



Text independent writer recognition using redundant writing patterns with contour-based orientation and curvature features

Imran Siddiqi^{a,b}, Nicole Vincent^{a,*}

^a LIAPDE – SIP, University Paris Descartes, 45 Rue des Saint Pères, 75006 Paris, France

^b College of Telecommunication Engineering, MCS-NUST, Rawalpindi, Pakistan

ARTICLE INFO

Article history:

Received 25 November 2009

Received in revised form

16 April 2010

Accepted 13 May 2010

Keywords:

Handwritten documents

Writer identification

Writer verification

Clustering

Freeman chain code

Polygonization

ABSTRACT

We propose an effective method for automatic writer recognition from unconstrained handwritten text images. Our method relies on two different aspects of writing: the presence of redundant patterns in the writing and its visual attributes. Analyzing small writing fragments, we seek to extract the patterns that an individual employs frequently as he writes. We also exploit two important visual attributes of writing, orientation and curvature, by computing a set of features from writing samples at different levels of observation. Finally we combine the two facets of handwriting to characterize the writer of a handwritten sample. The proposed methodology evaluated on two different data sets exhibits promising results on writer identification and verification.

© 2010 Elsevier Ltd. All rights reserved.

1. Introduction

The need to recognize the writer of a handwritten document is a recurrent problem with potential applications for document examiners, paleographers and forensic analysts. These experts might need to identify the authorship or authenticity of a questioned document (e.g. an ancient manuscript, a will, a ransom note, a threatening letter), verify signatures, identify forgeries, detect alterations or identify indented writings. These tasks require the document examiners to recognize the writer of a questioned document which naturally is a tiresome procedure. Thus developing computerized systems for automatic writer recognition could serve as valuable tools in ancient and forensic document analysis. These systems take as input the scanned document images and employ image processing and pattern classification techniques to solve a given problem. Of course the results of these systems cannot be accepted as evidence in a court case and the intervention of human experts is inevitable; nevertheless, they could still be quite helpful in speeding up the examination process considerably.

Writer recognition is generally distinguished into writer identification and verification. Writer identification involves finding the author of a query document given a reference base with documents of known writers. Writer verification on the

other hand determines whether two samples are written by the same person or not. Writer identification is generally carried out by calculating a similarity index between the questioned writing and all the writings of known writers and sorting the retrieved results in a hit list with an increasing distance from the query. Choosing appropriate acceptance thresholds on these similarity indices can then be used to perform writer verification.

The problem of writer recognition is related to that of handwriting recognition. Handwriting recognition aims at eliminating the writer-dependent variations between writings and thus identifying the individual characters and words. Writer recognition on the other hand relies on these writer-specific variations between character shapes which allow to characterize its writer's hand [1]. Despite this contradiction between the two approaches, writer recognition can be handy in handwriting recognition by exploiting the principle of adaptation of the system to the type of writer [2].

This paper proposes a framework for automatic writer recognition based on two facets of handwriting. As an individual writes, he draws similar characters using the same basic shapes which eventually leads to a certain redundancy in the writing of an individual with some patterns more frequent than others. Based on the learned copy book style, personal preferences, and certain neurological and psychological factors, these frequent patterns will vary from one individual to another and hence could be exploited to characterize the writer of a handwritten sample. Unlike classical approaches which consider the patterns having an interpretation such as characters or graphemes, we think that writer of a document is characterized by writing gestures rather

* Corresponding author.

E-mail addresses: imran.siddiqi@mi.parisdescartes.fr (I. Siddiqi), nicole.vincent@mi.parisdescartes.fr, nicole.vincent@math-info.univ-paris5.fr (N. Vincent).

than the alphabet used. We therefore introduce much elementary writing fragments. This allows to extract the frequent writing patterns that might be common across different characters or graphemes. We also focus on two important visual attributes of writing namely orientation and curvature that enable distinguishing one writing from another. These properties are captured by a set of features extracted from the contours of handwriting images. Finally we explore the effect of combining the two aspects in characterizing the writer of a given sample.

This paper is organized as follows: in the first section we give a brief account of some significant recent contributions to writer identification. In the next part we introduce the data sets used in our study followed by the description of our proposed methodology. The following section presents the experimental results and their analysis and finally we give a conclusion with some future research directions.

2. Background

Despite the development of electronic documents and predictions of a paperless world, the importance of handwritten documents has retained its place and the problems of identification and authentication of writers have been an active area of research over the past few years. A wide variety of systems that are based on the use of computer image processing and pattern recognition techniques have been proposed to solve the problems encountered in automatic analysis of handwriting and recognition of the writer of a questioned document. A comprehensive survey of the work in writer recognition until 1989 has been presented in [3]. Here, we will be more interested in discussing the approaches developed in the last several years, thanks to the renewed interest of the document analysis community for this domain.

Over the recent years, a wide variety of features, local or global and structural or statistical, have been proposed that serve to distinguish the writing of an individual from another. Techniques based on these features are traditionally categorized into two broad classes: text-dependent and text-independent methods.

In text-dependent methods the writing samples to be compared require to contain the same fixed text (e.g. signature verification). These methods are generally based on comparison of individual characters or words and therefore require recognition or segmentation of text prior to writer recognition. The text-independent methods on the other hand identify the writer of a document independent of its semantic content. These methods use features extracted from the entire image of a text or from a region of interest.

The text-dependent studies on writer identification are mainly motivated by forensic applications and aim to design computational algorithms to extract the features used by forensic document examiners [4]. One of the most comprehensive studies in this area has been presented in [5] where the authors validate the hypothesis of the individuality of handwriting by extracting a set of macro and micro features from handwriting samples (comprising a letter of 156 words) of 1000 individuals. The same research was extended to analyze the discriminative powers of individual characters [6] and numerals [7] offering a general guidance for selecting the most informative characters (numerals) in examining forensic documents. In a separate study, the authors also compared the individuality of characters with that of words [8] using the same set of features, with features extracted from words naturally outperforming the ones extracted from characters.

In another recent work [9], a set of structural micro features (meant to represent a subset of the features used by forensic experts) was extracted from three characters ‘d’, ‘y’ and ‘f’ and

grapheme ‘th’. The usefulness of these features was then analyzed using a wrapper method.

In addition to the forensic features discussed above, [10] introduced a new feature based on a morphological transform of the vertical histogram function of a binarized one-pixel thinned word. The discrimination capabilities of the proposed feature were evaluated on a database of 50 writers each having 45 samples of two different words of the same length.

Text-dependent automatic writer identification might achieve very good performance but is very constrained and is mostly not applicable in many practical applications like the identification of the writers of archived handwritten documents and crime suspect identification from large sized forensic data sets etc. [11]. Therefore, most of the methods developed lately fall in the text-independent analysis of handwriting where a set of features, global or local, is extracted from the handwritten image.

Among the well-known text-independent approaches, Said et al. [11] present a global approach and consider each handwriting as a different texture employing Gabor filters and co-occurrence matrices. Marti et al. [12] analyze the difference between handwritings by extracting a set of 12 structural features from each line of the available text samples. The approaches based on the fractal analysis of handwriting include identification by the fractal dimension [13] and by fractal compression/decompression [14]. The edge-based directional probability distributions are used as features in [15] and the identification performance is compared to a number of non-angular features.

Schlapbach and Bunke [16] employ HMM-based handwriting recognition systems for the purpose of writer identification with a separate system for each writer, recognizer output score being used to characterize the writer. This work was followed by the use of Gaussian mixture models to model an individual's handwriting [17]. A comparative study of the two revealed that while the GMM based system is conceptually much simpler and faster to train than the HMM based system, it achieves significantly higher identification rates [18].

Lately, the codebook based writer identification techniques are gaining more and more attention. Among these methods, Bensefia et al. [19] exploit the morphological redundancy of graphemes that are produced by a segmentation algorithm based on the analysis of the minima of the upper contour [2]. These graphemes are clustered to produce a codebook termed as ‘writer invariants’ by the authors. Writer identification is carried out based on the direct correspondence of the invariant patterns extracted from the two texts to be compared. Extending the same idea in a later study [20], instead of determining the invariant clusters from the graphemes of each writer, the authors cluster all the graphemes of the database thus defining a common feature space over the entire data set. Writer identification is then performed by a textual based information retrieval model.

Table 1

Performance comparison of recent studies on writer identification.

Author	Year	Writers	Data set	Performance (%)
Said et al.	2000	20	–	96
Zois and Anastassopoulos	2000	50	One word	96.5
Marti et al.	2001	20	IAM	90.7
Srihari et al.	2002	100	CEDAR letter	94
		900	CEDAR letter	87
Bensefia et al.	2005	150	IAM	86
		88	PSI	96
Schlapbach and Bunke	2006	100	IAM	98.46
Bulacu and Schomaker	2007	250	Firemaker	83
		650	IAM	89
		900	IAM+Firemaker	87

The identification rates might not be truly comparable in all cases.

Using a similar approach, Schomaker and Bulacu computed a codebook from an independent training set and employed the probability-density function of the patterns in unknown writing to identify its author. The approach was first applied to isolated uppercase handwriting [21] and later extended to lowercase cursive writing [22]. An extensive analysis of combining these codebook based features with the texture based features [15] has been presented in [1].

