

6

Camera Calibration

For the ancient Egyptians, exactitude was symbolized by a feather that served as a weight on scales used for the weighing of souls.
Italo Calvino, Six Memos for the Next Millennium

This chapter tackles the problem of camera calibration; that is, determining the value of the extrinsic and intrinsic parameters of the camera.

Chapter Overview

Section 6.1 defines and motivates the problem of camera calibration and its main issues.

Section 6.2 discusses a method based on simple geometric properties for estimating the camera parameters, given a number of correspondences between scene and image points.

Section 6.3 describes an alternative, simpler method, which recovers the projection matrix first, then computes the camera parameters as functions of the entries of the matrix.

What You Need to Know to Understand this Chapter

- Working knowledge of geometric camera models (Chapter 2).
- SVD and constrained least-squares (Appendix, section A.6).
- Familiarity with line extraction methods (Chapter 5).

6.1 Introduction

We learned in Chapters 4 and 5 how to identify and locate image features, and we are therefore fully equipped to deal with the important problem of *camera calibration*; that

be estimating the values of the intrinsic and extrinsic parameters of the camera model, which was introduced in Chapter 2.

The key idea behind calibration is to write the projection equations linking the known coordinates of a set of 3-D points and their projections, and solve for the camera parameters. In order to get to know the coordinates of some 3-D points, camera calibration methods rely on one or more images of a calibration pattern: that is, a 3-D object of known geometry, possibly located in a known position in space and generating image locations which can be located accurately. Figure 6.1 shows a typical calibration pattern, consisting of two planar grids of black squares on a white background. It is easy to know the 3-D position of the vertices of each square once the position of the two planes has been measured, and locate the vertices on the images, for instance as intersection of image lines, thanks to the high contrast and simple geometry of the pattern.

Problem Statement

Given one or more images of a calibration pattern, estimate

1. the intrinsic parameters,
2. the extrinsic parameters, or
3. both.

** The accuracy of calibration depends on the accuracy of the measurements of the calibration pattern; that is, its construction tolerances. To be on the safe side, the calibration pattern should be built with tolerances one or two orders of magnitude smaller than the desired accuracy of calibration. For example, if the desired accuracy of calibration is 0.1mm, the calibration pattern should be built with tolerances smaller than 0.01mm.

Although there are techniques inferring 3-D information about the scene from uncalibrated cameras, some of which will be described in the next two chapters, effective camera calibration procedures open up the possibility of using a wide range of existing algorithms for 3-D reconstruction and recognition, all relying on the knowledge of the camera parameters.

This chapter discusses two algorithms for camera calibration. The first method recovers directly the intrinsic and extrinsic camera parameters; the second method estimates the projection matrix first, without solving explicitly for the various parameters, which are then computed as closed-form functions of the entries of the projection matrix. The choice of which method to adopt depends largely on which algorithms are to be applied next (Section 6.4).

¹ Recall that the projection matrix links world and image coordinates, and its entries are functions of the intrinsic and extrinsic parameters.

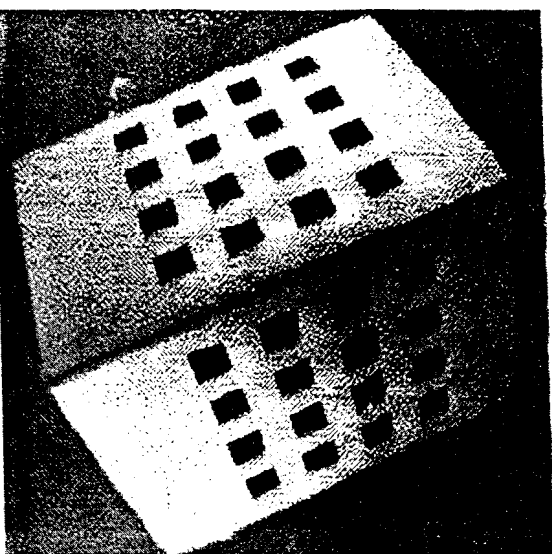


Figure 6.1 The typical calibration pattern used in this chapter.

6.2 Direct Parameter Calibration

We start by identifying the parameters to be estimated, and cast the problem in geometric terms.

6.2.1 Basic Equations

Consider a 3-D point, P , defined by its coordinates $[X^w, Y^w, Z^w]^T$ in the world reference frame. As usual in calibration, the world reference frame is known.

