

A Factorization Approach to Structure from Motion with Shape Priors*

Alessio Del Bue

Institute for Systems and Robotics – Instituto Superior Técnico
Av. Rovisco Pais 1 – 1049-001 Lisboa – Portugal

<http://www.isr.ist.utl.pt/~adb/>

Abstract

This paper presents an approach for including 3D prior models into a factorization framework for structure from motion. The proposed method computes a closed-form affine fit which mixes the information from the data and the 3D prior on the shape structure. Moreover, it is general in regards to different classes of objects treated: rigid, articulated and deformable. The inclusion of the shape prior may aid the inference of camera motion and 3D structure components whenever the data is degenerate (i.e. nearly planar motion of the projected shape). A final non-linear optimization stage, which includes the shape priors as a quadratic cost, upgrades the affine fit to metric. Results on real and synthetic image sequences, which present predominant degenerate motion, make clear the improvements over the 3D reconstruction.

1. Introduction

Extracting 3D data from monocular image sequences is a problem extensively studied in Computer Vision. The idea embracing various scenarios is relatively similar: to infer both the 3D structure parametrization and the camera parameters from a set of 2D points extracted from a sequence which depicts a moving object. In this context, methods based on the bilinear factorization of the image data have been very successful in both proposing simple and closed-form approaches with the least assumptions as possible.

Back in the early '90s, Tomasi and Kanade [14] introduced the first factorization algorithm dealing with a rigid object viewed by a simple orthographic camera. Later on, studies were mostly focused on upgrading the approach to more complex viewing conditions such as paraperspective [11] and projective [13]. Only recently, the framework was extended to deal with objects which may also vary their shape. Deformable [5] and articulated [17, 20] factoriza-

tions are among these examples.

However, one of the main problems of Structure from Motion (SfM) consists in the higher complexity in which shapes may vary their 3D structure. This may add nonlinearities and strong dependencies between both shape and motion components resulting in non-trivial solutions of the problem. Likewise, the higher number of degrees of freedom and the possible motion degeneracies in the measured data may lead to a local solution which correctly minimizes the 2D reprojection error but, however, results in a poor or even meaningless 3D reconstruction.

In order to counter this effect, prior information may be included to obtain reliable 3D reconstructions. Different priors have been shown to improve performances in rigid and non-rigid SfM. Forsyth et al. [7] first proposed priors over the specific camera constraints in a consistent Bayesian framework. Xiao et al. [19] use the prior information over a set of independent shapes to compute a closed-form solution. In a face modelling context, Solem and Kahl [12] used a learned shape model to aid the 3D inference over regions with no 2D information available. Del Bue et al. [6] enforce priors over the rigidity of some points to obtain reliable estimations of the object rigid components. Torresani et al. [15] propose the use of gaussian priors over the deformation parameters in order to avoid arbitrarily variations. Finally, Olsen and Bartoli [10] impose a prior over temporal variations of the camera parameters combined with constraints over the proximity of projected 2D points and reconstructed 3D points.

Differently from the mentioned solutions, priors are here introduced in the form of previously computed 3D shapes. In such way, it is possible to obtain a description of the object shape jointly given by the measured data and the prior information available. This approach especially supports the computation of reliable 3D shapes whenever the image sequence contains strong degeneracies. For instance, a common case is a talking head in front of a camera performing tiny pose changes. In this case depth would be lost if not irremediably mixed with the ongoing deformations. A known 3D metric description of the subject face may how-

¹To appear in CVPR08, Proc. of IEEE Conf. on Computer Vision and Pattern Recognition, Anchorage, Alaska, June 2008.

ever recover the lost depth and disambiguate better motion and shape components. At this end, the proposed computational methods directly include the 3D metric priors in a factorization framework for SfM. First, the approach uses the data and the prior to compute an initial affine solution which is then upgraded to metric with an iterative non-linear optimization procedure.

The next section provides an introduction to SfM methods showing how different problems can be treated with a unique approach. Section 3 shows how to compute a metric solution for the motion and structure components. The proposed algorithm is then explained in Section 4 showing computational tools to compute 3D models given prior information. The experiments in Section 5 show real and synthetic examples of 3D reconstruction in such cases.

2. Factorization for SfM

The key idea in SfM is to gather all the 2D image coordinates lying on a generic shape at each frame in a single measurement matrix W . The location of a point j in a certain frame i can be defined by a vector $\mathbf{w}_{ij} = (u_{ij} \ v_{ij})^T$ where u_{ij} and v_{ij} are the horizontal and vertical image coordinates respectively. A compact matrix representation can be expressed as:

$$W = \begin{bmatrix} \mathbf{w}_{11} & \dots & \mathbf{w}_{1p} \\ \vdots & \ddots & \vdots \\ \mathbf{w}_{f1} & \dots & \mathbf{w}_{fp} \end{bmatrix} = \begin{bmatrix} W_1 \\ \vdots \\ W_f \end{bmatrix} \quad (1)$$

