

A method of visual metrology from uncalibrated images

Ze zhi Chen ^{*}, Nick Pears, Bojian Liang

Department of Computer Science, University of York, York YO10 5DD, UK

Received 26 May 2005; received in revised form 28 October 2005

Available online 30 May 2006

Communicated by Prof. H.H.S. Ip

Abstract

A method of measuring the height of any feature above a reference plane from a pair of uncalibrated images, separated by a (near) pure translation is presented. The output of the algorithm is a feature height, expressed as a fraction of the height of the camera above the reference plane. There are three contributions. Firstly a robust method of computing the dual epipole or focus of expansion (FOE) under pure translation is presented. Secondly, a novel reciprocal-polar (RP) image rectification scheme is presented, which allows planar image motion, expressed as a planar homography, to be accurately detected and recovered by 1D correlation. The technique can work even when there are no corner features on the reference plane and even over large image distortions caused by large camera motion, which would cause correlation techniques in the original image space to fail. Thirdly, we present a projective construct to enable measurement of the relative (or affine) feature height. Results show that our algorithm performs very well against outliers and noise. The mean of absolute error is 1.8 mm, and the mean of relative error is only 0.13% with two outliers removed.
© 2006 Elsevier B.V. All rights reserved.

Keywords: Reciprocal-polar (RP) transform; Homography matrix; Fundamental matrix; Metrology; Focus of expansion (FOE); Uncalibrated image

1. Introduction

The goal of this paper is to be able to measure the height of features above a reference plane from two uncalibrated images separated by a (near) pure translation. We have used this to detect obstacles in a mobile robot application and measure the height of doorways, windows and furniture in our laboratory in a visual metrology application. In both of these applications we use the laboratory floor as our reference plane. In the mobile robot application, we know the height of the camera above the reference plane, whereas in the metrology experiments, we know the height of one object in the image. There are several different methods of visual metrology in the literature. For example, 3D world structure computed from uncalibrated views of a scene,

given sufficient correspondences in general position, has already been in use for answering specific metric questions about the scene. The approach used by Tomasi and Kanade (1992) is known as the factorization method; Triggs (1996) later extended the factorization method to a projective camera model using epipolar constraints to calculate depth scale factors; Heyden and Berthilsson (1999) upgraded the affine approximation method to projective results employing iterative optimisation techniques. Another method is the camera-centered approach (Pollefeys et al., 1998) in which the first image is used as a reference to determine the projection matrices of other images in a projective frame under multiple geometric constraints. The world-centered approach is suitable for long image sequences while the camera-centered approach on the other hand is suitable for short image sequences. The photogrammetry technique concentrates on accuracy problems, whereas its derived techniques produces highly accurate three-dimensional models (Atkinson, 1996). However, photogrammetry techniques generally require heavy human interaction.

^{*} Corresponding author. Tel.: +44 (0)131 451 4179; fax: +44 (0)131 451 3327.

E-mail addresses: chen@macs.hw.ac.uk (Z. Chen), nep@cs.york.ac.uk (N. Pears), bojian@cs.york.ac.uk (B. Liang).

Many researchers have produced a number of techniques for measuring objects from images. Criminisi et al. (1999, 2000) proposed several methods to make measurement based on world planes from their perspective images. Reid and Zisserman (1996) and Kim et al. (1998) gave some methods for locating the 3D position and calculating the height of a soccer ball from monocular image sequences of soccer games. Liang and Pears (2002a,b) presented a method for computing the height of obstacles for robot visual navigation.

It is well known that, for a plane in general position, the homography is determined uniquely by the plane and vice versa. A homography can be determined by four correspondences in general position (with no three collinear correspondences) and is described by a 3×3 matrix H . Under pure camera translation in some direction \mathbf{t} , image motion sources or sinks to a point \mathbf{v} in the direction of \mathbf{t} , which is termed the focus of expansion (FOE). It is evident that \mathbf{v} is also the epipole e and e' for both views and the epipolar lines are radial and centred on the FOE. In this situation, it is possible to carry out an affine reconstruction using two images. The geometric significance is that the vanishing line of the reference plane can be determined by two pairs of correspondences, and the normal of ground plane can be determined by orientation the vanishing line (Hartley and Zisserman, 2001; Faugeras et al., 2001). The recovered H-matrix, in conjunction with a cross-ratio construct, is further applied to measure the height of a visual feature to the reference plane. The measurement method developed here is independent of the camera's intrinsic parameters (focal length, aspect ratio, principle point, skew) and extrinsic parameter (translation \mathbf{t}). The method does not require any geometric constraints on the 3D scene and the only limiting aspect of the method is that it is only applicable to (near) pure translation. However, in many applications, such as mobile robot navigation, deliberative (near) translation motion can be executed to probe the environment in terms of height, and measured heights of near zero are obstacle-free, navigable zones.

Six salient aspects of this method are as follows: (i) A robust method for estimating the FOE was developed yielding a highly accurate result. (ii) A new reciprocal-polar (RP) image rectification renders all planar motion to a pure shift, when the camera translation direction is parallel to the plane. This allows correlation based planar image motion recovery even over large perspective image distortions induced by large camera motions. (iii) Since the image motion recovery is correlation based, no corner feature correspondences are required on the reference plane, just local intensity variation. All coplanar pixels with local intensity variation contribute to the reference plane homography estimation, not just a potentially low density of corner feature matches. (iv) We show that the magnitude of image motion in the $1/r$ dimension of the rectified image pair follows a sinusoidal form along the θ dimension over a maximum of π radians, for the four degrees-of-freedom (dof) class of planar homographies called elations. (v) Ran-

dom sample consensus (RANSAC) and least squares (LS) estimation of the phase of the sinusoid yields a highly accurate vanishing line orientation of the reference plane (and simultaneously a segmentation of the reference plane) and this, along with the FOE, allows an accurate reference plane homography relation to be obtained. (vi) Finally, we present a new projective construction to compute affine height via the plane-and-parallax cue.

