

Document Image Retrieval using Signatures as Queries

Sargur N. Srihari, Shravya Shetty,
Siyuan Chen, Harish Srinivasan, Chen Huang
CEDAR, University at Buffalo(SUNY)
Amherst, New York 14228

Gady Agam and Ophir Frieder
Illinois Institute of Technology,
Chicago, Illinois 60616

Abstract

In searching a repository of business documents, a task of interest is that of using a query signature image to retrieve from a database, other signatures matching the query. The signature retrieval task involves a two-step process of extracting all the signatures from the documents and then performing a match on these signatures. This paper presents a novel signature retrieval strategy, which includes a technique for noise and printed text removal from signature images, previously extracted from business documents. Signature matching is based on a normalized correlation similarity measure using global shape-based binary feature vectors. In a retrieval task involving a database of 447 signatures, on an average 4.43 out of the top 5 choices were signatures belonging to the writer of the queried signature. On considering the Top 10 ranks, a *F-measure* value of 76.3 was obtained and the precision and recall values at this *F-measure* were 74.5% and 78.28% respectively.

1. Introduction

In searching complex documents, such as a repository of archival office documents, a task of relevance is relating the signature in a given document to the closest matches within a database of documents; this is known as the signature retrieval task. This paper presents an effective signature retrieval technique and presents its accuracy using a system called CEDAR-FOX [8, 9] which attempts to provide the full range of functionalities for a digital library of handwritten/signed documents. Signature retrieval is the process of retrieving the closest matching signatures to the questioned signature from a database of known signatures, the objective is to identify other closely resembling signatures by the same writer. Given a database of signed documents, it would be of interest to relate a queried document to other documents in this database which have been signed by the same author. Figure 1 shows examples of 3 documents signed by the author "10". The retrieval task would

be to retrieve all the other documents signed by the same author. This would involve extracting all the signatures from the documents and then performing a match on these signatures. The task of extracting signatures from these documents has not been dealt with in this paper, the signatures were manually cropped from the documents and the retrieval was performed on these images. Figure 2 shows a few sample signature images after extraction.

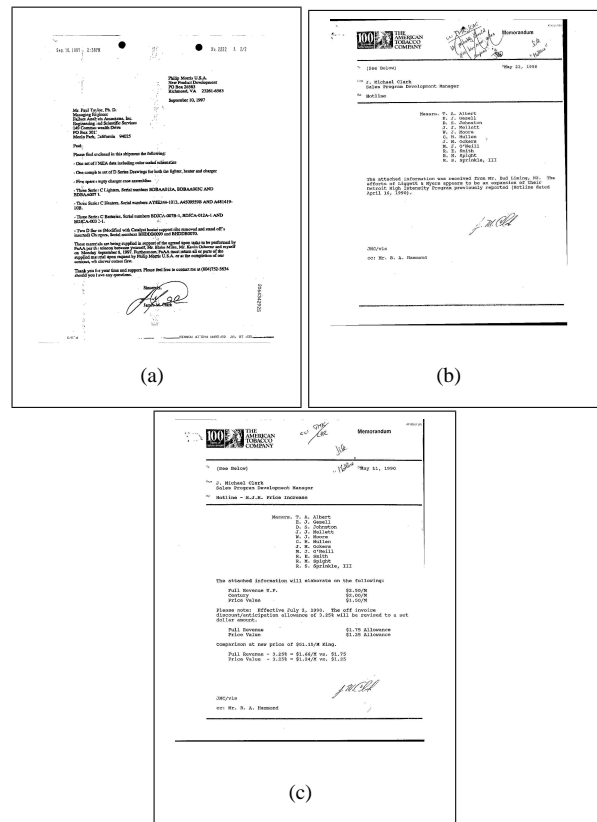


Figure 1. Three examples of business documents signed by the same author in an archival database: (a) document 1, (b) document 2 and (c) document 3.

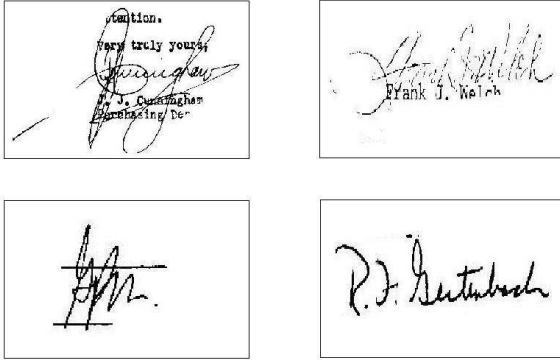


Figure 2. Sample signature images cropped from business documents.