A quantitative performance comparison (based on the results reported in the literature) of some of the methods discussed earlier has been summarized in Table 1. Although identification rates of as high as 98% have been reported, they are based on a smaller number of writers. Only the studies [1,5] have been evaluated on significantly large data sets (900 writers). The data set in the first one, however, comprises the same text written three times by each of the authors hence the effectiveness of the system on text-independent analysis is yet to be explored. Therefore, we can conclude that Bulacu and Schomaker [1] currently hold the best performance results reading around 87% on 900 writers.

From the discussion of text-dependent and text-independent methods, one can conclude that in general, higher identification rates are achievable with the former type but at the cost of the requirement to have same fixed text or human intervention to extract the elements (characters or words) to be compared. Text-independent methods are much more useful and applicable. These methods, however, require a certain minimum amount of text to produce acceptable results.

Resuming, we could say that the research on writer recognition that started with the analysis of very constrained writings and very few writers has matured really well over time. Regarding the methods developed, in addition to the structural and statistical features, codebook generation has emerged as a very popular as well as effective method for writer identification. These codebooks could be computed universally for the entire set of writers or for each of the writers separately. The methods based on a universal codebook are generally efficient in terms of computational cost, however, a new codebook is to be generated if the script changes. On the other hand, writer specific codebooks have high computational costs but they could present a generic framework independent of the alphabet under study. We will discuss both these aspects of codebook in Section 4. Before that, in the next section, we will give a brief account of the data sets used in our study.

3. Data sets

We mainly carried out our experimental evaluations on two data sets, IAM and RIMES.

The IAM data set [23] is one of the most well known and widely used data sets in problems such as handwriting recognition and writer identification. It comprises forms with handwritten English text of variable content. The images have been scanned at 300 dpi, 8 bits/pixel, gray-scale. A total of 650 writers have contributed to the data set with 350 writers having only one page, 300 writers with at least two and 125 writers with at least four pages. In order to fit in all the 650 writers in our experiments, we keep only the first two images for the writers having more than two pages and divide the image into two parts for writers who contributed a single page thus ensuring two images per writer, one used in training while the other in testing. Furthermore, in order to have an independent validation set, we used the pages ‘three’ and ‘four’ of the 125 writers having at least four samples in the actual data set.

RIMES, a relatively new data set [24], comprises handwritten letters in French text representing the content sent by individuals to companies or administrations. More than 1300 writers contributed to the data collection by writing up to five letters, making a total of 5600 letters. This data set has also been used in the ICDAR 2009 handwriting recognition contest. In our experimental set up, we have used the RIMES writing samples contributed by 375 different writers.

After having discussed the data sets, we will present in the next section, the proposed methodology that is based on a complementary approach. We will first discuss the codebook features followed by a presentation of the contour features that are based on the visual attributes of writing. We will then see how these two facets of handwriting could be combined together as illustrated in Fig. 1.

4. Codebook based features

In Section 2, we discussed the significant contributions to the field of writer recognition over the last two decades. Lately, the focus of research in the domain has shifted towards the extraction of writer-specific patterns in writing. These redundant writing patterns, generally termed as ‘codebook’ could either be writer specific where the frequent writing forms are extracted separately for each writer [19] or universal where these forms are extracted globally (either from the data set under study [20] or from an independent data set [1]). In each of these studies, the codebook is computed from graphemes which could represent over, under or well segmented characters [25].

The approach that we present is inspired by the same idea of frequent writing shapes. We, however, chose to work on a much smaller scale of observation. The methods [1,19,20], although text-independent, are more linked to the way the characters are drawn and segmented and the extracted graphemes might also

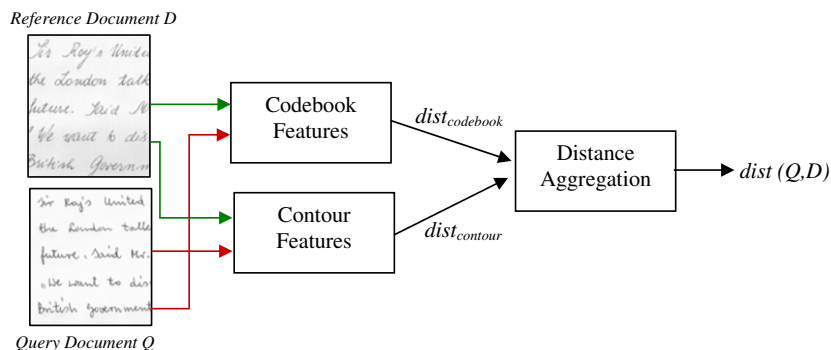


Fig. 1. Overview of the proposed methodology.

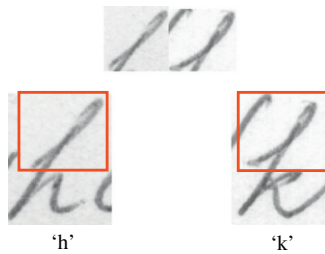


Fig. 2. Presence of similar loops in two different characters.

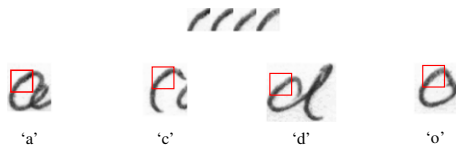


Fig. 3. Same pattern repeated across different characters.

carry some semantic information. We think the recognition of the writer is more linked to the physical way the lines or loops are produced and hence the observation scale may be inferior to that of a grapheme. The fragments that we consider are small parts of handwritten text that do not carry any semantic information and are obtained by a division of handwriting into small windows. The redundant patterns characterizing an individual are then extracted by grouping similar patterns into clusters.

The main advantage of working on small observation windows is that it allows studying the frequent shapes that might be parts of different characters/graphemes. An individual who draws a particular shape (e.g. loops) in a specific way is expected to always employ the same (similar) patterns when drawing that shape, irrespective of the character being written. As an example, Fig. 2 shows two loops that are very much similar but come from two different characters ('h' and 'k', respectively) written by the same author. This redundancy may or may not be captured at the grapheme level depending upon the segmentation scheme employed as well as the characters these loops belong to.

Developing the same idea and reducing the observation scale further, one can identify the presence of redundancy at the level of even smaller fragments. This is explained by the fact that a writer might use the same gesture of hand and hence the same pattern while writing the characters that share similar basic forms. A writer-specific curved fragment for example, is likely to be produced in a similar form for a number of characters as illustrated in Fig. 3 that shows four small fragments which are extracted from four different characters but they are all very close to one another. This redundancy of fragments hence is different from the one at the grapheme level and will be the main focus of our study. We will show how these frequent forms can be determined from a handwritten sample. We will first introduce the proposed division (segmentation) scheme and then discuss how the segmented fragments (represented by a set of features) are grouped into clusters to generate a codebook (writer-specific or universal) which would eventually characterize the writer of a given text sample.

4.1. Division of handwriting image

The division of handwriting is an important step as a 'good' division would allow to exploit most of the redundancy in writing. For our problem, a 'good' division is the one that produces writing

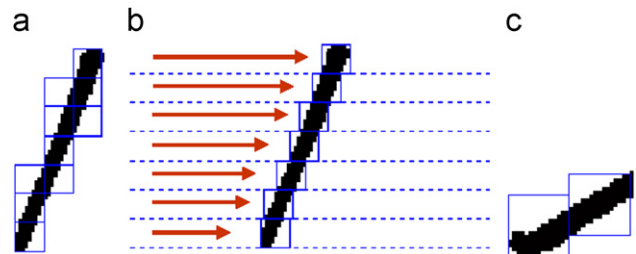


Fig. 4. Division of writing: (a) regular division; (b) horizontal sliding; and (c) adaptive division.

fragments in a way that allows meaningful comparison of these fragments. We have chosen to carry out this division employing square windows of size n . This size should be large enough to contain ample information about the style of the author and small enough to ensure a good identification performance. For our system, the window size was empirically fixed to 13×13 [26].

An important aspect with window positioning is the normalization of writings (character size and ink thickness). If the images under study are all scanned at the same resolution (as in our case) the size of characters in the writing is a writer-dependent attribute and should not be normalized (as opposed to the case of character recognition). For writing instrument independency, we studied the distribution of ink widths in the validation data set and normalized each writing to a fixed thickness using a thinning algorithm. This, however, resulted in a degradation of performance as some writer-specific information is also lost during the normalization. We therefore decided not to normalize the ink thickness. In any case the standard deviation of the ink thickness distribution was not very significant. This observation is in fact in accordance with real world documents where in general, the commonly used writing instruments do not vary too much in their ink widths (of course the images are assumed to be scanned at the same resolution).

In our previous work [26] we introduced a window positioning algorithm where for each connected component in the text, we divide the component into equally spaced lines separated by a distance of n (window size) pixels and slide a window on each of the lines, from left to right, to find the first text (black) pixel that is not already contained within a window. This leads to better positioning as compared to the one achieved by dividing the text regularly in a top-bottom left-right fashion but still suffers from the problems that are linked to the lack of invariance of line position within the window. To overcome these issues, we seek to follow the ink trace in our window positioning algorithm. Since the images are offline, it is not possible to follow the stroke trajectory that was followed by the writer. Nevertheless, we try to follow the ink trace with the objective of achieving an optimal window positioning. The main idea of different window positioning methods has been illustrated in Fig. 4 while the algorithmic details of these methods could be found in [27].

