

Handwritten Document Security System with Inner Product and Shape based Feature Extraction Method

Dian Pratiwi
Trisakti University
Jl Kyai Tapa No.1, Jakarta-Indonesia

Syaifudin Abdullah
Trisakti University
Jl Kyai Tapa No.1, Jakarta-Indonesia

ABSTRACT

Security of handwritten document in present is one of important things, because many crimes against the falsification of document is growing. For example, in the case of signature forgery or falsification of land certificate. This research was conducted with the aim of designing an application system which can recognize handwritten of the owner of the document through the characteristics of shape, so that the falsification of documents can be prevented. The method applied in this study consisted of writing stage analog to digital conversion, pre-processing, automatic segmentation into the size of 500x200 pixels every data word by word and 10 grids, feature extraction stage, and the percentage calculation of similarity through the similarity measures and inner product method. From the research that has been conducted on 100 documents of 20 owners handwriting, 72 documents identified the owner managed appropriately through matching the features of the 12 words in each document, namely "The", "You", "Will", "To", "He", "And", "It", "Is", "Are", "His", "Have", "For". So that the percentage of accuracy, precision, and recall obtained against the document security system that is equal to 72%, 72.8%, and 67.4%.

General Terms

Data Mining and Pattern Recognition.

Keywords

Handwritten, Falsification, Document Security, Similarity Measures, Inner Product

1. INTRODUCTION

Currently, the development of various types of computer technology has developed rapidly following the various needs that exist. Computer technology is no longer foreigners and almost every individual is able to use it in everyday, such as in the use of notebook or laptop. These advances also have an impact on developments in software device such as the handwriting recognition device. Handwriting recognition tool is pretty much up to now has been applied in a handheld device or mobile phone based touchscreen. In function, the device is generally used only as a reader of results of user's handwriting on the screen to be known the meaning of the sentence. Handwriting recognition function then can be developed further to be known owner of the handwriting of the results of the analysis of shapes or patterns that exist in it. Handwriting patterns can be known through a series of image processing and feature extraction of texts that will produce features or special characteristics that can be used to identify each owner handwriting tested. With this handwriting recognition software, in the future is not just to be able to identify the owner of the handwriting pattern, but can also be used as a protection system such as in terms of securing

important documents from counterfeiting writings and the system is maintained through the input of pattern article.

2. THEORETICAL

2.1 Pre-processing

Pre-processing is an early stage that needs to be done to get the post data in digital form with the same size of the pixel and the same greylevel of a set of analog handwriting that has been digitized by means of scanner. This stage consists of RGB color conversion into grayscale and thresholding.

- RGB to greyscale color conversion

RGB to greyscale color conversion is a stage to change the color value of 24 bits to 8 bits, so the size of the resulting color will be smaller with the interval between 0 to 255 [1].

- Thresholding

Thresholding is a process to separate the object region (foreground) to the background area through a certain threshold value [2]. In this study, threshold value is also determined by trial and error

2.2 ROI Formation

ROI (Region of Interest) formation is a technique that is commonly performed to assist the analysis of the object to be observed, such as fMRI image analysis conducted by researchers from the UCLA - Los Angeles, Russel A Poldrack in 2007. This technique can improve the success of the introduction phase, due to the information of ROI, feature extraction process to be performed is limited to a specific region or area that has been restricted [3]

2.3 Feature Extraction

Feature extraction is an important stage of the pattern recognition application. This stage will give results in the form of values of the feature to be measured or recognized as a pattern. With feature or traits extraction, important information of data (which in this study is the form of image data) will be taken and stored in the feature vector [2]. Features that can be extracted in the form of image data including color features, shapes, and textures. And in this study, which will be extracted feature is based on the representation of handwritten form. The values of the handwriting feature extraction forms-based will be binary values (worth "0" and "1") for each grid for each image, where the value of "0" will be given if the representation of the grid is background object, and the value "1" if the grid is representation of foreground objects with a minimum of 15% of total pixels of each grid is a foreground object [4]

2.4 Pattern Recognition

Pattern recognition is one of the artificial intelligence techniques that aims to recognize the features or specific characteristics of data set (both text and image document) and classify [5]. Pattern recognition can be done in several ways, one of which is by using the method of *similarity measures*.

Similarity measures is a method that can be used to look for similarities from one object to another, by calculating the distance of which [6].

