

A Multi-body Factorization Method for Motion Analysis

João Costeira * Takeo Kanade
School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213

Abstract

The structure-from-motion problem has been extensively studied in the field of computer vision. Yet, the bulk of the existing work assumes that the scene contains only a single moving object. The more realistic case where an unknown number of objects move in the scene has received little attention, especially for its theoretical treatment. In this paper we present a new method for separating and recovering the motion and shape of multiple independently moving objects in a sequence of images. The method does not require prior knowledge of the number of objects, nor is dependent on any grouping of features into an object at the image level. For this purpose, we introduce a mathematical construct of object shapes, called the shape interaction matrix, which is invariant to both the object motions and the selection of coordinate systems. This invariant structure is computable solely from the observed trajectories of image features without grouping them into individual objects. Once the structure is computed, it allows for segmenting features into objects by the process of transforming it into a canonical form, as well as recovering the shape and motion of each object.

1 Introduction

A motion image sequence allows for the recovery of the three-dimensional structure of a scene. While a large amount of literature exists about this structure-from-motion problem, most previous theoretical work is based on the assumption that only a single motion is included in the image sequence; either the environment is static and the observer moves, or the observer is static and only one object in the scene is moving. More difficult and less studied is the general case of an unknown number of objects moving independently. Suppose that a set of features has been extracted and tracked in an image sequence, but it is not known which feature belongs to which object. Given a set of such feature trajectories, the question is whether we can segment and recover the motion and shape of multiple objects contained in the image sequence.

The previous approaches to the structure-from-motion problem for multiple objects can be grouped into two classes: image motion-based (2D) and three-dimensional (3D) modeling. The image-motion based approach relies mostly on spatio-temporal properties

of an image sequence. For example, regions corresponding to different velocity fields are extracted by using Fourier domain analysis [1] or scale-space and space-time filters [2, 6, 7]. These image-based methods have limited applicability either because object motions are restricted to a certain type, such as translation only, or because image-level properties, such as locality, need to be used for segmentation without assuring consistent segmentation into 3D objects.

To overcome these limitations, models of motion and scene can be introduced which provide more constraints. Representative constraints include rigidity of an object [12] and smoothness (or similarity) of motion [10, 3]. Then the problem becomes segmenting image events, such as feature trajectories, into objects so that the recovered motion and shape satisfy those constraints. It is now a clustering problem with constraints derived from a physical model. Though sound in theory, the practical difficulty is the cyclic dilemma: to check the constraints it is necessary to segment features and to segment it is necessary to compute constraints. So, developed methods tend to be of a "generate-and-test" nature, or require prior knowledge of the number of objects (clusters). Ullman [12] describes a computational scheme to recursively recover shape from the tracks of image features. A model of the object's shape is matched to the current position of the features, and a new model that maximizes rigidity is computed to update the shape. He suggests that this scheme could be used to segment multi-body scenes by local application of the rigidity principle. Since a single rigid body model does not fit the whole data, collections of points that could be explained by a rigid transformation would be searched and grouped into an object. Under the framework of the factorization method [11], this view of the problem is followed by Boulton and Brown [3] and Gear [5], where the role of rigidity is replaced by linear dependence between feature tracks. Since the factorization produces a matrix that is related with shape, segmentation is obtained by recursively clustering columns of feature trajectories into linearly dependent groups.

This paper presents a new method for segmenting and recovering the motion and shape of multiple independently moving objects from a set of feature trajectories tracked in a sequence of images. Developed by using the framework of the factorization by Tomasi and Kanade [11], the method does not require any grouping of features into an object at the

*Also affiliated with Instituto Sup. Técnico-ISR. Partially funded by JNICT-Portugal

image level or prior knowledge of the number of objects. It directly computes shape information and allows segmentation into objects. This has been made possible by introducing a linear-algebraic construct of object shapes, called the shape interaction matrix. The entries of this matrix are invariant to individual object motions and yet is computable only from tracked feature trajectories without knowing their object identities (ie, segmentation). Once the matrix is computed, transforming it into the canonical form results in segmenting features as well as recovering the shape and motion of each object. We will present our theory by using the orthographic camera model. It is, however, easily seen that the theory, and thus the method, works under a broader projection model including weak perspective (scaled orthography) and paraperspective [9] up to an affine camera [8]

2 Factorization Method: A New Formulation Including Translation

The factorization method was originally introduced by Tomasi and Kanade[11] for the case of single static object viewed by a moving camera. Here we will reformulate the method in such a way that a static camera observes a scene with a moving object. Also, whereas the translation component of motion is first eliminated in the Tomasi-Kanade formulation, we will retain that component in our formulation.

