

**LANGUAGE MODELS AND AUTOMATIC TOPIC  
CATEGORIZATION FOR INFORMATION RETRIEVAL  
IN HANDWRITTEN DOCUMENTS.**

By

**Faisal Farooq**

May 2008

A DISSERTATION SUBMITTED TO THE  
FACULTY OF THE GRADUATE SCHOOL OF STATE  
UNIVERSITY OF NEW YORK AT BUFFALO  
IN PARTIAL FULFILLMENT OF THE REQUIREMENTS  
FOR THE DEGREE OF  
DOCTOR OF PHILOSOPHY

© Copyright 2008

by

Faisal Farooq

*To my grandfather Prof. Mohammad Sidiq*

# Acknowledgments

I thank the Almighty for providing me with this opportunity to serve Him and make a contribution through His infinite wisdom. I thank my parents for their perseverance and unconditional support, without which I could never have accomplished this endeavor. I would also like to thank other members of my family including my cousin Muneer who has been watching my back from day one.

I want to extend my deep appreciation to Dr. Venu Govindaraju, the chair of my dissertation committee. He has been an advisor and a mentor. His persistent guidance, omnipresent motivation and overall support have been the foundation of this thesis. He introduced me to the area of handwriting recognition and encouraged me to address the open challenge of retrieval from handwritten documents.

I want to show my gratitude to Dr. Peter Scott, member of my dissertation committee. His course Computer Vision and Image Processing indeed laid a solid foundation for

this research. His guidance and advise has been always helpful. In addition, I had the opportunity to be his Teaching Assistant for three semesters and his passion for teaching was a great motivation.

Special thanks go to Dr. Hung Ngo who served as a member of my committee and provided full support and assistance. I would also like to thank Dr. Matthew Beal whose extraordinary guidance in Machine Learning enabled me to apply and modify the state of the art methods for the purpose of my dissertation. I thank Dr. Salim Roukos, manager of the multilingual natural language processing group at the IBM T.J. Watson Research Center who introduced me to Statistical Language Modelling which proved very helpful in my dissertation.

I would like give my special thanks to Eugenia Smith, who has been extremely supportive and helpful during my stay at CEDAR/CUBS. In addition I would also like to thank Mary Jane Gallo for her administrative help.

Last but not the least, I acknowledge my mates at CEDAR/CUBS for their collaboration, discussions and assistance. To name a few, they are Suryaprakash Kompalli, Karthik Sridharan, Huaigu Cao, Gaurav Chandalia, Praveer Mansukhani, Damien Jose and Anurag Bhardwaj.

# Abstract

People have become accustomed to accessing (searching, retrieving, reading) information online with ease. This information is largely limited to ASCII documents, though access to printed and spoken documents has significantly improved in recent years. Despite several decades of research in handwriting recognition, the goal of having computers access handwritten information from unconstrained document images is still elusive. Current handwriting recognition systems are only capable of recognizing words that are present in a restricted lexicon typically comprised of 10-1000 words. As the size of the lexicon grows, the recognition accuracy falls sharply and is reported to be around 30% for a 10K word lexicon. It is generally believed in the Information Retrieval (IR) community that OCR accuracy needs to be around 70%-80% for an acceptable user experience. The objective of this research is to raise the accuracy levels on unconstrained handwritten documents so that they can be accessed with the same ease as printed or ASCII documents.

We have used statistical language modeling techniques to advance the state-of-the-art in handwriting recognition. We use an innovative ‘noisy-channel’ model to correct the errors made by OCR systems and use automatic topic categorization methods based on state-of-the art statistical models like Maximum Entropy and Latent Dirichlet Allocation to reduce the effective size of the lexicon. A beam-search technique is used to prune and re-rank the multiple hypotheses using the confidences returned by the recognizer and the topic categories. We have adapted the Vector Space Model and the Latent Semantic Analysis technique to account for the multiple (top  $N$ ) recognition hypotheses. In the retrieval phase we use convolution and relevance models to improve the precision of the retrieval results.

The experimental evaluation was performed on a publicly available IAM dataset of over 1500 pages written by more than 600 writers covering a wide variety of topics and a separate dataset of 5000 hand-filled medical forms. Our topic categorization approach improves the top choice accuracy of word recognition from 30% to 45% (on average) and the top 10 accuracy is improved from 50% to the “sweet spot” range of 70%. Also, the experiments on these improved top 10 results in conjunction with our improved retrieval models show a performance at par with the retrieval on ground truth of the handwritten documents.

# Contents

	iii
Acknowledgments	iv
Abstract	vi
1 Introduction	6
1.1 Outline of the Dissertation . . . . .	11
1.2 Background . . . . .	12
1.3 Previous Work . . . . .	17
1.4 Proposed Approach . . . . .	20



<b>2</b>	<b>Handwriting Recognition</b>	<b>23</b>
2.1	Existing Recognizers . . . . .	23
2.2	Correction Models . . . . .	25
2.2.1	A phrase based correction model . . . . .	29
2.2.2	Training and Test Data . . . . .	30
2.2.3	Experimental Settings . . . . .	32
2.2.4	Character-based Model . . . . .	33
2.2.5	Word-based Model . . . . .	40
2.2.6	Evaluation . . . . .	44
2.3	Statistical Language Models . . . . .	45
2.3.1	Interpolated Language Model . . . . .	47
2.4	From OCR to Handwriting Recognizer . . . . .	48
2.5	Related Work . . . . .	49
2.6	Training and Test Data . . . . .	50

2.7	Phrase-based Direct Model . . . . .	52
2.7.1	Decoding . . . . .	53
2.8	Results . . . . .	57
<b>3</b>	<b>Lexicon Pruning</b>	<b>58</b>
3.1	Effect of Lexicon Size on Accuracy . . . . .	58
3.2	Lexicon Pruning . . . . .	59
3.3	Proposed Method . . . . .	63
3.4	Models for Topic Categorization . . . . .	67
3.4.1	Naïve Bayes . . . . .	67
3.4.2	Maximum Entropy . . . . .	69
3.5	Latent Dirichlet Allocation . . . . .	72
3.5.1	LDA Model . . . . .	73
3.5.2	Inference and Parameter Estimation . . . . .	76

3.6	Experiments . . . . .	77
3.6.1	Dataset . . . . .	77
3.6.2	Results . . . . .	78
3.7	Combination of Recognizers . . . . .	79
3.8	Conclusions . . . . .	80
<b>4</b>	<b>Information Retrieval</b>	<b>83</b>
4.1	Introduction . . . . .	83
4.2	Multiple-candidate Vector Model . . . . .	86
4.2.1	Benefit from multiple-candidates . . . . .	86
4.2.2	A Vector Model for multiple word candidates . . . . .	87
4.3	Test results and discussions . . . . .	89
4.4	Overall System for Information Retrieval . . . . .	94
<b>5</b>	<b>Conclusions</b>	<b>98</b>

5.1	Summary . . . . .	98
5.2	Contributions . . . . .	100
5.3	Future Work . . . . .	101

# List of Figures

1.0.1 Information retrieval in handwritten documents. . . . .	10
1.2.1 (a) Page of Newton’s handwritten manuscript, (b) line and word seg- mentation for the first 10 lines of text, (c) top ten choices using a random lexicon of size 1000. . . . .	16
1.3.1 Block diagram of a generic word spotting system. . . . .	18
1.4.1 Block Diagram of the proposed system. Gray blocks represent topics addressed in this thesis. . . . .	20
2.2.1 a) Typical NLP architecture b) NLP architecture with proposed cor- rection model. . . . .	27
2.2.2 Conversion of input into a character stream. . . . .	35

<i>LIST OF FIGURES</i>	2
2.2.3 Examples of merge/split errors. . . . .	36
2.2.4 Phrases extracted from training data after conversion to character stream. . . . .	38
2.2.5 Distribution of the correct word in the candidate list. . . . .	41
2.6.1 Samples generated for the same word 'word'. . . . .	51
2.7.1 Examples of merge/split errors in the handwriting context. . . . .	53
2.7.2 Phrases extracted from handwritten training data. . . . .	54
2.7.3 Example of Sentence Decoding. . . . .	55
3.3.1 Generating topic specific lexicons. . . . .	64
3.3.2 Topic Categorization in handwritten document images. . . . .	65
3.3.3 Example of N-best choices from an HR for an IAM sentence. . . . .	66
3.4.1 Plate representation of the LDA. . . . .	71
3.6.1 Example image from IAM Database. . . . .	81

3.6.2 Topic Classification Results (NB - Naives Bayes, ME - MaxEnt, suffix 10 implies top-10 choices for each word image. . . . .	82
4.3.1 An example of NYS Pre-hospital Care Report form. . . . .	90
4.3.2 Performances without pseudo relevance feedback. . . . .	92
4.3.3 Performances with pseudo relevance feedback. . . . .	93
4.3.4 MAP values using different values of $c_t \in \{0.02, 0.05, 0.075, 0.1, 0.2, 0.3\}$ are tested. . . . .	94
4.4.1 Precision-Recall comparison for the overall IR systems. . . . .	96
5.3.1 Written examples from George Washington. Worth noticing is the style of writing ‘ $p$ ’ . . . . .	102

# List of Tables

1.0.1 Effect of OCR Accuracy on Information Retrieval [14]. . . . .	8
1.2.1 Performance of line separation, word segmentation and word recognition modules. . . . .	15
2.2.1 OCR Error Rates for the three datasets. . . . .	28
2.2.2 Reduction in WER of development test set using word-based model before and after normalization. . . . .	39
2.2.3 Reduction in WER of <i>MT03</i> test set using word-based model before and after normalization. . . . .	39
2.2.4 BLEU scores for the <i>MT03</i> test set( $\pm 0.0164$ for statistical significance). . . . .	45



<i>LIST OF TABLES</i>	5
2.6.1 Details of the training and test data. . . . .	49
2.7.1 Accuracy of the HR with and without correction. . . . .	56
3.2.1 Recognition Accuracy (%) for the system described in [31] . . . . .	60
3.3.1 Top- $N$ Recognition Accuracies of the Word Recognizer. . . . .	65
3.6.1 Topic Categories of the IAM database. . . . .	78
3.6.2 Recognition Accuracies before and after lexicon reduction. . . . .	79
3.7.1 Recognition Accuracies before and after lexicon reduction using combination of two HR systems. . . . .	80
4.2.1 Recognition rates of PCR forms for different $n$ . . . . .	86
4.3.1 MAP values of retrievals using different $tf-idf$ schemes (with feedback and $c_t=0.1$ ). . . . .	91
4.4.1 Retrieval performance for the different configurations . . . . .	97

# Chapter 1

## Introduction

People have come to rely upon being able to access (organize, search, retrieve and read) information with ease. While much of this information is limited to accessing ASCII documents, recent advances in optical character recognition technology and speech recognition have allowed us to significantly extend access to printed and spoken documents. There have been substantial gains in on-line handwriting recognition as well where the process of writing is captured and advantage can be taken of pen stroke and velocity information. The problem of off-line handwriting recognition, though, is considerably more difficult. The current state-of-the-art in handwriting technology is such that when large (10000 words or more) vocabularies are referenced, recognition accuracy sharply drops. Access to handwritten documents is difficult and usually

requires manual transcription that is laborious and expensive before content can be automatically searched and organized. Such transcriptions are limited.

