

Factorization as a Rank 1 Problem

Pedro M. Q. Aguiar*

Instituto de Sistemas e Robótica
IST, Lisboa, Portugal
aguiar@isr.ist.utl.pt

José M. F. Moura

Carnegie Mellon University
Pittsburgh, PA, USA
moura@ece.cmu.edu

Abstract

Tomasi and Kanade [1] introduced the factorization method for recovering 3D structure from 2D video. In their formulation, the 3D shape and 3D motion are computed by using an SVD to approximate a matrix that is rank 3 in a noiseless situation. In this paper we reformulate the problem using the fact that the x and y coordinates of each feature are known from their projection onto the image plane in frame 1. We show how to compute the 3D shape, i.e., the relative depths z , and the 3D motion by a simple factorization of a matrix that is rank 1 in a noiseless situation. This allows the use of very fast algorithms even when using a large number of features and large number of frames. We also show how to accommodate confidence weights for the feature trajectories. This is done without additional computational cost by rewriting the problem as the factorization of a modified matrix.

1 Introduction

The recovery of the 3D structure (3D shape and 3D motion) from a 2D video sequence has been widely considered by the computer vision community. If no prior knowledge about the scene is available, the cue to estimating the 3D structure is the 2D motion of the brightness pattern in the image plane. For this reason, this problem is generally referred to as *structure from motion*. The two major steps in *structure from motion* are usually the following: compute the 2D motion in the image plane; and estimate the 3D shape and the 3D motion from the computed 2D motion. Early approaches to *structure from motion* processed a single pair of consecutive frames and provided existence and uniqueness results to the problem of estimating 3D motion and absolute depth from the 2D motion in the camera plane between two frames. Two-frame based algorithms are highly sensitive to image noise, and, when the object is far from the camera, i.e., at a large distance when compared to the object depth, they fail

even at low level image noise. More recent research has been oriented towards the use of longer image sequences, leading to filtering algorithms that integrate along time a set of two-frame depth estimates.

Tomasi and Kanade [1] introduced an elegant method to recover rigid structure from an image sequence. Instead of representing the 3D positions of feature points by their image coordinates and their depths, they adopt Ullman's original formulation of the *structure from motion* problem [2]. In [1], as in [2], the 3D positions of the feature points are expressed in terms of cartesian coordinates in a world-centered coordinate system, and the images are modeled as orthographic projections. In [1], the 2D projection of each feature point is tracked along the image sequence. The 3D shape and motion are then estimated by factorizing a measurement matrix whose entries are the set of trajectories of the feature point projections. This work was later extended to the scaled-orthography and paraperspective projections [3, 4].

In this paper, we adopt the scenario of [1] and [2], i.e., we also consider a world-centered coordinate system and the orthographic projection model. In our formulation the unknowns are the 3D motion and the relative depths of the set of features, not their 3D positions as considered in [1]. The coordinates of the features along the camera plane are given by their image positions in the first frame. The knowledge of the x and y coordinates of the features enables us to solve the *structure from motion* problem by factorizing a rank 1 matrix instead of a rank 3 matrix as in [1]. This simplifies the decomposition and normalization stages involved in the factorization approach. Our approach is easily extended to the scaled-orthography and paraperspective projections.

We avoid the computation of the *Singular Value Decomposition* (SVD) by using a fast iterative method to compute the rank 1 matrix that best matches the data. Reference [5] introduced a recursive formulation for the original rank 3 factorization that also does not

*The work of P. Aguiar was partially supported by INVOTAN.

use SVD. The method in [5] stores and updates a matrix that becomes very large with the increasing of the number of features. In our implementation, it is not necessary to store or compute any other matrix than the matrix to be factorized.

In the paper we also show how to include confidence weights for the feature trajectories, i.e., how to weight more the trajectory of a "sharp" feature than the trajectory of a "smooth" feature. Reference [3] proposes an iterative method to solve a more general problem, where the weights are time-variant. As reported in [3] the iterative method may fail to converge. We show that when the weights are time-invariant, the problem is rewritten as the non-weighted factorization of a modified matrix. Then, any method can be used to factorize this matrix.

In [6] we extended the factorization approach to recover 3D structure from a sequence of optical flow parameters. Instead of tracking pointwise features, we track regions where the optical flow is described by a single set of parameters. The reliability of the estimates of the optical flow parameters depends not only on the spatial variability of the brightness pattern but also on the size of the region. The weighted factorization method we describe below is general enough to be applied to that framework.

Paper overview In section 2 we formulate the *structure from motion* problem. Section 3 shows how to recover structure from motion by factorizing a rank 1 matrix. In section 4 we extend the approach to accommodate confidence weights. The algorithm to compute the best rank 1 approximation is described in section 5. Experimental results are in section 6. Section 7 concludes the paper.

