

Indexing and Retrieval of Degraded Handwritten Medical Forms

Huaigu Cao, Faisal Farooq and Venu Govindaraju

Center for Unified Biometrics and Sensors (CUBS)

Dept. of Computer Science and Engineering

University at Buffalo, Amherst, NY

{hcao3,ffarooq2,govind}@cubs.buffalo.edu

Abstract

The tasks of indexing and retrieval are specifically challenging for the erroneous output of handwriting recognition (*HR*) systems. This paper proposes an approach of indexing and retrieving degraded documents with very low recognition rates. We present a modified version of the popular Vector Model in information retrieval (*IR*). Our model incorporates *top n* candidates from a HR system into the scheme of calculating the term frequency (*tf*) and the inverted document frequency (*idf*). Standardized IR Tests show that the proposed approach outperforms the retrieval of ordinary HR text in terms of mean average precision (MAP) and R-Precision.

1 Introduction

Decades of the development in optical character recognition (OCR) techniques has made it possible to convert great volumes of documents into digital form such as text or structured forms of PDF and XML. In recent years, the achievement in information retrieval has provided a powerful tool for indexing and searching large scale on-line database of documents. Although OCR has been successful in applications of machine-printed document recognition and handwriting recognition (HR) with small lexicon, unconstrained handwriting with large lexicon in general still has very high error rate. Researchers [Croft *et al.*, 1994; Beitzel *et al.*, 2003] have shown that the performance of OCR text retrieval is badly affected when dealing with short or low quality documents. Our work is motivated by the objective of building a system for recognizing and retrieving New York State Pre-hospital Care Reports (PCR forms). In New York State all patients who enter the Emergency Medical System (EMS) are tracked through their pre-hospital care to the emergency room using the PCR. The PCR is used to gather vital patient information and has manifold purpose: (i) it is used by the emergency room staff in providing appropriate care, (ii) it is to facilitate record keeping about the patient and the care provided, and (iii) to afford large regional EMS (such as WREMS, Western Regional Emergency Medical Systems.) to gain

trend knowledge through macro analysis. Our goal is to automate the collection of data from the PCR and enable efficient maintenance, and dissemination of information. The task proposed is quite challenging for several reasons: (i) handwritten responses are very loosely constrained in terms of writing style, format of response, and choice of text, (ii) medical lexicons of words are very large (about 50,000 entries), (iii) forms are often very noisy due to irrepressible emergency situations. This leads to the automatic transcription of the forms very difficult. Our dataset consists of over 6000 scanned PCR forms containing medicine and healthcare information. Each PCR contains about 100 handwritten words on average so the content is very short. The word recognition rate of the forms using Word Model Recognizer (WMR) [Kim and Govindaraju, 1997] is as low as around 20% as reported in [Milewski and Govindaraju, 2006].

The problem of indexing and retrieving documents with low recognition rate has recently been addressed by researchers. In [Mittendorf and Schauble, 1996; Ohta *et al.*, 1997; Jing, 2002] different approaches modeling typical recognition errors were proposed. In [Mittendorf and Schauble, 1996] a probabilistic model for common recognition errors was proposed and this model was used to design the term-weighting scheme of information retrieval. The approach that generates candidate terms for each “true” search term and adds the retrieval results of candidate terms into the final result was studied in [Ohta *et al.*, 1997]. In [Jing, 2002], a language model that took common recognition errors into account was built. This language model can then be used to approximate an “uncorrupted” version of a particular document, and it can be used for retrieval in a language modeling approach. Due to low recognition accuracies, the task of retrieval and indexing has been elusive. Keyword spotting as an alternative approach of indexing and retrieving handwritten documents has been proposed in [Manmatha and Croft, 1997; Rath and Manmatha, 2006; Zhang *et al.*, 2004]. The idea is to search the document for a certain keyword by feature matching instead of recognition. Rath and Manmatha developed a word matching algorithm that compares word images using Dynamic Time Warping (DTW) [Rath *et al.*, 2004]. They address a “single-

