

Multi-stage Optimization for Multi-body Motion Segmentation

Kenichi Kanatani and Yasuyuki Sugaya

Department of Information Technology, Okayama University, Okayama 700-8530 Japan

{kanatani, sugaya}@suri.it.okayama-u.ac.jp

Many techniques have been proposed for separating feature point trajectories tracked through a video sequence into independent motions, but objects are usually assumed to undergo general 3-D motions. As a result, the separation accuracy considerably deteriorates in realistic video sequences in which object motions are nearly degenerate. In this paper, we introduce unsupervised learning assuming degenerate motions followed by unsupervised learning assuming general 3-D motions. This multi-stage optimization allows us to not only separate simple motions that we frequently encounter with high precision but also preserve the high performance for considerably general 3-D motions. Doing simulations and real video experiments, we demonstrate that our method is superior to all existing methods.

1. Introduction

Separating feature point trajectories tracked through a video sequence into independent motions is the first step of many video processing applications. Already, many techniques have been proposed for this task.

Costeira and Kanade [1] proposed a segmentation algorithm based on the shape interaction matrix. Gear [3] used the reduced row echelon form and graph matching. Ichimura [4] used the discrimination criterion of Otsu [11]. He also used the QR decomposition [5]. Inoue and Urahama [6] introduced fuzzy clustering. Kanatani [8, 9, 10] incorporated model selection using the geometric AIC [7]. Wu et al. [18] introduced orthogonal subspace decomposition.

However, all these methods assume that the objects undergo general 3-D motions relative to the camera. As a result, segmentation fails when the motions are degenerate, e.g., all the objects are simply translating independently (not necessarily along straight lines). This type of degeneracy frequently occurs in practical applications. Though strict degeneracy may be rare, the segmentation accuracy considerably deteriorates if the motions are nearly degenerate.

At first sight, segmenting simple motions may seem easier than segmenting complicated motions. In reality, however, the opposite is the case, because complicated motions have sufficient cues for mutual discrimination. In fact, we have found through our experiments that many methods that exhibit high accuracy for complicated simulations perform very poorly for real video sequences.

To cope with this, we have presented a method for automatically selecting the best motion model using the geometric AIC, but the accuracy improvement was very much limited [14, 15].

In this paper, we introduce unsupervised learning [13] assuming degenerate motions followed by unsupervised learning assuming general 3-D motions. This multi-stage optimization allows us to not only separate simple motions with high precision but also preserve the high performance for considerably general 3-D motions.

In Sec. 2, we describe the geometric constraints that underlie our method. In Sec. 3, we introduce unsupervised learning of the non-Bayesian and Bayesian types. Our multi-stage optimization scheme is described in Sec. 4. In Sec. 5, we show synthetic and real video examples and demonstrate that our method is superior to all existing methods. Section 6 concludes this paper.

2. Geometric Constraints

2.1 Trajectory of feature points

Suppose we track N feature points over M frames. Let $(x_{\kappa\alpha}, y_{\kappa\alpha})$ be the coordinates of the α th point in the κ th frame. Stacking all the coordinates vertically, we represent the entire trajectory by the following $2M$ -dimensional *trajectory vector*:

$$\mathbf{p}_\alpha = (x_{1\alpha} \ y_{1\alpha} \ x_{2\alpha} \ y_{2\alpha} \ \cdots \ x_{M\alpha} \ y_{M\alpha})^\top. \quad (1)$$

For convenience, we identify the frame number κ with “time” and refer to the κ th frame as “time κ ”.

We identify the XYZ camera coordinate system with the world frame, relative to which multiple objects (including the background) are moving. Consider a 3-D coordinate system fixed to one moving object, and let \mathbf{t}_κ and $\{\mathbf{i}_\kappa, \mathbf{j}_\kappa, \mathbf{k}_\kappa\}$ be, respectively, its origin and basis vectors at time κ . If the α th point has coordinates $(a_\alpha, b_\alpha, c_\alpha)$ with respect to this coordinate system, its position with respect to the world frame at time κ is

$$\mathbf{r}_{\kappa\alpha} = \mathbf{t}_\kappa + a_\alpha \mathbf{i}_\kappa + b_\alpha \mathbf{j}_\kappa + c_\alpha \mathbf{k}_\kappa. \quad (2)$$

2.2 Affine camera model

We assume an affine camera, which generalizes orthographic, weak perspective, and paraperspective projections [12]: the 3-D point $\mathbf{r}_{\kappa\alpha}$ is projected onto the image position

$$\begin{pmatrix} x_{\kappa\alpha} \\ y_{\kappa\alpha} \end{pmatrix} = \mathbf{A}_\kappa \mathbf{r}_{\kappa\alpha} + \mathbf{b}_\kappa, \quad (3)$$

where \mathbf{A}_κ and \mathbf{b}_κ are, respectively, a 2×3 matrix and a 2-dimensional vector determined by the position and orientation of the camera and its internal