2 Problem Formulation

We consider the same scenario of [1]. Attach a coordinate system to the object. The object coordinate system has axes labeled by (x, y, z) . The shape of the object is described by the 3D position of a set of N feature points. The 3D position of feature n is expressed in terms of the object coordinate system by (x_n, y_n, z_n) . To describe the 3D motion of the object, we attach a different coordinate system to the camera. The camera coordinate system has axes labeled by (u, v, w) . The 3D motion of the object is described by the position of the object coordinate system with respect to the camera coordinate system.

A set of N feature points are tracked along an image sequence of F frames. Under orthography, the projection of feature n in frame f , denoted by (u_{fn}, v_{fn}) , is

$$u_{fn} = i_{xf}x_n + i_{yf}y_n + i_{zf}z_n + t_{uf} \quad (1)$$

$$v_{fn} = j_{xf}x_n + j_{yf}y_n + j_{zf}z_n + t_{vf} \quad (2)$$

where $i_{xf}, i_{yf}, i_{zf}, j_{xf}, j_{yf},$ and j_{zf} are entries of the well known 3D rotation matrix, and t_{uf} and t_{vf} are the components of the object translation along the camera plane. We make the object coordinate system and camera coordinate system coincide in the first frame, so we have

$$u_{1n} = x_n \quad \text{and} \quad v_{1n} = y_n. \quad (3)$$

The coordinates of the feature points along the camera plane $\{x_n, y_n, 1 \leq n \leq N\}$ are given by (3). We formulate the *structure from motion* problem as solving the overconstrained equation system (1,2) with respect to the following set of unknowns: the 3D positions of the object for $2 \leq f \leq F$, and the relative depths $\{z_n, 1 \leq n \leq N\}$.

By choosing the origin of the object coordinate system to coincide with the centroid of the set of feature points, we get the estimate for the translation as the centroid of the feature point projections

$$\hat{t}_{uf} = \frac{1}{N} \sum_{n=1}^N u_{fn} \quad \text{and} \quad \hat{t}_{vf} = \frac{1}{N} \sum_{n=1}^N v_{fn}. \quad (4)$$

Replacing the translation estimates in equation system (1,2), and defining

$$\tilde{u}_{fn} = u_{fn} - \frac{1}{N} \sum_{n=1}^N u_{fn}, \quad \tilde{v}_{fn} = v_{fn} - \frac{1}{N} \sum_{n=1}^N v_{fn}, \quad (5)$$

$$\mathbf{R} = \begin{bmatrix} \tilde{u}_{21} & \tilde{u}_{22} & \cdots & \tilde{u}_{2N} \\ \cdots & \cdots & \cdots & \cdots \\ \tilde{u}_{F1} & \tilde{u}_{F2} & \cdots & \tilde{u}_{FN} \\ \tilde{v}_{21} & \tilde{v}_{22} & \cdots & \tilde{v}_{2N} \\ \cdots & \cdots & \cdots & \cdots \\ \tilde{v}_{F1} & \tilde{v}_{F2} & \cdots & \tilde{v}_{FN} \end{bmatrix}, \quad (6)$$

$$\mathbf{M} = \begin{bmatrix} i_{x2} & \cdots & i_{xF} & j_{x2} & \cdots & j_{xF} \\ i_{y2} & \cdots & i_{yF} & j_{y2} & \cdots & j_{yF} \\ i_{z2} & \cdots & i_{zF} & j_{z2} & \cdots & j_{zF} \end{bmatrix}^T, \quad (7)$$

$$\text{and} \quad \mathbf{S}^T = \begin{bmatrix} x_1 & x_2 & \cdots & x_N \\ y_1 & y_2 & \cdots & y_N \\ z_1 & z_2 & \cdots & z_N \end{bmatrix}, \quad (8)$$

we rewrite (1,2) in matrix format as

$$\mathbf{R} = \mathbf{M}\mathbf{S}^T. \quad (9)$$

Matrix \mathbf{R} is $2(F-1) \times N$ but it is rank deficient. In a noiseless situation, \mathbf{R} is rank 3 reflecting the high redundancy in the data, due to the 3D rigidity of the object. The relation expressed in matrix format as in (9) was first introduced in [1]. In our formulation, the rows corresponding to frame 1 are not included in \mathbf{R} and \mathbf{M} , and we use the fact that the first two rows of \mathbf{S}^T are known.

3 Rank 1 Factorization

The factorization approach [1] finds a suboptimal solution to the problem

$$\min_{\mathbf{M}, \mathbf{S}} \left\| \mathbf{R} - \mathbf{M}\mathbf{S}^T \right\|_F \quad (10)$$

where the solution space is constrained by the orthonormality of the rows of the matrix \mathbf{M} (7). $\|\cdot\|_F$ denotes the Frobenius norm. In [1], the nonlinear minimization above was solved in two stages. The first stage, *decomposition stage*, solves $\mathbf{R} = \mathbf{M}\mathbf{S}^T$ in the *Least Squares* (LS) sense by computing the SVD of the matrix \mathbf{R} and selecting the 3 largest singular values. From $\mathbf{R} \simeq \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, the solution is $\mathbf{M} = \mathbf{U}\mathbf{\Sigma}^{\frac{1}{2}}\mathbf{A}$ and $\mathbf{S}^T = \mathbf{A}^{-1}\mathbf{\Sigma}^{\frac{1}{2}}\mathbf{V}^T$ where \mathbf{A} is a non-singular 3×3 matrix. The second stage, *normalization stage*, computes \mathbf{A} by approximating the constraints imposed by the structure of the matrix \mathbf{M} .

