

Recognition of Degraded Machine-Printed Characters Using a Complementary Similarity Measure and Error-Correction Learning

Minako SAWAKI[†] and Norihiro HAGITA[†], Members

SUMMARY Most conventional methods used in character recognition extract geometrical features, such as stroke direction and connectivity, and compare them with reference patterns in a stored dictionary. Unfortunately, geometrical features are easily degraded by blurs and stains, and by the graphical designs such as used in Japanese newspaper headlines. This noise must be removed before recognition commences, but no preprocessing method is perfectly accurate. This paper proposes a method for recognizing degraded characters as well as characters printed on graphical designs. This method extracts features from binary images, and a new similarity measure, the *complementary similarity measure*, is used as a discriminant function; it compares the similarity and dissimilarity of binary patterns with reference dictionary patterns. Experiments are conducted using the standard character database ETL-2, which consists of machine-printed Kanji, Hiragana, Katakana, alphanumeric, and special characters. The results show that our method is much more robust against noise than the conventional geometrical-feature method. It also achieves high recognition rates of over 97% for characters with textured foregrounds, over 99% for characters with textured backgrounds, over 98% for outline fonts and over 99% for reverse contrast characters. The experiments for recognizing both the fontstyles and character category show that it also achieves high recognition rates against noise.

key words: character recognition, error-correction learning, similarity measure, noise model, fontstyle recognition

1. Introduction

In mainstream of character recognition, geometrical-feature methods that extract structural features, such as stroke direction and connectivity, have been studied for many years [1]. These methods, however, are easily influenced by degradation (or noise) such as blurs or stains. In order to apply these methods, preprocessing is required to remove the noise before recognition begins. Existing preprocessing techniques, however, have difficulty in distinguishing character parts from noise for degraded characters such as those found in faxed messages and for characters with the graphical designs used, for instance, in Japanese newspaper headlines (Fig. 1).

Sakou et al. [2] attempted to remove the textured backgrounds using a textured segregation approach, while Okamoto et al. [3] proposed a morphological approach for removing them. They assumed that the stroke widths in the characters are greater than line widths in the textured backgrounds and the connected area of each character image

政策に生活者の実感 必要

田中 真紀子・科技庁長官に聞く

(a)

住宅公庫 金利 12月半ばにも上げ

(b)

Fig. 1 Examples of textlines with graphical designs in Japanese newspapers. (a) Upper: textured characters; lower: reverse contrast characters. (b) Left: black-plain characters; right: outline font characters with textured backgrounds.

could be separated from the background. These assumptions, however, do not always hold because some strokes in complicated Kanji characters such as 設 and 証 are often narrower than the lines in the textured backgrounds, and these strokes are erroneously removed. In addition, character parts may not be correctly extracted when multiple characters are connected due to underlines or smears. On the other hand, Ozawa et al. [4] focused on the gray-scale images of headlines rather than on binary images to enhance the contrast between character parts and backgrounds. With all of these methods, however, the removal of textured backgrounds through image processing often extracts incomplete areas of the character parts. Moreover, headlines usually include not only black-plain characters but also other kinds of degraded characters: reverse contrast characters, outline characters, textured characters, and various mixtures of each, as shown in Fig. 1. Therefore, a segmentation and recognition method without removing graphical designs will be needed for successful recognition.

This paper focuses on a method for recognizing these degraded characters without removing graphical designs. We assume that either character height or width is detected before recognition, since the proposed method uses whole binary image as features to avoid the drawback of geometrical features for degradation. This assumption could be satisfied by using features that reflect the difference between character parts and background parts, since only detecting the character height or width would be easier than

Manuscript received November 1, 1995.

Manuscript revised January 22, 1996.

[†] The authors are with the Basic Research Laboratories of NTT, Atsugi-shi, 243-01 Japan.

removing graphical designs. A new similarity measure, the *complementary similarity measure*, is proposed and is applied to the error-correction learning to make reference patterns. Experimental results confirm the method's robustness. It also achieves a high accuracy for fontstyle and category recognition.

