are given.

On the other hand, this result also provides us with the true dimensionality of the solution space of Q_k , and note that Q_k is precisely what we are after. However, in their paper, this practical implication had not been explicitly exploited.

In the following, we will show how one can take advantage of this result, and derive a practical algorithm that directly leads to a parametrization of this solution space. More precisely, we will prove that the solution space of Q_k is actually the *null-space* of a certain matrix A which can be directly obtained from the input image data. Practical usefulness of this representation is obvious.

3.1. Null-space Representation

First, denote $\text{vec}()$ as the vectorization operator, and $\mathbf{q}_k = \text{vec}(Q_k)$. Using $\text{vec}(AXB^T) = (B \otimes A)\text{vec}(X)$, we rewrite the linear system Eq.-(3) as:

$$\begin{bmatrix} (\hat{\Pi}_i \otimes \hat{\Pi}_i)(1, :) - (\hat{\Pi}_i \otimes \hat{\Pi}_i)(4, :) \\ (\hat{\Pi}_i \otimes \hat{\Pi}_i)(2, :) \end{bmatrix} \mathbf{q}_k \doteq A_i \mathbf{q}_k = \mathbf{0}, \quad (4)$$

where $(\hat{\Pi}_i \otimes \hat{\Pi}_i)(j, :)$ denotes the j -th row of $(\hat{\Pi}_i \otimes \hat{\Pi}_i)$.

Stacking all such equations from all frames ($i = 1, \dots, F$), we then have

$$\text{Avec}(Q_k) = A \mathbf{q}_k = \mathbf{0}, \quad (5)$$

where $A = [A_1^T, A_2^T, \dots, A_F^T]^T$. This is a linear system of equations over the unknown $9K^2$ -vector \mathbf{q}_k .

Note that the $9K^2$ -vector \mathbf{q}_k has $(3K)(3K+1)/2$ independent entries. It may appear that, given enough frames, *i.e.* when $2F \geq (3K)(3K+1)/2$, \mathbf{q}_k should be able to be solved via linear least squares. However, this is not the case, because all valid solutions reside in a $2K^2 - K$ dimensional space as shown in Theorem 3.1. Moreover, Eq.-(5) shows, the solution space is nothing but the null-space of A . In addition, it is easy to verify that the minimum required number of frames for computing the null-space linearly is $F \geq (5K^2 + 5K)/4$.

3.2. The Intersection Theorem

Combining all the proceeding results, we now arrive at the central theorem of this paper:

Theorem 3.2 (Intersection Theorem). *Under non-degenerate and noise-free conditions, any correct solution of Q_k (i.e. the Gram matrix of a column-triplet of the true Euclidean corrective matrix G_k) must lie in the intersection of the $(2K^2 - K)$ -dimensional null-space of A and a rank-3 positive semi-definite matrix cone, i.e., Q_k belongs to*

$$\{A \text{ vec}(Q_k) = \mathbf{0}\} \cap \{Q_k \succeq 0\} \cap \{\text{rank}(Q_k) = 3\}. \quad (6)$$

Proof. Denote G_k as column triplet of a correct rectifying transform, and $Q_k = G_k G_k^T$, then $\text{rank}(Q_k) = \text{rank}(G_k) = 3$, $Q_k \succeq 0$. Additionally, $\text{vec}(Q_k)$ lies in the null space of A as $\text{vec}(Q_k)$ gives correct rectification with zero error. Thus G_k is a solution to the equation system which means that the equation system is well-defined. Denote \tilde{G}_k as solution to the equation system Eq.-(6), then $\text{rank}(\tilde{Q}_k) = 3$, $\tilde{Q}_k \succeq 0$ and $\text{vec}(\tilde{Q}_k)$ lies in the null space of the system of linear constraints on the elements \tilde{Q}_k . Thus all the solutions to the equation system Eq.-(6) satisfy the condition for correct rectifying transforms. There is no difference between these solutions. \square

4. Algorithm Solution

Armed with the above results (in particular Theorem 3.2), we are now ready to present our simple algorithm to the non-rigid factorization problem.

Recall that the goal is to recover the true motion matrix R and the true non-rigid shape matrix S from image measurement W , such that $W = RS = R(C \otimes I_3)B$. Note that due to the inherent basis ambiguity, it is hopeless to recover a unique B or C . While in previous work many researchers chose to use a pre-selected special shape bases B (or enforce arbitrary priors on the shape bases or shape coefficients) to pin down the undetermined degrees-of-freedom, in this work we will show how one can directly estimate the S without fixing B or C .

Our algorithm consists of three steps to be applied in sequel: (1) Estimate the (Gram of) corrective matrix G_k , (2) Estimate camera rotations R and (3) Estimate the nonrigid shape S . We now explain the three steps in order.