Our formulation takes advantage of the fact that the first two rows of \mathbf{S} are known ($\{x_n, y_n\}$ are known from the position of the features in the first frame). The problem is now reduced to, given \mathbf{R} , compute \mathbf{M} and $\{z_n\}$. This problem, although nonlinear, has a specific structure: it is a bilinear constrained LS problem. The bilinear relation comes from (9) and the constraints are imposed by the orthonormality of the rows of the matrix \mathbf{M} (7). We also perform in sequence the decomposition stage (that solves the unconstrained bilinear problem) and the normalization stage.

Decomposition Because the first two rows of \mathbf{S}^T are known, we show that the unconstrained bilinear problem $\mathbf{R} = \mathbf{M}\mathbf{S}^T$ is solved by the factorization of a rank 1 matrix, rather than a rank 3 matrix like in [1]. Define $\mathbf{M} = [\mathbf{M}_0, \mathbf{m}_3]$ and $\mathbf{S} = [\mathbf{S}_0, \mathbf{z}]$. \mathbf{M}_0 and \mathbf{S}_0 contain the first two columns of \mathbf{M} and \mathbf{S} , respectively, \mathbf{m}_3 is the third column of \mathbf{M} , and \mathbf{z} is the third column of \mathbf{S} . We decompose the relative depth vector \mathbf{z} into the component that belongs to the space spanned by the columns of \mathbf{S}_0 and the component orthogonal to this space as

$$\mathbf{z} = \mathbf{S}_0\mathbf{b} + \mathbf{a}, \quad \text{with } \mathbf{a}^T\mathbf{S}_0 = [0 \ 0]. \quad (11)$$

We rewrite \mathbf{R} by inserting (11) in (9), as

$$\mathbf{R} = \mathbf{M}_0\mathbf{S}_0^T + \mathbf{m}_3\mathbf{b}^T\mathbf{S}_0^T + \mathbf{m}_3\mathbf{a}^T. \quad (12)$$

The decomposition stage is formulated as

$$\min_{\mathbf{M}_0, \mathbf{m}_3, \mathbf{b}, \mathbf{a}} \left\| \mathbf{R} - \mathbf{M}_0\mathbf{S}_0^T - \mathbf{m}_3\mathbf{b}^T\mathbf{S}_0^T - \mathbf{m}_3\mathbf{a}^T \right\|_F. \quad (13)$$

Since we know \mathbf{S}_0 , we eliminate the dependence of the expression above on \mathbf{M}_0 by solving the linear LS for \mathbf{M}_0 in terms of the other variables. We get

$$\widehat{\mathbf{M}}_0 = \mathbf{R}\mathbf{S}_0 \left(\mathbf{S}_0^T\mathbf{S}_0 \right)^{-1} - \mathbf{m}_3\mathbf{b}^T \quad (14)$$

where we used the Moore-Penrose pseudoinverse and the orthogonality between \mathbf{a} and \mathbf{S}_0 . By replacing $\widehat{\mathbf{M}}_0$ in (13), we get

$$\min_{\mathbf{m}_3, \mathbf{a}} \left\| \widetilde{\mathbf{R}} - \mathbf{m}_3\mathbf{a}^T \right\|_F, \quad (15)$$

$$\text{where } \widetilde{\mathbf{R}} = \mathbf{R} \left[\mathbf{I} - \mathbf{S}_0 \left(\mathbf{S}_0^T\mathbf{S}_0 \right)^{-1} \mathbf{S}_0^T \right]. \quad (16)$$

We see that the decomposition stage does not determine the vector \mathbf{b} . This is because the component of \mathbf{z} that lives in the space spanned by the columns of \mathbf{S}_0 does not affect the space spanned by the columns of the entire matrix \mathbf{S} and the decomposition stage restricts only this last space.

