

Automatic Thresholding for Correspondence Detection

Kenichi Kanatani* and Yasushi Kanazawa†

*Department of Information Technology, Okayama University, Okayama 700-8530 Japan

†Department of Knowledge-based Information Engineering

Toyohashi University of Technology, Toyohashi, Aichi 441-8580 Japan

kanatani@suri.it.okayama-u.ac.jp, kanazawa@tutkie.tut.ac.jp

Summary

We study the problem of thresholding the residual of template matching for selecting correct matches between feature points detected in two images. In order to determine the threshold dynamically, we introduce a statistical model of the residual and compute an optimal threshold according to that model. The model parameters are estimated from the histogram of the residuals of candidate matches. We demonstrate the effectiveness of our scheme using real images.

1. Introduction

Establishing point correspondences over multiple images is the first step of many video processing applications. Two approaches exist for this purpose: tracking correspondences over successive frames, and direct matching between separate frames. This paper focuses on the latter.

The basic principle is local correlation measurement by template matching. Detecting feature points in the first and second images separately using a corner detector [3, 8], we measure the correlation between the neighborhoods of the two points for each candidate pair and match those that have a high correlation.

To do so, we need to set an appropriate threshold for distinguishing correct matches from incorrect ones. This problem has not been considered fully in the past, chiefly because template matching alone is insufficient for establishing point correspondences; an outlier removal technique, such as LMedS [7] and RANSAC [2], needs to be applied afterwards. The thresholding task is usually passed on to the outlier removal stage [1, 9].

However, most outlier removal techniques do not work if the outlier ratio is as high as 50%. Hence, setting a good threshold at the template matching stage for removing incorrect matches is essential for the subsequent outlier removal procedure to be effective. If the threshold is too high, however, a lot of correct matches are lost. This reduces the number of final matches, making subsequent computations less reliable. Hence, we need a good balance between removing incorrect matches and retaining correct ones.

For this purpose, we introduce a statistical model of the template matching residual and compute an op-

timal threshold according to that model. The model parameters are estimated from the histogram of the residuals of candidate matches. We demonstrate the effectiveness of our method using real images.

2. Template matching

After feature points are detected from the two images using a corner detector, the similarity between point P in the first image and point Q in the second is measured by the following *residual* (sum of squares)

$$J(P, Q) = \sum_{(i,j) \in \mathcal{N}} |T_P(i, j) - T_Q(i, j)|^2, \quad (1)$$

where $T_P(i, j)$ and $T_Q(i, j)$ are the intensity values of the templates obtained by cutting out an $n \times n$ pixel region \mathcal{N} centered on P and Q , respectively. The following argument can be extended to color images and other measures such as the normalized correlation.

Basically, each point P in the first image is matched to the point Q in the second for which $J(P, Q)$ is the smallest, but overlaps and conflicts must be resolved. So, we apply the following uniqueness enforcing procedure. We first choose the pair (P, Q) that has the smallest residual $J(P, Q)$. Then, we remove from the candidate pairs those that involve P and Q . From the remaining pairs, we choose the pair (P', Q') that has the smallest residual $J(P', Q')$. We repeat this procedure until all pairs are exhausted.

3. Thresholding the residual

If the above uniqueness enforcing procedure is applied to all the pairs, we may end up with matching points for which no counterparts exist. So, we need to remove beforehand those pairs for which the residual is very large. For this, the threshold is usually set empirically [1, 9]. For example, Zhang et al. [9] accepted those pairs for which the normalized correlation is larger than 0.8.

However, the threshold cannot be fixed, because the residual $J(P, Q)$ is determined not only by the image intensity fluctuations but also by the relative distortion of the two images. For example, but if rotation and

scale change exist between the two images, the residual $J(P, Q)$ is not zero even when no image noise exists. It follows that the threshold should depend on the degree of the rotation and scale change, which is unknown and different from image to image.

Our strategy is that we introduce a statistical model of the residual and compute an optimal threshold according to that model. The model parameters are estimated from the histogram of the residual.

