

A Performance Criterion for the Depth Estimation with Application to Robot Visual Servo Control

Chang-Jia Fang,* Shir-Kuan Lin†

*Institute of Electrical and Control Engineering
National Chiao Tung University
Hsinchu 30010, Taiwan*

Received 23 February 2001; accepted 18 April 2001

This article deals with the depth observability problem of a robot visual system with a moving camera. In the visual system, the unknown depth of a feature point is estimated from the input of the camera velocity and the output of the image of the feature point. Although it is well known that the linear velocity of the camera must satisfy some constraints for successful depth estimation, this proposes a criterion to measure the performance of the depth estimation, which is a heuristic extension from an estimation result of a linear system. This performance criterion depends on both the image position and the linear velocity of the camera. Some simulation and experiment examples demonstrate and verify the proposed performance criterion. Furthermore, this criterion is used to develop a new visual servo control scheme that has good performance in both the depth estimation and the visual control. This control scheme is also verified by a simulation example. © 2001 John Wiley & Sons, Inc.

1. INTRODUCTION

Depth estimation is a very important problem for three-dimensional (3-D) vision application, e.g., in object tracking,¹⁻³ motion estimation,⁴⁻⁷ and obstacle detection. Unfortunately, the value of depth cannot be directly measured. In the machine vision field, the stereo vision system extracts the 3-D space data from multiview images by matching the fea-

tures in different images.⁹⁻¹¹ Some works research the stereo visual system.¹² However, the stereo vision system is more expensive and the feature-corresponding problem needs to be solved. Many hand-eye robot systems still employ simple active vision systems that estimate the depth from a sequence of images through the motion of the camera.¹³ A typical estimator in this system is the extended Kalman filter (EKF).¹⁴

Nowadays, visual sensors are widely used in the motion control of the robot manipulator. Visual servo robot control overcomes the difficulties of uncertain models and unknown environments and broadens the domain of application of current

Contract grant sponsor: National Science Council, Taiwan.

Contract number: NSC89-2213-E0009-216.

*To whom all correspondence should be addressed; e-mail: conga.ece83g@nctu.edu.tw.

†e-mail: sklin@cc.nctu.edu.tw.

robots. In the literature, there are two types of visual servo controllers: one is the feature-based method^{10,15–20}; the other is the position-based method.^{21,22} The main distinction is the input. The former uses the image features as an input command, while the latter takes the position in space calculated from the images as an input. The feature-based control has the input command described directly in the feature space; it is then easy to generate the input trajectory by video-aid, computer-aided design, or virtual reality techniques. Therefore, most of visual servo controllers are feature-based.

Jang and Bien¹⁸ investigated the mathematical descriptions of various image features and introduced the redundant features to improve the control performance for noise. Corke and Paul¹⁹ designed a Single-Input Single-Output featured-based controller for each degree of freedom of the camera motion to track some particular objects. Weiss et al.¹⁶ first proposed the feature Jacobian matrix that establishes the relationship between the differential change in the feature vector and that in the relative pose of the camera. However, the singularity of the feature Jacobian matrix is an inherent problem in the feature-based control.

To alleviate the singularity problem, some authors use the singular value decomposition method to determine the input of the camera velocity.²³ Feddema et al.²⁴ deal with selection and weighting of features for the condition and the sensitivity of the feature Jacobian matrix. Furthermore, a trade-off of measure of lens focus, field of view, robot configuration, and resolvability is also introduced into the visual control object for special tasks.²⁵ In these works, good performance in several visual servo control schemes was reported. Until now, the unknown depths were in the feature Jacobian matrix of the feature-based control; the depth estimation is needed in visual servo control. In some cases, the exact value of depth may not be necessary for the stability of the visual servo control scheme.³ The depth is still important information for handling the object. The depth observability is the vital problem in the depth estimation, which determines the success of the estimation. Although there have been several methods proposed for the depth estimation,^{24,26} the popular one is the EKF.^{13,22,27} However, little attention has been paid to the effect of the velocity of the camera on the depth observability. The depth estimation incorporated in visual servo control has not been deeply investigated.

Recently, the observability problem of the 3-D structure from motion of a visual system was discussed.²⁸ Two different initial states may be undistinguishable by the scale ambiguity or by a special kind of translational velocity. The effects of the velocity of the camera on the observability are investigated. Dayawansa et al.²⁹ proposed a necessary and sufficient condition for the perspective observability problem. This condition is derived from the generalized Popov–Belevitch–Hautus test for the camera with a constant velocity. Sharma and Hutchinson^{23,30} proposed a measure of motion perceptibility and some cost functions to determine the optimal camera placement (or trajectory) for the observability of robot motion.

This article considers a hand-eye robot system. The feature-based visual servo control is used to move the end-effector to the desired position relative to an object by inputting the image features of the object. In contrast to earlier works, this article tries to improve the depth estimation incorporated in visual servo control. According to previous research on the observability of 3-D structure from visual motion,²⁸ a necessary and sufficient condition of depth observability for a moving camera is found, which states what types of camera velocities can make the trajectory of the depth different for different initial depths. Although this condition can make the depth estimation successful, it provides no information about the convergence rate of the depth estimation. If the nonlinear depth estimator could be virtually seen as a linear unbiased estimator, the variance of the reciprocal of the depth estimate could be explicitly found. We then suggest use of this variance as an index for the performance of the depth estimation. This index is also verified by simulations and experiments as a rule of thumb for the performance evaluation of the EKF, especially for slow camera velocity. We then try to develop a new visual servo control scheme by making the index as small as possible while having little effect on the control performance, so that the depth estimation is improved. Finally, a simulation example shows that the resulting control scheme reaches this goal.

This article is organized as follows. Section 2 introduces the camera model and discusses the observability of the nonlinear optic flow model. The result of Section 2 is used to develop a new visual servo control scheme in Section 3, which tries to improve the depth estimation during control. Finally, we conclude our results in Section 4.

2. DEPTH ESTIMATION

Consider a pinhole camera model and assign a body-fixed coordinate frame E_{XYZ} on it. The origin of the camera frame E_{XYZ} is at the center of the camera lens and the Z -axis of E_{XYZ} is in alignment with the optical axis of the lens. Another frame is fixed on the image plane and denoted by E_{xyz} whose x - and y -axes are parallel to the X - and Y -axes of E_{XYZ} and whose origin is at the intersection point of the image plane and the optical axis. The distance between the origin of E_{XYZ} and that of origin of E_{xyz} is named the effective focal length and is denoted by f_e .

Consider a point P (e.g., a corner of a 3-D object) with coordinates (X, Y, Z) with respect to the camera frame E_{XYZ} . The value of Z is referred to as the depth of point P to the camera lens. The image of point P projected onto the image plane is denoted by p with coordinates $(x, y, 0)$ with respect to E_{xyz} . The projection relationships state

$$x = \gamma_x \frac{X}{Z}, \quad y = \gamma_y \frac{Y}{Z} \quad (1)$$

where $\gamma_x = f_e/S_x$ and $\gamma_y = f_e/S_y$, in which S_x, S_y , are, respectively, the horizontal and vertical lengths per pixel on the camera sensing array. Equation (1) is the so-called perspective projection equation.³¹

In an active vision system, the camera is capable of moving. Assume that the linear velocity and the angular velocity of the camera are \mathbf{v} and $\boldsymbol{\omega}$, respectively, with respect to E_{XYZ} . For a stationary point P , the relative velocity of point P with respect to the camera is

$$\begin{bmatrix} \dot{X} \\ \dot{Y} \\ \dot{Z} \end{bmatrix} = -\boldsymbol{\omega}(t) \times \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} - \mathbf{v}(t) \quad (2)$$

Let $\boldsymbol{\xi} \triangleq [x, y, Z]^T$ and $\mathbf{u} \triangleq [\mathbf{v}^T, \boldsymbol{\omega}^T]^T$. Differentiating (1) and using (2) yields the so-called optic flow-motion equation³¹:

$$\begin{bmatrix} \dot{x} \\ \dot{y} \end{bmatrix} = \bar{\mathbf{J}}(\boldsymbol{\xi})\mathbf{u} \quad (3)$$

where

$$\bar{\mathbf{J}}(\boldsymbol{\xi}) = \begin{bmatrix} -\frac{\gamma_x}{Z} & 0 & \frac{x}{Z} & \frac{xy}{\gamma_y} & -\frac{\gamma_x^2 + x^2}{\gamma_x} & \frac{\gamma_x}{\gamma_y} y \\ 0 & -\frac{\gamma_y}{Z} & \frac{y}{Z} & \frac{\gamma_y^2 + y^2}{\gamma_y} & -\frac{xy}{\gamma_x} & -\frac{\gamma_y}{\gamma_x} x \end{bmatrix} \quad (4)$$

Since the unknown depth Z is involved in (3), the dynamic equation of Z in (2) needs to be considered together. Note that (3) is physically valid only when $Z > 0$.