where f is the total number of image frames and p the number of image points. The image trajectories stored in W can be expressed as a bilinear product as $W = M_{2f \times r} S_{r \times p}$. The matrices M and S refer to the motion and shape subspaces respectively with dimension r where $r \ll \min\{2f, p\}$. As a result, the rank of W is constrained to be $\text{rank}\{W\} \leq r$. The matrices M and S can be further decomposed in:

$$M = \begin{bmatrix} M_1 \\ \vdots \\ M_f \end{bmatrix} \quad S = [S_1 \ \dots \ S_p] \quad (2)$$

where M_i with $i = 1 \dots f$ is a $2 \times r$ matrix projecting the shape onto the image frame i . The size of M_i directly depends on the type of camera and motion that appears in the scene. The component S_j with $j = 1 \dots p$ is a r -vector that defines the 3D parametrization for each point j and its size depends on the shape properties (i.e. rigid or non-rigid).

2.1. Rigid factorization

The projection of a 3D rigid shape by means of an orthographic camera is historically the first factorization problem studied [14]. In this case the camera motion M_i and the 3D

point S_j can be expressed as:

$$M_i = \begin{bmatrix} r_{i1} & r_{i2} & r_{i3} & t_{ui} \\ r_{i4} & r_{i5} & r_{i6} & t_{vi} \end{bmatrix} = [R_i \mid \mathbf{t}_i] \quad (3)$$

$$S_j = [X_j \ Y_j \ Z_j \ 1]^T = \begin{bmatrix} \mathbf{X}_j \\ 1 \end{bmatrix}$$

where R_i contains the first two rows of a rotation matrix (i.e. $R_i R_i^T = I_{2 \times 2}$), S_j is a 4-vector containing the homogeneous metric coordinates of the 3D point \mathbf{X}_j , and \mathbf{t}_i is a 2-vector representing a translation into the image plane. Every point belonging to the rigid structure shares the same rotation and translation, thus we can compact 3D points in a single $4 \times p$ matrix S giving:

$$W_i = [R_i \mid \mathbf{t}_i] [S_1 \ \dots \ S_p] = M_i S. \quad (4)$$

Stacking the rows of W_i for every frame, we obtain the full measurement matrix as:

$$W = MS = [R \mid \mathbf{t}] S \quad (5)$$

where R is the $2f \times 3$ collection of f rotation matrices, \mathbf{t} is a $2f$ -vector which contains the translation for every frame. The dimension of M and S is fixed to $r = 4$. If the 2D points in W are registered to the shape centroid (i.e. $W \mathbf{1}^T = \mathbf{0}$), the maximum dimensionality of each subspace is $r = 3$ and equation (5) can be written as:

$$W = MS = R [X_1 \ \dots \ X_p]. \quad (6)$$

2.2. Articulated factorization

If the measurements in W belong to two independent moving objects, the overall rank sums to eight since it is possible to write motion and shape components as:

$$M_i = \begin{bmatrix} R_i^{(1)} \mid \mathbf{t}_i^{(1)} \mid R_i^{(2)} \mid \mathbf{t}_i^{(2)} \end{bmatrix} \quad (7)$$

$$S = \begin{bmatrix} S^{(1)} & \mathbf{0} \\ \mathbf{0} & S^{(2)} \end{bmatrix}$$

such that:

$$W = \begin{bmatrix} W_i^{(1)} \mid W_i^{(2)} \end{bmatrix} = M_i S \quad (8)$$

where $W_i^{(1)}$ and $W_i^{(2)}$ are the measured data at frame i for the first and second shape respectively. The components of M_i and S for each shape are in the form as shown by equation (4) for a single rigid object. However, the rank r may decrease if the moving objects show a dependency such as a common rotational axis.

In articulated SfM [17, 20], this dependency is given by the joints which constrain the degrees of freedom of the moving objects. In the case of a *universal joint* [17] the distance between the center of the shapes is constant (for instance, the head and the torso of a human body) but they

show independent rotation components. At each frame the shapes connected by a joint satisfy:

$$\mathbf{t}^{(1)} + \mathbf{R}^{(1)}\mathbf{d}^{(1)} = \mathbf{t}^{(2)} + \mathbf{R}^{(2)}\mathbf{d}^{(2)} \quad (9)$$

where $\mathbf{t}^{(1)}$ and $\mathbf{t}^{(2)}$ are the 2D image centroid of the two objects, $\mathbf{R}^{(1)}$ and $\mathbf{R}^{(2)}$ the 2×3 orthographic camera matrices and $\mathbf{d}^{(1)}$ and $\mathbf{d}^{(2)}$ the 3D displacement vectors of each shape from the joint. The relation in equation (9) gives the reduced dimensionality in the motion and shape subspaces. Thus, the shape matrix \mathbf{S} can be written as:

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}^{(1)} & \mathbf{d}^{(1)} \\ 0 & \mathbf{S}^{(2)} - \mathbf{d}^{(2)} \\ 1 & 1 \end{bmatrix} \quad (10)$$

where \mathbf{S} is a full rank-7 matrix. The motion for a frame i has to be accordingly arranged to satisfy equation (9) as:

$$\mathbf{M}_i = \begin{bmatrix} \mathbf{R}_i^{(1)} & \mathbf{R}_i^{(2)} & \mathbf{t}_i^{(1)} \end{bmatrix}. \quad (11)$$

Further details can be found in [17], alongside a description of additional joint models. Notice that is necessary to solve an assignment problem in order to form \mathbf{W} as in equation (8). At this end, a recent approach able to detect articulated parts from images is presented in [20].