It has been shown that for the task of Information Retrieval (IR), in case of printed documents, even “moderate” recognition accuracies can yield encouraging results [14]. Thus current OCRs are well suited for processing printed documents for subsequent IR. It has been observed (in various Text Retrieval conferences and symposia) that when OCR accuracy is in the range of 70 – 80%, specific retrieval methods have performed well 1. For 80 – 95% accuracy, sophisticated text retrieval methods are sufficient whereas above 95%, most IR techniques are unaffected by OCR errors[14]. The ability to tolerate OCR errors stems from the redundancy of words in the target documents [46]. In addition, advanced retrieval models on noisy data (fuzzy [36], approx. string matching [66]) have been proposed to overcome these errors. The same would hold for handwritten documents and handwritten word recognizers as well. As long as the recognition rates are “moderate”, i.e., in the range 60 – 80%, a usable retrieval system can be envisioned. Unfortunately, with any realistic lexicon size of (say) above 10,000 words, the accuracy of state-of-the art handwritten word recognizers is only around 30%.

The motivation of this research arises from the following facts. A plethora of information exists in large volumes of historic and modern handwritten documents. This

OCR Accuracy	Effect on IR
95% – 100%	completely unaffected
80% – 95%	sophisticated text methods sufficient
70% – 80%	specific retrieval methods perform well

Table 1.0.1: Effect of OCR Accuracy on Information Retrieval [14].

information is, however, not accessible to the masses. Existing digital libraries (*e.g. Citeseer, IEEE, ACM*) and search engines (*e.g. Google, Altavista*) have only addressed text and printed documents. Libraries and museums possessing handwritten scholarly and historic articles intend to make them available widely. Handwritten patient information available in large volumes is a critical resource for mining demographic disease data [43]. Forms currently available in handwritten formats can be made searchable. Academic, industrial and personal notes can be indexed and made accessible. Thus, advancement in this area would be a boon for a wide range of users.

Decades of development in optical character recognition (OCR) techniques has made it possible to convert large 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 [12, 7] have shown that the performance of OCR text retrieval is severely affected when dealing with short or low quality documents. Our work is motivated by the objective of building a system for recognizing and retrieving handwritten documents.

Our goal is to automate the collection of data from the handwritten documents and enable efficient maintenance, and dissemination of information. Figure 1.0.1 describes in brief the task of information retrieval in handwritten documents. The research 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) lexicons of words are very large and iii) handwriting is often very noisy.

The most obvious methodology for searching handwritten documents is to automatically recognize the handwriting followed by retrieval methods applied on the resulting noisy text. However, recognition of large vocabulary and unconstrained offline handwritten documents has remained a very challenging task. It has been successful only in applications of limited vocabulary and constrained domains like postal address recognition and bank check processing. Due to the varying styles and forms of writing, character based recognition systems, otherwise successful in recognizing printed text, have proved to be inapplicable in the Handwriting Recognition (HR) domain.

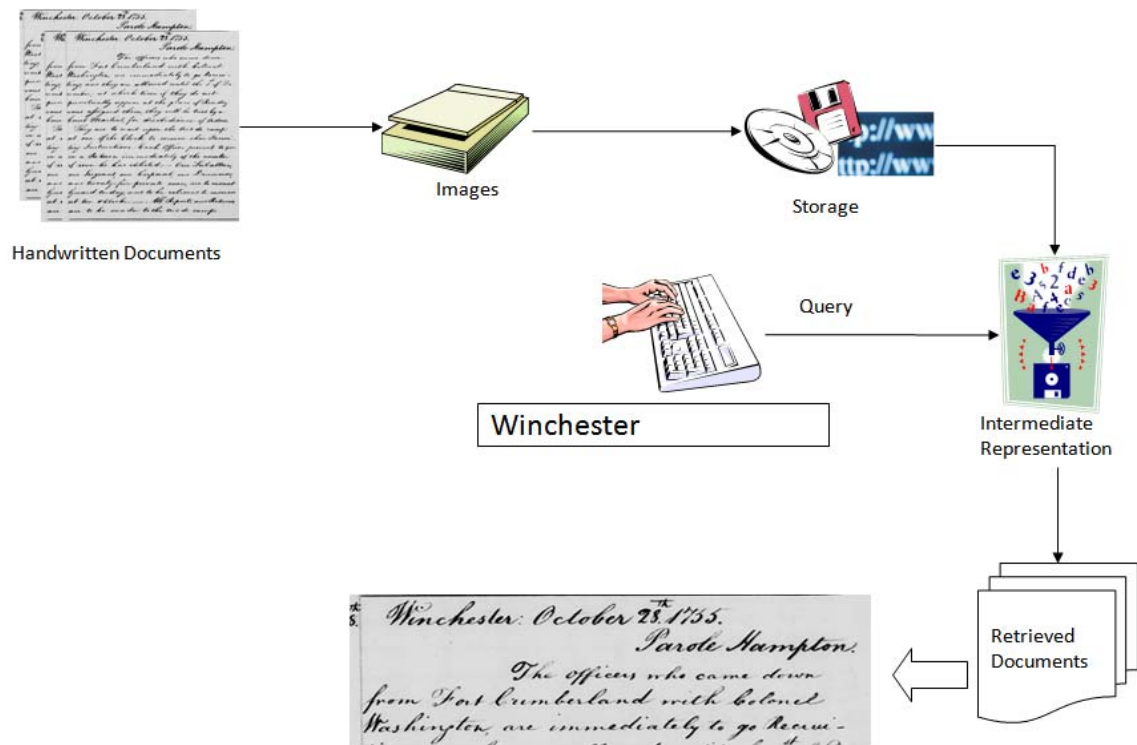


Figure 1.0.1: Information retrieval in handwritten documents.

Thus, most practical HR systems are word based i.e., the smallest unit of recognition is an isolated dictionary word. This makes the system dependent on a vocabulary, also known as the lexicon. It has been established that the recognition accuracies show an inverse relationship with the size of the lexicons. High word error rates in large vocabulary systems render the output unfit for current indexing and retrieval models. Also, given the growing volume of the documents, manual transcription to make them searchable, is economically infeasible and almost intractable. Thus, automated indexing and retrieval of handwritten documents is an open problem and

has gained impetus very recently. More recently a ‘*recognition free*’ strategy was adopted by Rath et al [53]. The following sections discuss in more detail the existing recognition based and recognition free methods for this task. The Newton Project ([www.newtonproject.ic.ac.uk](http://www.newtonproject.ic.ac.uk)) and Washington’s manuscripts [53] are only a few examples that have been undertaken due to the complexity of the problem. While manual annotation and indexing in the former proved inefficient and expensive, the latter was the first attempt at automating the process. As we shall discuss, the models are still not sufficiently mature to solve the problem.

## 1.1 Outline of the Dissertation

In this work, we propose a framework for indexing and retrieval of handwritten documents. We propose a methodology for document categorization which is not only an essential first step in indexing, but will also help reduce the vocabulary size and constrain the problem for better recognition. In addition we propose convolution and relevance based retrieval models on raw recognition results. We start with a literature survey in section 1.3 where we describe the existing methods and their drawbacks. In section 1.4 we present methods of overcoming these drawbacks. In section 4.2 we describe our methodology for indexing and retrieval. As we shall see, the challenge is to use techniques of topic categorization, query processing, and statistical language

modeling to boost the recognition rates above 60% to enable IR. This must be accomplished while treating the word recognizers as black boxes thus keeping the approach general and also flexible to accommodate different types of word recognizers within the proposed framework.

## 1.2 Background

Handwriting Recognition (HR) has been successfully used in several applications such as postal address interpretation [64], bank check reading [21], and forms reading [16]. These applications are characterized by small or fixed lexicons afforded by contextual knowledge. Recognition of unconstrained handwritten document images is difficult because of the diversity in writing styles, inconsistent spacing between words and lines, and uncertainty of the number of lines in a page as well as the number of words in a line [62]. Also, most current handwritten word recognition approaches depend on the availability of a lexicon of words for matching, making the recognition accuracy dependent upon the size of the lexicon. Accuracy dramatically decreases with increase in lexicon size and therefore domain constraints are important for solving the problem. In check processing, a small, limited lexicon of 33 is used, and with postal address domains the street names, city names and zip codes mutually constrain each other. However, for general handwritten manuscripts, these constraints do not apply and



therefore the lexicon is essentially the entire English dictionary causing the accuracy of recognition to be very low.

The inputs to a Handwriting Recognition (HR) system are a word (text block) image and a lexicon that includes the truth of the image. The handwritten text must first be located, extracted, made free of artifacts stemming from the medium (for example in the check domain these include underlines and background from the check leaf), separated into lines if necessary and finally into individual words before it can be recognized. These steps are generally nontrivial and research issues in their own right [32, 45]. Subsequently, the recognizer extracts features from the word image and matches the features against internal word models. Based on the matches, lexicon words are assigned scores or confidence values to indicate how close they are to the truth, ordered accordingly and output by the recognizer. If the truth is ranked at the top of the output, then the recognition is considered successful.

From the earliest days of research in HR, two approaches to the problem have been identified. The first approach, often called the analytical approach, treats a word as a collection of simpler subunits such as characters and proceeds by segmenting the word into these units, identifying the units and building a word-level interpretation using the lexicon. The other approach treats the word as a single, indivisible entity and attempts to recognize it using features of the word as whole. The latter approach

is referred to as the word-based or holistic approach. Because analytical approaches decompose HR into the problem of identifying a sequence of smaller sub-units, the chief problems they face are (i) segmentation ambiguity: deciding where to segment the word image and (ii) variability of segment shape: determining the identity of each segment. Holistic approaches circumvent these problems by making no attempt to segment the word into sub-units. Instead, they rely on features at the word-level to determine the identity of the word.

