

Visual Servoing

Vision allows a robotic system to obtain geometrical and qualitative information on the surrounding environment to be used both for motion planning and control. In particular, control based on feedback of visual measurements is termed *visual servoing*. In the first part of this chapter, some basic algorithms for image processing, aimed at extracting numerical information referred to as *image feature parameters*, are presented. These parameters, relative to images of objects present in the scene observed by a camera, can be used to estimate the pose of the camera with respect to the objects and vice versa. To this end, analytical *pose estimation methods*, based on the measurement of a certain number of points or *correspondences* are presented. Also, numerical pose estimation methods, based on the integration of the linear mapping between the camera velocity in the operational space and the time derivative of the feature parameters in the *image plane*, are introduced. In cases in which multiple images of the same scene, taken from different viewpoints, are available, additional information can be obtained using *stereo vision* techniques and *epipolar geometry*. A fundamental operation is also *camera calibration*; to this end, a calibration method based on the measurement of a certain number of correspondences is presented. Then, the two main approaches to visual servoing are introduced, namely *position-based visual servoing* and *image-based visual servoing*, as well as a scheme, termed *hybrid visual servoing*, which combines the benefits of both approaches.

10.1 Vision for Control

Vision plays a key role in a robotic system, as it can be used to obtain geometrical and qualitative information on the environment where the robot operates, without physical interaction. Such information may be employed by the control system at different levels, for the sole task planning and also for feedback control.

As an example, consider the case of a robot manipulator, equipped with a camera, which has to grasp an object using a gripper. Through vision the robot may acquire information capable of identifying the relative pose of the object with respect to the gripper. This information allows the control system to plan a trajectory leading the manipulator in an appropriate grasping configuration, computed on the basis of the pose and of the shape of the object, from which the closure of the gripper can be commanded.

The planned trajectory can be executed using a simple motion controller. In this approach, termed *look-and-move*, visual measurements are used in open loop, making the system very sensitive to uncertainties due, for instance, to poor positioning accuracy of the manipulator or to the fact that the object may have moved while the gripper reaches the grasp position.

On the other hand, in *vision-based control* or *visual servoing*, the visual measurements are fed back to the control to compute an appropriate error vector defined between the current pose of the object and the pose of the manipulator's end-effector.

A key characteristic of visual servoing, compared to motion and force control, is the fact that the controlled variables are not directly measured by the sensor, but are obtained from the measured quantities through complex elaborations, based on algorithms of *image processing* and *computational vision*.

In Sect. 5.4.3 it was shown that a monochrome camera simply provides a two-dimensional matrix of values of light intensity. From this matrix, the so-called *image feature parameters* are to be extracted in real time. The geometric relationships between one or more two-dimensional views of a scene and the corresponding 3D space are the basis of techniques of *pose estimation* of objects in the manipulator workspace or of the end-effector with respect to the surrounding objects. In this regard, of fundamental importance is the operation of *camera calibration*, which is necessary for calculating the *intrinsic parameters*, relating the quantities measured in the image plane to those referred to the camera frame, and the *extrinsic parameters*, relating the latter to quantities defined in a frame attached to the manipulator.

The vision-based control schemes can be divided into two categories, namely, those that realize *visual servoing in operational space*, also termed *position-based visual servoing*, and those that realize *visual servoing in the image space*, also known as *image-based visual servoing*. The main difference lies in the fact that the schemes of the first category use visual measurements to reconstruct the relative pose of the object with respect to the robot, or vice versa, while the schemes of the second category are based on the comparison of the feature parameters of the image of the object between the current and the desired pose. There are also schemes combining characteristics common to both categories, that can be classified as *hybrid visual servoing*.

Another aspect to be considered for vision-based control is the type of camera (colour or monochrome, resolution, fixed or variable focal length, CCD or CMOS technology). In this chapter, only the case of monochrome cameras with fixed focal length will be considered.

Equally important is the choice of the number of cameras composing the visual system and their location; this issue is briefly discussed in the following.

10.1.1 Configuration of the Visual System

A visual system may consist of only one camera, or two or more cameras. If more cameras are used to observe the same object of a scene, it is possible to retrieve information about its depth by evaluating its distance with respect to the visual system. This situation is referred to as *3D vision* or *stereo vision*, where the term *stereo* derives from the Greek and means solid. The human capability of perceiving objects in three dimensions relies on the fact that the brain receives the same images from two eyes, observing the same scene from slightly different angles.

It is clear that 3D vision can be achieved even with one camera, provided that two images of the same object, taken from two different poses, are available. If only a single image is available, the depth can be estimated on the basis of certain geometrical characteristics of the object known in advance. This means that, in many applications, mono-camera systems are often preferred to multi-camera systems, because they are cheaper and easier to calibrate, although characterized by lower accuracy.

Another feature that distinguishes visual systems for robot manipulators is the placement of cameras. For mono-camera systems there are two options: the *fixed configuration*, often referred to as *eye-to-hand*, where the camera is mounted in a fixed location, and the *mobile configuration*, or *eye-in-hand*, with the camera attached to the robot. For multi-camera systems, in addition to the mentioned solutions, it is also possible to consider the *hybrid configuration*, consisting of one or more cameras in eye-to-hand configuration, and one or more cameras in eye-in-hand configuration.

In the eye-to-hand configuration, the visual system observes the objects to be manipulated by a fixed pose with respect to the base frame of the manipulator. The advantage is that the camera field of view does not change during the execution of the task, implying that the accuracy of such measurements is, in principle, constant. However, in certain applications, such as assembly, it is difficult to prevent that the manipulator, moving in the camera field of view, occludes, in part or in whole, the view of the objects.

In the eye-in-hand configuration, the camera is placed on the manipulator and can be mounted both before and after the wrist. In the first case, the camera can observe the end-effector by a favourable pose and without occlusions caused by the manipulator arm; in the latter case, the camera is attached to the end-effector and typically observes only the object. In both situations, the camera field of view changes significantly during the motion and this produces a high variability in the accuracy of measurements. However, when the end-effector is close to the object, the accuracy becomes almost constant and is usually higher than that achievable with eye-to-hand cameras, with the advantage that occlusions are virtually absent.

Finally, hybrid configuration combines the benefits of the other two configurations, namely, ensures a good accuracy throughout the workspace, while avoiding the problems of occlusions.

A separate category is represented by robotic heads, which are typically equipped with a stereo vision system consisting of two cameras mounted on motorized mechanisms that allow for yaw motion, or *pan*, and pitch motion, or *tilt*, hence the name of *pan-tilt* cameras.

In this chapter, only schemes based on a single eye-in-hand camera will be considered. The extension of the algorithms to the case of a eye-to-hand camera, or to the case of multiple cameras, requires only minor modifications.

10.2 Image Processing

Visual information, unlike the information provided by other types of sensors, is very rich and varied and thus requires complex and computational expensive transformations before it can be used for controlling a robotic system. The objective of these transformations is the extraction of numerical information from the image, which provides a synthetic and robust description of the objects of interest in the scene, through the so-called *image feature parameters*.

To this end, two basic operations are required. The first is so-called *segmentation*, which aims at obtaining a representation suitable for the identification of measurable features of the image. The subsequent operation, termed *interpretation* is concerned with the measurement of the feature parameters of the image.

The source information is contained in a framestore, namely the two-dimensional memory array representing the spatial sample of the image. On the set of pixels the so-called *image function* is defined which, in general, is a vector function whose components represent the values of one or more physical quantities related to the pixel, in a sampled and quantized form.

For example, in the case of color images, the image function defined on a pixel of coordinates (X_I, Y_I) has three components $I_r(X_I, Y_I)$, $I_g(X_I, Y_I)$ and $I_b(X_I, Y_I)$, corresponding to the light intensity in the wavelengths of red, green and blue. For a monochrome black-and-white image, the image function is scalar and coincides with the light intensity in shades of gray $I(X_I, Y_I)$, also referred to as *gray level*. In the following, for simplicity, only monochrome images will be considered.

The number of gray levels depends on the adopted grey-scale resolution. In all cases, the gray scale is bounded by two gray levels, black and white, corresponding to the minimum and maximum measurable light intensity respectively. Most current acquisition equipments adopt a scale consisting of 256 gray levels, that can be represented by a single byte of memory.

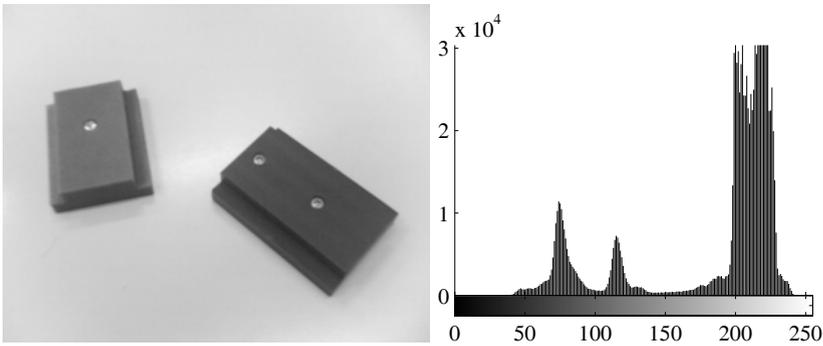


Fig. 10.1. Black-and-white image and corresponding gray-level histogram on the right

A representation of the framestore which is particularly useful for subsequent processing is the gray-level histogram, which provides the frequency of occurrence of each gray level in the image.

Where the gray levels are quantized from 0 to 255, the value $h(p)$ of the histogram at a particular gray level $p \in [0, 255]$ is the number of image pixels with gray level p . If this value is divided by the total number of pixels, the histogram is termed *normalized histogram*.

Figure 10.1 shows a black-and-white image and the corresponding gray-level histogram. Proceeding from left to right, three main peaks can be observed — from left to right — corresponding to the darkest object, the lightest object, and the background.

10.2.1 Image Segmentation

Segmentation consists of a grouping process, by which the image is divided into a certain number of groups, referred to as *segments*, so that the component of each group are similar with respect to one or more characteristics. Typically, distinct segments of the image correspond to distinct objects of the environment, or homogeneous object parts.

There are two complementary approaches to the problem of image segmentation: one is based on finding connected *regions* of the image, the other is concerned with detection of *boundaries*. The objective of region-based segmentation is that of grouping sets of pixels sharing common features into two-dimensional connected areas, with the implicit assumption that the resulting regions correspond to real-world surfaces or objects. On the other hand, boundary-based segmentation is aimed at identifying the pixels corresponding to object contours and isolating them from the rest of the image. The boundary of an object, once extracted, can be used to define the position and shape of the object itself.

The complementarity of the two approaches relies on the fact that a boundary can be achieved by isolating the contours of a region and, conversely, a region can be achieved simply by considering the set of pixels contained within a closed boundary.

The problem of segmentation is not trivial and there exist many solutions, some of which are sketched below. From the point of view of memory usage, boundary-based segmentation is more convenient, since boundaries contain a reduced number of pixels. However, from the computational load point of view, region-based segmentation is faster because it requires a reduced number of memory accesses.

Region-based segmentation

The central idea underlying region-based segmentation techniques is that of obtaining connected regions by continuous merging of initially small groups of adjacent pixels into larger ones.

Merging of two adjacent regions may happen only if the pixels belonging to these regions satisfy a common property, termed *uniformity predicate*. Often the uniformity predicate requires that the gray level of the pixels of the region belongs to a given interval.

In many applications of practical interest a thresholding approach is adopted and a light intensity scale composed of only two values (0 and 1) is considered. This operation is referred to as *binary segmentation* or image *binarization*, and corresponds to separating one or more objects present in the image from the background by comparing the gray level of each pixel with a threshold l . For light objects against a dark background, all the pixels whose gray level is greater than the threshold are considered to belong to a set S_o , corresponding to the objects, while all the other pixels are considered to belong to a set S_b corresponding to the background. It is obvious that this operation can be reversed for dark objects against a light background. When only an object is present in the image, the segmentation ends with the detection of sets S_o and S_b , representing two regions; in the presence of multiple objects, a further elaboration is required to separate the connected regions corresponding to the single objects. The image obtained assigning a light intensity equal to 0 to all the pixels of set S_o , and a light intensity equal to 1 to all the pixels of set S_b , or vice versa, is termed *binary image*.

A crucial factor for the success of binary segmentation is the choice of the threshold. A widely adopted method for selecting the threshold is based on the gray-level histogram, under the assumption that it contains clearly distinguishable minimum and maximum values, corresponding to the gray levels of the objects and of the background; the peaks of the histogram are also termed *modes*. For dark objects against a light background, the background corresponds to the mode which is located further to the right — as, for example, in the case of Fig. 10.1 — and the threshold can be chosen at the closest minimum to the left. For light objects against a dark background, the

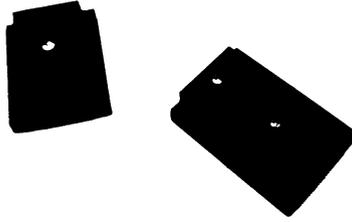


Fig. 10.2. Binary image corresponding to image of Fig. 10.1

background corresponds to the mode which is located further to the left and the threshold should be selected accordingly. With reference to Fig. 10.1, the threshold can be set to $l = 152$. The corresponding binary image is reported in Fig. 10.2.

In practice, the gray-scale histogram is noisy and the modes are difficult to identify. Often, there is no clear separation between the gray levels of the objects and those of the background. To this end, various techniques have been developed to increase the robustness of binary segmentation, which require appropriate filtering of the image before binarization and the adoption of algorithms for automatic selection of the threshold.

Boundary-based segmentation

Boundary-based segmentation techniques usually obtain a boundary by grouping many single local edges, corresponding to local discontinuities of image gray level. In other words, local edges are sets of pixels where the light intensity changes abruptly.

The algorithms for boundary detection first derive an intermediate image based on local edges from the original gray-scale image, then they construct short-curve segments by edge linking, and finally obtain the boundaries by joining these curve segments through geometric primitives often known in advance.

Boundary-based segmentation algorithms vary in the amount of a priori knowledge they incorporate in associating or linking the edges and their effectiveness clearly depends on the quality of the intermediate image based on local edges. The more reliable the local edges in terms of their position, orientation and ‘authenticity’, the easier the task of the boundary detection algorithm.

Notice that edge detection is essentially a filtering process and can often be implemented via hardware, whereas boundary detection is a higher level task usually requiring more sophisticated software. Therefore, the current trend

is that of using the most effective edge detector to simplify the boundary detection process. In the case of simple and well-defined shapes, boundary detection becomes straightforward and segmentation reduces to the sole edge detection.

Several edge detection techniques exist. Most of them require the calculation of the gradient or of the Laplacian of function $I(X_I, Y_I)$.

Since a local edge is defined as a transition between two regions of significantly different gray levels, it is obvious that the spatial gradient of function $I(X_I, Y_I)$, which measures the rate of change of the gray level, will have large magnitude close to these transitional boundary areas. Therefore, edge detection can be performed by grouping the pixels where the magnitude of the gradient is greater than a threshold. Moreover, the direction of the gradient vector will be the direction of maximum variation of the gray level.

Again, the choice of the value of the threshold is extremely important; in the presence of noise, the threshold is the result of a trade-off between the possibility of losing valid edges and that of detecting false edges.

For gradient computation, it suffices to evaluate the directional derivatives of function $I(X_I, Y_I)$ along two orthogonal directions. Since this function is defined on a discrete set of pixels, the derivatives are computed in an approximate way. The essential differences between gradient-based edge detection techniques are the directions used for the computation of the derivatives and the manner in which they approximate these derivatives and compute the gradient magnitude.

The most common operator for the computation of the gradient is that approximating the derivative along directions X_I and Y_I with the first differences:

$$\begin{aligned}\Delta_1 &= I(X_I + 1, Y_I) - I(X_I, Y_I) \\ \Delta_2 &= I(X_I, Y_I + 1) - I(X_I, Y_I).\end{aligned}$$

Other operators, less sensitive to noise effects are, for example, the *Roberts operator*, based on the first differences computed along the diagonals of a (2×2) square of pixels:

$$\begin{aligned}\Delta_1 &= I(X_I + 1, Y_I + 1) - I(X_I, Y_I) \\ \Delta_2 &= I(X_I, Y_I + 1) - I(X_I + 1, Y_I),\end{aligned}$$

and the *Sobel operator*, defined on a (3×3) square of pixels:

$$\begin{aligned}\Delta_1 &= (I(X_I + 1, Y_I - 1) + 2I(X_I + 1, Y_I) + I(X_I + 1, Y_I + 1)) \\ &\quad - (I(X_I - 1, Y_I - 1) + 2I(X_I - 1, Y_I) + I(X_I - 1, Y_I + 1)) \\ \Delta_2 &= (I(X_I - 1, Y_I + 1) + 2I(X_I, Y_I + 1) + I(X_I + 1, Y_I + 1)) \\ &\quad - (I(X_I - 1, Y_I - 1) + 2I(X_I, Y_I - 1) + I(X_I + 1, Y_I - 1)).\end{aligned}$$

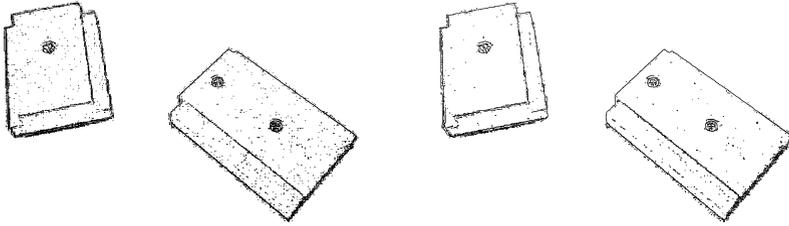


Fig. 10.3. Contours of image of Fig. 10.1 obtained using Roberts (*left*) and Sobel (*right*) operators

Then, the approximated magnitude, or norm, of gradient $G(X_I, Y_I)$ can be evaluated using one of the following two expressions:

$$G(X_I, Y_I) = \sqrt{\Delta_1^2 + \Delta_2^2}$$

$$G(X_I, Y_I) = |\Delta_1| + |\Delta_2|,$$

and direction $\theta(X_I, Y_I)$ with the relationship

$$\theta(X_I, Y_I) = \text{Atan2}(\Delta_2, \Delta_1).$$

Figure 10.3 shows the images obtained from that of Fig. 10.1 by applying the gradient operators of Sobel and Roberts and binarization; the thresholds have been set to $l = 0.02$ and $l = 0.0146$, respectively.

An alternative edge detection method is based on the *Laplacian operator*, which requires the computation of the second derivatives of function $I(X_I, Y_I)$ along two orthogonal directions. Also in this case, suitable operators are used to discretize the computation of derivatives. One of the most common approximations is the following:

$$L(X_I, Y_I) = I(X_I, Y_I) - \frac{1}{4} (I(X_I, Y_I + 1) + I(X_I, Y_I - 1) \\ + I(X_I + 1, Y_I) + I(X_I - 1, Y_I)).$$

In this case, the pixels of the contour are those where the Laplacian is lower than a threshold. The reason is that the Laplacian is null at the points of maximum magnitude of the gradient. The Laplacian, unlike the gradient, does not provide directional information; moreover, being based on the calculation of second derivatives, it is more sensitive to noise than the gradient.

10.2.2 Image Interpretation

Image interpretation is the process of calculating the image feature parameters from the segments, whether they are represented in terms of boundaries or in terms of regions.

The feature parameters used in visual servoing applications sometimes require the computation of the so-called *moments*. These parameters are defined on a region \mathcal{R} of the image and can be used to characterize the position, orientation and shape of the two-dimensional object corresponding to the region itself.

The general definition of moment $m_{i,j}$ of a region \mathcal{R} of a framestore, with $i, j = 0, 1, 2, \dots$, is the following:

$$m_{i,j} = \sum_{X_I, Y_I \in \mathcal{R}} I(X_I, Y_I) X_I^i Y_I^j.$$

In the case of binary images, by assuming the light intensity equal to one for all the points of region \mathcal{R} , and equal to zero for all the points not belonging to \mathcal{R} , the following simplified definition of moment is obtained:

$$m_{i,j} = \sum_{X_I, Y_I \in \mathcal{R}} X_I^i Y_I^j. \quad (10.1)$$

In view of this definition, moment $m_{0,0}$ coincides with the area of the region, computed in terms of the total number of pixels of region \mathcal{R} .