The paper is structured as follows. In the following section, we describe the relation between the fundamental matrix and the homography under pure translation. In Section 3, we make an incremental contribution to the estimation of the FOE or epipoles, by suggesting that a linear estimate should be refined using an optimisation based on the fundamental matrix. Section 4 presents the main result of the paper, where the analysis suggests the use of a reciprocal-polar rectification to recover the homography (H-matrix) representing reference plane motion across the image pair. Section 5 gives the result for computing affine height via the plane-and-parallax cue. Section 6 validates the method through several experiments, first using simple point correspondences and then using 1D intensity correlation, which is where the power of the method lies. Also, it is shown that, in terms of accuracy of the height, the method performs as well as the state-of-the-art. In the final section, conclusions are drawn.

2. The relation of the F and H matrix under pure translation

Images of points on a plane are related to corresponding image points in the second view by a planar homography (Hartley and Zisserman, 2001; Faugeras et al., 2001; Forsyth and Ponce, 2003). This is a projective relation since it depends only on the intersections of planes with lines. The homography map transfers points from one view to the other as if they were images of points on the plane. There are then two relations between the two views: first, through the epipolar geometry a point in one view determines a line in the other which is the image of the ray through that point; and second, through the homography a point in one view determines a point in the other which is the image of the intersection of the ray with a plane.

In pure translation of the camera, Eq. (1) gives the relation of the fundamental matrix, F , and the epipoles e and e' .

$$F = [e']_{\times} = [e]_{\times} = [v]_{\times} \quad (1)$$

where $[\cdot]_{\times}$ is a skew-symmetric matrix corresponding to the vector and prime denotes the second view in the image pair.

Let X_i be a set of points, which are coplanar in 3D Euclidean space. The images of X_i from two view points are related by a plane to plane projectivity or homography H , such that,

$$\lambda x'_i = Hx_i \quad (2)$$

where λ is a scalar, \mathbf{x}_i and \mathbf{x}'_i are homogenous image coordinates of the images of point X_i , and \mathbf{H} is a 3×3 matrix representing the homography. As homogenous coordinates are defined up to a scale factor, the H-matrix has only eight degrees-of-freedom (dof), and it can be determined by standard linear methods of four corresponding point pairs in general position (no three collinear). When the number of point pairs is more than four, a standard least square method or SVD (singular value decomposition) method can be used.

However, two corresponding point pairs fully determine the H-matrix under pure translation parallel to the plane. Suppose the cameras are calibrated with the origin of world coordinate system at the first camera, and the intrinsic parameters, \mathbf{K} , constant. The camera matrices, \mathbf{P} , \mathbf{P}' , for the two views are

$$\mathbf{P} = \mathbf{K}[\mathbf{I}|\mathbf{0}] \quad \mathbf{P}' = \mathbf{K}[\mathbf{R}|\mathbf{T}] \quad (3)$$

where \mathbf{I} is the identity matrix and \mathbf{R} and \mathbf{T} represent the rotational and translation displacement between the two views respectively. If the world plane π_E has normal, \mathbf{n} , and distance to origin, d , so that its coordinates are $\pi_E = (\mathbf{n}^T, d)^T$, then

$$\mathbf{H} = \mathbf{K}(\mathbf{R} - \lambda \mathbf{t}\mathbf{n}^T)\mathbf{K}^{-1} \quad (4)$$

where $\lambda = \frac{1}{d}$. For a pure translation, $\mathbf{R} = \mathbf{I}$, and so \mathbf{H} has the form

$$\mathbf{H} = \mathbf{I} - \lambda(\mathbf{K}\mathbf{t})(\mathbf{K}^{-T}\mathbf{n})^T \quad (5)$$

We note that $\mathbf{K}\mathbf{t}$ is the FOE $\mathbf{v} = \eta(x_f \ y_f \ 1)^T$, and $\mathbf{K}^{-T}\mathbf{n}$ is the vanishing line $\mathbf{l} = v(a_v \ b_v \ 1)^T$ corresponding to the plane π_E . Thus, we have

$$\mathbf{H} = \mathbf{I} - k\mathbf{v}\mathbf{l}^T \quad (6)$$

where k is a constant scalar. Since two corresponding point pairs fully define the FOE and vanishing line, so the

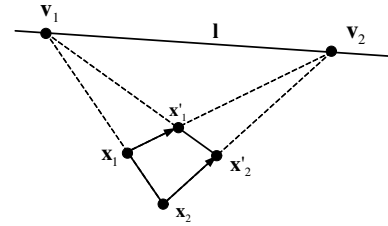


Fig. 1. Two corresponding point pairs fully define the vanishing point v_1 , v_2 and vanishing line l .

H-matrix can, in theory, also be fully determined by two corresponding matches (Fig. 1).

3. Accurate estimation of the FOE

We can detect (near) pure translation by intersecting all lines defined by all corner correspondences from the image pair and if great percentage of intersections lies in a small area (such as 85% of intersection should lie within a 50 pixels radius), then pure translation is assumed and the mean of the intersections is an initial value of FOE. The question now is: how can we calculate the FOE with high accuracy and high stability? A robust algorithm for estimating the FOE is summarized in Table 1.