** This means to pick an accessible object defining three mutually orthogonal directions intersecting in a common point. In this chapter, this object is the calibration pattern: in vision systems for indoor robot navigation, for instance, it can be a corner of a room.

Let $[X^c, Y^c, Z^c]^T$ be the coordinates of P in the camera reference frame (with $Z^c > 0$ if P is visible). As usual, the origin of the camera frame is the center of projection, and the Z axis is the optical axis. The position and orientation of the camera frame is unknown, since, unlike the image and world reference frames, the camera frame is inaccessible directly. This is equivalent to saying that we do not know the extrinsic parameters; that is, the 3×3 rotation matrix R and 3-D translation vector T such that

$$\begin{bmatrix} X^c \\ Y^c \\ Z^c \end{bmatrix} = R \begin{bmatrix} X^w \\ Y^w \\ Z^w \end{bmatrix} + T. \tag{6.1}$$

In components, (6.1) can be written as

$$\begin{aligned} X^c &= r_{11}X^w + r_{12}Y^w + r_{13}Z^w + T_x \\ Y^c &= r_{21}X^w + r_{22}Y^w + r_{23}Z^w + T_y \\ Z^c &= r_{31}X^w + r_{32}Y^w + r_{33}Z^w + T_z. \end{aligned} \tag{6.2}$$

Note the slight but important change of notation with respect to Chapter 2. In that chapter, the transformation between the world and camera reference frames was defined by translation *followed* by rotation. Here, the order is reversed and rotation *precedes* translation. While the rotation matrix is the same in both cases, the translation vectors differ (see Question 6.2).

Assuming that radial distortions (Section 2.4.3) can be neglected² we can write the image of $[X^c, Y^c, Z^c]^T$ in the image reference frame as (see (2.14) and (2.20))

$$x_{im} = -\frac{f}{Z^c} X^c + o_x \tag{6.3}$$

$$y_{im} = -\frac{f}{Z^c} Y^c + o_y \tag{6.4}$$

For simplicity, and since there is no risk of confusion, we drop the subscript $_{im}$ indicating image (pixel) coordinates, and write (x, y) for (x_{im}, y_{im}) . As we know from Chapter 2, (6.3) and (6.4) depend on the five intrinsic parameters f (focal length), s_x and s_y (horizontal and vertical effective pixel size), and o_x and o_y (coordinates of the image center), and, owing to the particular form of (6.3) and (6.4), the five parameters are not independent. However, if we let $f_x = f/s_x$ and $\alpha = s_y/s_x$, we may consider a new set of four intrinsic parameters, o_x, o_y, f_x , and α , all *independent* of one another. The parameter f_x is simply the focal length expressed in the effective horizontal pixel size (the *focal length in horizontal pixels*), while α , usually called *aspect ratio*, specifies the pixel deformation induced by the acquisition process defined in Chapter 2. Let us now summarize all the parameters to be calibrated in a box (see also the discussion in sections 2.4.2 and 2.4.3).

Extrinsic Parameters

- R , the 3×3 rotation matrix
- T , the 3-D translation vector

²We shall reconsider this assumption in Exercise 6.1, which suggests how to calibrate the parameters of the radial distortion model of Chapter 2.

Intrinsic Parameters

- $f_x = f/s_x$, length in effective horizontal pixel size units
- $\alpha = s_y/s_x$, aspect ratio
- (o_x, o_y) , image center coordinates
- k_1 , radial distortion coefficient

Plugging (6.2) into (6.3) and (6.4) gives

$$x - o_x = -f_x \frac{r_{11}X^w + r_{12}Y^w + r_{13}Z^w + T_x}{r_{31}X^w + r_{32}Y^w + r_{33}Z^w + T_z} \tag{6.5}$$

$$y - o_y = -\alpha f_x \frac{r_{21}X^w + r_{22}Y^w + r_{23}Z^w + T_y}{r_{31}X^w + r_{32}Y^w + r_{33}Z^w + T_z} \tag{6.6}$$

Notice that (6.5) and (6.6) bypass the inaccessible camera reference frame and link *directly* the world coordinates $[X^w, Y^w, Z^w]^T$ with the coordinates (x, y) of the corresponding image point. If we use a known calibration pattern, both vectors are measurable. This suggests that, given a sufficient number of points on the calibration pattern, we can try to solve (6.5) and (6.6) for the unknown parameters. This is the idea behind the first calibration method, which is articulated in two parts:

1. assuming the coordinates of the image center are known, estimate all the remaining parameters
2. find the coordinates of the image center

6.2.2 Focal Length, Aspect Ratio, and Extrinsic Parameters

We assume that the coordinates of the image center are known. Thus, with no loss of generality, we can consider the translated coordinates $(x, y) = (x - o_x, y - o_y)$. In other words, we assume that the image center is the origin of the image reference frame. As we said in the introduction, the key idea is to exploit the known coordinates of a sufficient number of corresponding image and world points.