As in the research conducted Anna Huang in 2008, this study also used the technique of Similarity Measures to recognizing handwriting patterns by calculating the distance between the patterns by using the Euclidean Distance formula [6] :

$$d(p, q) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \dots + (q_n - p_n)^2}$$

$$= \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Where d is the distance between the total value of the handwriting image pixel with each other, q and p are the pixel image

The formula has been used by researchers in a previous paper [7], and the results are not too good. Therefore, in this study, the researchers adapted the use of weights from inner product method in order to obtain a more accurate distance value again :

$$d(p, q) = \sqrt{(q_1 \cdot w - p_1 \cdot w)^2 + (q_2 \cdot w - p_2 \cdot w)^2 + \dots + (q_n \cdot w - p_n \cdot w)^2}$$

$$= \sqrt{\sum_{i=1}^n (q_i \cdot w - p_i \cdot w)^2}$$

Where w is a random weights 0 to 1.

2.5 Calculation of Accuracy, Precision, and Recall

This step is the final stage, where each image handwriting tested will be the level of success of its introduction. The formula used is :

$$Accuracy = \frac{\text{Total of Image Handwriting Successfully Identified}}{\text{Overall Count of Handwriting Tested}} \times 100\%$$

$$Precision = \frac{\text{Total of relevant images} \cap \text{total of retrieved images}}{\text{Total of retrieved images}} \times 100\%$$

$$Recall = \frac{\text{Total of relevant images} \cap \text{Total of retrived images}}{\text{Total of relevant images}} \times 100\%$$

3. PROCEDURE & IMPLEMENTATION

3.1 Collecting Data

Collecting data in this study is done through direct searches of a number of resources (randomly close person with researcher) by a certain time interval. They are recorded in a book in white background, with intervals of three days in a row one times, and intervals of one week later one times, and intervals of one month later one times. This is done to see in future studies if there's any change in the shape of a person's handwriting in a certain period. So that the number of documents collected for each resource amount to five documents.

Data in the form of handwritten documents obtained by asking authors to write a sentence like the following :

"The man who loves you more will allow you to grow as a person without taking space. He will be patient and kind. It is because you are his priority. He will always have a reason for seeing you."

These sentences chosen by the researchers because it has a number of conjunction and common words in the English language each article, such as "The", "Will", "And", and so on. In this study, the handwritten data are tested using a sentence written with the English language, but can also use the handwritten with other languages. The entire document then will be scanned and stored as digital data.

3.2 Design and Testing

After the document that containing the digitized handwriting, the next step is pre-processing, where each article/handwritten document will be converted into a greyscale color that has a little bit more color. After that, every file or document will be taken 12 words in it through assimilation and cropping stage by using this system and stored in the folder of trial data. The word assimilation is done in accordance with the words to be taken as "The", "You", "Will", "To", "He", "And", "It", "Is", "Are", "His", "Have", "For". Cropping phase will produce an image with a size of 500x200 pixels automatically, and each document will produce 12 pictures. From 12 pictures, each will be extracted feature shape after forming grids sized 100x100 pixels and the values of the features that contain a 10-digits binary value is then stored in the form of a text (.txt) via the application

Overall, the data collected by researchers amounted to 100 documents from 20 sources (we call : A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, Q, R, S, T). All these documents resulted in a total of 1200 words in common, characteristics are then calculated through the method of similarity measures. The results of these calculations will show the distance value of each author/writer/source, where the owner of the real documents will be selected based on the value of the smallest distance generated. Because the smaller the distance writings produced, will be more similar to the tested article.

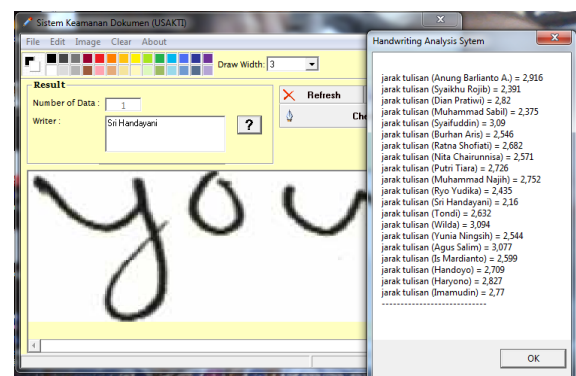


Fig. 1. Distance Calculation of Handwritten between each Writer

Based on the increase in the number of data as much as 80 documents, with the total number of handwritten being tested now as many as 100 data (20 authors) and the application of the weights of the inner product method (which in previous experiments [7] have not been applied). The results obtained can be seen in the tables below :