2.1 World and Observations

The object moves relative to the camera which acquires images. In the sequence we track feature points from frame to frame. The position of an object point $\mathbf{p}_i^T = [X_i Y_i Z_i]^T$ expressed in homogeneous coordinates in the camera frame, is given by:

$$\mathbf{s}_{fi}^C \equiv \begin{bmatrix} \mathbf{p}_{fi}^C \\ 1 \end{bmatrix} = \begin{bmatrix} \mathbf{R}_f & \mathbf{t}_f \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{p}_i \\ 1 \end{bmatrix} \quad (1)$$

$$= \begin{bmatrix} \mathbf{R}_f & \mathbf{t}_f \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \mathbf{s}_i \quad (2)$$

where R_f and t_f are, respectively, the rotation and translation components. Suppose that we track N feature points over F frames, and that we collect all these measurements into a single matrix:

$$\begin{bmatrix} u_{11} \dots & u_{1N} \\ \vdots & \vdots \\ u_{F1} \dots & u_{FN} \\ v_{11} \dots & v_{1N} \\ \vdots & \vdots \\ v_{F1} \dots & v_{FN} \end{bmatrix} = \begin{bmatrix} \mathbf{i}_1^T & t_{x_1} \\ \vdots & \vdots \\ \mathbf{i}_F^T & t_{x_F} \\ \mathbf{j}_1^T & t_{y_1} \\ \vdots & \vdots \\ \mathbf{j}_F^T & t_{y_F} \end{bmatrix} [\mathbf{s}_1 \dots \mathbf{s}_N] \quad (3)$$

$$\mathbf{W} = \mathbf{M}\mathbf{S}. \quad (4)$$

where (u_{fi}, v_{fi}) are the feature image position, vectors $\mathbf{i}_f^T = [i_{x_f} \ i_{y_f} \ i_{z_f}]$, $\mathbf{j}_f^T = [j_{x_f} \ j_{y_f} \ j_z]^T$, ($t_f = 1 \dots F$) are the first two rows of the rotation matrix at instant f , and $(\mathbf{t}_{x_f}, \mathbf{t}_{y_f})$ are the X and Y coordinates of the position of the object's coordinate frame, in the camera frame, at the same instant.

2.2 Solution for Shape and Motion by Factorization

Recovering the shape and motion is equivalent to start with a given matrix \mathbf{W} and obtain a factorization into motion matrix \mathbf{M} and shape matrix \mathbf{S} . By simple inspection of (4) we can see that since \mathbf{M} and \mathbf{S} can be at most rank 4, \mathbf{W} will be at most rank 4. Using Singular Value Decomposition (SVD), \mathbf{W} is decomposed and approximated as:

$$\mathbf{W} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T. \quad (5)$$

Matrix $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \sigma_3, \sigma_4)$ is a diagonal matrix made of the four biggest singular values which reveal the most important components in the data. Matrices $\mathbf{U} \in R^{2F \times 4}$ and $\mathbf{V} \in R^{N \times 4}$ are the left and right singular matrices respectively, such that $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathcal{I}$ (the 4×4 identity matrix).

