

Collinearity and Coplanarity Constraints for Structure from Motion

Gang Liu¹, Reinhard Klette², and Bodo Rosenhahn³

¹ Institute of Information Sciences and Technology, Massey University, New Zealand,
Department of Computer Science, The University of Auckland, New Zealand

² Max Planck Institute Saarbrücken, Germany

Abstract. Structure from motion (SfM) comprises techniques for estimating 3D structures from uncalibrated 2D image sequences. This work focuses on two contributions: Firstly, a stability analysis is performed and the error propagation of image noise is studied. Secondly, to stabilize SfM, we present two optimization schemes by using a priori knowledge about collinearity or coplanarity of feature points in the scene.

1 Introduction

Structure from motion (SfM) is an ongoing research topic in computer vision and photogrammetry, which has a number of applications in different areas, such as e-commerce, real estate, games and special effects. It aims at recovering 3D (shape) models of (usually rigid) objects from an (uncalibrated) sequence (or set) of 2D images.

The original approach [5] of SfM consists of the following steps: (1) extract corresponding points from pairs of images, (2) compute the fundamental matrix, (3) specify the projection matrix, (4) generate a dense depth map, and (5) build a 3D model. A brief introduction of some of those steps will be presented in Section 2.

Errors are inevitable to every highly complex procedure depending on real-world data, and this also holds for SfM. To improve the stabilization of SfM, two optimizations are proposed using information from the 3D scene; see Section 3. Section 4 presents experimental results, and Section 5 concludes the paper with a brief summary.

2 Modules of SfM

This section gives a brief introduction for some of the SfM steps (and related algorithms). For extracting correspondent points, we recall a method proposed in [14]. Then, three methods for computing the fundamental matrix are briefly introduced. To specify a projection matrix from a fundamental matrix, we describe two common methods based on [3, 4]. In this step we also use the knowledge of intrinsic camera parameters, which can be obtained through Tsai calibration [12]; this calibration is performed before or after taking the pictures for the used

camera. It allows to specify the effective focal length f , the size factors k_u and k_v of CCD cells (for calculating the physical size of pixels), and the coordinates u_0 and v_0 of the principal point (i.e., center point) in the image plane.

2.1 Corresponding points

We need a number of at least seven pairs of corresponding points to determine the geometric relationship between two images, caused by viewing the same object from different view points. One way to extract those points from a pair of images is as follows [14]:

(i) extract candidate points by using the Harris corner detector [2], (ii) utilize a correlation technique to find matching pairs, and (iii) remove outliers by using a LMedS (i.e., least-median-of-squares) method.

Due to the poor performance of the Harris corner detector on specular objects, this method is normally not suitable.

2.2 Fundamental and Essential matrix

A fundamental matrix is an algebraic representation of epipolar geometry [13]. It can be calculated if we have at least seven correspondences (i.e., pairs of corresponding points), for example using linear methods (such as the *8-Point Algorithm* of [8]) or nonlinear methods (such as the *RANSAC Algorithm* of [1], or the *LMedS Algorithm* of [14]).

In the case of a linear method, the fundamental matrix is specified through solving an overdetermined system of linear equations utilizing the given correspondences. In the case of a nonlinear method, subsets (at least seven) of correspondences are chosen randomly and used to compute candidate fundamental matrices, and then the best is selected, which causes the smallest error for all the detected correspondences.

According to our experiments, linear methods have a more time efficient and provide reasonably good results for large (say more than 13) numbers of correspondences. Nonlinear methods are more time consuming, but less sensible to noise, especially if correspondences also contain outliers.

For given intrinsic camera parameters K_1 and K_2 , the Essential matrix E can be derived from F by computing

$$E = K_2^T F K_1$$

2.3 Projection matrix

A projection matrix P can be expressed as follows:

$$P = K[R \mid -RT]$$

where K is a matrix of the intrinsic camera parameters, and R and T are the rotation matrix and translation vector (the extrinsic camera parameters). Since the

intrinsic parameters are specified by calibration, relative rotation and translation can be successfully extracted from the fundamental matrix F . When recovering the projection matrices in reference to the first camera position, the projection matrix of the first camera position is given as $P_1 = K_1[I \mid 0]$, and the projection matrix of the second camera position is given as $P_2 = K_2[R \mid -RT]$.

