

CAMERA CALIBRATION USING NEURAL NETWORKS

Márcio Mendonça

Department of Electrical Technology
Federal Center of Technology and Education
86300-000, Cornelio Procópio, PR
Brazil
marcio_mendonca@zipmail.com.br

Ivan N. da Silva, José E. C. Castanho

Department of Electrical Engineering
School of Engineering, São Paulo State University, CP. 473
17015-970, Bauru, SP
Brazil
ivan|castanho@feb.unesp.br <http://www.feb.unesp.br>

ABSTRACT

This work presents a procedure for camera calibration using artificial neural networks of the type back propagation perceptron. Camera calibration is employed in computer vision for pose determination and it requires a solution of non-linear system of equations. By employing neural network, it becomes unnecessary to know the parameters of the cameras, such as focus, distortions besides the geometry of the system. Camera simulations and real experiments are used to demonstrate and evaluate the procedure.

Keywords: Computer vision, camera calibration, neural networks, and stereo vision.

1. INTRODUCTION

The goal of camera calibration is to establish the relationship between the global 3D-coordinates of a point and the 2D-coordinates of the projected image [Echig90]. The process of camera calibration is a pre-requirement for most applications in computer vision to determine the position and orientation of an object with regard to a global coordinate system. Most methods of camera calibration usually require setting a careful and difficult procedure and include a complex mathematical model.

In this work, we present a method to establish the relationship between 3D coordinates and the image coordinates of a point using neural networks. The method is advantageous when compared to traditional ones since it does not need the information of the camera geometry and neither needs a complex experimental procedure to be settled. Therefore, the procedure can be more useful since it is easy to employ in common situations. There are many techniques to obtain 3-D information [Heikk00], [Andre99], including stereo

vision [Aguil96], and structured light [Guiss00]. In stereo vision, two or more cameras are used to visualize the same point in the space inside the field of vision. So, each camera obtains a different coordinate for this point in their respective images, that is, each point in the space can generate only a pair of coordinates, (x_1, y_1) and (x_2, y_2) , respectively in the each camera. In that way, using triangulation it is possible to recover the position of the point in the space starting from the coordinates of the point projected in both cameras. However, this implies in the solving a non-linear system of equations, which is usually done through optimization techniques. With the proposed method, a neural network is used to learn the relationship among global coordinates and the camera coordinates. Moreover, it is straight up to recover the 3D information without knowing camera parameters explicitly if there are two images of the same scene. The method is also robust enough for dealing with different cameras - with different focal lengths - because it is capable of recognizing the actual focus of the cameras once the neural networks have been trained.

2. THE CAMERA MODEL

Each point in the space can have only one corresponding projection in images - (x_1, y_1) and (x_2, y_2) - for each camera. Using the camera model, it is possible to determine the coordinate x and y of a given point in the field of vision for both images. This model, that is used to simulate the camera, includes the rotation, translation, perspective, and location transformations that are needed to get the image projection of a point in space [Gonza92]. Usually, the coordinates of cameras are expressed in homogeneous form to simplify the processing of matrices. The complete transformation from real world coordinates to image coordinates is given by the composition of all the transformations as shown in Eq. 1.

$$C_h = P \times C \times R \times G \times W_h \quad (1)$$

Where the matrix P , for perspective transformation, is given by:

$$P = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & \frac{-1}{\lambda} & 1 \end{bmatrix} \quad (2)$$

The matrix G , for translation of the camera related to the origin of the global coordinates, is given by:

$$G = \begin{bmatrix} 1 & 0 & 0 & -X_0 \\ 0 & 1 & 0 & -Y_0 \\ 0 & 0 & 1 & -Z_0 \\ 0 & 0 & 1 & 1 \end{bmatrix} \quad (3)$$

Matrix rotation R , with two degrees of freedom, is given by:

$$R = \begin{bmatrix} \cos \theta & \sin \theta & 0 & 0 \\ -\sin \theta \cdot \cos \alpha & \cos \theta \cdot \cos \alpha & \sin \alpha & 0 \\ \sin \theta \cdot \sin \alpha & -\cos \theta \cdot \sin \alpha & \cos \alpha & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (4)$$

The displacement of the structure to fix the camera in a particular position is given by:

$$C = \begin{bmatrix} 1 & 0 & 0 & -r_1 \\ 0 & 1 & 0 & -r_2 \\ 0 & 0 & 1 & -r_3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad (5)$$

W_h contains the global coordinates expressed in homogeneous form, and C_h the homogeneous image

coordinates after transformations. Image coordinates in both cameras are obtained from points (x_1, y_1) , and (x_2, y_2) by Eq. 1. More details of this formulation can be found in [Gonza92]. The model was used to simulate the experiments under controlled conditions. This is a quite simple representation of a real camera, since it does not consider the lens distortions. In several other works lens distortions are considered and suitable mathematical models are developed [Heikk00], [Weng92]. In this work, this problem is not considered since it is supposed that using neural networks the distortions problems will be automatically accounted for.