2.3. Deformable factorization

In the case of deformable objects, a single shape varies its 3D structure with respect to a set of deformation modes. The number of modes used to define the shape deformations results in a specific rank-constraint over the image trajectories in \mathbf{W} . The representation for the deformations is a simple model where any specific 3D configuration \mathbf{X} is approximated by a linear combination of a set of k basis shapes \mathbf{B}_d which represent the principal modes of deformation:

$$\mathbf{X} = \sum_{d=1}^k c_d \mathbf{B}_d \quad \mathbf{X}, \mathbf{B}_d \in \mathfrak{R}^{3 \times p} \quad c_d \in \mathfrak{R} \quad (12)$$

Bregler et al. [5] were the first to propose an extension of factorization algorithms able to deal with the case of deformable shapes assuming an orthographic camera model. In this case, the coordinates of the 2D image points observed at each frame i are related to the coordinates of the 3D points according to the following equation:

$$\mathbf{W}_i = \mathbf{R}_i \left(\sum_{d=1}^k c_{id} \mathbf{B}_d \right) + \mathbf{T}_i \quad (13)$$

where c_{id} is the configuration weight for basis d at frame i . When the image coordinates are registered to the object's centroid, equation (13) can be rewritten as:

$$\mathbf{W}_i = \begin{bmatrix} c_{i1}\mathbf{R}_i & \dots & c_{ik}\mathbf{R}_i \end{bmatrix} \begin{bmatrix} \mathbf{B}_1 \\ \vdots \\ \mathbf{B}_k \end{bmatrix} = \mathbf{M}_i \mathbf{S} \quad (14)$$

Again, by stacking the measurements at each frame we obtain the following compact matrix form:

$$\mathbf{W} = \begin{bmatrix} c_{11}\mathbf{R}_1 & \dots & c_{1k}\mathbf{R}_1 \\ \vdots & \ddots & \vdots \\ c_{f1}\mathbf{R}_f & \dots & c_{fk}\mathbf{R}_f \end{bmatrix} \begin{bmatrix} \mathbf{B}_1 \\ \vdots \\ \mathbf{B}_k \end{bmatrix} = \mathbf{M} \mathbf{S} \quad (15)$$

Since \mathbf{M} is a $2f \times 3k$ matrix and \mathbf{S} is a $3k \times p$ matrix, the rank of \mathbf{W} when no noise is present must be $r \leq 3k$.

3. Affine and metric reconstruction

In order to extract the motion and shape components, the classical procedure solves two separate problems:

1. Find an affine fit $\tilde{\mathbf{M}}$ and $\tilde{\mathbf{S}}$ to the measured data \mathbf{W} given the rank r .
2. Enforce the metric structure by computing a corrective transform $\mathbf{Q}_{r \times r}$ such that $\mathbf{W} = \tilde{\mathbf{M}} \mathbf{Q} \mathbf{Q}^{-1} \tilde{\mathbf{S}} = \mathbf{M} \mathbf{S}$.

The first step can be trivially solved using any rank revealing technique such as SVD. However this solution is merely one means of numerical computations, other approaches may be used as pointed out in [9]. Given the chosen rigid/non-rigid model, the numerical rank it is fixed to r and thus it is possible to approximate the decomposition as:

$$\mathbf{W} \xrightarrow{SVD} \sum_{i=1}^n \mathbf{u}_i \sigma_i \mathbf{v}_i^T = \sum_{i=1}^r \mathbf{u}_i \sigma_i \mathbf{v}_i^T = \mathbf{U}_r \Sigma_r \mathbf{V}_r^T \quad (16)$$

where \mathbf{U}_r is a $2f \times r$ orthogonal matrix, Σ_r a $r \times r$ diagonal matrix and \mathbf{V}_r^T a $p \times r$ orthogonal matrix. After simple operations the product can be arranged in the affine form $\mathbf{W} = \tilde{\mathbf{M}} \tilde{\mathbf{S}}$. Notice that this initial solution is independent from the problem considered, either the considered shape is rigid or non-rigid.

The second step, the computation of \mathbf{Q} , is where the metric properties are imposed. For instance, in the rigid orthographic case, this leads to the computation of the transformation $\mathbf{Q}_{3 \times 3}$ which renders each row of $\tilde{\mathbf{M}}_i$ orthonormal (i.e. $\tilde{\mathbf{M}}_i \rightarrow \mathbf{R}_i$). In the non-rigid case the computation is more complex since the relationships in \mathbf{M} are strongly non-linear. For instance, in the deformable case, a closed form solution can be found only if a correct set of independent bases is chosen [19]. A wrong set of bases may lead to inaccurate solutions as shown in [4]. On the other hand, non-linear optimization or hybrid approaches have been proposed to solve the problem [4, 15]. In the case of articulated shapes, solutions are available [17] but based on the correct knowledge of the type of joint connecting the bodies.