Table 1. The Test Result of First Sampling

Word	1	2	3	4	5	6	7	8	9	10	11	12	True	False	Author's Prediction
A	A	L	E	B	A	Q	H	S	M	H	A	P	3	9	A
B	H	I	B	O	M	P	H	B	T	E	C	T	2	10	B/H/T
C	S	C	C	T	M	Q	I	S	C	T	A	L	3	9	C
D	A	L	B	M	Q	D	D	B	I	S	D	P	3	9	D
E	A	L	G	T	E	P	G	S	H	J	E	T	2	10	E/T
F	A	F	M	R	F	Q	S	F	F	I	O	P	4	8	F
G	R	G	G	G	C	D	T	A	I	I	O	F	3	9	G
H	B	C	C	G	Q	H	I	M	I	H	R	1	11	C/G/I	
I	B	L	R	O	I	I	C	M	S	T	O	P	2	10	I/O
J	J	B	P	M	I	M	Q	S	C	T	J	J	3	9	J
K	A	B	P	G	F	I	K	I	K	T	O	K	3	9	K
L	S	H	L	F	P	L	T	K	D	J	N	E	2	10	L
M	A	M	P	S	A	M	O	M	D	H	A	P	1	11	A
N	B	L	T	N	A	S	O	M	I	T	N	P	2	10	N
O	O	C	L	N	M	M	O	O	I	T	E	Q	3	9	O
P	P	C	B	M	I	A	P	T	E	P	H	E	3	9	P
Q	Q	Q	Q	B	F	D	P	T	I	T	S	D	3	9	Q
R	R	R	G	B	D	A	I	O	S	T	Q	N	2	10	R
S	A	E	T	F	H	P	S	H	S	N	F	N	2	10	S
T	A	J	G	G	T	P	L	H	Q	I	Q	M	0	12	G

Note : 1 = The, 2 = You, 3 = Will, 4 = To, 5 = He, 6 = And, 7 = It, 8 = Is, 9 = Are, 10 = His, 11 = Have, 12 = For

In the table 1, shows that the number of correct predictions of writers as much as 14, and incorrect is 6.

Table 2. The Test Result of Second Sampling

Word	1	2	3	4	5	6	7	8	9	10	11	12	True	False	Author's Prediction
A	A	L	E	B	A	Q	H	S	M	H	A	P	3	9	A
B	H	I	B	O	M	P	H	F	T	E	C	T	1	11	H/T
C	S	C	C	T	M	Q	I	S	C	T	A	L	3	9	C
D	A	L	B	M	Q	D	D	B	I	S	D	P	3	9	D
E	A	L	G	T	E	P	G	S	H	J	E	B	2	10	E
F	A	F	M	R	N	Q	S	F	L	I	O	P	2	10	F
G	R	G	G	G	C	D	T	A	I	I	O	F	3	9	G
H	B	C	C	K	G	Q	H	I	M	I	H	R	2	10	C/H/I
I	B	L	R	O	I	I	C	M	S	T	O	P	2	10	I/O
J	J	B	P	M	I	M	Q	S	C	T	J	J	3	9	J
K	A	B	P	G	F	I	K	I	K	T	O	K	3	9	K
L	S	L	I	F	P	T	L	K	D	J	N	E	2	10	L
M	A	M	P	S	F	M	O	M	D	H	F	P	3	11	M
N	B	L	T	N	A	S	S	O	I	T	S	P	1	11	S
O	O	C	L	N	M	M	O	O	I	T	E	Q	3	9	O
P	P	C	B	M	I	A	P	T	E	P	H	E	3	9	P
Q	Q	Q	Q	B	F	D	P	T	I	T	S	D	3	9	Q
R	R	R	G	B	D	A	I	O	S	T	Q	N	2	10	R
S	A	E	T	N	H	T	S	H	S	N	F	N	2	10	N
T	A	J	G	G	T	P	L	H	Q	I	Q	M	0	12	G

Note : 1 = The, 2 = You, 3 = Will, 4 = To, 5 = He, 6 = And, 7 = It, 8 = Is, 9 = Are, 10 = His, 11 = Have, 12 = For

In the table 2, shows that the number of correct predictions of writers as much as 14, and incorrect is 6.