The solution for \mathbf{m}_3 and \mathbf{a} is given by the rank 1 matrix that best approximates $\widetilde{\mathbf{R}}$. In a noiseless situation, $\widetilde{\mathbf{R}}$ is rank 1 (we get $\widetilde{\mathbf{R}} = \mathbf{m}_3\mathbf{a}^T$ by replacing (12) in (16)). By computing the largest singular value of $\widetilde{\mathbf{R}}$ and the associated singular vectors, we get

$$\widetilde{\mathbf{R}} \simeq \mathbf{u}\sigma\mathbf{v}^T, \quad \widehat{\mathbf{m}}_3 = \alpha\mathbf{u}, \quad \widehat{\mathbf{a}}^T = \frac{\sigma}{\alpha}\mathbf{v}^T \quad (17)$$

where α is a normalizing scalar different from 0. To compute \mathbf{u} , σ , and \mathbf{v} we could use an SVD, but the rank deficiency of $\widetilde{\mathbf{R}}$ enables the use of less expensive algorithms to compute \mathbf{u} , σ , and \mathbf{v} , as detailed in section 5. This makes our decomposition stage simpler than the one in [1]. In fact, $\widetilde{\mathbf{R}}$ in (16) is \mathbf{R} multiplied by the orthogonal projector onto the orthogonal complement of the space spanned by the columns of \mathbf{S}_0 . This projection reduces the rank of the problem from 3 (matrix \mathbf{R}) to 1 (matrix $\widetilde{\mathbf{R}}$).

Normalization In this stage, we compute α and \mathbf{b} by imposing the constraints that come from the structure of \mathbf{M} . By replacing $\widehat{\mathbf{m}}_3$ in (14), we get

$$\widehat{\mathbf{M}} = \begin{bmatrix} \widehat{\mathbf{M}}_0 & \widehat{\mathbf{m}}_3 \end{bmatrix} = \mathbf{N} \begin{bmatrix} \mathbf{I}_{2 \times 2} & \mathbf{0}_{2 \times 1} \\ -\alpha\mathbf{b}^T & \alpha \end{bmatrix}, \quad (18)$$

$$\text{where } \mathbf{N} = \begin{bmatrix} \mathbf{R}\mathbf{S}_0 \left(\mathbf{S}_0^T\mathbf{S}_0 \right)^{-1} & \mathbf{u} \end{bmatrix}. \quad (19)$$

The constraints imposed by the structure of \mathbf{M} are the unit norm of each row and the orthogonality between row j and row $j + F - 1$. In terms of \mathbf{N} , α , and \mathbf{b} , the constraints are

$$\mathbf{n}_i^T \begin{bmatrix} \mathbf{I}_{2 \times 2} & -\alpha\mathbf{b} \\ -\alpha\mathbf{b}^T & \alpha^2(1 + \mathbf{b}^T\mathbf{b}) \end{bmatrix} \mathbf{n}_i = 1, \quad (20)$$

$$\mathbf{n}_j^T \begin{bmatrix} \mathbf{I}_{2 \times 2} & -\alpha \mathbf{b} \\ -\alpha \mathbf{b}^T & \alpha^2(1 + \mathbf{b}^T \mathbf{b}) \end{bmatrix} \mathbf{n}_{j+F-1} = 0 \quad (21)$$

where \mathbf{n}_i^T denotes the row i of matrix \mathbf{N} . We compute α and \mathbf{b} from the linear LS solution of the system above in a similar way to the one described in [1]. The normalization stage is also simpler than the one in [1] because the number of unknowns is 3 (α and $\mathbf{b} = [b_1, b_2]^T$) as opposed to the 9 entries of a generic 3×3 normalization matrix.

4 Weighted Factorization

Usually, feature tracking is achieved by matching the intensity pattern of a small block around each feature point across the frame sequence. The accuracy of the matching depends on the spatial variability of the intensity pattern around each feature. The model in the previous section, as well as the model in reference [1], weights equally the contribution of each feature to the final shape and motion estimates. A more robust estimate would weight more a trajectory corresponding to a “sharp” feature than a trajectory corresponding to a feature with more smooth texture. In this section, we introduce such a model and show that it leads to the factorization of a modified matrix. The computation of the weighted estimates is then done by using the method of the previous section, i.e., without additional computational cost.

The derivations in this section are valid for a more general factorization framework like the one we presented in [6]. In that paper, we parameterize the 3D shape and recover the 3D shape parameters and 3D motion parameters from the induced optical flow field. This is done by using the factorization of a matrix that collects optical flow parameters over a set of frames. The reliability of the estimates of the optical flow parameters depends on the size of the image region of analysis and on the spatial variability of the brightness pattern. The weighted factorization method we describe below is general enough to be applied to that framework.