4. Reciprocal-polar transform and recovery of H-matrix

Once the FOE has been computed, we shift image coordinates so that each image is centred on the FOE $\mathbf{v} = (x_f \ y_f \ 1)^T$. Let

Table 1
A robust method to estimate FOE

- (1) *Extract interest points*: Compute interest points in each image by using the Plessey–Harris corner detector (Harris and Stephens, 1988) or KLT algorithm (Shi and Tomasi, 1994) or SUSAN (Smith and Brady, 1997) method
- (2) *Putative correspondences*: Compute a set of interest point matches based on proximity and similarity of their intensity neighbourhood
- (3) *RANSAC robust estimation*: Repeat for m samples, where m is determined adaptively by using binning technique (Zhang, 1998)
 - (a) Select a random sample of at least 2 correspondences and compute the FOE by using the simultaneous equations:

$$\begin{cases} (\mathbf{x}_1 \times \mathbf{x}'_1) \cdot \mathbf{v} = 0 \\ (\mathbf{x}_2 \times \mathbf{x}'_2) \cdot \mathbf{v} = 0 \\ \vdots \\ (\mathbf{x}_n \times \mathbf{x}'_n) \cdot \mathbf{v} = 0 \end{cases}$$

where $\mathbf{x}_i \leftrightarrow \mathbf{x}'_i$ ($i = 1, 2, \dots, n$) is any pair of matching points in two images

- (b) Calculate the epipolar distance $f(\mathbf{v})$ for each putative correspondence

$$f(\mathbf{v}) = \left(\frac{1}{\sqrt{(\mathbf{F}\mathbf{x}_i)_1^2 + (\mathbf{F}\mathbf{x}_i)_2^2}} + \frac{1}{\sqrt{(\mathbf{F}^T\mathbf{x}'_i)_1^2 + (\mathbf{F}^T\mathbf{x}'_i)_2^2}} \right) |\mathbf{x}'_i{}^T \mathbf{F}\mathbf{x}_i|$$

where \mathbf{F} is the fundamental matrix and the subscripts 1, 2 denote vector components.

- (c) Compute the number of *inliers* consistent with \mathbf{v} by the number of correspondences for which $f(\mathbf{v}) < threshold$. Choose the FOE with the largest number of *inliers*

- (4) *Optimal estimation*: Re-estimate the FOE from all correspondences classified as *inliers*, by minimizing the object function $f(\mathbf{v})$
- (5) *Repeat* steps (3)–(4) until the number of correspondences are stable

$$\mathbf{T}_c = \begin{bmatrix} 1 & 0 & -x_f \\ 0 & 1 & -y_f \\ 0 & 0 & 1 \end{bmatrix} \quad (7)$$

After translation \mathbf{T}_c is applied, the FOE is at homogenous coordinates $\mathbf{v}_c = (0, 0, 1)^T$ and vanishing line becomes $\mathbf{l}_c = \mathbf{T}_c^{-T} \mathbf{l} = (a_v, b_v, \mathbf{v}^T \mathbf{l})^T$. The homography relating points in FOE centred coordinates is

$$\mathbf{H}_c = \mathbf{I} - k \mathbf{v}_c \mathbf{l}_c^T = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -ka_v & -kb_v & (1 - k \mathbf{v}^T \mathbf{l}_c) \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ s & \mu & q \end{bmatrix} \quad (8)$$

where $q = 1 - k \mathbf{v}^T \mathbf{l}$, $s = -ka_v$ and $\mu = -kb_v$. The homography \mathbf{H}_c has a very simple form and corresponding points, $\mathbf{x}_c, \mathbf{x}'_c$, in FOE centered images are related as follows:

$$\lambda \begin{bmatrix} x'_c \\ y'_c \\ 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ s & \mu & q \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ 1 \end{bmatrix} \quad (9)$$

If the robot's translation direction is parallel to the ground, $q = 1$ since the FOE lies on the vanishing line for this special motion. In this specialisation, the original H-matrix has four degrees-of-freedom, and is sometimes termed an elation (Hartley and Zisserman, 2001). Otherwise, the FOE is at a distance d_1 from the vanishing line and the five degree-of-freedom two-view planar relationship is termed a homology.

$$d_1 = \left| \frac{1 - q}{k \sqrt{a_v^2 + b_v^2}} \right| \quad (10)$$

Under the new homography, we have

$$(s x_c + \mu y_c + q)^2 (x_c'^2 + y_c'^2) = (x_c^2 + y_c^2) \quad (11)$$

If we define $\rho = \frac{1}{r}$, where $r = \sqrt{x_c^2 + y_c^2}$ is the Euclidean distance between an image point and the FOE in a frame, then taking square roots of Eq. (11)

$$f(\theta) = \rho' - q\rho = s \frac{x_c}{r} + \mu \frac{y_c}{r} = k_{s\mu} \sin(\theta + \alpha) \quad (12)$$

where θ is the angular position of a pixel in a frame centred on the FOE.

$$k_{s\mu} = \sqrt{s^2 + \mu^2}, \quad \sin \alpha = \frac{s}{k_{s\mu}}, \quad \cos \alpha = \frac{\mu}{k_{s\mu}}, \quad \tan \alpha = \frac{a_v}{b_v} \quad (13)$$

Eq. (12) indicates that we need to find three constants ($q, k_{s\mu}, \alpha$) in order to recover the homography and that the computation should be implemented in (ρ, θ) image space (note that a planar homology has five dof, but two have been recovered in the FOE computation). We call this reciprocal-polar (RP) image space. Thus, after computing the FOE, an interpolation procedure is used to generate a RP image for each image in the image pair.

Along any radius from the FOE, $f(\theta)$ is constant, so for any two pairs of correspondences (i, j)

$$q = \frac{\rho'_i - \rho'_j}{\rho_i - \rho_j} \quad (14)$$

Values of ρ in above equation can be determined by 1D windowed correlations between the two images in RP image space, i.e., along lines of constant θ in the RP image space. q can be obtained by using all the strong correlation results.

Given a value q , the RP image of the first image may be scaled along the ρ dimension. This allows correlations to be made along the ρ dimension at each angle θ_i . For each pixel in the first image, its position in RP image space is computed, and a 1D window is created around this position along the $q\rho$ dimension. We then correlate this window along the ρ dimension in the second RP image 2, at the same value of θ . This correlation process is possible because of the 'pure-shift' relation between ρ' and $q\rho$, expressed in Eq. (12) and the position of the maximum value of the correlation is related as a value of $f(\theta)$.