The quantities

$$\bar{x} = \frac{m_{1,0}}{m_{0,0}} \quad \bar{y} = \frac{m_{0,1}}{m_{0,0}}$$

define the coordinates of the so-called *centroid* of the region. These coordinates can be used to detect uniquely the position of region \mathcal{R} on the image plane.

Using an analogy from mechanics, region \mathcal{R} can be seen as a two-dimensional rigid body of density equal to light intensity. Hence, moment $m_{0,0}$ corresponds to the mass of the body and the centroid corresponds to the centre of mass.

The value of moment $m_{i,j}$ in (10.1) depends on the position of region \mathcal{R} in the image plane. Therefore, the so-called *central moments* are often considered, defined as

$$\mu_{i,j} = \sum_{X_I, Y_I \in \mathcal{R}} (X_I - \bar{x})^i (Y_I - \bar{y})^j,$$

which are invariant with respect to translation.

According to the mechanical analogy, it is easy to recognize that the central moments of second order $\mu_{2,0}$ and $\mu_{0,2}$ have the meaning of inertia moments with respect to axes X_I and Y_I respectively, while $\mu_{1,1}$ is an inertia product, and the matrix

$$\mathcal{I} = \begin{bmatrix} \mu_{2,0} & \mu_{1,1} \\ \mu_{1,1} & \mu_{0,2} \end{bmatrix},$$

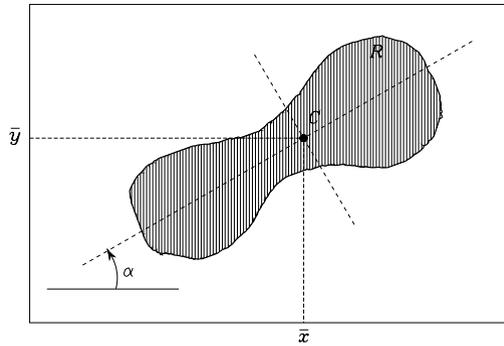


Fig. 10.4. Region of a binary image and some feature parameters

has the meaning of inertia tensor relative to the centre of mass. The eigenvalues of matrix \mathcal{I} define the principal moments of inertia, termed *principal moments* of the region and the corresponding eigenvectors define the principal axes of inertia, termed *principal axes* of the region.

If region \mathcal{R} is asymmetric, the principal moments of \mathcal{I} are different and it is possible to characterize the orientation of \mathcal{R} in terms of the angle α between the principal axis corresponding to the maximum moment and axis X . This angle can be computed with the equation (see Problem 10.1)

$$\alpha = \frac{1}{2} \tan^{-1} \left(\frac{2\mu_{1,1}}{\mu_{2,0} - \mu_{0,2}} \right). \quad (10.2)$$

As an example, in Fig. 10.4, the region of a binary image is shown; centroid C , the principal axes, and the angle α are evidenced.

Notice that the moments and the corresponding parameters can also be computed from the boundaries of the objects; moreover, these quantities are especially useful to characterize objects of generic form. Often, however, the objects present in the scene, especially those manufactured, have geometric characteristics which are useful to take into account for image interpretation.

For example, many objects have edges that, in the image plane, correspond to the intersection of linear parts of contour or to contour points of high curvature. The coordinates of these points on the image plane can be detected using algorithms robust against the noise, and therefore can be used as feature parameters of the image. They are usually termed *feature points*.

In other cases, it is possible to identify true geometric primitives such as lines or line segments, which are projections of linear edges or solids of revolution (cones, cylinders), or ellipses, obtained as projections of circles or spheres. These primitives can be characterized on the image plane in terms of a minimum set of parameters. For example, a line segment can be characterized by the coordinates of its endpoints, or alternatively, by the coordinates of its midpoint (centroid), its length (moment $m_{0,0}$) and its orientation (angle α); in both cases, the characterization of the line segment requires four parameters.

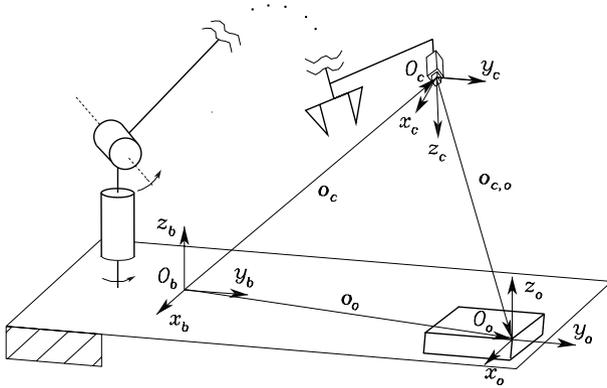


Fig. 10.5. Reference frames for an eye-in-hand camera

10.3 Pose Estimation

Visual servoing is based on the mapping between the feature parameters of an object measured in the image plane of the camera and the operational space variables defining the relative pose of the object with respect to the camera or, equivalently, of the camera with respect to the object. Often, it is sufficient to derive a differential mapping in terms of velocity. As for the computation of the inverse kinematics of a manipulator, the differential problem is easier to solve because the velocity mapping is linear; moreover, the solution to the differential problem can be used to compute the pose by using numerical integration algorithms.

The set of feature parameters of an image defines a $(k \times 1)$ vector \mathbf{s} , termed *feature vector*. In the following, to simplify notation, normalized coordinates (X, Y) defined in (5.44) will be used in place of pixel coordinates (X_I, Y_I) to define the feature vector. Since only pixel coordinates can be directly measured, the normalized coordinates should be computed from pixel coordinates using the inverse of mapping (5.45), provided that the intrinsic parameters of the camera are known.

The feature vector \mathbf{s} of a point is defined as

$$\mathbf{s} = \begin{bmatrix} X \\ Y \end{bmatrix}, \tag{10.3}$$

while

$$\tilde{\mathbf{s}} = \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

denotes its representation in homogeneous coordinates.

10.3.1 Analytical Solution

Consider a camera, for example an eye-in-hand camera, and a reference frame $O_c-x_c y_c z_c$ attached to the camera; consider also a reference frame $O_o-x_o y_o z_o$ attached to the object, supposed to be rigid, and let \mathbf{T}_o^c be the homogeneous transformation matrix corresponding to the relative pose of the object with respect to the camera, defined as

$$\mathbf{T}_o^c = \begin{bmatrix} \mathbf{R}_o^c & \mathbf{o}_{c,o}^c \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (10.4)$$

with $\mathbf{o}_{c,o}^c = \mathbf{o}_o^c - \mathbf{o}_c^c$, where \mathbf{o}_c^c is the position vector of the origin of the camera frame with respect to the base frame, expressed in camera frame, \mathbf{o}_o^c is the position vector of the origin of the object frame with respect to the base frame, expressed in the camera frame, and \mathbf{R}_o^c is the rotation matrix of the object frame with respect to the camera frame (Fig. 10.5).

The problem to solve is that of computing the elements of matrix \mathbf{T}_o^c from the measurements of object feature parameters in the camera image plane. To this end, consider n points of the object and let $\mathbf{r}_{o,i}^o = \mathbf{p}_i^o - \mathbf{o}_o^o$, $i = 1, \dots, n$, denote the corresponding position vectors with respect to the object frame. These quantities are assumed to be known, for example, from a CAD model of the object. The projections of these points on the image plane have coordinates

$$\mathbf{s}_i = \begin{bmatrix} X_i \\ Y_i \end{bmatrix},$$

and define the feature vector

$$\mathbf{s} = \begin{bmatrix} \mathbf{s}_1 \\ \vdots \\ \mathbf{s}_n \end{bmatrix}. \quad (10.5)$$

The homogeneous coordinates of the points of the object with respect to the camera frame can be expressed as

$$\tilde{\mathbf{r}}_{o,i}^c = \mathbf{T}_o^c \tilde{\mathbf{r}}_{o,i}^o.$$

Therefore, using (5.44), the homogeneous coordinates of the projections of these points on the image plane are given by

$$\lambda_i \tilde{\mathbf{s}}_i = \mathbf{I} \mathbf{T}_o^c \tilde{\mathbf{r}}_{o,i}^o, \quad (10.6)$$

with $\lambda_i > 0$.

Assume that n correspondences are available, namely n equations of the form (10.6) for n points of the object, whose coordinates are known both in the object frame and in the image plane. These correspondences define a system of equations to be solved for the unknown elements of matrix \mathbf{T}_o^c .

Computing the solution is a difficult task because, depending on the type and on the number of correspondences, multiple solutions may exist. This problem, in photogrammetry, is known as PnP (Perspective- n -Point) problem. In particular, it can be shown that:

- P3P problem has four solutions, in the case of three non-collinear points.
- P4P and P5P problems each have at least two solutions, in the case of non-coplanar points, while the solution is unique in the case of at least four coplanar points and no triplets of collinear points.
- PnP problem, with $n \geq 6$ non-coplanar points, has only one solution.

The analytical solution to PnP problem is rather laborious. However, the derivation becomes simpler in some particular cases as, for example, in the case of coplanar points.

Without loss of generality, assume that the plane containing the points of the object coincides with one of the three coordinate planes of the object frame, for instance, with the plane of equation $z_o = 0$; this implies that all the points of the plane have the third coordinate equal to zero. Multiplying both sides of (10.6) by the skew-symmetric matrix $\mathbf{S}(\tilde{\mathbf{s}}_i)$, the product on the left-hand side is zero, leading to the homogeneous equation

$$\mathbf{S}(\tilde{\mathbf{s}}_i)\mathbf{H} \begin{bmatrix} r_{x,i} & r_{y,i} & 1 \end{bmatrix}^T = \mathbf{0}, \quad (10.7)$$

where $r_{x,i}$ and $r_{y,i}$ are the two non-null components of vector $\mathbf{r}_{o,i}^o$ and \mathbf{H} is the (3×3) matrix

$$\mathbf{H} = [\mathbf{r}_1 \quad \mathbf{r}_2 \quad \mathbf{o}_{c,o}^c], \quad (10.8)$$

\mathbf{r}_1 and \mathbf{r}_2 being the first and the second column of rotation matrix \mathbf{R}_o^c , respectively.

Vector equation (10.7), defined on homogeneous coordinates of points belonging to two planes, is known as *planar homography*; this denomination is used also for matrix \mathbf{H} .

Notice that Eq. (10.7) is linear with respect to \mathbf{H} and, therefore, can be rewritten in the form

$$\mathbf{A}_i(\mathbf{s}_i)\mathbf{h} = \mathbf{0},$$

where \mathbf{h} is the (9×1) column vector obtained by staking the columns of matrix \mathbf{H} , while \mathbf{A}_i is the (3×9) matrix

$$\mathbf{A}_i(\mathbf{s}_i) = [r_{x,i}\mathbf{S}(\tilde{\mathbf{s}}_i) \quad r_{y,i}\mathbf{S}(\tilde{\mathbf{s}}_i) \quad \mathbf{S}(\tilde{\mathbf{s}}_i)]. \quad (10.9)$$

Since the rank of $\mathbf{S}(\cdot)$ is at most 2, then the rank of matrix \mathbf{A}_i is also at most 2; therefore, to compute \mathbf{h} (up to a scale factor), it is necessary to consider at least 4 equations of the form (10.9) written for 4 points of the plane, leading to the system of 12 equations with 9 unknowns

$$\begin{bmatrix} \mathbf{A}_1(\mathbf{s}_1) \\ \mathbf{A}_2(\mathbf{s}_2) \\ \mathbf{A}_3(\mathbf{s}_3) \\ \mathbf{A}_4(\mathbf{s}_4) \end{bmatrix} \mathbf{h} = \mathbf{A}(\mathbf{s})\mathbf{h} = \mathbf{0}, \quad (10.10)$$

with \mathbf{s} defined in (10.5).

It can be shown that, considering a set of four points with no triplets of collinear points, matrix \mathbf{A} has rank 8 and the system of equations in (10.10) admits a non-null solution $\zeta\mathbf{h}$, defined up to a scaling factor ζ (see Problem 10.2). As a result, matrix $\zeta\mathbf{H}$ can be computed up to a scaling factor ζ . The presented derivation is general and can be applied to any kind of planar homography defined by an equation of the form (10.7).

In view of (10.8), it is

$$\begin{aligned}\mathbf{r}_1 &= \zeta\mathbf{h}_1 \\ \mathbf{r}_2 &= \zeta\mathbf{h}_2 \\ \mathbf{o}_{c,o}^c &= \zeta\mathbf{h}_3\end{aligned}$$

where \mathbf{h}_i denotes the i -th column of matrix \mathbf{H} . The absolute value of constant ζ can be computed by imposing the unit norm constraint to vectors \mathbf{r}_1 and \mathbf{r}_2 :

$$|\zeta| = \frac{1}{\|\mathbf{h}_1\|} = \frac{1}{\|\mathbf{h}_2\|},$$

while the sign of ζ can be determined by choosing the solution corresponding to the object in front of the camera. Finally, the third column \mathbf{r}_3 of matrix \mathbf{R}_o^c can be computed as

$$\mathbf{r}_3 = \mathbf{r}_1 \times \mathbf{r}_2.$$

Notice that, because of the noise affecting the measurements of the coordinates in the image plane, the results of this derivation are affected by errors that can be reduced by considering a number $n > 4$ of correspondences and computing the solution $\zeta\mathbf{h}$ to the $3n$ equations in (10.10), up to a scaling factor ζ , using least-squares techniques. This, however, does not guarantee that the resulting matrix $\mathbf{Q} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$ is a rotation matrix.

A possible solution overcoming this problem consists of computing the rotation matrix ‘closest’ to \mathbf{Q} with respect to a given norm such as the matrix which minimizes the Frobenius norm¹

$$\|\mathbf{R}_o^c - \mathbf{Q}\|_F = \left(\text{Tr} \left((\mathbf{R}_o^c - \mathbf{Q})^T (\mathbf{R}_o^c - \mathbf{Q}) \right) \right)^{1/2}, \quad (10.11)$$

with the constraint that \mathbf{R}_o^c is a rotation matrix. The problem of minimizing norm (10.11) is equivalent to that of maximizing the trace of matrix $\mathbf{R}_o^c{}^T \mathbf{Q}$. It can be shown that the solution to this problem is

$$\mathbf{R}_o^c = \mathbf{U} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \sigma \end{bmatrix} \mathbf{V}^T \quad (10.12)$$

where \mathbf{U} and \mathbf{V}^T are, respectively, the left and right orthogonal matrices of the singular value decomposition of $\mathbf{Q} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$. The choice $\sigma = \det(\mathbf{U}\mathbf{V}^T)$ ensures that the determinant of \mathbf{R}_o^c is equal to one (see Problem 10.3).

¹ The Frobenius norm is defined in Sect. A.4.

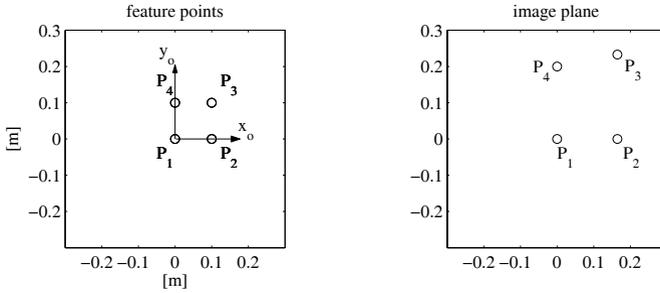


Fig. 10.6. Planar object; *left*: object frame and feature points; *right*: feature points projections on the normalized image plane of a camera in the pose of Example 10.1

Example 10.1

Consider a planar object characterized by four feature points shown in Fig. 10.6, where the object frame is represented. The feature points P_1, P_2, P_3, P_4 are the vertices of a square with side length $l = 0.1$ m. In Fig. 10.6, the images of the projections of the four points of the object on the normalized image plane of the camera are shown as well, under the assumption that the relative pose of the object frame with respect to the camera frame is characterized by rotation matrix

$$\mathbf{R}_o^c = \mathbf{R}_z(0)\mathbf{R}_y(\pi/4)\mathbf{R}_x(0) = \begin{bmatrix} 0.7071 & 0 & 0.7071 \\ 0 & 1 & 0 \\ -0.7071 & 0 & 0.7071 \end{bmatrix}$$

and position vector $\mathbf{o}_{c,o}^c = [0 \ 0 \ 0.5]^T$ m. The normalized coordinates of the four points of the object can be computed from the position vectors in the object frame $\mathbf{r}_{o,1}^o = [0 \ 0 \ 0]^T$ m, $\mathbf{r}_{o,2}^o = [0.1 \ 0 \ 0]^T$ m, $\mathbf{r}_{o,3}^o = [0.1 \ 0.1 \ 0]^T$ m, $\mathbf{r}_{o,4}^o = [0 \ 0.1 \ 0]^T$ m, using (10.6), which gives

$$\mathbf{s}_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \mathbf{s}_2 = \begin{bmatrix} 0.1647 \\ 0 \end{bmatrix} \quad \mathbf{s}_3 = \begin{bmatrix} 0.1647 \\ 0.2329 \end{bmatrix} \quad \mathbf{s}_4 = \begin{bmatrix} 0 \\ 0.2 \end{bmatrix}.$$

To solve the inverse problem, namely that of computing matrix \mathbf{T}_o^c from the coordinates of the four points both in the image plane and in the object frame, it is necessary to build matrix $\mathbf{A}(\mathbf{s})$ from four matrices $\mathbf{A}_i(\mathbf{s}_i)$ defined in (10.9). It is easy to verify that matrix $\mathbf{A}(\mathbf{s})$ has rank 8; moreover, a non-null solution to the system of equations in (10.10) can be computed using the singular value decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

and coincides with the last column of matrix \mathbf{V} , namely, with the right eigenvector corresponding to the null singular value of \mathbf{A} . From this computation, matrix

$$\zeta\mathbf{H} = \begin{bmatrix} -0.4714 & 0 & 0 \\ 0 & -0.6667 & 0 \\ 0.4714 & 0 & -0.3333 \end{bmatrix}$$

can be obtained. Normalization of the first column yields $|\zeta| = 1.5$. It is possible to verify that, with the choice $\zeta = -1.5$, the exact solution for $\mathbf{o}_{c,o}^c$ is obtained; moreover, matrix \mathbf{Q} coincides numerically with the sought rotation matrix \mathbf{R}_o^c , without using any kind of approximate solution. This result was expected, due to the absence of measurement noise affecting image plane coordinates \mathbf{s}_i .

The above derivation is a particular case of a method, termed *direct linear transformation*, which is aimed at computing the elements of matrix \mathbf{T}_o^c by solving a system of linear equations obtained using n correspondences relative to points in generic configuration. In detail, from equalities

$$\mathbf{S}(\tilde{\mathbf{s}}_i) [\mathbf{R}_o^c \quad \mathbf{o}_{c,o}^c] \tilde{\mathbf{r}}_{o,i}^o = \mathbf{0}, \quad (10.13)$$

which coincide with (10.7) in the case that points $\mathbf{r}_{o,i}^o$ belong to plane $z_0 = 0$, two independent linear equations with 12 unknowns are obtained, taking into account that matrix $\mathbf{S}(\cdot)$ is, at most, of rank 2. Therefore, n correspondences produce $2n$ equations.

It can be shown that, considering a set of 6 points not all coplanar, the matrix of coefficients of the corresponding system of 12 equations with 12 unknowns has rank 11; therefore, the solution is defined up to a scaling factor. Once this solution has been computed, the elements of rotation matrix \mathbf{R}_o^c and of vector $\mathbf{o}_{c,o}^c$ can be obtained using a derivation similar to that presented above. Notice that, in practical applications, due to the presence of noise, the system of equations has rank 12 and admits only the null solution. In this case, it is necessary to consider a number $n > 6$ of correspondences and to compute the solution of the resulting system of equations, defined up to a scaling factor, using least-squares techniques.

In conclusion, the presented method permits the computation of matrix \mathbf{T}_o^c , characterizing the relative pose of the object frame with respect to the camera frame, from the projections of n points of the object on the camera image plane. To this end, it is necessary to know the position of these points with respect to the object frame, and thus the object geometry, besides the camera *intrinsic parameters*. The latter are required for computing normalized coordinates \mathbf{s}_i from pixel coordinates.