The method proposed by Hartley and Zisserman for computing rotation matrix R and translation vector T (from the essential matrix E) is as follows [3]:

1. compute E by using $E = K_2^T F K_1$, where

$$K_i = \begin{pmatrix} f k_u & 0 & u_0 \\ 0 & f k_v & v_0 \\ 0 & 0 & 1 \end{pmatrix}$$

(note: $K_1 = K_2$ if we use the same camera at view points 1 and 2),

2. perform a singular value decomposition (SVD) of E by following the template $E = U \text{diag}(1, 1, 0) V^T$,
3. compute R and T (for the second view point), where we have two options, namely

$$\begin{aligned} R_1 &= U W V^T & R_2 &= U W^T V^T \\ T_1 &= u_3 & T_2 &= -u_3 \end{aligned}$$

where u_3 is the third column of U and

$$W = \begin{pmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Another method for computing R and T from E (also only using elementary matrix operations) is given in [4], which leads to almost identical results as the method by Hartley and Zisserman.

2.4 Dense depth map

At this point, the given correspondences allow only a few points to be reconstructed in 3D. A satisfactory 3D model of a pictured object requires a dense map of correspondences. The epipolar constraint (as calculated above) allows that correspondence search can be restricted to one-dimensional epipolar lines, it supports that images are at first rectified following the method in [10], and that correspondence matching is then done by searching along a corresponding scan line in the rectified image. We also require a recovered base line between both camera positions to calculate a dense depth map.

3 Optimization with Prior Knowledge

Since computations of fundamental and projection matrix are sensitive to noise, it is necessary to apply a method for reducing the effect of noise (to stabilize SfM). We utilize information about the given 3D scene, such as knowledge about collinearity or coplanarity.

3.1 Knowledge about collinearity

It is not hard to detect collinear points on man-made objects, such as buildings or furniture. Assuming ideal central projection (i.e., no lens distortion or noise), then collinear points in object space are mapped onto one line in the image plane. We assume that lens distortions are small enough to be ignored. Linearizing points which are supposed to be collinear can then be seen as a way to remove noise.

Least-square line fitting (minimizing perpendicular offsets) is used to identify the approximating line for a set of “noisy collinear points”. Assume that we have such a set of points $P = \{(x_i, y_i) | i = 1, \dots, n\}$ which determines a line $l(\alpha, \beta, \gamma) = \alpha x + \beta y + \gamma$. The coefficients α, β and γ are calculated as follows [7]:

$$\alpha = \frac{\mu_{xy}}{\sqrt{\mu_{xy}^2 + (\lambda^* - \mu_{xx})^2}}$$

$$\beta = \frac{\lambda^* - \mu_{xx}}{\sqrt{\mu_{xy}^2 + (\lambda^* - \mu_{xx})^2}}$$

$$\gamma = -(\alpha \bar{x} + \beta \bar{y})$$

where

$$\lambda^* = \frac{1}{2}(\mu_{xx} + \mu_{yy} - \sqrt{(\mu_{xx} - \mu_{yy})^2 + 4\mu_{xy}})$$

$$\mu_{xx} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad \mu_{yy} = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$\mu_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}), \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

After specifying the line, the points’ positions are modified through perpendicular projection onto the line.

3.2 Knowledge about coplanarity

Coplanar points can be expected on rigid structures such as on walls or on a tabletop. For a set of points, all incident with the same plane, there is a 3×3 matrix H called *homography* which defines a perspective transform of those points into the image plane [11].