4. Statistical model of the residual

If the match (P, Q) is correct, the image intensity difference

$$\Delta T_{ij} = T_P(i, j) - T_Q(i, j) \quad (2)$$

in eq. (1) is due to the relative distortion in the neighborhoods of P and Q , as well as random fluctuations of image intensity. We model these by a Gaussian distribution of mean 0 and standard deviation σ_0 . Then, J/σ_0^2 for a correct match is subject to a χ^2 distribution with n^2 degrees of freedom if the intensity difference is independent of the pixel (n is the template size).

If the match (P, Q) is incorrect, the difference ΔT_{ij} is due to the inhomogeneity of the intensity within the image of that scene. We assume that ΔT_{ij} is subject to a Gaussian distribution of mean 0 and standard deviation σ_1 . Then, J/σ_1^2 for an incorrect match is subject to a χ^2 distribution with n^2 degrees of freedom if the intensity difference is independent of the pixel.

Let $f_0(J)$ be the probability density of the residual J for correct matches, and $f_1(J)$ that for incorrect ones. According to the above model, we have

$$f_0(J) = \frac{1}{\sigma_0^2} \phi_{n^2}\left(\frac{J}{\sigma_0^2}\right), \quad f_1(J) = \frac{1}{\sigma_1^2} \phi_{n^2}\left(\frac{J}{\sigma_1^2}\right), \quad (3)$$

where $\phi_d(x)$ denotes the probability density of the χ^2 distribution with d degrees of freedom.

5. Effective template size

The assumption that the image intensity difference is independent of the pixel is not realistic. However, exactly modeling the interpixel correlation is difficult, so we introduce the following approximation.

If there are N points in the first image and M points in the second, the number of correct matches is at most $\min(N, M)$, which is much smaller than the number NM of all the pairs. If most of the matches are incorrect, the probability density of the residual J for all the matches is approximately $f_1(J)$, which has an expectation of $n^2\sigma_1^2$ and a variance of $2n^2\sigma_1^4$. It follows that if we compute the mean μ_J and the variance σ_J^2 of J from the histogram of J , we should have

$$\mu_J \approx n^2\sigma_1^2, \quad \sigma_J^2 \approx 2n^2\sigma_1^4, \quad (4)$$

provided each pixel value is independent. Eliminating σ_1 from these, we obtain $n^2 \approx 2\mu_J^2/\sigma_J^2$. However, n^2

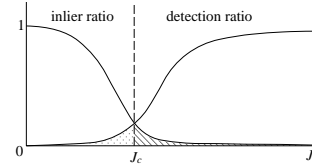


Figure 1: Determining the threshold that balances the inlier ratio and the detection ratio.

should be much smaller than this due to correlations. So, we define the *effective template size* by

$$n = \frac{\sqrt{2}\mu_J}{\sigma_J}. \quad (5)$$

In other words, we regard each pixel value as if independent within the template of this size, which need not be an integer.

6. Model parameter estimation

Let p and $q (= 1 - p)$ be the ratios of the correct and incorrect matches, respectively. The probability density of the residual J for all the matches is

$$f(J) = pf_0(J) + qf_1(J) = \frac{J^{n^2/2-1}}{2^{n^2/2}\Gamma(n^2/2)} \left(\frac{pe^{-J/2\sigma_0^2}}{\sigma_0^2} + \frac{qe^{-J/2\sigma_1^2}}{\sigma_1^2} \right). \quad (6)$$

We determine the model parameters σ_0 and σ_1 by maximum likelihood estimation. Let $J_1 \leq J_2 \leq \dots \leq J_{NM}$ be all the NM residual values sorted in ascending order. From eq. (6), their likelihood is

$$\prod_{i=1}^{NM} f(J_i) = \frac{\prod_{i=1}^{NM} J_i^{n^2/2-1}}{2^{n^2 NM/2} \Gamma(n^2/2)^{NM}} \times \prod_{i=1}^{NM} \left(\frac{pe^{-J_i/2\sigma_0^2}}{\sigma_0^2} + \frac{qe^{-J_i/2\sigma_1^2}}{\sigma_1^2} \right). \quad (7)$$