Employing the proposed segmentation scheme, we extract small writing fragments from the text and proceed to the generation of a codebook as discussed in the following.

4.2. Codebook generation

In order to generate a codebook of the characteristic (frequent) sub-images (writing fragments) which could eventually characterize the writer of a sample, we need to group similar fragments into classes. The comparison between the sub-images can be made either on the images directly (pattern matching) or by first extracting a set of features and representing the images in a

feature space. Directly comparing pixel values is simple but suffers from the disadvantage of keeping n^2 pixel values (for a window size of n) and in addition, the resulting comparisons might not be robust to noise and distortions. We therefore chose to represent the patterns by a set of features.

The features that we compute for each window include the horizontal and vertical histograms, upper and lower profiles and a set of well-known shape descriptors (orientation, eccentricity, rectangularity, elongation, perimeter and solidity). The feature values are normalized in the interval $[0\ 1]$ and hence each window is represented by a vector of dimension $d=4n+6$, where n is the window size. If S represents the set of vectors representing the sub-images, we have

$$S = \{S_i\} \quad \text{with each } S_i = (s_i^1, s_i^2, \dots, s_i^d) \quad (1)$$

The (dis)similarity between two patterns is computed by using a distance measure (Euclidean distance) defined on the feature space.

4.2.1. Writer specific codebook

Writing fragments of an individual may be grouped into clusters in a variety of ways. Methods like k-means, fuzzy c-means, learning vector quantization and the closely related self organizing maps have been successfully applied to similar problems of clustering allographs or graphemes [28–30,1]. Generating a writer-specific codebook, however, requires a clustering method that does not need to know a priori the number of clusters to retain as this number would vary from one writing to another. We therefore investigated a number of such algorithms including sequential clustering [31], iterative sequential clustering [2], hierarchical clustering [32] and minimum spanning tree clustering [33]. A comparative evaluation of these methods on the validation data set revealed that hierarchical clustering achieved the best performance on writer identification and verification. The subsequent evaluations were thus carried out using hierarchical clustering to group the writing fragments of a handwritten sample.

In order to find the classes (clusters) that correspond to the frequent writing patterns, we sort them with respect to their cardinality and keep only those having ‘sufficient’ number of elements. Since the number of elements per class depends upon the amount of text in the writing sample, the ‘sufficient’ number cannot be fixed value. We therefore compare the number of clusters against the corresponding area of text pixels covered and pick the top M classes which allow to cover 90% of text pixels in the image [34]. A document image D , having retained M^D classes is thus represented as

$$C^D = \{C_i | 1 \leq i \leq M^D\} \quad (2)$$

And for each class C_i in C^D , we have

$$C_i = \{S_{1,i}, S_{2,i}, \dots, S_{m,i}\}, \quad m = \text{card}(C_i) \quad (3)$$

We then estimate the probability of occurrence of each class $P(C_i)$ and also find the mean vector representing the class \bar{S}_i . The set of these occurrence probabilities could be considered as a probability distribution h^D , where each bin in h^D would represent the emission probability of the respective pattern by the author of document D . This distribution is then used to characterize the writer of a given sample.

4.2.2. Universal codebook

After having presented the extraction of frequent writing patterns for an individual leading to the generation of a writer specific codebook, we will now extend the same principle to a universal set of basic shapes as proposed at the grapheme level in [1,20]. It should be noted that the codebook is not meant to

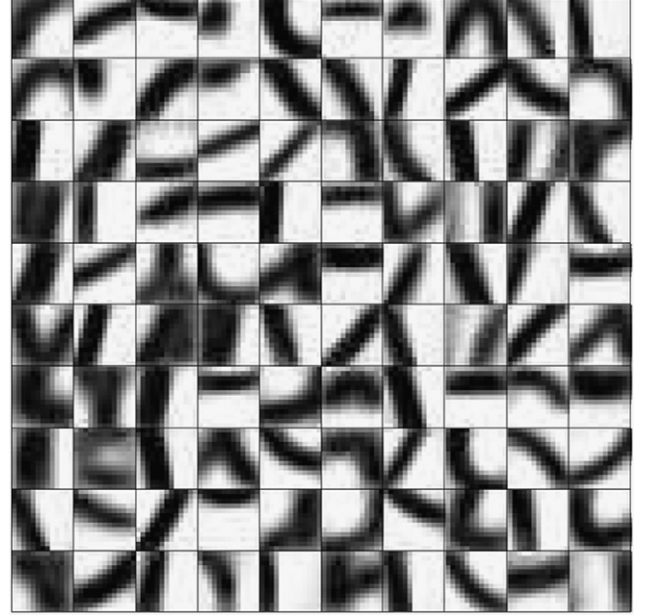


Fig. 5. Universal codebook of size 100 generated from 50 samples of the data set RIMES.

represent an exhaustive list of all possible fragment shapes of a particular script and alphabet. Rather, its objective as described in [1] is to ‘span a shape space and act as a set of nearest-neighbor attractors for the patterns extracted from a given writing sample’.

We apply the principle of universal codebook generation on the small writing fragments obtained by the segmentation method discussed earlier. Since we have chosen to work on a codebook that is not produced from the data set under study, the codebook is generated from handwritten samples of the RIMES data set to evaluate the writings of IAM data set and vice versa. A total of 50 writing samples are used to produce the codebook while the fragments are clustered using k-means algorithm in the feature space with k being varied from 50 to 750 and finally fixed to 100 after evaluations on the validation set. Fig. 5 shows an example codebook (gray levels correspond to intra-cluster variability) and comparing it with the one obtained in [1] where one can notice the presence of graphemes that correspond to complete characters, the shapes that we consider are much more elementary.

Once the codebook has been generated, we find for each writer, the frequencies of producing the patterns in the codebook, the distribution being characteristic of the writer. For a document image we initialize a histogram with all bins set to zero, the number of bins being equal to the size of the codebook. For each of the extracted sub-images (fragments) in the document, we find its nearest entry in the codebook and count it in the respective histogram bin. The distribution is finally normalized and is used to characterize the writer.

5. Contour based features

After having studied the redundant patterns of writing in characterizing the writer, we will now analyze the different visual attributes of writing that allow to capture the writing style of its author. The most important of these attributes is the overall writing orientation which is known to be a stable parameter [35] with the assumption that the writing under consideration represents the natural writing style of the writer and is not a

forged one. Besides orientation, curvature is known to be another fundamental characteristic of handwriting [36] and the features based on curvature have shown effective performance in characterizing writing styles [37] and writers [38] and, recognition of characters [39,40] and numerals [41]. Inspired by the power of these attributes of handwriting, we endeavored to design a set of features that would capture these characteristics and allow us to characterize the writer of a handwritten text sample. This section is devoted to the discussion of these features.

The proposed features are computed from the writing contours which encapsulate the writing style of the author and allow to preserve the writer-dependent variations between character shapes. We have chosen to represent the contours in two ways that correspond to two different observation scales and writing details. These representations include the Freeman chain codes and a set of polygons approximating the contours. The features computed from each of these representations are discussed in the sections to follow where we first introduce the chain code based features and then we present how we have modeled some loss of details keeping only a general and simple view of the writing that is enough to characterize its writer.

5.1. Chain code based features

Chain codes have shown effective performance on problems like shape registration [42] and object recognition [43]. In case of handwritten document images, these directional features and their improved variants have been successfully applied to character/word recognition [44–47] as well as classification of writing styles [48]. Since the handwritten characters issued by a particular writer can be regarded as having a specific shape/style, chain code based features are likely to work well on tasks like writer recognition as well. We therefore compute a set of features from the chain code sequence of the text contours first at a global level and then from the contour fragments within small observation windows [49].

5.1.1. Global features

At the global level, in order to capture the overall orientation information in the writing, we compute the well-known slope density function (histogram of all the chain codes/slopes f_1) of

the contours. The eight bins of the histogram represent the percentage contributions of the eight principal directions in an individual's writing.

In addition, we also find the histograms of the first (and second) order differential chain codes that are computed by subtracting each element of the chain code from the previous one and taking the result modulo connectivity (8 in our case). The differential chain code at pixel p_i represents the angle θ_i between the vectors $p_{i-1}p_i$ and p_ip_{i+1} and their histogram f_2 (also known as curvature density function) is used as the second feature to represent a handwritten text. Employing the same principle, we also compute the histogram f_3 from the second order derivative of the chain code, which is indicative of the variation of the angles as the contour progresses.

The distributions of chain codes and their differentials give a crude idea about the writing shapes but they might not be very effective in capturing the fine details in writing; we thus propose to count not only the occurrences of the individual chain code directions but also the chain code pairs, in a histogram f_4 , illustrated for two writings in Fig. 6. The bin (ij) of the (8×8) histogram represents the probability of finding the pair (ij) in the chain code sequence of the contours. Extending the same idea, we also compute the $(8 \times 8 \times 8)$ histogram of chain code triplets f_5 , the entry (ij,k) indicating the percentage contribution of the respective triplet in the chain code string of writing contours. It is important to note that all the 64 possible pairs and 512 possible triplets do not exist while tracing the contours and we can only have a total of 44 pairs and 236 triplets.

