

RECURSIVE ESTIMATION OF CAMERA MOTION FROM UNCALIBRATED IMAGE SEQUENCES

Stefano Soatto[†] and Pietro Perona[‡]

[†] California Institute of Technology 116-81, Pasadena-CA 91125

[‡] Università di Padova, Dipartimento di Elettronica ed Informatica, Padova-Italy
soatto@systems.caltech.edu

ABSTRACT

We describe a method for estimating the motion and structure of a scene from a sequence of images taken with a camera whose geometric calibration parameters are unknown.

The scheme is based upon a recursive motion estimation scheme, called the “essential filter” [16], extended according to the epipolar geometric representation presented in [7] in order to estimate the calibration parameters as well.

The motion estimates can then be fed into any “structure from motion” module that processes motion error, in order to recover the structure of the scene.

1. INTRODUCTION

Camera motion estimation is a key task in many applications ranging from image compression, to autonomous vehicle navigation, to recognition. Motion estimation from image sequences is usually performed in two steps: first the camera is calibrated, in order to establish metric relationships between world coordinates and image-plane measurements. The internal parameters (pixel size, optical center, focal length), are usually estimated *off-line*. Once calibration is performed, we can estimate camera motion and ambient structure recursively from the image sequence in a variety of ways [1, 14, 19, 12].

Most of the recursive motion estimation schemes rely upon the exact knowledge of internal camera parameters. However, experimental evidence shows that these can change drastically during a long sequence [4]

Research funded by the California Institute of Technology, a scholarship from the University of Padova, a fellowship from the “A. Gini” Foundation, an AT&T Foundation Special Purpose grant, ONR grant N0014-93-1-0990, grant ASI-RS-103 from the Italian Space Agency and the NSF National Young Investigator Award (P.P.). This work is registered as CDS Technical Report CIT-CDS 94-005, California Institute of Technology, January 1994 – revised February 1994.

due to zooming and changing of the aperture. Moreover, *often it is not possible to access the physical device* which produced the sequence. Therefore, a motion estimation scheme should be able to estimate camera calibration while processing the sequence and estimating motion and structure.

Many approaches for camera calibration are available in the literature; they can roughly be classified as:

1. Batch schemes, which rely on the knowledge of the *structure* by including a calibration rig in the field of view (see [10] and references).
2. Active devices, which rely on the knowledge of the camera *motion* by controlling the configuration (pose) of the camera [5, 4, 3].
3. *Arbitrary structure and motion*. Camera self calibration is performed along with motion estimation [7].

The first two approaches assume that the camera is available for measurements, by either controlling its motion or inserting a known object into the field of view. Therefore it seems that the third approach is the only feasible solution when the device which produced the sequence is not available, as for example in image compression applications or automation of image processing tasks for the movie industry.

Faugeras et al. [7] propose a *batch* scheme which reconstructs the epipolar transformation of the camera, and then imposes the structure of such a transformation by solving a set of polynomial equations, known as Kruppa’s equations. However, the scheme has some substantial drawbacks which make it unattractive for real world applications. In particular

- High sensitivity to pixel-noise
- Numerical instability

- Motion parameters and internal parameters are treated alike. While camera-motion can vary arbitrarily during a sequence, it is conceivable that some parameters (for example the pixel size or aspect ratio) are constant over long periods of time
- Not all the information coming from a sequence is exploited. The scheme processes 3 images at a time and does not use temporal coherence (recursion) or a-priori information (such as reference values for focal length, initial confidence in the position of the optical center etc.).

Hence we want a recursive scheme which, after each incoming image, updates the computation performed at the previous step. We also want the scheme to be *causal* so that it can be used for real-time implementations. Azarbayejani et al. [1] perform *partial calibration* by updating the focal length of the camera on-line together with camera motion. To our knowledge, the problem of estimating camera motion and calibration recursively from an image sequence has never been addressed in the literature before.

In this paper we present a scheme for performing ego-motion estimation and camera calibration recursively and causally for an image sequence. It does not need a calibration rig nor to control motion, while it exploits redundancy at each step and computations from each previous step by recursion. A priori information about calibration can be used, if available, as initial conditions for the estimation scheme. Internal parameter time constants are adjustable by tuning their random walk models.

The scheme is based upon a recent method for recursive motion estimation [16], extended to estimate camera parameters according to the representation of [7]. A key feature of our scheme is that *the structure of the epipolar geometry is imposed explicitly as the structure of the state-space of the filter*, so we do not need to solve *explicitly* complicated polynomial equations in order to enforce such a structure. From a different point of view, our filter can be viewed as a recursive differential scheme for solving Kruppa's equations.