Table 3. The Test Result of Third Sampling

Word	1	2	3	4	5	6	7	8	9	10	11	12	True	False	Author's Prediction
A	A	L	E	B	A	Q	H	S	M	H	A	P	3	9	A
B	B	I	B	O	M	P	B	F	T	E	C	T	3	9	B
C	S	C	C	T	M	Q	I	S	C	T	A	L	3	9	C
D	A	L	B	M	Q	D	D	B	I	S	D	P	3	9	D
E	A	L	G	T	E	P	G	S	H	J	E	B	2	10	E
F	A	F	M	R	N	Q	S	F	L	I	O	P	2	10	F
G	R	G	G	G	C	D	T	A	I	I	O	F	3	9	G
H	A	C	H	K	G	Q	H	L	M	I	H	R	3	9	H
I	B	L	R	O	I	I	C	N	L	T	O	I	3	9	I
J	J	B	P	M	I	K	Q	S	A	T	J	J	3	9	J
K	A	N	P	G	F	P	N	I	K	T	O	K	2	10	K/P/N
L	S	L	L	F	Q	T	O	K	D	J	B	E	2	10	L
M	M	M	P	K	F	M	O	M	D	H	O	A	4	8	M
N	B	H	T	N	A	I	S	S	I	T	S	P	1	11	S
O	O	C	L	N	M	M	O	O	I	T	E	Q	3	9	O
P	P	C	B	M	I	A	P	T	E	P	H	E	3	9	P
Q	Q	Q	Q	B	F	D	P	T	I	T	S	D	3	9	Q
R	R	R	G	B	D	A	I	O	S	T	Q	N	2	10	R
S	A	E	T	N	H	T	S	H	S	N	F	N	2	10	N
T	T	J	J	T	T	P	L	H	F	I	Q	M	3	9	T

Note : 1 = The, 2 = You, 3 = Will, 4 = To, 5 = He, 6 = And, 7 = It, 8 = Is, 9 = Are, 10 = His, 11 = Have, 12 = For

In the table 3, shows that the number of correct predictions of writers as much as 17, and incorrect is 3.

Table 4. The Test Result of Fourth Sampling

Word	1	2	3	4	5	6	7	8	9	10	11	12	True	False	Author's Prediction
A	I	L	E	A	Q	H	A	M	H	O	P	3	9	A	
B	B	I	N	O	M	P	B	F	B	E	L	I	3	9	B
C	S	C	S	T	M	Q	I	S	C	C	A	L	3	9	C/S
D	D	L	B	M	Q	N	G	B	I	S	D	P	2	10	B/D
E	A	D	G	T	H	P	G	M	H	J	E	J	1	11	J
F	A	F	M	R	L	Q	S	F	L	I	O	F	3	9	F
G	R	G	G	G	C	D	T	A	I	I	O	F	3	9	G
H	A	C	H	K	G	Q	H	L	M	I	H	R	3	9	H
I	B	L	R	O	I	I	C	N	L	T	O	I	3	9	I
J	J	B	S	M	I	K	Q	S	A	T	J	J	3	9	J
K	K	N	P	G	F	S	A	I	K	T	O	K	3	9	K
L	S	L	L	F	Q	T	O	K	D	J	B	E	2	10	L
M	M	A	P	K	D	M	O	M	S	H	J	A	3	9	M
N	B	H	T	N	A	I	S	S	I	T	S	P	1	11	S
O	O	C	L	N	M	M	O	O	I	T	E	Q	3	9	O
P	P	C	B	M	I	A	P	T	E	P	H	E	3	9	P
Q	Q	L	Q	O	F	B	P	T	I	L	S	L	2	10	L/Q
R	R	R	G	B	D	A	I	O	S	T	Q	N	2	10	R
S	A	E	T	N	H	T	S	H	S	N	F	N	2	10	N
T	T	J	J	T	T	P	L	H	F	I	Q	M	3	9	T

Note : 1 = The, 2 = You, 3 = Will, 4 = To, 5 = He, 6 = And, 7 = It, 8 = Is, 9 = Are, 10 = His, 11 = Have, 12 = For

In the table 4, shows that the number of correct predictions of writers as much as 14, and incorrect is 6.

Table 5. The Test Result of Fifth Sampling

Word	1	2	3	4	5	6	7	8	9	10	11	12	True	False	Author's Prediction
A	I	L	E	S	A	L	A	H	M	F	O	P	2	10	A/L
B	B	B	F	O	J	P	P	B	B	E	L	I	4	8	B
C	S	C	S	T	M	Q	S	S	N	C	A	L	2	10	S
D	A	R	Q	M	Q	N	G	I	I	S	D	D	2	10	D
E	A	D	G	T	H	P	G	M	H	J	E	J	1	11	J
F	A	F	M	R	L	Q	S	F	L	I	O	F	3	9	F
G	R	G	G	G	C	D	T	A	I	I	O	F	3	9	G
H	A	C	H	K	G	Q	H	L	M	I	H	R	3	9	H
I	B	L	R	O	I	I	C	N	L	T	O	I	3	9	I
J	J	B	S	M	I	K	Q	S	A	T	J	J	3	9	J
K	K	N	P	G	F	S	A	I	K	T	O	K	3	9	K
L	S	L	L	F	Q	T	O	L	D	J	B	E	3	9	L
M	M	A	P	K	D	M	O	M	S	H	J	A	3	9	M
N	B	I	T	J	A	I	S	S	I	T	S	P	1	11	I/S
O	A	C	L	N	K	M	O	F	I	T	E	Q	1	11	K
P	P	C	B	M	I	A	P	T	E	P	H	E	3	9	P
Q	Q	L	Q	O	F	B	P	T	I	M	S	L	2	10	Q
R	R	R	G	B	D	A	I	O	S	T	Q	N	2	10	R
S	A	E	T	N	H	T	S	H	S	N	F	N	2	10	N
T	A	J	J	A	T	P	D	P	P	S	Q	M	1	11	P