parameters at time κ . Substituting Eq. (2), we have

$$\begin{pmatrix} x_{\kappa\alpha} \\ y_{\kappa\alpha} \end{pmatrix} = \tilde{\mathbf{m}}_{0\kappa} + a_\alpha \tilde{\mathbf{m}}_{1\kappa} + b_\alpha \tilde{\mathbf{m}}_{2\kappa} + c_\alpha \tilde{\mathbf{m}}_{3\kappa}, \quad (4)$$

where $\tilde{\mathbf{m}}_{0\kappa}$, $\tilde{\mathbf{m}}_{1\kappa}$, $\tilde{\mathbf{m}}_{2\kappa}$, and $\tilde{\mathbf{m}}_{3\kappa}$ are 2-dimensional vectors determined by the position and orientation of the camera and its internal parameters at time κ . From Eq. (4), the trajectory vector \mathbf{p}_α in Eq. (1) can be written in the form

$$\mathbf{p}_\alpha = \mathbf{m}_0 + a_\alpha \mathbf{m}_1 + b_\alpha \mathbf{m}_2 + c_\alpha \mathbf{m}_3, \quad (5)$$

where \mathbf{m}_0 , \mathbf{m}_1 , \mathbf{m}_2 , and \mathbf{m}_3 are the $2M$ -dimensional vectors obtained by stacking $\tilde{\mathbf{m}}_{0\kappa}$, $\tilde{\mathbf{m}}_{1\kappa}$, $\tilde{\mathbf{m}}_{2\kappa}$, and $\tilde{\mathbf{m}}_{3\kappa}$ vertically over the M frames, respectively.

2.3 Constraints on image motion

Equation (5) implies that the trajectories of the feature points that belong to one object are constrained to be in the *4-dimensional subspace* spanned by $\{\mathbf{m}_0, \mathbf{m}_1, \mathbf{m}_2, \mathbf{m}_3\}$ in \mathcal{R}^{2M} . It follows that multiple moving objects can be segmented into individual motions by separating the trajectory vectors $\{\mathbf{p}_\alpha\}$ into distinct 4-dimensional subspaces. This is the principle of the method of *subspace separation* [8, 9].

In addition, the coefficient of \mathbf{m}_0 in Eq. (5) is identically 1 for all α . This means that the trajectories are in a *3-dimensional affine space* within that 4-dimensional subspace. It follows that multiple moving objects can be segmented into individual motions by separating the trajectory vectors $\{\mathbf{p}_\alpha\}$ into distinct 3-dimensional affine spaces. This is the principle of the method of *affine space separation* [10].

Theoretically, the segmentation accuracy should be higher if we use stronger constraints. In fact, according to simulations, the affine space separation performs better than the subspace separation except in the case in which perspective effects are very strong in the presence of small noise [10]. For real video sequences, however, the affine space separation accuracy is sometimes lower than that of the subspace separation [14, 15], which is inconsistent with the simulation results. The cause of this inconsistency will be clarified in the subsequent analysis.

3. Unsupervised Learning

3.1 Non-Bayesian type

Segmentation by the subspace separation and the affine space separation is not always correct. However, we can optimize the segmentation *a posteriori* by optimally fitting a 3-dimensional affine space (or a 4-dimensional subspace) to each trajectory class and reclassifying each trajectory to the closest affine space (or subspace) (Fig. 1(a)). This process is iterated until the classification converges.

If the noise in the coordinates of the feature points is an independent Gaussian random variable of mean 0 and a constant variance, this procedure can be

viewed as unsupervised learning based on maximum likelihood estimation, since minimizing the distance of points from the fitted space is equivalent to maximizing their likelihood under our noise model.

3.2 Bayesian type

We may also take into consideration the internal data distributions inside the fitted spaces (Fig. 1(b)). This is the standard approach to unsupervised learning for pattern recognition. However, the existence of geometric constraints somewhat complicates the likelihood computation. For the affine space constraint, the actual procedure is as follows (the procedure for the subspaces constraint goes similarly).

Let $n = 2M$. Suppose N n -dimensional trajectory vectors $\{\mathbf{p}_\alpha\}$ are initially classified into m classes. Define the weight $W_\alpha^{(k)}$ of the vector \mathbf{p}_α for the k th class by

$$W_\alpha^{(k)} = \begin{cases} 1 & \text{if } \mathbf{p}_\alpha \text{ belongs to the } k\text{th class} \\ 0 & \text{otherwise} \end{cases}. \quad (6)$$

Then, iterate the following procedures A and B in turn until all the weights $\{W_\alpha^{(k)}\}$ converge.