Homography Consider we have an image sequence (generalizing the two-image situation from before) and p_{ki} is the projection of 3D point P_i into the k th image, i.e. P_i is related to p_{ki} as follows:

$$p_{ki} = \omega_{ki} K_k R_k (P_i - T_k) \quad (1)$$

where ω_{ki} is an unknown scale factor, K_k denotes the intrinsic matrix (for the used camera), and R_k and T_k are the rotation matrix and translation vector. Following Equation (1), P_i can be expressed as follows:

$$P_i = \omega_{ki}^{-1} R_k^{-1} K_k^{-1} p_{ki} + T_k \quad (2)$$

Similarly, for point p_{li} lying on the l th image, we have

$$P_i = \omega_{li}^{-1} R_l^{-1} K_l^{-1} p_{li} + T_l \quad (3)$$

From Equations (2) and (3), we get

$$p_{ki} = \omega_{ki} K_k R_k (\omega_{li}^{-1} R_l^{-1} K_l^{-1} p_{li} + T_l - T_k) \quad (4)$$

With $R_{kl} = R_k R_l^{-1}$ we define $H_{kl}^\infty = K_k R_{kl} K_l^{-1}$. We also have epipole $e_{kl} = K_k R_k (T_l - T_k)$. Equation (4) can then be simplified to

$$p_{ki} = \omega_{ki} \omega_{li}^{-1} (H_{kl}^\infty p_{li} + \omega_{li} e_{kl}) \quad (5)$$

H_{kl}^∞ is what we call the homography which maps points at infinity ($\omega_{li} = 0$) from image l to image k . Consider a point P_i on plane $\hat{n}^T P_i - d = 0$. Then, from Equation (3), we have

$$\hat{n}^T P_i - d = \hat{n}^T \omega_{li}^{-1} R_l^{-1} K_l^{-1} p_{li} + \hat{n}^T T_l - d = 0$$

Then we have

$$\omega_{li} = \frac{\hat{n}^T R_l^{-1} K_l^{-1} p_{li}}{d - \hat{n}^T T_l}$$

what can be rewritten as follows:

$$\omega_{li} = d_l^{-1} \hat{n}^T R_l^{-1} K_l^{-1} p_{li}$$

where $d_l^{-1} = d - \hat{n}^T T_l$ is the distance from the camera center (principal point) of the l th image to the plane (\hat{n}, d) . Substituting ω_{li} into Equation (5), finally we have

$$p_{ki} = \omega_{ki} \omega_{li}^{-1} (H_{kl}^\infty + d_l^{-1} e_{kl} \hat{n}^T R_l^{-1} K_l^{-1}) p_{li}$$

Let

$$H = \omega_{ki} \omega_{li}^{-1} (H_{kl}^\infty + d_l^{-1} e_{kl} \hat{n}^T R_l^{-1} K_l^{-1})$$

This means: points lying in the same plane have identical H which can be utilized as coplanarity constraint; see [11].

Coplanarity optimization Coplanar points satisfy the relation described by homography. We use this relation for modifying ‘‘noisy coplanar points,’’ using the following equation:

$$p_{ki} = H_{kl} p_{li}$$

Here, H_{kl} is the homography between k th and l th image in the sequence, and p_{ki} , p_{li} are projections of point P_i on the k th and l th image, respectively.

4 Experiments and Analysis

To analyze the influence of noise, we perform SfM in a way as shown in Figure 1. At Step 1, Gaussian noise is introduced into coordinates of detected correspondences. At step 2, three different methods are compared to specify which one is the best to compute the fundamental matrix. At Step 3, a quantitative error analysis is performed.

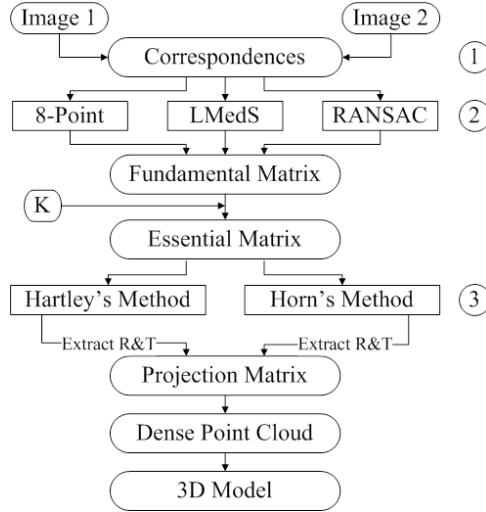


Fig. 1. Basic steps of SfM.

This section shows at first experiments of the performance of different methods for computing the fundamental matrix, and second the effect of those optimizations mentioned in the previous section.