We report some experiments on noisy synthetic image sequences, and are in the process of testing the scheme on real image sequences.

2. FORMULATION OF THE SCHEME

2.1. CAMERA MODEL: INTERNAL PARAMETERS AND EGO-MOTION

The camera may be modeled as a perspective projection map

$$\begin{aligned} M : \mathbb{R}^3 &\rightarrow \mathbb{R}^2 \\ \mathbf{X} &\mapsto \mathbf{x}. \end{aligned} \quad (1)$$

The simplest instance is the so called "ideal pinhole model":

$$\mathbf{X} \doteq \begin{bmatrix} X \\ Y \\ Z \end{bmatrix}^T \mapsto \begin{bmatrix} x \\ y \end{bmatrix}^T \doteq \begin{bmatrix} \frac{X}{Z} \\ \frac{Y}{Z} \end{bmatrix}^T \doteq \mathbf{x}. \quad (2)$$

It can also be represented as a *linear map* between real projective spaces, $\bar{M} : \mathbb{RP}^3 \rightarrow \mathbb{RP}^2$: in homogeneous coordinates it is represented by a 3×4 matrix $\begin{bmatrix} \mathbf{A} & | & 0 \end{bmatrix}$ where

$$\mathbf{A} \doteq \begin{bmatrix} f s_x & 0 & -i_0 \\ 0 & f s_y & -j_0 \\ 0 & 0 & 1 \end{bmatrix}$$

is the internal parameter matrix. f is the focal length, (i_0, j_0) the coordinates of the optical center and (s_x, s_y) the pixel sizes along the image plane coordinates. The deviation from 90° of the angle between the optical axis and the CCD surface is usually on the order of 1° , and we may therefore neglect it.

As the camera moves inside the (static) scene, the points move in its reference according to the rigid motion constraint:

$$\mathbf{X}(t+1) = R(t)\mathbf{X}(t) + T(t),$$

where $(T(t), R(t))$ represent the discrete camera motion between the time t and $t+1$. The goal of a self-calibrating motion scheme is to estimate the internal parameters and the camera motion from the time-varying projection $\mathbf{x}(t)$ of a number of feature points.

2.2. THE ESSENTIAL CONSTRAINT AND EPIPOLAR GEOMETRY

Longuet-Higgins [11] introduced a simple coplanarity constraint which links the projective coordinates $\bar{\mathbf{x}} \doteq \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$ of a point at time t , the corresponding $\bar{\mathbf{x}}'$ at $t+1$, and the motion $(T(t), R(t))$ undergone by the camera:

$$\bar{\mathbf{x}}_i'^T \mathbf{Q} \bar{\mathbf{x}}_i = 0 \quad \forall i = 1 \dots N. \quad (3)$$

where $\mathbf{Q} \doteq R(T\Lambda)$ is called the essential matrix. The essential matrices form a space that has the structure of an algebraic variety [6] and of a differentiable manifold as well [18]. Given a number of such constraints, it is possible to estimate the motion which generated it [11, 20, 13, 6]. It can be proved easily that a 3×3 matrix is essential if and only if it has two equal singular values and zero determinant [13].

In the case of an uncalibrated camera, a similar constraint can be derived based on the epipolar geometry: given $\bar{\mathbf{x}}(t)$ at time t , its correspondent at $t+1$, $\bar{\mathbf{x}}(t+1)$, must lie on the epipolar line ${}^t\mathbf{e}_{t+1}$. Such a line is described in projective coordinates by a linear function of $\bar{\mathbf{x}}(t)$. The representing matrix is called the *fundamental matrix* \mathbf{F} , which is defined by the relation ${}^t\mathbf{e}_{t+1} \doteq \mathbf{F}\bar{\mathbf{x}}(t)$. It can be shown [7] that $\mathbf{F} \doteq \mathbf{A}^{-T}\mathbf{Q}\mathbf{A}^{-1}$, where \mathbf{Q} is an essential matrix. From the definition of the epipolar line, one may derive a generalization of the essential constraint [7]:

$$\bar{\mathbf{x}}_i^T \mathbf{F} \bar{\mathbf{x}}_i = 0 \quad \forall i = 1 \dots N. \quad (4)$$

The scheme presented in [7] consists in first estimating \mathbf{F} from (4), and then imposing its structure a-posteriori by solving the Kruppa equations, which correspond to enforcing the fact that $\mathbf{A}^T \mathbf{F} \mathbf{A}$ (is essential and therefore) has two equal singular values and zero determinant.