Assumptions and Problem Statement

Assuming that the location of the image center (o_x, o_y) is known, and that radial distortion can be neglected, estimate f_x, α, R , and T from image points $(x_i, y_i), i = 1 \dots N$, projections of N known world points $[X_i^w, Y_i^w, Z_i^w]^T$ in the world reference frame.

The key observation is that (6.5) and (6.6) have the same denominator, therefore, from each corresponding pair of points $((X_i^w, Y_i^w, Z_i^w), (x_i, y_i))$ we can write an equation of the form

$$x_i f_x (r_{21}X_i^w + r_{22}Y_i^w + r_{23}Z_i^w + T_y) = y_i f_x (r_{11}X_i^w + r_{12}Y_i^w + r_{13}Z_i^w + T_x). \tag{6.7}$$

Since $\alpha = f_x/f_y$, (6.7) can be thought of as a linear equation for the 8 unknowns $\mathbf{v} = (v_1, v_2, \dots, v_8)$:

$$x_i X_i^w v_1 + y_i Y_i^w v_2 + x_i Z_i^w v_3 + x_i v_4 - y_i X_i^w v_5 - y_i Y_i^w v_6 - y_i Z_i^w v_7 - y_i v_8 = 0$$

where

$$\begin{aligned} v_1 &= r_{21} & v_5 &= \alpha r_{11} \\ v_2 &= r_{22} & v_6 &= \alpha r_{12} \\ v_3 &= r_{23} & v_7 &= \alpha r_{13} \\ v_4 &= T_x & v_8 &= \alpha T_x. \end{aligned}$$

Writing the last equation for the N corresponding pairs leads to the homogeneous system of N linear equations

$$A\mathbf{v} = 0 \tag{6.8}$$

where the $N \times 8$ matrix A is given by

$$A = \begin{bmatrix} x_1 X_1^w & x_1 Y_1^w & x_1 Z_1^w & x_1 & -y_1 X_1^w & -y_1 Y_1^w & -y_1 Z_1^w & -y_1 \\ x_2 X_2^w & x_2 Y_2^w & x_2 Z_2^w & x_2 & -y_2 X_2^w & -y_2 Y_2^w & -y_2 Z_2^w & -y_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ x_N X_N^w & x_N Y_N^w & x_N Z_N^w & x_N & -y_N X_N^w & -y_N Y_N^w & -y_N Z_N^w & -y_N \end{bmatrix}$$

If $N \geq 7$ and the N points are not coplanar, A has rank 7, and system (6.8) has a nontrivial solution (unique up to an unknown scale factor), which can be determined from the SVD of A , $A = UDV^T$, as the column of V corresponding to the only null singular value along the diagonal of D (Appendix, section A.6).

* The effects of the noise and the inaccurate localization of image and world points make the rank of A likely to be maximum (eight). In this case, the solution is the eigenvector corresponding to the *smallest* eigenvalue.

A rigorous proof of the fact that, in the ideal case (noise-free, perfectly known coordinates) the rank of A is 7 seems too involved to be presented here, and we refer you to the Further Readings section for details. Here, we just observe that, if the effective rank is larger than 7, system (6.8) would only have the trivial solution.