4.1 Computation of fundamental matrix

Three algorithms (8-Point, RANSAC and LMedS) are compared with each other in this section. To specify the most stable one in presence of noise, Gaussian Noise (with mean 0 and deviation $\delta = 1$ pixel) and one outlier are propagated to given correspondences. Performances of the three algorithms are characterized in Figure 2: due to the outlier, the 8-Point Algorithm is more sensible than the other two.

4.2 Optimizations

To test the effect of the optimizations mentioned in the previous section, the results of splitting essential matrices (rotation matrices and translation vectors)

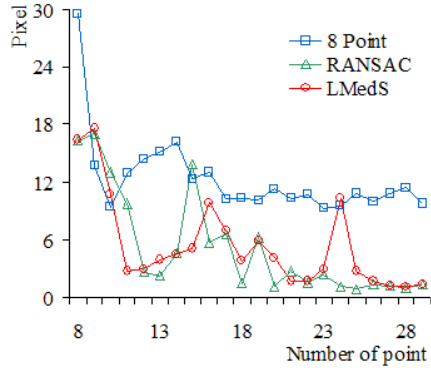


Fig. 2. Performance of three algorithms in presence of noise.

are utilized to compare with each other. Two images of a calibration object are used as test images (shown in Figure 3). The data got from calibration (intrinsic parameters and extrinsic parameters of camera) are used as the ground truth. Roll angle α , pitch angle β and yaw angle γ are used to compare the rotation matrices in a quantitative manner. These angles can be computed from a rotation matrix R by following equations [9]:

$$\begin{aligned}\alpha &= \text{atan2}\left(\frac{r_{23}}{\sin(\gamma)}, \frac{r_{13}}{\sin(\gamma)}\right) \\ \beta &= \text{atan2}\left(\frac{r_{32}}{\sin(\gamma)}, \frac{-r_{31}}{\sin(\gamma)}\right) \\ \gamma &= \text{atan2}\left(\sqrt{r_{31}^2 + r_{32}^2}, r_{33}\right)\end{aligned}$$

where r_{ij} is the element of R at i th row and j th column, and

$$\text{atan2}(y, x) = \begin{cases} \text{atan}\left(\frac{y}{x}\right) & (x > 0) \\ \frac{y}{|y|} \cdot (\pi - \text{atan}\left(|\frac{y}{x}\right|)) & (x < 0) \\ \frac{y}{|y|} \cdot \frac{\pi}{2} & (y \neq 0, x = 0) \\ \text{undefined} & (y = 0, x = 0) \end{cases}$$

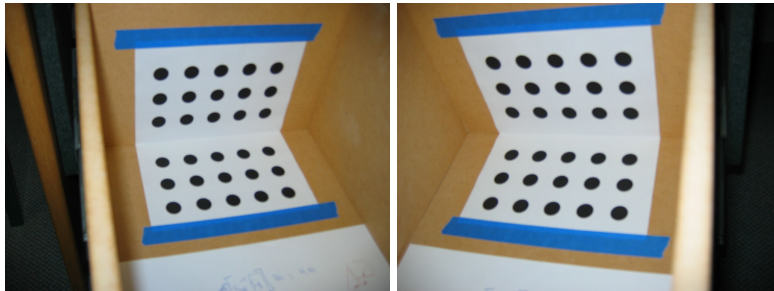


Fig. 3. The first (left) and second (right) candidate images.