Differentiating the logarithm of this with respect to σ_0^2 and σ_1^2 and letting the results be zero, we obtain

$$\sigma_0^2 = \frac{\sum_{i=1}^{NM} A_i J_i}{n^2 \sum_{i=1}^{NM} A_i}, \quad \sigma_1^2 = \frac{\sum_{i=1}^{NM} B_i J_i}{n^2 \sum_{i=1}^{NM} B_i}, \quad (8)$$

where we define

$$A_i = \frac{1}{1 + \frac{q}{p} \left(\frac{\sigma_0}{\sigma_1}\right)^{n^2} e^{\frac{J_i}{2} \left(\frac{1}{\sigma_0^2} - \frac{1}{\sigma_1^2}\right)}}, \quad B_i = \frac{1}{1 + \frac{p}{q} \left(\frac{\sigma_1}{\sigma_0}\right)^{n^2} e^{\frac{J_i}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2}\right)}}. \quad (9)$$

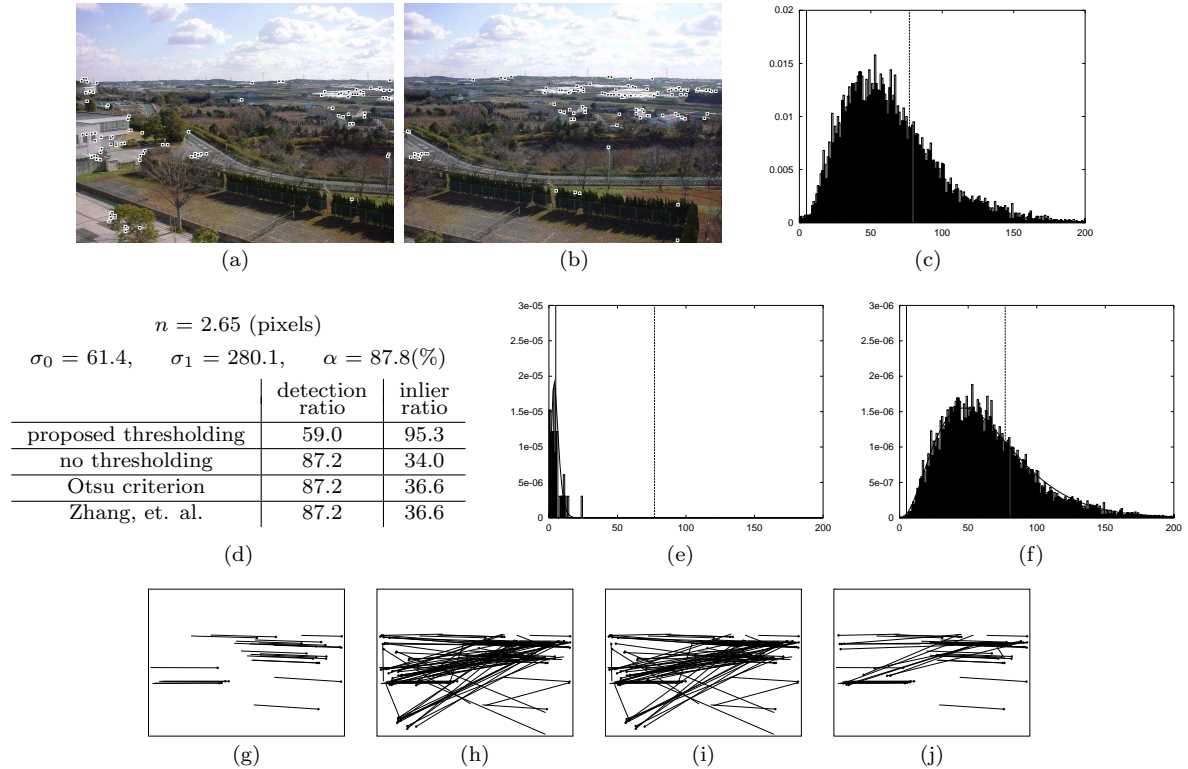


Figure 2: (a), (b) Input images and detected feature points. (c) The residual histogram of all the matches. (d) The model parameters, the detection ratio (%), and the inlier ratio (%). (e) The residual histogram of correct matches and the estimated density. (f) The residual histogram of incorrect matches and the estimated density. (g) Matches resulting from the proposed thresholding. (h) Matches without thresholding. (i) Matches resulting from the Otsu thresholding. (j) Matches resulting from the method of Zhang, et al. The vertical solid lines in (c), (e), and (f) indicate the threshold determined from the computed detection ratio α . The vertical dotted lines indicate thresholds obtained by the Otsu criterion.

