

Template-free word spotting in low-quality manuscripts

Huaigu Cao* and Venu Govindaraju

Center for Unified Biometrics and Sensors (CUBS)

Dept. of Computer Science and Engineering

University at Buffalo, Amherst, NY

E-mail: {hcao3, govind}@buffalo.edu*

As the OCR technique is not yet adequate for handwritten scripts with large lexicon, word spotting has been introduced as an alternative to OCR. This paper proposes a novel approach to word spotting that, instead of matching features of the word image to features extracted from predefined templates, uses the estimated posterior probability as the output of well trained classifier for spotting. Gabor features are extracted from gray scale image in order to yield higher performance on degraded, low quality document image.

Keywords: Keyword spotting; OCR.

1. Introduction

Nowadays, Optical Character Recognition (OCR) has advanced enough for a certain applications such as printed document, on-line handwritten document, and off-line hand-written documents with good image quality and relatively smaller lexicon. However, OCR is not adequate for the low-quality handwritten documents with large lexicon. Word spotting¹ has been proposed as an alternative to OCR for indexing and retrieving the keywords of low-quality handwritten scripts.

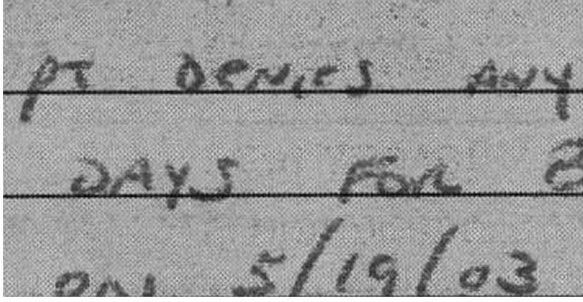
The idea of word spotting is to find all the occurrences of a given word in the document images. Here the input is a word which is usually provided as user input, and the output are all the coordinates of the corresponding word images. Most of the existing approaches to word spotting choose a certain number of word images containing the user query as templates and find the best matches in the dataset. Thus the most important step is to determine the similarity between two word images. Among the existing approaches, the similarity can be defined on either the intensity of the raw images or features extracted from the original images. Several similarities between raw image data using different definitions of distances such as XOR, SSD, and EDM are discussed in.² There are also some similarities defined over feature space such as SLH² (using Scott and Longuet-Higgins algorithm), SC³ (Shape Context matching), DTW⁴ (Dynamic time warping), CORR⁵ (recovering the correspondences between points of interest in two images). The performance of all the above similarities are compared in⁴ which shows that DTW

and CORR approaches have the best performance. B. Zhang, et al⁶ proposed a method based on word shape and claimed to have better performance than the DTW approach in.⁴

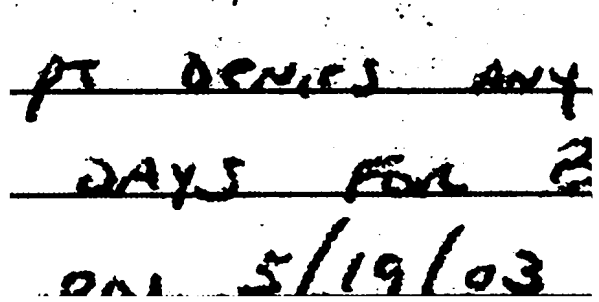
All of the above approaches assume that for any query, there are a few word images stored as templates in a training set. However, since it is not easy to get all the word images of possible user queries in advance, this assumption limits the application of word spotting to a small set of keywords. To solve this problem, in the approach proposed in this paper, matching is performed in character image level rather than in word image level. Since the possible number of characters is very limited, i.e., 26 letters for English text, it will be much easier to build a training set of character images.

Another novelty of our approach is that we use Gabor features extracted from grayscale images. For some document image with degraded quality, e.g., the noisy carbon medical form shown in figure 1, so much information was lost after binarisation that the binarised version is not even readable by human beings, whereas the grayscale version is still readable. Although there are some literatures^{8,9} on grayscale feature extraction for OCR, all the existing feature-matching based word spotting methods still extract features from binary image.

The following sections in this paper are organized as follows: in section 2 we will talk about how to extract features from grayscale word image; in section 3, two similarities are discussed; in section 4, the results of experiments are presented.



(a) a patch from the medical form



(b) the binarized image of the patch

Fig. 1. An example of the carbon copy of medical form.

