

# Jumping Manifolds: Geometry Aware Dense Non-Rigid Structure from Motion

## Supplementary Material

Suryansh Kumar  
Australian National University, CECS, Canberra, ACT, 2601  
suryansh.kumar@anu.edu.au

### Abstract

*This document provides mathematical derivation to the objective function proposed in the main paper [2]. Additionally, we provide some more qualitative results and statistical evaluations of our algorithm. Lastly, we made a brief comment on the challenges associated with handling temporal grassmannians for NRSfM problem.*

### 1. Mathematical derivation to the optimization of the objective function

In this section, we provide mathematical derivation of the following optimization proposed in the paper. Some of the derivations are similar to the Kumar *et al.* work [4, 5, 7]

$$\begin{aligned} & \underset{\mathbf{Z}, \tilde{\mathbf{C}}, \mathbf{S}, \mathbf{S}^\sharp}{\text{minimize}} \frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_{\mathbb{F}}^2 + \beta_1 \|\mathbf{X} - \chi \tilde{\mathbf{C}}\|_{\mathbb{F}}^2 + \beta_2 \|\mathbf{S}^\sharp\|_* + \\ & \quad \frac{\rho}{2} \|\mathbf{S}^\sharp - f(\mathbf{S})\|_{\mathbb{F}}^2 + \langle \mathbf{L}_1, \mathbf{S}^\sharp - f(\mathbf{S}) \rangle + \beta_3 \|\mathbf{Z}\|_* + \\ & \quad \frac{\rho}{2} \|\tilde{\mathbf{C}} - \mathbf{Z}\|_{\mathbb{F}}^2 + \langle \mathbf{L}_2, \tilde{\mathbf{C}} - \mathbf{Z} \rangle \\ & \text{subject to: } \xi = f_g(\mathbf{P}, \mathbf{S}), \tilde{\xi} = f_h(\Delta, \xi), \\ & \mathbf{S} = f_s(\xi, \Sigma, \xi_v), \mathbf{P} = f_p(\tilde{\xi}, \tilde{\mathbf{C}}, \mathbf{P}_o) \end{aligned} \quad (1)$$

The constraints in the Eq.(1) are invoked over iteration. The solution to each sub-problem is obtained by taking the derivative of the above ALM form w.r.t the concerned variable and equating it to zero.

#### 1.1. Solution to ‘S’

$$\begin{aligned} & \equiv \underset{\mathbf{S}}{\text{argmin}} \frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_{\mathbb{F}}^2 + \frac{\rho}{2} \|\mathbf{S}^\sharp - f(\mathbf{S})\|_{\mathbb{F}}^2 + \\ & \quad \langle \mathbf{L}_1, \mathbf{S}^\sharp - f(\mathbf{S}) \rangle \\ & \equiv \underset{\mathbf{S}}{\text{argmin}} \frac{1}{2} \|\mathbf{W} - \mathbf{R}\mathbf{S}\|_{\mathbb{F}}^2 + \frac{\rho}{2} \left\| \mathbf{S} - \left( f^{-1}(\mathbf{S}^\sharp) + \frac{f^{-1}(\mathbf{L}_1)}{\rho} \right) \right\|_{\mathbb{F}}^2 \end{aligned} \quad (2)$$

Taking the derivative of the above equation w.r.t ‘S’ and equating it to zero gives

$$(\mathbf{R}^T \mathbf{R} + \rho \mathbf{I}) \mathbf{S} = \mathbf{R}^T \mathbf{W} + \rho \left( f^{-1}(\mathbf{S}^\sharp) + \frac{f^{-1}(\mathbf{L}_1)}{\rho} \right) \quad (3)$$

We used MATLAB `mldivide()` function to solve it during our implementation. You may use any linear algebra package to solve the above well known form.

#### 1.2. Solution to ‘S<sup>♯</sup>’

Similar to previous derivation, we can write the ALM form for the variable S<sup>♯</sup>

$$\begin{aligned} & \equiv \underset{\mathbf{S}^\sharp}{\text{argmin}} \beta_2 \|\mathbf{S}^\sharp\|_* + \frac{\rho}{2} \|\mathbf{S}^\sharp - f(\mathbf{S})\|_{\mathbb{F}}^2 + \langle \mathbf{L}_1, \mathbf{S}^\sharp - f(\mathbf{S}) \rangle \\ & \equiv \underset{\mathbf{S}^\sharp}{\text{argmin}} \beta_2 \|\mathbf{S}^\sharp\|_* + \frac{\rho}{2} \left\| \mathbf{S}^\sharp - \left( f(\mathbf{S}) - \frac{\mathbf{L}_1}{\rho} \right) \right\|_{\mathbb{F}}^2 \end{aligned} \quad (4)$$