The values of σ_0 and σ_1 are obtained by iterations: guessing the initial values to be, for example,

$$\sigma_0 = \sqrt{\frac{\sum_{i=1}^{\lfloor pNM \rfloor} J_i}{n^2 \lfloor pNM \rfloor}}, \quad \sigma_1 = \frac{\sigma_J}{\sqrt{2\mu_J}}, \quad (10)$$

and substituting them into the right-hand sides of eqs. (8), we obtain their updated values. This process is repeated until σ_0 and σ_1 converge. The first of eqs. (10) is the value of σ_0 we would have if the $\lfloor pNM \rfloor$ matches with the smallest residuals were all correct. The second of eqs. (10) is obtained by eliminating n^2 from eqs. (4).

The ratios p and q ($= 1 - p$) are given empirically. Since the number of correct matches between N points and M points is at most $\min(N, M)$, we let $p_{\max} = \min(N, M)/NM$ and set, for example, $p = 0.6p_{\max}$ if no knowledge is available about the correctness of the matches. *This estimate of p need not be precise*, as we will show later.

7. Detection ratio vs. inlier ratio

Suppose we set a threshold J_c for the residual J and accept those matches with $J \leq J_c$ as correct. Let α be the ratio of the accepted correct matches among all the

correct ones; we call it the *detection ratio*. A correct match with residual J is accepted with the probability

$$\alpha = P_0[J < J_c] = P_0\left[\frac{J}{\sigma_0^2} < \frac{J_c}{\sigma_0^2}\right], \quad (11)$$

where $P_0[\cdot]$ denotes the probability for correct matches. Let $\chi_{n^2}^2(\alpha)$ be the α th percentile of the χ^2 distribution with n^2 degrees of freedom. Since J/σ_0^2 for a correct match is subject to a χ^2 distribution with n^2 degrees of freedom, eq. (11) implies that J_c/σ_0^2 equals $\chi_{n^2}^2(\alpha)$. Hence, the threshold J_c is given by

$$J_c = \sigma_0^2 \chi_{n^2}^2(\alpha). \quad (12)$$

Some incorrect matches are necessarily accepted by this thresholding. An incorrect match with residual J is accepted with the probability

$$\gamma = P_1[J \leq J_c] = P_1\left[\frac{J}{\sigma_1^2} \leq \left(\frac{\sigma_0}{\sigma_1}\right)^2 \chi_{n^2}^2(\alpha)\right], \quad (13)$$

where $P_1[\cdot]$ denotes the probability for incorrect matches. Let $\Phi_{n^2}(X)$ ($= \int_0^X \phi_{n^2}(x) dx$) be the accumulated probability function of the χ^2 distribution with