We would like to precise that there exists a partial redundancy between the distribution of first order differential chain codes (f_2) and the distribution of chain code pairs (f_4) as well as between the distributions of second order differential codes (f_3) and the code triplets (f_5). This will be discussed towards the end of this paper.

The estimate of curvature that we achieve from the differential chain codes and the code pairs and triplets is very local (span of few pixels only). Therefore, an approximation of curvature that is based on a relatively larger neighborhood of a contour pixel would be a better indicative of the stroke curves. We therefore use an estimate of curvature computed from the histograms of contour chain code, presented in [43] for object recognition. A correlation measure between the distribution of directions on both sides (forward and backward) of a contour pixel p_c is used to

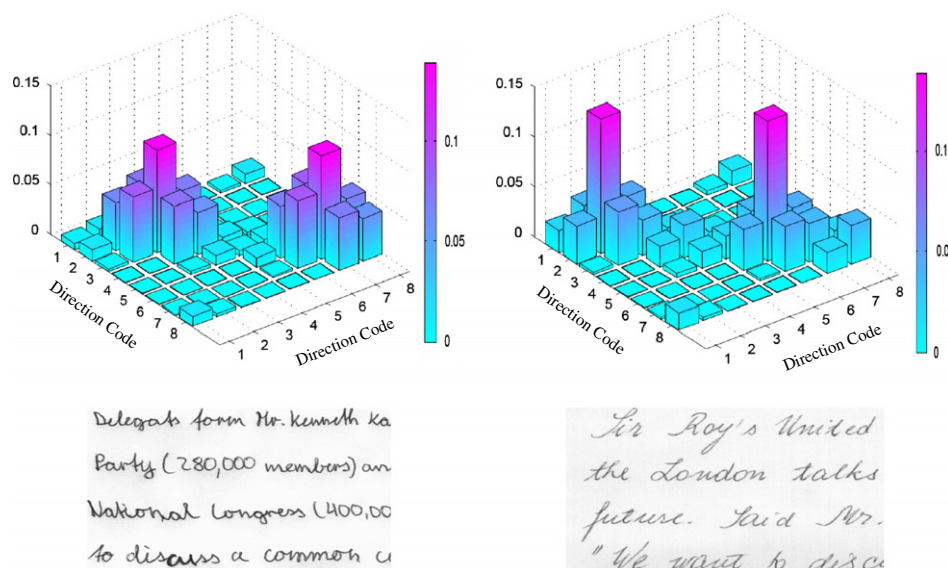


Fig. 6. Two writings and their respective distributions (normalized) of chain code pairs.

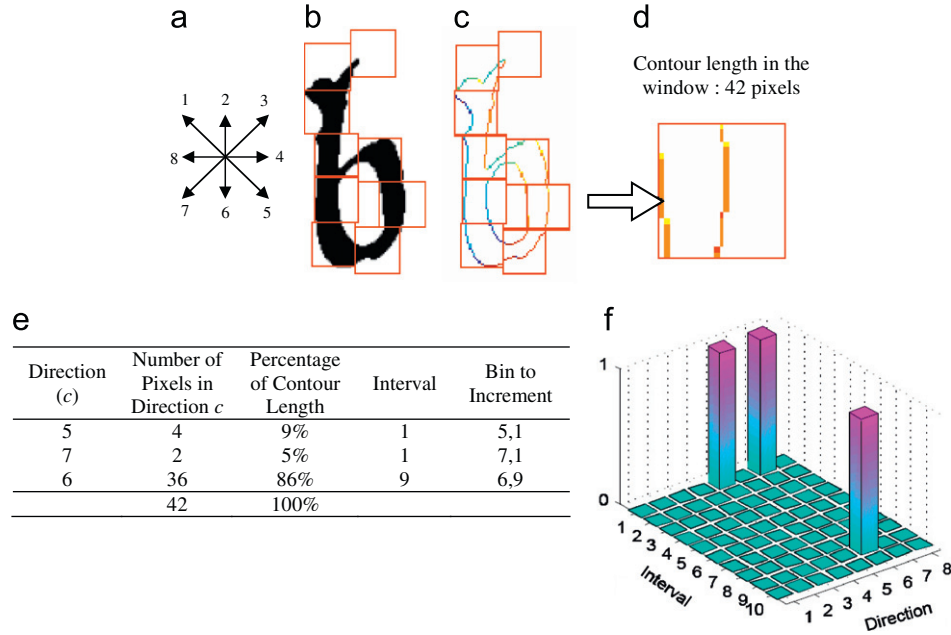


Fig. 7. Windows positioned over a text and respective contour image and contribution of one of the windows to local stroke distribution.

approximate the curvature at p_c and it is the distribution of these estimates which is used to characterize an author (f 6).

After having defined the features that bring some global information on the directions of the strokes, the evolution of directions and the curvature of drawings, we now introduce some local features.

5.1.2. Local features

The features $f 1-f 6$, although computed locally, capture the global aspects of writing and hence the relative stroke information is lost. We therefore chose to carry out an analysis of small stroke (contour) fragments as well. These fragments are extracted by dividing text into small windows in a similar fashion as presented in Section 4.1. An example of windows positioned over a text image and the corresponding contour image is illustrated in Figs. 7(b and c).

For each window, we determine how the eight directions are distributed with respect to the total contour length within the window, the percentages being divided into p intervals. We build an accumulator (local stroke direction distribution $f 7$) which is a two dimensional $d \times p$ array where d is the connectivity (eight directions). The accumulator is initialized with all bins set to zero. For each window w , containing the chain code sequence C^w , the bin (i, j) of the histogram (accumulator) is incremented by 1 if the frequency of direction i is represented in the j th interval, where j is given by

$$j = \text{ceil}\left(\frac{\text{card}(C_i^w)}{\text{card}(C^w) \times p} \times 100\right)$$

$$\text{With } C_i^w = \{c_k \in C^w | c_k = i\} \quad \text{and} \quad i = 0, 1, \dots, 7 \quad (4)$$

The number of intervals p should be large enough to capture the differences in the distributions of directions in the windows and small enough to allow marginally different distributions contribute to the same intervals of $f 7$. We have made p vary from 2 to 20 on the validation data set and finally chosen a value of p equal to 10 for our system.

Fig. 7 shows the contribution of one of the windows to the histogram where the three direction codes present in the window lead to three contributions to the histogram. The process is carried out for all the windows and the distribution is finally normalized.

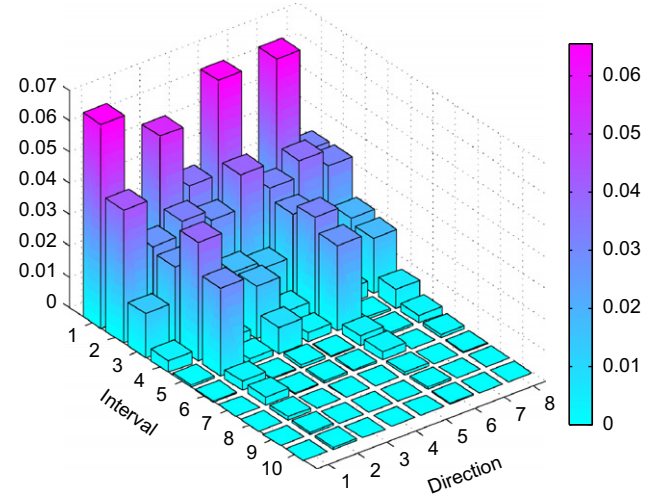


Fig. 8. Normalized local stroke direction distribution of a writing.

This distribution ($f 7$) as indicated in Fig. 8 can thus be considered as a window-based local variant of $f 1$. Using the same idea, we also compute the distributions $f 8$ and $f 9$ from the differential chain codes (local variants of $f 2$ and $f 3$, respectively).

Summarizing, these chain code based features compute the orientation and curvature information of writing. However, these estimates are computed from raw pixels and it would be interesting to carry out a similar analysis at a different observation level. We therefore propose to estimate the contours by a set of polygons and then proceed to feature extraction (a set of global features) which not only corresponds to a distant scale of observation but the computed features are also more robust to noise distortions. These features are discussed in the following section.

5.2. Polygon based features

These features are aimed at keeping only the significant characteristics of writing (enough to characterize its author)

discarding the minute details [50]. For this purpose, we carry out an estimation of the writing contours by a set of line segments, employing the sequential polygonization algorithm presented in [51]. The algorithm requires a user defined parameter T that controls the accuracy of approximation. Larger values of T create longer segments at the cost of character shape degradation and vice versa. Fig. 9 shows the polygon estimation of the contours of a handwritten word for different values of T . For our system, we have used a value of T equal to 2, chosen empirically on the validation set. We then extract a set of features from these line segments.

From the polygonized contours of the handwritten image, we first compute the slope of each of the line segments and employ their distribution (f_{10}) for characterizing the writer. Each line is identified as belonging to one of the bins (classes) illustrated in Fig. 10. These bins are chosen in such a way that the lines having nearly the same orientations as the principal directions (vertical, horizontal etc.) fall in their respective classes.