4.1. Step-1: Estimate G_k by Trace-Minimization

Our main intersection theorem (Theorem-3.2) naturally leads to an easy algorithm to solve for G_k , that is: to find the intersection of the aforementioned null-space and a rank-3 positive semi-definite matrix cone.

Because the rank-function itself is not very numerically-stably, measurements noise will increase the *numerical rank* of Q_k dramatically, we slightly relax the $\text{rank}(Q_k) = 3$ condition to a *rank-minimization* problem, *i.e.* $\min \text{rank}(Q_k)$. Note that however, rank-minimization is an NP-hard problem in general, and is very difficult to solve exactly. We therefore further relax it to a nuclear-norm minimization form, *i.e.*, $\min \|Q_k\|_*$. Moreover in our case, since Q_k is a symmetric positive definite matrix, the nuclear norm is simply its trace [23][4]. Thus we have $\|Q_k\|_* = \text{trace}(Q_k)$.

Then we arrive at the following trace-minimization to solve for the corrective matrix.

$$\begin{aligned} & \min \text{trace}(Q_k), \text{ such that,} \\ & Q_k \succeq 0, \\ & A \text{ vec}(Q_k) = \mathbf{0}. \end{aligned} \quad (7)$$

To avoid a trivial solution at $\mathbf{Q}_k = \mathbf{0}$, we express $\text{vec}(\mathbf{Q}_k)$ in explicit form of the null-space representation.

Easy to see that the above trace-minimization problem is a standard semi-definite programming (SDP). Also note that this SDP is actually of small and fixed size (of $2K^2 - K$) which is independent of the size of the measurement matrix. Thus this SDP can be solved very easily and very efficiently by any off-the-shelf SDP solvers. Once \mathbf{Q}_k is found, we use SVD to extract \mathbf{G}_k . This solved \mathbf{G}_k can be directly used to find \mathbf{R} and then \mathbf{S} . Alternatively, if higher accuracy is desired, one can further improve the numerical accuracy of \mathbf{G}_k by feeding it as an initial point to a non-linear refinement procedure, such as via the following unconstrained minimization:

$$\min_{\mathbf{G}_k} \sum_{i=1}^F \left[\left(1 - \frac{\hat{\Pi}_{2i} \mathbf{G}_k \mathbf{G}_k^T \hat{\Pi}_{2i}^T}{\hat{\Pi}_{2i-1} \mathbf{G}_k \mathbf{G}_k^T \hat{\Pi}_{2i-1}^T} \right)^2 + \left(2 \frac{\hat{\Pi}_{2i-1} \mathbf{G}_k \mathbf{G}_k^T \hat{\Pi}_{2i}^T}{\hat{\Pi}_{2i-1} \mathbf{G}_k \mathbf{G}_k^T \hat{\Pi}_{2i-1}^T} \right)^2 \right],$$

where the objective function is nothing but the orthonormality condition and this refinement is similar to a bundle-adjustment process.

4.2. Step-2: Compute Rotation Matrix \mathbf{R}

Conventionally, once $\mathbf{G}_k \in \mathbb{R}^{3K \times 3}$ is solved (w.l.o.g., let's denote it as \mathbf{G}_1), which is merely a single column-triplet in the full corrective matrix $\mathbf{G} \in \mathbb{R}^{3K \times 3K}$, the commonly-used next step is to solve for the other $K - 1$ of independent column-triplets $[\mathbf{G}_2, \dots, \mathbf{G}_K]$, and use them to populate the entire matrix \mathbf{G} . Brand [6] proposed a linear method to solve for the big \mathbf{G} (for affine case). Because these \mathbf{G}_k s always have rotation ambiguity, in order to align them, Procrustes method must be employed subsequently (c.f. [26][2]). Once the big \mathbf{G} is obtained, one then is allowed to compute the camera motion \mathbf{R} , the shape coefficients \mathbf{C} and the shape bases \mathbf{B} , and then reconstruct the non-rigid shape \mathbf{S} . However, the above approach is not only rather involved, but also not numerically stable. More importantly, it is *not necessary*, as shown in [3].

In this work, we adopt a simpler approach that directly computes the camera motion \mathbf{R} from a single column-triplet \mathbf{G}_k , without the need to fill in a big and full \mathbf{G} matrix. The method goes as follows. Once \mathbf{G}_k is solved, the rotation at every frame $i = 1, \dots, F$ can be solved by using:

$$\hat{\Pi}_{2i-1:2i} \mathbf{G}_k = c_{ik} \mathbf{R}_i, \quad i = 1, \dots, F. \quad (8)$$

Note that we do not need to care about the unknown value of c_{ik} , though its sign ambiguity must be taken care of (c.f. [2]). Finally, the full motion matrix \mathbf{R} is formed as $\mathbf{R} = \text{blkdiag}([\mathbf{R}_1, \mathbf{R}_2, \dots, \mathbf{R}_F])$.