According to (3), we describe the present system by the nonlinear system of equations

$$\dot{\boldsymbol{\xi}} = \mathbf{f}(\boldsymbol{\xi}, \mathbf{u}) \triangleq \mathbf{G}(\boldsymbol{\xi})\mathbf{u} \quad (5)$$

$$\boldsymbol{\psi} = \mathbf{h}(\boldsymbol{\xi}) \triangleq \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \boldsymbol{\xi} \quad (6)$$

where

$$\mathbf{G}(\boldsymbol{\xi}) = \begin{bmatrix} & \mathbf{J}(\boldsymbol{\xi}) & & & & \\ 0 & 0 & -1 & -\frac{yZ}{\gamma_y} & \frac{xZ}{\gamma_x} & 0 \end{bmatrix} \quad (7)$$

Equations (5) and (6) are the state and output equations, respectively. In the state vector $\boldsymbol{\xi}$, x , and y are the visual measurements, so the only unknown

is the depth Z . Therefore, the state estimation of the system of (5) and (6) is equivalent to the depth estimation for the given camera velocity.

Let the linear velocity of the camera be decomposed into

$$\mathbf{v}(t) = \alpha_1(t)\mathbf{v}_1(\xi) + \alpha_2(t)\mathbf{v}_2(\xi) + \alpha_3(t)\mathbf{v}_3(\xi) \quad (8)$$

where $\alpha_1(t)$, $\alpha_2(t)$, and $\alpha_3(t)$ are bounded functions of t , and

$$\mathbf{v}_1(\xi) = \begin{bmatrix} -\gamma_x \\ x \end{bmatrix}, \quad \mathbf{v}_2(\xi) = \begin{bmatrix} 0 \\ -\gamma_y \\ y \end{bmatrix}, \quad \mathbf{v}_3(\xi) = \begin{bmatrix} x \\ \gamma_x \\ y \\ \gamma_y \\ 1 \end{bmatrix} \quad (9)$$

Note that the dimension of the distribution span $\{\mathbf{v}_1(\xi), \mathbf{v}_2(\xi), \mathbf{v}_3(\xi)\}$ is always 3. It is well known that the unknown depth of the point cannot be estimated when its image is on the focus of expansion,²⁸ i.e., the camera always moves along the direction of $\mathbf{v}_3(\xi)$. The system (5) and (6) is locally observable for $Z > 0$, if and only if $\alpha_1^2(t) + \alpha_2^2(t)$ is nonzero for some time interval [i.e., $\mathbf{v}(t) \neq \alpha_3(t)\mathbf{v}_3(\xi)$], for some time interval). This restriction on the input is a necessary and sufficient condition for the success of the depth estimation. However, there is no information about the convergent rate of the depth estimation error. It is then interesting to know if there is a relationship of \mathbf{v}_1 and \mathbf{v}_2 to the convergence of the depth estimation.

In this article, the nonlinear time-varying system of (5) and (6) is estimated by the extended Kalman filter. Since the image data x and y are available, their estimation errors could be ignored. If the variation rate of the depth Z is also small enough to be negligible, then the extended Kalman filter (it is a nonlinear filter) can be reduced to the linear least-squares estimator. Thus, the performance of the linear least-squares estimate can be determined by the corresponding index function if only Gaussian noises are considered.

2.1. Performance Index for Depth Estimation

Consider the linear stochastic system

$$\mathbf{y}_o = \Phi\theta + \rho \quad (10)$$

where $\mathbf{y}_o \in \mathcal{R}^m$ is composed of the measurable signals, $\Phi \in \mathcal{R}^{m \times n}$ is known, $\theta \in \mathcal{R}^n$ is a parameter vector to be estimated, and $\rho \in \mathcal{R}^m$ is a zero-mean random vector with variance matrix \mathbf{S} . The best unbiased linear estimator for θ in (10) is $\hat{\theta}^* = \mathbf{S}^{-1}\Phi(\Phi^T\mathbf{S}^{-1}\Phi)^{-1}\mathbf{y}_o$.³² The corresponding covariance of $\hat{\theta}^*$ is then

$$\text{Cov}(\hat{\theta}^*) = (\Phi^T\mathbf{S}^{-1}\Phi)^{-1} \quad (11)$$

Since $\text{Tr}[\text{Cov}(\hat{\theta})] \triangleq E[(\theta - \hat{\theta})^T(\theta - \hat{\theta})]$, where $\text{Tr}(\cdot)$ is a trace operator, $\text{Tr}[\text{Cov}(\hat{\theta})]$ is often used as an accuracy index of the parameter estimation.

Now consider (3) and reform it in the form of (10) as

$$\begin{bmatrix} \dot{x} - a_x(\xi, \omega) \\ \dot{y} - a_y(\xi, \omega) \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1 \cdot \mathbf{v} \\ \mathbf{v}_2 \cdot \mathbf{v} \end{bmatrix} \frac{1}{Z} + \begin{bmatrix} \rho_1 \\ \rho_2 \end{bmatrix} \quad (12)$$

where $a_x(\xi, \omega) = (xy/\gamma_y)\omega_x - [(\gamma_x^2 + x^2)/\gamma_x]\omega_y + (\gamma_x y/\gamma_y)\omega_z$, $a_y(\xi, \omega) = [(\gamma_y^2 + y^2)/\gamma_y]\omega_x - (xy/\gamma_x)\omega_y - (\gamma_y x/\gamma_x)\omega_z$, and ρ_1 and ρ_2 are two zero-mean Gaussian noise terms. The variance matrix \mathbf{S} is then $\text{diag}(q_{11}, q_{22})$, where q_{ii} is the variance of ρ_i , $i = 1, 2$. In this case, $m = 2$, $n = 1$, $\hat{\theta} = 1/\hat{Z}$, $\mathbf{y}_o = [\dot{x} - a_x(\xi, \omega), \dot{y} - a_y(\xi, \omega)]^T$, $\Phi = [\mathbf{v} \cdot \mathbf{v}_1, \mathbf{v} \cdot \mathbf{v}_2]^T$, and $\rho = [\rho_1, \rho_2]^T$. It then follows from (11) that

$$\text{Cov}(1/\hat{Z}) = [\mathbf{v}^T\mathbf{A}(\xi)\mathbf{v}]^{-1} \quad (13)$$

where

$$\mathbf{A}(\xi) = \frac{\mathbf{v}_1(\xi)\mathbf{v}_1^T(\xi)}{q_{11}} + \frac{\mathbf{v}_2(\xi)\mathbf{v}_2^T(\xi)}{q_{22}} \quad (14)$$

Although $\text{Cov}(1/\hat{Z}) \neq \text{Cov}(\hat{Z})$, the accuracy of $1/\hat{Z}$ also reflects that of \hat{Z} . Since a smaller $\text{Cov}(1/\hat{Z})$ means a more accurate estimate $1/\hat{Z}$, we could then guess that this index function could still be a performance criterion of the extended Kalman filter for estimating the states of the nonlinear system of (5) and (6), when the variation rate Z is very small.

Proposition 1: Consider the system of (5) and (6) with input $\mathbf{u} = [\mathbf{v}^T, \omega^T]^T$. Suppose that the variation rate of the depth is small enough that the nonlinear depth estimator (such as the extended Kalman filter) can be approximated as a linear least-squares estimator. When the norm of the linear velocity of the camera $\|\mathbf{v}(t)\|$ is

fixed, a larger $\mathcal{F}_0(\boldsymbol{\xi}, \mathbf{v})$ guarantees a faster convergent rate of the depth estimation, where

$$\mathcal{F}_0(\boldsymbol{\xi}, \mathbf{v}) = \mathbf{v}^T \mathbf{A}(\boldsymbol{\xi}) \mathbf{v} \quad (15)$$

in which $\mathbf{A}(\boldsymbol{\xi})$ is defined in (14). Since \mathbf{v}_3 is orthogonal to \mathbf{v}_1 and \mathbf{v}_2 , \mathbf{v}_3 is then in the null space of $\mathbf{A}(\boldsymbol{\xi})$.

The two examples in the next subsection will show that it is acceptable for slow camera velocity.