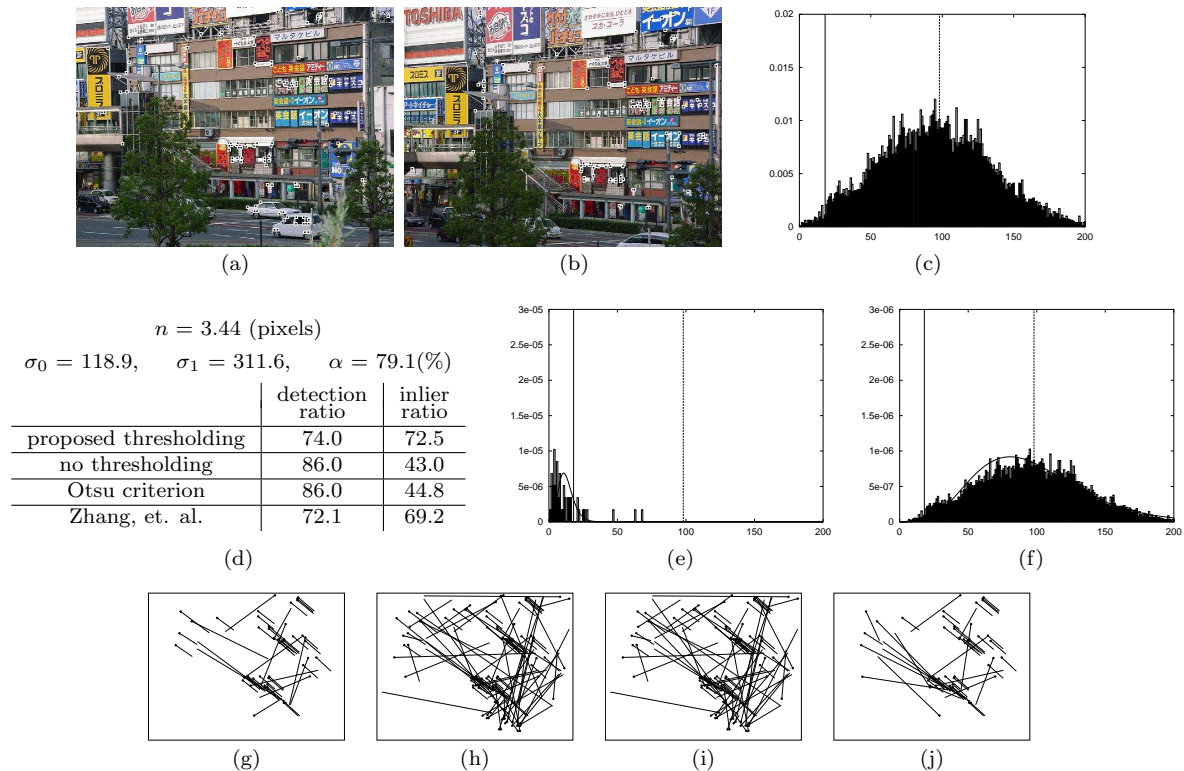


Figure 3: (a), (b) Input images and detected feature points. (c) The residual histogram of all the matches. (d) The model parameters, the detection ratio (%), and the inlier ratio (%). (e) The residual histogram of correct matches and the estimated density. (f) The residual histogram of incorrect matches and the estimated density. (g) Matches resulting from the proposed thresholding. (h) Matches without thresholding. (i) Matches resulting from the Otsu thresholding. (j) Matches resulting from the method of Zhang, et al. The vertical solid lines in (c), (e), and (f) indicate the threshold determined from the computed detection ratio α . The vertical dotted lines indicate thresholds obtained by the Otsu criterion.

n^2 degrees of freedom. Since J/σ_1^2 for an incorrect match is subject to a χ^2 distribution with n^2 degrees of freedom, eq. (13) implies

$$\gamma = \Phi_{n^2}\left(\left(\frac{\sigma_0}{\sigma_1}\right)^2 \chi_{n^2}^2(\alpha)\right). \quad (14)$$

Among the NM possible matches, the numbers of correct and incorrect matches are pNM and qNM , respectively. After the thresholding, we obtain αpNM correct matches and γqMN incorrect ones on average. Hence, the *inlier ratio*, i.e., the ratio of correct matches among the accepted matches, is approximately

$$\beta = \frac{\alpha pNM}{\alpha pNM + \gamma qMN} = \frac{\alpha p}{\alpha p + \gamma q}. \quad (15)$$

8. Threshold selection

The threshold J_c is determined by the detection ratio α in the form of eq. (12), but how should we set α ? It should be large if we want to collect many correct matches, but the number of incorrect matches also increases, lowering the inlier ratio β as a result.