Not only the number of slopes in a particular direction is important but their corresponding lengths as well, so in order to complement the distribution f_{10} , we also compute a length-weighted distribution of slopes (f_{11}), where for each segment at slope i , the bin i in f_{11} is incremented by the length of the segment. The distribution is finally normalized by the total length of segments in the image. The distributions f_{10} and f_{11} can thus

be viewed as the number and the length of segments belonging to the pre-defined segment classes, respectively.

A histogram based estimate of curvature, calculated from the chain code representation of the contours has been presented in Section 5.1.1. From the polygonized version of writing, we compute the angle measurement between two connected straight segments as

$$\alpha_i = \pi - \arccos \frac{V_i \cdot V_{i+1}}{|V_i| |V_{i+1}|} \tag{5}$$

With V_i and V_{i+1} being the vectors from (x_{i-1}, y_{i-1}) to (x_i, y_i) and from (x_i, y_i) to (x_{i+1}, y_{i+1}) , respectively, as illustrated in Fig. 11.

The distribution of these angles is then employed as our next feature (f_{12}). The angle bins (0° – 180°) are partitioned in a similar fashion as for the slopes. Similarly, in order to take into account the lengths of the segments forming a particular angle, a length-weighted version of f_{12} , f_{13} is also computed where each bin of f_{13} is incremented by the sum of lengths of the two segments forming a given angle. The distribution is then normalized to have a sum equal to 1.

Finally, irrespective of the orientation, it would be interesting to analyze the distribution of the lengths of segments in a writing. Generally, smooth strokes will lead to longer and fewer segments while shaky strokes will result in many small segments, thus the straight segment lengths could be useful in distinguishing the writings of different authors. We therefore use the distribution of these lengths (f_{14}) as a writer specific feature. An important issue here is how to determine the number and partitioning of the bins in f_{14} as the lengths are not normalized. This is done by studying the variation of lengths on the validation data set and setting an upper limit on the allowed segment length $maxL$, above which the segments are considered to be spurious and are discarded. In our case, this value was set to 100 pixels.

Summarizing, we extract a set of 14 (normalized) distributions to represent a document image. The distributions for which the

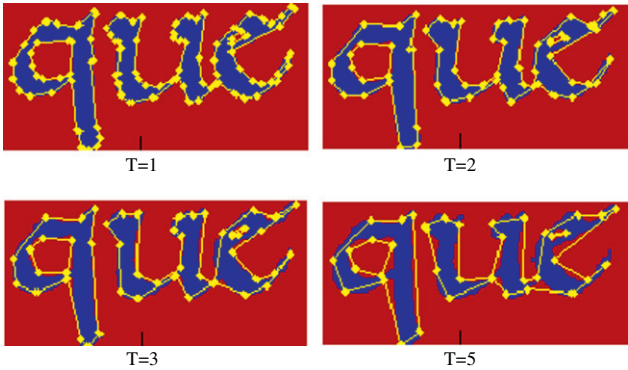


Fig. 9. Polygonization at different values of T .

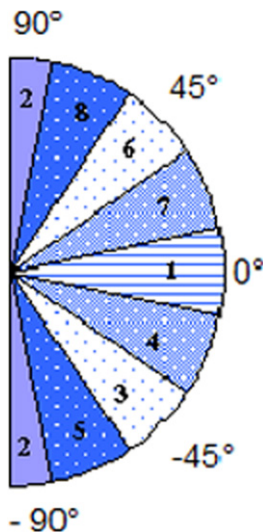


Fig. 10. Division of slopes (-90° to 90°) into bins and the corresponding segment classes.

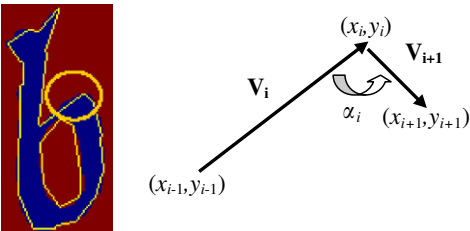


Fig. 11. Curvature (angle) between two connected segments.

Table 2
Summary of contour features.

Feature	Description	Dimension
f_1	Distribution of chain codes	8
f_2	Distribution of 1st order differential chain codes	7
f_3	Distribution of 2nd order differential chain codes	8
f_4	Distribution of chain code pairs	44
f_5	Distribution of chain code triplets	236
f_6	Distribution of curvature indices	11
f_7	Local stroke direction distribution	80
f_8	f_2 computed locally	70
f_9	f_3 computed locally	80
f_{10}	Distribution of segment slopes	8
f_{11}	Length-weighted distribution of segment slopes	8
f_{12}	Distribution of curvatures	8
f_{13}	Length-weighted distribution of curvatures	8
f_{14}	Distribution of segment lengths	10
Total		586

number of bins is not discussed explicitly are partitioned empirically on the validation set. Table 2 summarizes the proposed features with the dimensionalities of each.

6. Writer recognition

Once the handwriting samples have been represented by their respective features, we proceed to writer recognition. We will first define a (dis)similarity measure between two writing samples and then present how it is employed to perform writer identification and verification.

6.1. Dissimilarity measure

In order to compare two handwritten documents, we need to compute the distances between their respective feature values and hence define a dissimilarity measure. We evaluated a number of distance measures including: Minkowski (order 1–5), χ^2 , Bhattacharyya, (Non-)Intersection and Hamming distance. Among these, χ^2 distance performed the best in our series of experimentation hence the results that we will report in the subsequent sections will be based on χ^2 distance:

$$\chi^2(p, q) = \sum_{i=1}^{dim} \frac{(p_i - q_i)^2}{p_i + q_i} \quad (6)$$

where p and q represent the two histograms (distributions) to be compared, p_i and q_i are the entries in bin i of the histograms and dim represents the total number of bins in the histograms.

For contour based features, the computation of dissimilarity between two documents D and Q is quite straight forward and is given by the (χ^2) distance between features f_i^D and f_i^Q .

For codebook based features, we divide the text in the query image Q into small sub-images as discussed in Section 4.1. When using a writer-specific codebook, each of the segmented sub-images is assigned to one of the clusters in the reference document. To compare the questioned document Q with a reference document D having M^D clusters, we build a histogram h^{DQ} with one bin allocated to each cluster of D . For every pattern p (represented by its respective feature vector) in the test document, its nearest cluster is found using the Euclidean distance (with the mean of each cluster) and the occurrence is counted in the respective histogram bin:

$$b = \operatorname{argmin}_j (\operatorname{dist}(p, \bar{s}_j)); \quad h_b^{DQ} \leftarrow h_b^{DQ} + 1$$

$$j = 1, \dots, M^D; \quad M^D = \operatorname{card}(C^D) \text{ And } \bar{s}_j = \text{mean of cluster } j \quad (7)$$

Thus in fact, the questioned document is represented in the feature space of the reference document and the distance between the two documents is computed by calculating the (χ^2) distance between the respective distributions h^D and h^{DQ} .

For a universal codebook, we compare the sub-images (patterns) extracted from the questioned document to the shapes in the codebook and hence find the occurrence probabilities of the codebook shapes for a particular writer (similar to Section 4.2.2). Two writings are then compared by computing the distance between their respective probability distributions (h^D and h^Q) of producing the codebook patterns.

6.2. Writer identification

Writer identification is performed by computing the distance between the query image Q and all the images in the training data set, the writer of Q being identified as the writer of the document that reports the minimum distance. This corresponds to the

nearest neighbor classification (knn with $k=1$). For a query document, we not only find the nearest neighbor (Top-1) but a longer list up to a given rank (Top-K) thus increasing the chance of finding the correct writer in the retrieved list.

6.3. Writer verification

For writer verification, we compute the distance between two given samples and consider them to be written by the same person if the distance falls within a predefined decision threshold. Beyond the threshold value, we consider the samples to be written by different writers. Varying the acceptance threshold the ROC curves are computed and the verification performance is quantified by the Equal Error Rate (EER), the point on the curve where the False Acceptance Rate (FAR) equals the False Rejection Rate (FRR). The lower the equal error rate value, the higher the accuracy of the system. The identification and verification results are presented in the following section.

7. Experimental results

The experimental study was carried out on the writing samples from IAM and RIMES data sets presented in Section 3. For writer identification we report the Top-1 and Top-10 identification rates while for writer verification we present the Equal-Error-Rate (EER). We will first present the individual performances of codebook and contour features and then that of their combination. We will also analyze how stable the proposed features are with respect to certain parameters.

7.1. Performance of codebook features

We first performed a comparative evaluation of writer-specific and universal codebooks on the two data sets as summarized in Table 3 (numbers represent percentages). A comparative analysis of the performance of the two reveals that on both data sets, representing writings into a common codebook space leads to better results on writer identification and verification as compared to representing them in a writer-specific space. So in our subsequent experimentation we will report the results based on universal codebook only (which will be termed as f 15 from here onwards).

7.2. Performance of contour features

For contour based features, we first present the performance of individual features in Table 4. The results are grouped according to the feature type (global, local or polygon-based). Although the performance of the features varies significantly, it can be noticed that, for a chosen feature, the performance is more or less consistent across the two data sets. Another interesting observation is that the local variants (f 7, f 8 & f 9) of the chain code based features outperform their global counter parts (f 1, f 2 & f 3). Overall we achieved maximum identification rate of 79% (f 5) on IAM and 78% (f 7) on RIMES images. For writer verification,

Table 3

Performance (expressed in percentages) of codebook features on the two data sets.