2.2. Simulations and Experiments

According to the form of $\mathbf{v}_1(\boldsymbol{\xi})$, $\mathbf{v}_2(\boldsymbol{\xi})$, and $\mathbf{v}_3(\boldsymbol{\xi})$ in (9), we establish a distribution spanned by the orthonormal basis $\{\bar{\mathbf{b}}_1, \dots, \bar{\mathbf{b}}_6\}(\boldsymbol{\xi})$, where $\bar{\mathbf{b}}_i = \mathbf{b}_i / \|\mathbf{b}_i\|$ and

$$\begin{aligned} \mathbf{b}_1 &= \begin{bmatrix} -\gamma_x \\ 0 \\ x \\ 0 \\ 0 \\ 0 \end{bmatrix}, \\ \mathbf{b}_2 &= \begin{bmatrix} \mathbf{v}_3(\boldsymbol{\xi}) \times \mathbf{v}_1(\boldsymbol{\xi}) \\ 0 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{xy}{\gamma_y} \\ \gamma_y \\ -(\gamma_x^2 + x^2) \\ \gamma_x \\ \frac{\gamma_x}{\gamma_y} y \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad (16) \\ \mathbf{b}_3 &= \begin{bmatrix} \frac{x}{\gamma_x} \\ \gamma_x \\ \frac{y}{\gamma_y} \\ \gamma_y \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \\ \mathbf{b}_4 &= \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{b}_5 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{b}_6 = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} \quad (17) \end{aligned}$$

Although the basis varies, it is independent of the unknown depth Z . With the basis, we can express the camera velocity as

$$\mathbf{u}(t) = \sum_{i=1}^6 c_i(t) \bar{\mathbf{b}}_i(\boldsymbol{\xi}) = [\bar{\mathbf{b}}_1, \dots, \bar{\mathbf{b}}_6] \cdot [\mathbf{u}(t)]_{\mathbf{b}} \quad (18)$$

where $[\mathbf{u}(t)]_{\mathbf{b}} \triangleq [c_1(t), \dots, c_6(t)]^T$ is the representation of $\mathbf{u}(t)$ with respect to the basis $\{\bar{\mathbf{b}}_1, \dots, \bar{\mathbf{b}}_6\}(\boldsymbol{\xi})$. It is apparent that $c_1^2(t) + c_2^2(t) = 0$ is equivalent to $\mathcal{F}_0(\boldsymbol{\xi}, \mathbf{v}) = 0$, which occurs when \mathbf{v} is parallel to $\mathbf{v}_3(\boldsymbol{\xi})$.

2.2.1. Simulations

The model used in the first example is the visual system (5) and (6) with the effective focal length $f_e = 16.53$ mm, the horizontal length per pixel $S_x = 0.0161$ mm/pixel, and the vertical length per pixel $S_y = 0.0189$ mm/pixel. The image sampling period is 200 ms, the measurement noise covariance matrix $\mathbf{R} = \text{diag}(2.25 \text{ pixel}^2, 2.25 \text{ pixel}^2)$, and the system noise covariance matrix $\mathbf{Q} = \text{diag}(2.25 \text{ pixel}^2, 2.25 \text{ pixel}^2, 16 \text{ mm}^2)$. The states are estimated by the extended Kalman filter. The differential Eqs. (5) and (6) are numerically solved by the 4th order Adams–Bashforth method. Note that the velocities are expressed as (18). The units of the first three components of $\mathbf{u}(t)$ are mm/s and those of the others are rad/s.

Suppose that an interesting feature point on the screen is initially located at (100, 100). The guessed value of the initial depth of the point is 400 mm, while its true depth is 450 mm; i.e., the initial depth estimation error is 50 mm.

According to Proposition 1, it is recommended that \mathbf{v} be chosen such that $\mathcal{F}_0(\boldsymbol{\xi}, \mathbf{v})$ is as large as possible. Clearly, the value of the index (15) depends on the magnitude and the direction of \mathbf{v} , the image coordinates (x, y) , and the property of the system noise. For the convenience of comparison, we only consider $\|\mathbf{v}\|$ to be constant.

Consider the input in the form of

$$\mathbf{u} = \|\mathbf{v}\| \cdot [C_\beta (C_\gamma \bar{\mathbf{b}}_1 + S_\gamma \bar{\mathbf{b}}_2) + S_\beta \bar{\mathbf{b}}_3] + [\boldsymbol{\omega}]_{\mathbf{b}}, \quad (19)$$

where $S_\beta \triangleq \sin(\beta)$, $C_\beta \triangleq \cos(\beta)$, $S_\gamma \triangleq \sin(\gamma)$, $C_\gamma \triangleq \cos(\gamma)$, β and $\gamma \in \mathcal{R}$, and $[\boldsymbol{\omega}]_{\mathbf{b}} \triangleq [0, 0, 0, \boldsymbol{\omega}^T]^T$. When $\|\mathbf{v}\|$ is constant, it is known that any unit vector in \mathcal{R}^3 can be represented by $[C_\beta C_\gamma, C_\beta S_\gamma, S_\beta]^T$ with appropriate β and γ , where γ and β are the azimuth and elevation angles with respect to the basis

vectors $\bar{\mathbf{b}}_1$, $\bar{\mathbf{b}}_2$, and $\bar{\mathbf{b}}_3$. It follows from (15) that

$$\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v}) = \|\mathbf{v}\| \left[\frac{(C_\beta C_\gamma)^2}{q_{11}} + \frac{(C_\beta S_\gamma)^2}{q_{22}} \right] \quad (20)$$

$\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ is then proportional to C_β^2 when C_γ . In Example 1, the ratio of $|C_\beta|$ to $|S_\beta|$ will be changed to investigate the relationship between the index function $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ and the convergent rate of the depth estimation error. On the other hand, Example 2 will provide an experimental result.

Example 1: We consider the velocities of the camera with $\|\mathbf{v}(t)\| = 10$ mm/s for all cases. In this example, $c_1^2 + c_2^2$ varies with C_β and then $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ increases with $|C_\beta|$ increasing. Figure 1 shows the simulation results for $C_\gamma = 1/\sqrt{2}$ and $C_\beta = 1, 0.8, 0.5, 0.3, 0.1, 0, -0.4,$ and -0.7 , while $\boldsymbol{\omega} = 0.01[2, 1, 3]^T$ rad/s. It should be remarked that $\boldsymbol{\omega}$ is arbitrarily assigned and does not affect the property of the simulations. It can be seen from the depth estimation errors in Figure 1a and the cost functions $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ in Figure 1b that the depth estimation error converges more quickly when $\mathcal{J}_o(\boldsymbol{\xi}_0, \mathbf{v})$ is larger. We also show the simulation results of $\mathbf{v} = \pm 10\bar{\mathbf{v}}^*$ in Figure 1, where $\bar{\mathbf{v}}^*$ is the solution for the optimal problem:

$$\max_{\mathbf{v}} \mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v}) \text{ subject to } \|\mathbf{v}\| = 1 \text{ mm/s} \quad (21)$$

Figure 1b shows that the curves of $\mathcal{J}_o(\boldsymbol{\xi}, \pm 10\bar{\mathbf{v}}^*)$ are above all the others with the exception that the curve of $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ for $\mathbf{v} = 10\bar{\mathbf{v}}^*$ and that with $C_\beta = 1$ intersect about $t = 8.5$ s. This exception occurs because the image states for both are different for $t > 0$, and the optimization problem (21) is under the assumption that the image position (x, y) is given. Figure 1c shows that the image trajectories for different velocities are entirely different except that the initial image $(x_0, y_0) = (100, 100)$ in pixels. However, Proposition 1 still holds true, although the optimal solutions of (21) are different for different image states. The simulations in Figure 1 indicate that $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ is a good index for the convergent rate of the depth estimation.

We also changed the value of the fixed C_γ in the range $[-1, 1]$ and found that the simulations have the same property of the depth estimation performance as that in Figure 1.

