

# A Unified Framework for Quasi-Linear Bundle Adjustment

Adrien Bartoli, *Adrien.Bartoli@inria.fr*

INRIA Rhône-Alpes, 655, av. de l'Europe, 38334 St. Ismier cedex, France.

## Abstract

*Obtaining 3D models from long image sequences is a major issue in computer vision. One of the main tools used to obtain accurate structure and motion estimates is bundle adjustment. Bundle adjustment is usually performed using non-linear Newton-type optimizers such as Levenberg-Marquardt which might be quite slow when handling a large number of points or views.*

*We investigate an algorithm for bundle adjustment based on quasi-linear optimization. The method is straightforward to implement and relies on solving weighted linear systems obtained as simple functions of the input data. Important features are that (i) the original cost function is preserved, (ii) different projection models, either calibrated or not, are handled in a unified framework and (iii) any number of views and points as well as missing data can be handled.*

*Experimental results on simulated and real data show that the algorithm is as accurate as standard techniques while requiring less computational time to converge.*

## 1. Introduction

Recovering metric structure and motion from uncalibrated images is a central problem for computer vision. Most of the time, a sub-optimal solution is obtained in projective space for motion, then for structure [1], or jointly [6]. Self-calibration is then performed to upgrade the reconstruction to metric space.

In this paper, we investigate a very simple technique for bundle adjustment. Experimental results reveal that it is as accurate as traditional techniques while requiring less CPU time. We consider both the full perspective and the affine projection, either calibrated or uncalibrated. We may end up therefore with projective, affine or metric reconstructions.

Let us consider point features. Bundle adjustment basically consists in finding the minimum of a cost function defined as the sum of discrepancies between actual and predicted image points. These discrepancies are measured using the Euclidean distance. The optimization is performed over a set of parameters that represents structure and motion using most of the time non-linear Newton-type optimization techniques such as Levenberg-Marquardt [7]. This technique

has the following drawbacks. Firstly, handling the free scale of homogeneous entities such as vectors or matrices representing points or cameras may be non-trivial. Secondly, it requires the computation of at least the first order Jacobian matrix of the residuals with respect to structure and motion parameters which might be non-trivial. Thirdly, the computational cost may be non-negligible since the Hessian matrix has to be inverted, even if specific techniques have been proposed to speed up the process [7].

In this paper, we adapt the resection-intersection technique [2, 5, 7] to perform bundle adjustment using quasi-linear optimizations. Our contribution differs from previous ones in that (i) we use the original cost function of bundle adjustment, which preserves optimality and (ii) we handle a great variety of camera models in a unified manner. In particular, we deal with calibrated configurations where non-linear constraints hold on the recovered motion. The final algorithm is simplicity itself. The method relies on rewriting the non-linear Euclidean distance used in bundle adjustment as a weighted bilinear algebraic distance [4]. In practice, it consists in iteratively solving weighted least squares systems that are simple functions of image point positions. Experimental results reveal that it performs as well as standard techniques in terms of convergence accuracy while greatly reducing computational cost.

This paper is organized as follows. In §2, we give our notation and preliminaries. §3 introduces the principle and the algorithm for quasi-linear bundle adjustment. These results are independent of the camera model. The approach is specialized in §§4.1 and 4.2 to uncalibrated and calibrated perspective cameras respectively. Finally, §5 validates the approach on simulated and real data and §6 gives our conclusions and perspectives.

## 2. Notation and preliminaries.

Vectors are typeset using bold letters ( $\mathbf{q}$ ,  $\mathbf{Q}$ ), matrices using sans-serif fonts ( $P$ ,  $S$ ) and scalars in italics ( $x$ ,  $w$ ). We make no formal distinction between coordinate vectors and physical entities. Everything is expressed in homogeneous coordinates, e.g.  $\mathbf{q}^T \sim (q_1 \ q_2 \ q_3)$  is an image point, where  $\sim$  means equality up to scale and  $^T$  is vector/matrix transpo-

sition.