The above sub-problem is well-known form for nuclear norm minimization. By defining the soft-thresholding operator  $\mathcal{S}_\tau(v) = \text{sign}(v) \max(|v| - \tau, 0)$ , the solution of S<sup>♯</sup> can be obtained by

$$\mathbf{S}^\sharp = \mathbf{U}_s \mathcal{S}_{\frac{\beta_2}{\rho}}(\Sigma_s) \mathbf{V}_s \quad (5)$$

where,  $[\mathbf{U}_s, \Sigma_s, \mathbf{V}_s] = \text{svd}\left(f(\mathbf{S}) - \frac{\mathbf{L}_1}{\rho}\right)$

#### 1.3. Solution to ‘Z’

$$\begin{aligned} & \equiv \underset{\mathbf{Z}}{\text{argmin}} \beta_3 \|\mathbf{Z}\|_* + \frac{\rho}{2} \|\tilde{\mathbf{C}} - \mathbf{Z}\|_{\mathbb{F}}^2 + \langle \mathbf{L}_2, \tilde{\mathbf{C}} - \mathbf{Z} \rangle \\ & \equiv \underset{\mathbf{Z}}{\text{argmin}} \beta_3 \|\mathbf{Z}\|_* + \frac{\rho}{2} \left\| \mathbf{Z} - \left( \tilde{\mathbf{C}} + \frac{\mathbf{L}_2}{\rho} \right) \right\|_{\mathbb{F}}^2 \end{aligned} \quad (6)$$

Using the soft-thresholding function as mentioned before, the solution to Z is given by

$$\mathbf{Z} \equiv \mathbf{U}_z \mathcal{S}_{\frac{\beta_3}{\rho}}(\Sigma_z) \mathbf{V}_z \quad (7)$$

where  $[\mathbf{U}_z, \Sigma_z, \mathbf{V}_z] = \text{svd}\left(\tilde{\mathbf{C}} + \frac{\mathbf{L}_2}{\rho}\right)$

#### 1.4. Solution to ‘ $\tilde{\mathbf{C}}$ ’

Deriving the solution for ‘ $\tilde{\mathbf{C}}$ ’ from the sub-problem involving the variable ‘ $\tilde{\mathbf{C}}$ ’ is not straight forward rather, it’s a bit involved and therefore, we first derive an equivalent form for the error term that involves tensor  $\chi$ . The equivalent form is easy to handle and program on computers. Lets consider the following error term:

$$\|\chi - \chi\tilde{\mathbf{C}}\|_{\mathbb{F}}^2 \quad (8)$$

Using the notation from our paper, for any  $i^{\text{th}}$  Grassmann point this error term in Eq:(8) can be written as

$$\text{Tr}\left(\left((\Theta_i\Theta_i^T) - \sum_{j=1}^K c_{ij}(\Theta_j\Theta_j^T)\right)^T \left((\Theta_i\Theta_i^T) - \sum_{j=1}^K c_{ij}(\Theta_j\Theta_j^T)\right)\right) \quad (9)$$

Expanding the above form gives

$$\begin{aligned} &\equiv \text{Tr}((\Theta_i\Theta_i^T)^T(\Theta_i\Theta_i^T)) - 2\sum_{j=1}^K c_{ij}\text{Tr}((\Theta_i\Theta_i^T)^T(\Theta_j\Theta_j^T)) + \\ &\sum_{l=1}^K \sum_{m=1}^K c_{il}c_{im}\text{Tr}((\Theta_i\Theta_i^T)^T(\Theta_m\Theta_m^T)) \end{aligned} \quad (10)$$

From our definition  $\Theta \in \mathbb{R}^{\tilde{\mathbf{d}} \times p}$  as an orthonormal matrix. Using it simplifies the above equation to:

$$\begin{aligned} &\equiv p - 2\sum_{j=1}^K c_{ij}\Gamma_{ij} + \sum_{l=1}^K \sum_{m=1}^K c_{il}c_{im}\Gamma_{lm} \\ &\text{where, } \Gamma_{ij} = \text{Tr}((\Theta_i^T\Theta_i)(\Theta_j^T\Theta_j)) \quad \{\text{using trace cyclic property}\} \end{aligned} \quad (11)$$