2.2.2. Experiments

Example 2: In this example, a 3-DOF XYZ-type manipulator is used to drive the camera moving along three orthogonal directions. The images used in the experiment are 640×480 pixels in size. The parameters needed in (5) are $\gamma_x = 1663.4$ pixels and $\gamma_y = 1679.0$ pixels. The states are estimated by the EKF with the measurement noise covariance matrix $\mathbf{R} = \text{diag}(0.25 \text{ pixel}^2, 0.25 \text{ pixel}^2)$, the system noise covariance matrix $\mathbf{Q} = \text{diag}(0.25 \text{ pixel}^2, 0.25 \text{ pixel}^2, 100 \text{ mm}^2)$, and the image sampling period is 500 ms. The image coordinates are obtained from the center of the image of a black circle with diameter 6.4 mm. The initial image coordinates are located at $(4.90, -22.56)$ in pixels. The real initial depth is about 550 mm, while the estimate of the initial depth is 700 mm. The relationship between the camera frame and the base frame of the manipulator has been obtained after the camera calibration. Then the velocities can be described with respect to the camera frame. The camera velocities with $\|\mathbf{v}\| = 10$ mm/s are listed as follows [cf. (18)]:

$$\begin{aligned} [\mathbf{u}_1(\boldsymbol{\xi})]_b &= 10[1, 0, 0, 0, 0, 0]^T \\ [\mathbf{u}_2(\boldsymbol{\xi})]_b &= 10[0, 1, 0, 0, 0, 0]^T \\ [\mathbf{u}_3(\boldsymbol{\xi})]_b &= 10[0, 0, 1, 0, 0, 0]^T \\ [\mathbf{u}_4(\boldsymbol{\xi})]_b &= (10/3)[-2, 1, 2, 0, 0, 0]^T \\ [\mathbf{u}_5(\boldsymbol{\xi})]_b &= (10/\sqrt{14})[1, 2, 3, 0, 0, 0]^T \\ [\mathbf{u}_6(\boldsymbol{\xi})]_b &= (10/\sqrt{10})[-1, 0, 3, 0, 0, 0]^T \\ [\mathbf{u}_7(\boldsymbol{\xi})]_b &= (10/\sqrt{15})[0, 1, -4, 0, 0, 0]^T \\ [\mathbf{u}_8(\boldsymbol{\xi})]_b &= (10/\sqrt{27})[-1, 1, 5, 0, 0, 0]^T \end{aligned}$$

Note that \mathbf{u}_1 and \mathbf{u}_2 have zero components along \mathbf{v}_3 , while \mathbf{u}_3 has zero components along \mathbf{v}_1 and \mathbf{v}_2 . The angular velocities are zero for these eight inputs. Figure 2a shows that the depth estimation errors for these eight inputs, except for \mathbf{u}_3 , have a tendency to converge to zero. $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ curves in Figure 2b emphasize again that Proposition 1 is practical and useful.

Results of both examples support Proposition 1. However, the component ratio of \mathbf{v} is determined by a camera motion controller, not by a depth estimator. In next section, we attempt to increase $\mathcal{J}_o(\boldsymbol{\xi}, \mathbf{v})$ by a modified camera velocity controller for given feature points.

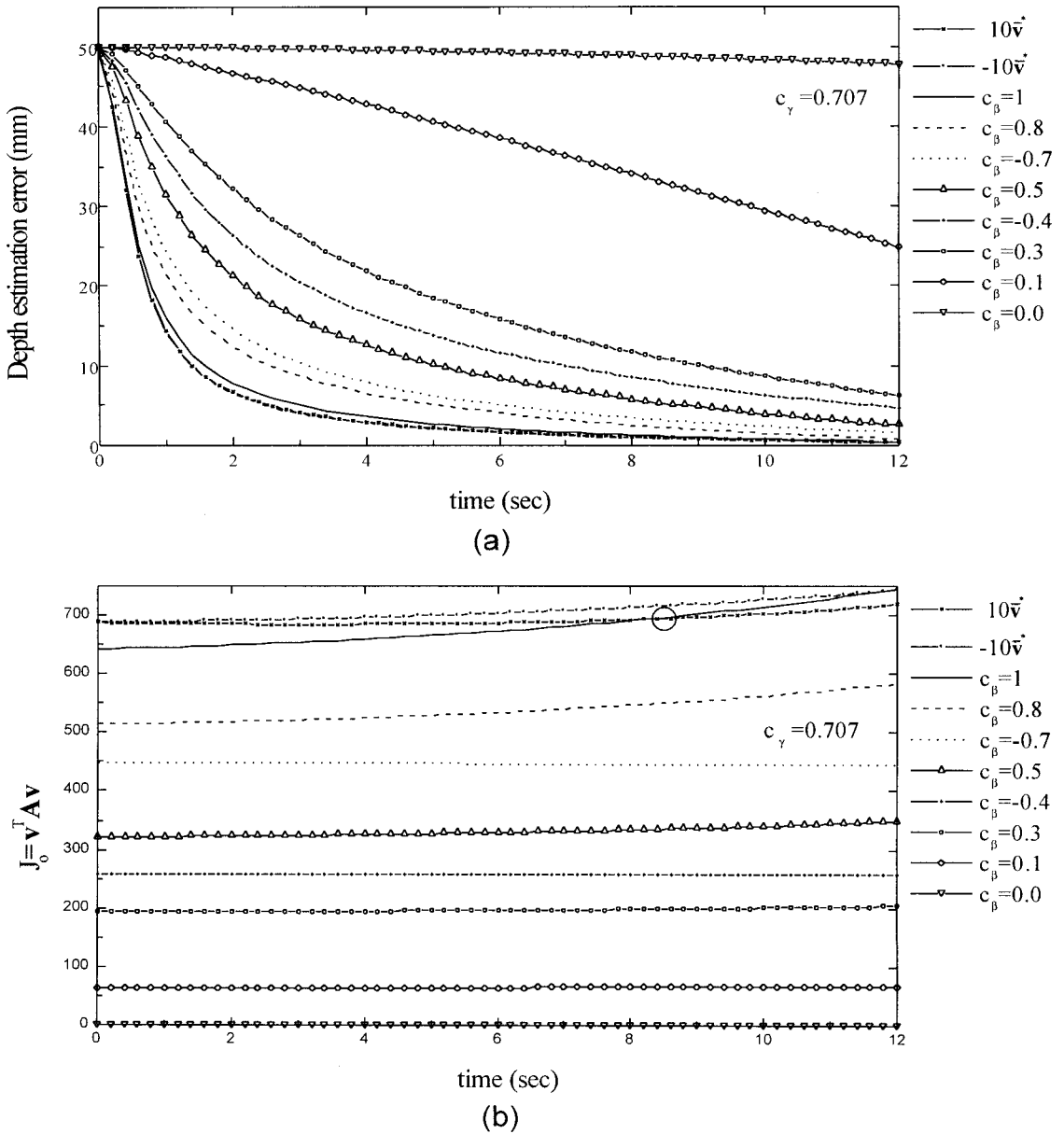


Figure 1. (a) Depth estimation errors. (b) The cost function $\mathcal{F}_0(\xi, \mathbf{v})$.

3. APPLICATION TO VISUAL SERVO CONTROL

The results of Section 2 will be used to design a visual servo control scheme. The concept is to correct the linear velocity \mathbf{v} of the camera by increasing $\mathcal{F}_0(\xi, \mathbf{v})$ as possible.

3.1. Control Scheme

First, we introduce the visual servo control scheme with the damped least-squares method (DLSM). Suppose that there are n feature points (x_i, y_i) .

Applying (3) to n feature points, we obtain

$$\dot{\mathbf{f}} = \mathbf{J} \mathbf{u} \tag{22}$$

where the feature vector \mathbf{f} and the visual Jacobian matrix \mathbf{J} are

$$\mathbf{f} = \begin{bmatrix} x_1 \\ y_1 \\ \vdots \\ x_n \\ y_n \end{bmatrix}, \quad \mathbf{J} = \begin{bmatrix} \bar{\mathbf{J}}(\xi_1) \\ \vdots \\ \bar{\mathbf{J}}(\xi_n) \end{bmatrix} \tag{23}$$

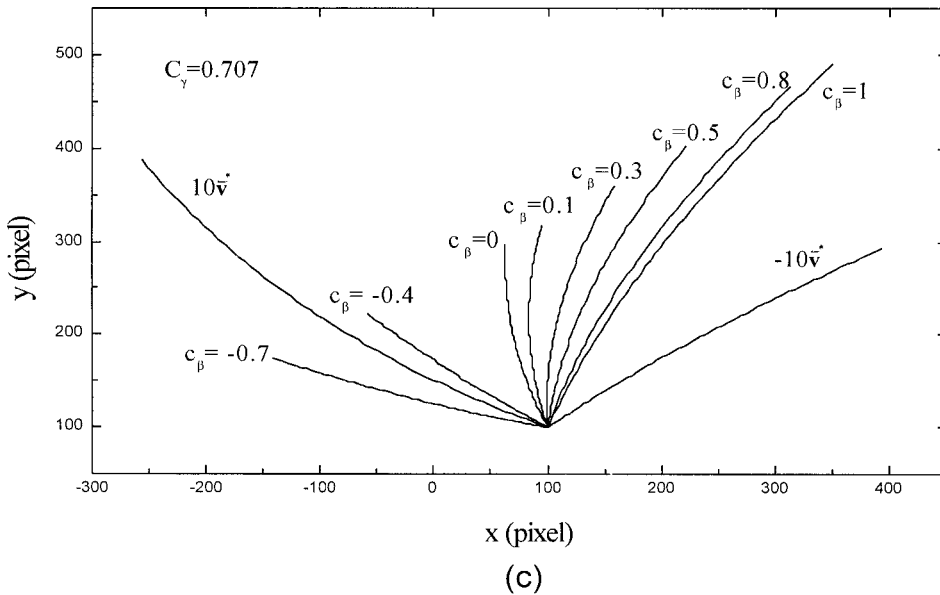


Figure 1. (continued) (c) The image trajectories for different inputs \mathbf{u} with $\|\mathbf{v}\| = 10$ mm/s for Example 1.