4. Factorization with shape priors

It is clear from the previous section that the metric upgrade is dependent on the initial affine decomposition. An

inaccurate affine fit of \tilde{M} and \tilde{S} may irremediably compromise the following metric upgrade. In this sense, the introduction of a prior on the values of the shape subspace S can bring the estimation close to the desired solution. Moreover, shape priors can support the computation in case of deficient components in the motion subspace M – a likely case when the image motion is weak. In such cases, it is always possible to compute a factorization which well fit W numerically but nonetheless the factors \tilde{M} and \tilde{S} may contain meaningless components. This problem is hard to solve unless prior information is included in the computation.

4.1. Measurements and priors

Shape priors are introduced as a matrix $L \in \mathfrak{R}^{l \times p}$ which holds a parametrization of each image trajectory stored in W . A point trajectory \mathbf{w}_j can be written as a $2f$ -vector such that $W = [\mathbf{w}_1 \cdots \mathbf{w}_p]$. Thus accordingly, we can write the shape prior¹ as $L = [L_1 \cdots L_p]$. The size of each l -vector L_j depends on the type of prior chosen. Often L may represent a rigid 3D shape leading to $l = 3$. Such prior may be used in deformable SfM to obtain a more reliable estimation of an object with complex shape variations. Alternatively, it may support the estimation of the object's depth when the shape is moving planarly. In general, L may as well store more complex shape descriptions (i.e. $l > 3$) such as a set of deformable basis shapes computed previously from a similar shape.

The main idea is to join the information stored in L with the available measurements in W in order to extract an affine fit which is dependent on both components. This can be formulated as two bilinear models for the data:

$$W = \tilde{M}_{2f \times t} \tilde{S}_{t \times p} = [M_J | M_I] \begin{bmatrix} S_J \\ S_I \end{bmatrix} \quad (17)$$

and for the shape prior:

$$L = N_{l \times l} S_J \quad (18)$$

where the J subscript refers to the components obtained by the joint space between prior and image measurements while the I refers to the remaining ones. Notice that we always consider L being full rank thus the following properties hold:

$$\text{rank}(W) = r, \text{rank}(L) = l \text{ and } \text{rank} \left(\begin{bmatrix} W \\ L \end{bmatrix} \right) = t \quad (19)$$

where $t = \max\{r, l\}$ is the overall rank for both prior and measurements.

¹Notice that both the prior and measurements are registered to the respective centroids i.e. $W \mathbf{1}_{p \times 1} = \mathbf{0}_{2f \times 1}$ and $L \mathbf{1}_{p \times 1} = \mathbf{0}_{l \times 1}$.

4.2. Generalized singular value decomposition

Once the shape prior and the data are defined, we seek a computational solution able to find the joint factorization for W and L . This can be obtained using a *Generalized Singular Value Decomposition* (GSVD) which can decompose the matrix pair $\{W, L\}$ as:

$$\begin{aligned} W &= U D_U X^T \\ L &= V D_V X^T \end{aligned} \quad (20)$$

where X^T is a $p \times p$ matrix which span the common row space of $\{W, L\}$, U is a $2f \times 2f$ matrix with orthonormal columns ($U^T U = I$) and V is a $l \times l$ matrix such that $V^T V = I$. Notice that the matrix X is rank deficient with t non-zero singular values (for more details and proofs on GSVD see [2], Sec. 4.2.2). The diagonal value matrices D_U and D_V of size $2f \times p$ and $l \times p$ are arranged as:

$$D_U = \begin{bmatrix} \Sigma_U & 0 \\ 0 & I \end{bmatrix} \text{ and } D_V = \begin{bmatrix} \Sigma_V & 0 \\ 0 & 0 \end{bmatrix}. \quad (21)$$

The diagonal matrices $\Sigma_U = \text{diag}(\sigma_1, \dots, \sigma_l)$ and $\Sigma_V = \text{diag}(\mu_1, \dots, \mu_l)$ of size $l \times l$ are constrained such that $\Sigma_U^2 + \Sigma_V^2 = I$ and they have the diagonal elements ordered as:

$$0 \leq \sigma_1 \leq \dots \leq \sigma_l \leq 1 \text{ and } 1 \geq \mu_1 \geq \dots \geq \mu_l > 0$$

The ratio between the diagonal values are called *Generalized Singular Values* (GSV) and they are defined as $\gamma_i = \sigma_i / \mu_i$. As a further note, the data and prior matrices are usually pre-scaled such that $\|W\|^2 = \|L\|^2$. This condition [8] guarantees a well-conditioning of the matrix X and it is the only data scaling performed in the decomposition.