4.3. Step-3: Estimate \mathbf{S} by Rank-Minimization

Now we show how to solve the non-rigid shape matrix \mathbf{S} . Most conventional methods do this *indirectly*, in the sense that they often start from solving the big \mathbf{G} matrix firstly, and then use pre-selected special shape bases \mathbf{B} (such as the

first K frame [26], or DCT bases in the dual space [3], or assume the the shape coefficients are also DCT-expandable [14]), then the corresponding coefficient matrix \mathbf{C} can be determined, and also the shape matrix $\mathbf{S} = (\mathbf{C} \otimes \mathbf{I}_3) \mathbf{B}$.

In the next two sub-sections, we will provide two simpler, more direct methods for solving \mathbf{S} .

4.3.1 Pseudo Inverse Method

Recall that our goal is to solve \mathbf{S} through the equation of $\mathbf{W} = \mathbf{R}\mathbf{S}$ given \mathbf{W} and \mathbf{R} . This equation is under-determined because \mathbf{R} is a short matrix of size $2F \times 3F$. There should be no unique but an infinite family of solutions to \mathbf{S} . However, we also notice that: the low-order linear model, *i.e.* $\mathbf{S} = (\mathbf{C} \otimes \mathbf{I}_3) \mathbf{B}$ immediately suggests that $\text{rank}(\mathbf{S}) \leq 3K$.

Taking into account of both of the above arguments, we reach: a valid solution to the shape matrix \mathbf{S} must lie in the intersection of low-rank matrix set of $\{\text{rank}(\mathbf{S}) \leq 3K\}$ and the solution space of equation $\mathbf{W} = \mathbf{R}\mathbf{S}$. As usual we relax the low-rank condition to rank-minimization. Now the shape matrix \mathbf{S} must be a solution to the following rank minimization problem:

$$\min \text{rank}(\mathbf{S}), \text{ such that, } \mathbf{W} = \mathbf{R}\mathbf{S}. \quad (9)$$

Remarks. We now make two important remarks: (1) the above rank-minimization problem (in the context of NRSFM) must have a unique solution; (2) this unique solution is nothing but the (unique) Moore-Penrose pseudo-inverse solution, *i.e.* $\mathbf{S} = \mathbf{R}^\dagger \mathbf{W} = (\mathbf{R}^T (\mathbf{R} \mathbf{R}^T)^{-1}) \mathbf{W}$.

Remark-1 is not surprising, as it is simply the main conclusion of [2], which states that: once \mathbf{R} is fixed, there is no ambiguity in finding the shape matrix \mathbf{S} , and the solution is *unique*.

Remark-2 was a bit surprising to the authors, as it seems to suggest a simple, linear, and closed-form solution to our rank-minimization problem of (9)—which otherwise would be NP-hard to solve in general.

Fortunately, recent progress in Compressive Sensing has confirmed the correctness of our Remark-2. In particular, we use the following result due to [19] that: *the Moore-Penrose pseudo-inverse solution $\mathbf{S} = \mathbf{R}^\dagger \mathbf{W}$ is the unique minimizer of the above rank minimization problem (9)*. For a detailed proof the reader is referred to [19]. From this result, we see that the pseudo-inverse solution is indeed a *unique* matrix that satisfies the necessary conditions which a low-rank shape matrix must all satisfy.

We have applied the pseudo-inverse method to both synthetic data and real data. In terms of shape-recovery accuracy, the pseudo-inverse solution outperforms Xiao *et al.*'s K -basis method by a large margin, and achieves comparable performance with Metric-projection [21] and EM-PPCA [25]. The reader is referred to the second last column of

Table-1. This is already very encouraging, as our method does not use any priors except for the low-rank condition. Note that however, our method is slightly inferior to the more recent DCT trajectory basis method [3] or the CSF method [14] (, both methods rely on strong smoothness prior)—which prompts us to think: can we do any better, without using any prior ?

4.3.2 Block Matrix Method

In the above pseudo-inverse method, we mainly make use of the rank- $3K$ condition that is $\text{rank}(\mathbf{S}) \leq 3K$. This $(3F \times P)$ matrix \mathbf{S} is simply a stack of P 3D points $[X_i, Y_i, Z_i]^T$ over F frames.