We consider the model in expressions (1,2). Each component of the trajectory of the feature n is observed with additive gaussian white noise of variance σ_n^2 . By choosing the origin of the object coordinate system in such a way that $\sum_n \frac{x_n}{\sigma_n^2} = \sum_n \frac{y_n}{\sigma_n^2} = \sum_n \frac{z_n}{\sigma_n^2} = 0$, we get the estimates

$$\hat{t}_{uf} = \frac{\sum_{n=1}^N \frac{u_{fn}}{\sigma_n^2}}{\sum_{n=1}^N \frac{1}{\sigma_n^2}} \quad \text{and} \quad \hat{t}_{vf} = \frac{\sum_{n=1}^N \frac{v_{fn}}{\sigma_n^2}}{\sum_{n=1}^N \frac{1}{\sigma_n^2}} \quad (22)$$

for the translation along the camera plane. Expressions (22) generalize (4). Replacing the translation

estimates in equation system (1,2), and defining

$$\tilde{u}_{fn} = u_{fn} - \frac{\sum_{n=1}^N \frac{u_{fn}}{\sigma_n^2}}{\sum_{n=1}^N \frac{1}{\sigma_n^2}}, \quad \tilde{v}_{fn} = v_{fn} - \frac{\sum_{n=1}^N \frac{v_{fn}}{\sigma_n^2}}{\sum_{n=1}^N \frac{1}{\sigma_n^2}} \quad (23)$$

and the matrices \mathbf{R} , \mathbf{M} , and \mathbf{S} as in (6,7,8), we obtain, as before, $\mathbf{R} = \mathbf{M}\mathbf{S}^T$.

To take into account the different variances $\{\sigma_n^2\}$, we define the $2(F-1) \times N$ matrix $\mathbf{\Sigma}$ as

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1^{-1} & \sigma_2^{-1} & \cdots & \sigma_N^{-1} \\ \sigma_1^{-1} & \sigma_2^{-1} & \cdots & \sigma_N^{-1} \\ \cdots & \cdots & \cdots & \cdots \\ \sigma_1^{-1} & \sigma_2^{-1} & \cdots & \sigma_N^{-1} \end{bmatrix} \quad (24)$$

where each entry of $\mathbf{\Sigma}$ represents the weight to be given to each entry of matrix \mathbf{R} . We formulate the weighted factorization by writing the *Maximum Likelihood* estimate. This estimate generalizes (10) as

$$\min_{\mathbf{M}, \mathbf{S}} \left\| \left(\mathbf{R} - \mathbf{M}\mathbf{S}^T \right) \odot \mathbf{\Sigma} \right\|_F \quad (25)$$

where \odot denotes the Hadamard product. Due to the structure of $\mathbf{\Sigma}$, we rewrite (25) as the factorization of a modified matrix \mathbf{R}_w ,

$$\min_{\mathbf{M}, \mathbf{S}_w} \left\| \mathbf{R}_w - \mathbf{M}\mathbf{S}_w^T \right\|_F \quad (26)$$

$$\text{where} \quad \mathbf{R}_w = \mathbf{R}\mathbf{W}, \quad \mathbf{S}_w = \mathbf{W}\mathbf{S}, \quad (27)$$

$$\mathbf{W} = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_N}\right) = [w_{ij}], \quad w_{ij} = \frac{1}{\sigma_i} \delta_{i-j}. \quad (28)$$

The factorization in (26) is similar to the one studied in the previous section. Note that \mathbf{R}_w and the first two rows of \mathbf{S}_w^T are known from (27). We minimize (26) by using the procedure described in the previous section to compute the estimates $\widehat{\mathbf{M}}$ and $\widehat{\mathbf{S}}_w$. The estimate $\widehat{\mathbf{S}}$ is obtained from $\widehat{\mathbf{S}}_w$ by $\widehat{\mathbf{S}} = \mathbf{W}^{-1}\widehat{\mathbf{S}}_w$.

Reference [3] also considers the reliability weights in estimating \mathbf{M} and \mathbf{S} . In [3], the author allows the matrix $\mathbf{\Sigma}$ to have a general structure, i.e., one can give each feature a weight that varies along the time. The solution is found by an iterative process that, as reported in [3], may fail to converge. In our formulation, this is not the case because we restrict the weight matrix $\mathbf{\Sigma}$ to have the structure of (24). For a general matrix $\mathbf{\Sigma}$, it is not possible to write the minimization (25) in the form of a factorization (26). In fact, the unconstrained bilinear weighted LS problem (25) has a single global minimum, up to a scale factor, when $\mathbf{\Sigma}$ has the rows all equal or the columns all equal. It can be shown that this is not true for a generic matrix $\mathbf{\Sigma}$. In this case, the existence local minima makes nontrivial the using of iterative numerical techniques.

5 Singular Value Computations

This section describes how we compute the rank 1 approximation of a given matrix $\tilde{\mathbf{R}}$. The problem is

$$\min_{\mathbf{u}, \mathbf{v}} \left\| \tilde{\mathbf{R}} - \mathbf{u}\mathbf{v}^T \right\|_F \quad (29)$$

which comes from expression (17) by including the scaling factor σ into the vector \mathbf{v} . The solution to (29) is known to be given by the SVD of $\tilde{\mathbf{R}}$ after selecting the largest singular value.