Our next task is to determine the unknown scale factor (and hence the various camera parameters) from the solution vector $\mathbf{v} = \bar{\mathbf{v}}$. If we call γ the scale factor, we have

$$\bar{\mathbf{v}} = \gamma(r_{21}, r_{22}, r_{23}, T_x, \alpha r_{11}, \alpha r_{12}, \alpha r_{13}, \alpha T_x). \tag{6.9}$$

Since $r_{21}^2 + r_{22}^2 + r_{23}^2 = 1$, from the first three components of $\bar{\mathbf{v}}$ we obtain

$$\sqrt{v_1^2 + v_2^2 + v_3^2} = \sqrt{\gamma^2(r_{21}^2 + r_{22}^2 + r_{23}^2)} = |\gamma|. \tag{6.10}$$

Similarly, since $r_{11}^2 + r_{12}^2 + r_{13}^2 = 1$ and $\alpha > 0$, from the fifth, sixth, and seventh component of $\bar{\mathbf{v}}$ we have

$$\sqrt{v_5^2 + v_6^2 + v_7^2} = \sqrt{\gamma^2 \alpha^2 (r_{11}^2 + r_{12}^2 + r_{13}^2)} = \alpha |\gamma|. \tag{6.11}$$

We can solve (6.10) and (6.11) for $|\gamma|$ as well as the aspect ratio α . We observe that the first two rows of the rotation matrix, \hat{R} , and the first two components of the translation vector, \mathbf{T} , can now be determined, up to an unknown common sign, from (6.9). Furthermore, the third row of the matrix \hat{R} can be obtained as the vector product of the first two estimated rows thought of as 3-D vectors. Interestingly, this implies that the sign of the third row is already fixed, as the entries of the third row remain unchanged if the signs of all the entries of the first two rows are reversed.

* Since the computation of the estimated rotation matrix, \hat{R} , does not take into account explicitly the orthogonality constraints, \hat{R} cannot be expected to be orthogonal ($\hat{R}\hat{R}^T = I$). In order to enforce orthogonality on \hat{R} , one can resort to the ubiquitous SVD decomposition. Assume the SVD of \hat{R} is $\hat{R} = UDV^T$. Since the three singular values of a 3×3 orthogonal matrix are all 1, we can simply replace D with the 3×3 identity matrix, I , so that the resulting matrix, UV^T , is exactly orthogonal (see Appendix, section A.6 for details).

Finally, we determine the unknown sign of the scale factor γ , and finalize the estimates of the parameters. To this purpose, we go back to (6.5), for example, with x instead of $x - o_x$, and recall that for every point $Z^w > 0$ and, therefore, x and $r_{11}X^w + r_{12}Y^w + r_{13}Z^w + T_x$ must have opposite sign. Consequently, it is sufficient to check the sign of $x(r_{11}X^w + r_{12}Y^w + r_{13}Z^w + T_x)$ for one of the points. If

$$x(r_{11}X^w + r_{12}Y^w + r_{13}Z^w + T_x) > 0, \tag{6.12}$$

the signs of the first two rows of \hat{R} and of the first two components of the estimated translation vector must be reversed. Otherwise, no further action is required. A similar argument can be applied to y and $r_{21}X^w + r_{22}Y^w + r_{23}Z^w + T_y$ in (6.6).

At this point, we have determined the rotation matrix, \hat{R} , the first two components of the translation vector, \mathbf{T} , and the aspect ratio, α . We are left with two, still undetermined parameters: T_z , the third component of the translation vector, and f_x , the focal length in horizontal pixel units. Both T_z and f_x can be obtained by least squares from a system of equations like (6.5) or (6.6), written for N points. To do this, for each point (x_i, y_i) we can write

$$x_i(r_{31}X_i^w + r_{32}Y_i^w + r_{33}Z_i^w + T_z) = -f_x(r_{11}X_i^w + r_{12}Y_i^w + r_{13}Z_i^w + T_x), \tag{6.13}$$

then solve the overconstrained system of N linear equations

$$A \begin{pmatrix} T_z \\ f_x \end{pmatrix} = \mathbf{b} \tag{6.14}$$

in the two unknowns T_z and f_x , where

$$A = \begin{bmatrix} x_1 & (r_{11}X_1^w + r_{12}Y_1^w + r_{13}Z_1^w + T_x) \\ x_2 & (r_{11}X_1^w + r_{12}Y_1^w + r_{13}Z_1^w + T_x) \\ \vdots & \vdots \\ x_N & (r_{11}X_N^w + r_{12}Y_N^w + r_{13}Z_N^w + T_x) \end{bmatrix}$$

and

$$b = \begin{pmatrix} -x_1(r_{31}X_1^w + r_{32}Y_1^w + r_{33}Z_1^w) \\ \vdots \\ -x_N(r_{31}X_N^w + r_{32}Y_N^w + r_{33}Z_N^w) \end{pmatrix}$$