author” problem by matching word images with each other to create equivalence classes. Each equivalence class consists of multiple instances of the same word. We propose an approach based on the idea of building index on the text of *top n* ($n > 1$) word recognition candidates rather than on the ordinary HR text containing the top choice candidate only. For retrieval, a variant of the popular Vector Model with a scheme of using a global probabilistic model of word recognition ranks to calculate TF and TDF is introduced. IR tests show that the proposed approach outperforms the standard Vector Model on OCR text with a gain of the Mean Average Precision (MAP) from 10.6% to 18.1% and the R-Precision from 10.0% to 22.1%. The following sections are organized as follows. The Top-n-Candidate Vector Model is introduced in section 2. Results of IR tests are presented and discussed in section 3. Finally, section 4 concludes the paper.

2 Multiple-candidate Vector Model

2.1 Benefit from multiple-candidates

The existence of errors in HR text decreases the performance of IR in the sense that a query word present in a certain position of the document, due to incorrect recognition, might not be present in the recognized text, but the count of occurrences of the query word is essential to the IR algorithm. On the one hand, as results in table 1 suggest, the use of *top n* ($n > 1$) candidates in terms of word recognition ranks will improve the count of occurrences.

Table 1: Recognition rates of PCR forms for different n .

| n | 1 | 2 | 5 | 10 |
|---------------------|------|------|------|------|
| Recognition rate(%) | 17.0 | 24.2 | 35.6 | 43.2 |

This is because the percentage chance of the correct choice being in the *top n* ($n > 1$) is significantly greater. On the other hand, use of too many candidates or many candidates with equal weights introduces more false occurrences of terms, thus effecting precision adversely.

2.2 A Vector Model for multiple word candidates

Suppose **WR** is a word recognizer, **I** is a collection of document images, and **O** is the output of **WR**. In order to retrieve multiple-candidate recognized text using the Vector Model, first we need to define the formula for calculating the term frequency for **WR** or for general purpose any OCR system. We will denote it by **OCR-tf**. Here we assume unigrams and ignore the effects of stemming. Then the **OCR-tf** of a given term (word) t_i in a document d_j is defined as the expectation of the *tf* in the ground truth of document d_j , denoted by

$$\begin{aligned} \mathbf{tf}_{i,j}^{\text{ocr}} &= E\{tf_{i,j}|\mathbf{O}\} \\ &= E\left\{\frac{freq_{i,j}}{\sum_l freq_{l,j}}|\mathbf{O}\right\} \\ &= \frac{E\{freq_{i,j}|\mathbf{O}\}}{\sum_l E\{freq_{l,j}|\mathbf{O}\}} \end{aligned} \quad (1)$$

where $freq_{i,j}$ is the raw frequency of term t_i in document d_j (i.e., the number of times term t_i is mentioned in the text of document d_j) and $\sum_l freq_{l,j}$ is a normalization factor that stands for the total number of terms in document j . In the second and third lines of Equation 1, $\sum_l freq_{l,j}$ is assumed to be independent of $tf_{i,j}$ so

$$\begin{aligned} E\{freq_{i,j}\} &= E\left\{tf_{i,j} \cdot \sum_l freq_{l,j}|\mathbf{O}\right\} \\ &= E\{tf_{i,j}|\mathbf{O}\} \cdot E\left\{\sum_l freq_{l,j}|\mathbf{O}\right\} \\ &= E\{tf_{i,j}|\mathbf{O}\} \cdot \sum_l E\{freq_{l,j}|\mathbf{O}\}. \end{aligned}$$

Denote the set of words (or terms, for simplicity) in a document d by $w(d)$, then

$$E\{freq_{i,j}|\mathbf{O}\} = \sum_{w \in w(d_j)} P(w = t_i) \quad (2)$$

Therefore,

$$\mathbf{tf}_{i,j}^{\text{ocr}} = E\{tf_{i,j}|\mathbf{O}\} = \frac{\sum_{\mathbf{w} \in \mathbf{w}(d_j)} \mathbf{P}(\mathbf{w} = \mathbf{t}_i)}{\sum_1 \sum_{\mathbf{w} \in \mathbf{w}(d_j)} \mathbf{P}(\mathbf{w} = \mathbf{t}_1)} \quad (3)$$

Equation 3 is an analytical expression of $tf_{i,j}^{\text{ocr}}$. Given probabilities $P(w = t_i) \forall t_i$, we will be able to calculate the **OCR-tf**.