Notice that, if it is required to compute the object pose with respect to the base frame (as usually happens in the case of eye-to-hand cameras) or to the end-effector frame (as usually happens in the case of eye-in-hand cameras), then it is also necessary to know the camera *extrinsic parameters*. In fact, in the first case, it is

$$\mathbf{T}_o^b = \mathbf{T}_c^b \mathbf{T}_o^c, \quad (10.14)$$

where the elements of matrix \mathbf{T}_c^b represent the extrinsic parameters of an eye-to-hand camera; on the other hand, in the case of eye-in-hand camera, it is

$$\mathbf{T}_o^e = \mathbf{T}_c^e \mathbf{T}_o^c, \quad (10.15)$$

where the extrinsic parameters matrix \mathbf{T}_c^e characterize the pose of the camera with respect to the end-effector frame.

10.3.2 Interaction Matrix

If the object is in motion with respect to the camera, the feature vector \mathbf{s} is, in general, time-varying. Therefore, it is possible to define a $(k \times 1)$ velocity vector in the image plane $\dot{\mathbf{s}}$.

The motion of the object with respect to the camera is characterized by the relative velocity

$$\mathbf{v}_{c,o}^c = \begin{bmatrix} \dot{\mathbf{o}}_{c,o}^c \\ \mathbf{R}_c^T(\boldsymbol{\omega}_o - \boldsymbol{\omega}_c) \end{bmatrix}, \quad (10.16)$$

where $\dot{\mathbf{o}}_{c,o}^c$ is the time derivative of vector $\mathbf{o}_{c,o}^c = \mathbf{R}_c^T(\mathbf{o}_o - \mathbf{o}_c)$, representing the relative position of the origin of the object frame with respect to the origin of the camera frame, while $\boldsymbol{\omega}_o$ and $\boldsymbol{\omega}_c$ are the angular velocities of the object frame and camera frame, respectively.

The equation relating $\dot{\mathbf{s}}$ to $\mathbf{v}_{c,o}^c$ is

$$\dot{\mathbf{s}} = \mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) \mathbf{v}_{c,o}^c, \quad (10.17)$$

where \mathbf{J}_s is a $(k \times 6)$ matrix termed *image Jacobian*. This equation is linear but \mathbf{J}_s depends, in general, on the current value of the feature vector \mathbf{s} and on the relative pose of the object with respect to the camera \mathbf{T}_o^c .

It is useful to consider also the mapping between the image plane velocity $\dot{\mathbf{s}}$, the absolute velocity of the camera frame

$$\mathbf{v}_c^c = \begin{bmatrix} \mathbf{R}_c^T \dot{\mathbf{o}}_c \\ \mathbf{R}_c^T \boldsymbol{\omega}_c \end{bmatrix}$$

and the absolute velocity of the object frame

$$\mathbf{v}_o^c = \begin{bmatrix} \mathbf{R}_c^T \dot{\mathbf{o}}_o \\ \mathbf{R}_c^T \boldsymbol{\omega}_o \end{bmatrix}.$$

To this end, vector $\dot{\mathbf{o}}_{c,o}^c$ can be expressed as

$$\dot{\mathbf{o}}_{c,o}^c = \mathbf{R}_c^T(\dot{\mathbf{o}}_o - \dot{\mathbf{o}}_c) + \mathbf{S}(\mathbf{o}_{c,o}^c) \mathbf{R}_c^T \boldsymbol{\omega}_c,$$

which allows equality (10.16) to be rewritten in the compact form

$$\mathbf{v}_{c,o}^c = \mathbf{v}_o^c + \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) \mathbf{v}_c^c, \quad (10.18)$$

with

$$\mathbf{\Gamma}(\cdot) = \begin{bmatrix} -\mathbf{I} & \mathbf{S}(\cdot) \\ \mathbf{O} & -\mathbf{I} \end{bmatrix}.$$

Therefore, Eq. (10.17) can be rewritten in the form

$$\dot{\mathbf{s}} = \mathbf{J}_s \mathbf{v}_o^c + \mathbf{L}_s \mathbf{v}_c^c, \quad (10.19)$$

where the $(k \times 6)$ matrix

$$\mathbf{L}_s = \mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) \quad (10.20)$$

is termed *interaction matrix*. This matrix, in view of (10.19), defines the linear mapping between the absolute velocity of the camera \mathbf{v}_c^c and the corresponding image plane velocity $\dot{\mathbf{s}}$, in the case that the object is fixed with respect to the base frame ($\mathbf{v}_o^c = \mathbf{0}$).

The analytical expression of the interaction matrix is, in general, simpler than that of the image Jacobian. The latter can be computed from the interaction matrix using the equation

$$\mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) = \mathbf{L}_s \mathbf{\Gamma}(-\mathbf{o}_{c,o}^c), \quad (10.21)$$

obtained from (10.20), with $\mathbf{\Gamma}^{-1}(\mathbf{o}_{c,o}^c) = \mathbf{\Gamma}(-\mathbf{o}_{c,o}^c)$. In the following, examples of computation of interaction matrix and image Jacobian for some of the most common cases in applications are provided.

Interaction matrix of a point

Consider a point P of the object characterized, with respect to the camera frame, by the vector of coordinates

$$\mathbf{r}_c^c = \mathbf{R}_c^T(\mathbf{p} - \mathbf{o}_c), \quad (10.22)$$

where \mathbf{p} is the position of point P with respect to the base frame. Choose vector \mathbf{s} of normalized coordinates (10.3) as the feature vector of the point. In view of (5.44), the following expression holds:

$$\mathbf{s} = \mathbf{s}(\mathbf{r}_c^c), \quad (10.23)$$

with

$$\mathbf{s}(\mathbf{r}_c^c) = \frac{1}{z_c} \begin{bmatrix} x_c \\ y_c \end{bmatrix} = \begin{bmatrix} X \\ Y \end{bmatrix}, \quad (10.24)$$

and $\mathbf{r}_c^c = [x_c \ y_c \ z_c]^T$. Computing the time derivative of (10.23) and using (10.24) yields

$$\dot{\mathbf{s}} = \frac{\partial \mathbf{s}(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} \dot{\mathbf{r}}_c^c, \quad (10.25)$$

with

$$\frac{\partial \mathbf{s}(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -x_c/z_c \\ 0 & 1 & -y_c/z_c \end{bmatrix} = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -X \\ 0 & 1 & -Y \end{bmatrix}.$$

To compute the interaction matrix, vector $\dot{\mathbf{r}}_c^c$ can be obtained from the time derivative of (10.22) under the assumption that \mathbf{p} is constant:

$$\dot{\mathbf{r}}_c^c = -\mathbf{R}_c^T \dot{\mathbf{o}}_c + \mathbf{S}(\mathbf{r}_c^c) \mathbf{R}_c^T \boldsymbol{\omega}_c = [-\mathbf{I} \quad \mathbf{S}(\mathbf{r}_c^c)] \mathbf{v}_c^c. \quad (10.26)$$

Combining Eqs. (10.25), (10.26), the following expression of interaction matrix of a point can be found:

$$\mathbf{L}_s(\mathbf{s}, z_c) = \begin{bmatrix} -\frac{1}{z_c} & 0 & \frac{X}{z_c} & XY & -(1+X^2) & Y \\ 0 & -\frac{1}{z_c} & \frac{Y}{z_c} & 1+Y^2 & -XY & -X \end{bmatrix}, \quad (10.27)$$

revealing that this matrix depends on the components of vector \mathbf{s} and the sole component z_c of vector \mathbf{r}_c^c .

The image Jacobian of a point can be computed using (10.21), (10.27) and has the expression

$$\mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -X & -r_{o,y}^c X & r_{o,z}^c + r_{o,x}^c X & -r_{o,y}^c \\ 0 & 1 & -Y & -(r_{o,z}^c + r_{o,y}^c Y) & r_{o,x}^c Y & r_{o,x}^c \end{bmatrix},$$

where $r_{o,x}^c, r_{o,y}^c, r_{o,z}^c$ are the components of vector $\mathbf{r}_o^c = \mathbf{r}_c^c - \mathbf{o}_{c,o}^c = \mathbf{R}_o^c \mathbf{r}_o^o$, \mathbf{r}_o^o being the constant vector expressing the position of point P with respect to the object frame.

Interaction matrix of a set of points

The interaction matrix of a set of n points of the object P_1, \dots, P_n can be built by considering the $(2n \times 1)$ feature vector (10.5). If $\mathbf{L}_{s_i}(\mathbf{s}_i, z_{c,i})$ denotes the interaction matrix corresponding to point P_i , then the interaction matrix of the set of points will be the $(2n \times 6)$ matrix

$$\mathbf{L}_s(\mathbf{s}, \mathbf{z}_c) = \begin{bmatrix} \mathbf{L}_{s_1}(\mathbf{s}_1, z_{c,1}) \\ \vdots \\ \mathbf{L}_{s_n}(\mathbf{s}_n, z_{c,n}) \end{bmatrix},$$

with $\mathbf{z}_c = [z_{c,1} \dots z_{c,n}]^T$.

The image Jacobian of a set of points can be easily computed from the interaction matrix, using (10.21).

Interaction matrix of a line segment

A line segment is the part of the line connecting two points P_1 and P_2 . The projection on the image plane is still a line segment that can be characterized in terms of the middle point coordinates \bar{x}, \bar{y} , the length L , and the angle α

formed by the line with respect to axis X . Therefore, the feature vector can be defined as

$$\mathbf{s} = \begin{bmatrix} \bar{x} \\ \bar{y} \\ L \\ \alpha \end{bmatrix} = \begin{bmatrix} (X_1 + X_2)/2 \\ (Y_1 + Y_2)/2 \\ \sqrt{\Delta X^2 + \Delta Y^2} \\ \tan^{-1}(\Delta Y/\Delta X) \end{bmatrix} = \mathbf{s}(\mathbf{s}_1, \mathbf{s}_2) \quad (10.28)$$

with $\Delta X = X_2 - X_1$, $\Delta Y = Y_2 - Y_1$ and $\mathbf{s}_i = [X_i \ Y_i]^T$, $i = 1, 2$. Computing the time derivative of this equation yields

$$\begin{aligned} \dot{\mathbf{s}} &= \frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \dot{\mathbf{s}}_1 + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \dot{\mathbf{s}}_2 \\ &= \left(\frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \mathbf{L}_{s_1}(\mathbf{s}_1, z_{c,1}) + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \mathbf{L}_{s_2}(\mathbf{s}_2, z_{c,2}) \right) \mathbf{v}_c^c, \end{aligned}$$

where \mathbf{L}_{s_i} is the interaction matrix of point P_i , under the assumption that the line segment is fixed with respect to the base frame. Therefore, the interaction matrix of a line segment is

$$\mathbf{L}_s(\mathbf{s}, z_c) = \frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \mathbf{L}_{s_1}(\mathbf{s}_1, z_{c,1}) + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \mathbf{L}_{s_2}(\mathbf{s}_2, z_{c,2}),$$

with

$$\frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \\ -\Delta X/L & -\Delta Y/L \\ \Delta Y/L^2 & -\Delta X/L^2 \end{bmatrix} \quad \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \\ \Delta X/L & \Delta Y/L \\ -\Delta Y/L^2 & \Delta X/L^2 \end{bmatrix}.$$

Notice that vectors \mathbf{s}_1 and \mathbf{s}_2 can be computed as functions of parameters \bar{x} , \bar{y} , L , α , using (10.28). Therefore, the interaction matrix can be expressed as a function of the feature vector $\mathbf{s} = [\bar{x} \ \bar{y} \ L \ \alpha]^T$, besides the components $z_{c,1}$ and $z_{c,2}$ of the endpoints P_1 and P_2 of the line segment.

The image Jacobian of a line segment can be easily computed from the interaction matrix using (10.21).

10.3.3 Algorithmic Solution

The interaction matrix \mathbf{L}_s is, in general, a matrix of dimension $(k \times m)$, where k is equal to the number of feature parameters of the image and m is the dimension of velocity vector \mathbf{v}_c^c . Usually $m = 6$, but it may happen that $m < 6$, when the relative motion of the object with respect to the camera is constrained.

The image Jacobian \mathbf{J}_s is also of dimension $(k \times m)$, being related to \mathbf{L}_s by mapping (10.21). Since this mapping is invertible, the rank of \mathbf{J}_s coincides with that of \mathbf{L}_s .

In the case that \mathbf{L}_s is full rank, by using (10.17), it is possible to compute $\mathbf{v}_{c,o}^c$ from $\dot{\mathbf{s}}$.

In particular, if $k = m$, the velocity $\mathbf{v}_{c,o}^c$ can be obtained using the expression

$$\mathbf{v}_{c,o}^c = \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) \mathbf{L}_s^{-1} \dot{\mathbf{s}}, \quad (10.29)$$

which requires the computation of the inverse of the interaction matrix.

In the case $k > m$, the interaction matrix has more rows than columns and Eq. (10.17) can be solved using a least-squares technique, whose solution can be written in the form

$$\mathbf{v}_{c,o}^c = \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) (\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T \dot{\mathbf{s}}, \quad (10.30)$$

where $(\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T$ is the left pseudo-inverse of \mathbf{L}_s . This situation is rather frequent in applications, because it permits using interaction matrices with good condition numbers.

Finally, in the case $k < m$, the interaction matrix has more columns than rows and Eq. (10.17) admits infinite solutions. This implies that the number of parameters of the observed image is not sufficient to determine uniquely the relative motion of the object with respect to the camera. Hence, there exist relative motions of the object with respect to the camera (or vice versa) that do not produce variations of the image feature parameters. The velocities associated with these relative motions belong to the null subspace of \mathbf{J}_s , which has the same dimension of the null subspace of \mathbf{L}_s . If the problem is that of computing uniquely the relative pose of the object with respect to the camera from feature parameters in the image plane, this case has no interest.

Example 10.2

The interaction matrix of a point P is a matrix with more columns than rows of dimension (2×6) and rank 2; therefore, the null subspace has dimension 4. It can be seen immediately that this subspace contains the velocity of the camera translational motion along the visual ray projecting point P on the image plane, proportional to vector

$$\mathbf{v}_1 = [X \ Y \ 1 \ 0 \ 0 \ 0]^T,$$

as well as the velocity of the camera rotational motion about this visual ray, proportional to vector

$$\mathbf{v}_2 = [0 \ 0 \ 0 \ X \ Y \ 1]^T.$$

Vectors \mathbf{v}_1 and \mathbf{v}_2 are independent and belong to a base of the null subspace of \mathbf{L}_s . The remaining base vectors are not easy to find geometrically, but can be easily computed analytically.

The pose estimation problem may be cast in a form analogous to that of inverse kinematics algorithms for robot manipulators. To this end, it is

necessary to represent the relative pose of the object with respect to the camera using a minimum number of coordinates, in terms of the $(m \times 1)$ vector

$$\mathbf{x}_{c,o} = \begin{bmatrix} \mathbf{o}_{c,o}^c \\ \phi_{c,o} \end{bmatrix}, \quad (10.31)$$

where $\mathbf{o}_{c,o}^c$ characterizes the position of the origin of the object frame with respect to the camera frame and $\phi_{c,o}$ characterizes the relative orientation. If Euler angles are used to represent orientation, then $\phi_{c,o}$ is the vector of the angles extracted from rotation matrix \mathbf{R}_o^c and the mapping between $\mathbf{v}_{c,o}^c$ and $\dot{\mathbf{x}}_{c,o}$ is expressed by

$$\mathbf{v}_{c,o}^c = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{T}(\phi_{c,o}) \end{bmatrix} \dot{\mathbf{x}}_{c,o} = \mathbf{T}_A(\phi_{c,o}) \dot{\mathbf{x}}_{c,o}. \quad (10.32)$$

Example 10.3

Consider a camera mounted on the end-effector of the SCARA manipulator of Fig. 2.36. Choose the camera frame parallel to the end-effector frame, with axis z_c pointing downward. Assume that the camera observes a fixed planar object, parallel to the image plane, and that axis z_o of the object frame is parallel to axis z_c and points downward.

The geometry of the problem suggests that the relative position of the object with respect to the camera can be represented by a vector $\mathbf{o}_{c,o}^c$, whereas the relative orientation can be defined by the angle α between the object frame and the camera frame about axis z_c . Therefore, $m = 4$ and

$$\mathbf{x}_{c,o} = \begin{bmatrix} \mathbf{o}_{c,o}^c \\ \alpha \end{bmatrix}. \quad (10.33)$$

Moreover, the time derivative $\dot{\alpha}$ coincides with the component of $\boldsymbol{\omega}_{c,o}^c$ along z_c , and this is the sole non-null component of the angular velocity of the relative motion of the object frame with respect to the camera frame. Hence, in (10.32), $\mathbf{T}_A(\phi_{c,o})$ is the (4×4) identity matrix.

Equation (10.17) can be rewritten in the form

$$\dot{\mathbf{s}} = \mathbf{J}_{A_s}(\mathbf{s}, \mathbf{x}_{c,o}) \dot{\mathbf{x}}_{c,o}, \quad (10.34)$$

where the matrix

$$\mathbf{J}_{A_s}(\mathbf{s}, \mathbf{x}_{c,o}) = \mathbf{L}_s \boldsymbol{\Gamma}(-\mathbf{o}_{c,o}^c) \mathbf{T}_A(\phi_{c,o}) \quad (10.35)$$

has a meaning analogous to that of the analytical Jacobian of a manipulator.

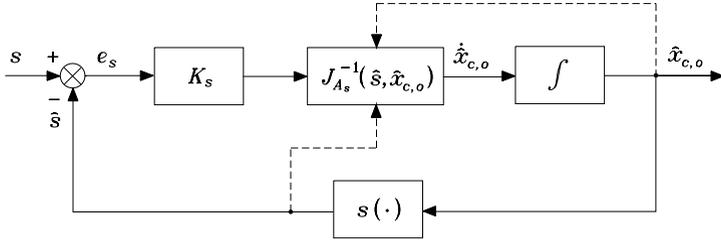


Fig. 10.7. Pose estimation algorithm based on the inverse of image Jacobian

Equation (10.34) is the starting point of a numeric integration algorithm for the computation of $\mathbf{x}_{c,o}$, similar to inverse kinematics algorithms. Let $\hat{\mathbf{x}}_{c,o}$ denote the current estimate of vector $\mathbf{x}_{c,o}$ and let

$$\hat{\mathbf{s}} = \mathbf{s}(\hat{\mathbf{x}}_{c,o})$$

be the corresponding vector of image feature parameters computed from the pose specified by $\hat{\mathbf{x}}_{c,o}$; the objective of this algorithm is the minimization of the error

$$\mathbf{e}_s = \mathbf{s} - \hat{\mathbf{s}}. \quad (10.36)$$

Notice that, for the purpose of numerical integration, vector \mathbf{s} is constant while the current estimate $\hat{\mathbf{s}}$ depends on the current integration time. Therefore, computing the time derivative of (10.36) yields

$$\dot{\mathbf{e}}_s = -\dot{\hat{\mathbf{s}}} = -\mathbf{J}_{A_s}(\hat{\mathbf{s}}, \hat{\mathbf{x}}_{c,o})\dot{\hat{\mathbf{x}}}_{c,o}. \quad (10.37)$$

Assumed that matrix \mathbf{J}_{A_s} is square and nonsingular, the choice

$$\dot{\hat{\mathbf{x}}}_{c,o} = \mathbf{J}_{A_s}^{-1}(\hat{\mathbf{s}}, \hat{\mathbf{x}}_{c,o})\mathbf{K}_s\mathbf{e}_s \quad (10.38)$$

leads to the equivalent linear system

$$\dot{\mathbf{e}}_s + \mathbf{K}_s\mathbf{e}_s = \mathbf{0}. \quad (10.39)$$

Therefore, if \mathbf{K}_s is a positive definite matrix (usually diagonal), system (10.39) is asymptotically stable and the error tends to zero with a convergence speed that depends on the eigenvalues of matrix \mathbf{K}_s . The convergence to zero of error \mathbf{e}_s ensures the asymptotic convergence of the estimate $\hat{\mathbf{x}}_{c,o}$ to the true value $\mathbf{x}_{c,o}$.

