

A Probabilistic Method for Keyword Retrieval in Handwritten Document Images

Huaigu Cao, Anurag Bhardwaj and Venu Govindaraju

Center for Unified Biometrics and Sensors (CUBS)

Dept. of Computer Science and Engineering

University at Buffalo, Amherst, NY

{hcao3, ab94, govind}@cubs.buffalo.edu

Abstract

Keyword retrieval in handwritten document images (word spotting) is very challenging given that OCR accuracy is not yet adequate for handwritten scripts, specially with large lexicons. Various proposed approaches build indices on information such as image features or OCR scores and have improved the performance of the traditional approach that builds index on OCR'ed text. In this paper, we improve existing keyword retrieval (word spotting) method by modeling imperfect word segmentation as probabilities and integrating the word segmentation probabilities into the word spotting algorithm. The word recognition scores are also converted into probabilities that are compatible with the probabilistic word spotting model.

Keywords: Word Spotting, Information Retrieval, Handwriting Recognition.

1. Introduction

Automatic indexing and searching in handwritten documents for keywords is a very challenging task. The OCR technique is adequate only for applications of printed documents or hand-written documents of good image quality and relatively small lexicons. Thus, instead of indexing the OCR'ed text of handwriting for keyword retrieval, approaches of indexing of word image features, intermediate OCR likelihoods or matching distances have been proposed.

Keyword retrieval in document image is generally referred to as *word spotting*. The index for word spotting can be built from images features [9, 13] or OCR scores (likelihoods/matching distances) [4, 10, 5]. In image feature based indexing approaches, after preprocessing of document images and word segmentation, feature vectors are extracted from word images and stored in a database. When a user provides a query word, the similarity between the query and the word image in the database is

computed, and word images are returned in the decreasing order of similarities. Similarity [9, 7] is achieved by Dynamic Time Warping (DTW) matching of profile features computed using various definitions of distances [9, 12, 7] in the feature space. Similarity [13] is based on bitwise matching of the GSC features of two word images. Image feature based indexing approaches requires on-line matching which is time-consuming. Trade-off between accuracy and speed has to be made in order to be scalable to large database. Thus Image feature based indexing approaches are limited in feature selection and the complexity of matching and training methods and are only applicable to constrained handwriting such as when dealing with a single writer or small lexicons.

OCR score based indexing approaches [4, 10, 5] conquered the speed problem. In these approaches, the indices are built from OCR scores such posterior probabilities or feature vector observational likelihoods (probability density) converted from distances returned by word recognizer. [4] performs word recognition on character level and search for optimal matchings in the series of character recognition scores. However searching in character recognition scores requires additional time and is time-consuming for large-scale data. The index for fast retrieval can be built on the results of word level recognition in lexicon-driven mode [10, 5]. In the lexicon driven mode, any word that is not in the lexicon can't be retrieved. Thus one need to select the set of keywords to meet the requirement of the application.

We improve the OCR score based indexing method by integrating word segmentation probabilities into the ranking scheme of word images. Word spotting methods mentioned above assume perfect word segmentation: word images are given by word segmentation algorithm, and the ranks of word images are obtained by ordering the word recognition scores. However it is difficult to get a nearly perfect word segmentation in unconstrained handwriting due to irregular variation of the word gaps. Thus, the performance of word spotting can be improved by mod-

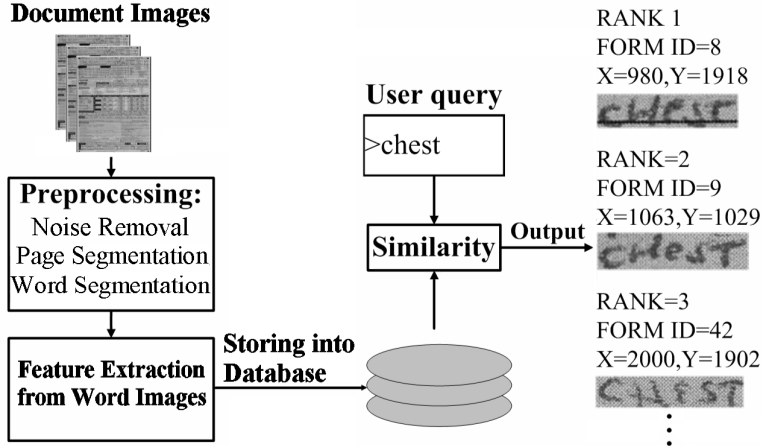


Figure 1. Diagram of the keyword spotting system.

eling word segmentation probabilities. In this paper, we describe a probabilistic model of word spotting that integrates word segmentation probabilities and word recognition probabilities. The word segmentation probabilities are obtained by modeling the conditional distribution of multivariate distance features of word gaps. The word recognition results are also represented by a probabilistic model so that they are compatible with the probabilistic word spotting model. The modeling of the word recognition probabilities is obtained from the distances returned by the word recognizer.