However, most word recognizers do not (or are not necessarily to) produce the posterior probability for each class. The only information readily available from recognition output is the rank of each candidate word. The rank of a candidate is correlated with the posterior probability so we can find some way of probability estimation when we know only the rank. A technique to map the ranks to a normalized probability distribution by fitting to a Zipfian distribution was presented in [Howe *et al.*, 2005]. We present a novel method based on the following motivation. Given an image η of word w , the recognizer lexicon L , and a sequence of entries in $L : w_1, w_2, \dots, w_n$ representing the recognition result,

where $L = \{w_1, w_2, \dots, w_n\}$ and the rank of term w_i w.r.t. word image η , $rank_\eta(w_i) = i (1 \leq i \leq n)$, we want to calculate $\Pr(w = w_i | w \in L)$, denoted by $P^{(i)}$, for each i . To simplify the analysis, assuming that for any large lexicon L , the top-1 probability is intrinsic of the recognizer and therefore is a constant p , define a sequence of nested lexicons

$$\begin{cases} L^{(1)} &= L \\ L^{(k)} &= L^{(k-1)} - \{w_{k-1}\}, 1 < k \leq n \end{cases}$$

then we have

$$\begin{aligned} P^{(1)} &= p \\ P^{(2)} &= P(w \neq w_1 | w \in L) \times P(w = w_2 | w \in L^{(2)}) \\ &= (1 - p)p \\ P^{(3)} &= P(w \neq w_1 \wedge w \neq w_2 | w \in L) \\ &\quad \times P(w = w_3 | w \in L^{(3)}) \\ &= (1 - p)^2 p \\ &\dots \end{aligned}$$

and finally

$$P^{(i)} = (1 - p)^{i-1} p, \quad i \ll n$$

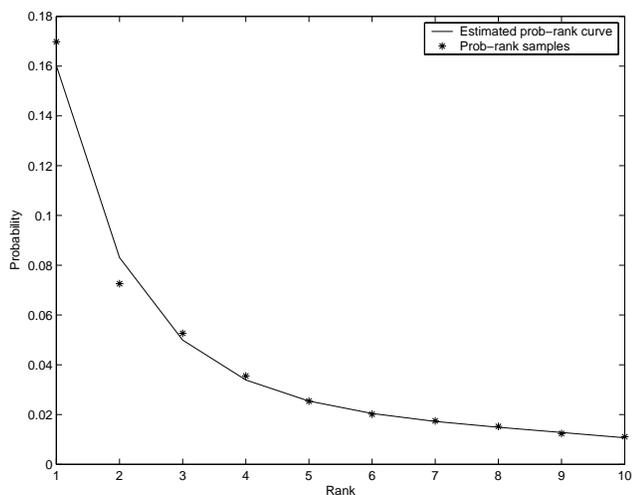


Figure 1: Top-10 recognition rates and the estimated posterior probability curve $P(w = w_i | w \in L) = e^{\theta(i)}$.

In other words, $P(w = w_i | w \in L)$ is an exponential function of i in this simplified scenario. To be more general, suppose $P^{(i)} = e^{\theta(i)}$ where $\theta(i)$ is a, possibly non-linear, function of i . In our test, $\theta(i)$ is assumed to be a third degree polynomial of i and is estimated by fitting a polynomial to logarithms of normalized word recognition rate samples of each rank i . The word recognition rate sample at rank i is calculated by counting the total number of true matches of rank i candidates and then divided by total number of words recognized during the test. Figure 1 shows the word

recognition rate samples of ranks 1 to 10 calculated from the recognition results of PCR forms and the estimated curve $P^{(i)} = e^{\theta(i)}$.