2. Binary Image Feature Method Using Complementary Similarity Measure

2.1 Basic Framework of Proposed Method

Let X be an n -pixel binary image normalized in terms of size and position. We assume that X retains a distinctive shape, such as 王 for discriminating 微 from 微 (namely, $n \geq 1024$ for Kanji characters). Pattern matching based on binary image features often uses m -multivalued features ($m \ll n$, namely $m = 8 \times 8 = 64$) in order to save computation time; a binary image is divided into m local areas and m -multivalued features are obtained by accumulating the number of black pixels in each local area. The downsampling of n to m , in general, deteriorates class separability because the shapes that distinguish similar characters, such as 微-微, 思-恩, 由-田-甲-旧-困-因, are not retained and results in lower recognition rates than geometrical feature methods [1]. We, therefore, attempt to use X as a feature vector to avoid the problem of downsampling and maintain a high degree of separability against noise. This technique is feasible with existing computers.

In pattern matching, the conventional similarity measure $S(X, M)$ is often used as the discriminant function:

$$s(X, M) = \frac{\langle X, M \rangle}{\|X\| \cdot \|M\|} \quad (1)$$

Here, $\langle \cdot, \cdot \rangle$ and $\|\cdot\|$ denote the inner product and norm, and X and M are the input feature vector and the mean vector in a category.

The measure is theoretically affected by additive noise. That is, since the increase in $\|X\|$ due to additive noise increases the value of the similarity between X and M of a category with larger $\|M\|$, X is always assigned to the category.

Moreover, binary image features commonly require a lot of reference patterns, so-called subcategories, to handle multiple fontstyles of characters. We, therefore, use *binary* reference pattern T in place of the multivalued M when generating the reference patterns for the dictionary. This simplifies the addition of new (sub)categories, correction of erroneous samples and ensures dictionary compactness. Note that M is rarely an element but T is always an element on a finite set of X .

The error-correction learning is also introduced to ensure the dictionary compactness.

2.2 Complementary Similarity Measure

Various similarity measures for n -Boolean feature space are summarized in [5]. The measures between an input pattern $X = (x_1, x_2, \dots, x_i, \dots, x_n)$ and a binary reference pattern $T = (t_1, t_2, \dots, t_i, \dots, t_n)$ are based on four parameters, a, b, c and e , which are given by

$$a = \sum_{i=1}^n x_i \cdot t_i, \quad b = \sum_{i=1}^n (1-x_i) \cdot t_i, \quad (2)$$

$$c = \sum_{i=1}^n x_i \cdot (1-t_i), \quad e = \sum_{i=1}^n (1-x_i) \cdot (1-t_i).$$

Note that the following relation always holds:

$$a + b + c + e = n,$$

$$X = \|X\|^2 = a + c, \quad T = \|T\|^2 = a + b. \quad (3)$$

where X and T denote the number of black pixels in X and T , respectively.

These four parameters have the following properties for noise:

- (1) When $X = T$, $a = T$, $b = c = 0$, $e = n - T$.
- (2) When X comes from T with additive noise, $a = T$, $b = 0$, $c > 0$ and $e < n - T$.
- (3) When X comes from T with deletion noise, $a < T$, $b > 0$, $c = 0$ and $e = n - T$.

Using the complementary relationship between a and e (or b and c) for noise, we define a similarity measure S_c , called the *complementary similarity measure*, as follows:

$$S_c(X, T) = \frac{a \cdot e - b \cdot c}{\sqrt{T \cdot (n - T)}} \quad (4)$$

$$= \frac{n \cdot a - T \cdot X}{\sqrt{T \cdot (n - T)}}. \quad (5)$$

When X comes from T with noise, term $b \cdot c$ is always equal to zero and S_c may have a high similarity value for the category of T . S_c in Eq. (4) can also be rewritten as Eq. (5) using n, T, X and a . S_c requires less computational time than the so-called four-fold point correlation efficient S_f in Eq. (6), since $\sqrt{X(n-X)}$ is a constant for a certain X .