Once the signatures have been extracted, the retrieval involves a 2 step process of 1) Selecting the questioned signature and 2) Selecting the list of known signatures to be compared against to obtain a ranking for each of the selected known signatures. Previous works [5, 4, 1] on signature verification and identification discuss methods for both writer dependent and writer independent cases. The organization of the rest of the paper is as follows. In Section 2, the various steps involved in the Signature Retrieval process are described. In Section 3, the dataset used for the signature retrieval experiments is described with a few examples. In Section 4, the various experiments conducted and the results obtained are discussed and Section 5, concludes the paper.

2. Signature Retrieval Strategy

The steps involved in signature retrieval are described here. For image analysis general nomenclature and approaches, see [7].

2.1. Printed Text Removal

In the preprocessing step the printed text is removed from the signature samples. To remove the printed text from the signature image, an image enhancement procedure based on chaincode information has been implemented. This procedure can be represented in three phases:

- 1) After generating the chaincode for the noisy image, x_1 (the ratio of the length and the width of the boundary box), x_2 (the ratio of the size of the contour and the largest one in the image) are recorded as features to separate the handwriting components and printed text components. Fisher linear discriminants are chosen to find the optimum direction W in a two-dimensional space to separate the two categories of

components using the ratio of average inter-class distance and intra-class distance. $y = W^T X$ which is the combined feature is further assumed to satisfy the Gaussian distribution and all the disconnected components of the image are separated into two classes (suspected or unsuspected) using a Bayesian approach.

- 2) The mutual spatial relation between the suspected components in the first phase is used to determine the category of the suspected components i.e. only those ranging in lines are classified and removed as printed text.



Figure 3. Example of printed text removal.

- 3) In the last phase, the neighborhood is scanned in a fixed-size window for the rest of the suspected components to remove the far and isolated stroke components. An example of the results obtained by this printed text removal procedure is shown in Figure 3.

2.2. Feature Extraction

The next step in signature retrieval involves converting each sample image into a set of binary feature vectors. The features used here are the Gradient, Structural and Concavity (GSC) features[2, 10] which measure the image characteristics at local, intermediate and large scales. The features for the signature images which are extracted under a 4 X 8 division, contain 384 bits of gradient features, 384 bits of structural features and 256 bits of concavity features, giving us a binary feature vector of length 1024. Each of these sets of binary features uniquely represents a given sample signature.

The gradient features detect local features of the image and provide a great deal of information about stroke shape on a small scale. The structural features extend the gradient features to longer distances and give useful information about stroke trajectories. The concavity features are used to detect stroke relationships at long distances which can span across the image. All the binary features are converted from their original floating number computations.

2.3. Retrieval

In the signature retrieval process, the first step is preprocessing where the printed text is removed from the questioned signature and from the set of known signatures as described above. Following this, the GSC binary feature vectors are extracted for the queried and the each of the known signatures and stored as a binary feature vectors. This is

	Writer 1				Writer 2				Writer 3			
	1_a	1_b	1_c	1_d	2_a	2_b	2_c	2_d	3_a	3_b	3_c	3_d
1_a	0	0.22	0.26	0.25	0.36	0.37	0.4	0.32	0.33	0.35	0.32	0.35
1_b	0.22	0	0.25	0.27	0.34	0.36	0.38	0.31	0.3	0.34	0.3	0.36
1_c	0.26	0.25	0	0.25	0.31	0.38	0.36	0.34	0.34	0.35	0.35	0.39
1_d	0.25	0.27	0.25	0	0.35	0.36	0.36	0.35	0.35	0.37	0.36	0.4
2_a	0.36	0.34	0.31	0.35	0	0.27	0.24	0.28	0.33	0.33	0.35	0.4
2_b	0.37	0.36	0.38	0.36	0.27	0	0.25	0.27	0.36	0.38	0.34	0.38
2_c	0.4	0.38	0.36	0.36	0.24	0.25	0	0.22	0.35	0.37	0.38	0.43
2_d	0.32	0.31	0.34	0.35	0.28	0.27	0.22	0	0.31	0.35	0.33	0.39
3_a	0.33	0.3	0.34	0.35	0.33	0.36	0.35	0.31	0	0.26	0.24	0.29
3_b	0.35	0.34	0.35	0.37	0.33	0.38	0.37	0.35	0.26	0	0.27	0.28
3_c	0.32	0.3	0.35	0.36	0.35	0.34	0.38	0.33	0.24	0.27	0	0.17
3_d	0.35	0.36	0.39	0.4	0.4	0.38	0.43	0.39	0.29	0.28	0.17	0

Table 1. Similarity distance values for 12 Signature images - Distance values between signatures belonging to the same author have been highlighted.