The rank deficiency of $\tilde{\mathbf{R}}$ enables the use of a less expansive iterative technique to compute the decomposition. It is based on the fact that only the singular vector \mathbf{v} that corresponds to the largest singular value has to be computed. Since the vector \mathbf{v} is the eigenvector of $\tilde{\mathbf{R}}^T \tilde{\mathbf{R}}$ that corresponds to the largest eigenvalue, we start with a random choice \mathbf{v}_0 and iterate,

$$\mathbf{v}_{i+1} = \frac{\mathbf{v}_i^T \mathbf{v}_i}{\mathbf{v}_i^T \tilde{\mathbf{R}}^T \tilde{\mathbf{R}} \mathbf{v}_i} \tilde{\mathbf{R}}^T \tilde{\mathbf{R}} \mathbf{v}_i \quad (30)$$

until convergence, see [7]. In each iteration, the component of \mathbf{v}_i along the vector \mathbf{v} is more magnified than the components along the other eigenvectors of $\tilde{\mathbf{R}}^T \tilde{\mathbf{R}}$. The fraction in (30) is a normalizing factor. The vector \mathbf{u} is then computed by solving (29) with \mathbf{v} given by the final value of the iterative process above,

$$\mathbf{u} = \frac{\tilde{\mathbf{R}} \mathbf{v}}{\mathbf{v}^T \mathbf{v}}. \quad (31)$$

This iterative procedure can be generalized to compute the best rank r approximation of a given matrix, for $r > 1$. This was done in reference [5] where the authors propose a recursive formulation to the original rank 3 factorization. Their method stores and updates the $N \times N$ matrix $\mathbf{R}^T \mathbf{R}$. Then, they use this matrix to compute, iteratively, the best rank 3 approximation of \mathbf{R} . When the number of features is large, the matrix $\mathbf{R}^T \mathbf{R}$ becomes very large. In our implementation of the decomposition stage, instead of computing $\tilde{\mathbf{R}}^T \tilde{\mathbf{R}}$, we split the computation of each iteration (30) by first computing $\tilde{\mathbf{R}} \mathbf{v}_i$,

$$\tilde{\mathbf{v}}_i = \tilde{\mathbf{R}} \mathbf{v}_i, \quad \mathbf{v}_{i+1} = \frac{\mathbf{v}_i^T \mathbf{v}_i}{\tilde{\mathbf{v}}_i^T \tilde{\mathbf{v}}_i} \tilde{\mathbf{R}}^T \tilde{\mathbf{v}}_i. \quad (32)$$

6 Experiments

We describe three experiments. The first experiment uses synthetic data to illustrate the properties of the rank 1 matrix $\tilde{\mathbf{R}}$. The second experiment compares the computational cost of our approach with the

one of the original factorization approach. In the third experiment, we recover the 3D shape and 3D motion from a real video sequence.

Experiment 1 This experiment illustrates the properties of the rank 1 matrix $\tilde{\mathbf{R}}$ by using synthetic data. We generated a set of 10 feature points randomly located inside a cube. The 3D rotational motion was simulated by synthesizing a smooth time evolution for the Euler angles. We used the perspective projection model to project the features onto the image plane. The lens focal length parameter was set to a value high enough such that the orthographic projection can be considered a valid approximation. Figure 1 shows the feature trajectories on the image plane obtained for a set of 50 frames, after adding noise. For each trajectory, the initial position is marked with 'o' and the final position is marked with '*'. The trajectories in figure 1 are a representation of the columns of matrix \mathbf{R} . The trajectory for feature n shows the n^{th} column of \mathbf{R} , i.e., it is the evolution of the image point $(\mathbf{R}_{f,n}, \mathbf{R}_{F-1+f,n})$ across the frame index f , see expression (6). The challenge in the *structure from motion* problem comes from the fact that the 3D shape and the 3D motion are observed in a coupled way through the 2D motion on the image plane (the trajectories in figure 1).

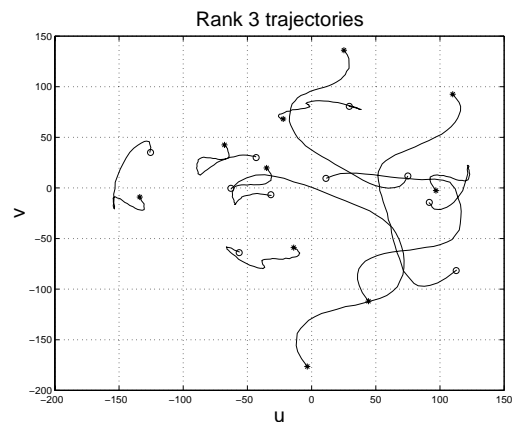


Figure 1: Feature trajectories on image plane.