By defining,

$$\hat{\mathbf{M}} \equiv \mathbf{U}\mathbf{\Sigma}^{\frac{1}{2}}, \quad \hat{\mathbf{S}} \equiv \mathbf{\Sigma}^{\frac{1}{2}}\mathbf{V}^T \quad (6)$$

we have the two matrices whose product can represent the bilinear system \mathbf{W} . However, this factorization is not unique, since for any invertible 4×4 matrix \mathbf{A} , $\mathbf{M} = \hat{\mathbf{M}}\mathbf{A}$ and $\mathbf{S} = \mathbf{A}^{-1}\hat{\mathbf{S}}$ are also a possible solution because

$$\mathbf{M}\mathbf{S} = (\hat{\mathbf{M}}\mathbf{A})(\mathbf{A}^{-1}\hat{\mathbf{S}}) = \hat{\mathbf{M}}\hat{\mathbf{S}} = \mathbf{W}. \quad (7)$$

The exact solution can be computed, using the fact that \mathbf{M} must have certain properties. Let us denote the 4×4 matrix \mathbf{A} as the concatenation of two blocks,

$$\mathbf{A} \equiv [\mathbf{A}_R | \mathbf{a}_t], \quad (8)$$

The first block \mathbf{A}_R is the first 4×3 submatrix related to the rotational component and the second block \mathbf{a}_t is a 4×1 vector related to translation. Now, since

$$\mathbf{M} = \hat{\mathbf{M}}\mathbf{A} = [\hat{\mathbf{M}}\mathbf{A}_R | \hat{\mathbf{M}}\mathbf{a}_t], \quad (9)$$

we can impose motion constraints, one on rotation and the other on translation, in order to solve for \mathbf{A} .

2.2.1 Rotation Constraints

Block \mathbf{A}_R of \mathbf{A} , which is related to rotational motion, is constrained by the orthonormality of axes vectors \mathbf{i}_f^T and \mathbf{j}_f^T : each of the $2F$ rows entries of matrix $\hat{\mathbf{M}}\mathbf{A}_R$ is a unit norm vector and the first and second set of F rows are pairwise orthogonal. This yields a set of constraints:

$$\hat{\mathbf{m}}_i \mathbf{A}_R \mathbf{A}_R^T \hat{\mathbf{m}}_i^T = 1 \quad \hat{\mathbf{m}}_j \mathbf{A}_R \mathbf{A}_R^T \hat{\mathbf{m}}_j^T = 1 \quad (10)$$

$$\hat{\mathbf{m}}_i \mathbf{A}_R \mathbf{A}_R^T \hat{\mathbf{m}}_j^T = 0 \quad (11)$$

where $\hat{\mathbf{m}}_i, \hat{\mathbf{m}}_j$ are rows i and j of matrix $\hat{\mathbf{M}}$ for $i = 1 \dots F$ and $j = F + 1 \dots 2F$. This is an overconstrained system which can be solved for the entries of $\mathbf{A}_R \mathbf{A}_R^T$ by using least squares techniques, and subsequently solving for \mathbf{A}_R . See [11] for a detailed solution procedure.

2.2.2 Translation Constraints

In orthography, the projection of the 3D centroid of an object features into the image plane is the centroid of the feature points. The X and Y position of the centroid of the feature points is the average of each row of \mathbf{W} :

$$\bar{\mathbf{w}} \equiv \begin{bmatrix} \frac{1}{N} \sum u_{1,i} \\ \vdots \\ \frac{1}{N} \sum v_{F,i} \end{bmatrix} = \mathbf{M}\bar{\mathbf{s}} = [\hat{\mathbf{M}}\mathbf{A}_R | \hat{\mathbf{M}}\mathbf{a}_t] \begin{bmatrix} \bar{\mathbf{p}} \\ 1 \end{bmatrix}, \quad (12)$$

where $\bar{\mathbf{p}} \equiv \frac{1}{N} \sum \mathbf{p}_i$ is the centroid of the object. The origin of the object's coordinate system is arbitrary, so we can choose to place it at the centroid of the object, that is $\bar{\mathbf{p}} = 0$. Then it follows immediately from (12) that

$$\bar{\mathbf{w}} = \hat{\mathbf{M}}\mathbf{a}_t \quad (13)$$

This expression is also an overconstrained system of equations, which can be solved for the entries of \mathbf{a}_t in the least square sense. The best estimate will be given by

$$\mathbf{a}_t = (\hat{\mathbf{M}}^T \hat{\mathbf{M}})^{-1} \hat{\mathbf{M}}^T \bar{\mathbf{w}} \quad (14)$$

$$= \Sigma^{-1/2} \mathbf{U}^T \bar{\mathbf{w}}, \quad (15)$$

which completes the computation of all the elements of matrix \mathbf{A} .