Eq. (12) indicates that correlation maxima and feature correlations in reciprocal-polar space, which are associated with a planar surface, lie on a sinusoid and the constants $(k_{s\mu}, \alpha)$ may be recovered by fitting a sinusoid to the data for $f(\theta)$.

Suppose that we have two values of $f(\theta), f_{i,j}$ measured at two angles, $\theta_{i,j}$, so that

$$f_i = k_{s\mu} \sin(\theta_i + \alpha), \quad f_j = k_{s\mu} \sin(\theta_j + \alpha) \quad (15)$$

$$\frac{f_i}{f_j} = \frac{\sin \theta_i + \cos \theta_i \tan \alpha}{\sin \theta_j + \cos \theta_j \tan \alpha} \quad (16)$$

collecting terms in $\tan \alpha$ and rearranging gives

$$\tan \alpha = \frac{f_j \sin \theta_i - f_i \sin \theta_j}{f_i \cos \theta_j - f_j \cos \theta_i} \quad (17)$$

Thus, in theory, a pair of f values, at different angular positions, for pixels belonging to the same plane, allows us to estimate the orientation of the vanishing line of that plane. Then, given the phase angle, α , corresponding to the orientation of the vanishing line, we can compute $k_{s\mu}$ from the Eq. (15).

In order to robustly and accurately estimate the vanishing line orientation from many correlation maxima and feature correspondences in reciprocal-polar space, many of which will not be associated with the ground plane, a RANdom SAMple Consensus (RANSAC) method (Fischer and Bolles, 1981) and iterated least-squares process are used. We define an optimization object function as:

$$\phi(\delta) = \sum_{i,j} \left| \delta - \frac{f_j \sin \theta_i - f_i \sin \theta_j}{f_i \cos \theta_j - f_j \cos \theta_i} \right| \quad (18)$$

where $\delta = \tan \alpha$, can be used to minimize $\phi(\delta)$, and determine the best set of inliers in the $f(\theta)$ data to a putative sinusoid. In this way, co-planar pixels may be grouped

without explicit construction of a homography matrix, although this is easily recovered in the FOE centred frame from the FOE and the two parameters of the sinusoid.

Let

$$f_i = k_{s\mu} \sin(\theta_i + \alpha) = s \cos \theta_i + \mu \sin \theta_i \quad (19)$$

Thus, for the all n inliers of the sinusoid model, we can write

$$\mathbf{AZ} = \mathbf{b} \quad (20)$$

where

$$\mathbf{A} = \begin{bmatrix} \cos \theta_1 & \sin \theta_1 \\ \vdots & \vdots \\ \cos \theta_n & \sin \theta_n \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} s \\ \mu \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix}$$

$$\mathbf{Z} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b} \quad (21)$$

The original homography matrix can be explicitly expressed as

$$\mathbf{H} = \mathbf{T}_c^{-1} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ s & \mu & q \end{bmatrix} \mathbf{T}_c \quad (22)$$

From Eq. (6), we have

$$k\mathbf{v}\mathbf{l}^T = \mathbf{I} - \mathbf{H} \quad (23)$$

If the j th row of the matrix $\mathbf{I} - \mathbf{H}$ is denoted by \mathbf{h}^T , then we may write

$$k \begin{pmatrix} x_j \mathbf{l}^T \\ y_j \mathbf{l}^T \\ \mathbf{l}^T \end{pmatrix} = \begin{pmatrix} \mathbf{h}^1 \mathbf{l}^T \\ \mathbf{h}^2 \mathbf{l}^T \\ \mathbf{h}^3 \mathbf{l}^T \end{pmatrix} \quad (24)$$

Then we can get the vanishing line

$$\mathbf{l}^T = \mathbf{h}^3 \mathbf{l}^T \quad (25)$$

5. Affine height measurements

The approach described above allows pixels to be classified as either belonging to the reference plane or not. For those non-reference plane regions, we would like to know their height above the reference plane. Our aim is to compute the relative height, h_r , of corner point A in Fig. 2, as a fraction of the height, h_c , of the camera optical centre O above the reference (ground) plane, when the camera undergoes pure translation t and the motion direction is parallel to the reference plane. Point A is the actual position of the corner point relative to the camera before the translation and the C is the position of the corner after the translation. Points A' and C' are the projections of these actual corner positions onto the ground plane. Points a and c are the image positions of the corner at positions A and C , respectively, and b is the predicted image position of the corner point by using H matrix of ground plane, if the corner point were to lie in the ground plane. Image point b is computed as $\mathbf{b} = \mathbf{H}\mathbf{a}$, using the recovered homography.

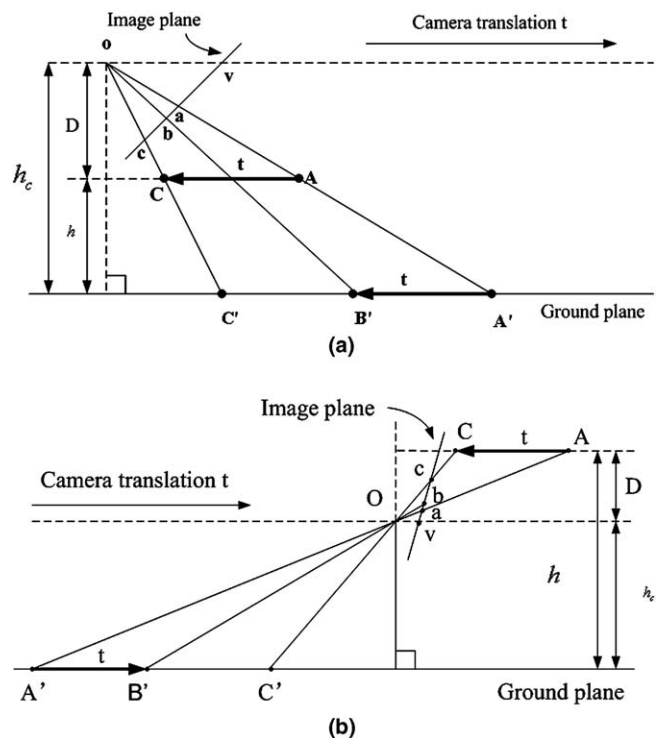


Fig. 2. Measuring the height of point A . (a) a lies below the vanishing line and (b) a lies above the vanishing line.