Here, we determine the threshold J_c so that the detection ratio α equals the inlier ratio β . This balances

the ratio $1 - \alpha$ of rejecting correct matches and the ratio $1 - \beta$ of accepting incorrect ones (Fig. 1). Letting $\beta = \alpha$ in eq. (14), we obtain

$$\alpha = 1 - \frac{q}{p} \Phi_{n^2}\left(\left(\frac{\sigma_0}{\sigma_1}\right)^2 \chi_{n^2}^2(\alpha)\right), \quad (16)$$

from which α is obtained by Newton iterations.

9. Real image examples

Figs. 2(a) and (b) show two real images of a distant scene. We detected 100 feature points from each image using the Harris operator [3], as marked there. Fig. 2(c) is the histogram of the residuals of all candidate matches we obtained using a 9×9 template. Letting $p/p_{\max} = 0.6$, we estimated the effective template size n , the model parameters σ_0 and σ_1 , and the optimal detection ratio α (= the inlier ratio β) as listed in Fig. 2(d). We see that n is much smaller than the actual size 9 due to correlations. The density $f(J)$ of the residual J estimated by eq. (6) is superimposed onto the histogram in Fig. 2(c). The estimated density agrees with the histogram very well.

Figs. 2(e) and (f) superimpose the densities $f_0(J)$ and $f_1(J)$ of correct and incorrect matches estimated