3. IMPLEMENTATION

In our approach, we have solved the calibration camera using neural networks. Neural networks are suitable for solving non-linear problems, and are capable of learning relationships between the camera coordinate and global coordinate [Linch99], [Qing99], without needing to determine a complex mathematical model. The layers of an artificial neural network are usually classified in three groups. Input layer, where the patterns are presented to the net. In this case the inputs are the cameras coordinates (x_{1i}, y_{1i}) and (x_{2i}, y_{2i}) . Hidden layers, where most of processing is carried out through weighted connections, can be considered as characteristics extractors. Out layer presents the results, which in this work are the global coordinates $Xg, Yg,$ and Zg of each point. Therefore, the inputs of the neural network are the points of image in the both cameras and the outputs are the corresponding 3D space coordinates of each point.

To implement the neural networks it was used the Neural Network Matlab Toolbox and the Levenberg method for training. Two different cameras were used in the experiments: a WebCam from Creative Labs and a Digital Camera from Sony. The experiments comprised comparisons among the results from simulations of camera model and those using real images. In addition, to have a reference for comparison the problem was solved using least square estimation method. Both results are presented. For purposes of experimentation, it were generated several grids of points in random space inside a defined area or volume. The experiments with real images were conducted with a square grid and volumes similar of those used in simulated images. When using simulated images two types of data were generated. One data set was generated without noise, and in the other, a Gaussian noise was added to represent the acquisition error for image coordinates. Experiments were arranged in six groups regarding net training in which the neural nets have generalized in all cases. These groups are: 1) training in which the net recognizes

camera intrinsic parameters as the focus with synthetic data set; 2) training in a surface (constant z coordinates) using synthetic data set; 3) training in a volume using synthetic data; 4) training with a surface using real images and points; 5) training with a volume with real images planes and real points; 6) training in which the net recognizes camera position in space.

4. RESULTS AND DISCUSSION

The generated and simulated data sets were used as inputs for the camera calibration problem for both neural networks and least square solution. The absolute (A.E.) and the relative error (R.E.) between the actual simulated coordinates and the solution inferred using neural network is show Table 1, as well as, the least square solution in Table 2. Table 3 shows several configurations parameters used in net experiments, such as topology, number of epochs, neurons, and input. In experiments 1, and 2, it was used neural network with a simulated input data set. The very small error obtained demonstrates the suitability of the approach.

Exp.	Points	A.E.	R.E.
01	300	0.00000	0%
02	300	0.00001	0%
03	250 ^E	0.00213	0.2%
04	250 ^E	0.00091	0.1%
06	200	-0.00011	0%
07	250	-0.33330	0%
08	250	0.00010	0%
09	250	0.00025	0%
13	300 ^A	0.00117	0.1%
14	100	-0.00024	0%
15	300	0.00000	0%
18 ^R	130	0.20560	5%
20 ^R	130	-0.28720	-2%
21 ^R	130	-0.01500	-4%
22 ^R	520	0.01895	1%
23 ^R	130 ^P	0.02560	1.5%

A - Random points.
R - Experiments with real points.
E - Simulations with Gaussian error.
F - Simulations with focus recognition.
P - Experience with position recognition.

Results using Neural Networks.
Table 1

In experiments 3, 4 and 5 it was added a Gaussian error of 5% in input data to get a closer simulation of the actual conditions of image acquisition. Experiments 3, and 4 were conducted with neural networks with good results. Experiment 5 was run using least square estimation with comparable results. Experiment 6 was carried out using cameras with different focus lengths, after a

previous training with two different cameras. This showed that the neural networks were capable of identifying which camera was used.

Exp.	Points	AE.	R.E.
05	150 ^E	-0.00678	-0.8%
10	100 ^A	0.01169	1.2%
11	200	-0.00267	-0.3%
12	300	0.00076	0.1%
17 ^R	130	0.23590	4%
19 ^R	130	-0.03695	1%

Results using Least Square Method.
Table 2

Experiments 7, 8, and 9 were accomplished with neural network and another set of synthetic data in a plane surface as input. In experiments 10, 11, and 12 the input data was generated in a volume delimited by a parallelepiped and the solution was got using least square method. The experiments 13, 14, and 15, were conducted with test data generated in a volume delimited by a parallelepiped and the solution was obtained using neural networks.

Num.	Error	Epochs	Topology	Input
01	10 ⁻⁷	7	[20-30-3]	300
02	10 ⁻⁷	9	[20-30-3]	300
03	10 ⁻⁶	252	[20-30-3]	250 ^E
04	10 ⁻⁷	989	[20-30-3]	250 ^E
06	10 ⁻⁸	43	[20-30-5] ^F	200
07	10 ⁻⁶	7	[20-30-3]	250
08	10 ⁻⁷	9	[20-30-3]	250
09	10 ⁻⁸	12	[20-30-3]	250
13	10 ⁻⁷	10	[20-30-3]	300 ^A
14	10 ⁻⁸	8	[20-30-3]	100
15	10 ⁻⁸	13	[20-30-3]	300
18 ^R	10 ⁻⁸	130	[20-30-3]	130
20 ^R	10 ⁻⁶	358	[50-3]	130
21 ^R	10 ⁻⁸	1359	[50-3]	130
22 ^R	10 ⁻⁵	5000	[55-3]	520
23 ^R	10 ⁻⁶	1256	[55-4]	130