Fig. 2(a) shows that point a lies below the vanishing line. Because $AC//A'C'$, using similar triangles, and denoting the distance between point x and y as $d(x,y)$, we can get

$$h_r = \frac{h}{h_c} = 1 - \frac{D}{h_c} = 1 - \frac{d(O,C)}{d(O,C')} = 1 - \frac{d(A,C)}{d(A',C')} \quad (26)$$

For pure translation, $d(A,C) = d(A',B')$, so that

$$h_r = 1 - \frac{d(A',B')}{d(A',C')} \quad (27)$$

Now, the four image points (a, b, c, v), where v is the FOE, and their corresponding four ground plane points (A', B', C', ∞) are collinear. The cross-ratio for this set of points remains invariant under perspective projection transform and so we can get

$$\frac{d(A',B')}{d(A',C')} = \frac{d(a,b)d(c,v)}{d(a,c)d(b,v)} \quad (28)$$

Hence, the height of the corner point relative to the height of the optical centre is

$$h_r = 1 - \frac{d(a,b)d(c,v)}{d(a,c)d(b,v)} \quad (29)$$

If the point a lies above the vanishing line (Fig. 2(b)), we have

$$h_r = \frac{h}{h_c} = 1 + \frac{D}{h_c} = 1 + \frac{d(A,C)}{d(A',C')} \quad (30)$$

$$h_r = 1 + \frac{d(a,b)d(c,v)}{d(a,c)d(b,v)} \quad (31)$$

If point a lies on the vanishing line, then $h_r = 1$.

The algorithm for computing the height of obstacles is summarised as follows:

$$h_r = \begin{cases} 1 - \frac{d(\mathbf{a}, \mathbf{b})d(\mathbf{c}, \mathbf{v})}{d(\mathbf{a}, \mathbf{c})d(\mathbf{b}, \mathbf{v})} & \mathbf{a} \text{ lies below the vanishing line} \\ 1 & \mathbf{a} \text{ lies on the vanishing line} \\ 1 + \frac{d(\mathbf{a}, \mathbf{b})d(\mathbf{c}, \mathbf{v})}{d(\mathbf{a}, \mathbf{c})d(\mathbf{b}, \mathbf{v})} & \mathbf{a} \text{ lies above the vanishing line} \end{cases} \quad (32)$$

Note that h_r can be interpreted as the height of point A in units of height h_c . The absolute distance can be obtained from this distance ratio once the camera's height h_c is specified. However, it is sometimes more practical to determine the distance via a second measurement in the image, that of a known reference length. Note that this approach only needs the H-matrix of ground plane, and the tracked image correspondences a and c of the feature to determine the height to the ground plane. The main advantages compared with other method (such as the method of Criminisi et al. (2000) and Wilczkowiak et al. (2001)) is that this method does not need any camera calibration or any geometry constraints in the scene, and it can be used to compute the height from any feature point to the reference plane.

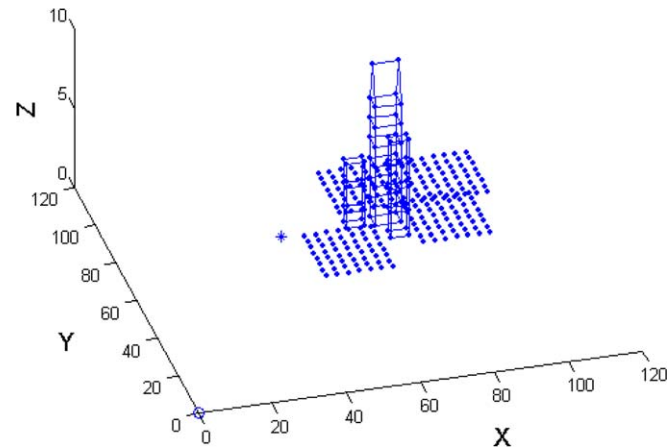


Fig. 3. A synthetic 3D scene.

6. Experimental results

A large number of synthetic data and real images were selected and intensive experimental work was carried out in order to test the robustness and accuracy of the method proposed in this paper.

6.1. Simulated experimental results

The simulated experiment was carried out on a 3D Euclidean model (Fig. 3). The intrinsic camera parameters were chosen as follows: the synthetic camera had a focal length which is invariant and an aspect ratio of one with no skew, $f/dx = f/dy = 500$, the principal point has a value of (319, 239) and the image size is 640×480 . Fig. 4 shows the images taken from the camera's translation direction parallel to the ground and includes their FOEs (o) and vanishing lines.

One hundred and fifty matches on the floor are used to compute the H matrix. Note that we compared our method with the standard least-squares estimation of the planar

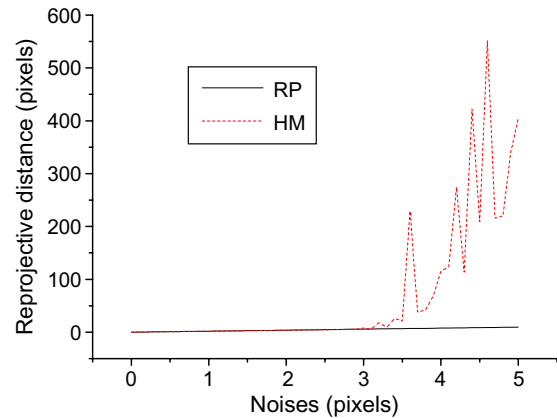


Fig. 5. Black line represents the reprojective errors of the RP method. Red dashed line represents the reprojective errors of the LS homography method (HM). (For interpretation of the references in colour in this figure legend, the reader is referred to the web version of this article.)