4.3. Generalized factorization for SfM

The GSVD decomposes the image measurements in W with a common row space which is dependent on both the measured data stored in W and the 3D prior information stored in L . However, to obtain a solution in the form of equation (17), we must reduce the decomposition of W given by GSVD into the two standard affine components $\tilde{M}_{2f \times t}$ and $\tilde{S}_{t \times p}$. Thus, it is convenient to split X in the components which are dependent on the prior (the first l) and the one dependent on the data (the remaining $p - l$) such that:

$$X = [X_J | X_I]. \quad (22)$$

In such way, we aim to preserve the common row space X_J^T component which was computed by GSVD.

Thus, equation (20) can be accordingly separated in two components W_J and W_I giving:

$$W = W_J + W_I = U_J \Sigma_U X_J^T + U_I X_I^T. \quad (23)$$

The matrix X_J^T of size $l \times p$ alone contains the row space components which mixes measurements and priors. The

row space in X_I of size $(p - l) \times l$ still entails the rank deficiency and it requires a further decomposition to extract the remaining $t - l$ components.

In order to obtain this further affine fit, the proposed method performs two projections of W_I : first over the subspace defined by X_J^T and then along its orthogonal complement. This is obtained by defining the orthogonal projector P such that:

$$P = X_J (X_J^T X_J)^{-1} X_J^T \quad (24)$$

giving:

$$W_I = W_I P + W_I P_{\perp} = W' + W'' \quad (25)$$

where $P_{\perp} = I - P$ which has the result to further split the remaining components in W' which still belongs to the subspace given the prior and W'' which is its orthogonal complement. The components in W'' can be reduced via SVD to obtain the remaining $t - l$ components giving:

$$W'' \xrightarrow{SVD} U_c D_c V_c^T \quad (26)$$

which can be re-arranged as:

$$M_I = U_c D_c \text{ and } S_I = V_c^T \quad (27)$$

where M_I is a $2f \times (t - l)$ matrix and S_I a $(t - l) \times p$ matrix. The remaining data W' projected along X_J^T is merged to the joint space obtaining the measurements W_g such that:

$$\begin{aligned} W_g &= W_J + W' = U_J \Sigma_U X_J^T + U_I X_I^T X_J (X_J^T X_J)^{-1} X_J^T \\ &= \left(U_J \Sigma_U + U_I X_I^T X_J (X_J^T X_J)^{-1} \right) X_J^T = M_J S_J \end{aligned}$$

where $S_J = X_J^T$ and M_J is given by the remaining factors. Given the factors M_J , M_I , S_J and S_I , it is possible now to form the bilinear decomposition as in equation (17).

4.4. Finding a metric solution

For the case of rigid shapes, the affine fit given by the priors can be finally upgraded to metric by forcing metric constraints in closed form for different type of cameras [9]. Differently for non-rigid shapes, we opt for a non-linear optimization stage based on bundle adjustment [18] where a prior on the rigid basis shape is included as an additional quadratic cost (i.e. L is a $3 \times p$ matrix). In principle, the prior may correspond to a full parametrization of a deformable shape, i.e. $l > 3$, however here we focus on priors which are describing the rigid component of a deformable object. Thus, the cost function minimized reflects the deformable model presented in Section 2.3 giving:

$$\min_{R_i, B_{dj}, c_{id}} \sum_{i,j} \| \mathbf{w}_{ij} - (R_i \sum_d c_{id} \mathbf{B}_{dj}) \|^2 + \sum_j \| \mathbf{B}_{1j} - \mathbf{C} \mathbf{l}_j \|^2$$

where \mathbf{B}_{dj} is the 3×1 basis component for the point j such that $\mathbf{B}_d = [\mathbf{B}_{d1} \cdots \mathbf{B}_{dp}]$. The matrix \mathbf{C} performs a metric



Figure 1. Three images sampled from a 160 frames sequence showing different facial expressions.

alignment of L to the first basis shape B_1 and it can be computed using standard Procrustes analysis. The minimization of the first sum of quadratic costs is equivalent to a Maximum Likelihood (ML) estimate of the model parameters if i.i.d gaussian noise is affecting the measurements. However, by including the second sum, we obtain a Maximum A Posteriori (MAP) estimate given the shape prior.

In order to initialize the non-linear optimization, the affine fit given by equation (17) is used to compute an initial metric solution for the prior constrained components M_J and S_J . This procedure has analogies with the approach first proposed in [16] where the rigid component was used to initialize a non-linear optimization procedure. The main difference here is that the rigid shape is given by a mixture of prior and measured data. To compute each rotation matrix R_i and the first configuration weight c_{i1} each frame-wise component of M_J can be decomposed using an orthonormal decomposition [3]. The remaining values c_{id} with $d = 2 \dots k$ are initialized close to zero. Finally notice that the non-linear optimization with $k = 1$ can be used to infer the 3D structure of a rigid object. For the articulated case, the algorithm presented in [17] can be extended to compute an affine fit using shape priors which can be then corrected by forcing specific metric constraints for the given joint. Non-linear optimization can be then applied by adding additional priors for each articulated part.