$$S_f(X, T) = \frac{a \cdot e - b \cdot c}{\sqrt{X \cdot (n - X) \cdot T \cdot (n - T)}} \quad (6)$$

For the reverse contrast image X^c of X , from Eqs. (3) and (4),

$$S_c(X^c, T) = -S_c(X, T) \quad (7)$$

When X and X^c should be assigned to the same category, the absolute value $|S_c|$ of S_c may be used as a discriminant function. Though a conventional similarity measure S also realizes $S(X, T) = S(X^c, T)$ by applying definite

canonicalization to X [6], the proposed measure $|S_c|$ requires more simple process than this.

3. Error-correction Learning in Designing Binary Reference Patterns

Binary reference patterns are obtained from a learning sample set in a (sub)category. Let m_i and m be the mean value at x_i and the average of all m_i in the learning sample set, respectively. The i -th pixel value t_i in T is determined as $t_i = 1$ when $m_i \geq m$ and $t_i = 0$ when $m_i < m$.

Now consider a learning algorithm that iteratively modifies the hypotheses subsets so as to correct error samples in the previous round. The algorithm adopted here is based on the “*self-corrective recognition algorithm*” reported by Nagy et al. [7] except that we do not take the rejection threshold into account. In the first iteration cycle, one binary reference pattern per category is formed using a learning sample set. The initial reference patterns are obtained from the mean vectors over all samples in each category. In each cycle, the similarity measure S_c is used as the discriminant function. Samples in a category are recognized and disjointly fall into cluster(s) of correctly recognized patterns which corresponds to each binary reference pattern and/or a cluster of misrecognized patterns of the category. New binary reference patterns are generated from the samples in each cluster of correctly recognized patterns. A cluster containing no sample is discarded. In the misrecognized cluster in the category, a sample that has the highest similarity to the binarized mean pattern is selected as a binary reference pattern. This ensures learning convergence. Iteration stops when no further error sample is found.

4. Recognition Experiments

4.1 Character Data

The machine-printed character database ETL-2 was used for the recognition experiments. Its contents are listed in Table 1. The database is open to the public and the characters in it are one of four fontstyles: the mincho font used in newspapers, the mincho font in patents, the gothic font in newspapers, and the gothic font in patents. Five characters per category were used as learning samples for designing binary reference patterns and the remaining five characters per category were test samples for each font style. A gray-scale character image consisting of 60×60 pixels was binarized and then the binary image was normalized in size and position as follows. The circumscribed rectangle of the black region in each 60×60 -pixel image was detected and the center of the rectangle was moved to the center of the 60×60 -pixel image. The shifted image was then expanded or shrunk to a 32×32 -pixel image X . Figure 2 shows all normalized binary images of a category “示” in the gothic font used in newspapers.

The error-correction learning algorithm mentioned in

Table 1 The ETL-2 database.

Font style	No. of categories	Kinds of character
Mincho newspaper	1895	Kanji, Hiragana, Katakana, Alphanumeric, Symbol
Mincho patent	2182	Kanji, Hiragana, Katakana, Alphanumeric, Symbol
Gothic newspaper	571	Kanji, Hiragana, Symbol
Gothic patent	571	Kanji, Hiragana, Symbol

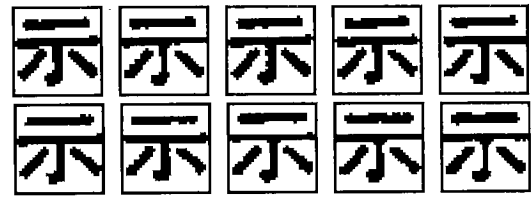


Fig. 2 Examples in the ETL-2 database (Kanji “示” in gothic fontstyle in newspapers).

chapter 3 was applied to the learning samples of noise-free characters for each font style. The similarity measure S_c was used as the discriminant function. The algorithm stopped after four iterations for the four font styles. The number of binary reference patterns formed was 1,901 for 1,895 categories in the mincho font in newspapers, 2,210 for 2,182 categories in the mincho font in patents, 578 for 571 categories in the gothic font in newspapers and 577 for 571 categories in the gothic font in patents. All learning samples were recognized correctly.