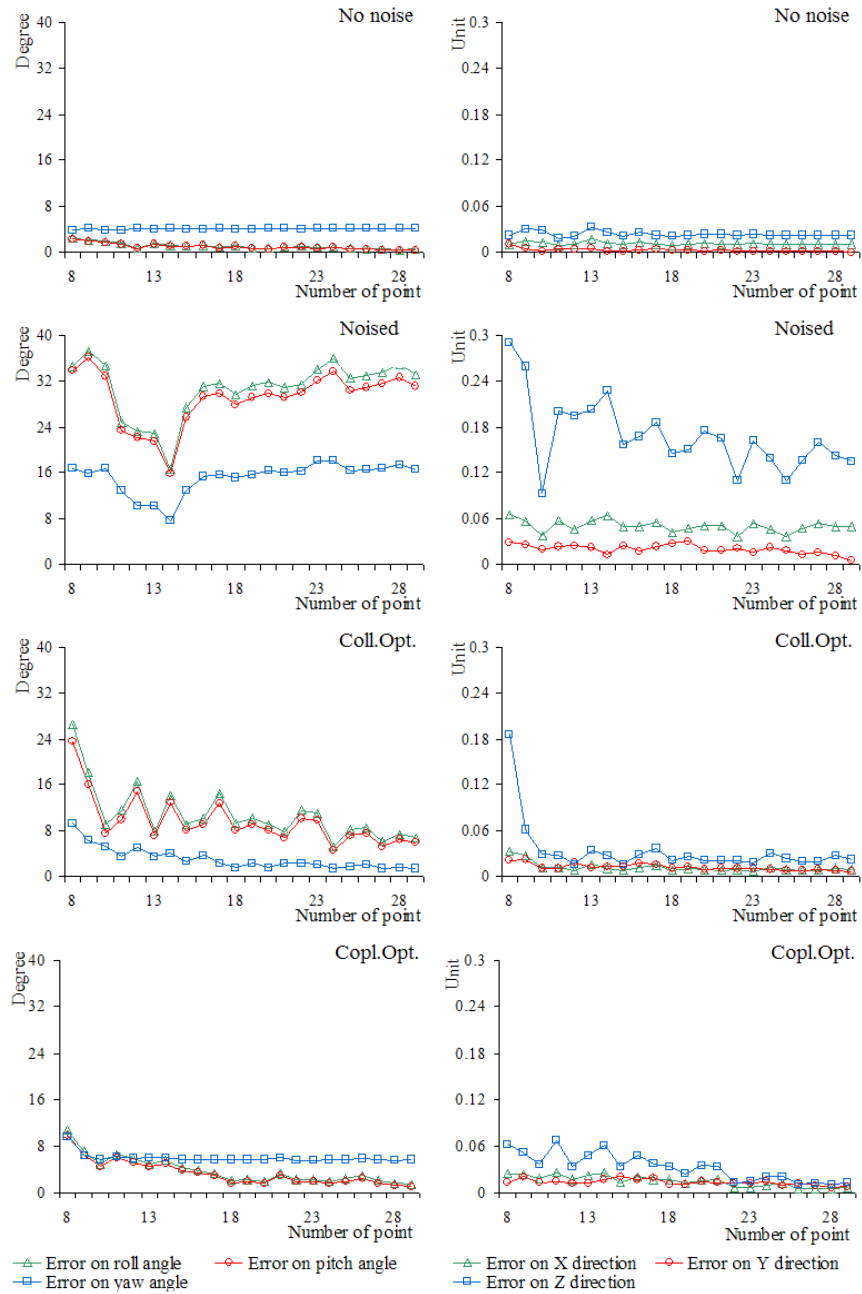


Fig. 4. Errors in rotation matrices (left) or translation vectors (right). First row: errors from non-noisy data. Second row: noisy data. Third or fourth row: errors from noisy data after optimization with collinearity or coplanarity knowledge, respectively.



Fig. 5. Epipolar lines result from different data sets. The green lines (dash lines) from data without generated noise; the red lines (straight lines) are from noisy data; the blue lines (dash-dot lines) and yellow lines (dot lines) are from noisy data which has been optimized with collinearity and coplanarity knowledge, respectively.

Since splitting the essential matrix only results in a translation vector up to a scale factor, all translation vectors (include the ground true one) are transformed into a normalized vector (length equal to one unit) to compare with each other in a quantitative manner. The comparison of rotation matrices and translation vectors are shown in Figure 4. The errors are mean error of ten times iteration when different number of correspondences are given. The noise propagated is Gaussian noise (with mean 0 and deviation $\delta = 1$ pixel). The method used to compute the fundamental matrix is the 8-Point Algorithm, which is more sensitive to noise than RANSAC and LMedS Algorithm. The method of Hartley and Zisserman is used to split essential matrix.