5. Experiments

The experiments are mainly focused on deformable and articulated shapes. First, synthetic tests are performed on a deformable face in order to verify the validity of the reconstruction using shape priors. The results from non-linear optimization are compared against ground truth obtained from a VICON motion capture system. Then, two further real tests show the method performances with real imaging condition for the deformable and articulated case.

5.1. Deformable face with ground truth

In this experiment 37 point tracks from a 160 frames long sequence were obtained with a VICON system which cap-

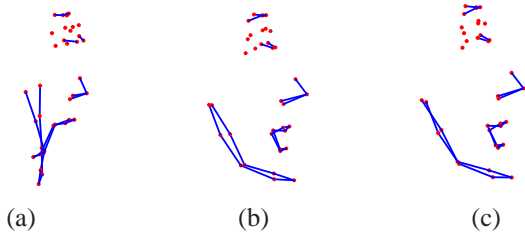


Figure 2. Shape prior and its effect over the estimated basis shapes. (a) The 1st deformable basis computed without priors. (b) The 3D prior used in the test. (c) The 1st basis computed using priors.

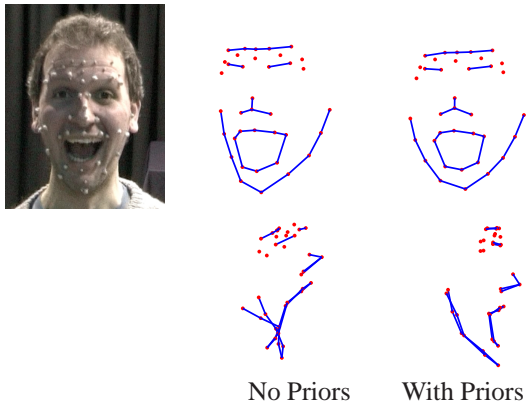


Figure 3. Comparison between the solution with and without priors. Even if the frontal view appears correct in both cases, only the solution with priors can properly estimate the shape depth.

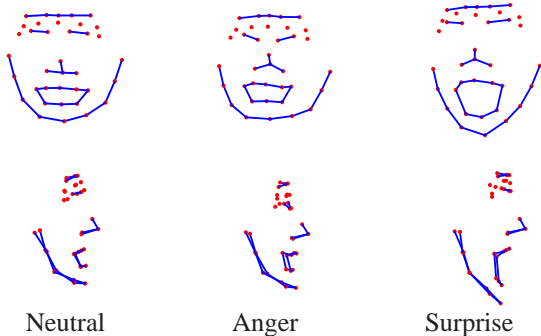


Figure 4. Front and side views of the 3D Reconstructions after non-linear optimization with shape priors of the neutral, anger and surprise facial expressions.

tured 2D and 3D locations of a set of markers overlaid on a deforming face (see Figure 1). The subject was performing tiny head pose changes which in turn affected the extraction of satisfactory 3D reconstructions from 2D trajectories. Figure 2(a) shows the first basis shape B_1 obtained using deformable non-linear optimization without shape priors. The basis shape was generally quite flat and the mouth bend inward into the head. The shape prior $L_{3 \times 37}$ is taken from a 3D reconstruction of the VICON system itself when the

subject was performing a neutral pose as shown in Figure 2(b). After performing non-linear optimization with priors, the basis shape B_1 much resemble the prior with some variations located over the temple area (Figure 2(c)).

Figure 3 shows a comparison between 3D reconstructions for the surprise expression. Notice that both the frontal views of the 3D reconstructions apparently estimated correctly the face shape. However, the inclusion of the prior is critical as shown in the side views. Given a tiny rigid motion, the ML solution alone is very ambiguous since strong variations in depth results in small displacements onto the image plane. The inclusion of a prior over the rigid shape component constrains the object depth and deformation estimates. Figure 4 presents front and side views of the final 3D reconstruction after 21 iterations of non-linear optimization. Facial symmetry is well preserved and generally the depth of the shape is correctly estimated.

In order to compute quantitatively the algorithm performances, we performed 100 trials for each test with different conditions. Gaussian noise of different levels was added to the image points while the prior accurateness was as well altered by adding gaussian noise to L . These tests were performed in order to show how much the algorithm is resilient to inaccurate priors and increasing image noise. Results are presented in Figure 5 showing that the algorithm can deliver satisfactory performances even with some degrees of inaccuracy on the shape prior.

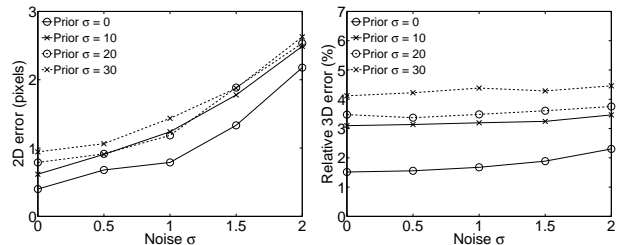


Figure 5. Synthetic experiments results for both 2D reprojection error in pixel (left) and 3D error in units (right). The image shape is approximately of size 114×84 pixels and zero-mean gaussian noise of variance $\sigma = \{0, 0.5, 1, 1.5, 2\}$ pixels is added to W . The 3D prior L is contained in a box of size $169 \times 193 \times 102$ and zero-mean gaussian noise of variance $\sigma = \{0, 10, 20, 30\}$ units is added to simulate inaccurate shape priors.