Let  $\Gamma = (\Gamma_{ij})_{i,j=1}^K \in \mathbb{R}^{K \times K}$ . Its easy to verify that  $\Gamma$  is symmetric positive semi-definite. Therefore, using cholsky factorization of  $\text{chol}(\Gamma) = \mathbf{L}\mathbf{L}^T$ , we can re-write the above equation as

$$\begin{aligned} &\equiv p - 2\text{Tr}(\tilde{\mathbf{C}}\mathbf{L}\mathbf{L}^T) + \text{Tr}(\tilde{\mathbf{C}}\mathbf{L}\mathbf{L}^T\tilde{\mathbf{C}}^T) \\ &\equiv \text{const} + \|\mathbf{L} - \mathbf{L}\tilde{\mathbf{C}}\|_{\mathbb{F}}^2 \end{aligned} \quad (12)$$

where, const. means constant w.r.t  $\tilde{\mathbf{C}}$

By substituting the result from Eq:(12) to the sub-problem w.r.t  $\tilde{\mathbf{C}}$ , we get the following form:

$$\begin{aligned} &\equiv \underset{\tilde{\mathbf{C}}}{\text{argmin}} \beta_1 \|\mathbf{L} - \mathbf{L}\tilde{\mathbf{C}}\|_{\mathbb{F}}^2 + \frac{\rho}{2} \|\tilde{\mathbf{C}} - \mathbf{Z}\|_{\mathbb{F}}^2 < \mathbf{L}_2, \tilde{\mathbf{C}} - \mathbf{Z} > \\ &\equiv \underset{\tilde{\mathbf{C}}}{\text{argmin}} \beta_1 \|\mathbf{L} - \mathbf{L}\tilde{\mathbf{C}}\|_{\mathbb{F}}^2 + \frac{\rho}{2} \|\tilde{\mathbf{C}} - \left(\mathbf{Z} - \frac{\mathbf{L}_2}{\rho}\right)\|_{\mathbb{F}}^2 \end{aligned} \quad (13)$$

Taking the derivative of the Eq:(13) w.r.t  $\tilde{\mathbf{C}}$  and equating it to zero.

$$(2\beta_1\mathbf{L}\mathbf{L}^T + \rho\mathbf{I})\tilde{\mathbf{C}} = 2\beta_1\mathbf{L}\mathbf{L}^T + \rho\left(\mathbf{Z} - \frac{\mathbf{L}_2}{\rho}\right) \quad (14)$$

#### 2. Solution to $\mathbf{E}(\Delta)$

$$\begin{aligned} \mathbf{E}(\Delta) &\equiv \underset{\Delta}{\text{minimize}} \sum_{(i,j)}^K w_{ij} \frac{1}{2} \|\Delta^T(\Lambda_{ij})\Delta\|_{\mathbb{F}}^2 \\ &\text{subject to:} \end{aligned} \quad (15)$$

$$\text{Tr}\left(\Delta^T \left(\sum_{i=1}^K \lambda_{ii} \Omega_i \Omega_i^T\right) \Delta\right) = 1$$

The optimization equation proposed for  $\mathbf{E}(\Delta)$  is a well-studied optimization form and Riemann Conjugate gradient toolbox can be employed to achieve the solution. Nevertheless, we can also derive augmented lagrangian form to solve the same problem. By letting  $\mathbf{X} = (\sum_{i=1}^K \lambda_{ii} \Omega_i \Omega_i^T)$  and expanding the Frobenius norm term, we can re-write the equation as:

$$\mathbf{E}(\Delta) \equiv \underset{\Delta}{\text{minimize}} \sum_{(i,j)}^K \frac{w_{ij}}{2} \text{Tr}(\Delta^T \Lambda_{ij} \Delta \Delta^T \Lambda_{ij} \Delta)$$

$$\mathbf{E}(\Delta) \equiv \underset{\Delta}{\text{minimize}} \text{Tr}\left(\Delta^T \sum_{(i,j)}^K \frac{w_{ij}}{2} \Lambda_{ij} \Delta^{t-1} \Delta^{(t-1)^T} \Lambda_{ij} \Delta\right)$$

subject to:

$$\text{Tr}(\Delta^T \mathbf{X} \Delta) = 1 \quad (16)$$

Here,  $t-1$  refers to its known value before the current iteration. Now, by assuming  $\mathbf{Y} = \frac{w_{ij}}{2} \Lambda_{ij} \Delta^{t-1} \Delta^{(t-1)^T} \Lambda_{ij}$ , the above equation simplifies to standard eigen value decomposition problem *i.e.*

$$\begin{aligned} \mathbf{E}(\Delta) &\equiv \underset{\Delta}{\text{minimize}} \text{Tr}(\Delta^T \mathbf{Y} \Delta) \\ &\text{subject to:} \\ &\text{Tr}(\Delta^T \mathbf{X} \Delta) = 1 \end{aligned} \quad (17)$$