3 The Multi-body Factorization Method

So far we have assumed that the scene contains a single moving object. If there is more than one moving object, the measurement matrix \mathbf{W} will contain features (columns) which originate from different motions. One may think that solving the problem requires first sorting the columns of the measurements matrix \mathbf{W} into submatrices, each of which contains features solely from one object, so that the factorization technique of the previous sections can be applied individually. We will show in this section that the multi-body problem can be solved without prior segmentation. For the sake of simplicity in presentation we will present the theory and method for the case of two bodies, but it will be clear that the method is applicable to the general case of an arbitrary unknown number of objects.

3.1 Multi-body Motion Recovery Problem: Its Difficulty

Suppose we have a scene in which two objects are moving and we take an image sequence of F frames. Suppose also that the set of features that we have observed and tracked in the image sequence actually consists of N_1 feature points from object 1 and N_2 from object 2 which are observed in an image sequence of F frames.

Imagine for the moment that somehow we know the classification of features and thus could permute the columns of \mathbf{W} in such a way that the first N_1 columns

belong to object 1 followed by the N_2 columns from object 2. Matrix \mathbf{W} would have the canonical form:

$$\mathbf{W}^* \equiv [\mathbf{W}_1 | \mathbf{W}_2]. \quad (16)$$

Each measurement submatrix can be factorized as

$$\mathbf{W}_l = \mathbf{U}_l \Sigma_l \mathbf{V}_l^T \quad (17)$$

$$= \mathbf{M}_l \mathbf{S}_l = (\hat{\mathbf{M}}_l \mathbf{A}_l) (\mathbf{A}_l^{-1} \hat{\mathbf{S}}_l) \quad (18)$$

with $l = 1$ and 2 for object 1 and 2 respectively. By denoting

$$\mathbf{M}^* \equiv [\mathbf{M}_1 | \mathbf{M}_2], \quad \mathbf{S}^* \equiv \begin{bmatrix} \mathbf{S}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_2 \end{bmatrix} \quad (19)$$

$$\mathbf{A}^* \equiv \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 \end{bmatrix}, \quad \mathbf{U}^* \equiv [\mathbf{U}_1 | \mathbf{U}_2] \quad (20)$$

$$\Sigma^* \equiv \begin{bmatrix} \Sigma_1 & \mathbf{0} \\ \mathbf{0} & \Sigma_2 \end{bmatrix}, \quad \mathbf{V}^{*T} \equiv \begin{bmatrix} \mathbf{V}_1^T & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_2^T \end{bmatrix} \quad (21)$$

we express a factorization in a similar way of the single object, that is, equation (16) now has the canonical factorization:

$$\mathbf{W}^* = \mathbf{M}^* \mathbf{S}^* \quad (22)$$

$$\mathbf{S}^* = \mathbf{A}^{*-1} \Sigma^{*\frac{1}{2}} \mathbf{V}^{*T}, \quad \mathbf{M}^* = \mathbf{U}^* \Sigma^{*\frac{1}{2}} \mathbf{A}^* \quad (23)$$

From equation (22), we see that \mathbf{W}^* (and therefore \mathbf{W}) will have at most rank 8, since \mathbf{W}_1 and \mathbf{W}_2 are at most rank 4. Let us consider for the remainder of this paper the non-degenerate case where the rank of \mathbf{W} is in fact equal to 8; that is, the object shape is actually three-dimensional (not planar or line) and the motion vectors span 3D for both objects. The degenerate cases will be briefly touched in the last section and are discussed in more detail in [4].

In reality, we do not know which features belong to which object, and thus the columns of the given measurement matrix \mathbf{W} are a mixture of features from object 1 and 2. We can still apply singular value decomposition (SVD) to the measurement matrix, and obtain

$$\mathbf{W} = \mathbf{U} \Sigma \mathbf{V}^T. \quad (24)$$

Then it may appear that the remaining task is to find the linear canonical transformation \mathbf{A}^* in (20) such that shape and motion will have the block structure of equations (23) and (23).