Data set	IAM (650 writers)			RIMES (375 writers)		
	Top1	Top10	EER	Top1	Top10	EER
Codebook						
Writer-specific	81	94	5.44	69	84	11.14
Universal	84	96	4.49	74	85	10.57

Table 4

Performance of individual contour features on the two data sets.

Data set		IAM (650 writers)			RIMES (375 writers)		
Feature class	Feature	Top1	Top10	EER	Top1	Top10	EER
Global	f_1	36	74	7.23	48	77	10.32
	f_2	34	76	6.89	50	74	11.40
	f_3	42	81	6.56	52	76	11.47
	f_4	67	88	5.67	68	85	8.05
	f_5	79	93	4.64	75	91	6.70
	f_6	43	77	6.96	55	77	10.87
Local	f_7	77	93	3.86	78	92	7.02
	f_8	46	83	7.10	58	82	10.25
	f_9	42	79	7.95	55	79	11.39
Polygon	f_{10}	55	86	5.82	64	86	8.21
	f_{11}	58	87	5.42	65	88	7.98
	f_{12}	37	75	6.97	52	76	11.12
	f_{13}	40	78	6.56	52	79	10.29
	f_{14}	31	72	7.51	51	75	12.04

we achieve equal error rates of as low as 3.86% and 6.70% on the two data sets, respectively.

After having evaluated the performance of individual contour features, we combine them by computing the distance between two writings as an average of the distances between the individual features. Since the individual distances d_i can be of quite different dynamic ranges and combining these distances the distance with larger magnitude might dominate the others, we first need to carry out a normalization. In order to be less sensitive to the outliers, we employ a Gaussian normalization with the assumption that the values in d_i are distributed normally.

In the training set of N writing samples, we compute the individual distances (d_i) between each pair of images I_m and I_n . For N images, there are $N \times (N-1)/2$ possible values for each of the distances. Treating each d_i as a data sequence we find its mean μ_{d_i} and standard deviation σ_{d_i} . These values are computed offline and are based on the assumption that N is large enough so that the calculated values are good estimates of the true mean and standard deviation of the distances (d_i) between images. When a questioned document Q is presented to the system, we first compute the raw distances (with respect to each of the features) between Q and all the documents in the database. These distances are then normalized to the range [0 1] using μ_{d_i} and σ_{d_i} [52].

Once the distances are normalized, we proceed to feature combination which is performed by averaging the distances of the features participating in the combination. Among the various combinations tested, we will report a subset of results only that corresponds to the natural combinations as summarized in Table 5. As with the individual features, the performance of combined features is more or less consistent across the two data sets. Combining all the contour features (f_1 – f_{14}) we achieve overall identification rates of 89% and 85% (Equal Error Rates of 2.46% and 4.87%) on IAM and RIMES, respectively.

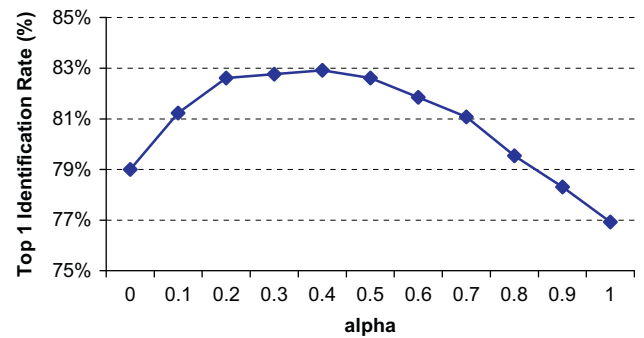
The results that we presented are based on simple distance averaging where each feature contributes equally to the sum. We also evaluated a weighted combination of distances, the weights being chosen with respect to the performance of individual features on the validation data set:

$$D(Q,D) = \sum_{i=1}^{No. \text{ of features}} w_i d_i(Q,D) \quad (8)$$

where d_i represents the distance between the two documents (Q and D) computed for feature i and is given by Eq. (6). This

Table 5Performance of combined contour features (f_1 – f_{14}) on the two data sets.

Data set		IAM (650 writers)			RIMES (375 writers)		
Feature combination		Top1	Top10	EER	Top1	Top10	EER
Global		81	93	4.08	77	92	6.18
Local		81	95	3.76	77	89	7.85
Polygon based		83	97	2.77	81	93	5.16
Global & Local		83	96	3.81	80	91	6.65
Global & Polygon		85	96	3.32	82	92	5.92
Local & Polygon		87	97	3.03	83	93	5.11
All features		89	97	2.46	85	93	4.87

**Fig. 12.** Writer identification performance for a weighted combination of f_5 and f_7 .**Table 6**Combined performance of contour and codebook features (f_1 – f_{15}).

Performance Features: f_1 – f_{15}			
Data set	Top1	Top10	EER
IAM (650 writers)	91	97	2.23
RIMES (375 writers)	84	93	4.90
IAM+RIMES (1025 writers)	88	95	2.86

weighted combination, however, resulted only in marginal improvements in the over all performance of the system. Fig. 12 illustrates the weighted combination of two best performing features (f_5 and f_7) where the final distance between two writings is computed as: $(1-\alpha)d_5 + \alpha d_7$ and it can be observed that the results exhibit only slight improvement in performance as compared to simple distance averaging ($\alpha = 0.5$).

7.3. Combining codebook and contour features

After having presented the codebook and contour features, it would be interesting if we could combine the feature sets. In addition to the two data sets, we will evaluate the performance of the combined feature set (f_1 – f_{15}) by merging the two data sets into one large set of (650+375) 1025 writers. This data set represents, in terms of number of writers, the largest data set used for text independent writer recognition until present. The results are presented in Table 6 where we achieve identification rates of 91%, 84% and 88% on the IAM, RIMES and the combined data sets, respectively. For the data set RIMES, there is a marginal decrease in the performance when combining the codebook and contour features. This is linked to a relatively weak performance of codebook features on the RIMES data set.

Finally, we would present a performance comparison of our method with some recent studies on writer identification. We will first compare the results of our codebook features with the ones achieved by generating a codebook at the grapheme level. For the writer-specific codebook [19] reports an identification rate of 98% on 88 writers but since the data set employed is not the same, the comparison would not be meaningful. Later studies [20] by the same authors achieve an identification rate of 86% on 150 writers employing a codebook generated from the entire data set. The best performance so far has been reported in [1] where the authors achieve an identification rate of 80% with a codebook generated from an independent data set. By changing the observation scale from grapheme to small fragments, we achieve an identification rate of 84%. It is, however, important to precise that the evaluation criterion in [1] is not the same as ours. We have distinguished the training and test sets while in [1] the authors have used a leave-one-out approach on the entire data set. Thus, in order to present an honest comparison we also carried out a similar experimentation and achieved an identical identification rate of 80%. Our method, however, relies on a much smaller codebook size (100 against 400) and the patterns contributing to the codebook are issued from a very generic segmentation scheme.

Comparing the overall system performance (f 1– f 15) on 650 writers of IAM data set, we achieve an identification rate of 91% (89% using leave-one-out method) which is comparable to the best results reported so far on this data set [1].

Two important parameters that influence the performance of the system are the number of writers and the amount of text available for each writer. It would therefore be interesting to analyze how these parameters affect the performance as presented in the following section.

7.4. Stability of features

The influence of amount of available text for each writer is studied on the 300 writers (having contributed at least two pages) of the IAM base. This allows varying the amount of text from one word to complete page (instead of half a page if we also use the 350 writers having contributed only one page).

The corresponding writer identification rates for codebook and contour features are summarized in Fig. 13. Naturally, the contour-based features are relatively more stable when the amount of text is varied where the identification rate rises from about 20% for a single word to 91% for the complete page (5–10

lines on the average). Comparing the two codebooks, the performance of universal codebook is relatively less sensitive to the amount of text. The performance (specially the Top10 identification rates) begins to stabilize a bit from three lines of text onwards.

Regarding the influence of number of writers on the identification rate, the evaluations are carried out on the IAM data set by varying the number of writers from 10 to 650. It can be seen from Fig. 14 that there is a consistent but not dramatic decrease in the identification rate as the number of writer increases starting with 100% for 10 writers and dropping to 89%, 84% and 81% for features computed from contours, universal and writer-specific codebooks, respectively, for a total of 650 writers. Unsurprisingly, the Top-10 performances are much more stable dropping to and 97%, 96% and 94%, respectively, when the entire data set is used.

Finally, we will be interested to study the evolution of universal codebook performance with respect to the number of writing samples used to generate the codebook as well as the size of codebook. This is carried out (on 300 writers of the IAM data set with two samples each) by fixing one of these two parameters and varying the other as illustrated in Fig. 15. It can be noticed that performance shows a consistence decrease as the size of the codebook increases beyond 250. This is quite natural as increasing the dimension of the representation space too much would imply that the writing in question would be compared to a large codebook with occurrence frequencies of codebook entries not being very characteristic of the writer. Another interesting observation is that the number of samples used to generate the codebook does not cause a dramatic change in the performance.