2 Model Definition

2.1 Word Spotting Model

Given a series of consecutive connected components c_1, c_2, \dots, c_n and a possible word image w represented by c_i, c_{i+1}, \dots, c_j ($1 \leq i, j \leq n$), then the similarity between w and the query word q can be represented by

$$\text{sim}(w, q) = \sigma_{i-1} \cdot (1 - \sigma_i) \cdot \dots \cdot (1 - \sigma_{j-1}) \cdot \sigma_j \cdot \text{Pr}(q|w) \quad (1)$$

where σ_k ($1 \leq k \leq n - 1$) is the probability that the gap between c_{k-1} and c_k is a valid word gap and $\text{Pr}(q|w)$ is the word recognition posterior probability. For the convenience of computation, we also treat the preceding blank of the left-most connected component and the succeeding blank of the right-most connected component as two “special” gaps (σ_0 and σ_n). The probability of validity of these two gaps is always 1, *i.e.*, $\sigma_0 = \sigma_n = 1$. Here we assume the gaps are independent and thus the word segmentation probability is $\sigma_{i-1} \cdot (1 - \sigma_i) \cdot \dots \cdot (1 - \sigma_{j-1}) \cdot \sigma_j$.

The size required to store the index can be reduced by applying constraints on the number of connected components within a word image and minimum similarity:

Given a series of consecutive connected components c_1, c_2, \dots, c_n ,

For i from 1 to n

For j from i to $\min(i + C_{max} - 1, n)$

If the similarity $\text{sim}(c_i \dots c_j, q) > \text{sim}_{min}$

Then store the document number, coordinates of the word image and the similarity into index.

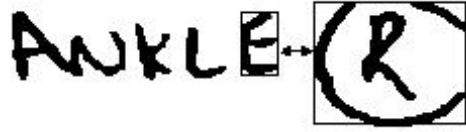
C_{max} is the maximum number of connected components within a word image. sim_{min} is the minimum similarity that can be stored into the index. We assume $C_{max} = 16$ and $\text{sim}_{min} = 0.1\%$ in our experiment.

2.2 Estimating Word Segmentation Probability

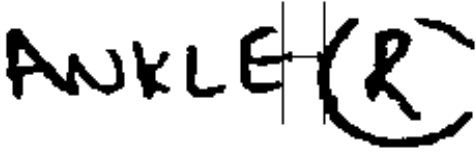
Word segmentation is defined as the process of segmenting a line into words. In handwritten lines, the space between words is uneven. Moreover, the same amount of space may be present between words, and between characters within a word. Such cases arise due to differences in writing styles, and space constraints.

In our word segmentation method [2], the probability of the gap between any two consecutive connected components being a valid word gap is estimated from feature space of distance measures. A gap between two connected components is represented by three distance features:

1. **Euclidean Distance.** This feature is defined as the horizontal distance between the bounding boxes of the two consecutive connected components of the line image (Figure 2(a)).
2. **Minimum Run Length.** This feature represents the minimum horizontal white run length distance between the two adjacent connected components of the line image.
3. **Convex Hull Distance.** We compute the convex hulls of two consecutive connected components and



(a) Euclidean distance.



(b) Run length distance.



(c) Convex hull distance.

Figure 2. Three feature representing a gap between two consecutive connected components.

draw a line connecting the mass centers of the two convex hulls. The Euclidean distance between points at which this line crosses the two convex hulls is defined as the Convex Hull distance of the two adjacent components.

To eliminate the effect of different text sizes, we compute the average height of all the components and normalize the extracted features by dividing them by the average height of all components in the same line.