However, we realize that, since in reality there are *in fact* only K shape bases (rather than $3K$), the shape matrix \mathbf{S} is not a fully-generic rank- $3K$ matrix, but has its special block structure.

In particular, we re-arrange the rows of \mathbf{S} that correspond to X, Y , and Z coordinate separately, in an $F \times 3P$ block matrix form, denoted by $\mathbf{S}^\#$ in below¹:

$$\mathbf{S}^\# = \begin{bmatrix} X_{11} & \dots & X_{1P} & Y_{11} & \dots & Y_{1P} & Z_{11} & \dots & Z_{1P} \\ \vdots & & \vdots & \vdots & & \vdots & \vdots & & \vdots \\ X_{F1} & \dots & X_{FP} & Y_{F1} & \dots & Y_{FP} & Z_{F1} & \dots & Z_{FP} \end{bmatrix}.$$

Then we must have: $\text{rank}(\mathbf{S}^\#) \leq K$. Note that this rank- K condition (on $\mathbf{S}^\#$) is stronger than the above rank- $3K$ condition (on \mathbf{S}), and the former captures the essence of the K -order linear combination model.

Now, to solve for this re-arranged shape matrix $\mathbf{S}^\#$, we use, again, a rank-minimization formulation:

$$\begin{aligned} \min \text{rank}(\mathbf{S}^\#), \text{ such that,} \\ \mathbf{W} = \mathbf{R}\mathbf{S}, \\ \mathbf{S}^\# = [\mathbf{P}_X \ \mathbf{P}_Y \ \mathbf{P}_Z](\mathbf{I}_3 \otimes \mathbf{S}). \end{aligned} \quad (10)$$

The last matrix equality condition is a compact (shorthand) representation of the re-arrangement relationship between $\mathbf{S}^\#$ and \mathbf{S} , where $\mathbf{P}_X, \mathbf{P}_Y, \mathbf{P}_Z \in \mathbb{R}^{F \times 3F}$ are some properly defined 0-1-valued “row-selection” matrices (similar to the “permutation matrix”).

Fast numerical implementation. We relax the above rank-minimization to nuclear-norm (*i.e.* trace-norm) minimization, *i.e.* $\min \|\mathbf{S}^\#\|_*$. In principle, this nuclear-norm minimization may be solved by a standard SDP solver. However, unlike the case of Eq.-7 where the resulted SDP has small and fixed size, here this SDP is of size $F \times 3P$, which renders the SDP technique very inefficient when either P or F is large.

Below, we give an efficient numerical implementation, based on *fixed point continuation* [20]. First, we re-cast the

above minimization Eq.- (10) in Lagrangian form as:

$$\begin{aligned} \min \mu \|\mathbf{S}^\#\|_* + \frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_F^2, \text{ such that,} \\ \mathbf{S}^\# = [\mathbf{P}_X \ \mathbf{P}_Y \ \mathbf{P}_Z](\mathbf{I}_3 \otimes \mathbf{S}), \end{aligned} \quad (11)$$

where μ is the continuation (homotopy) parameter which diminishes as the algorithm iterates. Next, the gradient of $\frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_F^2$ with respect to $\mathbf{S}^\#$ is obtained as:

$$g(\mathbf{S}^\#) = \frac{\partial \frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_F^2}{\partial \mathbf{S}^\#} = [\mathbf{P}_X \ \mathbf{P}_Y \ \mathbf{P}_Z](\mathbf{I}_3 \otimes (\mathbf{R}^T(\mathbf{R}\mathbf{S} - \mathbf{W}))). \quad (12)$$

Then, we solve the minimization of Eq.- (11) via the following two-line iteration update (cf. [20]):

$$\begin{cases} \mathbf{Y}^{(t)} = \mathbf{S}^{\#(t)} - \tau g(\mathbf{S}^{\#(t)}), \\ \mathbf{S}^{\#(t+1)} = \mathcal{S}_{\tau\mu}(\mathbf{Y}^{(t)}), \end{cases} \quad (13)$$

where τ is the step size of gradient descent, and $\mathcal{S}_v(\cdot)$ is the *matrix shrinkage* operator (cf. [20]). Once the iteration converges, we first project the solved $\mathbf{S}^\#$ to the nearest rank- K matrix (note: not $3K$), then rearrange it to \mathbf{S} .