HRs can be also classified depending on the type of model used. In word model based recognition, all lexicon entries are treated as word models and matched against the input. The entry with the best match is the top choice. In character model based recognition, segments are matched against individual characters without using any contextual information implied by the lexicon. Word hypotheses are generated by the character recognition results. If the best hypothesis is found in the lexicon, the recognition is done; otherwise, the second best hypothesis is generated and tested, and so on. Therefore, the lexicon plays an active role in the first strategy but a passive role in the second. In the strategy of over-segmentation, segments are marked at potential character boundaries. It is acceptable if a character is segmented into sub-characters (over segmentation) but no two characters should remain connected after the segmentation process. A graph is obtained where each segment is a node. Possible traversals of this graph through all the nodes constitute various interpretations of the

segments. Instead of passing on combinations of segments to a generic OCR, lexicon can be brought into play early in the process. Combinations of adjacent segments are compared to only those character choices which are possible at the position in the word being considered. The approach can be viewed as a process of accounting for all the segments generated by a given lexicon entry. Lexicon entries are ordered according to the "goodness" of match.

	writing style # of images	discrete 6	cursive 5	mixed 9	total 20
Line Separation	# of lines	118	101	210	429
	# of lines separated	114 96.6%	97 96.0%	198 94.3%	409 95.3%
Word Segmentation	# of words	641	692	1,427	2,760
	# of words segmented	597 93.1%	631 91.2%	1,379 96.6%	2,607 94.5%
Word Recognition	top 1	432 72.4%	268 42.5%	750 54.4%	1,450 55.6%
	top 10	541 90.6%	459 72.7%	1,081 78.4%	2,081 79.8%

Table 1.2.1: Performance of line separation, word segmentation and word recognition modules.

Table 1.2 shows performance of line separation, word segmentation, and word recognition on 20 (control) document images of different writing styles [32]. As can be seen from the illustration, word recognition is the most difficult piece of the solution. Figure 1.2.1 illustrates the results obtained by the same document analysis techniques on a page from Newton's notes. Word recognition for this example would involve the

entire English dictionary. Perhaps, if the context were known, a limited but still large dictionary could be used. For example, if the page is taken from Newton’s notes on Physics, the glossary of terms from a Physics textbook can be used. We propose to use this context knowledge to improve recognition accuracies that could be potentially used for retrieval. This task would require automatic classification of documents into semantic categories based on raw recognition results and then applying this prior knowledge in an iterative fashion for recognition.

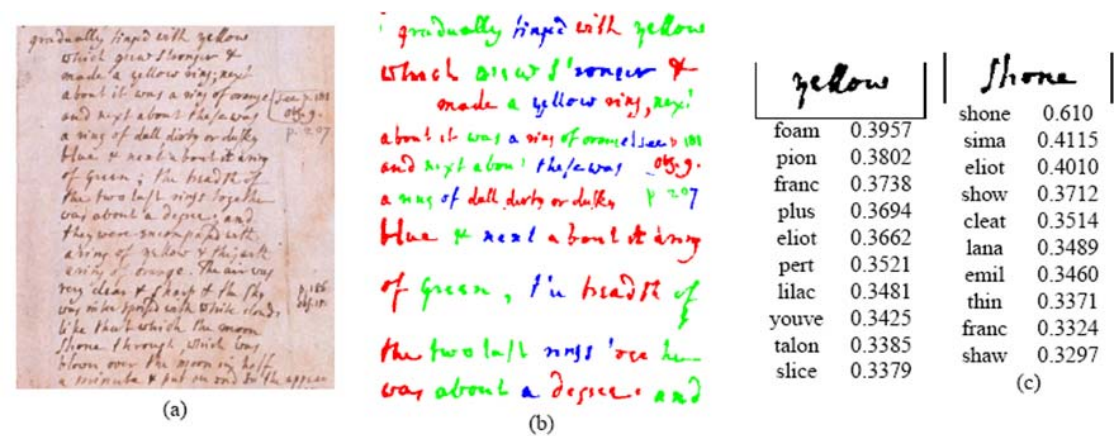


Figure 1.2.1: (a) Page of Newton’s handwritten manuscript, (b) line and word segmentation for the first 10 lines of text, (c) top ten choices using a random lexicon of size 1000.

### 1.3 Previous Work

For collections of handwritten manuscripts written by a single author, the images of multiple instances of the same word are likely to look similar. The Word Spotting idea [53] provides a different approach to building an index for such a collection. First a page is segmented into words. The different instances of a word are then clustered together using image matching. Then, a human can tag the  $n$  most interesting clusters for indexing with the appropriate ASCII-equivalent, which could be used to build a partial index for the analyzed collection.

In [29], Keaton and Goodman developed an alternative strategy based on learning a set of keyword signatures of particular words of interest. There is no page segmentation step, rather they consider a cross-correlation of the document image with a set of extracted keyword prototypes from a training set of documents. They report a classification accuracy of 91.7% for a set of 100 documents and two query keywords. In [53], Manmatha and others avoid OCR as well. 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.

Rath and Manmatha developed a word matching algorithm that compares word images using Dynamic Time Warping (DTW) [53]. DTW has been widely used in

speech processing, bio-informatics and also online handwriting communities to match 1-D signals. Although the matching of word images is in general a 2-dimensional problem, it is recast as a 1-dimensional problem since there is a loose association of image columns with the time that they were written over. By carefully preprocessing the image one can minimize the variations in the other dimension. A number of features are then extracted from each image column and the resulting feature sequences matched using Dynamic Time Warping (DTW). DTW can handle local distortions in word images and is not restricted to a single global transform.

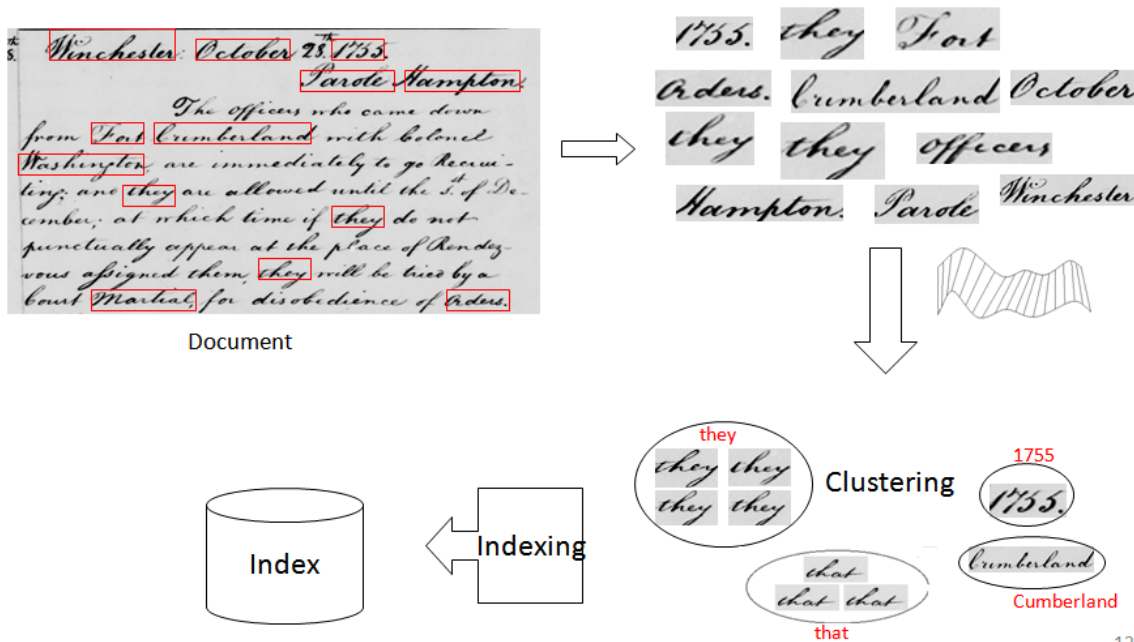


Figure 1.3.1: Block diagram of a generic word spotting system.

Figure 1.3.1 describes in brief a generic word spotting system. Where as the individual components may vary in different systems as proposed by researchers, the overall principle remains the same. One weakness of this approach is that it requires manual tagging of the clusters if they are to be queried using text. Even if one limits the word spotting technique to some small set of words it can prove valuable in a number of ways. It could be used to find multiple instances of some words. The word recognizer could then use the multiple instances to constrain its recognition and hence improve the recognition rate (see [30] for a similar idea in the context of print OCR). This could be done selectively so that one finds only multiple instances of words for which the recognizer is uncertain, or used as a verification procedure by the recognizer. While word spotting is likely to prove useful, it does not take advantage of language constraints. This method, however, has been shown to be successful in case of a single writer and also for Latin script. It performs poorly in multiple writer and unconstrained situations. It is based on the assumption that a query word has been encountered during training. In addition, conventional IR techniques such as relevance feedback and query expansion that have proven to increase retrieval performance are not applicable in this case. Further, word spotting is not suitable for scripts like Arabic where every character and word can take various canonical forms. Though, in this research, we have primarily concentrated on Latin script, our models are general and readily extend to other scripts like Arabic and Devanagari.