The block scheme of the pose estimation algorithm is shown in Fig. 10.7, where $\mathbf{s}(\cdot)$ denotes the function computing the feature vector of the ‘virtual’ image corresponding to the current estimate $\hat{\mathbf{x}}_{c,o}$ of the object pose with respect to the camera. This algorithm can be used as an alternative to the analytical methods for pose estimation illustrated in Sect. 10.3.1. Obviously, the convergence properties depend on the choice of the image feature parameters and on the initial value of estimate $\hat{\mathbf{x}}_{c,o}(0)$, which may produce instability problems related to the singularities of matrix \mathbf{J}_{A_s} .

Notice that, in view of (10.35), the singularities of matrix \mathbf{J}_{A_s} are both the representation singularities of the orientation and those of the interaction matrix. The most critical singularities are those of the interaction matrix, since they depend on the choice of the image feature parameters.

To separate the effects of the two types of singularities, it is convenient to compute (10.38) in two steps, evaluating first

$$\widehat{\mathbf{v}}_{c,o}^c = \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) \mathbf{L}_s^{-1} \mathbf{K}_s \mathbf{e}_s, \quad (10.40)$$

and then

$$\dot{\widehat{\mathbf{x}}}_{c,o} = \mathbf{T}_A^{-1}(\phi_{c,o}) \widehat{\mathbf{v}}_{c,o}^c. \quad (10.41)$$

Assumed to work far from representation singularities, the problem of singularities of \mathbf{L}_s can be overcome by using a number k of feature parameters greater than the minimum required m . This choice also allows a reduction of the effects of measurement noise. The resulting estimation algorithm requires the use of the left pseudo-inverse of \mathbf{L}_s in place of the inverse, namely

$$\widehat{\mathbf{v}}_{c,o}^c = \mathbf{\Gamma}(\mathbf{o}_{c,o}^c) (\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T \mathbf{K}_s \mathbf{e}_s \quad (10.42)$$

in place of (10.40). The convergence of error (10.36) can be shown using the direct Lyapunov method based on the positive definite function ²

$$V(\mathbf{e}_s) = \frac{1}{2} \mathbf{e}_s^T \mathbf{K}_s \mathbf{e}_s > 0 \quad \forall \mathbf{e}_s \neq \mathbf{0}.$$

Computing the time derivative of this function, and using (10.37), (10.35), (10.41), (10.42), yields

$$\dot{V} = -\mathbf{e}_s^T \mathbf{K}_s \mathbf{L}_s (\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T \mathbf{K}_s \mathbf{e}_s,$$

which is negative semi-definite because $\mathcal{N}(\mathbf{L}_s^T) \neq \emptyset$, \mathbf{L}_s^T being a matrix with more columns than rows. Therefore, the system is stable but not asymptotically stable. This implies that the error is bounded, but in some cases the algorithm can get stuck with $\mathbf{e}_s \neq \mathbf{0}$ and $\mathbf{K}_s \mathbf{e}_s \in \mathcal{N}(\mathbf{L}_s^T)$.

Notice that the pose estimation methods based on inverse Jacobian are as efficient in terms of accuracy, speed of convergence and computational load, as the initial estimate $\widehat{\mathbf{x}}_{c,o}(0)$ is close to the true value $\mathbf{x}_{c,o}$. Therefore, these methods are mainly adopted for real-time ‘visual tracking’ applications, where the estimate on an image taken at time \bar{t} is computed assuming as initial value the estimate computed on the image taken at time $\bar{t}-T$, T being the sampling time of the image (multiple of the sampling time of the numerical integration algorithm).

² See Sect. C.3 for the illustration of the direct Lyapunov method.

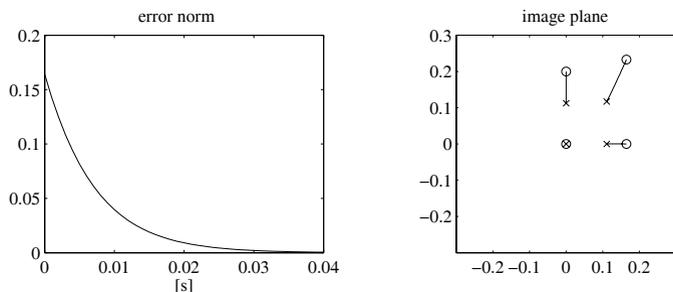


Fig. 10.8. Time history of the norm of the estimation error and corresponding paths of the feature points projections on image plane

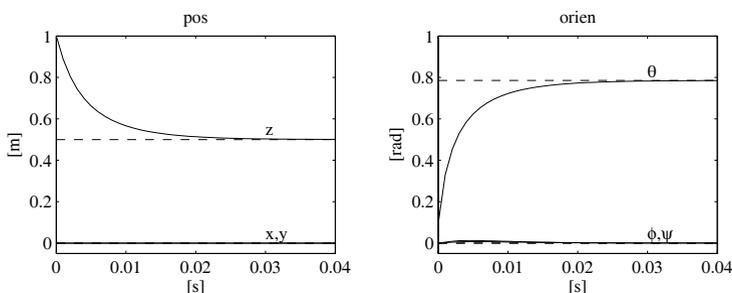


Fig. 10.9. Time history of camera pose estimate

Example 10.4

Consider the object of Example 10.1. Using the algorithmic solution, it is desired to compute the relative pose of the object frame, with respect to the camera frame, from the image plane coordinates of the projections of the four feature points of the object. The same numerical values of Example 10.1 are used.

Since the image Jacobian has dimension (6×8) , it is necessary to use the algorithm based on the pseudo-inverse of the interaction matrix. This algorithm was simulated on a computer by adopting the Euler numerical integration scheme with an integration time $\Delta t = 1$ ms, matrix gain $\mathbf{K}_s = 160\mathbf{I}_8$, and initial estimate $\hat{\mathbf{x}}_{c,o} = [0 \ 0 \ 1 \ 0 \ \pi/32 \ 0]^T$.

The results in Fig. 10.8 show that the norm of the estimation error of the feature parameters e_s tends to zero asymptotically with convergence of exponential type; moreover, due to the fact that matrix gain \mathbf{K}_s was chosen diagonal with equal elements, the paths of the projections of the feature points on the image plane (between the initial positions marked with crosses and the final positions marked with circles) are line segments.

The corresponding time histories of the components of vector $\hat{\mathbf{x}}_{c,o}$ for position and orientation are reported in Fig. 10.9, together with the time histories of the components of the true value $\mathbf{x}_{c,o} = [0 \ 0 \ 0.5 \ 0 \ \pi/4 \ 0]^T$ (represented with

dashed lines). It can be verified that, with the chosen value of \mathbf{K}_s , the algorithm converges to the true value in about 0.03 s, corresponding to 30 iterations.

The time histories of Fig. 10.9 can be interpreted as the position and orientation trajectories of a ‘virtual’ camera in motion between the initial pose $\hat{\mathbf{x}}_{c,o}(0)$ and the final pose $\mathbf{x}_{c,o}$.

Notice that, for the purpose of pose estimation, the illustrated algorithm may converge also in the case that only three feature points are used. In fact, in this case, \mathbf{J}_{A_s} is a square (6×6) matrix and the convergence is ensured provided that this matrix is nonsingular. However, since P3P problem has four solutions, the algorithm may converge to a solution different from that sought, unless, as for visual tracking applications, the initial estimate $\hat{\mathbf{x}}_{c,o}(0)$ is close enough to the true value $\mathbf{x}_{c,o}$.

10.4 Stereo Vision

The bidimensional image provided by a camera does not give any explicit information on depth, namely, the distance of the observed object from the camera. This information can be recovered in an indirect way from the geometric model of the object, assumed to be known.

On the other hand, the depth of a point can be directly computed in the case that two images of the same scene are available, taken from two different points of view. The two images can be obtained using two cameras, or sequentially, using a moving camera. These cases are referred to as *stereo vision*.

In the framework of stereo vision, two fundamental problems can be devised. The first is the *correspondence problem*, which consists of the identification of the points of the two images that are projections of the same point of the scene. These points are termed *conjugate* or *corresponding*. This problem is not easy to solve, and the solution is based on the existence of geometric constraints between two images of the same point, besides the fact that some details of the scene appear to be similar in the two images.

The second problem, illustrated below in some fundamental aspects, is that of *3D reconstruction* which, in general, consists of the computation of the relative pose of the cameras (calibrated and not) and thus, starting from this pose, the position in the 3D space of the points of the observed object.

10.4.1 Epipolar Geometry

Assume that two cameras are available, with respective reference frames, denoted as 1 and 2. Moreover, let $\mathbf{o}_{1,2}$ denote the position vector and \mathbf{R}_2^1 the rotation matrix of Frame 2 with respect to Frame 1 and let \mathbf{T}_2^1 be the corresponding homogeneous transformation matrix. The coordinates of a point P expressed in the two frames are related by equation

$$\mathbf{p}^1 = \mathbf{o}_{1,2}^1 + \mathbf{R}_2^1 \mathbf{p}^2. \quad (10.43)$$

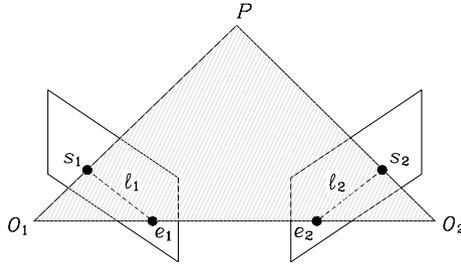


Fig. 10.10. Epipolar geometry

Let \mathbf{s}_1 and \mathbf{s}_2 be the coordinates of the projections of P on the image planes of the cameras; in view of (5.44) it is

$$\lambda_i \tilde{\mathbf{s}}_i = \mathbf{I} \tilde{\mathbf{p}}^i = \mathbf{p}^i, \quad i = 1, 2. \tag{10.44}$$

Substituting (10.44) into (10.43) yields

$$\lambda_1 \tilde{\mathbf{s}}_1 = \mathbf{o}_{1,2}^1 + \lambda_2 \mathbf{R}_2^1 \tilde{\mathbf{s}}_2. \tag{10.45}$$

Premultiplying both sides of (10.45) by $\mathbf{S}(\mathbf{o}_{1,2}^1)$ gives

$$\lambda_1 \mathbf{S}(\mathbf{o}_{1,2}^1) \tilde{\mathbf{s}}_1 = \lambda_2 \mathbf{S}(\mathbf{o}_{1,2}^1) \mathbf{R}_2^1 \tilde{\mathbf{s}}_2.$$

Hence, premultiplying both sides of the above equation by $\tilde{\mathbf{s}}_1^T$, the following equality is obtained

$$\lambda_2 \tilde{\mathbf{s}}_1^T \mathbf{S}(\mathbf{o}_{1,2}^1) \mathbf{R}_2^1 \tilde{\mathbf{s}}_2 = 0,$$

which has to be satisfied for any value of scalar λ_2 . Therefore, this equality is equivalent to the so-called *epipolar constraint* equation

$$\tilde{\mathbf{s}}_1^T \mathbf{E} \tilde{\mathbf{s}}_2 = 0, \tag{10.46}$$

where $\mathbf{E} = \mathbf{S}(\mathbf{o}_{1,2}^1) \mathbf{R}_2^1$ is a (3×3) matrix known as *essential matrix*. Equation (10.46) expresses in analytical form the geometric constraint existing between the projections of the same point on the image planes of the two cameras.

The geometric interpretation of the epipolar constraint can be derived from Fig. 10.10, where the projections of a point P on the image planes of the two cameras are reported as well as the respective optical centers O_1 and O_2 . Notice that points O_1 , O_2 and P are the vertices of a triangle whose sides O_1P and O_2P belong to the visual rays projecting point P into the points of coordinates \mathbf{s}_1 and \mathbf{s}_2 of the image plane, respectively. These rays lay along the directions of vectors $\tilde{\mathbf{s}}_1$ and $\mathbf{R}_2^1 \tilde{\mathbf{s}}_2$ respectively, expressed with respect to Frame 1. Line segment O_1O_2 , termed *base line*, is represented by vector $\mathbf{o}_{1,2}^1$. The epipolar constraint (10.46) corresponds to imposing that vectors $\tilde{\mathbf{s}}_1$, $\mathbf{R}_2^1 \tilde{\mathbf{s}}_2$

and $\mathbf{o}_{1,2}^1$ are coplanar. The plane containing these vectors is termed *epipolar plane*.

Notice that line segment O_1O_2 belongs to the visual ray projecting point O_2 on the image plane of Camera 1 and, at the same time, to the ray projecting point O_1 on the image plane of Camera 2. These projections, of coordinates \mathbf{e}_1 and \mathbf{e}_2 , respectively, are termed *epipoles*. Line ℓ_1 , passing through the points of coordinates \mathbf{s}_1 and \mathbf{e}_1 , and line ℓ_2 , passing through the points of coordinates \mathbf{s}_2 and \mathbf{e}_2 , are termed *epipolar lines*. The epipolar lines can also be obtained as the intersection of the epipolar plane with the image planes of the two cameras. Notice that, by varying point P , the epipolar plane describes a set of planes about the base line and the epipoles do not change.

For the purpose of computing the correspondences, the epipolar constraint can be exploited to reduce the complexity of the problem to find conjugate points. In fact, if \mathbf{s}_1 is the image of a point of the visual ray passing through O_1 and the point of coordinates \mathbf{s}_1 , the corresponding conjugate point of the image plane of Camera 2 must necessarily belong to the epipolar line ℓ_2 , which is known because the epipolar plane is uniquely defined by O_1 , O_2 and by the point of coordinates \mathbf{s}_1 . Finding correspondences, therefore, reduces to searching for a point along the epipolar line and not on the whole image plane.

In the framework of 3D reconstruction, different scenarios may arise, depending on the type of information which is available a priori.

10.4.2 Triangulation

In the case that both intrinsic and extrinsic parameters of the two cameras are known, the reconstruction problem consists of computing the position in the scene of the points projected on the two image planes using a geometric method known as *triangulation*. This method allows the computation of coordinates $\mathbf{p} = [p_x \ p_y \ p_z]^T$ of a point P with respect to the base frame, starting from normalized coordinates $\mathbf{s}_1 = [X_1 \ Y_1]^T$ and $\mathbf{s}_2 = [X_2 \ Y_2]^T$ of the projections of P on the image planes of the two cameras. Assume that, for simplicity, the base frame coincides with Frame 1, then $\mathbf{p}^1 = \mathbf{p}$, $\mathbf{o}_{1,2}^1 = \mathbf{o}$ and $\mathbf{R}_2^1 = \mathbf{R}$.

From (10.44), (10.45) the following equalities can be derived:

$$\mathbf{p} = \lambda_1 \tilde{\mathbf{s}}_1 \quad (10.47)$$

$$\mathbf{p} = \mathbf{o} + \lambda_2 \mathbf{R} \tilde{\mathbf{s}}_2, \quad (10.48)$$

where the first equality is the parametric equation of the visual ray passing through O_1 and the point of coordinates \mathbf{s}_1 , while the second represents the parametric equation of the visual ray passing through O_2 and the point of coordinates \mathbf{s}_2 ; both equations are expressed in the base frame.

Therefore, the coordinates of point P at the intersection of the two visual rays can be computed by solving the system of two Eqs. (10.47), (10.48) with

respect to \mathbf{p} . To this end, from (10.47), by computing λ_1 in the third equation and replacing its value into the other two equations, the following system is obtained:

$$\begin{bmatrix} 1 & 0 & -X_1 \\ 0 & 1 & -Y_1 \end{bmatrix} \mathbf{p} = \mathbf{0}. \quad (10.49)$$

Using a similar derivation for the equation obtained by premultiplying both sides of (10.48) by \mathbf{R}^T , the following system is obtained

$$\begin{bmatrix} \mathbf{r}_1^T - X_2 \mathbf{r}_3^T \\ \mathbf{r}_2^T - Y_2 \mathbf{r}_3^T \end{bmatrix} \mathbf{p} = \begin{bmatrix} o_x - o_z X_2 \\ o_y - o_z Y_2 \end{bmatrix}, \quad (10.50)$$

with $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$ and $\mathbf{R}^T \mathbf{o} = [o_x \ o_y \ o_z]^T$. Equations (10.49), (10.50) define a system of four equations and three unknowns, which is linear with respect to \mathbf{p} . Of these equations, only three are independent in the ideal case that the two visual rays intersect at point P . In practical applications, because of noise, these equations are all independent and the system has no solution; hence, suitable algorithms based on least-squares techniques have to be adopted to compute an approximate solution.

Computation of \mathbf{p} can be greatly simplified in the case, rather frequent in applications, that the two cameras have parallel and aligned image planes, with $\mathbf{R} = \mathbf{I}$ and $\mathbf{R}^T \mathbf{o} = [b \ 0 \ 0]^T$, $b > 0$ being the distance between the origins of the two camera frames. This implies that the solution to the system of Eqs. (10.49), (10.50) is

$$p_x = \frac{X_1 b}{X_1 - X_2} \quad (10.51)$$

$$p_y = \frac{Y_1 b}{X_1 - X_2} = \frac{Y_2 b}{X_1 - X_2} \quad (10.52)$$

$$p_z = \frac{b}{X_1 - X_2}. \quad (10.53)$$

10.4.3 Absolute Orientation

In the case of a calibrated system of two cameras observing a rigid object of unknown shape, triangulation can be used to compute the variation of pose of the object or of the system of cameras, due to the relative motion of the object with respect to the cameras. This problem, known as *absolute orientation*, requires the measurement of the positions of the projections of a certain number of feature points of the object.

If the stereo camera system is moving and the object is fixed, let $\mathbf{p}_1, \dots, \mathbf{p}_n$ denote the position vectors of n points of the rigid object measured at time t and $\mathbf{p}'_1, \dots, \mathbf{p}'_n$ the position vectors of the same points measured at time t' using triangulation. These vectors are all referred to Frame 1 and, under the assumption of rigid motion, satisfy the equations

$$\mathbf{p}_i = \mathbf{o} + \mathbf{R} \mathbf{p}'_i \quad i = 1, \dots, n \quad (10.54)$$

where vector \mathbf{o} and rotation matrix \mathbf{R} define the position and orientation displacement of Frame 1 between time t and time t' . The absolute orientation problem consists of the computation of \mathbf{R} and \mathbf{o} from \mathbf{p}_i and \mathbf{p}'_i .

From rigid body mechanics it is known that this problem has a unique solution in the case of three non-collinear points. In this case, nine nonlinear equations can be derived from (10.54), in terms of the nine independent parameters which characterize \mathbf{o} and \mathbf{R} . However, since the points are obtained using triangulation, the measurements are affected by error and the system may have no solution. In this case, it is convenient to consider a number $n > 3$ of points and compute \mathbf{o} and \mathbf{R} as the quantities which minimize the linear quadratic function

$$\sum_{i=1}^n \|\mathbf{p}_i - \mathbf{o} - \mathbf{R}\mathbf{p}'_i\|^2, \quad (10.55)$$

with the constraint that \mathbf{R} is a rotation matrix. The problem of computing \mathbf{o} can be separated from the problem of computing \mathbf{R} observing that the value of \mathbf{o} which minimizes function (10.55) is (see Problem 10.6)

$$\mathbf{o} = \bar{\mathbf{p}} - \mathbf{R}\bar{\mathbf{p}}' \quad (10.56)$$

where $\bar{\mathbf{p}}$ and $\bar{\mathbf{p}}'$ are the centroids of the set of points $\{\mathbf{p}_i\}$ and $\{\mathbf{p}'_i\}$, defined as

$$\bar{\mathbf{p}} = \frac{1}{n} \sum_{i=1}^n \mathbf{p}_i, \quad \bar{\mathbf{p}}' = \frac{1}{n} \sum_{i=1}^n \mathbf{p}'_i.$$

Hence the problem becomes that of computing the rotation matrix \mathbf{R} which minimizes the linear quadratic function

$$\sum_{i=1}^n \|\bar{\mathbf{p}}_i - \mathbf{R}\bar{\mathbf{p}}'_i\|^2, \quad (10.57)$$

where $\bar{\mathbf{p}}_i = \mathbf{p}_i - \bar{\mathbf{p}}$ and $\bar{\mathbf{p}}'_i = \mathbf{p}'_i - \bar{\mathbf{p}}'$ are the deviations with respect to the centroids.