The segmentation probability of a gap g is given by the Bayes' Rule

$$\sigma_g = \Pr(g|f_{1,g}, f_{2,g}, f_{3,g}) = \frac{\Pr(g)p(f_{1,g}, f_{2,g}, f_{3,g}|g)}{\Pr(g)p(f_{1,g}, f_{2,g}, f_{3,g}|g) + \Pr(\bar{g})p(f_{1,g}, f_{2,g}, f_{3,g}|\bar{g})} \quad (2)$$

where $\Pr(g)$ and $\Pr(\bar{g})$ are the prior probabilities of valid gaps and non-valid gaps, respectively. $f_{1,g}$, $f_{2,g}$ and $f_{3,g}$ are three features of g . $p(f_{1,g}, f_{2,g}, f_{3,g}|g)$ is the probability density of the features of valid gaps. $p(f_{1,g}, f_{2,g}, f_{3,g}|\bar{g})$ is the probability density of the features of non-valid gaps.

Given a set of gap features with the annotation of "valid" and "non-valid", we can estimate $\Pr(g)$, $\Pr(\bar{g})$, $p(f_{1,g}, f_{2,g}, f_{3,g}|g)$ and $p(f_{1,g}, f_{2,g}, f_{3,g}|\bar{g})$ as follows. $\Pr(g)$ and $\Pr(\bar{g})$ are estimated from the ratio of the numbers of valid and non-valid gaps in the training set.

$$\Pr(g) = \frac{\#\{\text{valid gaps}\}}{\#\{\text{valid gaps}\} + \#\{\text{non-valid gaps}\}} \quad (3)$$

$$\Pr(\bar{g}) = 1 - \Pr(g) \quad (4)$$

$p(f_{1,g}, f_{2,g}, f_{3,g}|g)$ and $p(f_{1,g}, f_{2,g}, f_{3,g}|\bar{g})$ are estimated non-parametrically using Parzen window technique with a Gaussian kernel function.

2.3 Estimating Word Recognition Probability

We use a lexicon-driven word recognition algorithm [6] that performs character segmentation and finds the best matching path using dynamic programming. First a word image is segmented into candidate character images. Then the directional features are extracted from the contours of character images and matched to every word in the lexicon by searching all possible segmentations for the minimum sum of Euclidean distances from the features of the test image and the character templates in the training set. The minimum Euclidean distance indicates the similarity between the word image and the term in the lexicon. The square of the distance associated with a pair of a word image w and a term t_i is denoted by $s(w, t_i)$.

The word recognition probability is estimated from the recognition score using a Universal Background Model (UBM) [11, 2]. In a Background Model, the posterior probability of the word recognition is given by Bayes' rule:

$$\Pr(w = t_i | s(w, t_i)) = \frac{\Pr(w = t_i)p_{t_i}(s(w, t_i)|w = t_i)}{\Pr(w = t_i)p_{t_i}(s(w, t_i)|w = t_i) + \Pr(w \neq t_i)p_{t_i}(s(w, t_i)|w \neq t_i)} \quad (5)$$

where $p_{t_i}(s(w, t_i)|w = t_i)$ is the likelihood of the correct matching score when the word is t_i , $p_{t_i}(s(w, t_i)|w \neq t_i)$ is the likelihood of the incorrect matching score when the word is t_i , and $\Pr(w = t_i)$, $\Pr(w \neq t_i)$ are the prior probabilities of correct and incorrect matchings of t_i , respectively.

We need a term specific training set for every term to learn the background model. This is a drawback in applications using large number of terms. The Universal Background Model is an alternative approach that solves this problem. In the UBM, we use a single Background Model for all of the terms. The correct matching probability is given by

$$\Pr(\text{Correct}|s) = \frac{\Pr(\text{Correct})p(s|\text{Correct})}{\Pr(\text{Correct})p(s|\text{Correct}) + \Pr(\text{Incorrect})p(s|\text{Incorrect})} \quad (6)$$

where s is a matching score, and $\Pr(\text{Correct})$ and $\Pr(\text{Incorrect})$ are the prior probabilities of correct and incorrect matchings, respectively, and $p(s|\text{Correct})$, $p(s|\text{Incorrect})$ are the likelihoods of the score of correct and incorrect matchings, respectively. $\Pr(\text{Correct})$, $\Pr(\text{Incorrect})$, $p(s|\text{Correct})$, and $p(s|\text{Incorrect})$ are estimated from the scores of all of the terms. We model

$p(s|Correct)$ and $p(s|Incorrect)$ as Gamma distributions.