While calculating OCR-TF using equation 3 with estimated probabilities, any probability that is below a certain threshold is omitted in order to suppress the number of false occurrences of terms. In our tests any $P^{(i)}$ smaller than $c_t \cdot P(1)$ is considered to be zero and hence omitted, where c_t is a truncation parameter and takes the value of 0.1. When calculating the IDF of a term, only those documents including the term recognized with $P^{(i)} \geq c_t \cdot P(1)$ are taken into account. In other words, the IDF of a given term t_i , $\text{idf}_i^{\text{ocr}}$, is calculated by

$$\text{idf}_i^{\text{ocr}} = \log \frac{|\{d_j\}|}{|\{d_j | \exists \eta \in d_j \text{ s.t. } P(\text{rank}_\eta(t_i)) \geq c_t \cdot P(1)\}|}$$

where $|\{d_j\}|$ denotes the total number of documents, η denotes a word image, and $\text{rank}_\eta(t_i)$ denotes the word recognition rank of term t_i w.r.t. word image η .

3 Test results and discussions

| | TIME | RESP | PULSE | B.P. | GLUCOSE | GCS | PUPILS | L | R | SKIN | STATUS |
|---|-------|----------|----------|--------|---------|-----|-------------|---|---|--------------|--------|
| V | 18:47 | 18 | 70 | 120/80 | 150 | 15 | 3 | 3 | 3 | Warm | 0 |
| I | | Regular | Regular | | | | Constricted | | | Wet | 0 |
| T | | Shallow | Shallow | | | | Sluggish | | | Mostly | 0 |
| A | | Labored | Labored | | | | No-Reflexes | | | Dry | 0 |
| L | | Rate: 18 | Rate: 70 | | | | Normal | | | Gold | 0 |
| S | | Regular | Regular | | | | Dieter | | | Warm | 0 |
| I | | Shallow | Shallow | | | | Constricted | | | Wet | 0 |
| N | | Labored | Labored | | | | Sluggish | | | Mostly | 0 |
| G | | Rate: 18 | Rate: 70 | | | | Normal | | | Unremarkable | 0 |
| N | | Regular | Regular | | | | Dieter | | | Warm | 0 |
| S | | Shallow | Shallow | | | | Constricted | | | Wet | 0 |
| | | Labored | Labored | | | | Sluggish | | | Mostly | 0 |
| | | Rate: 18 | Rate: 70 | | | | Normal | | | Unremarkable | 0 |
| | | Regular | Regular | | | | Dieter | | | Warm | 0 |
| | | Shallow | Shallow | | | | Constricted | | | Wet | 0 |
| | | Labored | Labored | | | | Sluggish | | | Mostly | 0 |

Figure 2: An example of NYS Pre-hospital Care Report form.

Our corpus contains 342 PCR forms with ground truth and the coordinates of each word. The test bed involves 21 queries, and the annotation of relevance of the 342 forms to these queries. An example of the handwritten region in a PCR form is shown in figure 2. Of all the 21 queries, thirteen contain 1 term each, six contain 2 terms each, and one contains 3 terms. The total number of relevant documents (with duplicates) is 379.

All the forms are first binarized and lines removed using the algorithm proposed in [Milewski and Govindaraju, 2006]. Then all the word images are extracted and fed into the word recognizer of [Kim and Govindaraju, 1997] with a lexicon of 4405 English words. The top 1, 2, 5, and 10 word recognition rates are shown in table 1. IR tests of all the 21 queries are performed on both ordinary recognition text (with top one output) and the recognition output of multiple candidates. The Vector Model

with standard *tf-idf* scheme is applied to ordinary recognition text, and the Vector Model with modified TF-IDF scheme is applied to multiple-candidate recognition text. All the IR tests are implemented on the basis of Lemur Toolkit (<http://www.lemurproject.org>), a benchmark IR system designed by Carnegie Mellon University and the University of Massachusetts, Amherst. In order to get better performance, a pseudo-feedback using *top 3* documents and maximum of 10 terms is applied. For the purpose of comparison, the following naive scheme of *tf-idf* is also tested

$$P^{(i)} = \begin{cases} \frac{1}{S}, & 1 \leq i \leq S \\ 0, & \text{otherwise} \end{cases}$$

where S is a truncation threshold.