Figure 4. Signature images used to compare correlation distance values: (a) writer 1, (b) writer 2 and (c) writer 3.

done by enclosing each image in a 4 x 8 grid and then computing gradient, structural and concavity features for each cell of the grid. The result is a binary vector of length 1024. [12] describes the use of these features in spotting words in document images.

The distance between the queried signature and each of the known signatures in the database is calculated using a normalized correlation similarity measure[11, 12]. Given the two binary feature vectors $X \in \Omega$ and $Y \in \Omega$, each similarity distance measure $S(X, Y)$ uses all or some of the four possible values, i.e. S_{00} ; S_{01} ; S_{10} ; S_{11} . Here S_{ij} , $(i,j) \in \{0,1\}$, is the number of occurrences where pattern i

occurs in the first binary vector and pattern j occurs in the second vector in the same position. The similarity measure $S(X, Y)$ is defined as follows.

$$S(X, Y) = \frac{1}{2} + \frac{S_{11}S_{00} - S_{10}S_{01}}{2((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}} \quad (1)$$

where

S_{00} = the first binary vector has a 0 and the second vector too has a 0 in the corresponding positions.

S_{11} = the first binary vector has a 1 and the second vector too has a 1 in the corresponding positions.

S_{01} = the first binary vector has a 0 while the second vector has a 1 in the corresponding positions.

S_{10} = the first binary vector has a 1 while the second vector has a 0 in the corresponding positions.

When constructing the similarity measure all possible matches $S_{ij} \in 0,1$ are considered for better classification. Also S_{00} has been weighted with a beta value of 0.5 in order to boost classification. The results are ranked in the increasing order of this similarity distance. In the signature retrieval process there is no prior knowledge of the writers signature, the goal is to identify the closest signatures and to identify other similar signatures by the writer of the queried signature.

Table 1 shows examples of the distance values obtained using this correlation similarity measure for the images shown in Figure 4. The table shows 12 signature images belonging to three authors, each of these signatures was matched against all the images in this set. It can be seen from the results that the scores obtained when matching a signature against signatures by the same author are considerably smaller than the values obtained when matched

against those belonging to other authors. A similarity score of 0 is always obtained when an image is matched against itself.

3. Dataset

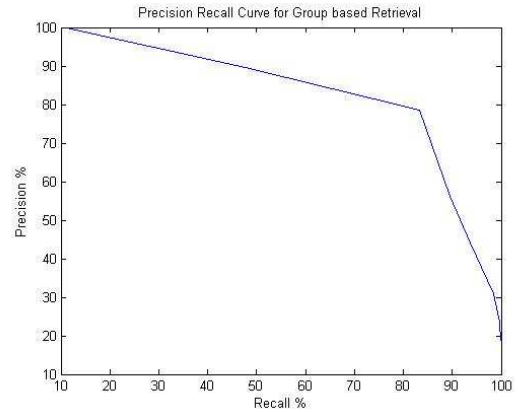
The dataset used for this experiment comprises of signature images which were cropped from signed documents collected from 40 different writers. The writers have been grouped into 6 groups consisting of around 6 authors per group. The total number of signature samples present is 447. It consists of a variety of sample images including signatures surrounded by printed text, signatures overlapping printed text, noisy images and just plain handwritten signature images. Figure 2 shows a few of these sample images. Table 2 displays all the writers and their corresponding number of samples, number of sample types and the approximate number of samples containing printed text. Some of the writers have several types of signatures like the writer's full name, initials, only first name, etc. But in most of these cases the majority of the images belonging to a particular writer were just of 1 or sometimes 2 sample types and just a few images belong to the remaining types. Figure 5 displays all the samples for the writer '20'. As seen in the figure, the samples b,d,e,f,h and i contain printed text and the sample b is of a different signature type as compared to the rest of the set.



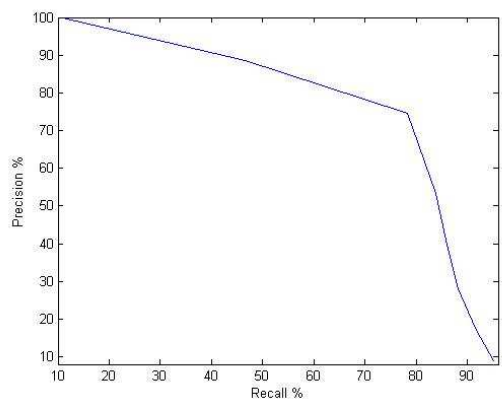
Figure 5. All the samples for writer Temko.