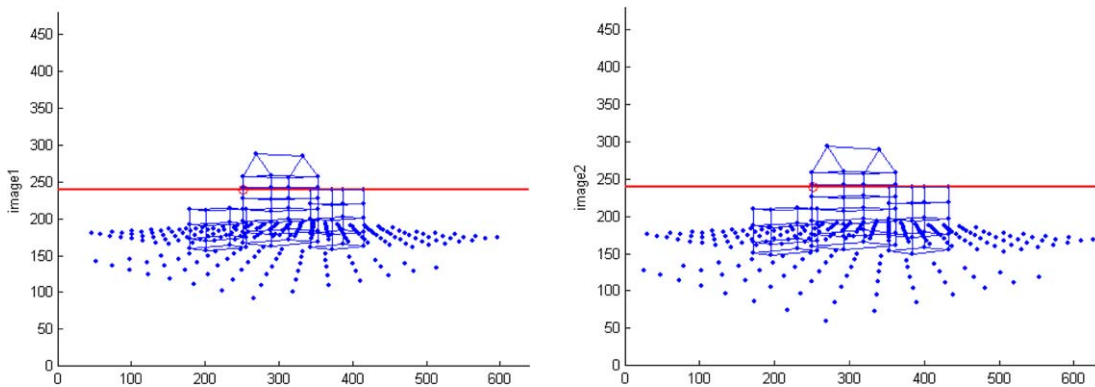
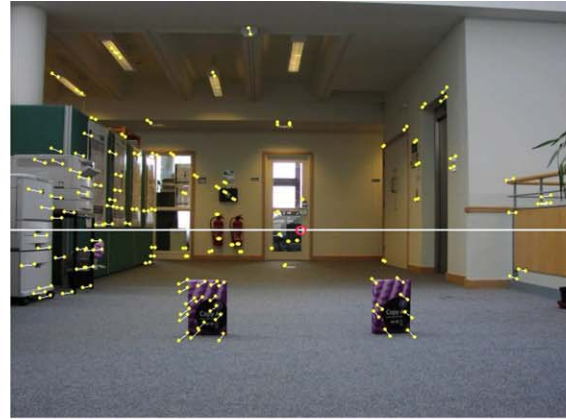


Fig. 4. Synthetic images used for simulation.



(a)



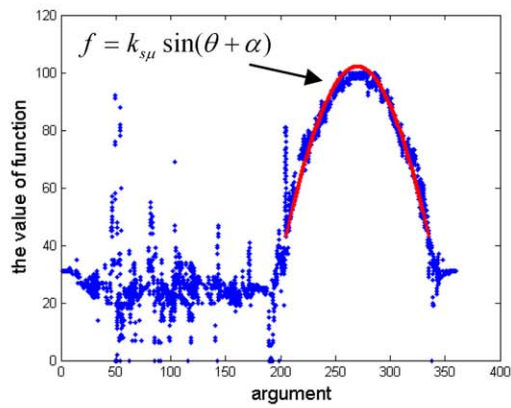
(b)



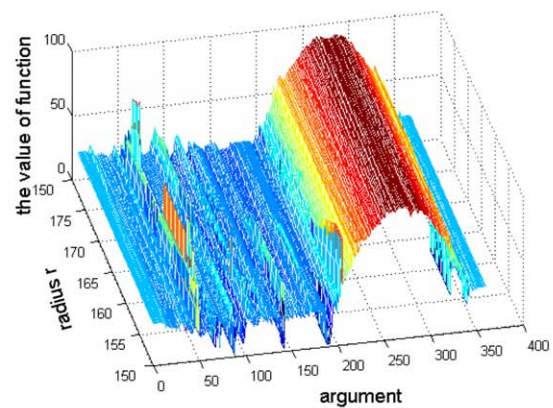
(c)



(d)



(e)



(f)

Fig. 6. (a) One of the original images. (b) Another image with their matching points (●), displacement of feature points, FOE (o) and vanishing line. (c) and (d) The reciprocal-polar images corresponding to (a) and (b), respectively. (e) The value of $f(\theta)$. (f) 3D sinusoid.

homography matrix (Hartley and Zisserman, 2001). Different Gaussian noises with mean from 0 to 5 pixels, variance one and standard deviation one were added to the image points to test the robustness of the reciprocal-polar (RP) method. The program executed 1000 times for every noise level. The means of the results are shown in Fig. 5. The results show that, if noise level is less than 3 pixels, then the two methods have almost the same robustness and accuracy. However, when noise level is increased to more than 3 pixels, our RP method is more robust than the standard least-squares (LS) homography matrix estimation method. After the H-matrix was computed, the height of obstacles can be recovered by using Eq. (32).

6.2. Real images experimental results

In the two visual metrology experiments (one indoor, one outdoor) presented here, $q = 1$ was assumed since two VGA frames (640×480 resolution) were captured with the translation direction parallel to the ground plane, and in both cases, the height AB , shown in Figs. 8 and 10, was used as the reference height.

The first experiment is conducted in an indoor scene, which resulted in 90.4% of intersections of all lines defined by all corner correspondences from the image pair to lie within a 50 pixels radius, which we classify as a (near) pure translation motion. The FOE (333.94, 264.41) was obtained using the robust estimation method presented here. Fig. 6(a) shows one of the sample image pair and Fig. 6(b) shows another image with their matching points (\cdot), displacement of feature points, FOE (o) and vanishing line. After the FOE had been obtained, the images were then converted to RP (ρ, θ) form. The RP images are shown in Fig. 6(c) and (d). These are quite severe distortions from the original, since radial lines centred on the FOE in the original image are now horizontal scan lines, and the intensity variations along these scan lines are flipped and non-linearly distorted due to the effect of the reciprocal operation. The relationship between the two figures is that reference plane pixels have undergone a pure shift along horizontal scan lines and the magnitude of this