The standard cost function for bundle adjustment is defined by  $\mathcal{C} = \sum_i \sum_j d^2(\mathbf{q}_{ij}, \hat{\mathbf{q}}_{ij})$ , where  $d^2$  is the squared Euclidean distance,  $\hat{\mathbf{q}}_{ij}$  the reprojection of the  $i$ -th estimated 3D point in the  $j$ -th image and  $\mathbf{q}_{ij}$  the corresponding image point. We drop the indices  $i$  and  $j$  for clarity. Bundle adjustment can be formulated as  $\min_{\mathbf{P}, \mathbf{Q}} \mathcal{C}$ , where  $\mathbf{P}$  and  $\mathbf{Q}$  respectively designate the estimates of the parameters of motion and structure.

**Distances in the image.** The Euclidean distance is employed for bundle adjustment and in many criteria since it is physically meaningful. Most of the time, its non-linearity induces that of the corresponding criterion which therefore can not be used to linearly estimate e.g. structure or motion. On the other hand, consider the squared *algebraic distance*  $d_a$  defined by [4]:

$$d_a^2(\mathbf{x}, \mathbf{y}) = w_{\mathbf{x}, \mathbf{y}}^2 d_a^2(\mathbf{x}, \mathbf{y}) \text{ where } w_{\mathbf{x}, \mathbf{y}} = (x_3 y_3)^{-2}. \quad (1)$$

The non-linearity is hidden in the weight factors  $w_{\mathbf{x}, \mathbf{y}}$  which can be thought of as a bias of the algebraic distance with respect to the Euclidean one. A convenient expression may be derived as follows:

$$d_a^2 = \|\mathbf{S}[\mathbf{x}]_{\times} \mathbf{y}\|^2 \text{ and } \mathbf{S} = \begin{pmatrix} -1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}, \quad (2)$$

where  $\|\cdot\|^2$  is the  $\mathcal{L}_2$  or squared Frobenius norm for vectors or matrices respectively and  $[\cdot]_{\times}$  the skew-symmetric  $3 \times 3$ -matrix associated with the cross product, i.e.  $[\mathbf{v}]_{\times} \mathbf{q} = \mathbf{v} \times \mathbf{q}$ . The algebraic distance does not have any physical meaning since it depends upon the relative scale of the homogeneous representation of  $\mathbf{x}$  and  $\mathbf{y}$ . Its advantage is that it is bilinear. Since in most estimation problems one of the points compared is known,  $d_a$  becomes linear for the other.

**Quasi-linear optimization.** The relationship (1) between the Euclidean and the algebraic distance is the basis for quasi-linear optimization. The principle is to initialize weight factors to unity and compute a biased estimate using the weighted algebraic distance. When such an estimate has been obtained, weight factors can be computed using (1). The process is then iterated until convergence. Typically, 3 or 4 iterations are enough.

### 3. Quasi-Linear Bundle Adjustment

We derive the bundle adjustment algorithm, which is valid for most camera models. We then detail camera-dependent elements.

Consider the initial bundle adjustment problem which involves optimizing simultaneously over the motion and structure parameters. It can be split, as indicated in [2, 5, 8] and becomes  $\min_{\mathbf{P}} (\min_{\mathbf{Q}} \mathcal{C})$ . This is strictly equivalent to the initial problem in that the same minima are kept. We split accordingly the optimization process into two distinct steps, for motion on the one hand and for structure on the other hand.

These optimization steps can be alternatively performed to refine structure and motion. This scheme is termed resection-intersection. Each step can then be conducted independently using non-linear optimizers such as Levenberg-Marquardt, as proposed in [8]. The advantages of doing so are numerous:

- it reduces the complexity, since the optimization consists of  $(m - 1)$  independent outer steps<sup>1</sup> over an 11-dimensional space and  $n$  independent 3-dimensional inner steps ( $m$  is the number of views of  $n$  is the number of points);
- it performs as well as directly optimizing over all the parameters in terms of convergence accuracy, as reported by [2, 5, 8] and confirmed by our results;
- each of the two alternated optimization steps can be performed quasi-linearly, see the next section.

The algorithm, given below, needs an initialization of either structure or motion whereas standard bundle adjusters require both:

1. *initialize weight factors* to unity:  $w_{\mathbf{q}, \hat{\mathbf{q}}} = 1$ . These are not necessary for the affine camera model;
2. *quasi-linearly solve for structure*;
3. *quasi-linearly solve for motion*;
4. *iterate* steps 2 and 3 until convergence (see below);

The estimation of weight factors, which corrects the bias of the algebraic distance with respect to the Euclidean one, is done inside steps 2 and 3.