It can be proven that the matrix \mathbf{R} which minimizes function (10.57) is the matrix which maximizes the trace of $\mathbf{R}^T \mathbf{K}$, with

$$\mathbf{K} = \sum_{i=1}^n \bar{\mathbf{p}}_i \bar{\mathbf{p}}'^T_i;$$

see Problem 10.7. Therefore, the solution has the form (10.12) where, for the purpose of this problem, \mathbf{U} and \mathbf{V} are respectively the left and right orthogonal matrices of the singular value decomposition of \mathbf{K} . Once that rotation matrix \mathbf{R} is known, vector \mathbf{o} can be computed using (10.56).

10.4.4 3D Reconstruction from Planar Homography

Another interesting case of application of 3D reconstruction occurs when the feature points of the observed object lie on the same plane. This geometric property represents an additional constraint between the projections of each point on the image plane of the two cameras, besides the epipolar constraint. This constraint is a *planar homography*.

Let \mathbf{p}^2 be the position vector of a point P of the object, expressed with respect to Frame 2. Moreover, let \mathbf{n}^2 denote the unit vector orthogonal to the plane containing the feature points and $d_2 > 0$ the distance of the plane from the origin of Frame 2. By virtue of simple geometric considerations, the following equation can be derived:

$$\frac{1}{d_2} \mathbf{n}^{2T} \mathbf{p}^2 = 1$$

which defines the set of points \mathbf{p}^2 belonging to the plane. In view of the above equality, Eq. (10.43) can be rewritten in the form

$$\mathbf{p}^1 = \mathbf{H} \mathbf{p}^2, \quad (10.58)$$

with

$$\mathbf{H} = \mathbf{R}_2^1 + \frac{1}{d_2} \mathbf{o}_{1,2}^1 \mathbf{n}^{2T}. \quad (10.59)$$

Replacing (10.44) into (10.58) yields

$$\tilde{\mathbf{s}}_1 = \lambda \mathbf{H} \tilde{\mathbf{s}}_2, \quad (10.60)$$

where $\lambda = \lambda_2/\lambda_1 > 0$ is an arbitrary constant. Premultiplication of both sides of (10.60) by $\mathbf{S}(\tilde{\mathbf{s}}_1)$ yields the equality

$$\mathbf{S}(\tilde{\mathbf{s}}_1) \mathbf{H} \tilde{\mathbf{s}}_2 = \mathbf{0}, \quad (10.61)$$

representing a planar homography defined by matrix \mathbf{H} .

Using a derivation similar to that presented in Sect. 10.3.1, it is possible to compute numerically matrix $\zeta \mathbf{H}$, up to a scaling factor ζ , starting from the coordinates of n points of the plane, with $n \geq 4$.

The value of the scaling factor ζ can be computed using a numerical algorithm based on expression (10.59) of matrix \mathbf{H} ; once \mathbf{H} is known, it is possible to compute quantities \mathbf{R}_2^1 , $\mathbf{o}_{1,2}^1/d_2$ and \mathbf{n}^2 in (10.59) — actually, it can be shown that two admissible solutions exist.

This result is of a certain relevance for visual servoing applications. For example, in the case of a camera in motion with respect to the object, if Frames 1 and 2 represent the poses of the camera in two different time instants, the computation of \mathbf{H} with decomposition (10.59) can be used to evaluate the orientation displacement of the camera frame and the position displacement of the origin, the latter defined up to a scaling factor d_2 . This information can be achieved without knowing the object geometry, as long as the feature points all belong to the same plane.

Example 10.5

Consider the SCARA manipulator of Example 10.3 and the planar object of Example 10.1. Assume that the feature vector of the four points of the object, measured by the camera at time t' , is

$$\mathbf{s}' = [0 \quad 0 \quad 0.2 \quad 0 \quad 0.2 \quad 0.2 \quad 0 \quad 0.2]^T,$$

while the feature vector, measured by the camera at time t'' , is

$$\mathbf{s}'' = [-0.1667 \quad 0.1667 \quad -0.0833 \quad 0.0223 \quad 0.0610 \quad 0.1057 \quad -0.0223 \quad 0.2500]^T.$$

It is desired to compute the quantities \mathbf{R}_2^1 , $(1/d_2)\mathbf{o}_{1,2}^1$ and \mathbf{n}^2 of the planar homography (10.59).

For simplicity, assume that the orientation of the planar object is known, namely

$$\mathbf{n}^2 = [0 \quad 0 \quad 1]^T.$$

A further simplification to the problem derives from the fact that, in this case, matrix \mathbf{R}_2^1 corresponds to a rotation about axis z of an angle β , namely $\mathbf{R}_2^1 = \mathbf{R}_z(\beta)$. Therefore, in view of (10.59), planar homography \mathbf{H} has the symbolic expression

$$\mathbf{H} = \begin{bmatrix} c_\beta & -s_\beta & o_x/d_2 \\ s_\beta & c_\beta & o_y/d_2 \\ 0 & 0 & 1 + o_z/d_2 \end{bmatrix},$$

with $\mathbf{o}_{1,2}^1 = [o_x \quad o_y \quad o_z]^T$.

On the other hand, starting from numerical values of \mathbf{s}' and \mathbf{s}'' and using a derivation similar to that used in Example 10.1, the following matrix can be obtained:

$$\zeta \mathbf{H} = \begin{bmatrix} 0.3015 & -0.5222 & 0.1373 \\ 0.5222 & 0.3015 & 0.0368 \\ 0 & 0 & 0.5025 \end{bmatrix},$$

which is the numerical value of the planar homography, up to a scaling factor ζ .

The symbolic expression of \mathbf{H} reveals that the first and second column of this matrix have unit norm. This property can be used to compute $|\zeta| = 1.6583$. The sign of ζ can be evaluated by imposing, for any of the feature points of the object, the constraint

$$\mathbf{p}^{1T} \mathbf{p}^1 = \mathbf{p}^{1T} \mathbf{H} \mathbf{p}^2 = \lambda_1 \lambda_2 \tilde{\mathbf{s}}_1^T \mathbf{H} \tilde{\mathbf{s}}_2 > 0.$$

Since scalars λ_i are positive, this inequality is equivalent to

$$\tilde{\mathbf{s}}_1^T \mathbf{H} \tilde{\mathbf{s}}_2 > 0,$$

hence, in this case, it is $\zeta > 0$. Therefore

$$\mathbf{H} = \begin{bmatrix} 0.5 & -0.8660 & 0.2277 \\ 0.8660 & 0.5 & 0.0610 \\ 0 & 0 & 0.8333 \end{bmatrix}.$$

At this point, the value of angle β can be computed from the elements h_{11} and h_{21} of matrix \mathbf{H} using equation

$$\beta = \text{Atan2}(h_{21}, h_{11}) = \frac{\pi}{3}.$$

Finally, vector $(1/d_2)\mathbf{o}_{1,2}^1$ can be computed from the elements of the last row of matrix \mathbf{H} using the equality

$$\frac{1}{d_2}\mathbf{o}_{1,2}^1 = \begin{bmatrix} h_{13} \\ h_{23} \\ h_{33} - 1 \end{bmatrix} = \begin{bmatrix} 0.2277 \\ 0.0610 \\ -0.0667 \end{bmatrix}.$$

Notice that the derivation illustrated in Example 10.5, in the case that \mathbf{n}^2 is unknown and \mathbf{R}_2^1 is a generic rotation matrix, becomes much more complex.

10.5 Camera Calibration

An important problem for visual servoing applications is the calibration of the camera, which is the sensor providing the information fed back to the controller. Calibration consists of the estimation of the *intrinsic parameters*, characterizing matrix $\mathbf{\Omega}$ defined in (5.41), and of the *extrinsic parameters*, characterizing the pose of the camera frame with respect to the base frame (for eye-to-end cameras) or to the end-effector frame (for eye-in-hand cameras). Various calibration techniques exist, which are based on algorithms similar to those used for the pose estimation of an object with respect to a camera.

In particular, the solution method of a PnP problem with n coplanar points illustrated in Sect. 10.3.1 can be directly used for the computation of the extrinsic parameters of the camera, if the intrinsic parameters are known.

In fact, in view of (10.14), extrinsic parameters of an eye-to-hand camera can be computed as

$$\mathbf{T}_c^b = \mathbf{T}_o^b(\mathbf{T}_o^c)^{-1}, \quad (10.62)$$

where matrix \mathbf{T}_o^c is the output of the algorithm solving a PnP planar problem, provided that matrix \mathbf{T}_o^b , expressing the position and the orientation of the object frame with respect to the base frame, is known. Similarly, in view of (10.15), the extrinsic parameters of an eye-in-hand camera can be computed as

$$\mathbf{T}_c^e = \mathbf{T}_o^e(\mathbf{T}_o^c)^{-1}, \quad (10.63)$$

provided that matrix \mathbf{T}_o^e , expressing the pose of the object frame with respect to the end-effector frame, is known.

If the intrinsic parameters are not known, the derivation of Sect. 10.3.1 has to be suitably extended and can be broken down in the three phases described below.

Phase 1

A planar homography can be computed starting from pixel coordinates

$$\mathbf{c}_i = \begin{bmatrix} X_{Ii} \\ Y_{Ii} \end{bmatrix}$$

in place of normalized coordinates \mathbf{s}_i . In detail, an equation formally identical to (10.7) can be obtained:

$$\mathbf{S}(\tilde{\mathbf{c}}_i)\mathbf{H}' [r_{x,i} \ r_{y,i} \ 1]^T = \mathbf{0}, \quad (10.64)$$

where \mathbf{H}' is the (3×3) matrix

$$\mathbf{H}' = \mathbf{\Omega}\mathbf{H}, \quad (10.65)$$

$\mathbf{\Omega}$ being the matrix of intrinsic parameters (5.41) and \mathbf{H} the matrix defined in (10.8). Using an algorithm similar to that presented in Sect. 10.3.1, planar homography $\zeta\mathbf{H}'$ can be computed, up to a scaling factor ζ , from the coordinates of n points of the plane, with $n \geq 4$.

Phase 2

Matrix $\mathbf{\Omega}$ can be computed from the elements of matrix $\zeta\mathbf{H}'$. In fact, taking into account (10.65), and the definition of \mathbf{H} in (10.8), yields

$$\zeta [\mathbf{h}'_1 \ \mathbf{h}'_2 \ \mathbf{h}'_3] = \mathbf{\Omega} [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{o}_{c,o}^c],$$

where \mathbf{h}'_i denotes the i -th column of matrix \mathbf{H}' . Computing \mathbf{r}_1 and \mathbf{r}_2 from this equation, and imposing the orthogonality and unit norm constraints on these vectors, the following two scalar equations are obtained:

$$\begin{aligned} \mathbf{h}'_1{}^T \mathbf{\Omega}^{-T} \mathbf{\Omega}^{-1} \mathbf{h}'_2 &= 0 \\ \mathbf{h}'_1{}^T \mathbf{\Omega}^{-T} \mathbf{\Omega}^{-1} \mathbf{h}'_1 &= \mathbf{h}'_2{}^T \mathbf{\Omega}^{-T} \mathbf{\Omega}^{-1} \mathbf{h}'_2 \end{aligned}$$

which, being linear, can be rewritten in the form

$$\mathbf{A}'\mathbf{b} = \mathbf{0}. \quad (10.66)$$

In the above equation, \mathbf{A}' is a (2×6) matrix of coefficients depending on $\mathbf{h}'_1, \mathbf{h}'_2$, while $\mathbf{b} = [b_{11} \ b_{12} \ b_{22} \ b_{13} \ b_{23} \ b_{33}]^T$, where b_{ij} is the generic element of the symmetric matrix

$$\mathbf{B} = \mathbf{\Omega}^{-T} \mathbf{\Omega}^{-1} = \begin{bmatrix} 1/\alpha_x^2 & 0 & -X_0/\alpha_x^2 \\ * & 1/\alpha_y^2 & -Y_0/\alpha_y^2 \\ * & * & 1 + X_0^2/\alpha_x^2 + Y_0^2/\alpha_y^2 \end{bmatrix}. \quad (10.67)$$

By repeating Phase 1 k times, with the same plane placed each time in a different pose, $2k$ equations of the form (10.66) are obtained. These equations, in the case $k \geq 3$, have a unique solution $\gamma\mathbf{b}$ defined up to a scaling factor γ . From matrix $\gamma\mathbf{B}$, in view of (10.67), the following expressions for the intrinsic parameters can be found:

$$X_0 = -b'_{13}/b'_{11}$$

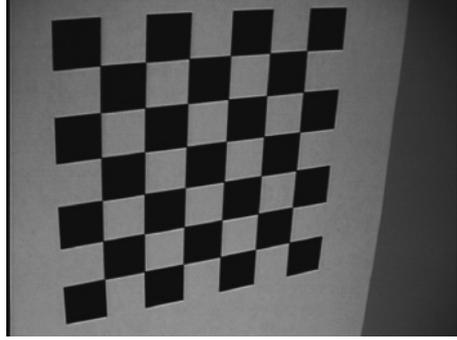


Fig. 10.11. Example of calibration plane

$$\begin{aligned}
 Y_0 &= -b'_{23}/b'_{22} \\
 \alpha_x &= \sqrt{\gamma/b'_{11}} \\
 \alpha_y &= \sqrt{\gamma/b'_{22}},
 \end{aligned}$$

where $b'_{i,j} = \gamma b_{i,j}$ and γ can be computed as

$$\gamma = b'_{13} + b'_{23} + b'_{33}.$$

Phase 3

Once the matrix $\mathbf{\Omega}$ of intrinsic parameters has been evaluated, it is possible to compute \mathbf{H} from \mathbf{H}' , up to a scaling factor ζ , using (10.65) for one of the k poses of the plane. Hence, matrix \mathbf{T}_o^c can be computed from \mathbf{H} as for the case of the solution of a PnP problem shown in Sect. 10.3.1. Finally, using equations (10.62), (10.63), the extrinsic parameters of the camera can be evaluated.

The method illustrated above is merely conceptual, because it does not provide satisfactory solutions in the presence of measurement noise or lens distortion — especially for the intrinsic parameters; however, the accuracy of the result can be improved using models that take into account the distortion phenomena of the lenses, together with nonlinear optimization techniques.

From the experimental point of view, the calibration method described above requires the use of calibration planes where a certain number of points can be easily detected; also, the position of these points with respect to a suitable reference frame must be known with high accuracy. An example of calibration plane with a chessboard pattern is shown in Fig. 10.11.

Finally, notice that a calibration method can also be set up starting from the solution of a nonplanar PnP problem, using the direct linear transformation method.

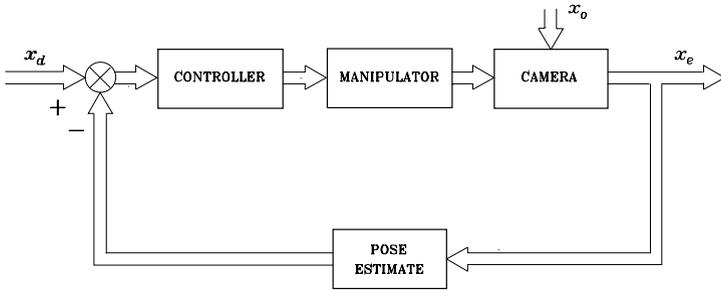


Fig. 10.12. General block scheme of position-based visual servoing

10.6 The Visual Servoing Problem

Visual measurements allow a robot to collect information on the surrounding environment. In the case of robot manipulators, such information is typically used to compute the end-effector pose with respect to an object observed by the camera. The objective of visual servoing is to ensure that the end-effector, on the basis of visual measurements elaborated in real time, reaches and keeps a (constant or time-varying) desired pose with respect to the observed object.

It is worth remarking that the direct measurements provided by the visual system are concerned with feature parameters in the image plane, while the robotic task is defined in the operational space, in terms of the relative pose of the end-effector with respect to the object. This fact naturally leads to considering two kinds of control approaches, illustrated by the block schemes of Figs. 10.12 and 10.13, namely *position-based visual servoing*, also termed *visual servoing in the operational space*, and *image-based visual servoing*, also termed *visual servoing in the image space*. In these schemes, the case of eye-in-hand camera is considered; for eye-to-hand cameras, similar schemes can be adopted.

The *position-based visual servoing* approach is conceptually similar to the operational space control illustrated in Fig. 8.2. The main difference is that feedback is based on the real-time estimation of the pose of the observed object with respect to the camera using visual measurements. Estimation can be performed analytically or using iterative numerical algorithms. Its conceptual advantage regards the possibility of acting directly on operational space variables. Therefore, the control parameters can be selected on the basis of suitable specifications imposed to the time response of the end-effector motion variables, both at steady state and during the transient. The drawback of this approach is that, due to the absence of direct control of the image features, the object may exit from the camera field of view during the transient or as a consequence of planning errors; hence, the feedback loop turns out to be open due to lack of visual measurements and instability may occur.

In the *image-space visual servoing* approach, the control action is computed on the basis of the error defined as the difference between the value of

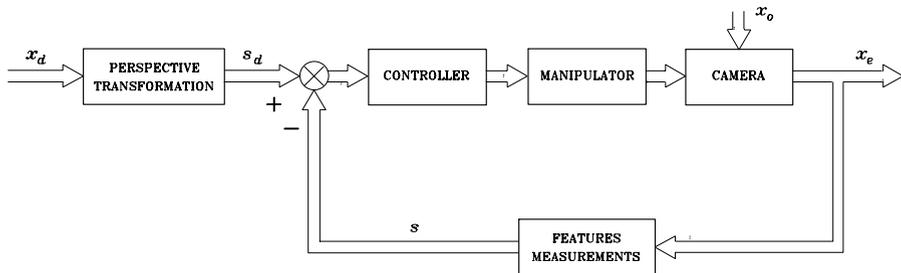


Fig. 10.13. General block scheme of image-based visual servoing

the image feature parameters in the desired configuration — computed using perspective transformation or directly measured with the camera in the desired pose — and the value of the parameters measured with the camera in the current pose. The conceptual advantage of this solution regards the fact that the real-time estimate of the pose of the object with respect to the camera is not required. Moreover, since the control acts directly in the image feature parameters, it is possible to keep the object within the camera field of view during the motion. However, due to the nonlinearity of the mapping between the image feature parameters and the operational space variables, singular configurations may occur, which cause instability or saturation of the control action. Also, the end-effector trajectories cannot be easily predicted in advance and may produce collisions with obstacles or joint limits violation.

To compare the two control strategies it is also worth considering the operating conditions. Of particular importance is the issue of camera calibration. It is easy to understand that position-based visual servoing is more sensitive to camera calibration errors compared to image-based visual servoing. In fact, for the first approach, the presence of uncertainties on calibration parameters, both intrinsic and extrinsic, produces errors on the estimate of operational space variables that may be seen as an external disturbance acting on the feedback path of the control loop, where disturbance rejection capability is low. On the other hand, in the image-based visual servoing approach, the quantities used for the computation of the control action are directly defined in the image plane and measured in pixel units; moreover, the desired value of the feature parameters is measured using the camera. This implies that the uncertainty affecting calibration parameters can be seen as a disturbance acting on the forward path of the control loop, where disturbance rejection capability is high.

A further aspect to analyze concerns knowledge of the geometric model of the object. It is evident that, for position-based visual servoing, the object geometry must be known if only one camera is used, because it is necessary for pose estimation, while it may be unknown when a stereo camera system is

employed. On the other hand, image-based visual servoing does not require, in principle, knowledge of the object geometry, even for mono-camera systems.

On the above premises, in the following, the main position-based and image-based visual servoing schemes are illustrated. For both approaches, the problem of regulation to a constant set-point is presented and the object is assumed to be fixed with respect to the base frame. Without loss of generality, the case of a single calibrated camera, mounted on the manipulator's end-effector, is considered (see Fig. 10.5); moreover, the end-effector frame is chosen so as to coincide with the camera frame.

10.7 Position-based Visual Servoing

In position-based visual servoing schemes, visual measurements are used to estimate in real time the homogeneous transformation matrix \mathbf{T}_o^c , representing the relative pose of the object frame with respect to the camera frame. From matrix \mathbf{T}_o^c , the $(m \times 1)$ vector of independent coordinates $\mathbf{x}_{c,o}$, defined in (10.31), can be extracted.

Assumed that the object is fixed with respect to the base frame, the position-based visual servoing problem can be formulated by imposing a desired value to the relative pose of the object frame with respect to the camera frame. This quantity can be specified in terms of homogeneous transformation matrix \mathbf{T}_o^d , where superscript d denotes the desired pose of the camera frame. From this matrix, the $(m \times 1)$ operational space vector $\mathbf{x}_{d,o}$ can be extracted.