## 1.4 Proposed Approach

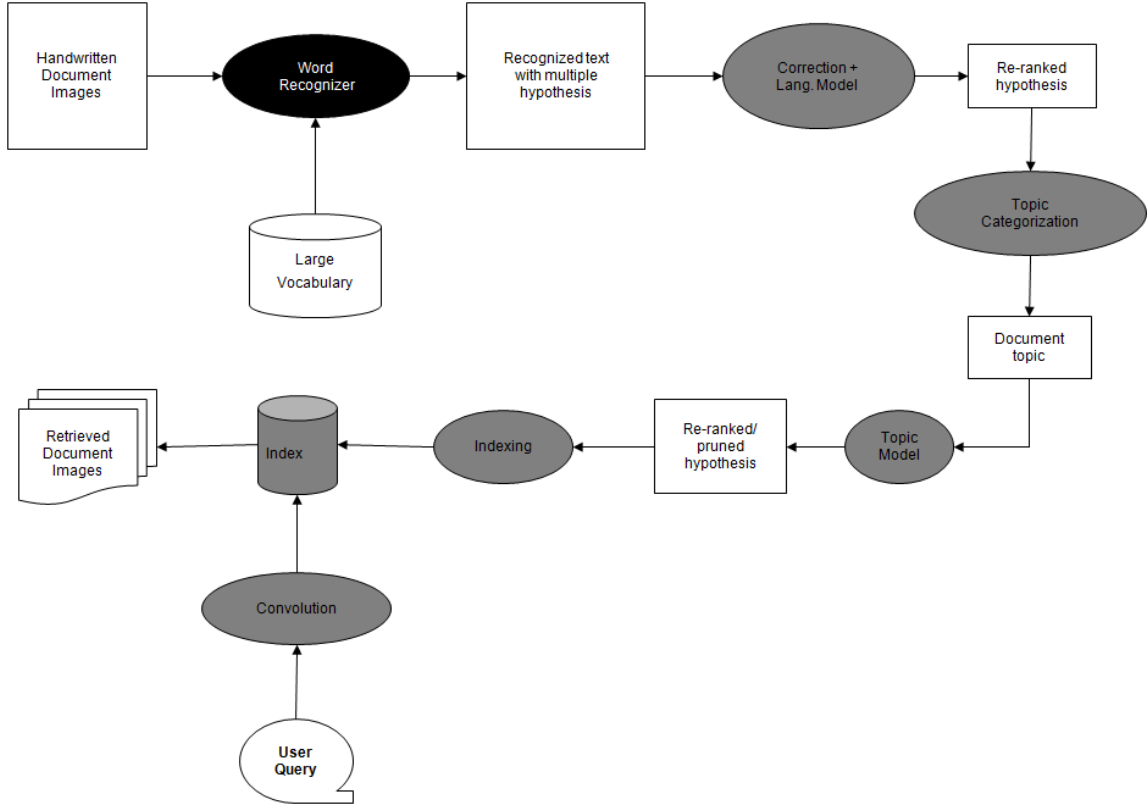


Figure 1.4.1: Block Diagram of the proposed system. Gray blocks represent topics addressed in this thesis.

Figure 1.4.1 describes in brief the overall proposed system. The blocks in *gray* will be addressed in this thesis. We will treat the handwriting recognizer as a black-box, such that there is no access to its internals and thus cannot be modified or retrained. We assume that these are supplied to us as *third-party* systems and we limit ourselves to the resources made available to us.



We have developed convolution and relevance based retrieval models on raw recognition results. Relevance feedback is an effective method for iteratively improving a query without increasing the computational requirements to perform the query. Relevance feedback uses the terms contained in relevant documents to supplement and enrich the user's initial keyword query, allowing greater retrieval performance. These methods are robust and can utilize multiple hypotheses from the handwriting recognition system. Each word of a handwritten document is recognized and mapped to a set of (word, score) pairs. For a text query, each word is statistically converted into an n-best list by assigning the maximum score to the query word and the minimum score to all the other n-best list entries. This n-best list for query generation is statistically learned from the output of the recognizer for a given set of training documents. This approach is robust when there are transcription errors and has the ability to retrieve words that would otherwise be missed by conventional transcription based methods. However, before these multiple hypotheses can be utilized, they need to be pruned and re-ranked so that the relevant ones appear higher up the list. We will describe how behavior of handwriting recognizers can be modeled in different domains and under different conditions. Behavior learning in conjunction with language models can improve the accuracy of character recognition systems significantly as described in our experiments. We show that correction models can be learned for different domains and even different writers. The *common mistakes* can be modeled using the

methodology described in our research and recognition accuracies can be improved. We have used language models and automatic topic categorization based feedback in order to accomplish this. Statistical language models describe the distribution of characters or words in a given language, thus ranking words that are more likely to follow certain words higher than the rest. Topic categorization provides an explicit context and reduces the search space as certain words are more likely to be used in a given topic and thus enable us to reduce the effective size of the lexicon.

We present research advances in five specific areas which will be described in this thesis

1. Learning the behavior of handwritten word recognizers
2. Use of contextual information and language models
3. Document categorization
4. Lexicon Pruning
5. Indexing and Retrieval

## Chapter 2

# Handwriting Recognition

### 2.1 Existing Recognizers

One of the most common constraints of current recognition systems is that they are only capable of recognizing words that are present in a restricted vocabulary, typically comprised of 10 to 1000 words. The restricted vocabulary, usually called a *lexicon*, is a list of all valid words that are expected to be recognised by the system. The lexicon is key to the success of such recognition systems, because it is a source of linguistic knowledge that helps disambiguate single characters by looking at the entire context. As the number of words in the lexicon grows, the difficulty of the recognition task

increases, because now similar words are more likely to be present in the lexicon. The computational complexity is also related to the lexicon, and it increases relatively to its size. Some open vocabulary systems, i.e. systems that are capable of dealing with any word presented at the input without relying on a lexicon, have also been proposed [63, 6], but their accuracy is still far below those relying on limited lexicons.

Most of the research efforts in the field have been devoted to improving the accuracy of constrained systems, notably small vocabulary systems. By increasing number of word candidates, the ambiguity increases due to the presence of more similar words in the vocabulary, thereby causing more confusion to the classifiers. Thus, a common characteristic of these systems is that the accuracy decreases as the number of words in the lexicon grows [70].

We have access to two kinds of recognizers. First is the Word Model Recognizer (WMR) which is described in [31]. The outputs of the recognizer are scores for a possible lexicon entry. The scores from this recognizer need to be converted into posterior probabilities as described in [11] or any other density estimation method. Another kind of recognizer is HMM based as described in [71] where posterior probabilities for each lexicon entry are output. It has been shown that when a large amount of training data is available, the performance of a word recognizer can generally be improved. However, our problem is two-fold: the lack of sufficient training data and the

availability of recognizers as black-boxes which cannot be retrained. In this chapter we will describe methods to improve the accuracy of OCRs in general using a statistical technique. We will first start with correcting the output of an Arabic machine print OCR. This is because Arabic OCRs share various common problems with HR viz., they are usually word based models due to the difficulty in character segmentation and the reported accuracies are lower than OCRs available for English. Since our model is statistical and general, it would be straight forward to extend it to the actual problem of improving HR accuracy.

## 2.2 Correction Models

Many NLP systems such as Information Retrieval or Speech-to-Speech Translation systems are multi-modal, where the input is often in different multimedia format such as scanned images of written text, audio, or video files. The building blocks of such systems are often other NLP applications pipelined together to form the wider NLP system. For instance a cross-lingual information retrieval system might be composed of an optical character recognition(OCR) system, a machine translation system, and a search engine. Therefore, such components must be robust against noisy input. Our problem here is similar. We have output from handwriting recognition systems in form of noisy text and our objective is information retrieval.

Component systems are often developed independently and possibly by different groups. In most cases, the internals of one component are not accessible to the developers of the subsequent components in the pipeline. In such cases, these components (e.g. OCR) must be treated as black boxes where only their output is observed (Figure 2.2.1a). Unfortunately, the output of these systems is error-prone, especially for non-Roman languages (e.g., Arabic) which have not been researched as extensively. Therefore, the downstream component must be robust enough to handle noisy input. The problem is compounded when the input to OCR systems, in the form of images, is also noisy.

There has been some work on post-processing of the OCR output in the machine print domain. A general survey of the research in this area can be found in [35]. Jones et al. [38] describe a multi-pass OCR post-processing system which carries out individual word corrections, combined edit distance corrections and bigram probability based correction in different passes. Perez-Cortes et al. [23] use a stochastic finite state machine to test hypotheses of words. If the machine accepts the word, then no correction is made, otherwise smallest set of transitions that could not be traversed show the most similar string in the model. Pal et al. [69] describe a method for OCR error correction of Devanagiri script using morphological parsing. Some of these techniques are highly dependent on the language and use features that are specific

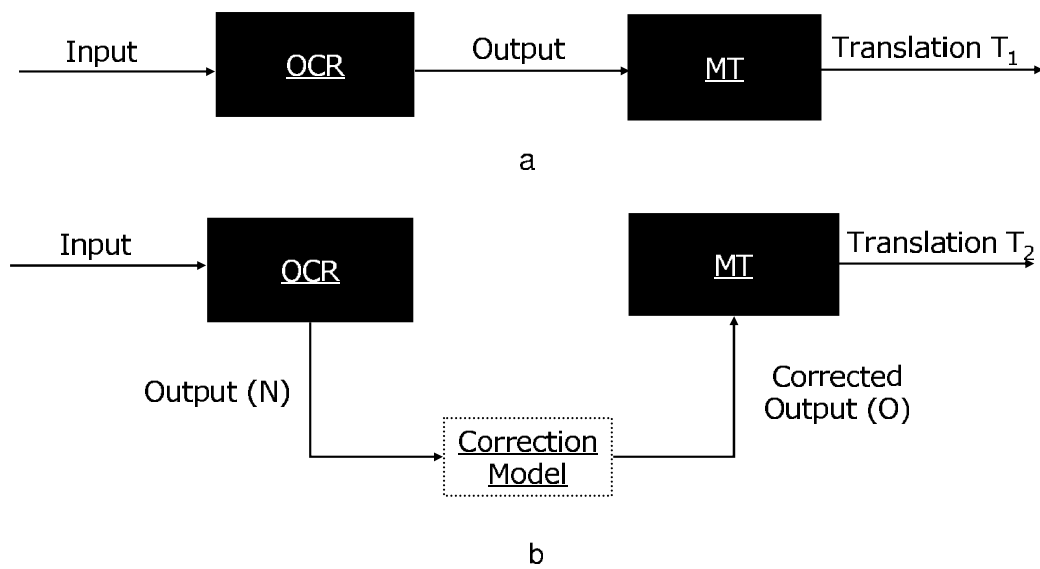


Figure 2.2.1: a) Typical NLP architecture b) NLP architecture with proposed correction model.

to the language in question, which makes such techniques difficult to adapt to a different language. The method we present here is trainable and hence scalable to other languages.

Although we show results on Arabic OCR and the effect on Arabic-to-English translation only, the technique is adaptable to other languages where training data is available. Okan et al. [48] introduce a generative probabilistic OCR model that describes an end-to-end process in a noisy channel framework. The technique they describe is implemented as a cascade of finite state transducers. We propose a different approach to this problem. We cast this problem as a simplified translation task where the source language is the error-full OCR output and the target language is the corrected output. Our model is statistical, easy to train and general. We later extend this model to handwriting recognizers.