5. Experiments

5.1. Setup

We compare our methods against the state-of-the-art methods, which include (1) Xiao *et al.*’s shape basis method (XCK) [26]; (2) Torresani *et al.*’s EM-PPCA [25]; (3) Metric projection [21]; (4) Trajectory basis method [3]; and (5) Column space fitting (CSF) [14].²

To facilitate the comparison, we use the same error metrics as reported in [3] and [14], that is: e_R measures the mean error in rotation estimation and $e_R = \frac{1}{F} \sum_{i=1}^F \|\mathbf{R}_i - \tilde{\mathbf{R}}_i\|_F$, where \mathbf{R}_i is the ground truth rotation at frame i and $\tilde{\mathbf{R}}_i$ the recovered rotation; e_{3D} measures the normalized mean 3D error in the reconstructed 3D points and $e_{3D} = \frac{1}{\sigma F P} \sum_{i=1}^F \sum_{p=1}^P e_{ip}$, $\sigma = \frac{1}{3F} \sum_{i=1}^F (\sigma_{ix} + \sigma_{iy} + \sigma_{iz})$, where σ_{ix}, σ_{iy} and σ_{iz} are the standard deviations in X, Y and Z coordinates of the original shape at frame i .

Extensive experiments are conducted to test the performance of the proposed methods, on both randomly synthetic data and on real motion capture data. The random synthetic data, which satisfy the low-rank nonrigid model perfectly, are used only for the purpose of algorithm validation, for which our methods have obtained nearly perfect result (with zero error) as we expect; the results are therefore omitted. Instead, only results on real sequences are reported below. The real sequences we have tested include the standard sequences of Drink (1102/41), Pickup (357/41), Yoga (307/41), Stretch (370/41), and Dance (264/75) used in [3], and Face (316/40), Shark (240/91) and Walking (260/55) in [25], where (F/P) denotes the number of frames (F) and points (P).

¹An identical rearrangement was used in [3] but with different motivation and for different purpose.

²We did not use the CSF variant of [15] as they are very similar.

5.2. Cumulative histograms of errors

Our first experiment is aimed to give a statistical comparison between the performance of our method and several existing methods. For this purpose we use a real motion capture sequence, here *e.g.* the Stretch sequence. From the ground-truth 3D point clouds of the sequence as well the true camera matrices, we re-synthesize F frames of image measurements with Gaussian random noise added in, where noise ratio is defined as $\|\text{Noise}\|_{\text{Fro}}/\|\mathbf{W}\|_{\text{Fro}}$. Use the obtained data, we test our methods, as well as several other existing methods. We repeat the random test 100 times. Then, we plot the cumulative histograms of the rotation estimation errors, and the 3D reconstruction errors, as shown in Fig. 1. This figure clearly reveals that: our block-matrix method outperforms most of the other methods, and our pseudo-inverse method also achieves better results compared with EM-PPCA and XCK.

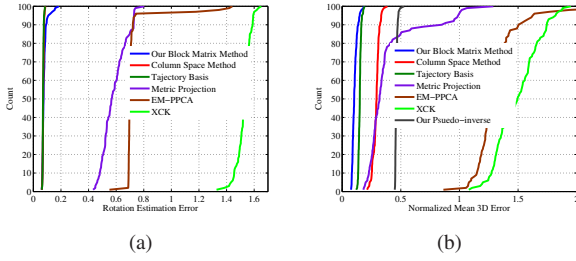


Figure 1. Cumulative histograms of errors tested on the “Stretch” sequence. The most left curve gives the best performance. **Left:** the rotation error; **Right:** the 3D reconstruction error. Our Block-Matrix method outperforms most of the other existing methods.

5.3. Noise performance

To analyze the behavior of our new methods under noise, we repeat the (above) first experiment at different noise ratios. Example results on the Stretch sequence are given in Fig. 2 which plots the estimation errors as a function of the noise ratio.

It is seen, our block matrix method achieves the best performance in terms of the accuracy for rotation estimation and for shape recovery, compared favorably with almost all the other state-of-the-art competitors. Note that our pseudo-inverse method also achieves better performance than EM-PPCA, Metric Projection and XCK.

5.4. Compare all methods on all real sequences

In this subsection, we provide experimental results of all the 5 methods we are benchmarking, on all the real sequences at hand. Table-1 summarizes our main results, where both the shape reconstruction error (mean 3D error) and the camera rotation error are provided (whenever the ground truth are available). Fig.-3 shows the comparison.

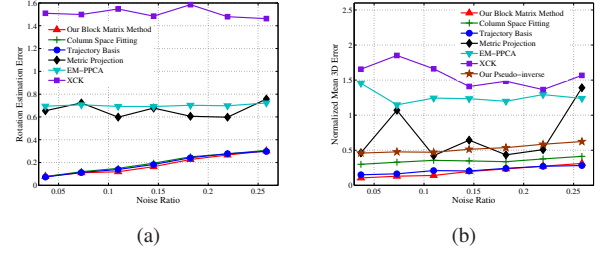


Figure 2. Noise performance. (a) Rotation estimation error. (b) Normalized mean 3D error.