5.2. Image data with motion degeneracy

The aim of this experiment is to enforce a shape prior belonging to a subject with measurements obtained from a different subject. A rigid 3D shape of a face is first extracted from the image sequence shown in Figure 6 using a rigid factorization approach. Then, the prior is used to infer the 3D structure of a 45 frames sequence with deformations mainly localized in the mouth region. In both sequences, the 65 image points were extracted using an AAM tracker.



Figure 6. The first two images show snapshots from a brief sequence of 75 frames showing dominant rigid motion. The right image shows the 3D rigid shape prior computed from the sequence.

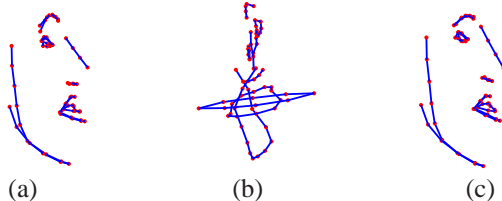
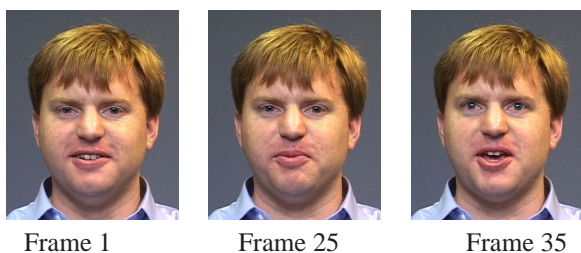


Figure 7. The first row shows three frames of a sequence with nearly planar motion and deformations mainly located in the mouth. The second row presents the 3D shape comparisons between the prior (a), the 1st basis extracted without prior (b) and the same basis computed using shape priors (c).

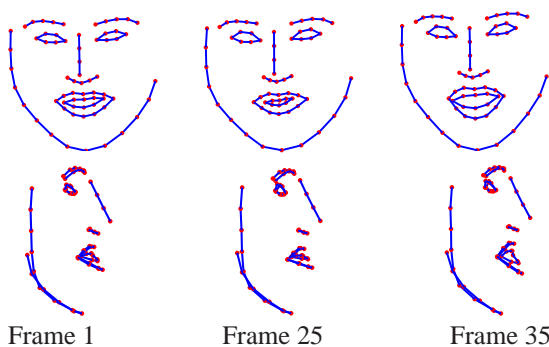


Figure 8. Front and side views after non-linear optimization. Deformations are mainly localized in the mouth region.

The second row of Figure 7 shows the priors and the difference between the computed rigid basis shape B_1 . From Figure 7(b) it is evident that the solution without priors presents depth estimates which were rather compromised. Differently, after the shape prior inclusion, depth is computed correctly and the face characteristics are adapted to

the new subject. It is possible to notice this in Figure 7(c) where the nose is more elongated and eyebrows are less bent. Finally, Figure 8 shows the 3D reconstruction for three frames after non-linear minimization.



Figure 9. The left image shows a frame from a sequence with two articulated bodies with a universal joint. The right image shows the overall 3D reconstruction which has wrong depth estimates.

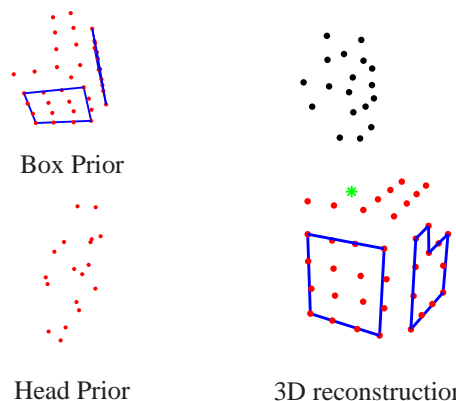


Figure 10. The two images on the left show the priors used to infer the 3D articulated structure. On the right it is shown the full reconstruction along with the estimation of the joint position (green).

5.3. Articulated shape with priors

This test was aimed to show the relevance of the inference with priors also in the case of articulated SfM. The experiment dealt with a universal joint between two shapes as shown by Figure 9. The 61 frames image sequence however contained motion degeneracies only in the box shape while the head was rotating enough to assure reliable 3D reconstructions. The shape priors were obtained from separated rigid factorizations of both objects from different sequences. Generalized factorization is then applied to obtain a better initial fit for the box shape. The overall affine structure is finally upgraded to metric with the closed-form

solution proposed in [17]. Figure 10 shows the box sides now preserving orthogonality and the depth was accurately estimated along with the position of the universal joint.