The equivalent Lagrangian function form is given by

$$\text{Tr}(\Delta^T \mathbf{Y} \Delta) + \lambda \left(1 - \text{Tr}(\Delta^T \mathbf{X} \Delta)\right) \quad (18)$$

The Eq:(18) is of the standard form to generalized eigen value problem. You may use any standard linear algebra package to solve it.

#### 3. More Results and Analysis

We provide some more qualitative results on Garg *et al.* benchmark dataset [1]. Figure (1) and Figure (2) show the 3D reconstruction results of our algorithm on real face and synthetic face sequence. Next, we provide one more statistical analysis of the our algorithm.

**1. Dependence of the algorithm on variable  $\tilde{\mathbf{d}}$ :** While reducing the dimension for grouping the grassmann points, one of the critical aspect is to determine the dimension to which we should project for better results. We used well-known procedure of cumulative energy of eigen vectors to get the value of  $\tilde{\mathbf{d}}$ . Mathematically, let  $\Omega$  be the set that stack all the Grassmannians and  $\sigma_i$

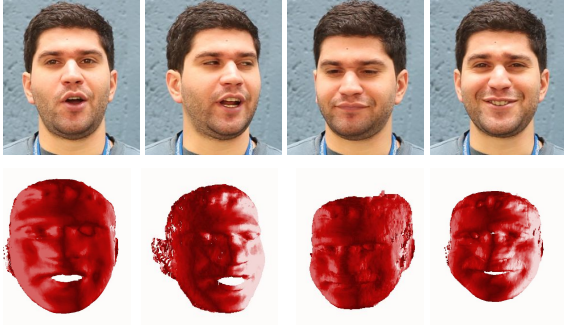


Figure 1: 3D reconstruction results on real face sequence [1]

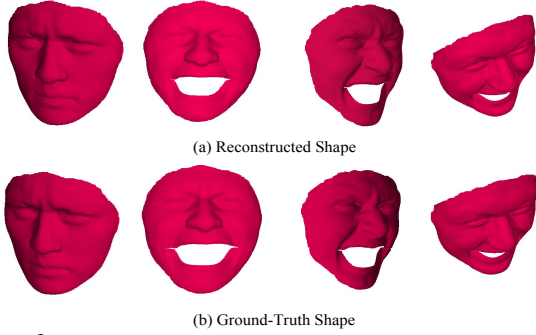


Figure 2: 3D reconstruction results on synthetic face sequence [1].

be the  $i^{\text{th}}$  singular value of  $\Omega\Omega^T$ , then

$$\tilde{d} = \underset{d}{\operatorname{argmin}} \frac{\sum_{i=1}^{d_{\text{opt}}} \sigma_i}{\sum_{i=1}^d \sigma_i} \geq \tau \quad (19)$$

where  $\tau$  can vary from 0 to 1 and  $d_{\text{opt}}$  (optimal dimension) is a positive integer. We put  $\tau = 0.97$  for all our experiment. Figure (3) show the variations in the reconstruction error with the value of  $\tau$ . It is observed that for different dataset the value of suitable  $\tilde{d}$  is different. The point to note is that if the reduced dimension is less than the intrinsic dimension, the samples may lose important information for better grouping of Grassmannians.