	Training Data	Development Test Set	<i>MT'</i> 03 Eval Set
WER(%)	25.20	21.21	28.87
CER(%)	10.60	8.16	11.00

Table 2.2.1: OCR Error Rates for the three datasets.



### 2.2.1 A phrase based correction model

The erroneous output from the OCR can be thought of as an output from a noisy channel problem[37]. It can be considered as a black-box through which the signal (truth) when passed gets corrupted and emerges out as the degraded output. We can learn the corruption models (correction model in the direct model framework) and then use this knowledge to correct the output. We leverage our experience in building SMT systems by casting this correction task as a simplified translation task. We propose a direct phrase-based translation approach that has proven effective in machine translation (e.g., [68], [50]).

As is done in a typical SMT system, given *sentence pairs* in the source (*foreign*) and the target (*english*) languages, we first align the source and target words. Word alignments can be obtained using one of several techniques such as HMM alignments [59], Maximum Entropy alignments [22], or Maximum-Posterior alignments [5]. Once a sentence pair is word-aligned, we extract phrase pairs (also referred to as block in the statistical machine translation literature) in a manner similar to [68]. These phrases are then used in combination with an  $n$ -gram language model to translate the source language into the target language.

We model the correction as a simplified translation task. The source language in our case is the OCR output (error-full) and the target language is the truth (correct).

Since there is no word re-ordering as is the case in a typical MT task, the alignments can be obtained easily by minimum edit distance alignments or monotone HMM alignments. We used a simple dynamic programming minimum edit distance alignment. We approached this task in two ways. First, in cases where the correction was achieved a character at a time( in terms of translation, the units of translation were taken to be characters instead of words as in a typical SMT system). The dynamic programming alignment between the truth and the corresponding OCR output was obtained at a character level. In the second method the same technique was applied, however, the units were words. Thus, we implemented two models - a character-based model and a word-based model. These models are explained further in sections 2.2.4 and 2.2.5.

## **2.2.2 Training and Test Data**

### **2.2.2.1 Training Data**

We use a commercially available Arabic OCR to recognize Arabic machine printed images for which we have the ground-truth. In order to create our correction training data, images of Arabic documents were needed. Since we did not have access to any publicly available large dataset of Arabic images and the corresponding text version,

we automatically generated images from the Arabic text of the LDC parallel corpus [3]. We used a Microsoft Word macro to open the *UTF-8* Arabic text files, format them using Bold Arabic Transparent Font at 14*pt.* and printing them into postscript files. The postscript files were then converted into jpg images at 300 dpi. This gave us a rich source of  $\sim 4000$  images containing 23205 segments or  $\sim 1.1M$  words. These images were recognized using the Arabic OCR and the OCR output, paired with the truth, was used as the training data for our models.

#### 2.2.2.2 Test Data

The development and blind test sets were created in the same manner described above. The development test set was created from editorials from an online Arabic newspaper ([www.asharqalawsat.com](http://www.asharqalawsat.com)). It consists of 390 segments,  $\sim 11K$  words. The blind test set was the NIST MT'03 [2] evaluation set. This test set is typically used to report the quality of machine translation systems. It consisted of 663 segments or  $\sim 15K$  words and was used in two-fold ways *i*) as a blind test set for the correction task, *ii*) to test the effect of input degradation on the SMT system. This was achieved by calculating the MT quality for the original set, then the quality of MT by translating the OCR output acquired from the corresponding images and finally for the corrected output.

### 2.2.3 Experimental Settings

We use the evaluation metrics reported by [48] for OCR evaluation. The evaluation metrics are the word error rate (WER) and the character error rate (CER) and are standard in measuring the accuracy of OCRs. The metrics are defined as follows:

$$WER = \frac{WordEditDistance(W_{truth}, W_{Output})}{|W_{truth}|}$$

$$CER = \frac{CharEditDistance(Truth, Output)}{|Truth|}$$

where the edit distance is the Levenshtein Edit Distance and is defined as the minimum number of point mutations required to change one unit(character or word) into another unit, and a point mutation is one of – substitute, insert or delete of a unit.  $|Truth|$  and  $|W_{truth}|$  is the number of characters and words respectively in the ground-truth. The error rates for the OCR for the three sets is given in Table 2.2.

Our correction model is evaluated first by measuring the reduction in recognition error rates. The error reduction is measured by comparing the error rate in  $N$  and  $O$  shown in Figure 2.2.1. Since MT is word-based, and any improvement in the WER would potentially affect the MT performance, we have focused solely on the WER of

the OCR or the corrected output. Secondly, we evaluate our correction model in the machine translation context by measuring the overall improvement in MT when the corrected OCR output is used instead of the original OCR output. This is achieved by comparing the translation quality (as measured by BLEU [25]) of  $T_1$  and  $T_2$  shown in Figure 2.2.1.

## 2.2.4 Character-based Model

Most OCR systems are conventionally segmentation based i.e. the text lines are segmented into words and words into their constituent characters prior to actual recognition. The OCR usually repeats the errors it makes in recognizing a certain character. In order to model the character-based nature of the OCR system, we defined the correction task as a translation of the character sequence that the OCR produces into the corresponding ground truth. Thus, the units of our input are characters instead of words in the case of a typical MT task.

### 2.2.4.1 Conversion to Character Stream

The first step for this task requires conversion of the input text into a stream of characters. This was achieved by separating individual characters by spaces. We

used a special symbol `<bs>` to denote the word boundaries in order to differentiate it from the space between the characters. Figure 2.2.2 shows the same example of conversion in Arabic (romanized for illustration). This was done for segments on both sides - the ground truth as well as the OCR output. Thus, our parallel corpus consisted of a character stream and its corresponding output from the OCR. It should be noted that the conversion on both sides into characters was performed after and not before recognition.

#### 2.2.4.2 Training

In order to extract the phrases from the training data, we align the truth and the OCR output character streams obtained above. Since there is no word re-ordering, we use a dynamic programming based minimum edit distance alignment. The edit distance was calculated using the regular Levenshtein edit distance formula with equal cost for deleting, inserting or substituting a unit. Using this, one-to-one alignments are obtained for the paired truth and output datasets. An unaligned character in the truth models a character deletion. Similarly an unaligned character in the OCR output models an insertion. An unaligned `<bs>` in the truth models a merge error where the OCR merges two words and an unaligned `<bs>` in the OCR output models a word split. Figure 2.2.3a shows an example where two words were merged by the

**Input String**

Arabic

وكان تم نقل التكريتي مع 17 سفيرا عراقيا

Romanized

wkAn tm nql Altkryty mE 17 sfyrA ErAqyA

**Output Character Stream**

Arabic

وكان&lt;bs&gt;تم&lt;bs&gt;نقل&lt;bs&gt;التكريتي&lt;bs&gt;مع&lt;bs&gt;17&lt;bs&gt;سفيرا&lt;bs&gt;عراقيا

Romanized

w k A n <bs> t m <bs> n q l <bs> A l t k r y t y <bs> m E <bs> 1 7 <bs> s f y r  
A <bs> E r A q y A

Figure 2.2.2: Conversion of input into a character stream.

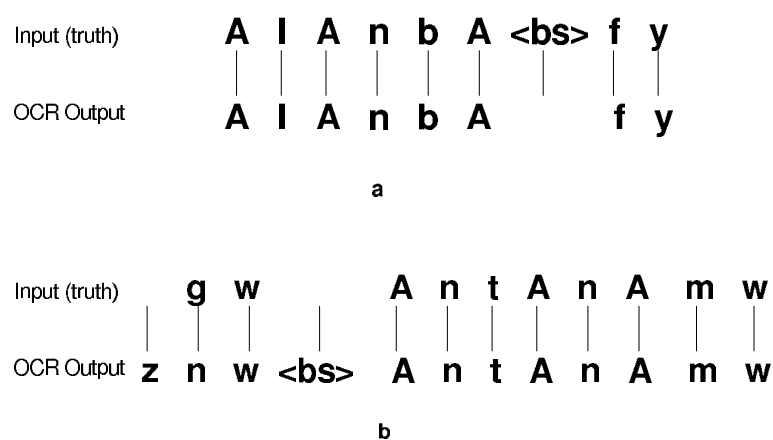


Figure 2.2.3: Examples of merge/split errors.

OCR into one and the **<bs>** is unaligned. Similarly, Figure 2.2.3b shows an example where a word was split and a character inserted (*z*). This time the **<bs>** on the OCR side is unaligned.

Using this alignment we extract all phrase pairs of lengths 1 through  $n$  ( $n = 6$  in our case) characters. Figure 2.2.4 shows example phrases extracted from the training data.



### 2.2.4.3 Decoding

To correct the OCR output for a given test sentence, we “translate” this sentence by decoding using two weighted components - the phrases obtained above and the language model. We trained a character trigram language model by converting the gigaword corpus [1] into characters using the same procedure as above. The original corpus has 319 files totaling  $\sim 4.3GB$  and  $\sim 39.1M$  words. The language model is an interpolated language model. The *lambdas* were calculated from held out data, which is the last 10% of the gigaword corpus.

The decoder is a beam search decoder. For the correction problem, the decoding is carried out in a monotone fashion such that characters at the beginning of the sentence must be decoded first. Decoding hypotheses are extended in cardinality-synchronous fashion such that at any given time (step  $n$ ), only those hypotheses that cover exactly  $n$  characters (i.e., of cardinality  $n$ ) are extended. For example, at time step 0 only hypotheses of cardinality 0 (in the special init hypothesis) are extended. The decoding is continued until all characters in the sentence are covered. The final decoding is the hypotheses of cardinality  $l$  with the least cost, where  $l$  is the length of the source text. More formally, the final decoding  $e$  for the source  $f$  is the one that satisfies the following equation:

Source (OCR)	Target (Truth)	Count
ي ن (n y)	ي ن (n y)	17437
ي ن (n y)	غ ي (g y)	22
ب ي ن و (n y b w)	ب ي ن و (n y b w)	2
ب ي ن و (n y b w)	ق ي ن و (n y q w)	2
ب ي ن و س (n y b w s)	ق ي ن و س (n y q w s)	35

Figure 2.2.4: Phrases extracted from training data after conversion to character stream.

$$\hat{e} = \underset{e}{\operatorname{argmax}} [-w_{ph} \times \log_{10} P(f|e) - w_{lm} \times \log_{10} P(e)]$$

where  $P(e)$  is the trigram character language model probability and  $P(f|e)$  is the phrase-based model. The weights  $w_{ph}$  and  $w_{lm}$  of these components were trained using the first 50 segments of the development test set. The weights are trained using maximum BLEU training such that the weights that resulted in the maximum BLEU score for these 50 segments are selected for correcting the rest of the test set and any subsequent test sets using the monotone decoding (as described earlier). Since,

	OCR Output	Corrected Output	Improvement (abs.)	Improvement (rel.)
WER % (Before Norm.)	21.21	20.08	1.13	5.32
WER % (After Norm.)	17.7	13.4	4.3	24.29

Table 2.2.2: Reduction in WER of development test set using word-based model before and after normalization.