3 Experimental Results

3.1 Data Collection

Our keyword retrieval algorithm has been tested on the New York State Pre-hospital Care Reports (PCR forms). The task is quite challenging for several reasons: (i) handwritten responses were very loosely constrained in terms of writing style, format of response, and choice of text due to irrepressible emergency situations; (ii) images are scanned from carbon copies and are very noisy (Figure 3). This leads to difficulties in the automatic transcription of forms. The word recognition rate of the forms using word recognizer [6] is about 20-30%.

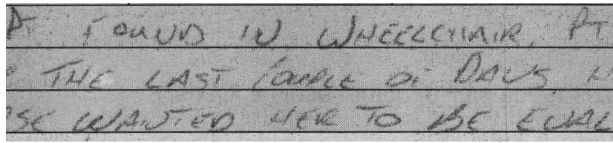


Figure 3. The text in a PCR form.

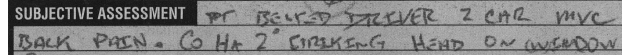
3.2 Preprocessing

First we detect and remove the skew of every PCR form image as follows.

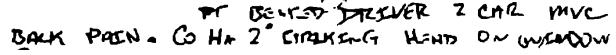
1. We manually de-skew a form and take it as a template. Two special regions are taken from the template as anchors.
2. The positions of two anchoring regions in any test image are found by cross-correlation.
3. The skew angle of the test image is obtained by the relative skewing between the test image and the template. We de-skew the image by rotating to the opposite direction.

By aligning the test image to the template image, we can also obtain the position of each form cell containing a line of text. The de-skewing and page segmentation method using template-matching works well on the PCR form images since they have a fixed layout and are scanned at the same resolution. Our approach is applicable to other types of forms as well.

We use the MRF based document image preprocessing algorithm [3] to binarize the form image and remove the grid lines from the image. Assuming the binarized objective image is x and the grayscale image is y , we solve the maximum a posteriori (MAP) estimation $\hat{x} = \underset{x}{\operatorname{argmax}} \operatorname{Pr}(x|y)$ using the Markov Random Fields



(a) The original grayscale image.



(b) The binarized image. Grid lines are removed and broken strokes are fixed.

Figure 4. An example of the binarization and line removal result.

(MRF). An example of binarization and line removal result is shown in Figure 4. The MRF based preprocessing method improves the word recognition accuracy from 18.7% (obtained by the PCR form preprocessing algorithm in [8]) to 28.6%.

3.3 Evaluation Metrics

The performance of word spotting can be evaluated using the precisions at standard recall levels (0, 0.1, ..., 1). We may also use single value measures such as the Mean Average Precision (MAP) [1] to evaluate the word spotting performance. The Mean Average Precision is computed as follows:

1. For each query, check the returned word images starting from rank 1. Whenever a relevant word image is found, record the precision of the word images from the one with rank 1 to the current one. The Average Precision (AP) of a given query q is weighted sum of the recorded precisions:

$$AP(q) = \frac{\sum_{1 \leq r \leq N_q, Rel(r) \text{ is true}} Prec(r)}{R_q} \quad (7)$$

where N_q is the number of word images returned, R_q is the number of relevant documents, $Prec(r)$ is the precision of top- r returned word images, and $Rel(r)$ is a Boolean function of the relevance of rank r .

2. The Mean Average Precision (MAP) of all the queries is:

$$MAP = \frac{\sum_q R_q \times AP(q)}{\sum_q R_q} \quad (8)$$

3.4 Keyword Retrieval Results

We performed two word spotting tests to show the improvement due to the use of word segmentation probability.

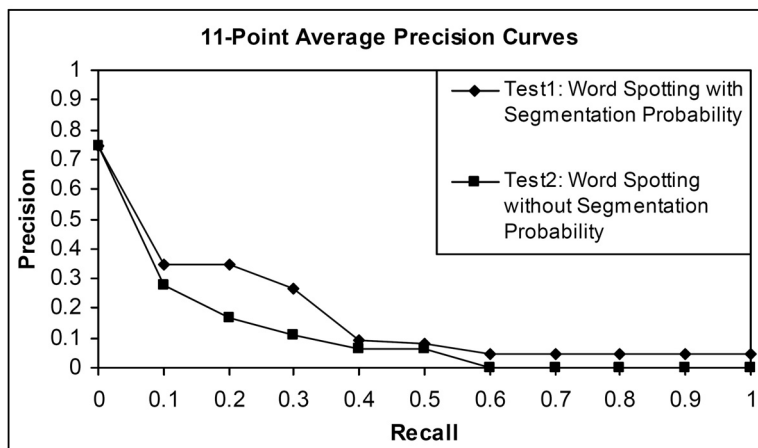


Figure 5. 11-point average precision curves of Tests 1-2.