Note that $\mathbf{J} \in \mathcal{R}^{2n \times 6}$. The control purpose is to design the velocity of the camera \mathbf{u} , so that the rate of change of the feature points in the image of the camera follows the desired one.

The damped least-squares control scheme is to minimize $\|\mathbf{J}\mathbf{u} - \dot{\mathbf{f}}^*\|^2 + \rho_s^2 \|\mathbf{u}\|^2$, where $\dot{\mathbf{f}}^*$ is the feature velocity command and $\rho_s \in \mathcal{R}$ is the damping factor which represents the weighting of $\|\mathbf{u}\|^2$ with respect to the feature velocity error $\|\dot{\mathbf{f}} - \dot{\mathbf{f}}^*\|$. The control command \mathbf{u}^* is the optimal solution

$$\mathbf{u}^* = (\mathbf{J}^T \mathbf{J} + \rho_s^2 \mathbf{I})^{-1} \mathbf{J}^T \dot{\mathbf{f}}^* \quad (24)$$

where \mathbf{I} is the identity matrix. A nonzero ρ_s makes $(\mathbf{J}^T \mathbf{J} + \rho_s^2 \mathbf{I})$ positive definite, even if $\mathbf{J}^T \mathbf{J}$ is singular. Although this control scheme can alleviate the singularity problem, the input velocity \mathbf{u}^* may not help the performance of the depth estimation Z_i , which is required by (4). To compensate for this drawback, Proposition 1 motivates us to minimize the following objective function $\mathcal{F}(\mathbf{u})$,

$$\mathcal{F}(\mathbf{u}) = \|\mathbf{J}\mathbf{u} - \dot{\mathbf{f}}^*\|^2 + \rho_s^2 \|\mathbf{u}\|^2 - \frac{\rho_o^2}{n} \sum_{i=1}^n \rho_{oi}^2 \mathcal{F}_o(x_i, y_i, \mathbf{u}) \quad (25)$$

where

$$\mathcal{F}_o(x_i, y_i, \mathbf{u}) = \mathbf{u}^T \bar{\mathbf{A}}_i \mathbf{u} \quad (26)$$

$$\bar{\mathbf{A}}_i = \begin{bmatrix} \mathbf{A}(\xi_i) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathcal{R}^{6 \times 6} \quad (27)$$

in which $\mathbf{A}(\xi_i)$ is a matrix function defined in (14) for the i th point. The first two terms on the right-hand side of (25) are those in the damped least-squares method. The term $\mathcal{F}_o(\xi, \mathbf{v})$ in (25) is for improving the depth estimation. The factor ρ_{oi} is a weighting factor to compromise the control error and the performance of the depth estimation.

Since matrix $\mathbf{A}(\xi_i)$ is symmetric, $\bar{\mathbf{A}}_i$ is orthogonally diagonalizable: $\bar{\mathbf{A}}_i = \mathbf{U}_i^T \Lambda_i \mathbf{U}_i$, where $\Lambda_i = \text{diag}(\sigma_{i1}, \sigma_{i2}, 0, 0, 0, 0) \in \mathcal{R}^{6 \times 6}$, σ_{i1} and σ_{i2} positive eigenvalues of $\mathbf{A}(\xi_i)$, and \mathbf{U}_i is an orthogonal matrix. Thus, (26) is rewritten as

$$\mathcal{F}_o(x_i, y_i, \mathbf{u}) = \mathbf{u}^T (\mathbf{U}_i^T \Lambda_i \mathbf{U}_i) \mathbf{u} \quad (28)$$

By calculus, we set $\partial \mathcal{F}(\mathbf{u}) / \partial \mathbf{u} = \mathbf{0}$ to obtain the optimal solution \mathbf{u}^* to (25) as

$$\mathbf{u}^* = \mathbf{W}^{-1} \mathbf{J}^T \dot{\mathbf{f}}^* \quad (29)$$

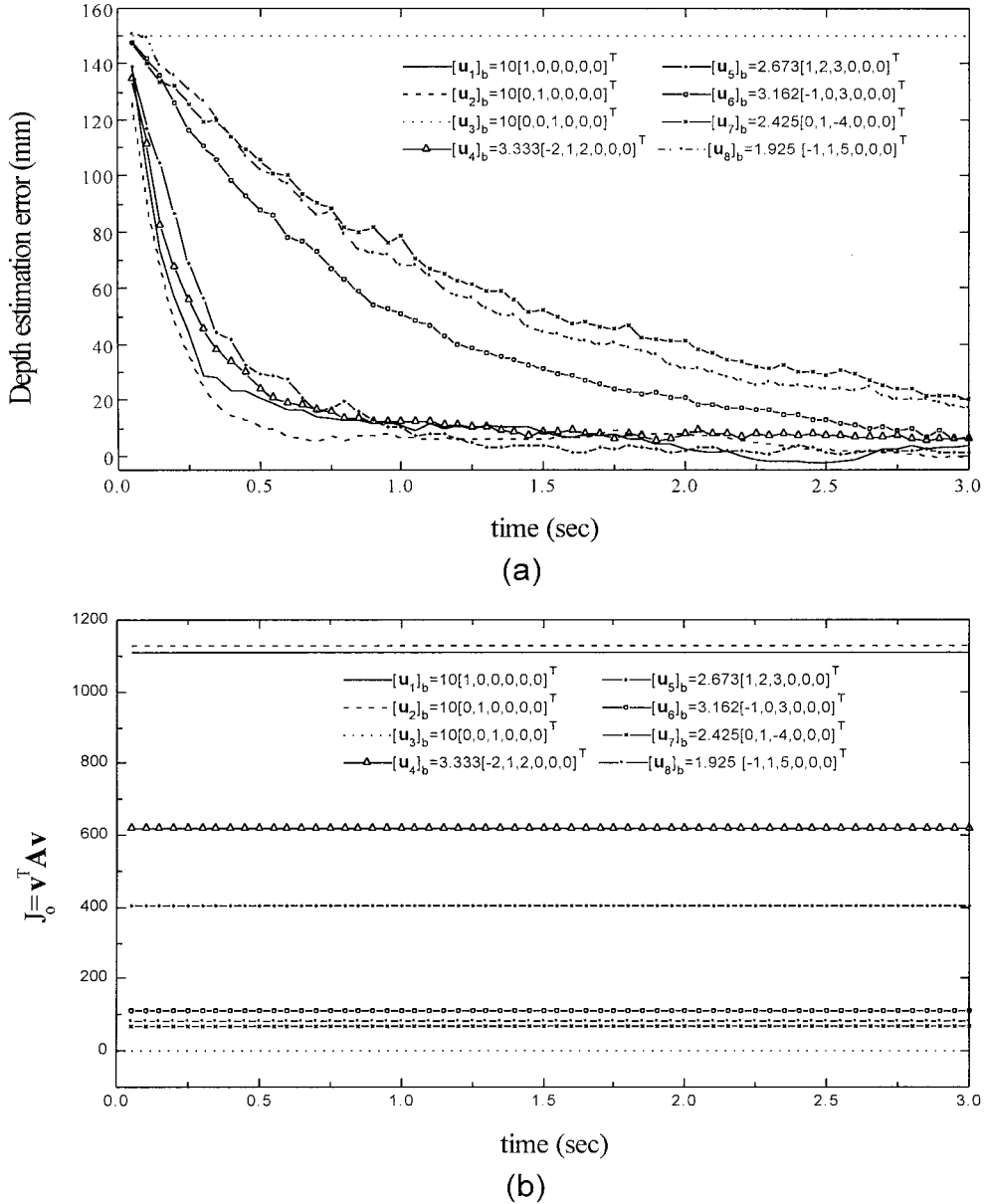


Figure 2. (a) Depth estimation errors. (b) The cost function $\mathcal{J}_0(\xi, \mathbf{v})$ for different inputs \mathbf{u} with $\|\mathbf{v}\| = 10$ mm/s for Example 2.