Note : 1 = The, 2 = You, 3 = Will, 4 = To, 5 = He, 6 = And, 7 = It, 8 = Is, 9 = Are, 10 = His, 11 = Have, 12 = For

In the table 5, shows that the number of correct predictions of writers as much as 13, and incorrect is 7.

Thus, of the five experiments may be summarized the results as follows :

Table 6. The Result of Precision, Recall, and Accuracy

Author's Prediction													Author Handwritten	Precision %	Recall %	Accuracy %							
A	B	C	D	E	F	G	H	I	J	K	L	M					N	O	P	Q	R	S	T
5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	A	83	83	
-	4	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	B	80	50	
-	-	2	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	C	66.7	66.7	
-	-	-	1	-	5	-	-	-	-	-	-	-	-	-	-	-	-	-	-	D	100	83	
-	-	-	-	3	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	E	100	50	
-	-	-	-	-	-	5	-	-	-	-	-	-	-	-	-	-	-	-	-	F	100	100	
-	-	-	-	-	-	-	5	-	-	-	-	-	-	-	-	-	-	-	-	G	62.5	100	
-	-	-	-	-	-	-	-	1	4	-	-	-	-	-	-	-	-	-	-	H	66.7	44.4	
-	-	-	-	-	-	-	-	-	-	5	-	-	-	-	-	-	-	-	-	I	62.5	71.4	
-																							

accuracy of 72%. Then to the value of precision obtained with an average of 72.8% and an average of recall value is 67.4%. Thus, this study quite successfully develop a system that can prevent the forgery of handwritten documents.

4. CONCLUSION

Based on the results of the data and document security applications, the overall researchers can provide the following conclusions :

1. The percentage of success rate on the document security system in this study reached 72%, where 72 of 100 documents successfully distinguished handwriting correctly through the features form. In addition, the percentage of precision and recall are also pretty good that is equal to 72.8% and 67.4%.
2. The number of features that are used relatively little that amounted to only 10 feature values, still causing some word used in writing the feature vector can have the same value even if the words a different.
3. Through weighting with the inner product method, can give results within a better value compared to just using the euclidean distance. Because by assigning weights, a small distance will be smaller, and the great distances will be greater. So that each document will be more clearly its differences and the better its match result.

5. REFERENCES

- [1] Pratiwi, D. 2012. The Use of Self Organizing Map Method and Feature Selection in Image Database Classification System. *International Journal of Computer Science Issues (IJCSI)*, Vol.9, Issue 3 No.2 ISSN : 1694-0814
- [2] Pratiwi, D. Santika, D.D, and Pardamean, B. 2011. An Application of Backpropagation Artificial Neural Network Method for Measuring The Severity of Osteoarthritis. *International Journal of Engineering & Technology (IJET-IJENS)*. Vol.11, No.3, ISSN: 117303-8585
- [3] Poldrak, R.A. 2007. Region of Interest Analysis for fMRI. *Oxford Journal*. Vol.2 Issue 1, pp: 67-70. Los Angeles-USA.
- [4] Lu, G. 1999. Multimedia Database Management Systems. Artech House Inc.
- [5] Absultanny, Y.A. 2003. Pattern Recognition using Multilayer Neural Genetic Algorithm. *Neurocomputing*. Pp.237-247. Elseiver Science.
- [6] Huang, A. 2008. Similarity Measures for Text Document Clustering. *New Zealand Computer Science Research Student Conference*. Christchurch. New Zealand
- [7] Pratiwi, D and Syaifudin. 2015. The Implementation of Shape Based Feature Extraction and Similarity Measures to Prevent Falsification of Handwritten Document. *Information Systems International Conference (ISICO)* ITS-Surabaya, Indonesia