shift follows a sinusoidal variation as we move down the image from one scan line to the next. Close inspection of the images does indeed reveal that this horizontal shifting process has occurred. Fig. 6(e) shows the reference plane image motion recovered using the 1D correlation method and the RANSAC method. The partial sinusoidal curve ($\theta \in [205^\circ, 335^\circ]$) is clearly shown which represents the motion of the ground plane in RP space. The phase of the curve is shown close to 180° rather than 0° because the direction of y in the image is directed upwards from the FOE rather than downwards. The three-dimensional sinusoidal form of RP ground plane motion, which is obtained when ρ is varied, is shown in Fig. 6(f). It clearly illustrates that the sinusoidal form is close to being constant irrespective of the radial distance of an image point from the FOE. In other words, planar image motion is a pure shift in RP space. Furthermore, we can get $\mu = -5.6374e - 004$, and $\tan \alpha = 0$ and hence

$$H = \begin{bmatrix} 1.0000 & -0.1883 & 49.7753 \\ 0 & 0.8509 & 39.4108 \\ 0 & -0.0006 & 1.1491 \end{bmatrix}$$

Fig. 7(a) and (b) show the reference plane (ground) recovered by the RP method and standard LS homography method, respectively. These results show the RP method performing better than the standard LS homography method.

Using the height of the box labelled “AB” (30 cm), we can recover the scale, after which we can correctly scale other heights in the scene, as shown in Fig. 8.

The second experiment implemented in an outdoor scene resulted in 89.9% of intersections of all lines defined by all corner correspondences from the image pair to lie within a 50 pixels radius after tracking the feature points. The FOE (261.17, 260.24) was obtained using the robust estimation method presented here. Fig. 9 shows one of the image pair with their matching points (\cdot), displacement of feature points, FOE and vanishing line. We computed $s = -8.0429e - 006$, $\mu = -2.9382e - 004$ and $\tan \alpha = 0.0274$ and hence

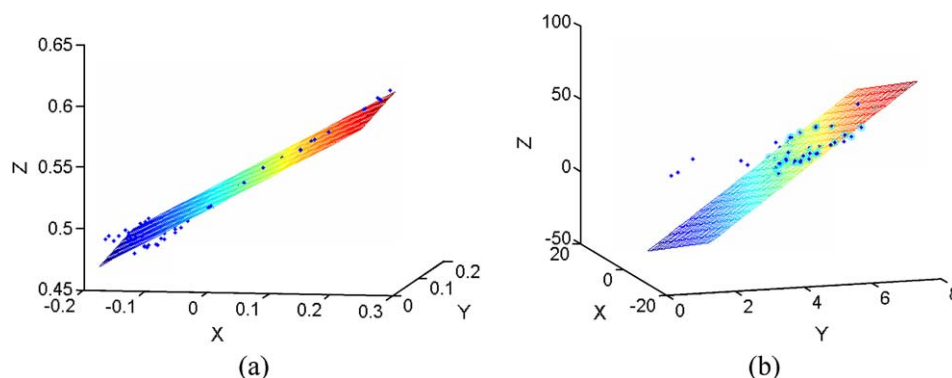


Fig. 7. (a) The reference plane (ground) recovered by RP and (b) the reference plane recovered by the standard LS homography method.

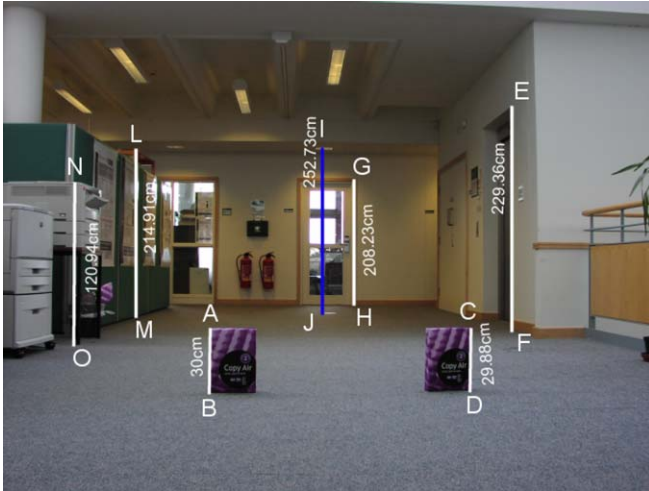


Fig. 8. Measuring the height of some interested point.

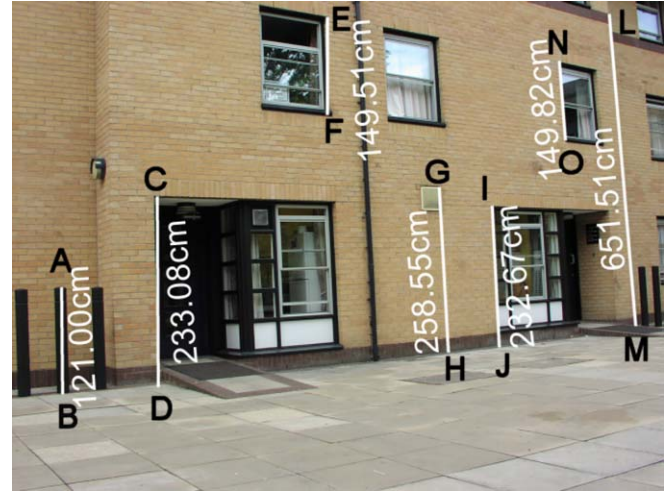


Fig. 10. Some measured results for outdoor scene.