where

$$\mathbf{W} = \mathbf{J}^T \mathbf{J} + \rho_s^2 \left[\mathbf{I} - \frac{1}{n} \sum_{i=1}^n \rho_{oi}^2 (\mathbf{U}_i^T \Lambda_i \mathbf{U}_i) \right] \quad (30)$$

which is a positive definite symmetrical matrix if $\rho_{oi}^2 \max\{\sigma_{i1}, \sigma_{i2}\} \leq 1, \forall i = 1, \dots, n$. We shall call the control law (29) the observabilized damped least-

squares method (abbreviated as ODLSM). When $\rho_{oi}^2 < 1/\max\{\sigma_{i1}, \sigma_{i2}\}, \forall i = 1, \dots, n$, it follows from (25) that $\mathcal{A}(\mathbf{u})$ is always positive.

Suppose that the true values of depths can be obtained by a depth estimator like the extended Kalman filter in a few seconds. Thereafter, \mathbf{J} in (23) is approximately calculated by the estimated values of depths. We can then expect that the overall system is asymptotically convergent according to Theo-

rem A1 in the Appendix. It is verified by the simulation described in the following.

3.2. Simulation Example

The simulation example considers three image points of the corners of a triangle pattern. The initial image coordinates of three corners are located at (10, 63), (90, 51), and (29, 153) in pixels and their real initial depths are all 550 mm but are unknown in this simulation. The desired image has three corner images at (−50, −68), (50, −68), and (−50, 43) in pixels. Note that the final depths are all 450 mm corresponding to the desired image feature. Suppose that the estimates of the initial depths are $\hat{Z}_1 = 595$ mm, $\hat{Z}_2 = 599$ mm, and $\hat{Z}_3 = 596$ mm; i.e., the initial estimation errors are about 50 mm. The intrinsic parameters of the camera are the same as those in Example 1 in Section 2.2. The noise covariance matrices in EKF are $\mathbf{R} = \text{diag}(0.25 \text{ pixel}^2, 0.25 \text{ pixel}^2)$ and $\mathbf{Q} = \text{diag}(0.25 \text{ pixel}^2, 0.25 \text{ pixel}^2, 25^2 \text{ mm}^2)$. Both DLSM and ODLSM controllers [see (24) and (29)] have the same following data: the sampling period $T_s = 200$ ms, the proportional gain $K_p = 0.65$, the damping factor $\rho_s = 0.003$, and the weighting factor $\rho_{oi} = 0.9/\max\{\sigma_{i1}, \sigma_{i2}\}$, $i = 1, 2, 3$, for each corner.

The history of the error norm of the estimated depths, i.e., $\sqrt{\sum_{i=1}^3 (Z_i - \hat{Z}_i)^2}/3$, in Figure 3a reveals that the depth estimation in the ODLSM controller is superior to that in the DLSM controller. The steady-state error of the depth estimate in DLSM is about 7 mm, while that in the ODLSM nearly vanishes for the same EKF estimator. The feature errors in both the ODLSM and DLSM are almost the same and converge to zero as is shown in Figure 3b. That means using a moderate small ρ_{oi} in ODLSM has little effect on the convergence performance of the visual servo control.

Figures 4 and 5 show the linear and angular velocities of the camera generated by the DLSM and ODLSM, respectively. The velocities are expressed with respect to the time-varying basis (16) by the component vector $[c_1, c_2, \dots, c_6]^T$. Comparing Figure 4 with Figure 5, we find that c_1 and c_2 are enlarged (while c_3 is reduced) by the ODLSM in the beginning of the control process. This velocity history in the ODLSM then provides a better depth estimation as was expected by Proposition 1. Consequently, the advantage of the ODLSM is that a good depth estimate can be achieved while the convergence performance is retained.

4. CONCLUSION

This article presents a performance criterion for the depth estimation problem of the visual servo system of a robot visual system. The convergence rate of the depth estimator increases with the value of the performance criterion stated in Proposition 1. Although this result is a heuristic extension of that of a linear system, some simulations and experiments verify the validation of this extension to the nonlinear visual system, provided that the camera velocity is not fast. Finally, we apply the performance criterion of the depth estimation to the design of an estimator-based visual servo control scheme for a manipulator, named the observabilized damped least-squares method, to improve the depth estimation. A simulation shows that the improvement of the depth estimation is achieved without any sacrifice of the convergence performance. Future work will be the extension of the control for a moving object.

APPENDIX

Theorem A1: Consider the visual servo system (22) with the ODLSM controller (29) and

$$\dot{\mathbf{f}}^* = K_p(\mathbf{f}_d - \mathbf{f}) \quad (\text{A1})$$

where $K_p > 0$ is a constant proportional gain. Suppose that the depth estimates of the feature points are bounded and will converge to the real values. Let the feature reference input \mathbf{f}_d be a step input. Then the feature vector \mathbf{f} asymptotically converges to \mathbf{f}_d or to some \mathbf{f}_s such that $\mathbf{J}^T \mathbf{f}_s = 0$, if the weighting factors ρ_{oi} , $i = 1, \dots, n$, in (30) all satisfy $\rho_{oi}^2 \max\{\sigma_{i1}, \sigma_{i2}\} \leq 1$.

Proof: We consider only the time after t_0 at which the depth estimates have converged to the true values. \mathbf{f} is the state of the system (22). So define the Lyapunov function candidate as

$$V(\tilde{\mathbf{f}}) = \frac{1}{2} \tilde{\mathbf{f}}^T \tilde{\mathbf{f}} \quad (\text{A2})$$

where $\tilde{\mathbf{f}} = \mathbf{f} - \mathbf{f}_d$. $V(\tilde{\mathbf{f}})$ is positive definite and radially unbounded. The time derivative of $V(\tilde{\mathbf{f}})$ along the solution trajectory of the closed-loop system

$$\dot{\mathbf{f}} = K_p \mathbf{J} \mathbf{W}^{-1} \mathbf{J}^T (\mathbf{f}_d - \mathbf{f}) \quad (\text{A3})$$

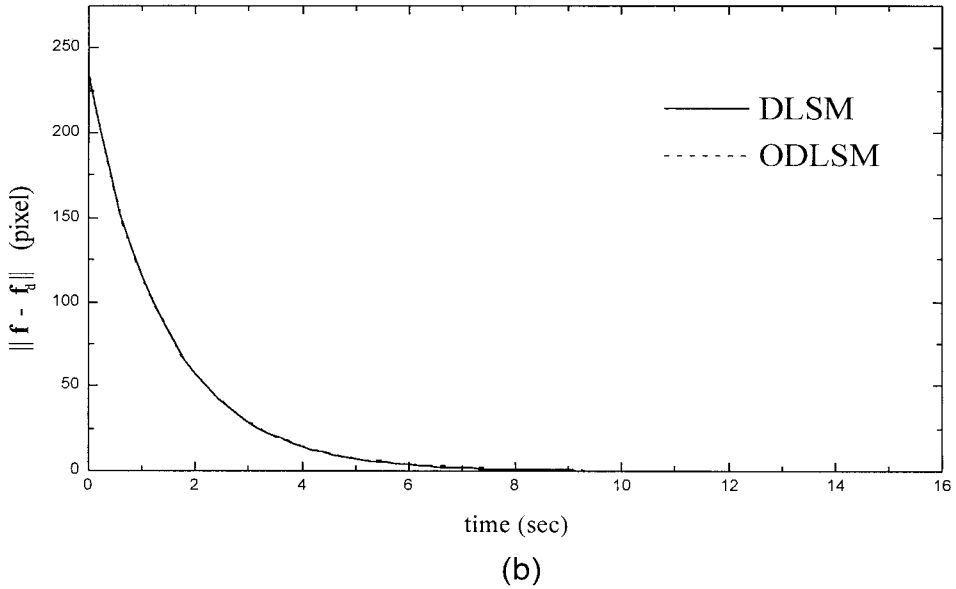
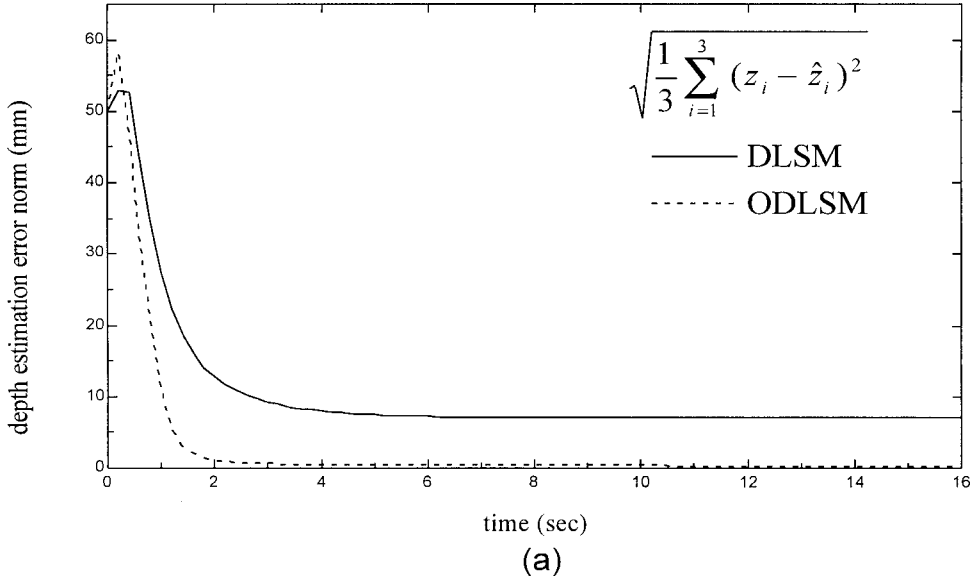