Test 1: Word spotting with segmentation probability

We searched 342 PCR forms for 33 keywords using the similarity function in Equation (1) and the estimation methods described in Section 2.

Test 2: Word spotting without segmentation probability

In this test we evaluate the performance of word spotting using connected component clustering based word segmentation method without probability annotation. This test is based on the same idea used by [5] but we implemented different word segmentation and recognition method. We performed word spotting test on the same PCR forms and keywords that we performed on in Test 1. The word images in test 2 were obtained by grouping adjacent connected components with gap segmentation probability $\sigma_k < C_{gap}$. The cutoff threshold $C_{gap} = 0.297$ in our test. The recall rate of the gap classification reached maximum (0.394) when $C_{gap} = 0.297$. We used the same word recognition method that we used in Test 1. The word recognition probability returned by the Universal Background Model (UBM) is taken as the similarity between the word image and the query.

The 11-point average precision curves of Tests 1-2 are shown in Figure. The Mean Average Precision (MAP) scores of Test 1 and Test 2 are 4.7% and 2.8%, respectively. The test results show the improvement obtained by using word segmentation probabilities.

4 Conclusion

This paper describes a method to search handwritten document images for keywords. The image/text similarity of the proposed method is defined as the product of word segmentation probability and word recognition probabil-

ity. Test results show the improvement of integrating the probabilistic annotation of word segmentation on handwritten document images with word segmentation errors.

References

- [1] R. A. Baeza-Yates and B. A. Ribeiro-Neto. *Modern Information Retrieval*. ACM Press / Addison-Wesley, 1999.
- [2] H. Cao, A. Bhardwaj, and V. Govindaraju. Handwritten document retrieval using mmse estimation of raw term frequencies. *submitted to IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [3] H. Cao and V. Govindaraju. Handwritten carbon form pre-processing based on markov random field. In *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007.
- [4] H. Cao and V. Govindaraju. Template-free word spotting in low-quality manuscripts. In *the Sixth International Conference on Advances in Pattern Recognition (ICAPR)*, volume 5296, pages 45–53, Calcutta, India, 2007.
- [5] N. R. Howe, T. M. Rath, and R. Manmatha. Boosted decision trees for word recognition in handwritten document retrievals. In *Proceedings of the SIGIR*, pages 377–383, 2005.
- [6] 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*, 19:366–379, April 1997.
- [7] R. Manmatha and T. M. Rath. Indexing of handwritten historical documents-recent progress. In *Symposium on Document Image Understanding Technology (SDIUT)*, pages 77–85, 2003.
- [8] R. Milewski and V. Govindaraju. Extraction of handwritten text from carbon copy medical form images. In *Document Analysis Systems*, pages 106–116, 2006.
- [9] T. M. Rath and R. Manmatha. Word image matching using dynamic time warping. In *Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 521–527, Center for Intelligent Information Retrieval, University of Massachusetts Amherst, 2003.
- [10] T. M. Rath, R. Manmatha, and V. Lavrenko. A search engine for historical manuscript images. In *Proceedings of the 27th annual international ACM SIGIR conference on research and development in information retrieval*, 2004.
- [11] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn. Speaker verification using adapted gaussian mixture models. *Digital Signal Processing*, 10(1-3):19–41, 2000.

- [12] J. L. Rothfeder, S. Feng, and T. M. Rath. Using corner feature correspondences to rank word images by similarity. In *CIIR Technical Report MM-44*, Center for Intelligent Information Retrieval, University of Massachusetts Amherst, 2003.
- [13] B. Zhang, S. N. Srihari, and C. Huang. Word image retrieval using binary features. In *Document Recognition and Retrieval XI, SPIE*, volume 5296, pages 45–53, Greenbelt, Maryland, April 2004.