There is, however, a fundamental difficulty in doing this. The metric (rotation and translation) constraints (eq.(10)-(10) and (13)-(15)) were obtained in section 2.2 by considering that the motion matrix for one object, that is, by assuming that the measurement matrix consists of features from a single object. Those constraints are therefore applicable only after knowing the segmentation. This is exactly the mathematical manifestation of the cyclic dilemma mentioned earlier.

Faced with this difficulty, a usual approach would be to group features bit by bit so that we segment \mathbf{W}

into two rank-4 matrices and obtain the factorization of the form (22). For example, a most simplistic procedure would be like the following. Pick the first four columns of \mathbf{W} and form a rank-4 subspace. If the fifth column belongs to the subspace (ie. is linear dependent on the first four, or "almost" linear dependent in the case of noisy measurement), then classify it to the same object as the first four columns and update the subspace representation. Otherwise, it belongs to a new object. Apply this procedure recursively to all the remaining columns. This approach is in fact essentially the one used by [3] and [5] to split matrix \mathbf{W} , and similar to what was suggested by Ullman [12], whose criteria for merging was local rigidity.

There are a few disadvantages in this cluster-and-test approach. First, there is no guarantee that the first four columns, which always form a rank-4 subspace, are from the same object. Second, if we use a sequential procedure like the one above or its variation, the final result is dependent on where we start the procedure, and alternatively, the search for the globally optimal segmentation most likely will be computational very expensive. Finally, the prior knowledge of the number of objects becomes very critical, since depending on the decision criterium of subspace inclusion the final number of objects may vary arbitrarily.¹

3.2 Mathematical Construct of Shapes Invariant to Motions

The main difficulty in the multi-body structure-from-motion problem revealed above is that shape and motion interact. Mathematically, the equation (22) indicates that the rank-8 measurement space is originally generated by the two subspaces of rank 4 each, represented by the block-diagonal shape matrix \mathbf{S}^* . However, the recovered shape space \mathbf{V}^T , obtained by the singular value decomposition of the non-canonical \mathbf{W} , is in general a linear combination of the two subspaces and has lost the block-diagonal structure.

There is however a mathematical construct that preserves the original subspace structure. Let us define \mathbf{Q} as $(N_1 + N_2) \times (N_1 + N_2)$ square matrix

$$\mathbf{Q} \equiv \mathbf{V}\mathbf{V}^T. \quad (25)$$

We will call this matrix the *shape interaction matrix*. Mathematically, it is the orthogonal operator that projects $N = (N_1 + N_2)$ dimensional vectors to the subspace spanned by the columns of \mathbf{V} . This matrix \mathbf{Q} has several interesting and useful properties. First, by definition it is uniquely computable only from the measurements \mathbf{W} without knowing the segmentation, since \mathbf{V} is uniquely obtained by the singular value decomposition of \mathbf{W} .

Secondly, permuting columns of \mathbf{W} does not change the set of values $\{Q_{ij}\}$ that appear in \mathbf{Q} though their

arrangement in \mathbf{Q} does; swapping columns l and m of \mathbf{W} results in swapping columns l and m of \mathbf{V}^T . Therefore it results in simultaneously swapping columns l and m and rows l and m in \mathbf{Q} , but not their entry values.

Thirdly, each element of \mathbf{Q} provides important information about whether a pair of features belong to the same object. Since the set of values do not change, let us compute \mathbf{Q}^* , the shape interaction matrix for the canonical measurement matrix \mathbf{W}^* . By substituting (23) into (25), we obtain

$$\mathbf{Q}^* = \mathbf{V}^* \mathbf{V}^{*T} \quad (26)$$

$$= \mathbf{S}^{*T} \mathbf{A}^{*T} \mathbf{\Sigma}^* \mathbf{A}^* \mathbf{S}^* \quad (27)$$

$$= \mathbf{S}^{*T} (\mathbf{A}^{*-1} \mathbf{\Sigma}^{*-1} \mathbf{A}^{*-T})^{-1} \mathbf{S}^* \quad (28)$$

$$= \mathbf{S}^{*T} \left[(\mathbf{A}^{*-1} \mathbf{\Sigma}^{*-1/2} \mathbf{V}^{*T}) (\mathbf{V}^* \mathbf{\Sigma}^{*-1/2} \mathbf{A}^{*-T}) \right]^{-1} \mathbf{S}^* \quad (29)$$