	OCR Output	Corrected Output	Improvement (abs.)	Improvement (rel.)
WER % (Before Norm.)	28.87	15.28	13.59	47.07
WER % (After Norm.)	13.71	11.54	2.17	15.83

Table 2.2.3: Reduction in WER of *MT03* test set using word-based model before and after normalization.

there were only two weights to be trained, we also adopted a grid-search technique. The translation accuracy(measured by WER) was calculated by varying the phrase component weight from 1.0 to 0.5 in steps of 0.1 and correspondingly the language model weight was varied from 0.0 to 0.5.

#### 2.2.4.4 Results

The character-based model resulted in a 1% absolute (12% relative) reduction in the CER of the development test set. There was no improvement in the WER, which remained at 21.21%. This could be attributed to a number of reasons. In addition to the character trigram language model being weak, the phrase features were also not helpful. It was noted that the OCR had a post-processing step where the recognized output was converted to the closest possible dictionary word. This led to the fact that besides correctly recognizing the character, the rest of the confusion matrix was smooth rather than sharp (as was expected), thus making the phrase features less reliable. Due to this post-processing, the actual error that the OCR would make on a character was hidden. By the creation of a dictionary word, the OCR output character could be replaced by another character thus making it difficult to model the frequent OCR errors.

#### 2.2.5 Word-based Model

As noted in section 2.2.4.4, the character-based model proved to be inadequate for correcting OCR errors. To overcome this, a word-based model is pursued, which proved to be useful. The word trigram language model would be much more beneficial

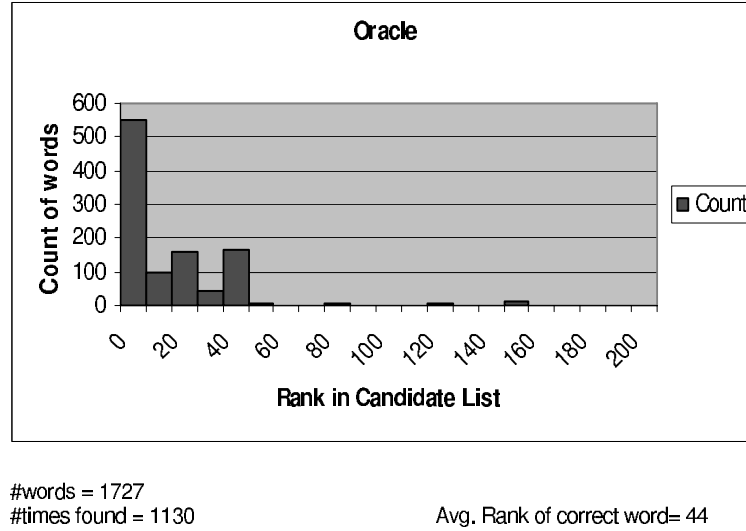


Figure 2.2.5: Distribution of the correct word in the candidate list.

since it provides more context.

### 2.2.5.1 Training

The training data was again aligned using the minimum edit distance, however, this time it was based on words rather than characters. After aligning the truth and the OCR output, we extracted the word phrases just as we extracted character phrases as described earlier in section 2.2.4.2. The drawback with this, however, would be that in order to get features for all possible words, we need much more data.

During training, the candidates for a word that the OCR produced were either found

from the training data or from  $V_{lm}$  based on the minimum edit distance and sorted on basis of counts in the training data or increasing order of edit distance. During decoding, candidates were selected from the list based on the formula as mentioned earlier. Selecting the correct word would be based on two factors - whether the correct word is a part of this list and whether the cost according to the formula above is the best. Figure 2.2.5 shows the oracle experiment as to how many times the correct word was found in our candidate list, irrespective of whether the decoding phase selected it or not. This is sampled from 1727 errors that the OCR makes in the development test set. As noted 1130 words out of 1727 were found in the candidate list at an average rank of 44. Thus, the best we can achieve, if all correct candidates were to be selected, is a 65.4% improvement. In addition to this, it was observed that 98% of the true words in the test set were in our vocabulary  $V_{lm}$ .

In order to overcome the above described problem, we adopted the following procedure:

1. Extract phrase pairs of length=1 pairs  $P = \bigcup p_w$ , where  $p_w$  is a phrase pair containing source word  $w$  and all target word(s) from the training data.
2. Get the vocabulary  $V_t$  of the test set.
3.  $\forall w_i \in V_t$ , if  $\exists p_{w_i} \in P$  such that  $count(p_{w_i}) < T$  then  $P = P \cup \{w_i, neighbour(w_i)\}$ , i.e. add  $w_i$  and  $neighbour(w_i)$  to  $P$ .

Given a very large language model vocabulary  $V_{lm}$ , the function  $neighbour(w_i)$  yields all words  $w'_i \in V_{lm}$  such that the minimum edit distance of  $w'_i$  from  $w_i$  is less than a threshold  $T'$ . For example, if the OCR output word was *AntAnAmw* (antanamo),  $neighbour(AntAnAmw)$  would yield  $\{AntAnAmw, gwAntAnAmw, ntAnyAhu, \dots\}$  (antanamo, Guantanamo, Netanyahu, ...). The thresholds  $T$  and  $T'$  are chosen empirically.

### 2.2.5.2 Decoding

Translation was carried out in the same manner as in case of a character-based model, however, the units this time were actual words. We trained a word trigram language model from the gigaword corpus[1] as described in section 2.2.4.3. However, this time there was no conversion into character streams. The weights for the phrase and the language model components were obtained by the grid search technique described in section 2.2.4.3.

### 2.2.5.3 Results

The word model performed better than the character model in general: 1.13% for the development test set and 13.71% for the blind test set absolute improvement in WER is observed, as compared to no improvement in the character model. Since Arabic is

highly inflectional, many of the errors were prefix/suffix errors. Our trigram language model had difficulty predicting the right set of prefix/suffix combination to choose for Arabic words. Other errors were less severe. For example, confusing various forms of *alef* (i.e., *alef* without *hamza* or with *hamza* above or below it). Since various forms of *alef* are mapped to one canonical form by the MT engine (in a process we call *normalization*), these errors were ignored by normalizing the OCR output as well as the corrected output that our model generated. The WER of the OCR after normalization was 17.7%. However, after normalizing the corrected output, the WER dropped from 20.08% to 13.4%. Thus, we were able to get a 4.3% absolute reduction in the WER after correction compared to 1.13% absolute reduction before normalization as shown in Table 2.2.4.3.

A similar improvement is observed for our blind test set as well. The WER before normalization was reduced from 28.87% to 15.28%, and from 13.71% to 11.54% after normalization as shown in Table 2.2.4.3.

### 2.2.6 Evaluation

In order to evaluate the effect of errors introduced by the OCR in an Arabic to English Machine Translation system, we calculated the quality of MT before correction and after correction of the OCR output for the *MT03* dataset. The translation is



	Ground Truth	OCR Output	Corrected Output
BLEU Score	0.4651	0.3776	0.4013

Table 2.2.4: BLEU scores for the *MT03* test set( $\pm 0.0164$  for statistical significance).

performed on the normalized text and the BLEU scores are calculated with the NIST provided reference[2]. Normalization is carried out as described in section 2.2.5.3. Table 2.2.6 shows the BLEU score for translation of the *MT03* test set before and after correction. The BLEU score for the actual *MT03* test set(truth) is also reported for comparison. We achieve a statistically significant improvement in the BLEU scores.

## 2.3 Statistical Language Models

The goal of Statistical Language Modeling is to build a statistical model that can estimate the distribution of natural language as accurate as possible. A statistical language model (SLM) is a probability distribution  $P(s)$  over strings  $s \in S$  that attempts to reflect how frequently a string  $S$  occurs as a sentence. The original (and is still the most important) application of SLMs is speech recognition, but SLMs also play a vital role in various other natural language applications as diverse as machine translation, part-of-speech tagging, intelligent input method and Text To

Speech system. N-gram model is the most widely used SLM today.

We can express the probability of a string  $s$  :  $P(s)$  as

$$P(s) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)...P(w_l|w_1, ..., w_{l-1}) = \prod_i^l (P(w_i|w_1, ..., w_{i-1})) \quad (2.3.1)$$

Depending on the size  $n$  of context, it is referred to as an n-gram model. For example, in bigram models, we make the approximation that the probability of a word only depends on the identity of the immediately preceding word, hence we can approximate  $P(s)$  as:

$$P(s) = \prod_i^l P(w_i|w_{i-1}) \quad (2.3.2)$$

Similarly, in tri-gram, it can be approximated as

$$P(s) = \prod_i^l P(w_i|w_{i-1}, w_{i-2}) \quad (2.3.3)$$

Whereas larger values of  $n$  provide more context, they are often computationally infeasible to manipulate. Traditionally, tri-gram models have been extensively used as a trade-off between complexity and accuracy.

Use of Statistical Language Models in handwriting recognition has recently gained impetus motivated by speech recognition. Where in some cases it has been used simply as post-processing independently from the core handwriting recognition unit, in others it has been used in conjunction with the recognition engine. In our case we use the handwriting recognizer as a black-box which may or may not have language modeling built into the system and we concern ourselves only with the final outcome. As mentioned earlier, the correction model has two components, and since the second component is the language model on the “correct” text, it is easy to train the model on any regular english (or corresponding language) text and use it in conjunction with the phrase model. In our experiments we have used an interpolated tri-gram language model.