The least squares solution (\hat{T}_z, \hat{f}_x) of system (6.14) is

$$\begin{pmatrix} \hat{T}_z \\ \hat{f}_x \end{pmatrix} = (A^T A)^{-1} A^T b. \quad (6.15)$$

It remains to be discussed how can we actually acquire an image of N points of known *world* coordinates, and locate the N corresponding image points accurately. One possible solution of this problem can be obtained with the pattern shown in Figure 6.1. The pattern consists of two orthogonal grids of equally spaced black squares drawn on white, perpendicular planes. We let the world reference frame be the 3-D reference frame centered in the lower left corner of the left grid and with the axes parallel to the three directions identified by the calibration pattern. If the horizontal and vertical size of the surfaces and the angle between the surfaces are known with high accuracy (from construction), then the 3-D coordinates of the vertices of each of the square in the *world* reference frame can be easily and accurately determined through simple trigonometry. Finally, the location of the vertices on the image plane can be found by intersecting the edge lines of the corresponding square sides. We now summarize the method:

Algorithm EXPL_PARS_CAL

The input is an image of the calibration pattern described in the text (Figure 6.1) and the location of the image center.

1. Measure the 3-D coordinates of each vertex of the n squares on the calibration pattern in the world reference frame. Let $N = 4n$.
2. In order to find the coordinates in the *image* reference frame of each of the N vertices:
 - locate the image lines defined by the sides of the squares (e.g., using procedures EDGE_COMP and HOUGH_LINES of Chapters 3 and 4).
 - estimate the image coordinates of all the vertices of the imaged squares by intersecting the lines found.

3. Having established the N correspondences between image and world points, compute the SVD of A in (6.8). The solution is the column of V corresponding to the smallest singular value of A .
4. Determine $|v|$ and α from (6.10) and (6.11).
5. Recover the first two rows of R and the first two components of T from (6.9).
6. Compute the third row of R as the vector product of the first two rows estimated in the previous step, and enforce the orthogonality constraint on the estimate of R through SVD decomposition.
7. Pick a point for which $(\alpha - \alpha_1)$ is noticeably different from 0. If inequality (6.12) is satisfied, reverse the sign of the first two rows of R and of the first two components of T .
8. Set up A and b of system (6.14), and use (6.15) to estimate T_z and f_x .

The output is formed by α , f_x , and the extrinsic parameters of the viewing camera.

When using a calibration pattern like the one in Figure 6.1, the 3-D squares lie on two different planes. The intersections of lines defined by squares from different world planes do not correspond to any image vertices. You must therefore ensure that your implementation considers only the intersections of pairs of lines associated to the *same* plane of the calibration pattern.

The line equations in image coordinates are computed by least squares, using as many collinear edge points as possible. This process improves the accuracy of the estimates of line parameters and vertex location on the image plane.

6.2.3 Estimating the Image Center

In what follows we describe a simple procedure for the computation of the image center. As a preliminary step, we recall the definition of *vanishing points* from projective geometry, and state a simple theorem suggesting how to determine the image center through the orthocenter³ of a triangle in the image.

Definition: Vanishing Points

Let $L_i, i = 1, \dots, N$ be parallel lines in 3-D space, and l_i the corresponding image lines. Due to the perspective projection, the lines l_i appear to meet in a point p_i , called *vanishing point*, defined as the common intersection of all the image lines l_i .

Orthocenter Theorem: Image Center from Vanishing Points

Let T be the triangle on the image plane defined by the three vanishing points of three mutually orthogonal sets of parallel lines in space. The image center is the orthocenter of T .

The proof of the theorem is left as an exercise (Exercise 6.2). The important fact is, the theorem reduces the problem of locating the image center to one of **interweaving**

³The orthocenter of a triangle is the common intersection of the three altitudes.

Image lines, and can be created easily on a suitable calibration pattern. In fact, we can use the same calibration pattern (actually, the same image!) of Figure 6.1, already used for EXPL_PARS_CAL, so that EXPL_PARS_CAL and the new algorithm, IMAGE_CENTER_CAL, fit nicely together.

Algorithm IMAGE_CENTER_CAL

The input is an image of the calibration pattern in Figure 6.1, and the output of the first two steps of algorithm EXPL_PARS_CAL.