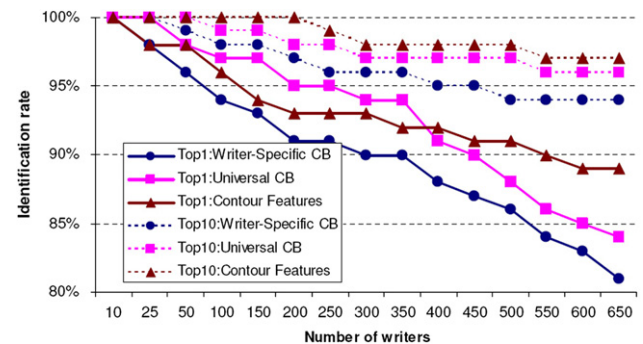


Fig. 14. Performance stability of features with respect to number of writers.

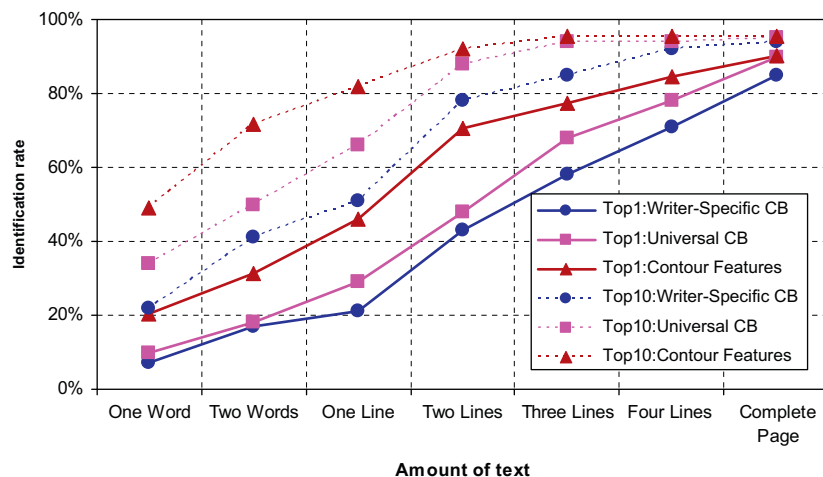


Fig. 13. Performance stability of features with respect to amount of text (on 300 writers of IAM data set).

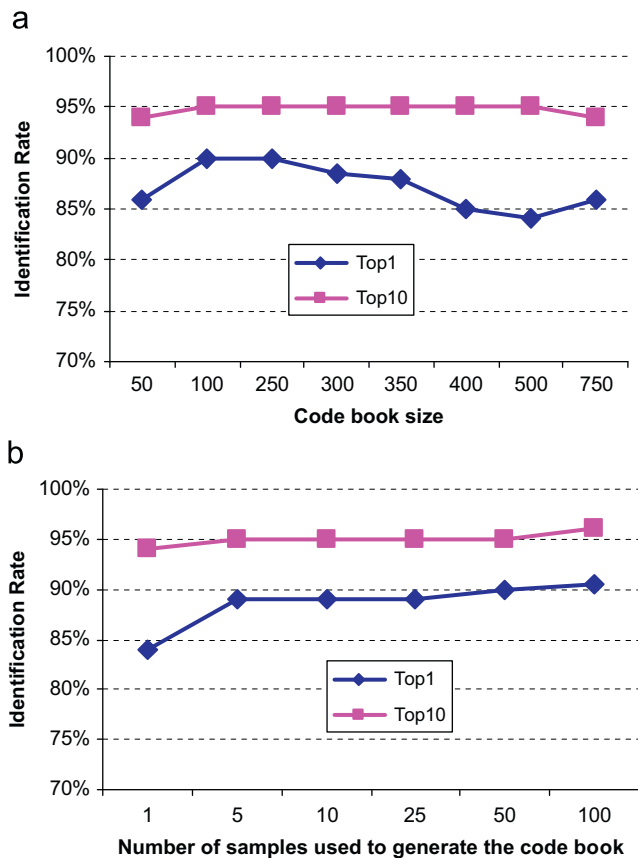


Fig. 15. Writer identification performance (on 300 writers) as a function of (a) codebook size and (b) the number of samples used to generate the codebook. (a) Codebook generated from 50 samples and (b) codebook size fixed to 100.

Employing only one writing sample to produce the codebook results in an identification rate of 84% which further validates the argument that the codebook needs not to be exhaustive.

8. Conclusion

In this paper we presented an original method for automatic writer identification and verification from offline scanned images of handwriting. Our method exploits two different facets of handwriting: the existence of certain redundant patterns in writing and the visual attributes of orientation and curvature characterizing the writer of a handwritten text. Contrary to the classical approaches which analyze the redundancy of writing shapes at the grapheme level, we exploited it at a much smaller scale of observation that corresponds to small writing fragments and allows to capture the redundancy of writing gestures which might be common across different graphemes. We next used the orientation and curvature information in writing which is extracted by computing a set of features from writing contours at different observation levels. Finally we combined the two feature sets and achieved very promising results on tasks of writer identification and verification which are comparable to the best results achieved so far in the domain. We also evaluated our system on the largest data set (in terms of number of writers) used so far in a study on text-independent writer recognition and the system was able to perform reasonably well without degrading the recognition rates too much.

The proposed features are simple to compute but are very effective and generic. We also evaluated them on Arabic hand-

written documents from the IFN/ENIT database [53] without changing any of the parameters and obtained very encouraging results (92% identification rate and 2.94% EER on 100 writers).

Concluding, we will enumerate some interesting research directions laid down by our work. First of all, for the extraction of frequent writing shapes, it would be better if one could determine the window size (for division of text) automatically depending upon the writing details rather than using a fixed size for all writings since in real world problems, one may need to compare writing images that are scanned at different resolutions. This would in turn require a representation that would allow a comparison of writing fragments at different sizes. The idea of redundant writing patterns may be developed further to study whether an individual who writes more than one script (for examples someone who writes French & Arabic) share some basic patterns/shapes across the two scripts. This could be realized by comparing the writer-specific codebooks computed from the texts written in the two scripts. It would be very interesting to analyze if one could recognize the writer of a text written in one script from the samples of same writer written using a different script or, identify the script of a given text based on the basic writing shapes present in the text, comparing them to a universal codebook of the respective alphabet.

For the contour-based features we tried to capture the direction and curvature information in writing by a set of distributions and it would be a good idea to extract this information at different image resolutions which in fact will correspond to different levels of observation. The features at different resolutions can then be combined and their effect on the performance could be studied to identify the observations scales which might be more effective for the problem of writer recognition.

Finally, as we discussed earlier there exists a partial redundancy among certain features of the proposed set f_1 – f_{15} . We therefore conducted a study on the relevance of the proposed features using a feature selection process (Genetic algorithm with a wrapper). A complete description of this study is out of the scope of this paper. However, we would like to summarize by stating that keeping approximately 50% of the features only we were able to maintain similar identification rates.

References

- [1] M. Bulacu, L. Schomaker, Text-independent writer identification and verification using textural and allographic features, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 29 (4) (2007) 701–717.
- [2] A. Nosary, L. Heutte, T. Paquet, Y. Lecourtier, Defining writer's invariants to adapt the recognition task, in: *ICDAR '99: Proceedings of the Fifth International Conference on Document Analysis and Recognition*, 1999, pp. 765–768.
- [3] R. Plamondon, G. Lorette, Automatic signature verification and writer identification—the state of the art, *Pattern Recognition* 22 (2) (1989) 107–131.
- [4] R.A. Huber, A. Headrick, *Handwriting Identification: Facts and Fundamentals*, CRC Press, Florida, USA, 1999.
- [5] S.N. Srihari, S.-H. Cha, H. Arora, S. Lee, Individuality of handwriting, *Journal of Forensic Sciences* 47 (4) (2002).
- [6] B. Zhang, S.N. Srihari, S. Lee, Individuality of handwritten characters, in: *ICDAR '03: Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 1086–1090.
- [7] S.N. Srihari, C.I. Tomai, B. Zhang, S. Lee, Individuality of numerals, in: *ICDAR '03: Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 1096–1100.
- [8] B. Zhang, S.N. Srihari, Analysis of handwriting individuality using word features, in: *ICDAR '03: Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 1142–1146.
- [9] V. Pervouchine, G. Leedham, Extraction and analysis of forensic document examiner features used for writer identification, *Pattern Recognition* 40 (3) (2007) 1004–1013.
- [10] E.N. Zois, V. Anastassopoulos, Morphological waveform coding for writer identification, *Pattern Recognition* 33 (3) (2000) 385–398.