4. Experiments and Results

In this section, the test setup and the experimental results obtained for the signature retrieval task are described. In the test setup for Signature Retrieval, the images were divided into 2 groups per writer. One group consisting of known signature images and the second group consisting of the questioned signatures for testing. The image formats supported are png, jpeg and tiff. The signature image



(a)



(b)

Figure 6. Precision-Recall Curves: (a) Group based retrieval - Precision of 78.61% at recall of 83.15% , (b) Entire Database retrieval - Precision of 74.5% at recall of 78.28%.

in question is first selected and this queried image is displayed both before and after the preprocessing step. Following this, the set of known signatures are selected and the signature retrieval process is carried out against this set of known signatures. In each case the precision and recall values [6, 3] are calculated. The precision and recall measures for a rank 'R' where the author of the questioned signature is represented by 'A' are defined as follows

Recall = Number of signatures belonging to A retrieved with rank $\leq R$ / Total number of signatures belonging to A

Precision = Number of signatures belonging to A retrieved with rank $\leq R$ / R

The testing for Signature Retrieval was conducted in 2 stages. First the retrieval was tested for signature images within each group and then the retrieval was tested against

Writer ID	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Total Samples	15	19	9	10	12	12	10	9	15	7	14	13	10	12	14	10	10	3	11	5
Types of Signatures	5	6	2	1	1	1	2	2	2	2	3	1	2	2	3	1	1	1	3	2
Printed Text	9	15	0	0	0	0	9	6	11	1	4	0	0	5	1	2	0	0	0	0
Writer ID	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40
Total Samples	13	10	6	11	13	10	13	16	10	9	12	11	13	10	11	11	10	12	15	11
Types of Signatures	3	1	2	2	4	1	2	2	2	1	3	1	2	1	2	1	1	3	3	1
Printed Text	5	0	0	2	1	0	0	0	3	2	4	0	12	5	11	11	0	8	4	6

Table 2. Signature samples for each writer.

Group	Rank 1	Within Rank 5	Within Rank 10	Within Rank 15	Within Rank 30	Within Rank 50
Group_01	10.95%	47.6%	76.19%	84.76%	99.04%	100%
Group_02	11.6%	45.8%	78.06%	85.16%	97.4%	100%
Group_06	10.75%	48.9%	83.8%	91.93%	100%	100%
Group_07	10.74%	50%	88.3%	91.12%	95.79%	98.13%
Group_08	9.3%	44.9%	87.1%	92.18%	98.23%	99.6%
Total	10.6%	47.4%	83.15%	89.32%	98.23%	99.6%

Table 3. Recall measures for group based signature retrieval.

the entire set of 447 images.

In the group based testing, each of the test images was tested against all the other signature images in the group to which the queried signature image belongs. The ranks of the signatures belonging to the correct authors of the questioned signatures were noted in each case. The results for this is displayed in Table 3. The results show that on an average 4.48 out of the Top 5 results are signatures by the writer of the queried signature. In the top 15 results a precision of 56.29% is obtained at a recall of 89.32%. This implies that out of all the images in a group, on an average 89.32% of all known images of the correct author of a queried image were returned in the Top 15 choices. The Precision-Recall graph is displayed in Figure 6(a).

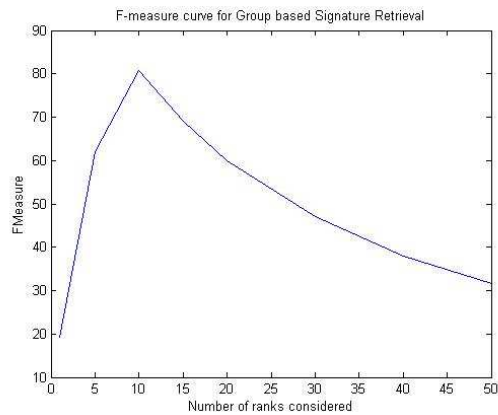
Similarly, in the second part of the signature retrieval experiment, this process was repeated for the entire database of signatures. The testing was then done for 60 randomly selected questioned signatures and the ranks of the signatures belonging to the correct authors of the questioned signatures were noted in each case. The results for this is displayed in Table 4. The results show that on an average 4.43 out of the Top 5 results are signatures by the writer of the queried signature. In the top 15 results a precision rate of

Number of Ranks Considered	Recall Measure
Rank 1	10.5%
Rank 5	46.5%
Rank 10	78.28%
Rank 15	83.88%
Rank 30	88.26%
Rank 50	91.59%