From the data in matrix \mathbf{R} , we computed matrix $\tilde{\mathbf{R}}$ given by expression (16). Figure 2 plots the columns of $\tilde{\mathbf{R}}$ in the same way as figure 1 plots \mathbf{R} , i.e., it shows the evolution of $(\tilde{\mathbf{R}}_{f,n}, \tilde{\mathbf{R}}_{F-1+f,n})$ across the frame index f , for each feature n . We see that all trajectories in figure 2 have equal shape, unlike the ones in figure 1. This is because we have eliminated the dependence of the trajectories on the x and y coordinates of the features, by making the subspace projection (16). Each trajectory in figure 2 is a scaled version of a fixed

trajectory that does not depend on the object shape. This fixed trajectory is determined uniquely by the 3D motion of the object; it corresponds to the third column of matrix \mathbf{M} , \mathbf{m}_3 . The scaling factor for each trajectory in figure 2 is the relative depth z of the corresponding feature point. This way we decouple the influence on the trajectories of the 3D motion from the influence of the 3D shape.

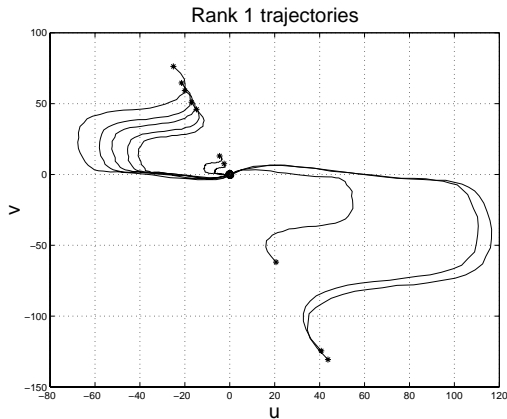


Figure 2: Feature trajectories from matrix $\tilde{\mathbf{R}}$.

Figure 3 plots the 10 larger singular values of matrices \mathbf{R} , marked with 'o', and $\tilde{\mathbf{R}}$, marked with '*'. While the 3 larger singular values of \mathbf{R} contain the most of the energy, $\tilde{\mathbf{R}}$ is well described by only the largest singular value. The plots in figure 3 confirm the analysis of sections 2 and 3 and the comment in the previous paragraph.

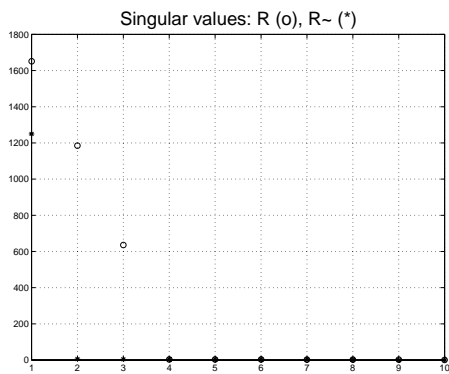


Figure 3: Singular values of matrices \mathbf{R} and $\tilde{\mathbf{R}}$.

Experiment 2 To compare the computational cost of the algorithms, we ran the experiment described for a fixed number of $F = 50$ frames and a number of N feature points varying from 10 to 100; and for a fixed number of $N = 50$ feature points and a number of F frames varying from 10 to 100. We computed the

average number of MATLAB floating point operations (FLOPS) over 1000 tests for each experiment.

For each experiment, we estimate the 3D shape and 3D motion by using three methods: i) the original factorization method [1] that computes the SVD, ii) the same method but computing the factorization of the rank 3 matrix \mathbf{R} by using an iterative procedure similar to the one described in section 5, expression (32), and iii) our formulation of the factorization as a rank 1 problem. The reason why we include method ii) in the experiment is because it is the fastest way to compute the rank 3 factorization, making fair the comparison with the method iii). It is worth mentioning that the method ii) is not the recursive formulation presented in [5]. The method in [5] is well suited for a fixed small number of points and an increasing number of frames. On the other hand, with the increasing of the number of feature points, it becomes useless because it uses a matrix that becomes prohibitively large. As discussed in section 5, the storage space and computational power required by the method [5] are then much higher than the ones required by method ii) above. Figures 4 and 5 plot the average number of FLOPS as a function of the number of frames and the number of feature points, for each of the three methods. The number of FLOPS are marked with dotted lines for method i), dashdotted lines for method ii), and solid lines for method iii). The left plots show the three curves, while the right plots show only the curves for methods ii) and iii) using a different vertical scale, for better visualization.

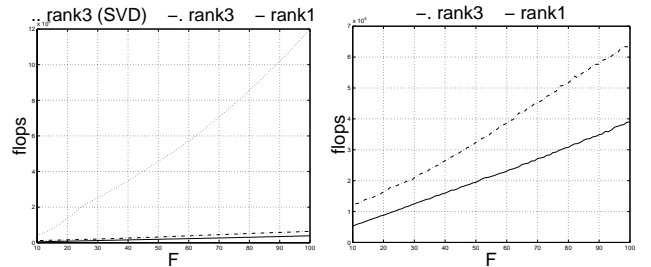


Figure 4: MATLAB floating point operations count.