Matrices \mathbf{T}_o^c and \mathbf{T}_o^d can be used to obtain the homogeneous transformation matrix

$$\mathbf{T}_c^d = \mathbf{T}_o^d (\mathbf{T}_o^c)^{-1} = \begin{bmatrix} \mathbf{R}_c^d & \mathbf{o}_{d,c}^d \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (10.68)$$

expressing the position and orientation displacement of the camera frame in the current pose with respect to the desired pose. From this matrix, a suitable error vector in the operational space can be computed, defined as

$$\tilde{\mathbf{x}} = - \begin{bmatrix} \mathbf{o}_{d,c}^d \\ \phi_{d,c} \end{bmatrix}, \quad (10.69)$$

where $\phi_{d,c}$ is the vector of the Euler angles extracted from rotation matrix \mathbf{R}_c^d . Vector $\tilde{\mathbf{x}}$ does not depend on the object pose and represents the error between the desired pose and the current pose of the camera frame. It is worth observing that this vector does not coincide with the difference between $\mathbf{x}_{d,o}$ and $\mathbf{x}_{c,o}$, but it can be computed from the corresponding homogeneous transformation matrices, using (10.68), (10.69).

The control has to be designed so that the operational space error $\tilde{\mathbf{x}}$ tends to zero asymptotically.

Notice that, for the choice of the set point $\mathbf{x}_{d,o}$, the knowledge of the object pose is not required. However, the control objective can be satisfied provided

that the desired pose of the camera frame with respect to the base frame, corresponding to the homogeneous transformation matrix

$$\mathbf{T}_d = \mathbf{T}_c(\mathbf{T}_c^d)^{-1} = \begin{bmatrix} \mathbf{R}_d & \mathbf{o}_d \\ \mathbf{0}^T & 1 \end{bmatrix}, \quad (10.70)$$

belongs to the dexterous workspace of the manipulator. If the object is fixed with respect to the base frame, this matrix is constant.

10.7.1 PD Control with Gravity Compensation

Position-based visual servoing can be implemented using PD control with gravity compensation, suitably modified with respect to that used for motion control.

Computing the time derivative of (10.69), for the position part, gives

$$\dot{\mathbf{o}}_{d,c}^d = \dot{\mathbf{o}}_c^d - \dot{\mathbf{o}}_d^d = \mathbf{R}_d^T \dot{\mathbf{o}}_c,$$

while, for the orientation part, it gives

$$\dot{\phi}_{d,c} = \mathbf{T}^{-1}(\phi_{d,c})\boldsymbol{\omega}_{d,c}^d = \mathbf{T}^{-1}(\phi_{d,c})\mathbf{R}_d^T\boldsymbol{\omega}_c.$$

To compute the above expressions, equalities $\dot{\mathbf{o}}_d^d = \mathbf{0}$ and $\boldsymbol{\omega}_d^d = \mathbf{0}$ have been taken into account, observing that \mathbf{o}_d and \mathbf{R}_d are constant. Therefore, $\dot{\tilde{\mathbf{x}}}$ has the expression

$$\dot{\tilde{\mathbf{x}}} = -\mathbf{T}_A^{-1}(\phi_{d,c}) \begin{bmatrix} \mathbf{R}_d^T & \mathbf{O} \\ \mathbf{O} & \mathbf{R}_d^T \end{bmatrix} \mathbf{v}_c. \quad (10.71)$$

Since the end-effector frame and the camera frame coincide, the following equality holds:

$$\dot{\tilde{\mathbf{x}}} = -\mathbf{J}_{A_d}(\mathbf{q}, \tilde{\mathbf{x}})\dot{\mathbf{q}}, \quad (10.72)$$

where the matrix

$$\mathbf{J}_{A_d}(\mathbf{q}, \tilde{\mathbf{x}}) = \mathbf{T}_A^{-1}(\phi_{d,c}) \begin{bmatrix} \mathbf{R}_d^T & \mathbf{O} \\ \mathbf{O} & \mathbf{R}_d^T \end{bmatrix} \mathbf{J}(\mathbf{q}) \quad (10.73)$$

has the meaning of analytical Jacobian of the manipulator in the operational space, as for the Jacobian in (9.14).

Position-based visual servoing of PD type with gravity compensation has the expression

$$\mathbf{u} = \mathbf{g}(\mathbf{q}) + \mathbf{J}_{A_d}^T(\mathbf{q}, \tilde{\mathbf{x}})(\mathbf{K}_P\tilde{\mathbf{x}} - \mathbf{K}_D\mathbf{J}_{A_d}(\mathbf{q}, \tilde{\mathbf{x}})\dot{\mathbf{q}}), \quad (10.74)$$

analogous to motion control law (8.110), but using a different definition of operational space error. The asymptotic stability of the equilibrium pose corresponding to $\tilde{\mathbf{x}} = \mathbf{0}$, under the assumption of symmetric and positive definite matrices \mathbf{K}_P , \mathbf{K}_D , can be proven using the Lyapunov function

$$V(\dot{\mathbf{q}}, \tilde{\mathbf{x}}) = \frac{1}{2}\dot{\mathbf{q}}^T \mathbf{B}(\mathbf{q})\dot{\mathbf{q}} + \frac{1}{2}\tilde{\mathbf{x}}^T \mathbf{K}_P\tilde{\mathbf{x}} > 0 \quad \forall \dot{\mathbf{q}}, \tilde{\mathbf{x}} \neq \mathbf{0},$$

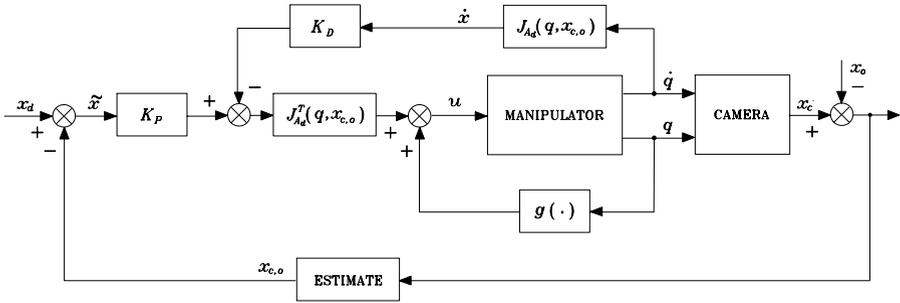


Fig. 10.14. Block scheme of position-based visual servoing of PD type with gravity compensation

similarly to the case of control law (8.110).

Notice that, for the computation of control law (10.74), the estimation of $\mathbf{x}_{c,o}$ and the measurements of \mathbf{q} and $\dot{\mathbf{q}}$ are required. Moreover, the derivative term can also be chosen as $-\mathbf{K}_D \dot{\mathbf{q}}$.

The block scheme of position-based visual servoing of PD type with gravity compensation is shown in Fig. 10.14. Notice that the sum block computing error $\tilde{\mathbf{x}}$ and that computing the output of the controlled system have a purely conceptual meaning and do not correspond to algebraic sums.

10.7.2 Resolved-velocity Control

The information deriving from visual measurements is computed at a frequency lower or equal to the camera frame rate. This quantity, especially for CCD cameras, is at least of one order of magnitude lower than the typical frequencies used for motion control of robot manipulators. As a consequence, in the digital implementation of control law (10.74), to preserve stability of the closed-loop system, the control gains must be set to values much lower than those typically used for motion control; therefore, the performance of the closed-loop system in terms of speed of convergence and disturbance rejection capability turns out to be poor.

This problem can be avoided assuming that the manipulator is equipped with a high-gain motion controller in the joint space or in the operational space. Neglecting the effects on the tracking errors deriving from manipulator dynamics and disturbances, the controlled manipulator can be considered as an ideal positioning device. This implies that, in the case of joint space motion control, the following equality holds:

$$\mathbf{q}(t) \approx \mathbf{q}_r(t), \tag{10.75}$$

$\mathbf{q}_r(t)$ being the imposed reference trajectory for the joint variables.

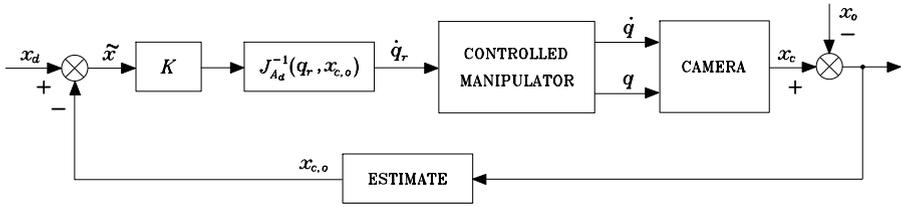


Fig. 10.15. Block scheme of resolved-velocity position-based visual servoing

Therefore, visual servoing can be achieved by computing the trajectory $\mathbf{q}_r(t)$ on the basis of visual measurements, so that the operational space tracking error (10.69) goes asymptotically to zero.

To this end, Eq. (10.72) suggests the following choice for the joint space reference velocity:

$$\dot{\mathbf{q}}_r = \mathbf{J}_{A_d}^{-1}(\mathbf{q}_r, \tilde{\mathbf{x}}) \mathbf{K} \tilde{\mathbf{x}} \tag{10.76}$$

which, replaced in (10.72), by virtue of equality (10.75), yields the linear equation

$$\dot{\tilde{\mathbf{x}}} + \mathbf{K} \tilde{\mathbf{x}} = \mathbf{0}. \tag{10.77}$$

This equality, for a positive definite matrix \mathbf{K} , implies that the operational space error tends to zero asymptotically with a convergence of exponential type and speed depending on the eigenvalues of matrix \mathbf{K} ; the larger the eigenvalues, the faster the convergence.

The above scheme is termed *resolved-velocity control* in the operational space, because it is based on the computation of velocity $\dot{\mathbf{q}}_r$ from the operational space error. Trajectory $\mathbf{q}_r(t)$ is computed from (10.76) via a simple integration.

The block scheme of resolved-velocity position-based visual servoing is reported in Fig. 10.15. Also in this case, the sum block computing the error $\tilde{\mathbf{x}}$ and that computing the output of the scheme have a purely conceptual meaning and do not correspond to algebraic sums.

Notice that the choice of \mathbf{K} influences the transient behaviour of the trajectory of the camera frame, which is the solution to the differential equation (10.77). If \mathbf{K} is a diagonal matrix with the same gains for the positional part, the origin of the camera frame follows the line segment connecting the initial position to the desired position. On the other hand, the orientation trajectory depends on the particular choice of Euler angles and, more in general, of the orientation error. The possible choices of the orientation error are those presented in Sect. 3.7.3, with the appropriate definition of Jacobian (10.73). The possibility of knowing in advance the trajectory of the camera is important because, during the motion, the object may exit from the camera field of view, making visual measurements unavailable.

10.8 Image-based Visual Servoing

If the object is fixed with respect to the base frame, image-based visual servoing can be formulated by stipulating that the vector of the object feature parameters has a desired constant value \mathbf{s}_d corresponding to the desired pose of the camera. Therefore, it is implicitly assumed that a desired pose $\mathbf{x}_{d,o}$ exists so that the camera pose belongs to the dexterous workspace of the manipulator and

$$\mathbf{s}_d = \mathbf{s}(\mathbf{x}_{d,o}). \quad (10.78)$$

Moreover, $\mathbf{x}_{d,o}$ is supposed to be unique. To this end, the feature parameters can be chosen as the coordinates of n points of the object, with $n \geq 4$ for coplanar points (and no triplets of collinear points) or $n \geq 6$ in case of non-coplanar points. Notice that, if the operational space dimension is $m < 6$, as for the case of SCARA manipulator, a reduced number of points can be used.

The interaction matrix $\mathbf{L}_s(\mathbf{s}, \mathbf{z}_c)$ depends on variables \mathbf{s} and \mathbf{z}_c with $\mathbf{z}_c = [z_{c,1} \dots z_{c,n}]^T$, $z_{c,i}$ being the third coordinate of the generic feature point of the object.

It is worth noticing that the task is assigned directly in terms of feature vector \mathbf{s}_d , while pose $\mathbf{x}_{d,o}$ does not need to be known. In fact, \mathbf{s}_d can be computed by measuring the feature parameters when the object is in the desired pose with respect to the camera.

The control law must be designed so as to guarantee that the image space error

$$\mathbf{e}_s = \mathbf{s}_d - \mathbf{s} \quad (10.79)$$

tends asymptotically to zero.

10.8.1 PD Control with Gravity Compensation

Image-based visual servoing can be implemented using a PD control with gravity compensation defined on the basis of the image space error.

To this end, consider the following positive definite quadratic form as Lyapunov function candidate:

$$\mathbf{V}(\dot{\mathbf{q}}, \mathbf{e}_s) = \frac{1}{2} \dot{\mathbf{q}}^T \mathbf{B}(\mathbf{q}) \dot{\mathbf{q}} + \frac{1}{2} \mathbf{e}_s^T \mathbf{K}_{Ps} \mathbf{e}_s > 0 \quad \forall \dot{\mathbf{q}}, \mathbf{e}_s \neq \mathbf{0}, \quad (10.80)$$

with \mathbf{K}_{Ps} symmetric and positive definite ($k \times k$) matrix.

Computing the time derivative of (10.80) and taking into account the expression (8.7) of the joint space dynamic model of the manipulator and property (7.49) yields

$$\dot{\mathbf{V}} = -\dot{\mathbf{q}}^T \mathbf{F} \dot{\mathbf{q}} + \dot{\mathbf{q}}^T (\mathbf{u} - \mathbf{g}(\mathbf{q})) + \dot{\mathbf{e}}_s^T \mathbf{K}_{Ps} \mathbf{e}_s. \quad (10.81)$$

Since $\dot{\mathbf{s}}_d = \mathbf{0}$ and the object is fixed with respect to the base frame, the following equality holds:

$$\dot{\mathbf{e}}_s = -\dot{\mathbf{s}} = -\mathbf{J}_L(\mathbf{s}, \mathbf{z}_c, \mathbf{q}) \dot{\mathbf{q}}, \quad (10.82)$$

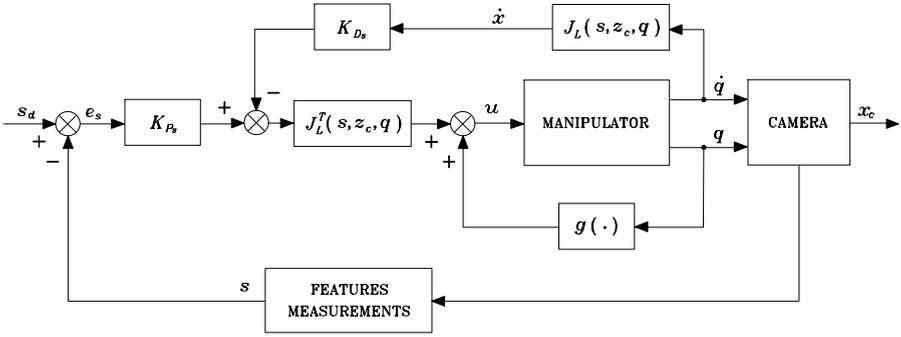


Fig. 10.16. Block scheme of image-based visual servoing of PD type with gravity compensation

where

$$J_L(s, z_c, q) = L_s(s, z_c) \begin{bmatrix} R_c^T & O \\ O & R_c^T \end{bmatrix} J(q), \quad (10.83)$$

the camera frame and the end-effector frame being coincident.

Therefore, with the choice

$$u = g(q) + J_L^T(s, z_c, q) (K_{Ps} e_s - K_{Ds} J_L(s, z_c, q) \dot{q}), \quad (10.84)$$

where K_{Ds} is a symmetric and positive definite ($k \times k$) matrix, Eq. (10.81) becomes

$$\dot{V} = -\dot{q}^T F \dot{q} - \dot{q}^T J_L^T K_{Ds} J_L \dot{q}. \quad (10.85)$$

Control law (10.84) includes a nonlinear compensation action of gravitational forces in the joint space and a linear PD action in the image space. The last term, in view of (10.82), corresponds to a derivative action in the image space and has been added to increase damping. The resulting block scheme is reported in Fig. 10.16.

The direct measurement of \dot{s} would permit the computation of the derivative term as $-K_{Ds} \dot{s}$; this measurement, however, is not available. As an alternative, the derivative term can simply be set as $-K_D \dot{q}$, with K_D symmetric and positive definite ($n \times n$) matrix.

Equation (10.85) reveals that, for all trajectories of the system, the Lyapunov function decreases until $\dot{q} \neq 0$. Therefore the system reaches an equilibrium state, characterized by

$$J_L^T(s, z_c, q) K_{Ps} e_s = 0. \quad (10.86)$$

Equations (10.86), (10.83) show that, if the interaction matrix and the geometric Jacobian of the manipulator are full rank, then $e_s = 0$, which is the sought result.

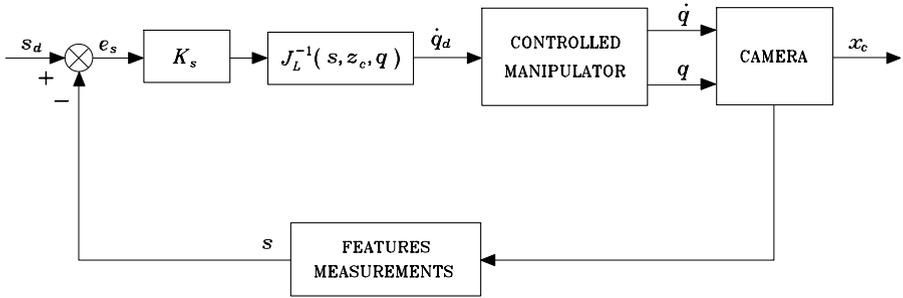


Fig. 10.17. Block scheme of resolved-velocity image-based visual servoing

Notice that control law (10.84) requires not only the measurement of \mathbf{s} but also the computation of vector \mathbf{z}_c which, in the image-based visual servoing philosophy, should be avoided. In some applications \mathbf{z}_c is known with good approximation, as in the case that the relative motion of the camera with respect to the object belongs to a plane. Alternatively, estimated or constant values can be used for \mathbf{z}_c , as the value in the initial configuration or that in the desired configuration. This is equivalent to using an estimate $\hat{\mathbf{L}}_s$ of the interaction matrix. In such cases, however, the stability proof becomes much more complex.

10.8.2 Resolved-velocity Control

The concept of resolved-velocity control can easily be extended to the image space. In such a case, Eq. (10.82) suggests the following choice of the reference velocity in joint space

$$\dot{\mathbf{q}}_r = \mathbf{J}_L^{-1}(\mathbf{s}, \mathbf{z}_c, \mathbf{q}_r) \mathbf{K}_s \mathbf{e}_s, \quad (10.87)$$

under the assumption of invertible matrix \mathbf{J}_L . This control law, replaced in (10.82), yields the linear equation

$$\dot{\mathbf{e}}_s + \mathbf{K}_s \mathbf{e}_s = \mathbf{0}. \quad (10.88)$$

Therefore, if \mathbf{K}_s is a positive definite matrix, Eq. (10.88) is asymptotically stable and error \mathbf{e}_s tends asymptotically to zero with convergence of exponential type and speed depending on the eigenvalues of matrix \mathbf{K}_s . The convergence to zero of the image space error \mathbf{e}_s ensures the asymptotic convergence of $\mathbf{x}_{c,o}$ to the desired pose $\mathbf{x}_{d,o}$.

The block scheme of resolved-velocity image-based visual servoing is shown in Fig. 10.17.

Notice that this control scheme requires the computation of the inverse of matrix \mathbf{J}_L ; therefore it is affected by problems related to the singularities of this matrix which, in view of (10.83), are both those of the geometric Jacobian and those of the interaction matrix. The most critical singularities are those

of the interaction matrix, since they depend on the choice of the image feature parameters.