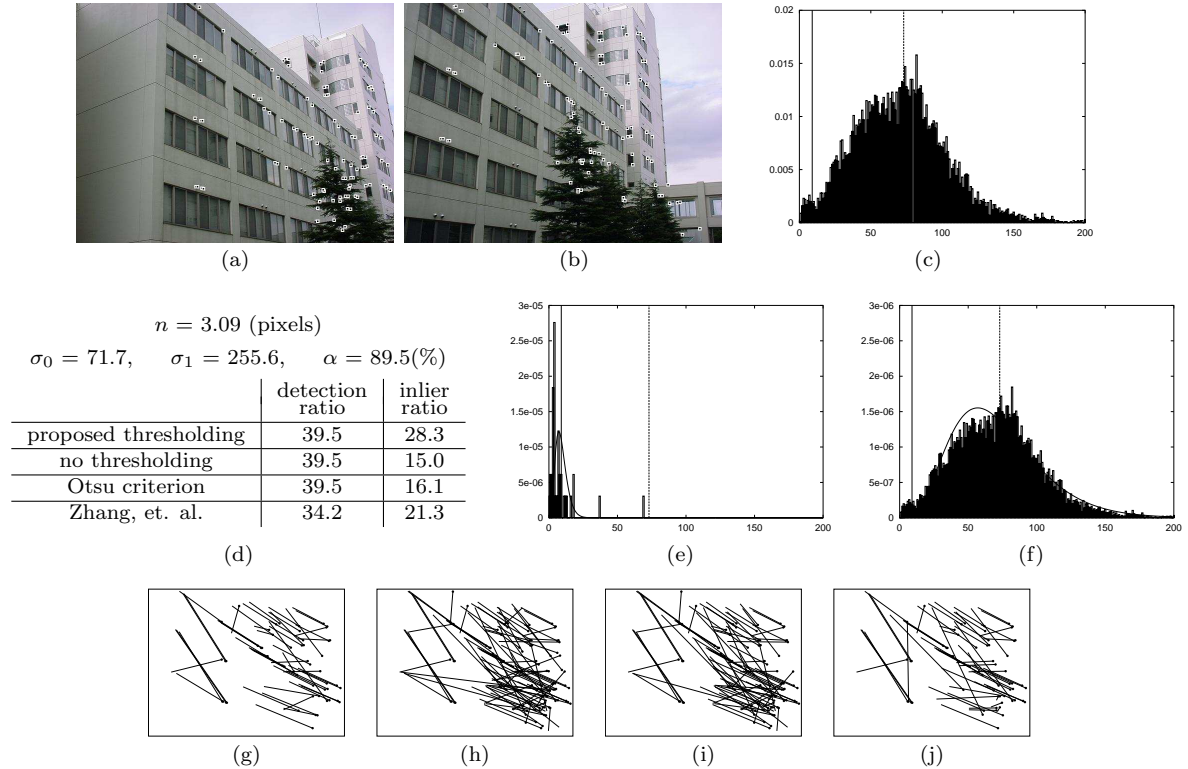


Figure 4: (a), (b) Input images and detected feature points. (c) The residual histogram of all the matches. (d) The model parameters, the detection ratio (%), and the inlier ratio (%). (e) The residual histogram of correct matches and the estimated density. (f) The residual histogram of incorrect matches and the estimated density. (g) Matches resulting from the proposed thresholding. (h) Matches without thresholding. (i) Matches resulting from the Otsu thresholding. (j) Matches resulting from the method of Zhang, et al. The vertical solid lines in (c), (e), and (f) indicate the threshold determined from the computed detection ratio α . The vertical dotted lines indicate thresholds obtained by the Otsu criterion.

by eqs. (3) onto the histograms of correct and incorrect matches, separately. Here, we checked the correctness of the matches as follows.

Since two images of a distant scene are related by a *homography*, we optimally computed the homography \mathcal{H} by *renormalization*¹ [4] from a large number of corresponding points selected by hand. For each candidate match (P, Q) , we mapped the point P to the second image by the computed homography \mathcal{H} and judged the match as correct if the point Q is within three pixels from its ideal position $\mathcal{H}P$. The result agrees very well with our prediction.

The vertical solid lines in Figs. 2(c), (e), and (f) indicate the computed threshold J_c . A well known scheme for automatic thresholding is the *Otsu discrimination criterion* [6]. The vertical dotted lines in Figs. 2(c), (e), and (f) indicate the corresponding threshold.

Fig. 2(g) shows the final matches obtained by applying the computed threshold J_c followed by the uniqueness enforcing procedure; they are displayed as “optical flow” (line segments connecting the matching positions). For comparison, Fig. 2(h) shows the result

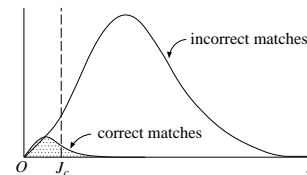


Figure 5: The residual distribution of correct matches is included in the residual distribution of incorrect matches to a large extent.

without thresholding; Fig. 2(i) is the result using the Otsu criterion; Fig. 2(j) is the result thresholded by the normalized correlation 0.8 according to Zhang et al. [9]. The actual detection ratio α and the inlier ratio β for these results are listed in Fig. 2(d).

We can see that without thresholding we can collect many correct matches but we also pick out many incorrect ones. As a result, the inlier ratio significantly drops. Our thresholding scheme balances the conflicting goals of collecting as many correct matches as possible and rejecting as many outliers as possible, using the knowledge of the residual distribution of the images in question. The Otsu criterion shows little effect for this, and the method of Zhang, et al. [9] gives an intermediate result between the Otsu criterion and our

¹The program is publicly available from <http://www.ail.cs.gunma-u.ac.jp/Labo/e-programs.html>.

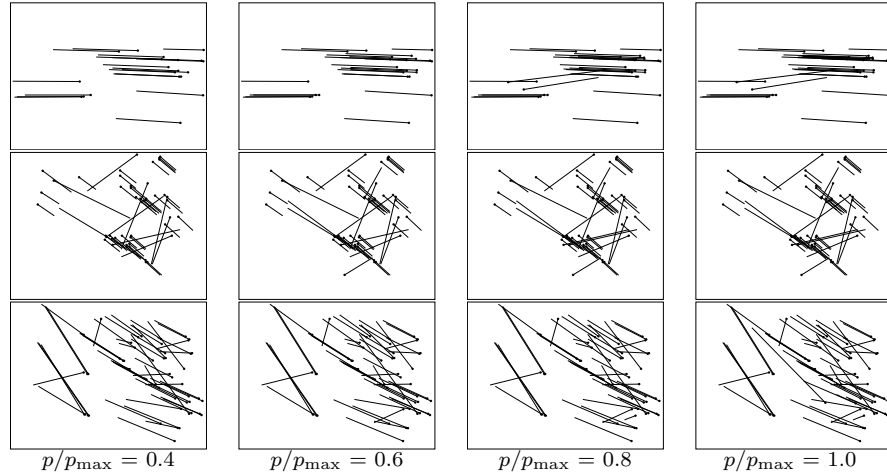


Figure 6: Matches resulting from different estimates of the ratio p of correct matches for the images in Figs. 2, 3, and 4 (from above).

method.