- [11] H.E.S. Said, T.N. Tan, K.D. Baker, Personal identification based on handwriting, *Pattern Recognition* 33 (2000) 149–160.
- [12] U.-V. Marti, R. Messerli, H. Bunke, Writer identification using text line based features, in: *ICDAR '01: Proceedings of the Sixth International Conference on Document Analysis and Recognition*, 2001, pp. 101–105.
- [13] V. Boulétreau, N. Vincent, R. Sabourin, H. Emptoz, Handwriting and signature: one or two personality identifiers?, in: *ICPR '98: Proceedings of the 14th International Conference on Pattern Recognition*, 1998, pp. 1758–1760.
- [14] A. Seropian, M. Grimaldi, N. Vincent, Writer identification based on the fractal construction of a reference base, in: *Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003.
- [15] M. Bulacu, L. Schomaker, L. Vuurpijl, Writer identification using edge-based directional features, in: *ICDAR '03: Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 937–941.
- [16] A. Schlappbach, H. Bunke, Using HMM based recognizers for writer identification and verification, in: *IWFHR'04: Proceedings of the Ninth International Workshop on Frontiers in Handwriting Recognition*, 2004, pp. 167–172.
- [17] A. Schlappbach, H. Bunke, Off-line writer identification using Gaussian mixture models, in: *ICPR '06: Proceedings of the 18th International Conference on Pattern Recognition*, 2006, pp. 992–995.
- [18] A. Schlappbach, *Writer Identification and Verification*, IOS Press, Amsterdam, The Netherlands, 2007.
- [19] A. Bensefia, A. Nosary, T. Paquet, L. Heutte, Writer identification by writer's invariants, in: *IWFHR '02: Proceedings of the Eighth International Workshop on Frontiers in Handwriting Recognition*, 2002, pp. 274–279.
- [20] A. Bensefia, T. Paquet, L. Heutte, A writer identification and verification system, *Pattern Recognition Letters* 26 (13) (2005) 2080–2092.
- [21] L. Schomaker, M. Bulacu, Automatic writer identification using connected-component contours and edge-based features of uppercase western script, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 26 (6) (2004) 787–798.
- [22] M. Bulacu, L. Schomaker, A comparison of clustering methods for writer identification and verification, in: *ICDAR '05: Proceedings of the Eighth International Conference on Document Analysis and Recognition*, 2005, pp. 1275–1279.
- [23] U.-V. Marti, H. Bunke, The IAM-database: an English sentence database for offline handwriting recognition, *International Journal on Document Analysis and Recognition* 5 (2002) 1433–2825.
- [24] E. Grosicki, M. Carré, J.-M. Brodin, E. Geoffrois, RIMES evaluation campaign for handwritten mail processing, in: *ICFHR'08: Proceedings of 11th International Conference on Frontiers in Handwriting Recognition*, 2008.
- [25] L. Heutte, T. Paquet, J.V. Moreau, Y. Lecourtier, C. Olivier, A structural/statistical feature based vector for handwritten character recognition, *Pattern Recognition Letters* 19 (1998) 629–641.
- [26] I. Siddiqi, N. Vincent, Writer identification in handwritten documents, in: *ICDAR '07: Proceedings of the Ninth International Conference on Document Analysis and Recognition*, vol. 1, 2007, pp. 108–112.
- [27] I. Siddiqi, N. Vincent, Combining global and local features for writer identification, in: *ICFHR '08: Proceedings of the Eleventh International Conference on Frontiers in Handwriting Recognition*, 2008, pp. 48–53.
- [28] F. Chang, C.-H. Chou, C.-C. Lin, C.-J. Chen, A prototype classification method and its application to handwritten character recognition, in: *Proceedings of IEEE International Conference on Systems, Man and Cybernetics*, 2004.
- [29] G.X. Tan, C. Viard-Gaudin, A.C. Kot, Automatic writer identification framework for online hand written documents using character prototypes, *Pattern Recognition* 42 (2009) 3313–3323.
- [30] F. Camastra, A. Vinciarelli, Combining neural gas and learning vector quantization for cursive character recognition, *Neurocomputing* 51 (2003) 146–159.
- [31] M. Friedman, A. Kandel, *Introduction to Pattern Recognition: Statistical, Structural, Neural and Fuzzy Logic Approaches*, World Scientific Publishing Company, Singapore, 1999.
- [32] S.C. Johnson, Hierarchical clustering schemes, *Psychometrika* 32 (3) (1967) 241–254.
- [33] C. Zahn, Graph-theoretical methods for detecting and describing gestalt clusters, *IEEE Transactions on Computers* C 20 (1971) 68–86.
- [34] I. Siddiqi, N. Vincent, How to define local shape descriptors for writer identification and verification, in: *PRIS'08: Proceedings of the Eighth International workshop on Pattern Recognition in Information Systems*, 2008, pp. 199–204.
- [35] F.J. Maarse, A.J.W.M. Thomassen, Produced and perceived writing slant: differences between up and down strokes, *Acta Psychologica* 54 (1–3) (1983) 131–147.
- [36] S.-W. Lee, *Advances in Handwriting Recognition*, World Scientific Publishing Company, Singapore, 1999.
- [37] G. Joutel, V. Eglin, S. Bres, H. Emptoz, Curvelets based queries for cbr application in handwriting collections, in: *ICDAR '07: Proceedings of the Ninth International Conference on Document Analysis and Recognition*, vol. 2, 2007, pp. 649–653.
- [38] M. Bulacu, L. Schomaker, Writer style from oriented edge fragments, in: *Proceedings of the 10th International Conference on Computer Analysis of Images and Patterns*, 2003, pp. 460–469.
- [39] R. Legault, C. Suen, A comparison of methods of extracting curvature features, in: *Proceedings of the 11th IAPR International Conference on Pattern Recognition*, 1992.
- [40] K.T. Miura, R. Sato, S. Mori, A method of extracting curvature features and its application to handwritten character recognition, in: *ICDAR '97: Proceedings of the Fourth International Conference on Document Analysis and Recognition*, 1997, pp. 450–454.
- [41] L. Yang, C.Y. Suen, T.D. Bui, P. Zhang, Discrimination of similar handwritten numerals based on invariant curvature features, *Pattern recognition* 38 (7) (2005) 947–963.
- [42] M.B. Ahmad, J.-A. Park, M.H. Chang, Y.-S. Shim, T.-S. Choi, Shape registration based on modified chain codes, in: *Advanced Parallel Processing Technologies*, 2003, pp. 600–607.
- [43] A. Bandera, C. Urdiales, F. Arrebola, F. Sandoval, 2D object recognition based on curvature functions obtained from local histograms of the contour chain code, *Pattern Recognition Letters* 20 (1999) 49–55.
- [44] H. Yamada, Y. Nakano, Cursive handwritten word recognition using multiple segmentation determined by contour analysis, *IEICE Transactions on Information and Systems* E79-D (1996) 464–470.
- [45] F. Kimura, N. Kayahara, Y. Miyake, M. Shridhar, Machine and human recognition of segmented characters from handwritten words, in: *ICDAR '97: Proceedings of the Fourth International Conference on Document Analysis and Recognition*, 1997, pp. 866–869.
- [46] M. Blumenstein, B. Verma, H. Basli, A novel feature extraction technique for the recognition of segmented handwritten characters, in: *ICDAR '03: Proceedings of the Seventh International Conference on Document Analysis and Recognition*, 2003, pp. 137–141.
- [47] M. Blumenstein, X.Y. Liu, B. Verma, An investigation of the modified direction feature for cursive character recognition, *Pattern Recognition* 40 (2) (2007) 376–388.
- [48] M.E. Dehkordi, N. Sherkat, T. Allen, Handwriting style classification, *International Journal on Document Analysis and Recognition* 6 (2003) 55–74.
- [49] I. Siddiqi, N. Vincent, A set of chain code based features for writer recognition, in: *ICDAR '09: Proceedings of the 10th International Conference on Document Analysis and Recognition*, 2009, pp. 981–985.
- [50] I. Siddiqi, N. Vincent, Combining contour based orientation and curvature features for writer recognition, in: *CAIP'09: Proceedings of the 13th International Conference on Computer Analysis of Images and Patterns*, 2009, pp. 245–252.
- [51] K. Wall, P.-E. Danielsson, A fast sequential method for polygonal approximation of digitized curves, *Computer Vision, Graphics, and Image Processing* 28 (3) (1984) 220–227.
- [52] Y. Rui, T.S. Huang, M. Ortega, S. Mehrotra, Relevance feedback: a power tool for interactive content-based image retrieval, *IEEE Transactions on Circuits and Systems for Video Technology* 8 (5) (1998) 644–655.
- [53] M. Pechwitz, S.S. Maddouri, V. Märgner, N. Ellouze, H. Amiri, IFN/ENIT-database of handwritten arabic words, in: *CIFED '02: Proceedings of the Seventh Colloque International Francophone sur l'Ecrit et le Document*, 2002.

About the Author—IMRAN SIDDIQI received his Masters degree in Image Processing in 2006 from University Paris Descartes, France. He completed his Ph.D. in 2009 and is currently working on writer identification and verification from offline handwritten documents images. His research interests include document analysis and pattern recognition and their applications.

About the Author—NICOLE VINCENT is full Professor since 1996. She presently heads the research group Systèmes Intelligents de Perception (SIP) at the Laboratoire d'Informatique Paris Descartes (LIPADE) in the university Paris Descartes—Paris 5. After studying in Ecole Normale Supérieure and graduation in Mathematics, Nicole Vincent received a Ph.D. in Computer Science in 1988 from Lyon Insa. She has been involved with several projects in pattern recognition, signal and image processing and video analysis. Her research interest concerns document image analysis, image retrieval and video sequence analysis.