A. Do the following computation for each class $k = 1, \dots, m$.

1. Compute

$$w^{(k)} = \frac{1}{N} \sum_{\alpha=1}^N W_\alpha^{(k)}. \quad (7)$$

2. Compute the centroid $\mathbf{p}_C^{(k)}$ of the k th class:

$$\mathbf{p}_C^{(k)} = \frac{\sum_{\alpha=1}^N W_\alpha^{(k)} \mathbf{p}_\alpha}{\sum_{\alpha=1}^N W_\alpha^{(k)}}. \quad (8)$$

3. Compute the $n \times n$ moment matrix of the k th class:

$$\mathbf{M}^{(k)} = \frac{\sum_{\alpha=1}^N W_\alpha^{(k)} (\mathbf{p}_\alpha - \mathbf{p}_C^{(k)}) (\mathbf{p}_\alpha - \mathbf{p}_C^{(k)})^\top}{\sum_{\alpha=1}^N W_\alpha^{(k)}}. \quad (9)$$

4. Let $\lambda_1 \geq \lambda_2 \geq \lambda_3$ be the three largest eigenvalues of the matrix $\mathbf{M}^{(k)}$, and $\mathbf{u}_1^{(k)}$, $\mathbf{u}_2^{(k)}$, and $\mathbf{u}_3^{(k)}$ the corresponding unit eigenvectors.
5. Compute the $n \times n$ projection matrices

$$\mathbf{P}^{(k)} = \sum_{i=1}^3 \mathbf{u}_i^{(k)} \mathbf{u}_i^{(k)\top}, \quad \mathbf{P}_\perp^{(k)} = \mathbf{I} - \mathbf{P}^{(k)}, \quad (10)$$

where \mathbf{I} denotes the $n \times n$ unit matrix.

6. Estimate the noise variance in the direction orthogonal to the k th affine space by

$$\hat{\sigma}_k^2 = \max\left[\frac{\text{tr}[\mathbf{P}_\perp^{(k)} \mathbf{M}^{(k)} \mathbf{P}_\perp^{(k)}]}{n-3}, \sigma^2\right], \quad (11)$$

where $\text{tr}[\cdot]$ denotes the trace and σ is an estimate of the tracking accuracy¹.

¹We found $\sigma = 0.5$ (pixels) a reasonable value [16].

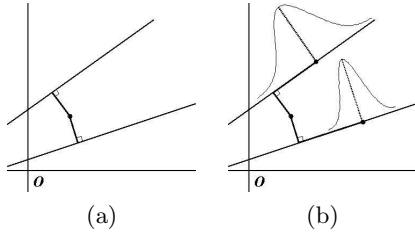


Figure 1: Segmentation criterion: (a) non-Bayesian type; (b) Bayesian type.

7. Compute the $n \times n$ covariance matrix of the k th class by

$$\mathbf{V}^{(k)} = \mathbf{P}^{(k)} \mathbf{M}^{(k)} \mathbf{P}^{(k)} + \hat{\sigma}_k^2 \mathbf{P}_\perp^{(k)}. \quad (12)$$

B. Do the following computation for each trajectory vector \mathbf{p}_α , $\alpha = 1, \dots, m$.

1. Compute the conditional likelihood $P(\alpha|k)$, $k = 1, \dots, m$, by

$$P(\alpha|k) = \frac{e^{-\langle \mathbf{p}_\alpha - \mathbf{p}_C^{(k)}, \mathbf{V}^{(k)-1} (\mathbf{p}_\alpha - \mathbf{p}_C^{(k)}) \rangle / 2}}{\sqrt{\det \mathbf{V}^{(k)}}}. \quad (13)$$

2. Recompute the weights $W_\alpha^{(k)}$, $k = 1, \dots, m$, by

$$W_\alpha^{(k)} = \frac{w^{(k)} P(\alpha|k)}{\sum_{l=1}^m w^{(l)} P(\alpha|l)}. \quad (14)$$

After the iterations of A, B, and C have converged, the α th trajectory is classified into the k th class that maximizes $W_\alpha^{(k)}$, $k = 1, \dots, N$.

3.3 Interpretation

In the above iterations, we fit a Gaussian distribution of mean $\mathbf{p}_C^{(k)}$ and the rank 3 covariance matrix $\mathbf{P}^{(k)} \mathbf{M}^{(k)} \mathbf{P}^{(k)}$ to the internal distribution of the trajectories inside the 3-dimensional affine spaces. For the deviations outside the fitted spaces, we fit a Gaussian distribution of mean 0 and a constant variance $\hat{\sigma}_k^2$.

Using these distributions, we compute the probability $P(\alpha|k)$ of the trajectory vector \mathbf{p}_α conditioned to be in the k th class. Regarding $w^{(k)}$ in Eq. (7) as the a priori probability of the k th class, we compute the a posterior probability $W_\alpha^{(k)}$ by Eq. (14) using Bayes' theorem. Then, we reclassify all the trajectories according to $W_\alpha^{(k)}$, which are fractions in general (i.e., one trajectory belongs to multiple classes with fractional weights). This procedure is iterated until all the weights $W_\alpha^{(k)}$ converge. Finally, we associate the α th trajectory with the k th class that maximizes $W_\alpha^{(k)}$.

If we consider only the deviations outside the fitted spaces, the above procedure reduces to the non-Bayesian type. This type of unsupervised learning²

²This scheme is often referred to as the *EM algorithm* [2], because the mathematical structure is the same as estimating parameters from incomplete data by maximizing the logarithmic likelihood marginalized by the posterior of the missing data given by Bayes' theorem.