Therefore, it is convenient to compute control law (10.87) in two steps. The first step is the computation of vector

$$\mathbf{v}_r^c = \mathbf{L}_s^{-1}(\mathbf{s}, \mathbf{z}_c) \mathbf{K}_s \mathbf{e}_s. \quad (10.89)$$

The second step is the computation of the joint space reference velocity using the relationship

$$\dot{\mathbf{q}}_r = \mathbf{J}^{-1}(\mathbf{q}) \begin{bmatrix} \mathbf{R}_c & \mathbf{O} \\ \mathbf{O} & \mathbf{R}_c \end{bmatrix} \mathbf{v}_r^c. \quad (10.90)$$

Far from the kinematic singularities of the manipulator, the problem of the singularities of the interaction matrix can be overcome by using a number k of feature parameters greater than the minimum required m , similarly to the case considered in Sect. 10.3.3. The control law can be modified by using the left pseudo-inverse of interaction matrix \mathbf{L}_s in place of the inverse, namely

$$\mathbf{v}_r^c = (\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T \mathbf{K}_s \mathbf{e}_s \quad (10.91)$$

in place of (10.89). Stability of the closed-loop system with control law (10.90), (10.91) can be shown using the direct Lyapunov method based on the positive definite function

$$V(\mathbf{e}_s) = \frac{1}{2} \mathbf{e}_s^T \mathbf{K}_s \mathbf{e}_s > 0 \quad \forall \mathbf{e}_s \neq \mathbf{0}.$$

Computing the time derivative of this function and taking into account (10.82), (10.83), (10.90), (10.91), yields

$$\dot{V} = -\mathbf{e}_s^T \mathbf{K}_s \mathbf{L}_s (\mathbf{L}_s^T \mathbf{L}_s)^{-1} \mathbf{L}_s^T \mathbf{K}_s \mathbf{e}_s$$

which is negative semi-definite because $\mathcal{N}(\mathbf{L}_s^T) \neq \emptyset$, \mathbf{L}_s^T being a matrix with more columns than rows. Therefore, the closed-loop system is stable but not asymptotically stable. This implies that the error is bounded, but in some cases the system may reach an equilibrium with $\mathbf{e}_s \neq \mathbf{0}$ and $\mathbf{K}_s \mathbf{e}_s \in \mathcal{N}(\mathbf{L}_s^T)$.

Another problem connected with the implementation of control law (10.89) or (10.91) and (10.90) depends on the fact that the computation of interaction matrix \mathbf{L}_s requires knowledge of \mathbf{z}_c . Similar to Sect. 10.8.1, this problem can be solved by using an estimate of matrix $\widehat{\mathbf{L}}_s^{-1}$ (or of its pseudo-inverse). In this case, the Lyapunov method can be used to prove that the control scheme remains stable provided that matrix $\mathbf{L}_s \widehat{\mathbf{L}}_s^{-1}$ is positive definite. Notice that \mathbf{z}_c is the only information depending on object geometry. Therefore, it can also be seen that image-based visual servoing, in the case that only one camera is used, does not require exact knowledge of object geometry.

The choice of the elements of matrix \mathbf{K}_s influences the trajectories of the feature parameters, which are solution to differential equation (10.88). In the

case of feature points, if a diagonal matrix \mathbf{K}_s with equal elements is set, the projections of these points on the image plane will follow line segments. The corresponding camera motion, however, cannot be easily predicted, because of the nonlinearity of the mapping between image plane variables and operational space variables.

10.9 Comparison Among Various Control Schemes

In order to make a comparison between the various control schemes presented, consider the SCARA manipulator of Example 10.3 and the planar object of Example 10.1. The base frame of the SCARA manipulator of Fig. 2.36 is set with the origin at the intersection of the axis of joint 1 with the horizontal plane containing the origin of the end-effector frame when $d_3 = 0$, d_3 being the prismatic joint variable; the axis z of the base frame points downward. The operational space, of dimension $m = 4$, is characterized by vector $\mathbf{x}_{c,o}$ in (10.33).

With reference to the dynamic model of Problem 7.3, the same data of Example 7.2 are considered; in addition, $m_{\ell_3} = 2 \text{ kg}$ and $I_{\ell_4} = 1 \text{ kg}\cdot\text{m}^2$, while the contribution of the motors of the last two links are neglected.

For position-based visual servoing schemes, the real-time estimation of vector $\mathbf{x}_{c,o}$ from a suitable feature vector is required. To this end, the algorithm of Sect. 10.3.3 can be used, based on the inverse of the image Jacobian. This is a classical visual tracking application that can be accomplished using only two feature points, because the corresponding Jacobian \mathbf{J}_{A_s} is a (4×4) matrix. The selected points are P_1 and P_2 of the object of Fig. 10.6.

The same points can also be used for image-based visual servoing, because the corresponding Jacobian \mathbf{J}_L is a (4×4) matrix.

It is assumed that at time $t = 0$ the pose of the camera frame, with respect to the base frame, is defined by the operational space vector

$$\mathbf{x}_c(0) = [1 \quad 1 \quad 0.5 \quad \pi/4]^T$$

and the pose of the object frame, with respect to the camera frame, is defined by the operational space vector

$$\mathbf{x}_{c,o}(0) = [0 \quad 0 \quad 0.5 \quad 0]^T.$$

The desired pose of the object frame with respect to the camera frame is defined by vector

$$\mathbf{x}_{d,o} = [-0.1 \quad 0.1 \quad 0.6 \quad -\pi/3]^T.$$

This quantity is assumed as the initial value of the pose estimation algorithm used by position-based visual servoing schemes.

For image-based visual servoing schemes, the desired value of the feature parameters of points P_1 and P_2 of the object, in the desired pose $\mathbf{x}_{d,o}$, is

$$\mathbf{s}_d = [-0.1667 \quad 0.1667 \quad -0.0833 \quad 0.0223]^T.$$

For all the schemes, a discrete-time implementation of the controller with sampling time of 0.04 s has been adopted, corresponding to a 25 Hz frequency. This value coincides with the minimum frame rate of analog cameras and allows the use of visual measurements also in the worst case.

In the numerical simulations, the following control schemes have been utilized:

- A.** Position-based visual servoing of PD type with gravity compensation with the following data:

$$\mathbf{K}_P = \text{diag}\{500, 500, 10, 10\}$$

$$\mathbf{K}_D = \text{diag}\{500, 500, 10, 10\}.$$

- B.** Resolved-velocity position-based visual servoing with the following data:

$$\mathbf{K} = \text{diag}\{1, 1, 1, 2\},$$

corresponding to a time constant of 1 s for the three position variables and of 0.5 s for the orientation variable.

- C.** Image-based visual servoing of PD type with gravity compensation with the following data:

$$\mathbf{K}_{Ps} = 300\mathbf{I}_4 \quad \mathbf{K}_{Ds} = 330\mathbf{I}_4.$$

- D.** Resolved-velocity image-based visual servoing with the following data:

$$\mathbf{K}_s = \mathbf{I}_4,$$

corresponding to a time constant of 1 s for the feature parameters.

For the simulation of resolved-velocity control schemes, the dynamics of the velocity controlled manipulator has been neglected. Therefore, a pure kinematic model has been considered, based on the analytical Jacobian of the manipulator.

For position-based control schemes, the pose estimation algorithm based on the inverse of the image Jacobian has been adopted, with integration step $\Delta t = 1$ ms and gain $\mathbf{K}_s = 160\mathbf{I}_4$. As shown in Example 10.4, this implies that the algorithm converges in a time of about 0.03 s, which is lower than the sampling time of the control, as required for a correct operation of position-based control.

For image-based control schemes, matrix $\mathbf{L}_s(\mathbf{s}, \mathbf{z}_c)$ has been approximated with matrix $\hat{\mathbf{L}}_s = \mathbf{L}_s(\mathbf{s}, z_d)$, where z_d is the third component of vector $\mathbf{x}_{d,o}$.

The parameters of the various controllers have been chosen in such a way as to show the particular features of the different control laws and, at the same time, to allow a significant comparison of the performance of each scheme in response to congruent control actions. In particular, it can be observed what follows:

- The gains in schemes **A** and **C** have been tuned in simulation so as to obtain transient behaviors similar to those of schemes **B** and **D**.
- In control scheme **B** the gains of the position variables have been intentionally chosen equal to one another but different from the gain of the orientation variable to show that a desired dynamics can be imposed to each operational space variable.
- In control scheme **D** the gains have all been chosen equal to one another, since imposing different dynamics to different coordinates of the projections of the feature points on the image plane is not significant.

The results obtained with the various control schemes are illustrated in Figs. 10.18–10.25 in terms of:

- The time history of position and orientation of the camera frame with respect to the base frame and corresponding desired values (represented with dashed lines).
- The time history of feature parameters and corresponding desired values (represented with dashed lines).
- The path of feature points projections on the camera image plane, from initial positions (marked with crosses) to final positions (marked with circles).

Regarding performance of the various control schemes, the following considerations can be drawn from the obtained results.

In principle, if position-based visual servoing is adopted, a desired transient behaviour can be assigned to the operational space variables. This is only partially true for control scheme **A**, because the dynamics of the closed-loop system, for PD control with gravity compensation, is nonlinear and coupled. Therefore, the transient behaviour shown in Fig. 10.18 may be different if the manipulator starts from a different initial pose or has to reach a different desired pose. Vice versa, for control scheme **B**, the time history of operational space variables reported in Fig. 10.20 shows transient behaviours of exponential type, whose characteristics depend only on the choice of matrix \mathbf{K} .

For both schemes **A** and **B**, the trajectories of the projections of the feature points and the corresponding paths in the image plane (Figs. 10.19 and 10.21 respectively) have evolutions that cannot be predicted in advance. This implies that, although the feature points projections are inside the image plane both in the initial and in the desired configuration, they may exit from the image plane during the transient, thus causing problems of convergence to the controlled system.

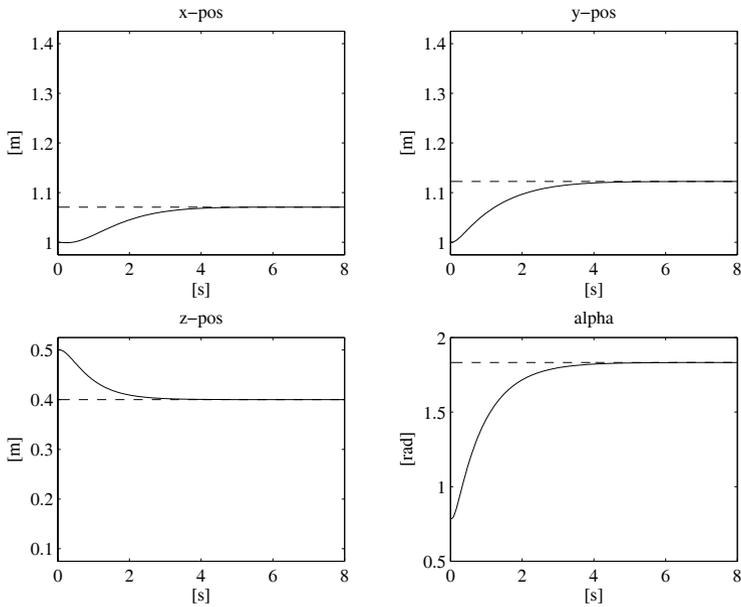


Fig. 10.18. Time history of camera frame position and orientation with control **A**

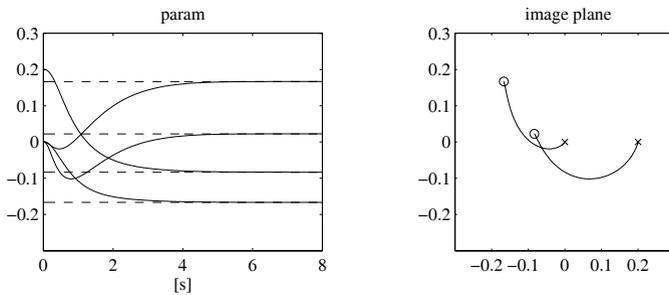


Fig. 10.19. Time history of feature parameters and corresponding path of feature points projections on image plane with control **A**

If image-based control is adopted, a desired transient behaviour can be assigned to the time histories of feature parameters and not to the operational space variables, in a dual fashion with respect to position-based control. This is confirmed by the results shown in Figs. 10.22–10.25, relative to control schemes **C** and **D**, respectively. In detail, especially in the case of control **C**, the time histories of the operational space variables are quite different from those reported in Figs. 10.18 and 10.20, despite the same initial and final configuration and a similar transient duration. This implies that the camera path, being unpredictable, may lead to joint limits violation or to collision of

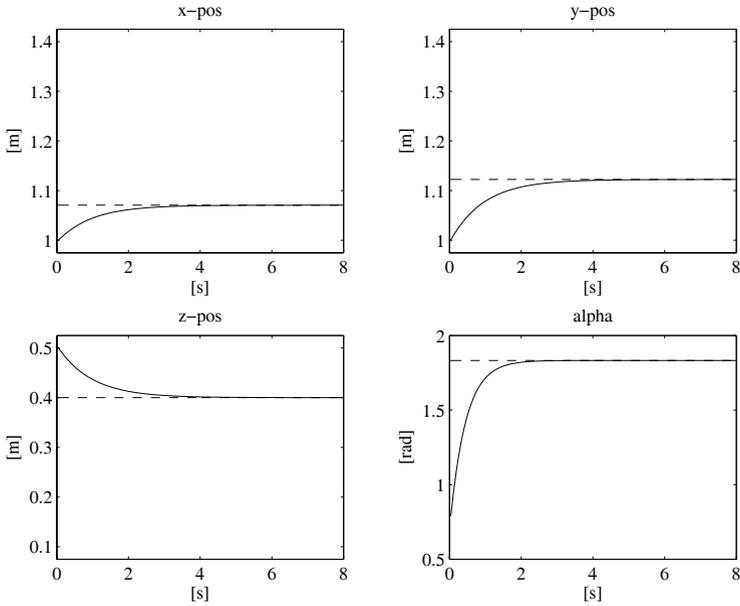


Fig. 10.20. Time history of camera frame position and orientation with control **B**

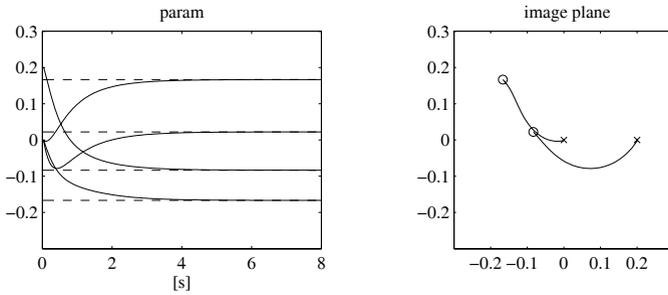


Fig. 10.21. Time history of feature parameters and corresponding path of feature points projections on image plane with control **B**

the manipulator with an obstacle. In the specific case of control scheme **C**, a 300% overshoot is present on the z component of the camera trajectory (see Fig. 10.22), which corresponds to a camera retreat movement with respect to the object of amplitude much higher than the distance reached at the end of the transient. The overshoot on the z component is present also for control scheme **D**, but is ‘only’ of 50% (see Fig. 10.24).

Notice that, for control scheme **C**, the presence of large displacements on some operational space variables does not correspond to significant deviations of the feature parameters with respect to their final values during the transient

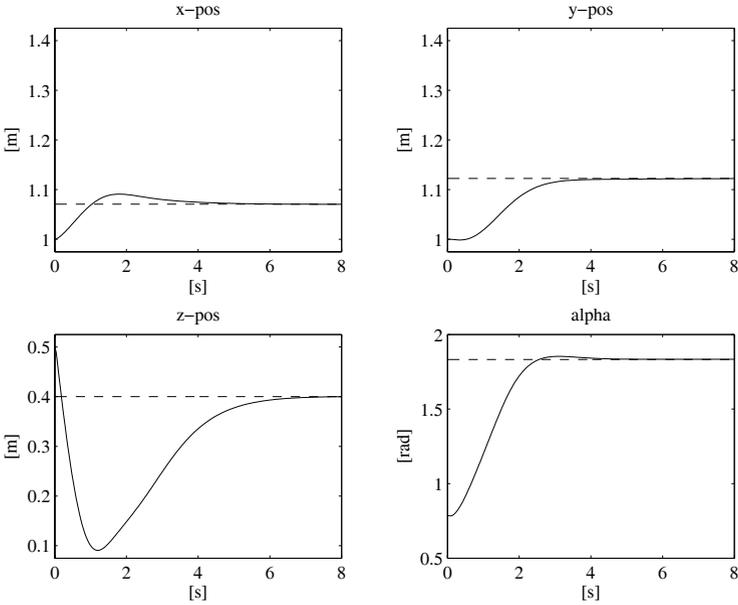


Fig. 10.22. Time history of camera frame position and orientation with control **C**

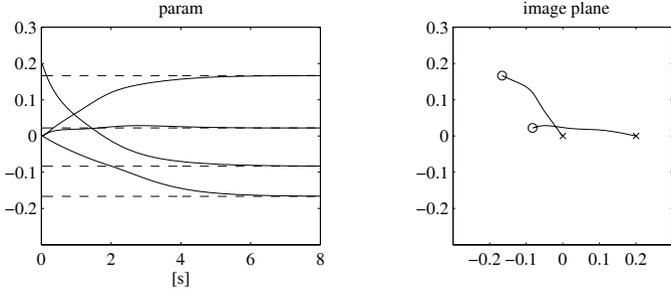


Fig. 10.23. Time history of feature parameters and corresponding path of feature points projections on image plane with control **C**

(see Fig. 10.23). Indeed, the paths of the feature points projections do not deviate much from the line segments connecting these points.

Figure 10.25, relative to control scheme **D**, reveals that the trajectories of the feature parameters are of exponential type. In this case, the transient behaviour depends only on matrix \mathbf{K}_s ; choosing a diagonal matrix with equal elements implies that the paths of the feature points projections are linear. In the case at hand, in view of the approximation $\mathbf{L}_s(\mathbf{s}, \mathbf{z}_c) \approx \mathbf{L}_s(\mathbf{s}, \mathbf{z}_d)$, the paths of the feature points projections shown in Fig. 10.25 are not perfectly linear.

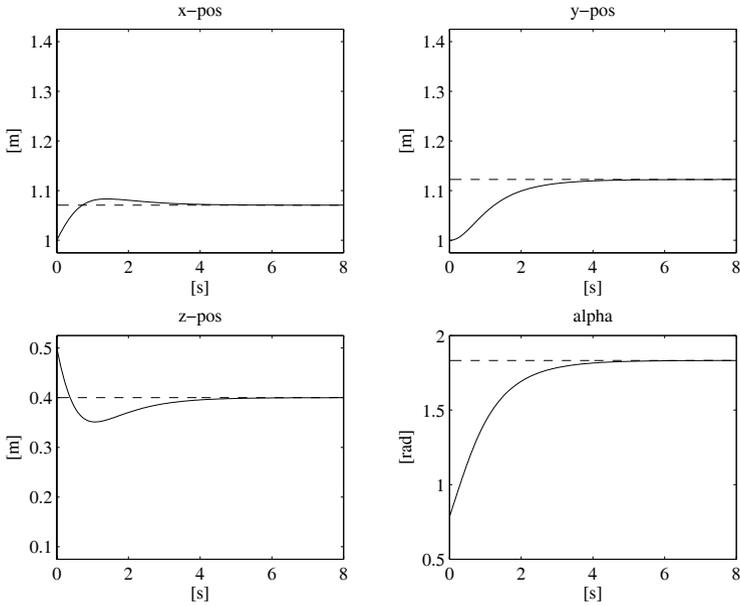


Fig. 10.24. Time history of camera frame position and orientation with control **D**

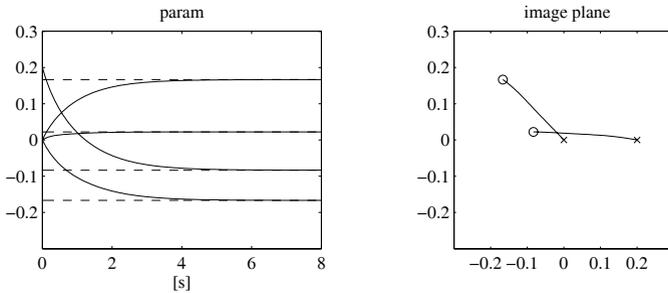


Fig. 10.25. Time history of feature parameters and corresponding path of feature points projections on image plane with control **D**

To conclude, Fig. 10.26 shows the paths of the origin of the camera frame obtained using the four control schemes. It can be observed that, with control scheme **B**, a perfectly linear path is obtained, thanks to the choice of a diagonal gain matrix \mathbf{K}_s with equal weights for the positional part. Using control scheme **A**, the path is almost linear, because, unlike case **B**, this type of control does not guarantee a decoupled dynamics for each operational space variable. Vice versa, using control schemes **C** and **D**, the path of the origin of the camera frame is far from being linear. In both cases, the phenomenon of camera retreat with respect to the object can be observed. To this end, notice

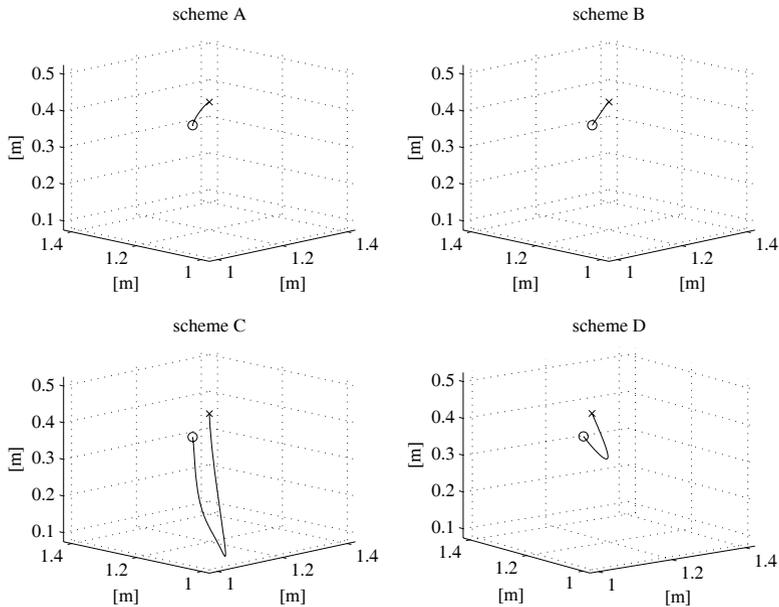


Fig. 10.26. Paths of camera frame origin in Cartesian space with the four control schemes

the axis z of the base frame of the SCARA manipulator and the axis z_c of the camera frame are aligned to the vertical direction and point downward; therefore, with respect to the reference frames of Fig. 10.26, the object is on the top and the camera points upward.