Figs. 3 and 4 show other examples with corresponding results. From these, we can confirm the effectiveness of our method.

10. Validation of our model

From the above experiments, we can see that the residual distribution of correct matches is included in the residual distribution of incorrect matches to a large extent with a long tail (Fig. 5). So, if we want to pick out a large number of correct matches, we must set a high threshold, which inevitably accepts many incorrect matches. Our analysis is for determining an optimal threshold by analyzing the distribution shapes.

For this purpose, we need to approximate the distributions by simple functions. Here, we fitted the χ^2 distribution density by adjusting the effective template size n and the parameters σ_0 and σ_1 and obtained a satisfactory fit. *As long as the fit is good, the underlying statistical hypotheses (independent Gaussian distributions, etc.) are not essential.*

We set the ratio p/p_{\max} to be 0.6, but this estimate need not be precise. Fig. 6 shows the final results for $p/p_{\max} = 0.4, 0.6, 0.8, 1.0$ for the images in Figs. 2, 3, and 4. The results are not so different, so we may set p/p_{\max} to be, for example, around 0.6 if no prior information is given.

11. Conclusion

We have studied the problem of thresholding the residual of template matching for selecting correct matches between feature points detected in two separate images. We dynamically determined the threshold by introducing a statistical model of the residual and computed an optimal threshold according to that model. The model parameters were estimated from the histogram of the residual of candidate matches. We

demonstrated the effectiveness of our scheme using real images.

Of course, template matching alone is not sufficient for practical applications. We must also incorporate outlier techniques such as LMedS [7] and RANSAC [2], exploiting geometric constraints such as the homography relationship and the epipolar equation. Our scheme is very effective as a preprocess for such outlier removal techniques (see, e.g., [5] for an image mosaicing application).

Acknowledgements: This work was supported in part by the Ministry of Education, Culture, Sports, Science and Technology, Japan, under a Grant in Aid for Scientific Research C(2) (No. 13680432), the Support Center for Advanced Telecommunications Technology Research, and Kayamori Foundation of Informational Science Advancement.

References

- [1] P. Beardsley, P. Torr and A. Zisserman, 3D model acquisition from extended image sequences, *Proc. 4th Euro. Conf. Comput. Vision*, April 1996, Cambridge, U.K., Vol. 2, pp. 683–695.
- [2] M. A. Fischler and R. C. Bolles, Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography, *Comm. ACM*, **24-6** (1981), 381–395.
- [3] C. Harris and M. Stephens, A combined corner and edge detector, *Proc. 4th Alvey Vision Conf.*, August 1988, Manchester, U.K., pp.147–151.
- [4] K. Kanatani, N. Ohta, and Y. Kanazawa, Optimal homography computation with a reliability measure, *IEICE Trans. Inf. & Syst.*, **E83-D-7** (2000), 1369–1374.
- [5] Y. Kanazawa and K. Kanatani, Image mosaicing by stratified matching, *Workshop on Statistical Methods in Video Processing*, June 2002, Denmark, Copenhagen.
- [6] N. Otsu, A threshold selection method from gray-level histograms, *IEEE Trans. Sys. Man Cyber.*, **9-1** (1979), 62–66.
- [7] P. J. Rousseeuw and A. M. Leroy, *Robust Regression and Outlier Detection*, Wiley, New York, 1987.
- [8] S. M. Smith and J. M. Brady, SUSAN—A new approach to low level image processing, *Int. J. Comput. Vision*, **23-1** (1997), 45–78.
- [9] Z. Zhang, R. Deriche, O. Faugeras and Q.-T. Luong, A robust technique for matching two uncalibrated images through the recovery of the unknown epipolar geometry, *Artif. Intell.*, **78** (1995), 87–119.