### 2.3.1 Interpolated Language Model

Linear interpolation, the most popular model combination technique, first trains individual  $n$ -gram LMs from each training corpus. Given the resulting set of  $n$ -gram language models, it computes the weighted average of the component model probabilities  $p^{LI}(w|h) = \sum_i \lambda_i p_i(w|h)$ , where  $p_i(w|h)$  is the smoothed probability of word  $w$  following  $n$ -gram history  $h$  in model  $i$ . The interpolation weight  $\lambda_i$ , satisfying

$\sum_i \lambda_i = 1$ , is typically tuned to optimize the development set perplexity. Any language model can be used as a component. In fact, once a commonset of held out data is selected for weight optimization, the component models need no longer be maintained explicitly. Instead, they can be represented in terms of the probabilities they assign to the held out data. Each model is represented as an array of probabilities. In our experiments we used the last 10% of the language modelling data to calculate the  $\lambda_i$ s. We experimented with the very useful publicly available language modeling toolkits [56, 65].

## 2.4 From OCR to Handwriting Recognizer

The state-of-art in handwriting recognition for a moderate lexicon of size  $\sim 10k$  is  $\sim 40\%$  [70]. Any improvement in the performance of the recognizer will improve the overall user experience of the systems which have HR as one of the components. We describe a novel “phrase-based” system for the correction of the output of a word recognizer that was developed in CEDAR [71]. We extend our correction model to the handwriting recognition paradigm and treat the recognizer as a black box that requires three inputs viz. i) a word image, ii) a lexicon and iii) an integer  $n \leq 10$ . As outputs we can acquire top- $n$  hypotheses from the lexicon for that word image. We process the output from the HR and correct it using a model trained for that particular

	Documents	Lines	Words
Training	1055	7878	127690
Test	340	2540	40789
Total	1395	10418	168479

Table 2.6.1: Details of the training and test data.

HR system. Thus, we add an additional module in the pipeline. We were able to extend our method for Arabic to English handwriting without any modifications.

## 2.5 Related Work

The work mentioned above has been applied on machine print OCR that conventionally uses “character-models” as opposed to the handwriting recognition systems that follow the word-based “multiple choice” paradigm. Thus, training the character confusion matrices is not straight forward. The method we propose here is trainable. Although we only show results on the in-house recognizer, the technique is adaptable to other recognizers and even other scripts where training data is available.

## 2.6 Training and Test Data

Due to lack of large truthed handwritten datasets, we automatically generated images from the english text used in the Conference on Computational Natural Language Learning (CONLL 2003) shared task (<http://www.cnts.ua.ac.be/conll2003/ner/>). Since our model required learning the behavior of the recognizer and the most common corrections for the mistakes that it makes, it is critical that the same input to the recognizer is provided in various styles and forms. This was not possible from any existing sentence database. Thus, we had to rely on generating data that would simulate handwriting. In case such databases are available, they can easily be utilized or even added to the current model. The handwriting is generated from text by concatenating character templates generated by the Blums Medial Axis Transform [10] from actual handwritten samples. This is followed by character auto-scaling, automatic baseline determination, ligature modeling, ligature joining, skeleton thickening and smoothing. The visual analysis and accuracies of the recognizers on this data support our simulation technique for generating handwriting samples. Details of this method can be found in [67]. In order to model the multiple writer scenario we sampled the characters at random from a set. Thus, during generation of a word e.g. 'word' at different positions in the data we sampled uniformly for  $i$  from different handwritten character samples  $\{w_i, o_i, r_i, d_i\}$ . Figure 2.6.1 shows the different samples generated

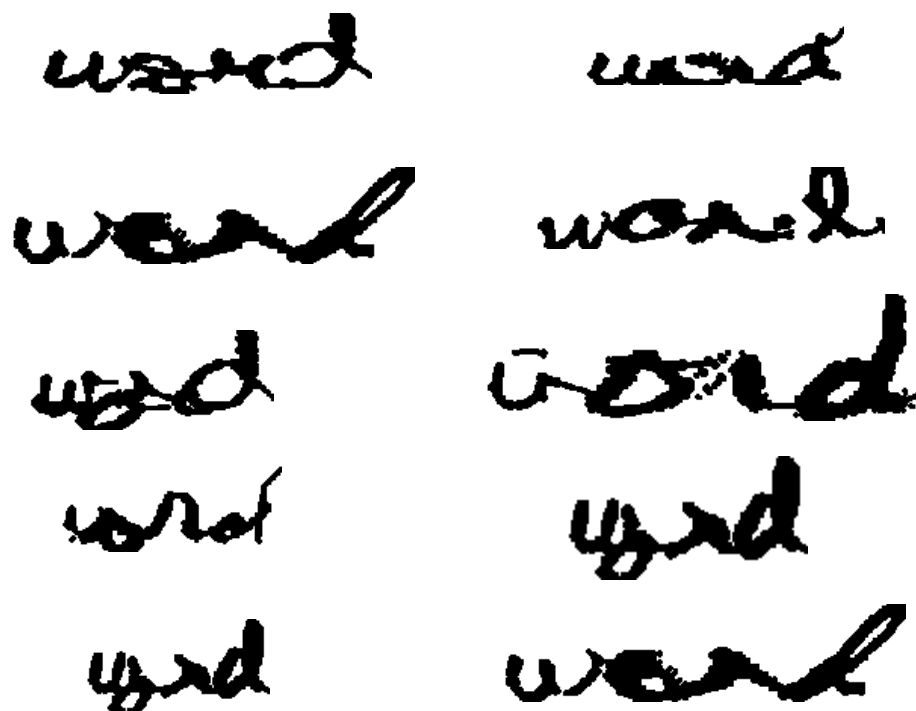


Figure 2.6.1: Samples generated for the same word 'word'.

for the same word.

We split the data into a training set (75%) and a test set (25%). The training data consisted of 127690 words from 1055 documents. Similarly, the test data consisted 40789 words from 340 documents. Table 2.6 gives the details of the two sets.

The in-house handwriting recognizer used for recognition of these images is a lexicon driven HMM based word model recognizer [71]. The recognizer takes as input a

lexicon and the word image. A lexicon of size  $\sim 21K$  was generated from the training data. Thus, during the recognition of the training data we used a closed lexicon such that the truth of the image being recognized was present in the lexicon. We used the same lexicon for recognizing the test data and observed that it covered 95% of the words and the rest were treated as OOVs (Out Of Vocabulary).

## 2.7 Phrase-based Direct Model

In order to extract the phrases from the training data we align the truth and the OCR output character streams obtained above. Since there is no word re-ordering, we use a dynamic programming based minimum edit distance alignment. The edit distance was calculated using the regular Levenshtein edit distance formula with equal cost for deleting, inserting or substituting a unit. Using this, one-to-one alignments are obtained for the paired truth and output datasets. An unaligned character in the truth models a character deletion. Similarly an unaligned character in the OCR output models an insertion. An unaligned *blank* in the truth models a merge error where the OCR merges two words and an unaligned *blank* in the OCR output models a word split. Figure 2.7.1a shows an example where two words were merged by the OCR into one and the *blank* is unaligned. Similarly, Figure 2.7.1b shows an example where a word was split and a word inserted (*el*). This time the *el* on the OCR side



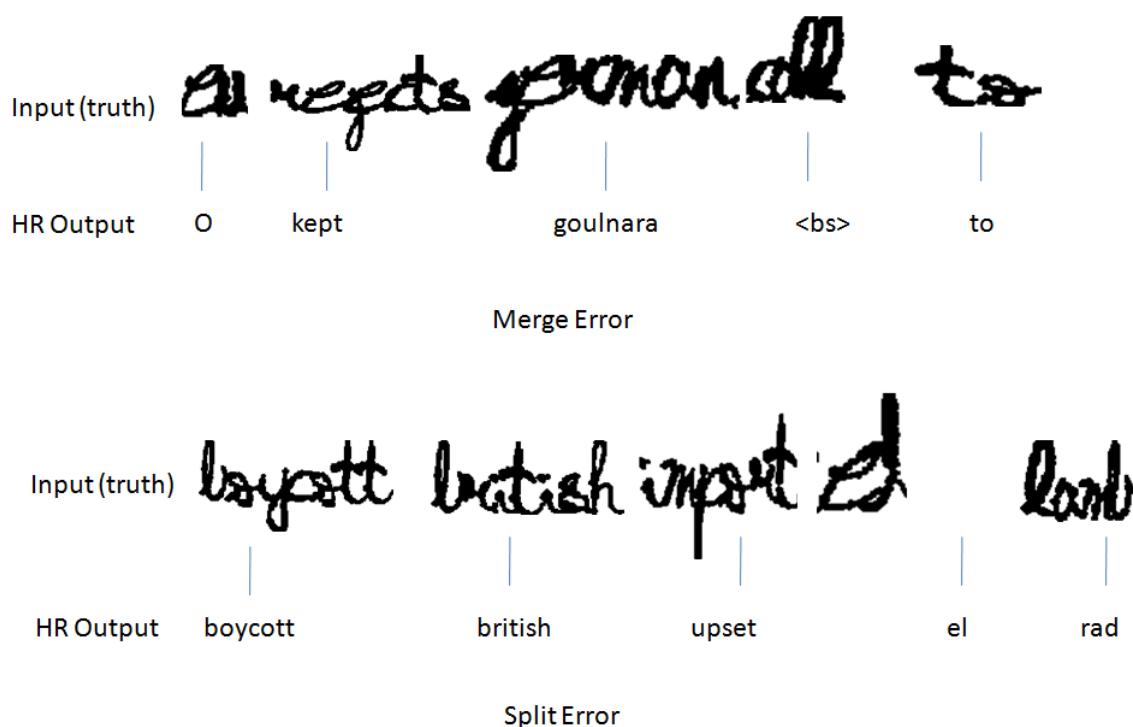



Figure 2.7.1: Examples of merge/split errors in the handwriting context.

is unaligned.

Using the alignment we extract all phrase pairs of lengths 1 through  $n$  words. Figure 2.7.2 shows example phrases extracted from the training data with the possible corrections along with their correction probabilities for a given word in the left column.

### 2.7.1 Decoding

To correct the OCR output for a given test sentence, we “translate” this sentence by decoding using two weighted components - the phrases (obtained above) and the



The diagram illustrates the relationship between recognized words and hidden words. On the left, the text "Recognized words" is accompanied by a downward-pointing orange arrow. Above the table, the text "Hidden words" is followed by a rightward-pointing blue arrow. A callout box labeled "Probability" points to the numerical values in the table cells.