4.2 Noise Models

The analysis of noise produced by random bit-switching in a binary image is not important in developing a recognition method for degraded characters, since such noise is rarely found in practical character recognition [8],[9]. On the other hand, additive and deletion noise should be investigated. Noise pattern Z_α is produced by randomly setting a given percentage ($|\alpha|$ percent, $-100 \leq \alpha \leq 100$) of all black (or white) pixels for $\alpha < 0$ (or $\alpha > 0$) to white (or black). The deletion noise pattern X_α is formed by the AND-operation of X and Z_α when $\alpha \leq 0$, as shown in Fig. 3 (a), while the additive noise pattern X_α is yielded by the OR-operation of X and Z_α when $\alpha > 0$, as shown in Fig. 3 (b). That is,

$$x_{ai} = x_i \wedge z_{ai} \text{ for } \alpha \leq 0, \quad x_{ai} = x_i \vee z_{ai} \text{ for } \alpha > 0 \quad (8)$$

where x_{ai} and z_{ai} denote the Boolean values at the i -th pixel on X_α and Z_α , respectively.

4.3 Recognition Results for Noisy Images

Recognition experiments were performed to compare the robustness to noise of our method, a conventional geometrical-feature method, and a binary image feature method having conventional similarity measures. We used the PDC (Peripheral Direction Contributivity) feature vector [10] as geometrical features. The 1536-dimensional feature vector represents the stroke direction, stroke connectivity, stroke complexity and relative stroke location, and has been used in the recognition of Kanji characters. Noise was reduced using 3×3-pixel mask patterns before extracting the PDC feature vector from noisy images. Error-correction learning was also applied to the PDC feature vector except that the multivalued reference patterns were formed using the Euclidean distance. For a conventional similarity measure S , error-correction learning was also applied to make binary reference patterns using a similarity measure S .

Noise patterns X_α were obtained from X for each font style in the range of $-90 \leq \alpha \leq 90$. Figure 4 shows degradation characteristics of noise for test samples in the gothic font used in newspapers, and includes the recognition performance of one human. Each degraded character was observed for several seconds by one subject. A total of 579 multivalued reference patterns for the PDC feature vector were obtained from noise-free binary images of the learning samples. Also, a total of 586 binary reference patterns were obtained for a conventional similarity measure S . No pattern was rejected. Figure 4 shows that the complementary similarity measure S_c produces the most robust hypothesis subsets for noise out of these three methods. The recognition rate for the geometrical feature method rapidly decreases even when the deletion or additive noise is less than 30% and noise reduction is carried out. The performance of the binary image feature method with conventional similarity measure S also decreases rapidly when additive noise exceeds 40%. Several similarity measures in [5] were also tested in primary experiments. However, they achieve lower recognition accuracy than a conventional similarity measure S in all range of α . The recognition rate for our method exceeds 99% in the noise range of $-40 \leq \alpha \leq 80$. Particularly noteworthy is that our method achieves better recognition performance than the human subject in the noise ranges of $-90 \leq \alpha \leq -40$ and $60 \leq \alpha \leq 90$. However, the recognition rates in these two ranges may not be important in practical use, because most degraded characters and the graphical design patterns printed in the headlines of Japanese newspapers usually yield a noise range of $-50 \leq \alpha \leq 50$. Figure 5 shows several examples misrecognized by our method.

The proposed method is also the most robust of the three methods for the other three fontstyles in ETL-2. For example, it achieves over 99% recognition rates in the noise range of $-50 \leq \alpha \leq 70$ for the mincho font in patents, constituting 2,182 categories; $-70 \leq \alpha \leq 80$ for the mincho font in newspapers, constituting 1,895 categories; and $-50 \leq \alpha \leq 80$ for the gothic font in patents, constituting 571 categories.