$$= \begin{bmatrix} \mathbf{S}_1^T & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_2^T \end{bmatrix} \begin{bmatrix} \mathbf{\Lambda}_1^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{\Lambda}_2^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{S}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_2 \end{bmatrix} \quad (30)$$

$$= \begin{bmatrix} \mathbf{S}_1^T \mathbf{\Lambda}_1^{-1} \mathbf{S}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{S}_2^T \mathbf{\Lambda}_2^{-1} \mathbf{S}_2 \end{bmatrix}. \quad (31)$$

where $\mathbf{\Lambda}_1$ and $\mathbf{\Lambda}_2$ are the 4×4 matrices of the moments of inertia of each object. This means that the canonical \mathbf{Q}^* matrix for the sorted \mathbf{W}^* has a very defined block-diagonal structure. Moreover, each entry has the value

$$Q_{ij}^* = \begin{cases} \mathbf{s}_{1i}^T \mathbf{\Lambda}_1^{-1} \mathbf{s}_{1j} & \text{feat. } i \text{ and } j \text{ belong to obj. 1} \\ \mathbf{s}_{2i}^T \mathbf{\Lambda}_2^{-1} \mathbf{s}_{2j} & \text{feat. } i \text{ and } j \text{ belong to obj. 2} \\ 0 & \text{feat } i \text{ and } j \text{ belong to diff. obj.} \end{cases} \quad (32)$$

Finally and most importantly, the set of values $\{Q_{ij}^*\}$, which is the same as $\{Q_{ij}\}$ are invariant to motion. This is true since equations (32) include only \mathbf{S} 's, and not \mathbf{M} . In other words, in whatever way the objects move they will produce the same set of entries in matrix \mathbf{Q} .

In summary, we have shown that without knowing the segmentation of features we can compute matrix \mathbf{Q} whose element Q_{ij} can be interpreted as a measure of the interaction between feature i and j : if the value is non zero, they belong to the same object, and if they don't belong to the same object, the value is zero. Also, if the features are sorted correctly into the canonical form of the measurement matrix \mathbf{W}^* , then the corresponding canonical shape interaction matrix \mathbf{Q}^* must be block diagonal.

3.3 Sorting Matrix \mathbf{Q} into Canonical Form

The problem of segmenting and recovering motion of multiple objects now has reduced to sorting the entries of matrix \mathbf{Q} by swapping pairs of rows and columns until it becomes block diagonal. Once achieved, the corresponding permutations of columns

¹While this is beyond the scope of the assumption in this section, this cluster-and-test approach also requires the prior knowledge of the ranks of objects as well. Since for example a rank-8 measurement matrix might have been generated by two line (rank-2) objects and one full 3D (rank 4) object instead of two full 3D objects, and therefore committing to find two rank-4 subspaces might be wrong.

of \mathbf{W} will transform it into its canonical form where features from one object are grouped into contiguous columns. This relationship between sorting \mathbf{Q} and permuting \mathbf{W} is illustrated in figure 1.

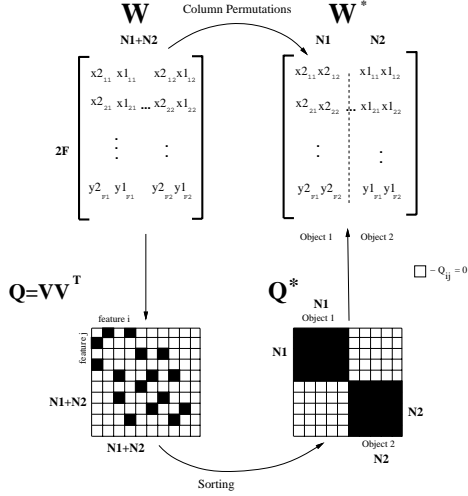


Figure 1: Segmentation process

With noisy measurements, a pair of features from different objects may exhibit a small non-zero entry in \mathbf{Q} . We can regard Q_{ij}^2 as representing the energy of the shape interaction, and the block diagonalization of \mathbf{Q} can be achieved by minimizing the total energy of all possible off-diagonal blocks over all set of permutations of rows and columns of \mathbf{Q} . We found that a simple iterative minimization procedure suffices for our purpose. Alternatively, we can regard matrix $\{Q_{ij}^2\}$ as defining a graph of $N_1 + N_2$ nodes, where the Q_{ij}^2 indicates the weight of the link (i, j) . We also found that graph-theoretical algorithms, such as the minimum spanning tree, can be used to achieve the block diagonalization more efficiently than the energy minimization. The detailed procedures are presented in [4].