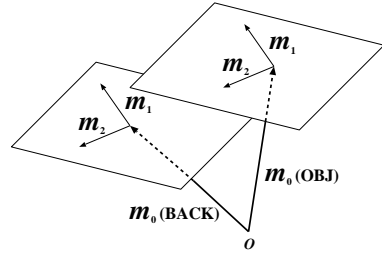


Figure 2: If the motions of the objects and the background are degenerate, their trajectory vectors belong to mutually parallel 2-dimensional affine spaces.

is widely used for pattern recognition, and the likelihood is known to increase monotonously in the course of iterations [13]. However, it is also well known that the iterations are very likely to be trapped at a local maximum. It is almost impossible to do correct segmentation by the above procedure alone unless we start from a very good initial value.

4. Degenerate Motion Model

Now, we model degenerate motions and derive an associated unsupervised learning procedure, from which we construct our multi-stage optimization procedure.

4.1 Degenerate motions

The motions we frequently encounter are such that the objects and the background are translating and rotating 2-dimensionally in the image frame with varying sizes.

For such a motion, we can choose the basis vector \mathbf{k}_κ in Eq. (2) in the Z direction (the camera optical axis is identified with the Z -axis). Under the affine camera model, motions in the Z direction do not affect the projected image except for its size. Hence, the vector $\tilde{\mathbf{m}}_{3\kappa}$ in Eq. (4) can be taken to be $\mathbf{0}$; the scale changes of the projected image are absorbed by the scale changes of $\tilde{\mathbf{m}}_{1\kappa}$ and $\tilde{\mathbf{m}}_{2\kappa}$ over time κ . It follows that the trajectory vector \mathbf{p}_α in Eq. (5) belongs to the 2-dimensional affine space passing through \mathbf{m}_0 and spanned by \mathbf{m}_1 and \mathbf{m}_2 .

All existing segmentation methods based on the shape interaction matrix of Costeira and Kanade [1] assume that the trajectories of different motions belong to independent 3-dimensional subspaces [8, 9]. Hence, degenerate motions cannot be correctly segmented.

If, in addition, the objects and the background do not rotate, we can fix the basis vectors \mathbf{i}_κ and \mathbf{j}_κ in Eq. (2) to be in the X and Y directions, respectively. Since the basis vectors \mathbf{i}_κ and \mathbf{j}_κ are common to the objects and the background, the vectors \mathbf{m}_1 and \mathbf{m}_2 in Eq. (5) are also common. Thus, the 2-dimensional affine spaces of all the motions are *parallel* (Fig. 2).

Note that two parallel 2-dimensional affine spaces can be included in a 3-dimensional affine space. Since the affine space separation method attempts to segment the trajectories into different 3-dimensional

affine spaces, it does not work if the objects and the background undergo such degenerate motions. This explains why the accuracy of the affine space separation is not as high as expected for real video sequences.

4.2 Unsupervised learning for degenerate motions

Since most of the motions we encounter in practice are degenerate, we can expect that the segmentation accuracy increases by unsupervised learning assuming such degenerate motions. The actual procedure goes as follows:

First, we set the weight $W_\alpha^{(k)}$ of \mathbf{p}_α for the k th class by Eq. (6). Next, we iterate the following procedures A, B, and C in turn until all the weights $\{W_\alpha^{(k)}\}$ converge:

A. Do the following computation for each class $k = 1, \dots, m$.

1. Compute $w^{(k)}$ by Eq. (7).
2. Compute the centroid $\mathbf{p}_C^{(k)}$ of the k th class by Eq. (8).
3. Compute the $n \times n$ moment matrix $\mathbf{M}^{(k)}$ by Eq. (9).

B. Do the following computation.

1. Compute the total $n \times n$ moment matrix

$$\mathbf{M} = \sum_{k=1}^m w^{(k)} \mathbf{M}^{(k)}. \quad (15)$$

2. Let $\lambda_1 \geq \lambda_2$ be the two largest eigenvalues of the matrix \mathbf{M} , and \mathbf{u}_1 and \mathbf{u}_2 the corresponding unit eigenvectors.
3. Compute the $n \times n$ projection matrices

$$\mathbf{P} = \sum_{i=1}^2 \mathbf{u}_i \mathbf{u}_i^\top, \quad \mathbf{P}_\perp = \mathbf{I} - \mathbf{P}. \quad (16)$$

4. Estimate the noise variance in the direction orthogonal to all the affine spaces by

$$\hat{\sigma}^2 = \max\left[\frac{\text{tr}[\mathbf{P}_\perp \mathbf{M} \mathbf{P}_\perp]}{n-2}, \sigma^2\right]. \quad (17)$$

5. Compute the $n \times n$ covariance matrix of the k th class by

$$\mathbf{V}^{(k)} = \mathbf{P} \mathbf{M}^{(k)} \mathbf{P} + \hat{\sigma}^2 \mathbf{P}_\perp. \quad (18)$$

C. Do the following computation for each trajectory vector \mathbf{p}_α , $\alpha = 1, \dots, N$.

1. Compute the conditional likelihood $P(\alpha|k)$, $k = 1, \dots, m$, by Eq. (13).
2. Recompute the weights $\{W_\alpha^{(k)}\}$, $k = 1, \dots, m$, by Eq. (14).