(a) Deletion noise image X_α



(b) Additive noise image X_α

Fig. 3 Deletion and additive noise image models.

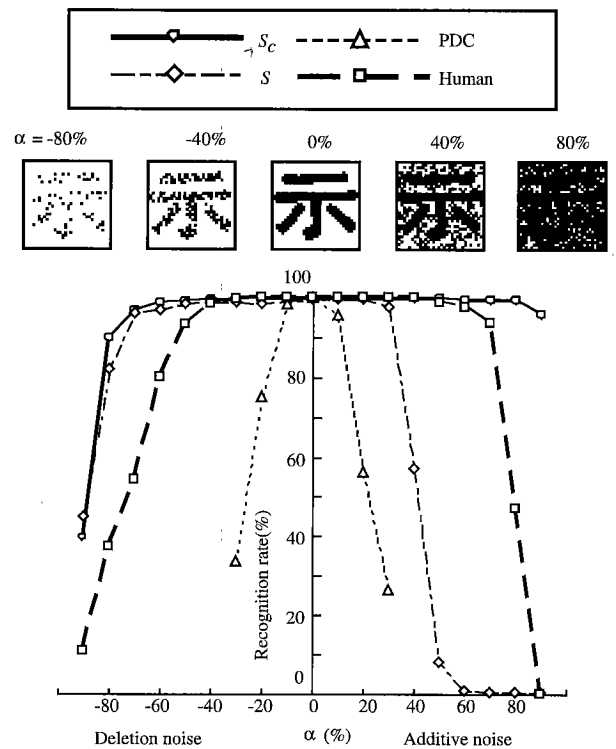


Fig. 4 Recognition rates for deletion and additive noise images.



Fig. 5 Examples of misrecognized noisy images.

4.4 Recognition Results for Characters with Textured Foregrounds and Backgrounds

The method was applied to the recognition of textured characters and characters with textured backgrounds, both of which are common in the headlines of Japanese newspapers. Textured character X_t was formed by the AND-operation of X and texture pattern Z , as shown in Fig. 6 (a), while a character with texture background X_b was formed by the

OR-operation of X and Z , as shown in Fig. 6 (b). Five kinds of textures were used in the experiment. The textures labeled T1 through T5 correspond to $|\alpha| = 50$ (%). Table 2 shows the recognition results of textured characters and characters with texture backgrounds obtained from test samples in the gothic font in newspapers. S_c was used as the discriminant function. Table 2 shows that our method achieves very high recognition rates of 97.97% to 99.79% for textured characters with the five textures and 99.26% to 99.75% for characters with textured backgrounds. Further error-correction learning against these error samples may reduce the error rate even further.

4.5 Recognition Results for Reverse Contrast Characters and Outline Font Characters

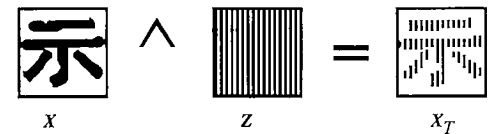
Recognition experiments for reverse contrast characters and outline font characters were also done. The contour pattern of a noise-free character image was used as an outline font character. Reverse contrast images were also obtained by reversing noise-free characters, outline font characters, textured characters, characters with textured backgrounds in terms of contrast. The absolute value $|S_c|$ of S_c was used as the discriminant function. The binary reference patterns obtained from noise-free characters were also used. Table 3 shows the recognition results for outline-font characters and the reverse contrast of those characters for test samples in the gothic font in newspapers. $|S_c|$ achieves high recognition rates of over 98% for reverse contrast images of original characters, synthesized outline font characters and their reverse contrast characters. Table 4 also shows the recognition results for reverse contrast images with textured foregrounds and backgrounds. The reverse contrast images were yielded by inverting the characters used in 4.4.