Clearly, our block matrix method achieves the best performance in shape recovery, on almost all of the benchmark sequences (the Shark sequence is an exception, but that possibly due to the fact that this Shark sequence is in fact degenerate [25]).

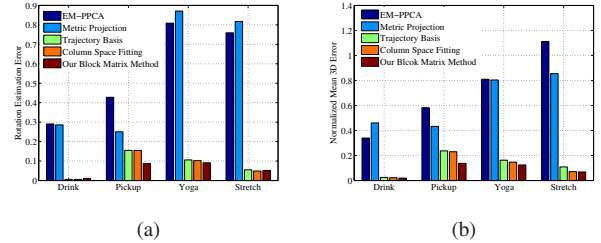


Figure 3. Motion capture data experimental results. **Left:** Rotation estimation error; **Right:** 3D reconstruction error.

5.5. Test on frame-reshuffled data

A highlight of this work is that our method does not assume any prior knowledge about the problem. For instance, we do not assume the trajectories are smooth across frames, while many other method do make this assumption, either explicitly or implicitly. So, we expect (predict) that our method is immune to random frame-order reshuffle (permutation).

To verify this point, we redo the experiments but on frame-reshuffled data, and obtain the following results in Fig.-4. It is seen that both the trajectory basis method and the CSF method perform very badly on the reshuffled sequence, while our method remains unaffected.

Fig.-5 gives a close inspection. The top row compares the trajectories recovered by using the original sequence, and using a frame-reshuffled sequence, while the bottom row shows the normalized mean 3D error for both sequences. From this figure, the trajectory basis method fails to output acceptable results on the frame-permuted sequence, but our method leads to identical results on both cases. This is not surprising as permutating a matrix will not change its rank, and our methods does not assume any frame order or temporal smoothness.

Table 1. Quantitative comparison of our proposed methods versus the state-of-the-art methods on benchmark video sequences. $e_{3D}(P)$ and $e_{3D}(B)$ denote the 3D errors of our Pseudo-inverse method and Block matrix method respectively.

Dataset	XCK		EM-PPCA		Metric Projection		Trajectory Basis		Column Space Fitting		Proposed Methods		
	e_R	e_{3D}	e_R	e_{3D}	e_R	e_{3D}	e_R	$e_{3D}(K)$	e_R	$e_{3D}(K)$	e_R	$e_{3D}(P)$	$e_{3D}(B)$
Drink	0.336	3.519	0.291	0.339	0.286	0.460	0.006	0.025(13)	0.006	0.022(6)	0.011	0.451(4)	0.019(4)
Pick-up	0.469	3.372	0.428	0.582	0.251	0.433	0.155	0.237(12)	0.155	0.230(6)	0.087	0.580(7)	0.138(7)
Yoga	1.201	7.494	0.809	0.810	0.871	0.804	0.106	0.162(11)	0.102	0.147(7)	0.091	0.659(9)	0.125(9)
Stretch	0.949	4.242	0.759	1.111	0.817	0.855	0.055	0.109(12)	0.049	0.071(8)	0.052	0.468(8)	0.069(8)
Dance	-	2.996	-	0.984	-	0.264	-	0.296(5)	-	0.271(2)	-	0.575(10)	0.171(10)
Face	-	-	-	0.033	-	0.036	-	0.044(5)	-	0.036(3)	-	0.485(7)	0.030(7)
Walking	-	-	-	0.492	-	0.561	-	0.395(2)	-	0.186(2)	-	0.471(6)	0.132(6)
Shark	-	-	-	0.050	-	0.157	-	0.180(9)	-	0.008(3)	-	0.902(3)	0.242(3)

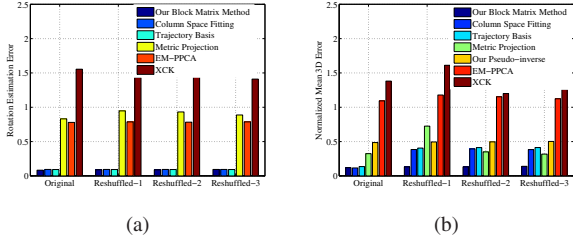


Figure 4. Performance on frame order reshuffled sequences. (a) Rotation error. (b) Normalized mean 3D error. (Better viewed in color).