1. Compute the three vanishing points p_1 , p_2 , and p_3 , determined by the three bundles of lines obtained in step 2 of EXPL_PARS_CAL.
2. Compute the orthocenter, O , of the triangle $p_1p_2p_3$.

The output are the image coordinates of the image center, O .

^{} It is essential that the calibration pattern is imaged from a viewpoint guaranteeing that no vanishing point lies much farther than the others from the image center; otherwise, the image lines become nearly parallel, and small inaccuracies in the location of the lines result in large errors in the coordinates of the vanishing point. This can happen if one of the three mutually orthogonal directions is nearly parallel to the image plane, a situation to be definitely avoided. Even with a good viewpoint, it is best to determine the vanishing points using several lines and least squares.

^{} To improve the accuracy of the image center estimate, you should run IMAGE_CENTER_CAL with several views of the calibration patterns, and average the results.

Experience shows that an accurate location of the image center is not crucial for obtaining precise estimates of the other camera parameters (see Further Readings). Be careful, however, as accurate knowledge of the image center is required to determine the ray in space identified by an image point (as we shall see, for example, in Chapter 8).

6.3 Camera Parameters from the Projection Matrix

We now move on to the description of a second method for camera calibration. The new method consists in two sequential stages:

1. estimate the projection matrix linking world and image coordinates;
2. compute the camera parameters as closed-form functions of the entries of the projection matrix.

6.3.1 Estimation of the Projection Matrix

As we have seen in Chapter 2, the relation between the 3-D coordinates (X_i^w, Y_i^w, Z_i^w) of a point in space and the 2-D coordinates (x, y) of its projection on the image plane

can be written by means of a 3×4 projection matrix, M , according to the equation

$$\begin{pmatrix} u_i \\ v_i \\ w_i \end{pmatrix} = M \begin{pmatrix} X_i^w \\ Y_i^w \\ Z_i^w \\ 1 \end{pmatrix},$$

with

$$\begin{aligned} x &= \frac{u_i}{w_i} = \frac{m_{11}X_i^w + m_{12}Y_i^w + m_{13}Z_i^w + m_{14}}{m_{31}X_i^w + m_{32}Y_i^w + m_{33}Z_i^w + m_{34}} \\ y &= \frac{v_i}{w_i} = \frac{m_{21}X_i^w + m_{22}Y_i^w + m_{23}Z_i^w + m_{24}}{m_{31}X_i^w + m_{32}Y_i^w + m_{33}Z_i^w + m_{34}} \end{aligned} \tag{6.10}$$

The matrix M is defined up to an arbitrary scale factor and has therefore only 11 independent entries, which can be determined through a homogeneous linear system formed by writing (6.16) for at least 6 world-image point matches. However, through the use of calibration patterns like the one in Figure 6.1, many more correspondences and equations can be obtained and M can be estimated through least squares techniques. If we assume we are given N matches for the homogeneous linear system we have

$$Am = 0, \tag{6.17}$$

with

$$A = \begin{bmatrix} X_1 & Y_1 & Z_1 & 1 & 0 & 0 & 0 & -x_1X_1 & -x_1Y_1 & x_1Z_1 & x_1 \\ 0 & 0 & 0 & 0 & X_1 & Y_1 & Z_1 & -y_1X_1 & -y_1Y_1 & y_1Z_1 & y_1 \\ X_2 & Y_2 & Z_2 & 1 & 0 & 0 & 0 & -x_2X_2 & -x_2Y_2 & x_2Z_2 & x_2 \\ 0 & 0 & 0 & 0 & 0 & X_2 & Y_2 & -y_2X_2 & -y_2Y_2 & y_2Z_2 & y_2 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ X_N & Y_N & Z_N & 1 & 0 & 0 & 0 & -x_NX_N & -x_NY_N & x_NZ_N & x_N \\ 0 & 0 & 0 & 0 & X_N & Y_N & Z_N & -y_NX_N & -y_NY_N & y_NZ_N & y_N \end{bmatrix}$$

and

$$m = [m_{11}, m_{12}, \dots, m_{33}, m_{34}]^T.$$

Since A has rank 11, the vector m can be recovered from SVD-related techniques as the column of V corresponding to the zero (in practice the smallest) singular value of A , with $A = UDV^T$ (see Appendix, section A.6). In agreement with the above definition of M , this means that the entries of M are obtained up to an unknown scale factor. The following is the detailed implementation of a method for estimating the matrix M :