2. Extracting Gabor features from word image

2.1. Gabor filter and Gabor wavelet

The two dimensional special function of Gabor filter and its two dimensional Fourier transform can be written as follows:⁷

$$g(x, y, F, \sigma_x, \sigma_y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left[-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right) + 2\pi jFx\right] \quad (1)$$

$$G(u, v) = \exp\left\{-\frac{1}{2}\left[\frac{(u - F)^2}{\sigma_u^2} + \frac{v^2}{\sigma_v^2}\right]\right\} \quad (2)$$

where $\sigma_u = 1/(2\pi\sigma_x)$, $\sigma_v = 1/(2\pi\sigma_y)$, and parameter F specifies the central frequency of interest. Gabor wavelet derived from Gabor filter in equation (1) is defined as follows

$$g_{mn}(x, y) = a^{-m}g(x', y', F, n\pi/K, \sigma_x, \sigma_y), \quad (3)$$

where

$$x' = a^{-m}(x \cos \theta + y \sin \theta),$$

and

$$x' = a^{-m}(-x \sin \theta + y \cos \theta).$$

In equation (3), $m = 0, 1, \dots, S-1$ where S is the total number of scales, $\theta = n\pi/K$, and $n = 0, 1, \dots, K-1$ where K is the total number of orientations. The central frequency of interest is F . By varying m and n , we can apply filter $g_{mn}(x, y)$ to the input image (using 2-D convolution) to get features at different scale and orientation. The Gabor wavelet is non-orthogonal, thus there is redundant information in

the filtered image. The following strategy is used to reduce the redundancy.⁷ Let U_l and U_h denote the lower and upper center frequencies of interest. Then the filter parameters a , σ_u and σ_v are determined by

$$\begin{cases} a = (U_h/U_l)^{\frac{1}{S-1}} \\ \sigma_u = \frac{(a-1)U_l}{(a+1)\sqrt{2\ln 2}} \\ \sigma_v = \tan\left(\frac{\pi}{2K}\right) \left[U_l - \frac{2\ln 2\sigma_u^2}{U_l} \right] \left[2\ln 2 - \frac{(2\ln 2)^2\sigma_u^2}{U_l^2} \right]^{-\frac{1}{2}} \end{cases} \quad (4)$$

In order to make the Gabor filter be sensitive to the strokes of the characters, the central frequency of interest should be set to $\frac{1}{2W}$, where W is the stroke width. This is because the stroke width is the half-period of the signal of interest, so the period is $2W$, and the frequency is $\frac{1}{2W}$. In PCR medical forms, W varies from 5 to 8 pixels which correspond to frequencies 0.1 and 0.0625, respectively. Our experiment is different from related work^{8,9} in that instead of applying Gabor filter of only one scale, we applied a Gabor wavelet of two different scales. The upper and lower frequencies of interest are $U_h = 0.1$, and $U_l = 0.05$ so the range of stroke width is completely covered. This seems more reasonable in that the stroke width is not fixed. Experiment of comparison indicates a better performance of using 2-scale Gabor wavelet over single scale Gabor filter.

2.2. Feature extraction from character image

In our dataset, the width and height of character images are between 30 and 40 pixels. So we take 64×64 pixels character image with the character at the center and apply a Gabor wavelet of two scales ($U_h =$

0.1, and $U_l = 0.05$) and four orientations to the images. Then a set of histogram features proposed in^{8,9} are taken from the 48×48 sub-image at the center of the 64×64 transformed image. Specifically, the 48 by 48 output of Gabor filter is divided into N by N blocks. Each block $r(x, y)$ is of size M by M , where $M=48/N$. In each block, histogram features are calculated separately by positive and negative real part outputs weighted by a Gaussian function⁸(See the following equations).

$$F_{x,y}^+ = \sum_{(m,n) \in r(x,y)} G(m-x, n-y) \cdot \max(0, F_K(m, n)) \quad (5)$$

and

$$F_{x,y}^- = \sum_{(m,n) \in r(x,y)} G(m-x, n-y) \cdot \min(0, F_K(m, n)) \quad (6)$$

where $G(x, y) = \exp\{-(x^2 + y^2)/(2\tau^2)\}/(2\pi)$, and F_K is the real part output of Gabor filter. For the scale $m=0$, N , M , τ are set to 4, 12, 6 respectively; for the scale $m=0$, N , M , τ are set to 2, 24, 12, respectively. As a result, totally $2 \times 4 \times (4 \times 4 + 2 \times 2) = 160$ features are extracted.