2.3. THE ESSENTIAL FILTER FOR FUNDAMENTAL MATRICES

The essential filter is a motion estimation paradigm recently presented in [16]. It solves motion estimation as identification of the exterior differential system determined by the essential constraint:

$$\begin{cases} \bar{\mathbf{x}}_i^T(t+1)\mathbf{Q}(t)\bar{\mathbf{x}}_i(t) = 0 \\ \bar{\mathbf{x}}_i(t) = \bar{\mathbf{x}}_i(t) + \mathbf{n}_i(t) \end{cases} \quad \forall i = 1 : N. \quad (5)$$

We propose to extend the essential filter to estimate fundamental matrices, and *impose the structure of the fundamental matrix explicitly* by writing the estimator in local coordinates: the estimate at each step determines a matrix which is *fundamental by construction*, and we do not need to enforce the structure by solving explicitly ill-conditioned polynomial equations. The structure of resulting update is very similar to the essential filter [16]:

$$\begin{bmatrix} \hat{\xi} \\ \hat{T} \\ \hat{R} \end{bmatrix} (t+1) = \begin{bmatrix} \hat{\xi} \\ \hat{T} \\ \hat{R} \end{bmatrix} (t) +$$

$$+L(t) \begin{bmatrix} \vdots \\ \bar{\mathbf{x}}_i^T(t)\mathbf{A}^{-T}(\hat{\xi})\mathbf{Q}(\hat{T},\hat{R})\mathbf{A}^{-1}(\hat{\xi})\bar{\mathbf{x}}_i(t-1) \\ \vdots \end{bmatrix} \quad (6)$$

where $\xi \doteq [fs_x, fs_y, i_0, j_0]^T$; L has the structure of the gain of an Implicit Extended Kalman Filter (IEKF) [9, 8, 16].

If we call $\alpha \doteq [\xi \ \theta \ \phi \ \Omega]^T \in \mathbf{R}^9$, where Ω are the exponential coordinates of $R = e^{\Omega\Lambda}$, and (θ, ϕ) are the spherical coordinates of T , then we can write the complete set of equations for the filter:

Prediction step

$$\begin{cases} \hat{\alpha}(t+1|t) = \hat{\alpha}(t|t) & \hat{\alpha}(0|0) = \alpha_0 \\ P(t+1|t) = P(t|t) + R_\alpha(t) & P(0|0) = P_0 \end{cases}$$

Update step

$$\begin{cases} \hat{\alpha}(t+1|t+1) = \hat{\alpha}(t+1|t) + \\ L(t+1)\bar{\mathbf{x}}_i^T(t)\mathbf{A}^{-T}\mathbf{Q}(\hat{\alpha}(t+1|t))\mathbf{A}^{-1}\bar{\mathbf{x}}_i(t-1) \\ P(t+1|t+1) = \\ \Gamma(t+1)P(t+1|t)\Gamma^T(t+1) + \\ L(t+1)D_+(t)R_n(t+1)D_+^T(t)L^T(t+1) \end{cases}$$

where

$$\begin{cases} L(t+1) = P(t+1|t)C^T(t+1)\Lambda^{-1}(t+1) \\ \Lambda(t+1) = C(t+1)P(t+1|t)C^T(t+1) + \\ D_+(t+1)R_n(t+1)D_+^T(t+1) \\ \Gamma(t+1) = I - L(t+1)C(t+1) \\ D_+(t+1) \doteq \left(\frac{\partial \bar{\mathbf{x}}_i^T(t+1)\mathbf{Q}(t)\bar{\mathbf{x}}_i(t)}{\partial \mathbf{x}(t+1)} \right)_{|\bar{\mathbf{x}}(t), \hat{\alpha}(t)} \\ C(t+1) \doteq \left(\frac{\partial \bar{\mathbf{x}}_i^T(t+1)\mathbf{Q}(t)\bar{\mathbf{x}}_i(t)}{\partial \alpha(t)} \right)_{|\bar{\mathbf{x}}(t), \hat{\alpha}(t)} \end{cases}$$

where R_α and R_n denote the variance of the noises $\alpha(t)$ and $n(t)$ respectively; the interested reader may find the detailed derivation in [17]¹.