**Convergence.** Convergence is achieved when steps 2 and 3 become idempotent for structure and motion, i.e. when the difference between two consecutive values of the cost function  $\mathcal{C}$  is lower than a threshold, typically  $10e^{-8}$  pixel, or when the residual error begins to increase.

## 4. Resection-Intersection

### 4.1. Uncalibrated Pin-Hole Cameras

We consider uncalibrated full perspective projection to model cameras. This yields projective reconstructions. We derive the equations that lead to structure, then motion, quasi-linear estimation. Consider one term  $d^2(\mathbf{q}, \hat{\mathbf{q}})$  of the double sum of the cost function  $\mathcal{C}$ . Using the link (1) between the Euclidean and the algebraic distances  $d$  and  $d_a$ , then the expression (2) of the algebraic distance, each of these terms can be written  $d^2(\mathbf{q}, \hat{\mathbf{q}}) = w_{\mathbf{q}, \hat{\mathbf{q}}}^2 \|\mathbf{S}[\mathbf{q}]_{\times} \hat{\mathbf{q}}\|^2$ . Using perspective projection  $\hat{\mathbf{q}} \sim \mathbf{P}\mathbf{Q}$  and some minor algebraic manipulations:

$$d^2(\mathbf{q}, \hat{\mathbf{q}}) = w_{\mathbf{q}, \hat{\mathbf{q}}}^2 \left\| \begin{pmatrix} \mathbf{p}_2^T \mathbf{Q} - q_2 \mathbf{p}_3^T \mathbf{Q} \\ \mathbf{p}_1^T \mathbf{Q} - q_1 \mathbf{p}_3^T \mathbf{Q} \end{pmatrix} \right\|^2, \quad (3)$$

where  $\mathbf{p}_i^T$  are the rows of  $\mathbf{P}$ . Using (1), we obtain the expression for the weight factors as:

$$w_{\mathbf{q}, \hat{\mathbf{q}}} = (\mathbf{p}_3^T \mathbf{Q})^{-2}. \quad (4)$$

<sup>1</sup>we do not optimize the first camera parameters to (partially) fix the coordinate frame.

**Intersection — Solving for the structure.** Using equation (3), and if weight factors  $w_{\mathbf{q}, \hat{\mathbf{q}}}$  are assumed to be known, one can recover the structure by solving the linear system  $\mathbf{M}^{\text{us}} \mathbf{Q} = \mathbf{0}$  for each point feature  $\mathbf{Q}$  considered where:

$$(\mathbf{M}^{\text{us}})_{2m \times 4} = \left( w_{\mathbf{q}, \hat{\mathbf{q}}} \begin{pmatrix} \dots & \dots & \dots & \dots \\ \mathbf{p}_2^\top - \mathbf{q}_2 \mathbf{p}_3^\top & & & \\ \mathbf{p}_1^\top - \mathbf{q}_1 \mathbf{p}_3^\top & & & \\ \dots & \dots & \dots & \dots \end{pmatrix} \right). \quad (5)$$

Each of the  $m$  views where the point  $\mathbf{Q}$  is visible provides 2 equations. The solution is the unit singular vector corresponding to the smallest singular value of  $\mathbf{M}^{\text{us}}$  which can be obtained using e.g. SVD. This naturally ensures the free scale of the recovered homogeneous vector since  $\|\mathbf{Q}\|^2 = 1$ .

The process is then iterated while re-estimating weight factors  $w_{\mathbf{q}, \hat{\mathbf{q}}}$  to reduce the bias of the algebraic distance until convergence, which is assessed as in §3.