2.3. Feature extraction from word image

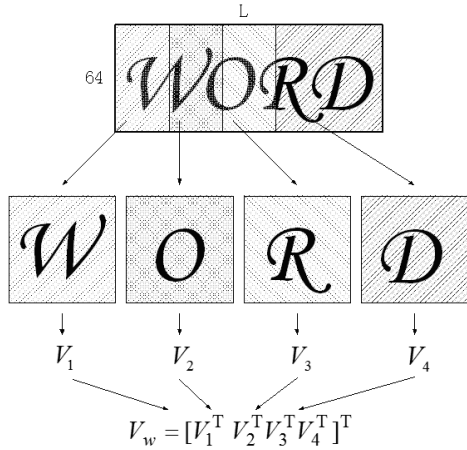


Fig. 2. The diagram indicating how a word image is split into four character images for feature extraction.

Let $I_w[1 : 64, 1 : L]$ denote a 64 by L ($L \geq 64$) word image, where the word is centered vertically in I_w . Suppose I_w contains n characters. Then a simplified definition of the feature vector of I_w is $V_w = [V_1^T V_2^T \dots V_n^T]^T$, where V_i ($1 \leq i \leq n$) is the 160-dimensional Gabor feature vector of the 64 by 64

sub-image

$$I_w[x_0^i : x_0^i + 63, 1 : 64]$$

where $x_0^i = \frac{(n-i)+(i-1)(L-63)}{n-1}$. In other words, a word image is divided evenly into n square blocks of character images allowing overlap. Figure 2 illustrates this feature mapping process. However, characters within a word image are seldom distributed evenly. The solution is to generate three candidates of feature vectors of the following sub-images:

$$I_w[x_0^i : x_0^i + 63, 1 : 64],$$

$$I_w[(x_0^i - 8) : (x_0^i - 8) + 63, 1 : 64], \text{ and}$$

$$I_w[(x_0^i + 8) : (x_0^i + 8) + 63, 1 : 64]$$

for the i 'th character and choose the one with minimum cost using criteria discussed in section 3.

3. Similarity between a word image and the query word

Suppose the query $w = c_1 c_2 \dots c_n$ is word consisting of n characters, and the n -character feature vector of the word image I_w is $V_w = [V_1^T V_2^T \dots V_n^T]^T$. We need to measure the similarity between w and V_w . In our experiments, we tested two different similarities that are based on Euclidian distance and posterior probability, respectively.

3.1. Euclidian distance similarity

We can define the Euclidian distance based similarity between a word image and the query word. First we extract the feature vectors of all the character images in the training set, and calculate the mean feature vector $\mu(c)$ of each character c . Then the squared root of squared sum of Euclidian distances between each V_i^T ($1 \leq i \leq n$) and $\mu(c_i)$ is calculated as a cost function, i.e.,

$$C_E(w, V_w) = \left(\sum_{i=1}^n \|V_i - \mu(c_i)\|^2 \right)^{\frac{1}{2}} \quad (7)$$

The less the cost is, the more similar the word image is to the query. In our experiment, since the training set is not very large, we reduce the feature vector of character image to 40-dimensional by PCA and get

better performance than the method using the original 160-dimensional feature vectors. All the 64 by 64 sub-images are taken from the word images, and their distances to each character are calculated and stored in the dataset so that the execution of query skips the step of feature extraction and is extremely fast.

3.2. Probabilistic similarity

Suppose that the posterior probabilities for character classification, i.e., $P(C_i|V_i)$ ($1 \leq i \leq n$) are known, and that all the characters are independent, then the posterior probability of the word $P(w|V_w) = \prod_{i=1}^n P(c_i|V_i)$. Then we can normalize the probability for different value of n by calculating the cost

$$C_P(w, V_w) = \ln \left[\frac{1}{P(w|V_w)} \right]^{\frac{1}{n}} = -\frac{1}{n} \sum_{i=1}^n \ln(P(c_i|V_i)) \quad (8)$$

The less the cost is, the more similar the word image is to the query. In our experiment, the posterior probabilities of character classification are estimated using the method proposed in¹⁰ from the output of LIBSVM¹¹ (an implementation of SVM classifier) using RBF kernel. All estimated probabilities are stored in our database to speed up the execution of queries.

3.3. WMR similarity

For the purpose of comparison, we also performed spotting test using the Word Model Recognizer (WMR),¹² a word recognizer that uses chain-code features from binarized images and performs word recognition under a scheme of over-segmentation and recognition. In spotting test with WMR, the word recognition distance produced by WMR is taken as a cost function.

4. Results and discussions

Our experiments are done on 12 medical form images. We take 5295 character images from the first ten document images of our dataset as a training set, and do spotting tests on the last two images using as queries all the 101 different words occurred 127 times in the tested document image.