3. EXPERIMENTAL ASSESSMENT

We report a set of simulations on a noisy synthetic sequence. In figure 1 we show the estimates of the translation and rotation parameters. In figure 2 we show the estimates of the internal parameters. The noise on the image-plane was one tenth of a pixel, according to the performance of the best optical flow/feature tracking techniques [2]. Convergence is reached in about 100 frames. Each iteration consists of about 100 Kflops: an implementation using Matlab (not optimized) runs at .6Hz on a Sparc 10-20. We are currently experimenting

¹This paper can be obtained via the Worldwide Net Mosaic (<http://avalon.caltech.edu/cds/techreports/>)

on real image sequences and higher noise levels. More detailed experiments are reported in [15].

Note that, once the motion has been reconstructed, we may feed the estimates onto any Structure-From-Motion module that processes motion error [14, 19]. However, the motion configurations that allow estimating accurately the scene structure, as for example fronto-parallel translation, are often not sufficiently exciting for estimating the camera parameters. Vice-versa, motions that allow a good estimation of the camera calibration are often ill-conditioned for estimating depth, as for example a spiral along the optical axis. Therefore there is an intrinsic conflict between the estimation of the camera parameters and the structure of the scene.

ACKNOWLEDGEMENTS

We wish to thank Prof. Ruggero Frezza and Prof. Giorgio Picci for their constant support and suggestions.

4. REFERENCES

- [1] A. Azarbayejani, B. Horowitz, and A. Pentland. Recursive estimation of structure and motion using relative orientation constraints. *Proc. CVPR*, New York, 1993.
- [2] J. Barron, D. Fleet, and S. Beauchemin. Performance of optical flow techniques. RPL-TR 9107, Queen's University Kingston, Ontario, Robotics and perception laboratory, 1992. Also in *Proc. CVPR 1992*, pp 236-242.
- [3] A. Basu. Active calibration: alternative strategy and analysis. *Trans. of the IEEE conf. CVPR*, 1993.
- [4] J.L. Crowley, P. Bobet, and C. Schmidt. Maintaining stereo calibration by tracking image points. *Trans. of the IEEE conf. CVPR*, 1993.
- [5] F. Du and M. Brady. Self-calibration of the intrinsic parameters of cameras for active vision systems. *Trans. of the IEEE conf. CVPR*, 1993.
- [6] O. Faugeras. *Three dimensional vision, a geometric viewpoint*. MIT Press, 1993.
- [7] O.D. Faugeras, Q.T. Luong, and S.J. Maybank. Camera self-calibration: theory and experiments. *Proc. of the ECCV92, Vol. 588 of LNCS, Springer Verlag*, 1992.
- [8] A.H. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, 1970.
- [9] R.E. Kalman. A new approach to linear filtering and prediction problems. *Trans. of the ASME-Journal of basic engineering.*, 35-45, 1960.
- [10] R.K. Lenz and R.Y. Tsai. Techniques for calibration of the image center and scale factor for high accuracy 3d vision. *IEEE Trans. Pattern Anal. Mach. Intell.*, Sept. 1988.
- [11] H. C. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, 293:133-135, 1981.
- [12] L. Matthies, R. Szeliski, and T. Kanade. Kalman filter-based algorithms for estimating depth from image sequences. *Int. J. of computer vision*, 1989.

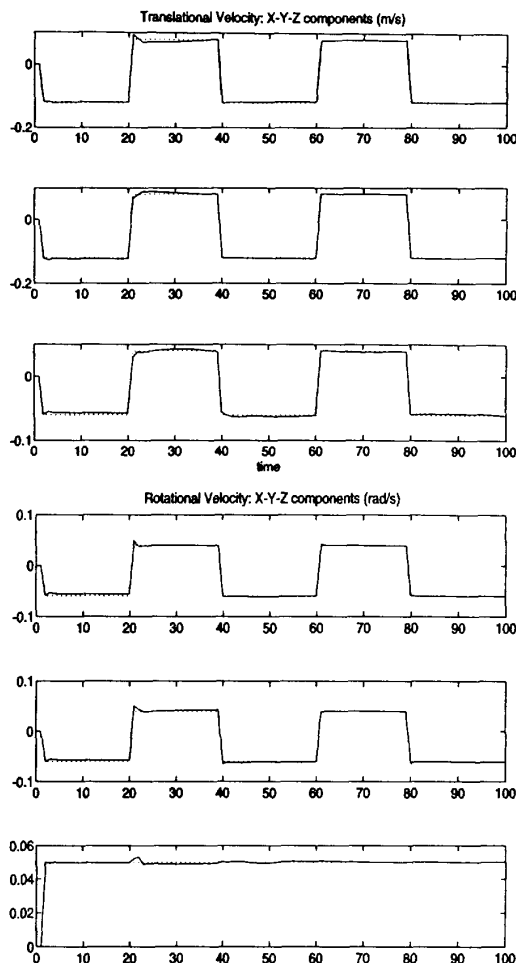


Figure 1: (Top) Translational velocity: filter estimates (solid) vs. true values (dotted). (Bottom) Components of rotational velocity.