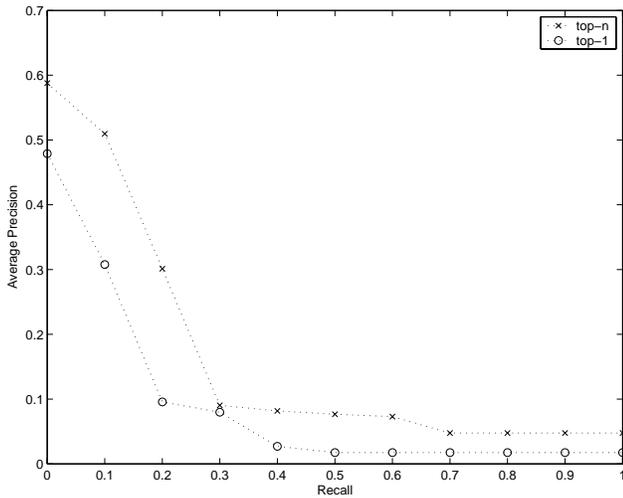


Figure 3: Performances without pseudo relevance feedback.

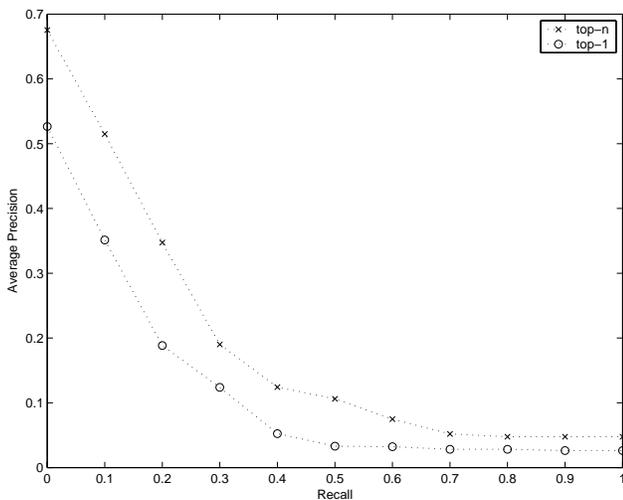


Figure 4: Performances with pseudo relevance feedback.

The results of interpolated average precisions at 11 different recall values are compared in figures 3 and 4. The Mean Average Precisions (MAP) and R-Precisions of the results of ordinary recognition text, multi-candidate text with naive *tf-idf* scheme and the proposed TF-IDF scheme are compared in table 2. The performances shown in figures 3 and 4 and table 2 indicate that the proposed method outperforms the Vector Models using both standard and naive *tf-idf* schemes. A gain of 7.5% of MAP and 12.1% of R-Precision is obtained as compared to the standard scheme for ordinary recognition text. We also learn from table 2 that the scheme of exponential probability is better than the scheme of assigning equal weight to each candidate.

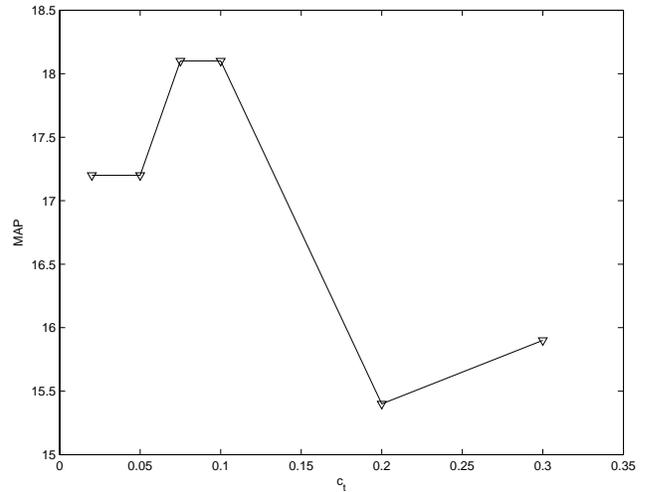


Figure 5: MAP values using different values of $c_t \in \{0.02, 0.05, 0.075, 0.1, 0.2, 0.3\}$ are tested.