The results in Tables 2 to 4 seem to be almost similar to the recognition rates for deletion and additive noise images, as shown in Fig. 4. The recognition characteristics for deletion and additive noise in Fig. 4 would be utilized for estimating recognition rates of unknown texture patterns. Experiments with a conventional similarity measure S also showed similar characteristics.

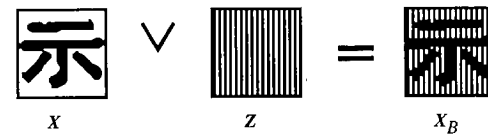
The high recognition rates shown in Tables 2 to 4 were also achieved for the three other font styles in ETL-2. These results, therefore, suggest that if character position is obtained in advance using the distinctive properties and knowledge of graphical backgrounds and characters, characters with textured backgrounds such as those in the headlines of Japanese newspapers can be directly recognized using our method. Such properties and common knowledge have already been used to segment sets of characters from textured backgrounds in document layout recognition [3],[4].

4.6 Recognition Results for Fontstyle and Category

Experiments for recognizing both the fontstyle and character category were also conducted. In order to eliminate effects



(a) Textured character



(b) Character with textured background

Fig. 6 Textured character and character with textured background.

Table 2 Recognition rates for textured characters and characters with textured backgrounds.

Texture	T1	T2	T3	T4	T5
X_T					
Recognition S_c rate (%)	99.79%	99.05	99.47	98.18	97.97
X_B					
Recognition S_c rate (%)	99.75%	99.72	99.26	99.58	99.61

Table 3 Recognition rates for outline characters and reverse contrast characters.

	Black-plain character	Reverse contrast image	Outline font	Reverse contrast outline font
Recognition $ S_c $ rate (%)	99.68%	99.82	98.74	98.74

Table 4 Recognition rates for reverse contrast characters with textured foreground and backgrounds.

Texture	T1	T2	T3	T4	T5
Reverse contrast image					
Recognition $ S_c $ rate (%)	99.68%	99.72	99.65	99.65	99.54
Reverse contrast image					
Recognition $ S_c $ rate (%)	99.61%	99.58	99.37	99.61	99.47

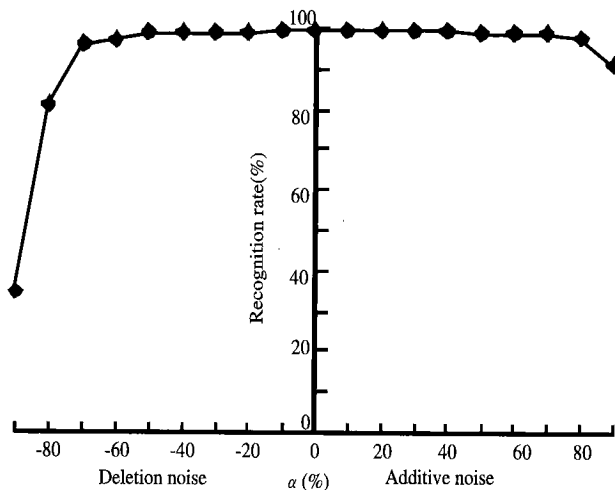


Fig. 7 Recognition rates for fontstyle and category.

of the total category number, 571 categories for each of 4 fontstyles were used in the experiments. Subsets of reference patterns corresponding to the 571 categories for each fontstyle were selected and then the subsets were concatenated as a reference pattern dictionary for these experiments. The dictionary consists of 574 (Mincho newspaper), 579 (Mincho patent), 578 (Gothic newspaper) and 577 (Gothic patent) reference patterns. Recognition experiments were carried out with the noise range of $-90 \leq \alpha \leq 90$. Experimental results are shown in Fig. 7. In these experiments, an input pattern was regarded as "correct" when its category and fontstyle were correctly recognized. The results show that over 99% of the patterns were correctly recognized in the noise range of $-40 \leq \alpha \leq 60$. The recognition rate of the category only was 99.93% at $\alpha = 0\%$. We also examined the recognition accuracy for the multiple-font recognition. Learning samples of these four fontstyles for each category were used in the error-correction learning. Experimental results show that over 98% of the patterns were correctly recognized in the noise range of $-30 \leq \alpha \leq 40$. These results suggest that this method may be applied to fontstyle identification and/or category recognition, if learning samples of fontstyles and categories are given.