- [13] S. Maybank. *Theory of reconstruction from image motion*. Springer Verlag, 1992.
- [14] J. Oliensis and J. Inigo-Thomas. Recursive multi-frame structure from motion incorporating motion error. *Proc. DARPA Image Understanding Workshop*, 1992.
- [15] S. Soatto, R. Frezza, and P. Perona. Recursive estimation of camera motion from uncalibrated image sequences. *Technical Report CIT-CDS 94-003*, California Institute of Technology, 1994.
- [16] S. Soatto, R. Frezza, and P. Perona. Motion estimation on the essential manifold. In *"Computer Vision ECCV 94, Lecture Notes in Computer Sciences vol. 801"*, Springer Verlag, May 1994.
- [17] S. Soatto, R. Frezza, P. Perona, and G. Picci. Motion estimation via dynamic vision. *Technical Report CIT-CDS 94-004*, California Institute of Technology. Reduced version to appear in the proceeding of the 33rd IEEE conf. on Decision and Control. Submitted to the IEEE transactions on Automatic Control, Feb. 1994.
- [18] S. Soatto and P. Perona. Structure-independent visual motion control on the essential manifold. *Technical Report CIT-CDS-94-013*, California Institute of Technology. Short version in *Proc. of the IFAC Symposium on Robot Control, Capri-Italy*, 1994.
- [19] S. Soatto, P. Perona, R. Frezza, and G. Picci. Recursive motion and structure estimation with complete error characterization. In *Proc. IEEE Comput. Soc. Conf. Comput. Vision and Pattern Recogn.*, pages 428–433, New York, June 1993.
- [20] J. Weng, T. Huang, and N. Ahuja. Motion and structure from two perspective views: algorithms, error analysis and error estimation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 11(5):451–476, 1989.

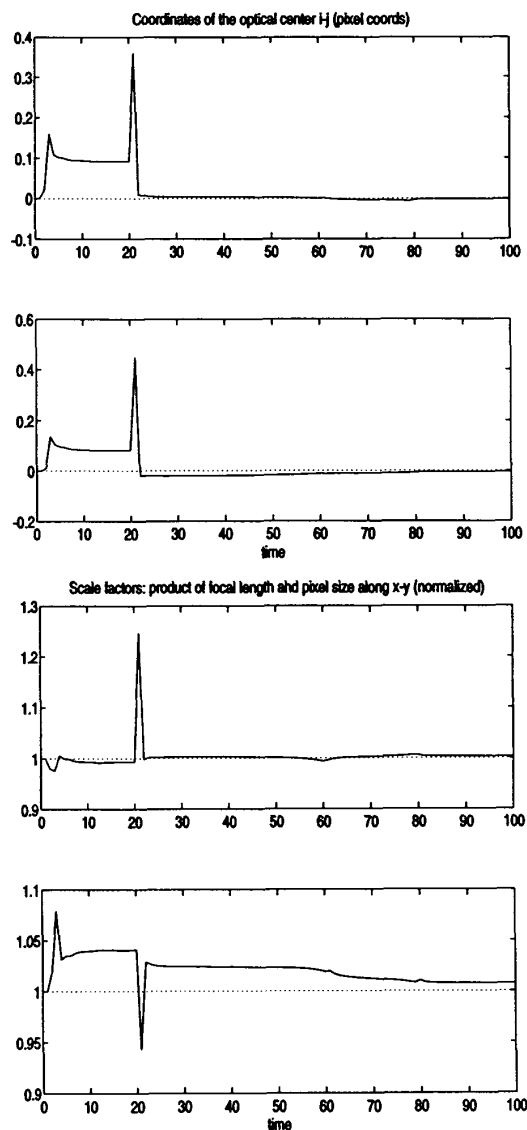


Figure 2: (Top) Coordinates of the center of projection: filter estimates vs. true values. (Bottom) Pixel size along image coordinates