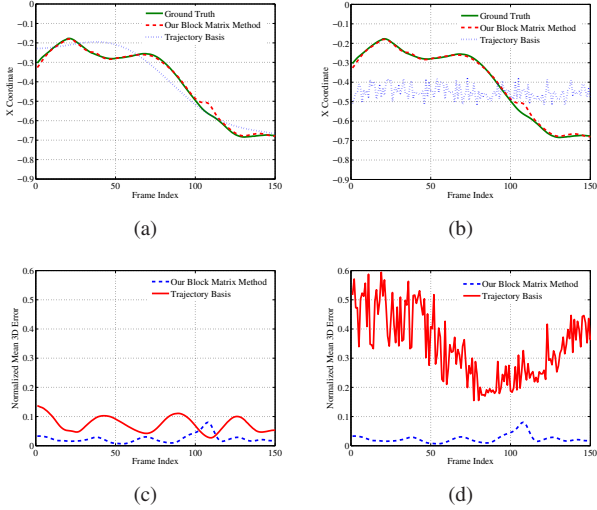


Figure 5. Comparison of our block matrix method versus the trajectory basis method, on an original input video sequence, as well as on a frame-reshuffled version. (a) Recovered trajectory of one point in the X coordinate by both methods on the original sequence. (b) Recovered trajectory of one point in the X coordinate by both methods on the frame-reshuffled sequence. (c) Normalized mean 3D error for both methods on the original sequence. (d) Normalized mean 3D error for both methods on the frame-reshuffled sequence.

5.6. Sample shape reconstruction results

For visual evaluation, we give result comparison between our block matrix method and the trajectory basis method on a more complex sequence, Dance, see Fig.-6.

We also test the Talking Face video³, using 500 frames and 68 feature tracks. Fig.-7 shows 3 frames of the original images and the resulted 3D points, where the reprojection error is 0.9222 pixels.

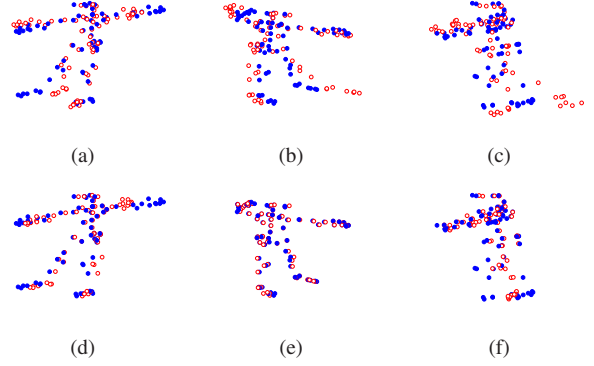


Figure 6. Comparison of the 3D reconstruction results on the Dance sequence. The blue dots are the ground truth 3D points, and the red circles show the reconstructed points. **Top row:** results by the trajectory basis method [3], where the 3D errors are **0.3011, 0.2827, 0.2814** for the 3 frames. **Bottom row:** our result by the block matrix method, where the 3D errors are **0.2228, 0.0355, 0.1389** for the 3 frames.

6. Closing remarks

This paper advocates a novel prior-free approach to non-rigid factorization. Our method is purely convex, very easy to implement, and is guaranteed to converge to an optimal solution (at least approximately up to certain relaxation). It shows that, contrary to common belief, the NRSFM factorization problem can be solved unambiguously, efficiently and accurately, without using extra priors. This said, however, from a practical point of view, we do not against the use of available prior, as long as the prior is sensible and

³<http://www.prima.inrialpes.fr/FGnet/data/01-TalkingFace/>

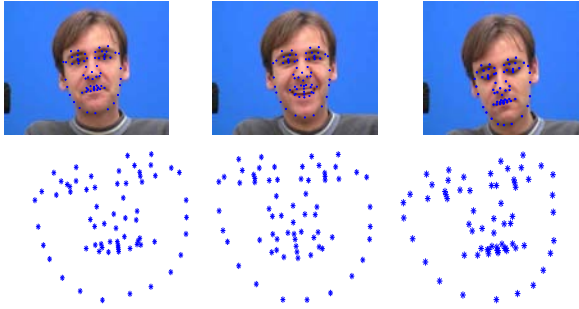


Figure 7. Example 3D deformable shape reconstruction results. **Top row:** sample frames of the input Face sequence. **Bottom row:** recovered 3D face shapes by using our block matrix method.

reflects the physical nature of the problem at hand. It is expected that, using good prior will further improve our solution, and make our method more applicable to complex scenarios.

In the present paper, we have concentrated on complete measurement case under orthographic camera model. Thanks to recent progress in SFM and Compressive Sensing, the proposed method can be easily adapted to handling missing-data case (e.g. [8, 10, 12]), outlier case (e.g. [9, 18, 4]), multibody motion case [17], as well as perspective camera case [10].

Acknowledgement. This work is funded, in part by Natural Science Foundation of China (60736007), and by Australia Research Council. The authors wish to thank Richard (Prof. Hartley) for invaluable discussions.