Fig. 9. One of the image pair with their matching points (●), displacement of feature points, FOE (o) and vanishing line.

$$H = \begin{bmatrix} 0.9979 & -0.0767 & 20.5190 \\ -0.0021 & 0.9235 & 20.4462 \\ 0 & -0.0003 & 1.0786 \end{bmatrix}$$

Fig. 10 shows the measured results. The quantitative results are given in Table 2, where ‘TM’ are the manual (tape measure) measurements, ‘VM’ are the visual metrology results, ‘AE’ are the absolute errors and ‘RE’ are the relative errors. We find a mean absolute error of 7.1 mm

Table 2

Visual metrology results in centimetres

Seg.	Indoor				Outdoor			
	TM	VM	AE	RE (%)	TM	VM	AE	RE (%)
CD	30.0	29.88	0.12	0.400	233.1	233.08	0.02	0.0086
EF	227.7	229.36	1.66	0.7290	149.8	149.51	0.29	0.1936
GH	208.4	208.23	0.17	0.0816	258.7	258.55	0.15	0.0580
IJ	252.5	252.73	0.23	0.0911	233.1	232.67	0.43	0.1845
NO	121.1	120.94	0.16	0.1321	149.8	149.82	0.02	0.0134
LM	210.3	214.91	4.61	2.1921	None	651.51		

and mean relative error of 0.37%. If we remove the two rather inaccurate measurements EF and LM (indoor), the remaining measurements have a mean absolute error of 1.8 mm and a 0.13% mean relative error. We have also used the method presented in this paper to detect potential obstacles for robot navigation. In this application, we measure the height of contours around segmented image regions, and label them as obstacles if they are above a height threshold, set at a fraction of the robot wheel radius. Some of our obstacle detection results are shown in Fig. 11. Our system is coded in MATLAB and can run at around five frames per second on a PC with a Pentium IV processor running at 2.4 GHz and 512MB of memory. Finally note that our system performance is largely unaffected by the magnitude of the translation, with the proviso that sufficient image motion (more than a few pixels) is generated in order to detect that the camera motion is near pure translation.



Fig. 11. Obstacles detection results.

7. Conclusions

The main contributions in this work are (in the order of algorithm execution) (1) robust FOE estimation, (2) RP rectification, (3) planar image motion estimation, planar homography estimation and plane segmentation by robust estimation of a sinusoid in RP image motion space, and (4) a projective construction allowing affine height to a plane to be measured from an uncalibrated image pair. Intensive experimental work was carried out in order to test the accuracy of the method proposed in this paper. The results show that our algorithm performs very well to outliers and noise and provides a practical method of visual metrology.

Acknowledgement

The authors acknowledge the support of the UK DTI Aeronautics Research Programme.

References

- Atkinson, K.B. (Ed.), 1996. *Close Range Photogrammetry and Machine Vision*. Whittles Publishing.
- Criminisi, A., Reid, I., Zisserman, A., 1999. A plane measuring device. *Image Vision Comput.* 17 (8), 625–634.
- Criminisi, A., Reid, I., Zisserman, A., 2000. Single view metrology. *Internat. J. Comput. Vision* 40 (2), 123–148.
- Faugeras, O., Luong, Q.-T., Papadopoulos, T., 2001. *The Geometry of Multiple Images: The Laws that Govern the Formation of Multiple Images of a Scene and Some of their Applications*. The MIT Press, Cambridge, Massachusetts, London, England.
- Fischer, M.A., Bolles, R.C., 1981. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Comm. Assoc. Comput. Machine* 24 (6), 381–395.
- Forsyth, D.A., Ponce, Jean, 2003. *Computer Vision: A Modern Approach*. Prentice Hall, Upper Saddle River, NJ.
- Harris, C.J., Stephens, M., 1988. A combined corner and edge detector. In 4th Alvey Vision Conference, Manchester, pp. 147–151.
- Hartley, R., Zisserman, A., 2001. *Multiple View Geometry in Computer Vision*. Cambridge University Press (reprinted).
- Heyden, A., Berthilsson, R., 1999. An iterative factorization method for projective structure and motion from image sequences. *Image Vision Comput.* 17 (13), 981–991.
- Kim, T., Seo, Y., Hong, K., 1998. Physics-based 3D position analysis of a soccer ball from monocular image sequences. *Proc. ICCV*, 721–726.
- Liang, B., Pears, N.E., 2002a. Ground plane segmentation from multiple visual cues. *Second Internat. Conf. on Image and Graphics*, Hefei, China, *Proc. of SPIE*, vol. 4875, pp. 822–829.
- Liang, B., Pears, N.E., 2002b. Visual navigation using planar homographies. In: *Proc. of IEEE Internat. Conf. on Robotics and Automation*. Washington, DC, USA, vol. 1, pp. 205–210.
- Pollefeys, M., Koch, R., Vergauwen, M., Gool, L.Van., 1998. Metric 3D surface reconstruction from uncalibrated image sequences. In: *Proc. SMILE Workshop (Post-ECCV'98)*, LNCS1506. Springer-Verlag, pp. 138–153.
- Reid, I., Zisserman, A., 1996. Goal-directed Video Metrology. In: Cipolla, R., Buxton, B. (Eds.), *Proc. ECCV*, vol. II. Springer, pp. 647–658.
- Shi, Jianbo, Tomasi, Carlo, 1994. Good features to track. *IEEE Conf. on Computer and Pattern Recognition*, Seattle, USA, pp. 593–600.
- Smith, S., Brady, J., 1997. SUSAN—a new approach to low level image processing. *Internat. J. Comput. Vision* 23 (1), 45–78.
- Tomasi, C., Kanade, T., 1992. Shape and motion from image streams under orthography: A factorization approach. *Internat. J. Comput. Vision* 9 (2), 137–154.
- Triggs, B., 1996. Factorization methods for projective structure and motion. In: *Proc. Conf. on Computer Vision and Pattern Recognition*, pp. 845–851.
- Wilczkowiak, M., Boyer, E., Sturm, P.F., 2001. Camera calibration and 3d reconstruction from single images using parallelepipeds. *Internat. Conf. on Computer Vision*, Vancouver, pp. 142–148.
- Zhang, Z., 1998. Determining the epipolar geometry and its uncertainty: a review. *Internat. J. Comput. Vision* 27 (2), 161–195.