First, all the word images are manually segmented, which can also be implemented automatically later. Then spotting tests using cost functions C_E , C_P and WMR distance are done respectively.

Binarization and line removal algorithms proposed in¹³ are performed before running WMR.

To evaluate the performance of our algorithm the recall-precision curves under all three cost functions are drawn in figure 3. For each cost function, a series of threshold values of the cost function are taken. The values of precision and recall are calculated for each threshold. We can see that the method using cost C_P outperforms that using cost C_E and WMR with gains of equal recall and precision rates of 26% and 15%, respectively. Test results show that our approach has a high performance. Unlike existing approaches,⁴ our approach does not require any word template.

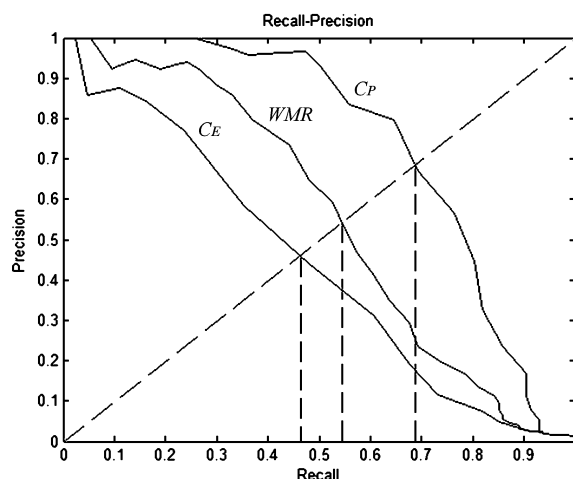


Fig. 3. Recall-precision curves of methods using WMR word recognizer (Precision=Recall=54%), Euclidian distance similarity (Precision=Recall=46%), and probabilistic similarity (Precision=Recall=69%), respectively.

References

1. R. Manmatha, C. Han, E. M. Riseman, and W. B. Croft, *Indexing handwriting using word matching*, in 1st ACM International Conference on Digital Libraries, Bethesda, MD, March 20-23, pp. 151-159 (1996).
2. S. Kane, A. Lehman, and E. Partridge, *Indexing George Washington's handwritten manuscripts*, Technical Report MM-34, Center for Intelligent Information Retrieval, University of Massachusetts Amherst (2001).
3. S. Belongie, J. Malik, and J. Puzicha, *Shape matching and object recognition using shape contexts*, IEEE Trans. on PAMI 24: 24 (2002) 509-522 (2002).
4. R. Manmatha and T. M. Rath, *Indexing of handwritten historical documents-recent progress*, in Sympo-

- sium on Document Image Understanding Technology (SDIUT), pp. 77-85 (2003).
5. J. L. Rothfeder, S. Feng, and T. M. Rath, *Using corner feature correspondences to rank word images by similarity*, CIIR Technical Report MM-44 (2003).
 6. B. Zhang, S. N. Srihari, and C. Huang, *Word image retrieval using binary features*, in Document Recognition and Retrieval XI, SPIE vol. 5296, pp. 45-53 (2004).
 7. B.S. Manjunath, W.Y. Ma, *Texture features for browsing and retrieval of image data*, IEEE Transactions on Pattern Analysis and Machine Intelligence, Volume: 18, Issue: 8, Pages: 837 - 842 (1996).
 8. Xuewen Wang, Xiaoqing Ding, Changsong Liu, *Optimized Gabor Filter Based Feature Extraction for Character Recognition*, ICPR (4) 2002: 223-226 (2002).
 9. Xuewen Wang, Xiaoqing Ding, Changsong Liu, *Gabor filters-based feature extraction for character recognition*, Pattern Recognition 38(3): 369-379 (2005).
 10. T.-F. Wu, C.-J. Lin, and R. C. Weng, *Probability estimates for multi-class classification by pairwise coupling*, Journal of Machine Learning Research, 5:975-1005 (2004).
URL: <http://www.csie.ntu.edu.tw/~cjlin/papers/svmprob/svmprob.pdf>.
 11. Chih-Chung Chang and Chih-Jen Lin, *LIBSVM: a library for support vector machines*, (2001).
Software available at
<http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
 12. G. Kim and V. Govindaraju, *A lexicon driven approach to handwritten word recognition for real-time applications*, IEEE Transactions on Pattern Analysis and Machine Intelligence, Volume: 19, Pages: 366 - 379 (1997).
 13. Robert Milewski, Venu Govindaraju, *Extraction of Handwritten Text from Carbon Copy Medical Form Images*, Document Analysis Systems pp. 106-116 (2006).