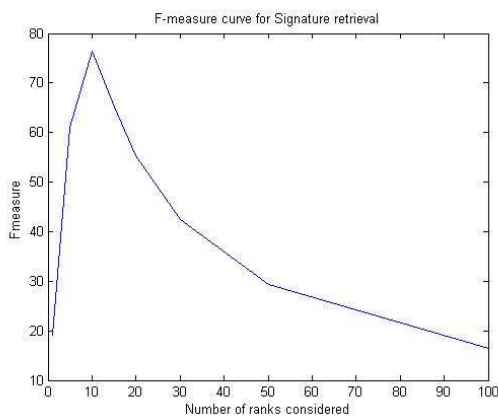
Table 4. Recall measures for signature retrieval from entire database.

53.22% is obtained at a recall of 83.88%. This implies that out of a total number of 447 possible ranks, on an average 83.88% of all known images of the correct author of a queried image were returned in the Top 15 choices. The Precision-Recall graph is displayed in Figure 6(b).

We also use a F-measure which combines the precision and recall values to provide a single measure of the retrieval accuracy, it is the harmonic mean of the precision and recall values. Figure 7 displays the F-measure curves for both group based and non-group based signature retrieval. In the



(a)



(b)

Figure 7. F-measure Curves: (a) Group based retrieval - Peak F-measure of 80.81 at top 10 ranks, (b) Entire database retrieval - Peak F-measure of 76.3 at top 10 ranks.

group based retrieval the peak for the F-measure was obtained for results in the Top 10 ranks, at this point the F-measure value obtained is 80.81, the recall value is 83.15% and the precision is 78.61%. In the non-group based retrieval the peak for the F-measure was again obtained for results in the Top 10 ranks. At this point the F-measure value is 76.3, the recall value is 78.28% and the precision is 74.5%. This implies that viewing the Top 10 results returned for a query supports high precision and recall measure.

5. Conclusions

Results from experiments performed for the problem of signature retrieval and its results were presented. The tests were conducted on a variety of signature samples including those with noise and printed text. The presence of noise and

the presence of multiple types of signature per author have led to a lower retrieval rate. In spite of this, our technique returned promising results for both group based and non-group based retrieval, we achieved a precision accuracy of 89.6% and 88.6% respectively when considering the Top 10 results. It was also seen that the retrieval results for group based testing yielded a slightly better precision and recall measure indicating that performing an unsupervised clustering before the retrieval process results in a better accuracy rate.

References

- [1] B. Fang, C. H. Leung, Y. Y. Tang, K. W. Tse, P. C. K. Kwok, and Y. K. Wong. Off-line signature verification by the tracking of feature and stroke positions. *Pattern recognition*, 36:91–101, 2003.
- [2] J. T. Favata and G. Srikantan. A multiple feature resolution approach for handprinted digit and character recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [3] D. Grossman and O. Frieder. *Information Retrieval: Algorithms and Heuristics*. Springer Publishers, Second Edition 2004.
- [4] M. K. Kalera, B. Zhang, and S. N. Srihari. Off-line signature verification and identification using distance statistics. *Proc. of the Int. Graphonomics society Conference*, pages 228–232, 2003.
- [5] R. Plamondon and G. Lorette. On-line and offline handwriting recognition: A comprehensive survey. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 22(1):63–84, 2000.
- [6] G. Salton and M. J. McGill. *Introduction to modern information retrieval*. McGraw-Hill, New York, 1983.
- [7] M. Sonka, V. Hlavac, and R. Boyle. *Image Processing : Analysis and Machine Vision*. Thomson-Engineering, 1998.
- [8] S. Srihari, B. Zhang, C. Tomai, S. Lee, Z. Shi, and Y. Shin. A system for handwriting matching and recognition. *Proceedings of the Symposium on Document Image Understanding Technology (SDIUT 03)*, pages 67–73, 2003.
- [9] S. N. Srihari, S. Cha, H. Arora, and S. Lee. Individuality of handwriting. *Journal of Forensic Sciences*, 47(4):856–872, 2002.
- [10] G. Srikantan, S. Lam, and S. Srihari. Gradient based contour encoding for character recognition. *Pattern Recognition*, 29(7):1147–1160, 1996.
- [11] B. Zhang and S. Srihari. Binary vector dissimilarity measures for handwriting identification. *Document Recognition and Retrieval X, Proceedings of Electronic Imaging, San Jose, CA*, pages 155–166, 2003.
- [12] B. Zhang, S. N. Srihari, and C. Huang. Word image retrieval using binary features. *Document Recognition and Retrieval XI, SPIE, San Jose, CA*, 2004.