The phenomenon of camera retreat appears whenever the camera is required to perform large rotations about the optical axis; this phenomenon can be intuitively explained through a simple example. Assume that a pure rotation of the camera about the axis z_c is required and control scheme **D** is adopted. Therefore, visual servoing imposes that the feature points projections on the image plane follow rectilinear paths from the initial to the desired positions, whereas simple camera rotation would have required circular paths. The constraint on the path in the image plane implies that, during rotation, the origin of the camera frame must move, with respect to the object, first backward and then forward to reach again, asymptotically, the initial position. It can be shown that, if the desired rotation tends to π , then the distance of the camera from the object tends to ∞ and the system becomes unstable.

10.10 Hybrid Visual Servoing

An approach which combines the benefits of position-based and image-based visual servoing is *hybrid visual servoing*. The name stems from the fact that

the control error is defined in the operational space for some components and in the image space for the others. This implies that a desired motion can be specified, at least partially, in the operational space so that the camera trajectory during visual servoing can be predicted in advance for some components. On the other hand, the presence of error components in the image space helps keep the image features in the camera field of view, which is a difficult task in position-based approaches.

Hybrid visual servoing requires the estimation of some operational space variables. Assume that the object has a planar surface where at least four feature points, and no triplets of collinear points, can be selected. Using the coordinates of these points in the camera image plane, both in the current and in the desired pose of the camera frame, it is possible to compute the planar homography \mathbf{H} as described in Sect. 10.4.4. Notice that, for this computation, knowledge of the current and the desired camera pose is not required, provided that the feature vectors \mathbf{s} and \mathbf{s}_d are known.

In view of (10.59), assuming that Frame 1 coincides with the camera frame in the current pose and Frame 2 coincides with the camera frame in the desired pose, the following equality holds:

$$\mathbf{H} = \mathbf{R}_d^c + \frac{1}{d_d} \mathbf{o}_{c,d}^c \mathbf{n}^{dT},$$

where \mathbf{R}_d^c is the rotation matrix between the desired orientation and the current orientation of the camera frame, $\mathbf{o}_{c,d}^c$ is the position vector of the origin of the camera frame in the desired pose with respect to the current pose, \mathbf{n}^d is the unit vector normal to the plane containing the feature points, and d_d is the distance between this plane and the origin of the camera frame in the desired pose. The quantities \mathbf{R}_d^c , \mathbf{n}^d , $(1/d_d)\mathbf{o}_{c,d}^c$, in the current camera pose, can be computed at each sampling time from matrix \mathbf{H} .

Adopting a resolved-velocity approach, the control objective consists of computing the reference absolute velocity of the camera frame

$$\mathbf{v}_r^c = \begin{bmatrix} \nu_r^c \\ \boldsymbol{\omega}_r^c \end{bmatrix}$$

from a suitably defined error vector.

To this end, the orientation error between the desired and the current camera pose can be computed from matrix \mathbf{R}_d^c , as for position-based visual servoing. If $\boldsymbol{\phi}_{c,d}$ denotes the vector of the Euler angles extracted from \mathbf{R}_d^c , the control vector $\boldsymbol{\omega}_r^c$ can be chosen as

$$\boldsymbol{\omega}_r^c = -\mathbf{T}(\boldsymbol{\phi}_{c,d}) \mathbf{K}_o \boldsymbol{\phi}_{c,d}, \quad (10.92)$$

where \mathbf{K}_o is a (3×3) matrix. With this choice, the equation of the orientation error has the form

$$\dot{\boldsymbol{\phi}}_{c,d} + \mathbf{K}_o \boldsymbol{\phi}_{c,d} = \mathbf{0}. \quad (10.93)$$

Equation (10.93), if \mathbf{K}_o is a symmetric and positive definite matrix, implies that the orientation error tends to zero asymptotically with convergence of exponential type and speed depending on the eigenvalues of matrix \mathbf{K}_o .

The control vector $\boldsymbol{\nu}_r^c$ should be selected so that the positional part of the error between the desired and the current camera pose converges to zero. The position error could be defined as the difference of the coordinates of a point of the object in the desired camera frame $\mathbf{r}_d^c = [x_d \ y_d \ z_d]^T$ and those in the current camera frame $\mathbf{r}_c^c = [x_c \ y_c \ z_c]^T$, namely $\mathbf{r}_d^c - \mathbf{r}_c^c$. These coordinates, however, cannot be directly measured, unlike the corresponding coordinates in the image plane, defining the feature vectors $\mathbf{s}_{p,d} = [X_d \ Y_d]^T = [x_d/z_d \ y_d/z_d]^T$ and $\mathbf{s}_p = [X \ Y]^T = [x_c/z_c \ y_c/z_c]^T$.

The information deriving from the computation of homography \mathbf{H} can be used to rewrite the ratio

$$\rho_z = z_c/z_d$$

in terms of known or measurable quantities in the form

$$\rho_z = \frac{d_c}{d_d} \frac{\mathbf{n}^{dT} \tilde{\mathbf{s}}_{p,d}}{\mathbf{n}^{cT} \tilde{\mathbf{s}}_p} \quad (10.94)$$

with

$$\frac{d_c}{d_d} = 1 + \mathbf{n}^{cT} \frac{\mathbf{O}_{c,d}^c}{d_d} = \det(\mathbf{H}), \quad (10.95)$$

and $\mathbf{n}^c = \mathbf{R}_d^c \mathbf{n}^d$, where vectors $\tilde{\mathbf{s}}_p$ and $\tilde{\mathbf{s}}_{p,d}$ denote the representations in homogeneous coordinates of \mathbf{s}_p and $\mathbf{s}_{p,d}$, respectively (see Problem 10.12).

The position error, expressed in terms of known or measurable quantities, can be defined as

$$\mathbf{e}_p(\mathbf{r}_d^c, \mathbf{r}_c^c) = \begin{bmatrix} X_d - X \\ Y_d - Y \\ \ln \rho_z \end{bmatrix}.$$

Notice that, in view of (5.44), convergence to zero of \mathbf{e}_p implies convergence to zero of $\mathbf{r}_d^c - \mathbf{r}_c^c$ and vice versa.

Computing the time derivative of \mathbf{e}_p yields

$$\dot{\mathbf{e}}_p = \frac{\partial \mathbf{e}_p(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} \dot{\mathbf{r}}_c^c,$$

\mathbf{r}_d^c being constant. By taking into account (10.26) and the decomposition

$$\mathbf{v}_c^c = \begin{bmatrix} \boldsymbol{\nu}_c^c \\ \boldsymbol{\omega}_c^c \end{bmatrix}$$

with $\boldsymbol{\nu}_c^c = \mathbf{R}_c^T \dot{\mathbf{o}}_c$, the above expression can be rewritten in the form

$$\dot{\mathbf{e}}_p = -\mathbf{J}_p \boldsymbol{\nu}_c^c - \mathbf{J}_o \boldsymbol{\omega}_c^c, \quad (10.96)$$

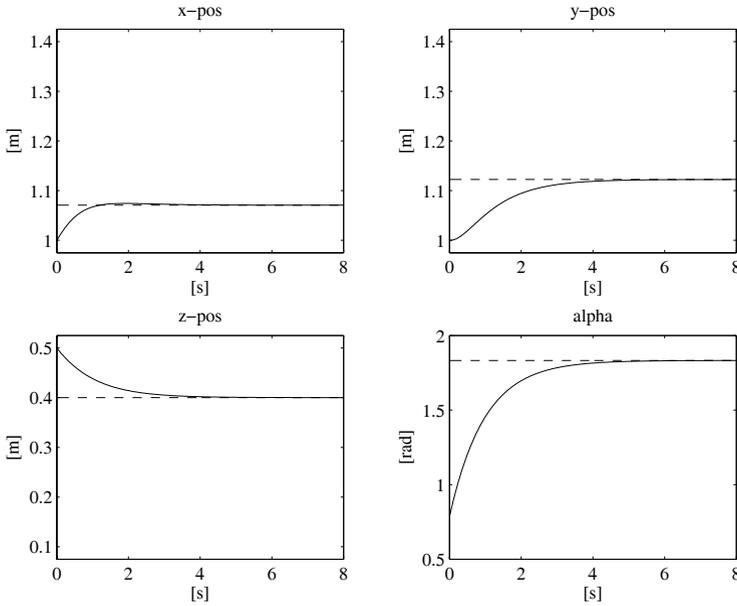


Fig. 10.27. Time history of position and orientation of camera frame with hybrid visual servoing

with (see Problem 10.13)

$$\mathbf{J}_p = \frac{1}{z_d \rho_z} \begin{bmatrix} -1 & 0 & X \\ 0 & -1 & Y \\ 0 & 0 & -1 \end{bmatrix}$$

and

$$\mathbf{J}_o = \begin{bmatrix} XY & -1 - X^2 & Y \\ 1 + Y^2 & -XY & -X \\ -Y & X & 0 \end{bmatrix}.$$

Equation (10.96) suggests the following choice of control vector $\boldsymbol{\nu}_r^c$

$$\boldsymbol{\nu}_r^c = \mathbf{J}_p^{-1}(\mathbf{K}_p \mathbf{e}_p - \mathbf{J}_o \boldsymbol{\omega}_r^c), \quad (10.97)$$

\mathbf{J}_p being a nonsingular matrix.

Notice that, for the computation of \mathbf{J}_p^{-1} , knowledge of the constant quantity z_d is required.

If z_d is known, control law (10.97), in view of assumptions $\dot{\boldsymbol{o}}_c^c \approx \boldsymbol{\nu}_r^c$ and $\boldsymbol{\omega}_c^c \approx \boldsymbol{\omega}_r^c$, yields the following error equation:

$$\dot{\mathbf{e}}_p + \mathbf{K}_p \mathbf{e}_p = \mathbf{0},$$

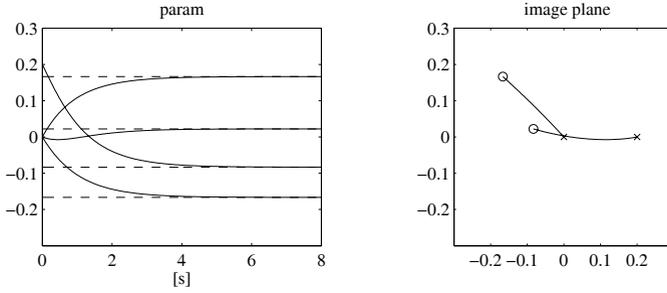


Fig. 10.28. Time history of feature parameters and corresponding path of feature points projections on image plane with hybrid visual servoing

which implies the exponential convergence of \mathbf{e}_p to zero, provided that \mathbf{K}_p is a positive definite matrix.

If z_d is not known, an estimate \hat{z}_d can be adopted. Therefore, in control law (10.97), matrix \mathbf{J}_p^{-1} can be replaced by an estimate $\hat{\mathbf{J}}_p^{-1}$. In view of the equality

$$\hat{\mathbf{J}}_p^{-1} = \frac{\hat{z}_d}{z_d} \mathbf{J}_p^{-1},$$

the following error equation is obtained:

$$\dot{\mathbf{e}}_p + \frac{\hat{z}_d}{z_d} \mathbf{K}_p \mathbf{e}_p = \left(1 - \frac{\hat{z}_d}{z_d} \right) \mathbf{J}_o \boldsymbol{\omega}_r^c.$$

This equation shows that the use of an estimate \hat{z}_d in place of the true value z_d implies a simple gain scaling in the error equation, and asymptotic stability is preserved. Moreover, due to the presence in the right-hand side of the above equation of a term depending on $\boldsymbol{\omega}_r^c$, the time history of \mathbf{e}_p is influenced by the orientation error, which evolves according to (10.93).

Example 10.6

For the SCARA manipulator and the task of Sect. 10.9, consider the hybrid visual servoing law with gains

$$\mathbf{K}_p = \mathbf{I}_3 \quad k_o = 2,$$

and compute the positional part of the error with respect to point P_1 . The planar homography and the corresponding parameters are estimated as in Example 10.5, using four points. The results are reported in Figs. 10.27 and 10.28, in terms of the same variables shown in Fig. 10.27. Notice that the time histories of the variables in the operational space of Fig. 10.27 are quite similar to that obtained with resolved-velocity position-based visual servoing (Fig. 10.20). On the other hand, the time histories of the feature parameters of Fig. 10.28 are substantially similar to those obtained with resolved-velocity image-based visual servoing (Fig. 10.25) except for

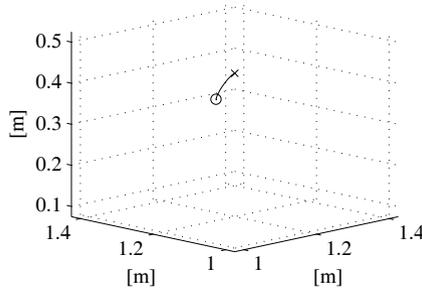


Fig. 10.29. Path of camera frame origin in Cartesian space with hybrid visual servoing

the fact that the path in the image plane of the projection of point P_1 is perfectly linear, as imposed by the control. The corresponding path of the camera frame origin, reported in Fig. 10.29, shows a substantial improvement with respect to those obtained with the image-based visual servoing schemes of Fig. 10.26.

The method illustrated above is only one of the possible visual servoing approaches based on the computation of the planar homography and on its decomposition. It is worth pointing out that knowledge of $(1/d_d)\sigma_{c,d}^c$ and \mathbf{R}_d^c allows the computation of the operational space error (10.69), up to a scaling factor for the positional part; therefore, it is also possible to use the position-based visual servoing schemes presented in Sect. 10.7. On the other hand, in hybrid visual servoing approaches, different choices are possible for the error components depending on feature parameters as well as for those depending on the operational space variables.

Bibliography

The literature dealing with computational vision is quite extensive and various. Image processing is treated, e.g., in [93], while geometrical issues are considered in [75] and [146]. The concept of interaction matrix was originally proposed in [239] under a different name, while the actual name was introduced in [71]. The pose estimation algorithm based on the inverse of the image Jacobian is also termed virtual visual servoing [47]. Position-based visual servoing was proposed in [239] and, more recently, in [244] and [134]. Several papers dealing with image-based visual servoing can be found, starting from early works of [239] and [79]. A rigorous stability analysis is reported in [108] for PD control with gravity compensation and in [71] for resolved-velocity control. Hybrid visual servoing, presented in [148], is only one of advanced control schemes based on decoupling of the camera DOFs, e.g., the partitioned approach [49] and the control based on image moments [35]. Finally,

visual servoing based on stereo vision is considered in [92]. An interesting review of the state of the art of visual servoing until mid 1990s is presented in [103], while a more recent review can be found in [36].

Problems

10.1. Derive Eq. (10.2).

10.2. Show that a non-null solution to (10.10) is the right eigenvector corresponding to the null singular value of matrix \mathbf{A} .

10.3. Show that (10.12) is the matrix which minimizes Frobenius norm (10.11) with the constraint that \mathbf{R}_o^c is a rotation matrix, in the case $\sigma > 0$. [*Hint*: consider that $\text{Tr}(\mathbf{R}_o^{cT} \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T) = \text{Tr}(\mathbf{V}^T \mathbf{R}_o^{cT} \mathbf{U} \boldsymbol{\Sigma})$; moreover, the absolute values of the diagonal elements of matrix $\mathbf{V}^T \mathbf{R}_o^{cT} \mathbf{U}$ are lower or equal to 1.]

10.4. Consider the SCARA manipulator of Example 10.3. Perform a computer implementation of the pose estimation algorithm based on the inverse of the image Jacobian considering the points P_1 and P_2 of the object of Example 10.1. Compute the homogeneous transformation matrix corresponding to the feature vector $\mathbf{s} = [-0.1667 \quad 0.1667 \quad -0.0833 \quad 0.0223]^T$. Assume that, in the initial pose, the axes of the camera frame are parallel to those of the object frame and the origin is at a distance of 0.5 m along the vertical axis.

10.5. Solve the previous problem using the feature parameters of the line segment $P_1 P_2$.

10.6. Show that the value of \mathbf{o} which minimizes function (10.55) has expression (10.56). [*Hint*: Use the equalities $\mathbf{p}_i = \bar{\mathbf{p}}_i + \bar{\mathbf{p}}$ and $\mathbf{p}'_i = \bar{\mathbf{p}}'_i + \bar{\mathbf{p}}'$ in (10.55), and the properties $\|\mathbf{a} + \mathbf{b}\|^2 = \|\mathbf{a}\|^2 + 2\mathbf{a}^T \mathbf{b} + \|\mathbf{b}\|^2$ and $\sum_{i=1}^n \bar{\mathbf{p}}_i = \sum_{i=1}^n \bar{\mathbf{p}}'_i = \mathbf{0}$.]

10.7. Show that the matrix \mathbf{R} which minimizes function (10.57) is the matrix which maximizes the trace of $\mathbf{R}^T \mathbf{K}$.

10.8. For the SCARA manipulator of Example 10.3, design a position-based visual servoing scheme of PD type with gravity compensation using the measurement of the feature parameters of line segment $P_1 P_2$ of Example 10.1. Perform a computer simulation in the same operating conditions of Sect. 10.9 and compare the results.

10.9. For the SCARA manipulator of Example 10.3, design a resolved-velocity position-based visual servoing scheme using the measurement of the feature parameters of line segment $P_1 P_2$ of Example 10.1. Perform a computer simulation in the same operating conditions of Sect. 10.9 and compare the results.

10.10. For the SCARA manipulator of Example 10.3, design an image-based visual servoing scheme of PD type with gravity compensation using the measurement of the feature parameters of the line segment P_1P_2 of Example 10.1. Perform a computer simulation in the same operating conditions of Sect. 10.9 and compare the results.

10.11. For the SCARA manipulator of Example 10.3, design a resolved-velocity image-based visual servoing scheme using the measurement of the feature parameters of line segment P_1P_2 of Example 10.1. Perform a computer simulation in the same operating conditions of Sect. 10.9 and compare the results.

10.12. Derive the expressions (10.94), (10.95).

10.13. Derive the expressions of \mathbf{J}_p and \mathbf{J}_o in (10.96).