According to the results shown in Figure 4, the coplanarity knowledge gives a better optimization than collinearity knowledge. One possible reason is that the collinearity optimization is performed on uncalibrated images, in which the true correlation of collinear points are not strictly lying in a straight line.

For arbitrary images, the effect of optimizations can be seen from Figure 5 through looking at relative positions of epipolar lines computed from different data sets. It shows that the two optimization strategies bring positive effects on

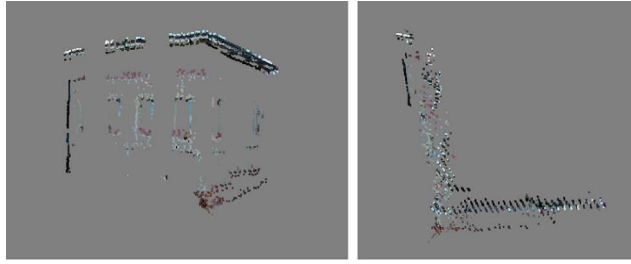


Fig. 6. Two different views of reconstructed points from the optimized SfM algorithm.



Fig. 7. Triangulated surface mesh with textures

reducing the influence of noise, and the coplanarity optimization performs better than the collinearity optimization. Figure 6 shows the reconstructed point cloud of the CITR-building in Auckland and Figure 7 visualizes the texture mapped surface model. The main edges of the building are reconstructed with near-perfect 90° angles. Slight image noise (less than 1 pixel) already leads to angles between 20° and 140° which indicates the sensitivity of classic SfM approaches. By incorporating the collinearity and coplanarity constraints, the reconstruction quality improved.

5 Summary

Modules relating to structure from motion have been discussed in this paper. According to experiments, structure from motion is sensitive to noise and it is necessary to improve its stability. Two optimizations, using collinearity and coplanarity knowledge, have been proposed, and the relating experiments show that the two proposed optimizations, especial the coplanarity one, bring positive effects on reducing influences of noise.

Acknowledgments

The authors thank Daniel Grest and Kevin Koeser from the University of Kiel for the BIAS-Library and useful hints.

References

1. M. Fischler and R. Bolles. Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography. *Comm. ACM*, **24**:381–385, 1981.
2. C. Harris and M. Stephen. A combined corner and edge detector. In Proc. *Alvey Vision Conf.*, pages 147–151, 1988.
3. R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, 2000.
4. B. K. P. Horn. Recovering baseline and orientation from essential matrix. <http://www.ai.mit.edu/people/bkph/papers/essential.pdf>, 1990.
5. T. Huang. Motion and structure from feature correspondences: a review. *Proc. IEEE*, **82**:252–268, 1994.
6. R. Klette, K. Schlüns, and A. Koschan. *Computer Vision – Three-dimensional Data from Images*. Springer, Singapore, 1998.
7. C. L. Lawson and R. J. Hanson. *Solving Least Squares Problems*. Prentice Hall, Englewood, 1974.
8. H. Longuet-Higgins. A computer algorithm for reconstructing a scene from two projections. *Nature*, **293**:133–135, 1981.
9. R. Murray, Z. Li, and S. Sastry. *A Mathematical Introduction to Robotic Manipulation*. CRC Press, Boca Raton, 1994.
10. M. Pollefeys, R. Koch, and L. van Gool. A simple and efficient rectification method for general motion. In Proc. *Int. Conf. Computer Vision*, pages 496–501, 1999.
11. R. Szeliski and P. H. S. Torr. Geometrically constrained structure from motion: points on planes. In Proc. *European Workshop 3D Structure Multiple Images Large-Scale Environments*, pages 171–186, 1998.
12. R. Tsai. A versatile camera calibration technique for high-accuracy 3d machine vision metrology using off-the-shelf tv cameras and lenses. *IEEE J. Robotics and Automation*, **3**:323–344, 1987.
13. Z. Zhang. Determining the epipolar geometry and its uncertainty: a review. *Int. J. Computer Vision*, **27**:161–198, 1998.
14. Z. Zhang, R. Deriche, O. Faugeras, and Q.-T. Luong. A robust technique for matching two uncalibrated images through the recovery of the unknown epipolar geometry. *Artificial Intelligence J.*, **78**:87–119, 1995.