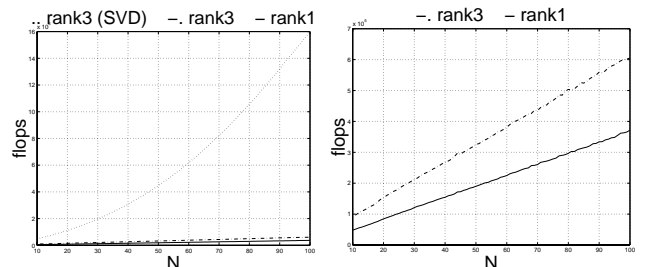


Figure 5: MATLAB floating point operations count.

From the left side plots, we see that that the number of FLOPS is much larger for the original factorization method than for the iterative approaches. This is due to the high computational cost of the SVD. From the right side plots, we see that the number of FLOPS increases approximately linearly with both the number of frames and the number of feature points, for both iterative methods ii) and iii). The increasing rate is lower for the factorization of the rank 1 matrix $\tilde{\mathbf{R}}$ than the rank 3 matrix \mathbf{R} . This is because both the decomposition and normalization stages in method iii) are simpler than the ones in method ii). In all the experiments, the performance of the three methods in terms of the accuracy of the estimates of the 3D structure are indistinct. However, it must be pointed out that the rank 1 factorization results can suffer a slight degradation if the x and y coordinates of the feature points are very inaccurate in frame 1.

Experiment 3 The first frame of a video sequence of 30 frames is on the left side of figure 6. It shows a building. Two planar walls meet along a smooth (round) edge. We marked with white squares 100 selected features. On the right side of figure 6, we represent the "trackability" of the feature candidates. The brighter is a point, the most reliable is the tracking. The method used to compute this image is described elsewhere. We choose the feature points by selecting the peaks of this image. We assign to each feature the confidence weight given by the value of the corresponding peak.

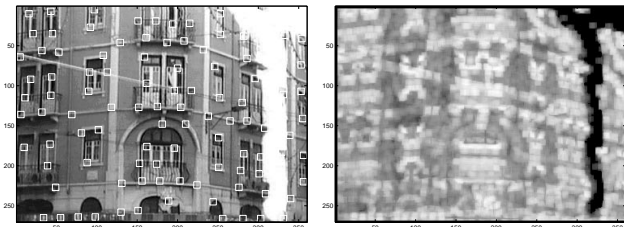


Figure 6: First frame and feature weights.

We tracked the feature points by matching the intensity pattern of each feature along the sequence. Using the rank 1 weighted factorization we recovered the 3D motion and the relative depth of the feature points from the set of feature trajectories, as described in sections 3 and 4. Figure 7 show two perspective views of the reconstructed 3D shape with the texture mapped on it. The angle between the walls is clearly seen and the round edge is also reconstructed.

7 Conclusion

We reformulate the factorization approach to the *structure from motion* problem. Since the coordinates

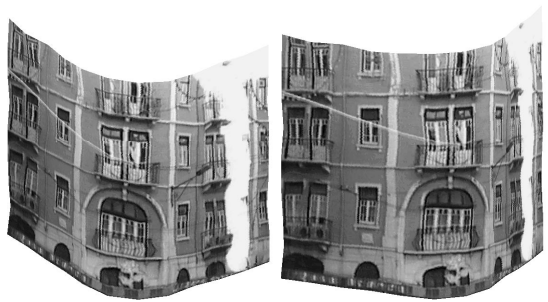


Figure 7: Reconstructed 3D shape and texture.

of the feature points along the camera plane are given by the coordinates of their projection onto the image plane in the first frame, we recover the 3D structure by estimating the relative depths of the feature points and the 3D motion. We show that this leads to the factorization of a matrix that is rank 1 in a noiseless situation. We use an efficient iterative algorithm to compute the factorization. This algorithm is very fast, even when processing a large number of features and large number of frames.

The paper also describes how to accommodate different weighting factors for the feature trajectories. This is done without additional computational cost by rewriting the problem as the factorization of a modified matrix.

References

- [1] C. Tomasi and T. Kanade, "Shape and motion from image streams under orthography: a factorization method," *IJCV*, vol. 9, no. 2, 1992.
- [2] S. Ullman, *The Interpretation of Visual Motion*, MIT Press, Cambridge, MA, USA, 1979.
- [3] C. J. Poelman, *A Paraperspective Factorization Method For Shape and Motion Recovery*, Ph.D. thesis, Carnegie Mellon University, USA, 1995.
- [4] C. Poelman and T. Kanade, "A paraperspective factorization method for shape and motion recovery," *IEEE Trans. on PAMI*, vol. 19, no. 3, 1997.
- [5] T. Morita and T. Kanade, "A sequential factorization method for recovering shape and motion from image streams," *PAMI*, vol. 19, no. 8, 1997.
- [6] Pedro M. Q. Aguiar and José M. F. Moura, "Video representation via 3D shaped mosaics," in *IEEE ICIP*, Chicago, USA, October 1998.
- [7] G. H. Golub and C. F. Van Loan, *Matrix Computations*, The Johns Hopkins Univ. Press, 1989.