The computation is the same as in Sec. 3.2 except that 2-dimensional affine spaces with the same orientation are fitted; the common basis vectors \mathbf{u}_1 and \mathbf{u}_2

and the common outside noise variance are estimated in the procedure B.

After the iterations of A, B, and C have converged, the α th trajectory is classified to the k th class that maximizes $W_\alpha^{(k)}$, $k = 1, \dots, N$. The corresponding non-Bayesian scheme can be obtained if we do not consider the internal distributions.

4.3 Multi-stage optimization

In order to start the above learning, we need a good initial value. Here, we use the affine space separation using 2-dimensional affine spaces, which effectively assumes planar motions with varying sizes. The resulting segmentation is then optimized by assuming non-rotational motions.

The solution should be very accurate if the motions are truly degenerate. In reality, however, rotations may be involved to some extent. So, we relax the constraint and optimize the solution by assuming general 3-D motions.

In sum, our scheme consists of the following three stages:

1. Initial segmentation by the affine space separation using 2-dimensional affine spaces.
2. Unsupervised learning of the Bayesian type assuming degenerate motions.
3. Unsupervised learning of the Bayesian type assuming general 3-D motions.

This multi-stage optimization allows us to not only separate degenerate motions that we frequently encounter with high precision but also preserve the high performance for general 3-D motions, as we now show.

5. Experiments

5.1 Simulations

Fig. 3 shows three sequences of five synthetic images (supposedly of 512×512 pixels) of 14 object points and 20 background points; the object points are connected by line segments for the ease of visualization. To simulate real circumstances better, all the points are perspectively projected onto each frame with 30° angle of view, although the underlying theory is based on the affine camera model without perspective effects.

In all the three sequences, the object moves toward the viewer in one direction (10° from the image plane), while the background moves away from the viewer in another direction (10° from the image plane). In (a), the object and the background are simply translating in different directions. In (b) and (c), they are additionally given rotations by 2° per frame in opposite senses around different axes; they make 10° from the optical axis in (b) and 60° in (c). Thus, all the three motions are not strictly degenerate (with perspective effects), but the motion is almost degenerate in (a), nearly degenerate in (b), and a general 3-D motion in (c).

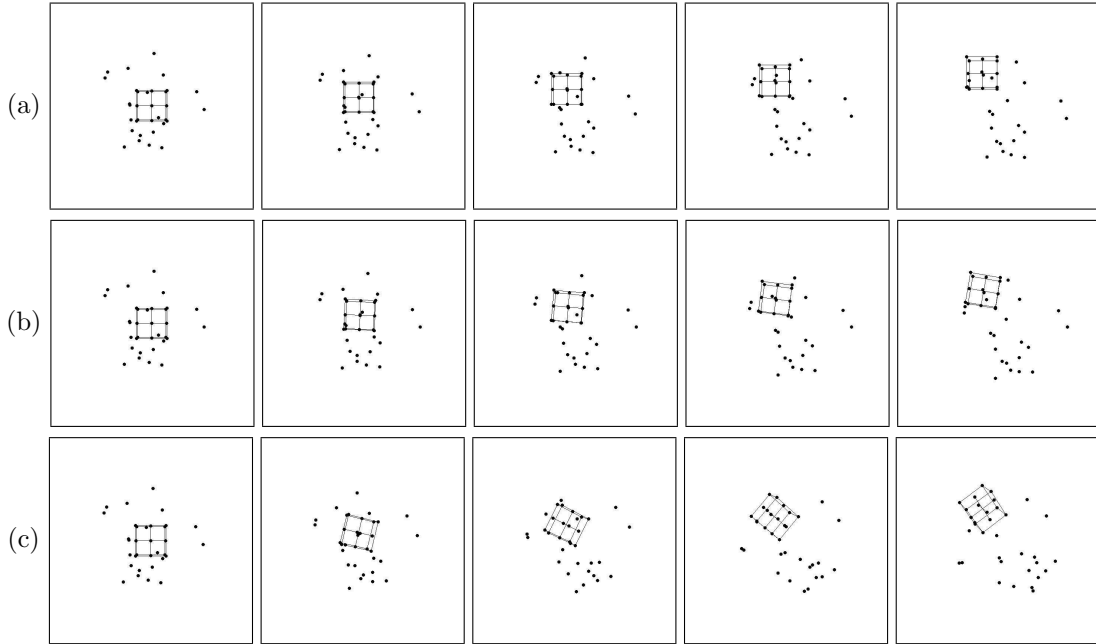


Figure 3: Simulated image sequences of 14 object points and 20 background points: (a) almost degenerate motion; (b) nearly degenerate motion; (c) general 3-D motion.