Figure 3. (a) The norm of the depth estimation errors. (b) The norm of the feature feedback errors by DLSM and ODLSM.

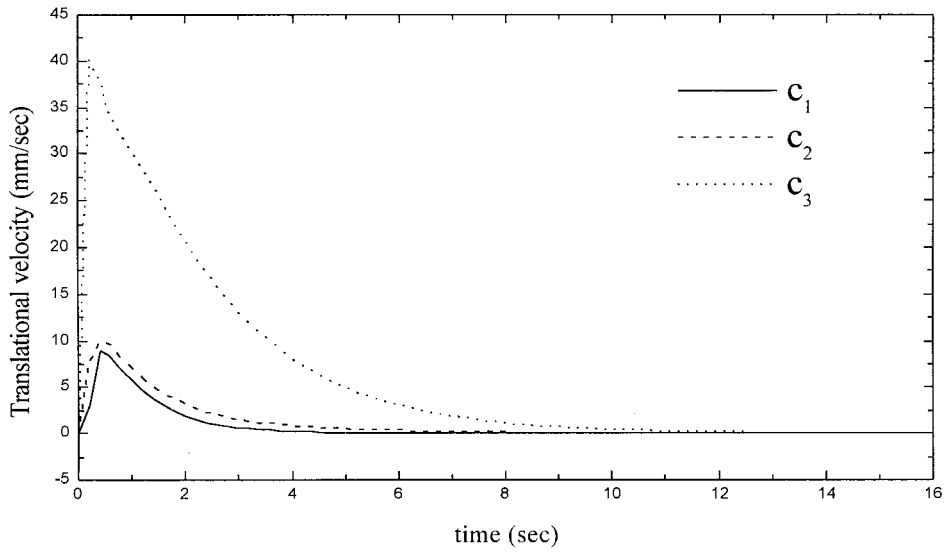
is

$$\dot{V}(\tilde{\mathbf{f}}) = -K_p \tilde{\mathbf{f}}^T \mathbf{J} \mathbf{W}^{-1} \mathbf{J}^T \tilde{\mathbf{f}} \leq 0 \quad (\text{A4})$$

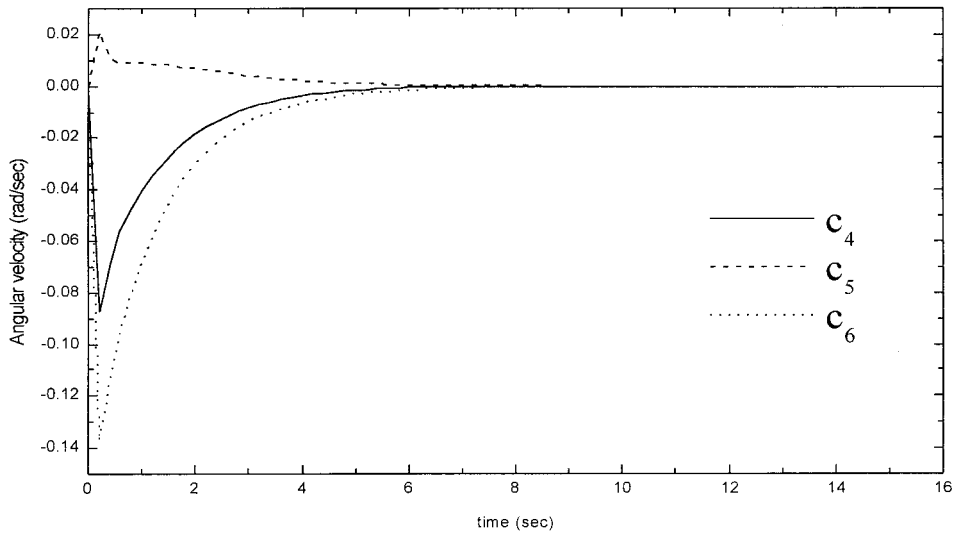
since \mathbf{W} is a positive definite matrix when

$$\rho_{oi}^2 \max\{\sigma_{i1}, \sigma_{i2}\} \leq 1.$$

Let $\mathcal{S} = \{\exists \mathbf{f}: t \geq 0 \text{ such that } \dot{V}(\mathbf{f} - \mathbf{f}_d) = 0\}$ and let \mathcal{M} denote the largest invariant set of (A3) contained in \mathcal{S} . LaSalle's invariance principle,^{33,34} states that all solution trajectories of (A3) globally asymptotically converge to \mathcal{M} as $t \rightarrow \infty$. Apparently, $\mathcal{S} = \{\mathbf{f}: \mathbf{f} = \mathbf{f}_d \text{ or } \mathbf{J}^T(\boldsymbol{\xi})(\mathbf{f} - \mathbf{f}_d) = 0\}$. It follows from (A3) that $\dot{\mathbf{f}} = 0$ and then \mathbf{f} is stationary when $\mathbf{J}^T(\boldsymbol{\xi})(\mathbf{f} - \mathbf{f}_d) = 0$. Thus, $\mathcal{M} = \mathcal{S}$. This completes the proof. Q.E.D.



(a)



(b)

Figure 4. Camera velocity by DLSSM (a) linear and (b) angular velocities.