Characteristics of Neural Network.
Table 3

Experiments 17 and 18 were carried out with acquisition of real points in images of planes using the WebCam, which present a high degree of lens distortion. These experiments made possible to evaluate the behavior of the method using low cost off the shelf cameras. The experiments 19, 20, and 21 were similar to 17 and 18 with the difference that they were got using a Sony camera. Both methods least square and neural network presented similar results. The experiment 22 was done using a neural network and real images with the training points inside a 3-D volume. In this experiment, a total of

520 points spread out in four planes were chosen in the digitized test image. In the experiment 23, results were obtained using a neural network with one additional neuron in the output layer. This was done to enable the neural network to produce the distance of the camera from a reference in space. This is useful in applications where the camera is not fixed and its position information is of interest. All the results were obtained using neural networks with one or two hidden layers. Observing the results it is possible to conclude that they are similar in both cases, confirming that the net was able to generalize and solve the problem. At least one hidden layer is necessary to solve any non-linear problem [Hayki99]. The need to use more than two hidden layers will occur only when the problem is in a discontinuous domain, which is not the present case. The experiments in which the neural network was trained with error 10^{-6} have produced good results. Increasing the training error to 10^{-7} and 10^{-8} does not produce significant improvement in the accuracy, and increases the processing time. In real images, the neural networks were trained with error 10^{-5} for volumes and 10^{-7} , 10^{-8} for surface when the processing time was not so long. Using less than 100 patterns for tests produced inaccurate results. Increasing to 200 patterns or more the error has decreased to an insignificant value. In cases in which Gaussian error was added in the input, the neural network still presented good accuracy demonstrating its robustness. Neural network generalized in all cases, presenting better precision than least square methods in simulated data. However, in experiments with real images the two methods are equivalent with a small advantage for least square. However, the difference in the accuracy for both methods is not too significant. Thus, using neural networks is simpler since the implementation does not require a complex mathematical model for the camera, and avoids much of the practical calibration details, mainly the geometric displacement of the system and modeling lens distortion.

5. CONCLUSIONS

In this paper a camera calibration method using neural networks in the context of 3D information recovery was presented, and the results were compared against the traditional least square estimation. The main advantages of the proposed method are: 1) one does not need to know complex mathematical models, that is, using the method does not require deep technical knowledge, and so it is useful for an inexperienced user; 2) an initial estimation of calibration parameters are not needed; 3) the method can be applied to several different cameras; 4) the neural network can sense and recognize when different cameras are used, and it will still produce

the correct outputs, once it was previously trained, without needing to enter an explicitly input parameter that warns the system of the change. 5) the neural network can also recognize some specific positions of camera, since it has been previously trained. This feature is useful in dynamic systems.

Therefore, the presented approach is very flexible and significantly more easy of being employed than those previously presented in the literature.

REFERENCES

- [Aguil96] Aguilar, J. J., F. Torres and M. A. Lopes *Stereo Vision for 3D measurement: accuracy analysis, calibration and industrial applications*. Elsevier Science. vol. 18, n.4, pp. 193-200, 1996.
- [Andre99] Andreff N., Horaud, R. and Espiau, B. *On-line Hand-Eye Calibration*. in Second International Conference on 3-D Digital Imaging and Modeling (3DIM'99), pp. 430-436, 1999.
- [Echig90] Echigo Tomio. *A Camera Calibration Technique using Three Sets of Parallel Lines*. Machine Vision and Applications. vol.3, pp.159-167, 1990.
- [Gonza92] Gonzalez, Rafael C., Richard E. Woods. *Digital Image Processing*. Addison-Wesley Publishing Company, Inc.1992.
- [Guiss00] Guisser, L., R. Payrissat, S. Castan. *An accurate 3D vision system using a projected grid for surface descriptions*. Image and Vision Computing, vol. 18, pp. 463-491, 2000.
- [Hayki99] Haykin, Simon. *Neural Networks, a Comprehensive Foundation*. Prentice Hall Inc., 2ed. 1999
- [Heikk00] Heikkilä, Janne. *Geometric Camera Calibration Using Circular Control Points*. IEEE Transactions on Pattern Analysis and Machine Intelligence. vol. 22, n.10, pp. 1066-1076, 2000.
- [Lynch99] Lynch, Mark B., Cihan H. Dagli , Mahesh Vallenki. *The use of feedforward neural networks for machine vision calibration*. Int. Journal of Production Economics, 60-61, pp. 479-489, 1999.
- [Qing99] Qing, Guo Wei and G Hirzinger *Multisensory Visual Servoing by a Neural Network*. IEEE Transactions on Systems, Man and Cybernetics - Part B, vol. 29, n.2, 1999.
- [Weng92] Weng, Juyang, Paul Cohen, and Marc Herniou. *Camera Calibration with Distortion Models and Accuracy Evaluation*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 14, n.10, 1992.