6. Conclusion

This paper presented a method capable of including shape priors in a factorization framework for SfM. These priors were in the form of previously computed 3D shapes which represented a close representation of the measured data. In such way it was possible to obtain reliable 3D reconstructions especially when the motion appearing in the image sequence was degenerate. At this end, a closed-form solution is used to compute an affine fit of motion and shape components which are then upgraded to metric using non-linear optimization with shape priors. Notice the strong relation of this approach to other methods based on factorization. For instance, in the case of photometric stereo [1], priors on the normals of the object may be directly applied using slight variations of our method. As ongoing work, in order to widen the approach applicability, it is necessary to develop an efficient algorithm for automatically matching 2D point trajectories and 3D shape priors.

Acknowledgments

This work was supported by Fundação para a Ciência e a Tecnologia (ISR/IST pluriannual funding) through the POS_Conhecimento Program (include FEDER funds) and grant PTDC/EEA-ACR/72201/2006, “MODI - 3D Models from 2D Images”. E. Muñoz, J. Xiao and P. Tresadern provided the sequences used in the experimental sections for the synthetic, deformable and articulated test respectively. Thanks to L. Agapito and X. Lladó for suggestions and for carefully reading this paper.

References

- [1] R. Basri, D. Jacobs, and I. Kemelmacher. Photometric stereo with general, unknown lighting. *International Journal of Computer Vision*, 72(3):239–257, 2007. 8
- [2] Å. Björck. *Numerical methods for least squares problems*. SIAM Philadelphia, 1996. 4
- [3] M. Brand. Morphable models from video. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, Kauai, Hawaii*, volume 2, pages 456–463, December 2001. 5
- [4] M. Brand. A direct method for 3d factorization of nonrigid motion observed in 2d. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, San Diego, California*, pages 122–128, 2005. 3
- [5] C. Bregler, A. Hertzmann, and H. Biermann. Recovering non-rigid 3d shape from image streams. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, Hilton Head, South Carolina*, pages 690–696, June 2000. 1, 3
- [6] A. Del Bue, X. Lladó, and L. Agapito. Non-rigid metric shape and motion recovery from uncalibrated images using priors. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, New York, NY*, pages 1191–1198, New York, June 2006. 1
- [7] D. Forsyth, S. Ioffe, and J. Haddon. Bayesian structure from motion. *Proc. 7th International Conference on Computer Vision, Kerkyra, Greece*, 01:660, 1999. 1
- [8] P. Hansen. *Rank-Deficient and Discrete Ill-Posed Problems: Numerical Aspects of Linear Inversion*. Society for Industrial Mathematics, 1998. 4
- [9] K. Kanatani and Y. Sugaya. Factorization without factorization: complete recipe. *Memories of the Faculty of Engineering, Okayama University*, 38(1–2):61–72, 2004. 3, 5
- [10] S. Olsen and A. Bartoli. Using priors for improving generalization in non-rigid structure-from-motion. *Proc. British Machine Vision Conference*, 2007. 1
- [11] C. J. Poelman and T. Kanade. A paraperspective factorization method for shape and motion recovery. In *Proc. 3rd European Conference on Computer Vision, Stockholm*, volume 2, pages 97–108, 1994. 1
- [12] J. Solem and F. Kahl. Surface reconstruction using learned shape models. *Advances in Neural Information Processing Systems*, 17, 2005. 1
- [13] P. Sturm and B. Triggs. A factorization based algorithm for multi-image projective structure and motion. In *Proc. 4th European Conference on Computer Vision, Cambridge*, pages 709–720, April 1996. 1
- [14] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography: A factorization approach. *International Journal of Computer Vision*, 9(2):137–154, 1992. 1, 2
- [15] L. Torresani, A. Hertzmann, and C. Bregler. Non-rigid structure-from-motion: Estimating shape and motion with hierarchical priors. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2008. 1, 3
- [16] L. Torresani, D. Yang, E. Alexander, and C. Bregler. Tracking and modeling non-rigid objects with rank constraints. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, Kauai, Hawaii*, 2001. 5
- [17] P. Tresadern and I. Reid. Articulated structure from motion by factorization. In *Proc. IEEE Conference on Computer Vision and Pattern Recognition, San Diego, California*, volume 2, pages 1110–1115, June 2005. 1, 2, 3, 5, 8
- [18] B. Triggs, P. McLauchlan, R. I. Hartley, and A. Fitzgibbon. Bundle adjustment – A modern synthesis. In W. Triggs, A. Zisserman, and R. Szeliski, editors, *Vision Algorithms: Theory and Practice*, LNCS, pages 298–375. Springer Verlag, 2000. 5
- [19] J. Xiao, J. Chai, and T. Kanade. A closed-form solution to non-rigid shape and motion recovery. *International Journal of Computer Vision*, 67(2):233–246, April 2006. 1, 3
- [20] J. Yan and M. Pollefeys. A factorization-based approach for articulated non-rigid shape, motion and kinematic chain recovery from video. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(5), May 2008. 1, 2, 3