### 3.1. Why we opt not to disturb the temporal continuity for this problem?

Although clustering of frames into smaller groups (Grassmannians) allows simpler model and be handy if prior informations about shapes/activities are available. However, its quite possible that there will be repeat of certain activities or expression in the video sequence (say facial expression). In such cases the Grassmannians at frame ‘f’ and frame ‘f+n’ will be assigned to same group (‘n’ is the time instant at which activities repeat or is similar). As a result, such representation procedure may disturb the overall time continuity of the sequence. Also, these group of frames may form high-dimensional grassmannians, in order to project it into low-dimension using neighboring Grassmannians will get extremely difficult, for example, how to decide neighboring grassmannians using temporal grassmann samples?. On the

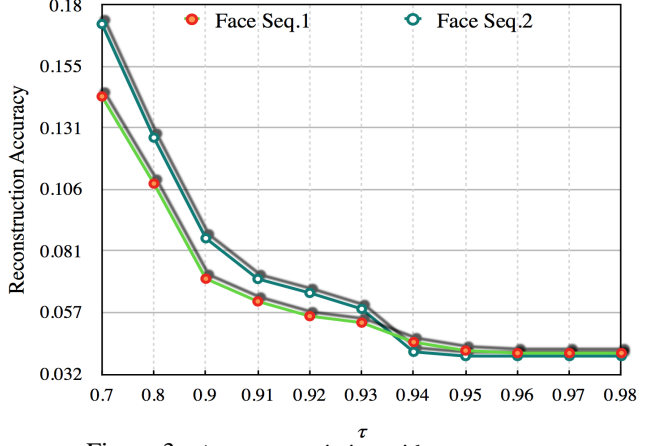


Figure 3: Accuracy variation with respect to  $\tau$ .

other hand, grouping of trajectories (spatial) does not disturb the temporal continuity of the trajectory and we can easily define the neighbors using spatial information *i.e.*, spatial neighbors tend to be neighbors throughout the sequence, for a single deforming object (unless breaks or disassociate, which is very rare). But in shape space, we don’t have any prior knowledge to define neighboring relation.

### 3.2. What is gained with the added extra complexity in the algorithm (Ablation Test)?

By adding Eq.(7) we can have more discriminative Grassmannian representation which is useful in practice. The performance on noisy sequence clearly demonstrates this. Also, by adding this local representation constraint we have better control over local shapes with very minimal increment in the processing time.

Data	WO_NC (Noiseless)	WO_NC (Noise5%)	W_NC (Noiseless)	W_NC (Noise5%)
Seq.1	0.1083	0.1592	0.0404	0.0482
Seq.2	0.0972	0.1491	0.0392	0.0429
Seq.3	0.0913	0.1354	0.0280	0.0365
Seq.4	0.0924	0.1433	0.0327	0.0402

Table 1: Ablation test showing the ( $e_{3D}$ ) comparison on dense synthetic face sequence. WO\_NC means WithOut Neighboring Constraint and W\_NC means With Neighboring Constraints.

### 3.3. Mean processing time per frame?

The mean processing time per frame is 10.57743s in comparison to [4] which is 9.53093s.

### 3.4. What we gained over [4]?

Temporal clustering does help improve reconstruction accuracy as shown in [4] but [4] algorithm does it by assuming that they know how many frames to select to represent a local temporal Grassmannian. Now such information is not available in practice. Our new algorithm show that without such temporal knowledge we can have an equivalent or better performance by using a better

representation of the grassmannians. This can be verified using the reconstruction accuracy table in the main paper.

## References

- [1] Ravi Garg, Anastasios Roussos, and Lourdes Agapito. Dense variational reconstruction of non-rigid surfaces from monocular video. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1272–1279, 2013.
- [2] Suryansh Kumar. Jumping manifolds: Geometry aware dense non-rigid structure from motion. *arXiv preprint arXiv:1902.01077*, 2019.
- [3] Suryansh Kumar. Non-rigid structure from motion: Prior-free factorization method revisited. *arXiv preprint arXiv:1902.10274*, 2019.
- [4] Suryansh Kumar, Anoop Cherian, Yuchao Dai, and Hongdong Li. Scalable dense non-rigid structure-from-motion: A grassmannian perspective. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [5] Suryansh Kumar, Yuchao Dai, and Hongdong Li. Multi-body non-rigid structure-from-motion. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 148–156. IEEE, 2016.
- [6] Suryansh Kumar, Yuchao Dai, and Hongdong Li. Monocular dense 3d reconstruction of a complex dynamic scene from two perspective frames. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4649–4657, Oct 2017.
- [7] Suryansh Kumar, Yuchao Dai, and Hongdong Li. Spatio-temporal union of subspaces for multi-body non-rigid structure-from-motion. *Pattern Recognition*, 71:428–443, May 2017.
- [8] Suryansh Kumar, Ram Srivatsav Ghorakavi, Yuchao Dai, and Hongdong Li. Dense depth estimation of a complex dynamic scene without explicit 3d motion estimation. *arXiv preprint arXiv:1902.03791*, 2019.