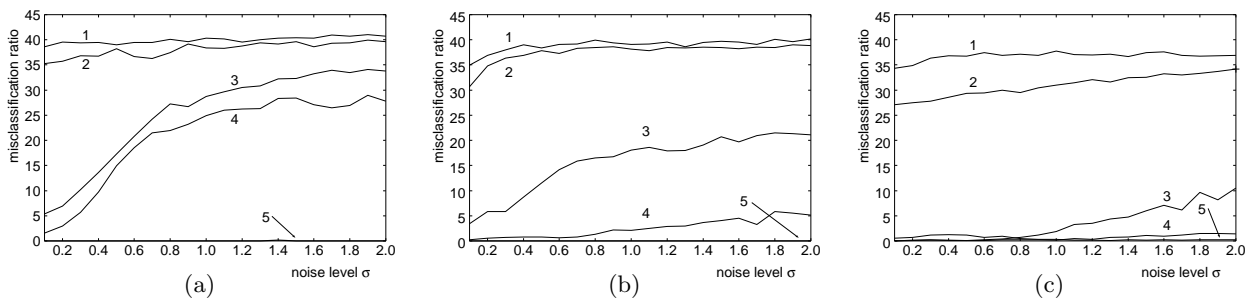


Figure 4: Misclassification ratio for the sequences (a), (b), and (c) in Fig. 3: 1) Costeira-Kanade; 2) Ichimura; 3) optimized subspace separation; 4) optimized affine space separation; 5) multi-stage optimization.

We added independent Gaussian random noise of mean 0 and standard deviation σ to the coordinates of all the points and segmented them into two groups. Fig. 4 plots the average misclassification ratio over 500 trials using different noise for different σ . We compared 1) the Costeira-Kanade method [1], 2) Ichimura’s method [4], 3) the subspace separation [8, 9] followed by unsupervised learning of the Bayesian type (we call this *optimized subspace separation* for short), 4) the affine space separation [10] followed by unsupervised learning of the Bayesian type (*optimized affine space separation* for short), and 5) our multi-stage optimization.

For the almost degenerate motion in Fig. 3(a), the optimized subspace separation and the optimized affine space separation do not work very well. Also, the affine space separation is not superior to the subspace separation (Fig. 4(a)). Since our multi-stage optimization is based on this type of degeneracy, it achieves 100% accuracy over all the noise range.

For the nearly degenerate motion in Fig. 3(b), the optimized subspace separation and the optimized affine space separation both work fairly well

(Fig. 4(b)). However, our method still attains almost 100% accuracy.

For the general 3-D motion in Fig. 3(c), the optimized subspace separation and the optimized affine space separation exhibit relatively high performance (Fig. 4(c)), but our method performs much better with nearly 100% accuracy again.

Although the same learning procedure is used in the end, the multi-stage optimization performs better than the optimal affine space separation, because the former starts from a better initial value than the latter. This is the reason why the multi-stage optimization achieves high performance even for considerably non-degenerate motions.

For all the motions, the Costeira-Kanade method performs very poorly. The accuracy is not 100% even in the absence of noise ($\sigma = 0$) because of the perspective effects. Ichimura’s method is not effective, either. It works to some extent for the general 3-D motion in Fig. 3(c), but it does not compare with the optimized subspace or affine space separation, much less with the multi-stage optimization method.

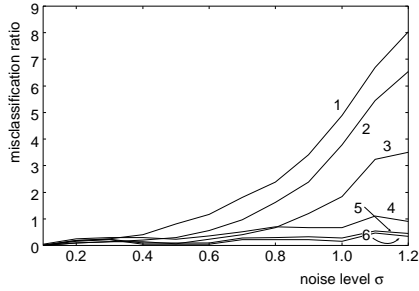


Figure 5: Effects of unsupervised learning for Fig. 3(c): 1) subspace separation; 2) subspace separation followed by unsupervised learning of the non-Bayesian type; 3) subspace separation followed by unsupervised learning of the Bayesian type; 4) affine space separation; 5) affine space separation followed by unsupervised learning of the non-Bayesian type 6) affine space separation followed by unsupervised learning of the Bayesian type.

5.2 Effects of learning

Fig. 5 shows the effects of unsupervised learning for Fig. 3(c). We plot the misclassification ratios of the subspace separation and the affine space separation with and without unsupervised learning of the non-Bayesian and Bayesian types. From Fig. 5, we can see that both the non-Bayesian and the Bayesian types work effectively but that the Bayesian type is slightly better. Yet, as we can see from Fig. 5, our multi-stage optimization is far superior to all other methods.

Fig. 6 shows the stage-wise effects of unsupervised learning of our multi-stage optimization for Fig. 3(c). For this general 3-D motion, the learning assuming degenerate motions does not perform so very well indeed, but the subsequent learning assuming general 3-D motions successfully restores the accuracy up to almost 100%.

The interesting fact is that the accuracy increases as the noise increases. This is perhaps because the discrepancy between the assumed degenerate motion and the actual non-degenerate motion is more conspicuous when the noise is smaller.

5.3 Real video examples

Fig. 7 shows five decimated frames from three video sequences A, B, and C (320×240 pixels). For each sequence, we detected feature points in the initial frame and tracked them using the Kanade-Lucas-Tomasi algorithm [17]. The marks \square indicate their positions. From the trajectories tracked throughout the sequence, we removed outlier trajectories using the method of Sugaya and Kanatani [16].