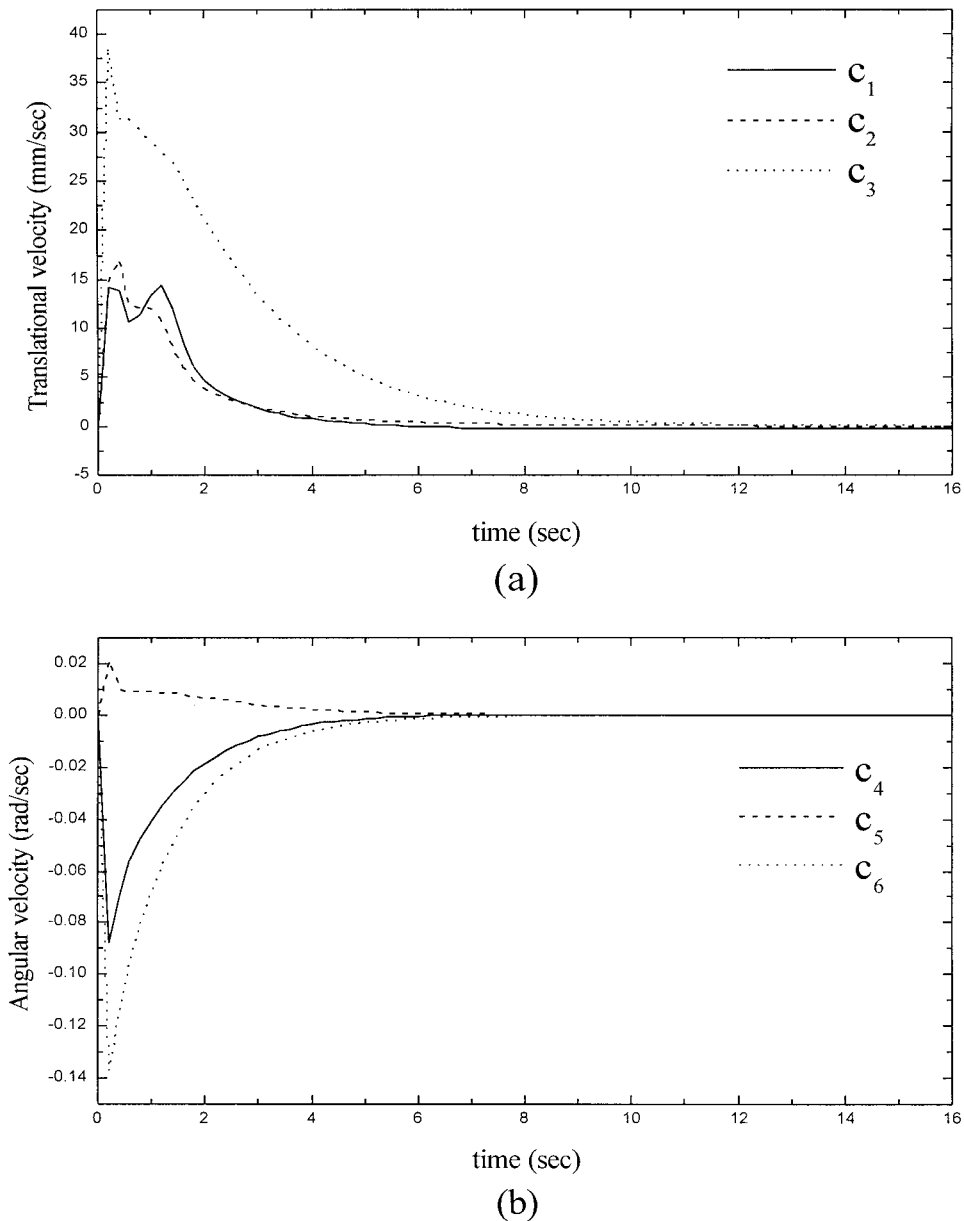


Figure 5. Camera velocity by ODLSM (a) linear and (b) angular velocities.

REFERENCES

1. N.P. Papanikolopoulos, P.K. Khosla, and T. Kanade, Visual tracking of a moving target by a camera mounted on a robot: a combination of control and vision, *IEEE Trans Robotics Automat* 9(1) (1993), 14–34.
2. P.K. Allen, A. Timcenko, B. Yoshimi, and P. Michelman, Automated tracking and grasping of a moving object with a robotic hand-eye systems, *IEEE Trans Robotics Automat* 9(2) (1993), 152–165.
3. A.A. Rizzi and D.E. Koditschek, An active visual estimator for dexterous manipulation, *IEEE Trans Robotics Automat* 12(5) (1996), 697–713.
4. T.S. Huang and A.N. Netravali, Motion and structure from feature correspondences: a review, *Proc IEEE* 82(2) (1994), 251–267.
5. K.S. Arun, T.S. Huang, and S.D. Blostein, Least-squares fitting of two 3-D point sets, *IEEE Trans Pattern Anal Machine Intell* PAMI-9(5) (1987), 698–700.
6. H.H. Chen and T.S. Huang, Maximal matching of 3-D points for multiple-object motion estimation, *Pattern Recognition* 21(2) (1988), 75–90.

7. P. Gros, Matching and clustering: two steps toward automatic object modeling in computer vision, *Int J Robotics Res* 14(6) (1995), 633–642.
8. B. Sridhar, R. Suorsa, P. Smith, and B. Hussien, Vision-based obstacle detection for rotorcraft flight, *J Robotic Syst* 9(6) (1992), 709–727.
9. G.D. Hager, A modular system for robust positioning using feedback from stereo vision, *IEEE Trans Robotics Automat* 13(4) (1997), 582–595.
10. S. Soatto, R. Frezza, and P. Perona, Motion estimation via dynamic vision, *IEEE Trans Automat Contr* 41(3) (1996), 393–413.
11. R.Y. Tsai and T.S. Huang, Uniqueness and estimation of three-dimensional motion parameters of rigid objects with curved surfaces, *IEEE Trans Pattern Anal Machine Intell PAMI-6*(1) (1984), 13–27.
12. J. Hespanha, Z. Dodds, G.D. Hager, and A.S. Morse, Decidability of robot positioning tasks using stereo vision systems, *Proc 37th IEEE Int Conf Decision Contr*, 1998, pp. 3736–3741.
13. L. Matthies and T. Kanade, Kalman filter-based algorithms for estimating depth from image sequences, *Int J Comput Vision* 3 (1989), 209–236.
14. C.K. Chui and G. Chen, *Kalman filtering with real-time applications*, Springer-Verlag, New York, 1991.
15. H. Sutanto and R. Sharma, Global performance evaluation of image features for visual servo control, *J Robotic Syst* 13(4) (1996), 243–258.
16. L.E. Weiss, A.C. Sanderson, and C.P. Neuman, Dynamic sensor-based control of robots with visual feedback, *IEEE J Robotics Automat RA-3*(5) (1987), 404–417.
17. J.T. Feddema and O.R. Mitchell, Vision guided servoing with feature-based trajectory generation, *IEEE Trans Robotics Automat* 5(5) (1989), 691–700.
18. W. Jang and Z. Bein, Feature-based visual servoing of an eye-in-hand robot with improved tracking performance, *Proc IEEE Int Conf Robotics Automat*, 1991, pp. 2254–2260.
19. P.I. Corke and R.P. Paul, Video-rate visual servoing for robots, Technical Report MS-CIS-89-18, GRASP Laboratory, University of Pennsylvania, 1989.
20. K. Hashimoto, T. Ebine, and H. Kimura, Visual servoing with hand-eye manipulator—optimal control approach, *IEEE Trans Robotics Automat* 12(5) (1996), 766–773.
21. A.J. Koivo and N. Houshangi, Real-time vision feedback for servoing robotic manipulator with self-tuning controller, *IEEE Trans Syst Man Cybernet* 21(1) (1991), 134–141.
22. W.J. Wilson, C.C.W. Hulls, and G.S. Bell, Relative end-effector control using cartesian position based visual servoing, *IEEE Trans Robotics Automat* 12(5) (1996), 684–696.
23. R. Sharma and S. Hutchinson, On the observability of robot motion under active camera control, *Proc IEEE Int Conf Robotics Automat*, 1994, pp. 162–167.
24. J.T. Feddema, C.S.G. Lee, and O.R. Mitchell, Weighted selection of image features for resolved rate visual feedback control, *IEEE Trans Robotics Automat* 7(1) (1991), 31–47.
25. B. Nelson and P.K. Khosla, Integrating sensor placement and visual tracking strategies, *Proc IEEE Int Conf Robotics Automat*, 1994, pp. 1351–1356.
26. N.P. Papanikolopoulos, P.K. Khosla, and T. Kanade, Adaptive robotic visual tracking: theory and experiments, *IEEE Trans Automat Contr* 38(3) (1993), 429–445.
27. C.E. Smith, S.A. Brandt, and N.P. Papanikolopoulos, Eye-in-hand robotic tasks in uncalibrated environments, *IEEE Trans Robotics Automat* 13(6) (1997), 903–914.
28. S. Soatto, 3-D structure from visual motion: modeling, representation and observability, *Automatica* 33(7) (1997), 1287–1312.
29. W.P. Dayawansa, B.K. Ghosh, C. Martin, and X. Wang, A necessary and sufficient condition for the perspective observability problem, *Syst Contr Lett* 25 (1995), 159–166.
30. R. Sharma and S. Hutchinson, Optimizing hand/eye configuration for visual-servo systems, *Proc IEEE Int Conf Robotics Automat*, 1995, pp. 172–177.
31. R.M. Haralick and L.G. Shapiro, *Computer and robot vision*, vol. II, Addison-Wesley, Reading, MA, 1993.
32. T. Söderström and P. Stoica, *System identification*. Prentice-Hall, Englewood Cliffs, NJ, 1989.
33. M. Vidyasagar, *Nonlinear system analysis*. Prentice-Hall, Englewood Cliffs, NJ, 1993.
34. H.K. Khalil, *Nonlinear systems*. Prentice-Hall, Englewood Cliffs, NJ, 1996.