**Resection — Solving for the motion.** We use equation (3) to form a linear system  $\mathbf{M}^{\text{um}} \mathbf{p}$  that can be solved using SVD where  $\mathbf{p} \sim (\mathbf{p}_1^\top \ \mathbf{p}_2^\top \ \mathbf{p}_3^\top)$  encapsulates the rows of  $\mathbf{P}$  and:

$$(\mathbf{M}^{\text{um}})_{2n \times 12} = \left( w_{\mathbf{q}, \hat{\mathbf{q}}} \begin{pmatrix} \dots & \dots & \dots & \dots & \dots & \dots \\ \mathbf{o}_4^\top & \mathbf{q}^\top & -\mathbf{q}_2 \mathbf{q}^\top & & & \\ \mathbf{q}^\top & \mathbf{o}_4^\top & -\mathbf{q}_1 \mathbf{q}^\top & & & \\ \dots & \dots & \dots & \dots & \dots & \dots \end{pmatrix} \right).$$

Each of the  $n$  points that project onto view  $\mathbf{P}$  provides 2 equations. The iteration is then conducted while re-estimating the weight factors, as for the structure. The free scale of  $\mathbf{P}$  is fixed since  $\|\mathbf{P}\|^2 = \|\mathbf{p}\|^2 = 1$ .

## 4.2. Calibrated Pin-Hole Cameras

We consider the calibrated full perspective camera model. This leads to metric reconstructions. In this case, resection is termed pose estimation and has been widely studied, see e.g. [3]. When the motion is frozen, solving for the structure can be conducted using the quasi-linear algorithm of §4.1.

Pose estimation from minimal data, i.e.  $n = 3$  points is non-linear by nature since there exist four solutions in this case. However, when an initial guess is known, it can be refined using linear algebra. Difficulties arise because of non-linear constraints on the orthonormality of the recovered rotation matrix. An algorithm for pose refinement with  $n \geq 3$  points that fits into our quasi-linear framework is described in the next section.

**Quasi-linear pose estimation.** Let  $\mathbf{P} \sim \mathbf{K}(\mathbf{R} \ \mathbf{t})$  be the projection matrix where  $\mathbf{K}$  is a  $3 \times 3$  matrix containing the known intrinsic parameters,  $\mathbf{R}$  a  $3 \times 3$  rotation matrix and  $\mathbf{t}$  a 3-vector. The unknowns are therefore  $\mathbf{R}$  and  $\mathbf{t}$ . Under the assumption of zero skew and unity aspect ratio which are valid for most modern cameras, distances measured onto the image plane are equal to distances measured onto the retina up to a scale factor as  $d^2(\mathbf{q}, \hat{\mathbf{q}}) = f^2 d^2(\mathbf{x}, \hat{\mathbf{x}})$  where  $\mathbf{x} = \mathbf{K}^{-1} \mathbf{q}$  and  $\hat{\mathbf{x}}^{-1} = \mathbf{K}^{-1} \hat{\mathbf{q}}$  are retinal points and  $f$  is the focal length (in number of pixels).

**Linearly update the motion.** Each term  $d^2(\mathbf{q}, \hat{\mathbf{q}})$  of the double sum of the cost function  $\mathcal{C}$  can be expanded as in the uncalibrated case, but using retinal points. This gives  $d^2(\mathbf{q}, \hat{\mathbf{q}}) = f^2 w_{\mathbf{x}, \hat{\mathbf{x}}}^2 \|\mathbf{S}[\mathbf{x}]_{\times} \hat{\mathbf{x}}\|^2$ . Using then perspective projection  $\hat{\mathbf{x}} \sim (\mathbf{R} \ \mathbf{t}) \mathbf{Q}$  of the point  $\mathbf{Q}^T \sim (\mathbf{Q}^T \ \alpha)$  on the retina, we obtain  $d^2(\mathbf{q}, \hat{\mathbf{q}}) = f^2 w_{\mathbf{x}, \hat{\mathbf{x}}}^2 \|\mathbf{S}[\mathbf{x}]_{\times} (\mathbf{R} \mathbf{Q} + \alpha \mathbf{t})\|^2$ . If we approximate the rotation between two iterations of the algorithm by a first order Taylor expansion,  $(\mathbf{I} + [\mathbf{w}]_{\times})$ , we obtain after some minor algebraic manipulations  $d^2(\mathbf{q}, \hat{\mathbf{q}}) = f^2 w_{\mathbf{x}, \hat{\mathbf{x}}}^2 \|\mathbf{S}[\mathbf{x}]_{\times} \mathbf{R} \mathbf{Q} - \mathbf{S}[\mathbf{x}]_{\times} [\hat{\mathbf{Q}}]_{\times} \mathbf{w} + \alpha \mathbf{S}[\mathbf{x}]_{\times} \mathbf{t}\|^2 + \mathcal{O}$ , where  $\mathcal{O}$  contains higher order terms. By neglecting  $\mathcal{O}$ , we can estimate the first order pose update  $\mathbf{m}^T = (\mathbf{w}^T \ \mathbf{t}^T)$  by solving the linear system  $\mathbf{M}^{\text{cm}} \mathbf{m} = \mathbf{d}^{\text{cm}}$ , where:

$$(\mathbf{M}^{\text{cm}})_{2n \times 6} = \begin{pmatrix} f^2 w_{\mathbf{x}, \hat{\mathbf{x}}}^2 ( \ \mathbf{S}[\mathbf{x}]_{\times} \mathbf{R} [\hat{\mathbf{Q}}]_{\times} & -\alpha \mathbf{S}[\mathbf{x}]_{\times} ) \\ \dots & \dots \end{pmatrix}$$

$$(\mathbf{d}^{\text{cm}})_{2n} = \begin{pmatrix} f^2 w_{\mathbf{x}, \hat{\mathbf{x}}}^2 ( \ \mathbf{S}[\mathbf{x}]_{\times} \mathbf{R} \mathbf{Q} ) \\ \dots \end{pmatrix}.$$

**Algorithm.** The algorithm for quasi-linear pose estimation combines first order updates of pose and re-estimation of calibrated weight factors  $w_{\mathbf{x}, \hat{\mathbf{x}}} = \frac{1}{f} w_{\mathbf{q}, \hat{\mathbf{q}}}$ :

1. *solve for pose update* using  $\mathbf{M}^{\text{cm}} \mathbf{m} = \mathbf{d}^{\text{cm}}$ ;
2. *update rotation:*  $\mathbf{R} \leftarrow \mathbf{R} \mathbf{R}_{\mathbf{w}}$ , where  $\mathbf{R}_{\mathbf{w}}$  gives the rotation corresponding to  $\mathbf{w}^T = (w_1 \ w_2 \ w_3)$  using e.g. Euler angles;
3. *re-estimate calibrated weight factors*  $w_{\mathbf{x}, \hat{\mathbf{x}}}$  using (4);
4. *iterate* steps 1-3 until convergence, see §3.

## 5. Experimental Results

In the following two sections, we compare previously described algorithms to existing ones. We use simulated data to compare them in the most general case, i.e. projective. We then consider real images and both projective and metric spaces. More detailed results, in particular on pose estimation and for the affine camera will be made available in a longer version of the paper. We compare the following algorithms. LIN is structure and motion initialization [1], QLIN is the proposed quasi-linear method and LM relies on the Levenberg-Marquardt non-linear optimizer. Note that our implementation uses analytical differentiation and a sparse inversed for the Hessian matrix, see [7]. We also compared a full Hessian inversion-based implementation that we denote as LM-full.

The threshold for convergence (see §3) is the same for QLIN and LM and is equal to  $10e^{-8}$ . Image point coordinates are standardized such that they lie in  $[-1 \dots 1]^2$ .

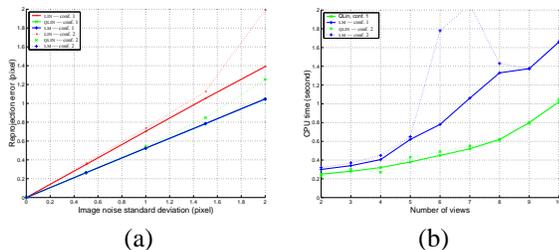
### 5.1. Simulated Data

The simulated scene consists of 3D point features chosen at random inside a cube of 2 meters side length. We simulate two camera configurations. The first one is stable and the second one contains nearly unreconstructable points. We measure the reprojection error to assess the convergence accuracy

and the CPU time to convergence to compare computational costs. In order to simulate realistic situations, we adopted the following parameters: the focal length of the cameras is 1000 (expressed as number of pixels) and the image size is  $512 \times 512$ . Values plotted have been averaged over 50 trials. Median values gave graphs similar to those for the means. The curves for 3D errors show expected results, see e.g. [6].