Table 1 lists the number of frames, the number of inlier trajectories, and the computation time for our multi-stage optimization. We reduced the computation time by compressing the trajectory data into 8-dimensional vectors [14]. We used Pentium 4 2.4B GHz for the CPU with 1 Gb main memory and Linux for the OS.

Table 2 lists the segmentation accuracies for different methods (“opt” stands for “optimized”). The accuracy is measured by (the number of correctly classi-

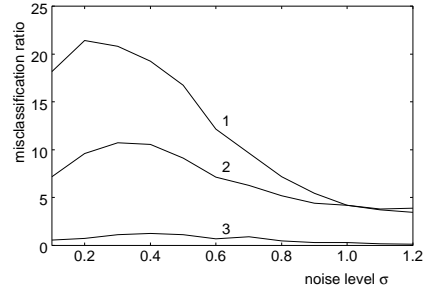


Figure 6: Effects of unsupervised learning for Fig.3(c): 1) affine space separation using 2-dimensional affine spaces; 2) unsupervised learning of the Bayesian type assuming degenerate motions; 3) unsupervised learning of the Bayesian type assuming general 3-D motions.

fied points)/(the total number of points) in percentage. For the methods other than Costeira-Kanade and Ichimura, this percentage is averaged over 50 trials, since both the subspace and the affine space separations internally use random sampling for robust estimation and hence the result is slightly different for each trial.

As we can see, the Costeira-Kanade method fails to produce meaningful segmentation. Ichimura’s method is effective for sequences A and B but not so very effective for sequence C. For sequence A, the affine space separation is superior to the subspace separation. For sequence B, the two methods have almost the same performance. For sequence C, in contrast, the subspace separation is superior to the affine space separation, strongly suggesting that the motion in sequence C is nearly degenerate.

The effect of learning is larger for sequence A than for sequences B and C, for which the accuracy is already high before the learning. Thus, the effect of unsupervised learning very much depends on the quality of the initial segmentation. For all the three sequences, our multi-stage optimization achieves 100% accuracy.

6. Concluding Remarks

In this paper, we have proposed multi-stage optimization by unsupervised learning assuming degenerate motions followed by unsupervised learning assuming general 3-D motions. Doing simulations and real video experiments, we have confirmed that our method is superior to all existing methods in realistic circumstances.

The reason for this superiority is that our method is tuned to realistic circumstances, where the motions of objects and backgrounds are almost degenerate, while existing methods mostly make use of the shape interaction matrix of Costeira and Kanade on the assumption that objects and backgrounds undergo general 3-D motions. As a result, they perform very poorly for simple motions such as in Fig. 7.

In contrast, our method³ has very high perfor-

³The source code is publicly available at: <http://www.suri.it.okayama-u.ac.jp/e-program.html>

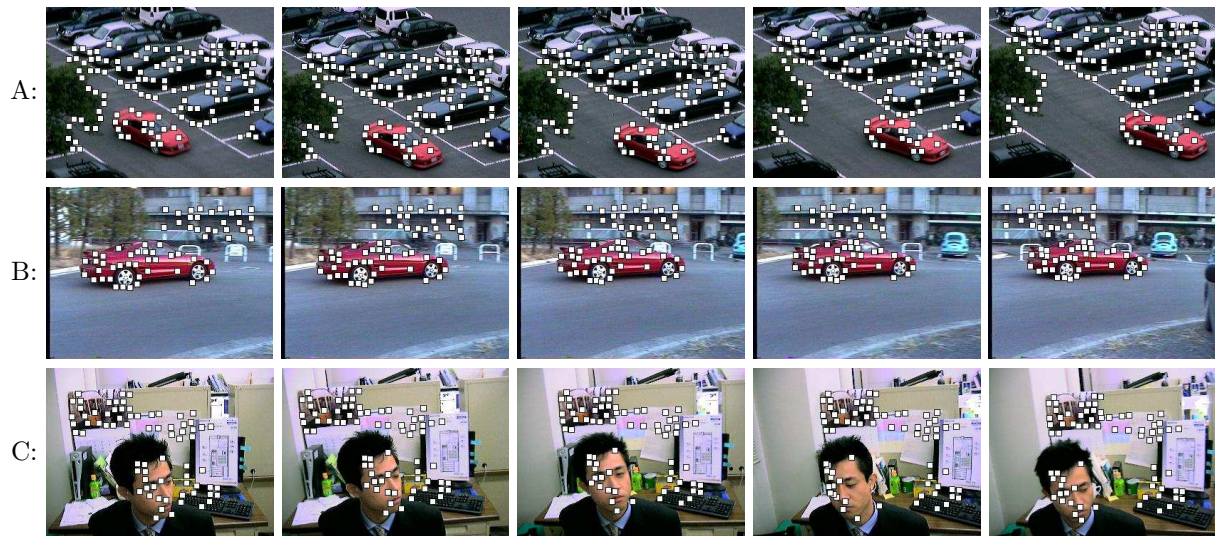


Figure 7: Three video sequences and successfully tracked feature points.