3.4 Summary of Algorithm

While we have presented the theory for the case of two full-3D objects, it is easy to see that its essential part holds for more general cases. First the matrix \mathbf{Q}^* has the block diagonal structure for an arbitrary number of moving objects, that is, an entry Q_{ij} of the \mathbf{Q} matrix equals to zero if features i and j belong to different objects. Furthermore, this property holds even when the shape matrix of the objects has rank less than 4 (planes and lines). The computation of \mathbf{Q} by (25) requires only the knowledge of the total rank of \mathbf{W} , which we can determine by SVD. Finally once \mathbf{Q}^* is obtained, instead of permuting columns of \mathbf{W} we can use the equivalent permutation of \mathbf{V}^T , since it is more computationally efficient.

The whole algorithm of the multi-body factorization method is now summarized as:

1. Extract and track features in the input image sequence and create matrix \mathbf{W}
2. Compute $r = \text{rank}(\mathbf{W})$
3. Decompose matrix \mathbf{W} using SVD
4. Compute shape interaction matrix \mathbf{Q} using the first r rows of \mathbf{V}^T
5. Block-diagonalize \mathbf{Q}
6. Permute matrix \mathbf{V}^T into submatrices, each corresponding to a single object
7. Compute \mathbf{A}_i for each object, and thus its shape and motion.

4 Experiments

The scene consists of two roughly cylindrical shapes covered by rolling a cardboard sheet and drawing dots on the surface. The cylinder on the right tilts and rotates in the plane parallel to the image plane while the cylinder on the left hand side rotates around its axis. In the sequence with 85 images, a total of 55 features are detected and tracked: 27 belonging to the left cylinder and 28 the other, while, of course, the algorithm was not given that information. Figure 2 shows the 85-th frame in the sequence with the tracks of the selected features superimposed. Figure 3(a) show the shape interaction matrix \mathbf{Q} for the unsorted input features. The sorted block diagonal matrix \mathbf{Q}^* is shown in figure 3(b), and the features are grouped accordingly for individual shape recovery. The resultant three-dimensional points are displayed in figure 4 with linearly interpolated surface in order to convey a better perception of their shape.

5 Discussion and Conclusion

In this paper we have shown that the problem of multi-body structure-from-motion problem can be solved systematically by using the shape interaction matrix. The striking fact is that the method allows for segmenting or grouping image features into separate objects *based on* their shape properties *without* explicitly computing the individual shapes themselves. Also, no prior knowledge of the number of moving objects in the scene is assumed in the algorithm.

This is due to the interesting and useful invariant properties of the shape-interaction matrix \mathbf{Q} . We have shown that \mathbf{Q} is motion invariant. Even when the matrix is computed from a different set of image-level measurements \mathbf{W} generated by a different set of motions of objects, its entries will remain invariant. The motion invariance property of \mathbf{Q} means also that the degree of complexity of the solution is dependent on the scene complexity, but not on the motion complexity.

The shape interaction matrix \mathbf{Q} is also invariant to the selection of individual object coordinate frames. We can easily see that by considering transforming the shape of object k , \mathbf{S}_k , by a 4×4 matrix \mathbf{T} ,

$$\mathbf{S}'_k = \mathbf{T}\mathbf{S}_k. \quad (33)$$



Figure 2: Image of the objects and feature tracks

The corresponding block-diagonal element matrix of \mathbf{Q}^* will be

$$\begin{aligned} \mathbf{S}'_k (\mathbf{S}'_k \mathbf{S}'_k^T)^{-1} \mathbf{S}'_k &= (\mathbf{T} \mathbf{S}_k)^T (\mathbf{T} \mathbf{S}_k \mathbf{S}_k^T \mathbf{T}^T)^{-1} (\mathbf{T} \mathbf{S}_k) \\ &= \mathbf{S}_k^T (\mathbf{S}_k \mathbf{S}_k^T)^{-1} \mathbf{S}_k \end{aligned} \quad (34)$$