**Convergence accuracy.** The standard configuration is 50 points, 10 views and an added image noise with a Gaussian distribution of 0.5 pixel standard deviation.

Figures 1 (a) shows the results when varying the image noise. We observe that for camera configuration 1, there is no difference between QLIN and LM, even when the starting point given by LIN lies quite far from the final result. The error for these two methods degrades gracefully as image noise increases. For camera configuration 2, figure 1 (a) shows that there is a slight difference between QLIN and LM only when the level of added image noise is beyond 1 pixel. LM-full, not shown on these graphs, gave exactly the same results as LM.



**Figure 1. (a): reprojection error at convergence versus added image noise. The curves for QLIN and LM and the first configuration, and for LM and the second configuration are nearly undistinguishable. (b): CPU time to convergence versus the number of views.**

**Computation time.** We observe on figure 1 (b) that the quasi-linear optimization requires less CPU time to convergence than the non-linear one whatever the number of views. Results from another experiment where added image noise is varied from 0 to 2 pixels, shows that CPU time to convergence linearly increases with the added image noise. The CPU time needed by LM is twice thus needed by QLIN. LM-full gave results out of range of these graphs.

## 5.2. Real Images

We use the 181 images of the hotel sequence<sup>2</sup> and 197 corner correspondences. These corners have been automatically detected and semi-automatically matched through the different views. Calibration data were approximately known for this sequence. The results given in table 1 show that QLIN

<sup>2</sup>These data have been provided by the Modeling by Videotaping group in the Robotics Institute, Carnegie Mellon University.

step	approach	first 5 images		full 181 images sequence	
		rep. error (pixel)	CPU time (second)	rep. error (pixel)	CPU time (second)
init.	-	9.4209	-	6.2524	-
ba	LM-full	0.0357	1056 $\approx$ 17 mn.	-	-
	LM	0.0357	2.55	2.2026	96
	QLIN	0.0391	0.95	2.2034	37

**Table 1. Reprojection errors (pixel) and CPU time to convergence (second) for various algorithms and sequence lengths.**

converges to the same residual error as LM and LM-full. In terms of computational cost, we observe that QLIN is roughly 2.5 times faster than LM.

## 6. Conclusions and Perspectives

We addressed the problem of structure and motion refinement using bundle adjustment. We perform this task based on quasi-linear optimization while keeping the original cost function of bundle adjustment. We split the problem into structure refinement on the one hand and motion refinement on the other hand. The principle is then to rewrite the Euclidean distance used in bundle adjustment as a weighted algebraic distance. Quasi-linear optimization iteratively updates these weights while refining structure or motion.

The result is a disarmingly simple algorithm that comes out to be a loop over two weighted linear systems constructed as simple functions as the input data.

We conducted numerous experiments on simulated and real data. In the light of these results, we observed that this algorithm is as accurate as standard Levenberg-Marquardt-based bundle adjusters in terms of convergence accuracy while being much faster in terms of computational cost.

**Acknowledgement.** I wish to thank Frederik Schaffalitzky for providing his implementation of LM used in §5.

## References

- [1] P. Beardsley, P. Torr, and A. Zisserman. 3D model acquisition from extended image sequences. In *ECCV*, 1996.
- [2] Q. Chen and G. Medioni. Efficient iterative solution to M-view projective reconstruction. In *CVPR*, 1999.
- [3] R. Haralick, H. Joo, C. Lee, X. Zhuang, V. Vaidya, and M. Kim. Pose estimation from corresponding point data. *IEEE SMC*, 6(19):1426–1446, 1989.
- [4] R. Hartley. Minimizing algebraic error. In *ICCV*, 1998.
- [5] S. Mahamud, M. Herbert, Y. Omori, and J. Ponce. Provably-convergent iterative methods for projective structure and motion. In *CVPR*, 2001.
- [6] P. Sturm and B. Triggs. A factorization based algorithm for multi-image projective structure and motion. In *ECCV*, 1996.
- [7] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon. Bundle adjustment — a modern synthesis. In *Vision Algorithms: Theory and Practice*, 2000.
- [8] Z. Zhang. Motion and structure from two perspective views: From essential parameters to euclidean motion via fundamental matrix. *Journal of the Optical Society of America A*, 1997.