Table 1: The computation time for the multi-stage optimization of the sequences in Fig. 7.

| | A | B | C |
|------------------------|------|------|------|
| number of frames | 30 | 17 | 100 |
| number of points | 136 | 63 | 73 |
| computation time (sec) | 2.50 | 0.51 | 1.49 |

Table 2: Segmentation accuracy (%) for the sequences in Fig. 7.

| | A | B | C |
|------------------------------|--------------|--------------|--------------|
| Costeira-Kanade | 60.3 | 71.3 | 58.8 |
| Ichimura | 92.6 | 80.1 | 68.3 |
| subspace separation | 59.3 | 99.5 | 98.9 |
| affine space separation | 81.8 | 99.7 | 67.5 |
| opt. subspace separation | 99.0 | 99.6 | 99.6 |
| opt. affine space separation | 99.0 | 99.8 | 69.3 |
| multi-stage optimization | 100.0 | 100.0 | 100.0 |

mance for degenerate motions, and the accuracy is preserved even for considerably non-degenerate motions due to the multi-stage optimization.

Acknowledgments: This work was supported in part by the Ministry of Education, Culture, Sports, Science and Technology, Japan, under a Grant in Aid for Scientific Research C(2) (No. 15500113), the Support Center for Advanced Telecommunications Technology Research, and Kayamori Foundation of Informational Science Advancement.

References

- [1] J. P. Costeira and T. Kanade, A multibody factorization method for independently moving objects, *Int. J. Comput. Vision*, vol.29, no.3, pp.159–179, Sept. 1998.
- [2] A. P. Dempster, N. M. Laird and D. B. Rubin, Maximum likelihood from incomplete data via the EM Algorithm, *J. Roy. Statist. Soc., ser.B*, vol.39, pp.1–38, 1977.
- [3] C. W. Gear, Multibody grouping from motion images, *Int. J. Comput. Vision*, vol.29, no.2, pp.133–150, Aug./Sept. 1998.
- [4] N. Ichimura, Motion segmentation based on factorization method and discriminant criterion, *Proc. 7th Int. Conf. Comput. Vision*, Kerkyra, Greece, vol.1, pp.600–605, Sept. 1999.
- [5] N. Ichimura, Motion segmentation using feature selection and subspace method based on shape space, *Proc. 15th Int. Conf. Pattern Recog.*, Barcelona, Spain, vol.3, pp.858–864, Sept. 2000.
- [6] K. Inoue and K. Urahama, Separation of multiple objects in motion images by clustering, *Proc. 8th Int. Conf. Comput. Vision*, Vancouver, Canada, vol.1, pp.219–224, July 2001.
- [7] K. Kanatani, Geometric information criterion for model selection, *Int. J. Comput. Vision*, vol.26, no.3, pp.171–189, Feb./March 1998.
- [8] K. Kanatani, Motion segmentation by subspace separation and model selection, *Proc. 8th Int. Conf. Comput. Vision*, Vancouver, Canada, vol.2, pp.301–306, July 2001.
- [9] K. Kanatani, Motion segmentation by subspace separation: Model selection and reliability evaluation, *Int. J. Image Graphics*, vol.2, no.2, pp.179–197, April 2002.
- [10] K. Kanatani, Evaluation and selection of models for motion segmentation, *Proc. 7th Euro. Conf. Comput. Vision*, Copenhagen, Denmark, vol. 3, pp. 335–349, June 2002.
- [11] N. Otsu, A threshold selection method from gray-level histograms, *IEEE Trans. Sys. Man Cyber.*, vol.9, no.1, pp.62–66, Jan. 1979.
- [12] C. J. Poelman and T. Kanade, A paraperspective factorization method for shape and motion recovery, *IEEE Trans. Pat. Anal. Mach. Intell.*, vol.19, no.3, pp.206–218, March 1997.
- [13] M. I. Schlesinger and V. Hlaváč, *Ten Lectures on Statistical and Structural Pattern Recognition*, Kluwer, Dordrecht, The Netherlands, 2002.
- [14] Y. Sugaya and K. Kanatani, Automatic camera model selection for multibody motion segmentation, *Proc. Workshop on Science of Computer Vision*, Okayama, Japan, pp.31–39, Sept. 2002.
- [15] Y. Sugaya and K. Kanatani, Automatic camera model selection for multibody motion segmentation, *IAPR Workshop on Machine Vision Applications*, Nara, Japan, pp.412–415, Dec. 2002.
- [16] Y. Sugaya and K. Kanatani, Outlier removal for motion tracking by subspace separation, *IEICE Trans. Inf. & Syst.*, vol.E86-D, no.6, pp.1095–1102, June 2003.
- [17] C. Tomasi and T. Kanade, Detection and Tracking of Point Features, *CMU Tech. Rep. CMU-CS-91-132*, Apr. 1991: <http://vision.stanford.edu/~birch/klt/>
- [18] Y. Wu, Z. Zhang, T. S. Huang and J. Y. Lin, Multibody grouping via orthogonal subspace decomposition, sequences under affine projection, *Proc. IEEE Conf. Computer Vision Pattern Recog.*, vol.2, pp.695–701, Kauai, Hawaii, U.S.A., Dec. 2001.