5. Conclusion

This paper proposed a method for recognizing degraded characters through a simple process assuming that either character heights or widths are detected before recognition. This method is based on the use of the binary image features where the whole binary image is used as a feature. A new similarity measure, the *complementary similarity measure*, was used for robust recognition against noise. A set of binary reference patterns is learned using an error-correction learning algorithm.

Recognition experiments were performed to compare the robustness of this method, a conventional geometrical feature method and a binary image feature method having

conventional similarity measure. The results show that our method is the most robust against noise out of these three methods. Next, the methods were applied to four character recognition tasks: the recognition of outline fonts, reverse contrast characters, textured characters and characters with textured backgrounds. It also achieved high recognition rates for these characters. These results show that degraded characters can be directly recognized by this method. Experiments for recognizing fontstyles and character category were also conducted. The results show that this method also achieves high recognition rates for fontstyle identification and character recognition. In the future, we will investigate a method for extracting character heights or widths in a text-line image with graphical designs.

Acknowledgments

We would like to acknowledge the encouragement and support of Dr. K. Ishii and Dr. K. Odaka. We also thank the researchers at ETL for permitting the use of the database.

References

- [1] S.Mori, C.Y.Suen, and K.Yamamoto, "Historical review of OCR research and development," Proc. IEEE, vol.80, no.7, pp.1029-1058, 1992.
- [2] H.Sakou, H.Matsushima, and M.Ejiri, "Texture discrimination using self-organized multiresolution filtering," IEICE Trans., vol.J73-D-II, no.4, pp.562-573, April 1990 (in Japanese).
- [3] M.Okamoto and H.Hayashi, "Character extraction from headlines with background patterns by using shrinking/expanding methods," IEICE Technical Report, PRU90-151, 1990 (in Japanese).
- [4] H.Ozawa and T.Nakagawa, "A character image enhancement method from characters with various background images," Proc. of Second ICDAR, Tsukuba, Japan, pp.58-61, Oct. 1993.
- [5] R.S.Michalski and E.Diday, "A recent advance in data analysis: clustering objects into classes characterized by conjunctive concepts," Progress in Pattern Recognition, eds. L.N. Kanal and A.Rosenfeld, pp. 33-56, North-Holland Publishing Company, 1981.
- [6] T.Iijima, "Theory of pattern recognition," pp.78-79, Morikita Publishing Company, 1989 (in Japanese).
- [7] G.Nagy and G.L.Shelton, Jr., "Self-corrective character recognition system," IEEE Trans. on Information Theory, vol.IT-12, no.2, pp.215-222, 1966.
- [8] G.Nagy, "At the frontiers of OCR," Proc. IEEE, pp.1093-1100, July 1992.
- [9] T.Bayer, J.Hull, and G.Nagy, "Character recognition: SSPR'90 working group report., Structured Document Image Analysis," eds. H.S.Baird, H.Bunke, and K.Yamamoto, pp.565-567, Springer-Verlag, Tokyo, 1992.
- [10] T.Akiyama and N.Hagita, "Automated entry system for printed documents," Pattern Recognition, vol.23, no.11, pp.1141-1154, 1990.



Minako Sawaki received the B.E. degree in electrical engineering from Keio University in 1989. From 1989, she has been with the Nippon Telegraph and Telephone Corporation (NTT). Currently, she is working on research of character recognition.



Norihiro Hagita received the B.E., M.E. and Ph.D. degrees in 1976, 1978 and 1986, respectively all in electrical engineering. In 1978, he joined Nippon Telegraph and Telephone Public Corporation (now NTT). He is currently a Senior Scientist, Supervisor in NTT Basic Research Laboratories.