References

- [1] H. Aanæs and F. Kahl. Estimation of deformable structure and motion. In *Workshop on Vision and Modelling of Dynamic Scenes, ECCV*, pages 1–4, 2002. 2
- [2] I. Akhter, Y. Sheikh, and S. Khan. In defense of orthonormality constraints for nonrigid structure from motion. In *CVPR*, pages 1534–1541, 2009. 2, 4
- [3] I. Akhter, Y. Sheikh, S. Khan, and T. Kanade. Nonrigid structure from motion in trajectory space. In *NIPS*, 2008. 1, 2, 4, 5, 7
- [4] R. Angst, C. Zach, and M. Pollefeys. The generalized trace-norm and its application to structure-from-motion problems. In *ICCV*, pages 1–8, 2011. 3, 8
- [5] A. Bartoli, V. Gay-Bellile, U. Castellani, J. Peyras, S. Olsen, and P. Sayd. Coarse-to-fine low-rank structure-from-motion. In *CVPR*, pages 1–8, 2008. 1, 2
- [6] M. Brand. A direct method for 3D factorization of nonrigid motion observed in 2D. In *CVPR*, 2005. 2, 4
- [7] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3D shape from image streams. In *CVPR*, pages 690–696, 2000. 1, 2
- [8] A. M. Buchanan and A. W. Fitzgibbon. Damped newton algorithms for matrix factorization with missing data. In *CVPR*, pages 316–322, 2005. 8
- [9] E. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis. *J. ACM*, 58(3):11:1–37, 2011. 8
- [10] Y. Dai, H. Li, and M. He. Element-wise factorization for n-view projective reconstruction. In *ECCV*, pages 396–409, 2010. 8
- [11] A. Del Bue. A factorization approach to structure from motion with shape priors. In *CVPR*, pages 1–8, 2008. 1
- [12] A. Eriksson and A. van den Hengel. Efficient computation of robust low-rank matrix approximations in the presence of missing data using the L_1 norm. In *CVPR*, pages 771–778, 2010. 8
- [13] J. Fayad, L. Agapito, and A. Del Bue. Piecewise quadratic reconstruction of non-rigid surfaces from monocular sequences. In *ECCV*, pages 297–310, 2010. 2
- [14] P. Gotardo and A. Martinez. Computing smooth time-trajectories for camera and deformable shape in structure from motion with occlusion. *PAMI*, 33(10):2051–2065, 2011. 1, 2, 4, 5
- [15] P. Gotardo and A. Martinez. Non-rigid structure from motion with complementary rank-3 spaces. In *CVPR*, pages 3065–3072, 2011. 5
- [16] R. Hartley and R. Vidal. Perspective nonrigid shape and motion recovery. In *ECCV*, pages 276–289, 2008. 2
- [17] H. Li. Two-view motion segmentation from linear programming relaxation. In *CVPR*, pages 1–8, june 2007. 8
- [18] H. Li. Consensus set maximization with guaranteed global optimality for robust geometry estimation. In *ICCV*, pages 1074–1080. IEEE, 2009. 8
- [19] G. Liu, Z. Lin, S. Yan, J. Sun, Y. Yu, and Y. Ma. Robust recovery of subspace structures by low-rank representation. *CoRR*, abs/1010.2955, 2010. 4
- [20] S. Ma, D. Goldfarb, and L. Chen. Fixed point and bregman iterative methods for matrix rank minimization. *Mathematical Programming, Series A*, 128(1,2):321–353, 2011. 5
- [21] M. Paladini, A. D. Bue, M. Stosic, M. Dodig, J. Xavier, and L. Agapito. Factorization for non-rigid and articulated structure using metric projections. In *CVPR*, pages 2898–2905, 2009. 1, 4, 5
- [22] V. Rabaud and S. Belongie. Linear embeddings in non-rigid structure from motion. In *CVPR*, pages 2427–2434, 2009. 2
- [23] B. Recht, M. Fazel, and P. A. Parrilo. Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Review*, 52(3):471–501, 2010. 3
- [24] J. Taylor, A. D. Jepson, and K. N. Kutulakos. Non-rigid structure from locally-rigid motion. In *CVPR*, pages 2761–2768, 2010. 2
- [25] L. Torresani, A. Hertzmann, and C. Bregler. Nonrigid structure-from-motion: Estimating shape and motion with hierarchical priors. *PAMI*, 30, 2008. 1, 4, 5, 6
- [26] J. Xiao, J.-x. Chai, and T. Kanade. A closed-form solution to non-rigid shape and motion recovery. In *ECCV*, pages 573–587, 2004. 1, 2, 3, 4, 5