and therefore the entries of matrix \mathbf{Q}^* remain the same. Another interesting fact is that the shape interaction matrix can handle many degenerate cases as well, where objects may be full 3-D object but also linear or planar. More research is required for the degenerate cases including the cases where the motions are degenerate. Also, in order to achieve robustness under the presence of noise we need to relate its level with the thresholds necessary in some of the decision making processes. They include the identification of the rank of the measurement matrix in the singular value decomposition, and the determination of block-diagonality in sorting the shape interaction matrix. The report [4] explores some of those issues.

Acknowledgments: The authors wish to thank Martial Hebert, José Moura and José Dias for useful suggestions and comments to this work.

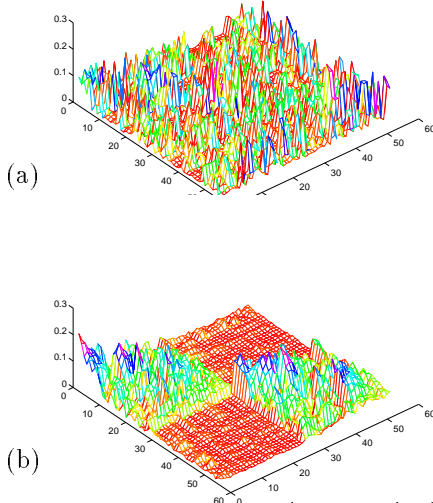


Figure 3: The shape interaction matrix for the lab scene: (a) Unsorted \mathbf{Q} ; (b) block-diagonalized \mathbf{Q}^*

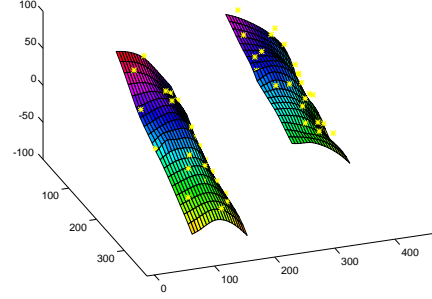


Figure 4: The recovered shape of the two cylinders

References

- [1] E. Adelson and J. Bergen. Spatiotemporal energy models for the perception of motion. *J. of the Opt. Soc. of America*, 2(2):284–299, 1985.
- [2] J. Bergen, P. Burt, R. Hingorani, and S. Peleg. Computing two motions from three frames. In *ICCV*, 1990.
- [3] T. Boulton and L. Brown. Factorization-based segmentation of motions. In *Proc. of the IEEE Workshop on Visual Motion*, 1991.
- [4] J. Costeira and T. Kanade. A multibody factorization method for motion analysis; Degenerate cases. Technical report, SCS, Carnegie Mellon University, to be printed.
- [5] C. W. Gear. Feature grouping in moving objects. In *Proc. of the workshop on motion of non-rigid and articulated objects*, Austin, Texas, 1994.
- [6] M. Irani, R. Benny, and S. Peleg. Computing occluding and transparent motions. *IJCV*, 12(1), 1994.
- [7] R. S. Jasinschi, A. Rosenfeld, and K. Sumi. Perceptual motion transparency: the role of geometrical information. *Journ. of the Opt. Soc. of America*, 9(11), 1992.
- [8] J. Koendering and A. van Doorn. Affine structure from motion. *Journ. of the Opt. Soc. of America*, 12(1), 1992.
- [9] C. Poelman and T. Kanade. A paraperspective factorization method for shape and motion recovery. In *Proc. ECCV*, 1994.
- [10] D. Sinclair. Motion segmentation and local structure. In *Proc. ICCV*, 1993.
- [11] C. Tomasi and T. Kanade. Shape from motion from image streams under orthography: A factorization method. *IJCV*, 9(2), 1992.
- [12] S. Ullman. Maximizing rigidity: The incremental recovery of 3-d structure from rigid and rubbery motion. Technical Report A.I. Memo No. 721, MIT.