The effect of choosing different values of c_t is tested and the result is shown in figure 5. We can see that the best performance is obtained when $c_t \in \{0.1, 0.075\}$. The c_t value as large as 0.2 or 0.3 does not involve enough top word candidates, whereas too small c_t value seems to involve too many irrelevant candidates which bring about poor performance.

4 Conclusion

This paper proposed a modified IR Vector Model for degraded documents with low recognition rate that incorporates the ranks of top word recognition candidates into the scheme of calculating *tf-idf*. The proposed approach improved the performance of IR as compared to the method using standard Vector Model. Future work will focus on two respects. Firstly, define and test a more general model that takes segmentation errors into account. Secondly, use language model such as n-gram. This may encounter the problem of learning the language model from multiple candidates of word recognition. Besides, we are still investigating methods of improving recognition performance on PCR forms.

Table 2: MAP values of retrievals using different *tf-idf* schemes (with feedback and $c_t=0.1$).

| TF-IDF scheme | Standard | S=3 | S=5 | S=10 | Proposed |
|----------------|----------|------|------|------|----------|
| MAP(%) | 10.6 | 16.3 | 16.1 | 15.5 | 18.1 |
| R-Precision(%) | 10.0 | 20.3 | 21.5 | 21.0 | 22.1 |

References

- [Beitzel *et al.*, 2003] Steven M. Beitzel, Eric C. Jensen, and David A. Grossman. A survey of retrieval strategies for ocr text collections. In *Proceedings of the Symposium on Document Image Understanding Technologies*, Greenbelt, Maryland, April 2003.
- [Croft *et al.*, 1994] W. B. Croft, S. M. Harding, K. Taghva, and J. Borsack. An evaluation of information retrieval accuracy with simulated ocr output. In *Proceedings of the Symposium on Document Analysis and Information Retrieval*, 1994.
- [Howe *et al.*, 2005] Nicholas R. Howe, Toni M. Rath, and R. Manmatha. Boosted decision trees for word recognition in handwritten document retrievals. In *Proceedings of the SIGIR*, pages 377–383, 2005.
- [Jing, 2002] Hongyan Jing. Using hidden markov modeling to decompose human-written summaries. *Computational Linguistics*, 28(4):527–543, 2002.
- [Kim and Govindaraju, 1997] 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.
- [Manmatha and Croft, 1997] R. Manmatha and W.B. Croft. Word spotting: Indexing handwritten manuscripts. In *Intelligent Multi-media Information Retrieval Collection*. AAAI press, May 1997.
- [Milewski and Govindaraju, 2006] Robert Milewski and Venu Govindaraju. Extraction of handwritten text from carbon copy medical form images. In *Document Analysis Systems*, pages 106–116, 2006.
- [Mittendorf and Schauble, 1996] E. Mittendorf and P. Schauble. Measuring the effects of data corruption on information retrieval. In *Proceedings of the Fifth Annual Symposium on Document Analysis and Information Retrieval*, 1996.
- [Ohta *et al.*, 1997] M. Ohta, A. Takasu, and J. Adachi. Retrieval methods for english text with misrecognized ocr characters. In *Proceedings of the International Conference on Document Analysis and Recognition*, 1997.
- [Rath and Manmatha, 2006] T. M. Rath and R. Manmatha. Word spotting for historical documents. *International Journal on Document Analysis and Recognition*, 2006.
- [Rath *et al.*, 2004] Toni M. Rath, R. Manmatha, and Victor 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.
- [Zhang *et al.*, 2004] B. Zhang, S. N. Srihari, and C. Huang. Word image retrieval using binary features. In *Document Recognition and Retrieval XI, SPIE*, pages 45–53. SPIE, 2004.