<b>pitcher</b>	<b>pitcher</b> 1.00			
<b>pascolo</b>	<b>financially</b> 0.40	<b>pascolo</b> 0.20	<b>speculates</b> 0.20	<b>poisonous</b> 0.20
<b>notation</b>	<b>protection</b> 0.50	<b>invitation</b> 0.40	<b>motivation</b> 0.05	<b>notation</b> 0.05
<b>experts</b>	<b>experts</b> 0.88	<b>expired</b> 0.13		
<b>updated</b>	<b>injunction</b> 0.40	<b>uprooted</b> 0.20	<b>infrastructural</b> 0.40	

Figure 2.7.2: Phrases extracted from handwritten training data.

language model. We trained a trigram language model on the Reuters data (<http://trec.nist.gov/data/reuters/reuters.html>) using the CMU language modeling toolkit [56].

The decoder is a trigram viterbi decoder. For the correction problem, the decoding is carried out in a monotone fashion such that words at the beginning of the sentence must be decoded first. Decoding hypotheses are extended in cardinality-synchronous fashion as in the case of the Arabic system. The final decoding  $e$  for the source  $f$  is again the one that satisfies the following equation:

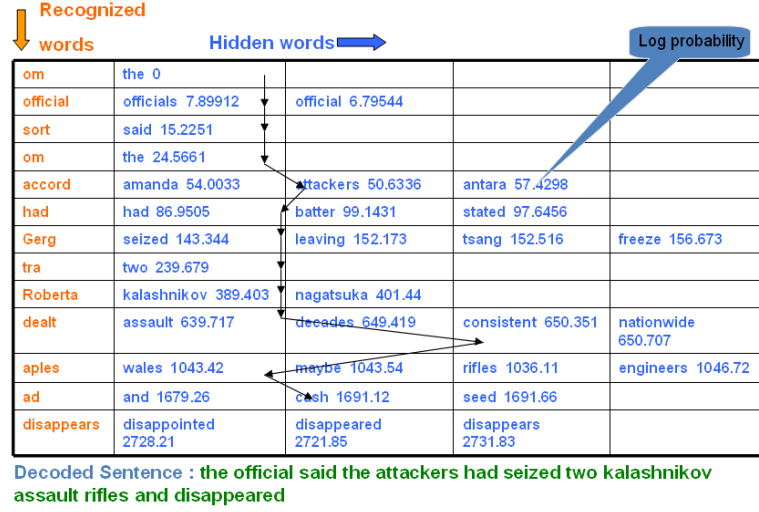


Figure 2.7.3: Example of Sentence Decoding.

$$\hat{e} = \arg \max_e [-w_{ph} \times \log_{10} P(f|e) - w_{lm} \times \log_{10} P(e)]$$

where  $P(e)$  is the trigram character language model probability and  $P(f|e)$  is the phrase-based model. In this situation, used equal weights  $w_{ph}$  and  $w_{lm}$  ( $w_{ph} + w_{lm} = 1$ ) for both mixture components. However, in general these weights need refinement (can also be trained) based on the prior evaluation of the word recognizer. For example, a recognizer with better accuracies would require a slightly larger  $w_{ph}$  and vice-versa.

	Raw	Top-10+3-gram LM	Corrected
Accuracy	8.1%	13.5%	71.3%

Table 2.7.1: Accuracy of the HR with and without correction.

The drawback of this approach, however, would be that in order to get features for all possible words, more data is needed. In order to overcome this problem, we adopted the following procedure:

1. Extract phrase pairs of length=1 pairs  $P = \bigcup p_w$ , where  $p_w$  is a phrase pair containing source word  $w$  and all target word(s) from the training data.
2. Get the vocabulary  $V_t$  of the test set.
3.  $\forall w_i \in V_t$ , if  $\exists p_{w_i} \in P$  such that  $count(p_{w_i}) < T$  then  $P = P \cup \{w_i, recog(w_i)\}$ ,  
i.e. add  $w_i$  and  $recog(w_i)$  to  $P$ .

Given the test lexicon  $V_t$ , the function  $recog(w_i)$  yields top-10 words  $w'_i \in V_t$  from the word recognizer. This technique was motivated by the back-off language modeling [27]. Thus, wherever, the unigram is not found in the training model ( $T < 1$ ), we simply utilize the top-10 unigram outputs from the word recognizer for that word image.

## 2.8 Results

Figure 2.7.3 shows an example of sentence decoding. The matrix shows the correction model options for the “observed” output from the recognizer with the corresponding probabilities. The arrows show the decoding path in conjunction with the language model. The final output sentence is shown at the bottom of the figure.

Table 2.7.1 shows the accuracies of the recognizer on the test set before and after the correction. Before correction, both raw accuracies as well as accuracy with post processing using top-10 choices in conjunction with a trigram language model are reported. We observe that there is a considerable increase in the accuracy after the correction. Since, there has been no related work in the handwriting domain, the increase in accuracies can only be compared with methods from machine print OCRs. The relative improvement from our method is significantly superior to other reported methods [48, 69, 23].

# Chapter 3

## Lexicon Pruning

### 3.1 Effect of Lexicon Size on Accuracy

Handwriting Recognition (HR) has been successfully used in several applications such as postal address interpretation [64], bank check reading [21], and forms reading [16]. These applications are characterized by (small) fixed lexicons afforded by contextual knowledge. Recognition of unconstrained handwritten document images is difficult because of the diversity in writing styles, inconsistent spacing between words and lines, and uncertainty of the number of lines in a page as well as the number of words in a line [62]. Also, most current handwritten word recognition approaches depend on

the availability of a lexicon of words for matching, making the recognition accuracy dependent upon the size of the lexicon. Accuracy dramatically decreases with an increase in lexicon size and therefore domain constraints are important for solving the problem. In check processing, a small (33), limited lexicon is used, and with postal address domains the street names, city names and zip codes mutually constrain each other. However, for general handwritten manuscripts, neither of these constraints apply and therefore the lexicon is the entire English dictionary and therefore the accuracy of recognition is very low.

Figure 1.2.1 illustrates the results obtained by the same document analysis techniques on a page from Newton's notes. Word recognition for this example would involve the entire English dictionary. Perhaps, if the context were known, a limited but large dictionary could be used. For example, if the page is taken from Newton's notes on Physics, the glossary of terms from a Physics textbook can be used. We propose to use such context knowledge to improve recognition accuracies.

## 3.2 Lexicon Pruning

One of the elements that contributes more to the complexity of the recognition task is the size of the lexicon. Table 3.2.1 describes how the increase in the size of the lexicon

reduces the recognition performance of a handwriting recognition system developed at CEDAR [31]. One problem with large lexicons is the number of times that the observation sequence extracted from the input image has to be matched against the words in the lexicon. So, a more intuitive approach attempts to limit the number of words to be compared during the recognition.

Rank	Lex. size 10	Lex. size 100	Lex. size 1000
Top 1	96.80	88.23	68.70
Top 2	98.63	93.36	81.40
Top 20		98.40	97.60
Top 50			98.70

Table 3.2.1: Recognition Accuracy (%) for the system described in [31]

Whereas, this is straight forward in a task that is geared towards a particular application such as postal automation or bank check recognition, it is non-trivial in large vocabulary and unconstrained systems. Researchers have proposed methods that are based on word length or shape.

Short words can be easily distinguished from long words by comparing only their lengths. Thus, the length is one simple criterion for lexicon reduction. Kaufmann et al [28] use a length classifier to eliminate (from the lexicon) those models which



differ significantly from the unknown pattern in the number of symbols. For each model, a minimal and a maximal length are determined. Based on this range, a distance between a word and the model class is defined and used during the recognition process to select only the pertinent models. Kaltenmeier et al [26] use the word length information given by a statistical classifier adapted to features derived from Fourier descriptors for the outer contours to reduce the number of entries in vocabulary of city names. Koerich et al [4] used the length of the feature vector and the topology of the HMMs that model the characters to limit the length of the words dynamically in a level building framework. Kimura et al [33] estimate the length of the possible word candidates using the segments resulting from the segmentation of the word image. Powalka et al [52] estimate the length of cursive words based on the number of times an imaginary horizontal line drawn through the middle of the word intersects the trace of the pen in its densest area. A similar approach by Guillevic et al [17] estimates word length and reduces the lexicon size. The number of characters is estimated using the counts of stroke crossings within the main body of a word.

The shape of the words is another important cue about the length and the style of the words. Zimmerman and Mao [74] use key characters in conjunction with word length estimation to limit the size of the lexicon. They attempt to identify some key characters in cursive handwritten words, and use them to generate a search string.

Kaufmann et al [28] proposed a method of reducing the size of vocabulary based on the combination of four classifiers: a length classifier, the profile range, an average profile and a transition classifier. All the classifiers use as input the same feature vectors used by the recognition system. Seni et al [61] extract a structural description of a word and use it to derive a set of matchable words. This set consists of entries from the system lexicon that are similar in shape or structure to the input word. The set of matchable words forms the reduced lexicon that is employed during the recognition process. Madhvanath et al [57, 58] use holistic features based on a coarse representation of the word shape for pruning large lexicons . A survey of various lexicon reduction techniques can be found in [34].

Most of these techniques rely on features extracted from the word image which are susceptible to errors. In addition to this the coverage of the lexicon (percentage of test words in lexicon) tends to reduce drastically using these techniques. Most of the features extracted for word shape based reduction are non-trivial to extract and prone to errors. Also, in large vocabulary systems the number of words that have same the length or global shape tend to be large thus leading to ineffective reduction in the size of the lexicon.

### 3.3 Proposed Method

The task of word recognition is a pattern recognition problem that can be simply stated as ‘given a word image what is the best possible corresponding word in a given lexicon’. Table 3.2.1 describes the effect of the size of this lexicon on the accuracies of a word recognizer. Thus, if we are able to reduce the size of the lexicon, the accuracies can increase significantly. However, lexicon reduction cannot be achieved trivially. It is either a painstaking manual effort or can be automated if we can constrain the lexicon by the domain we are interested in. This is not only important for indexing of documents, it can be exploited in the reduction of the input lexicon to a word recognizer as well. Milewski and Govindaraju [43] have described preliminary work in this area where the domain is restricted to medical forms. Each medical form is associated with a ‘concept’ e.g. leg injury, heart attack etc. They have shown that given this information, restricting the lexicon by concepts increases the efficiency of the recognizer and hence the retrieval efficiency.

We have explored automatic creation of small and representative lexicons  $Lex_i$  for topics of interest  $T_i$  from a large vocabulary  $V$  (Figure 3.3.1) where  $|Lex_i| \ll |V|$ . Given the automatically detected topic from Figure 3.3.2, we can then use this smaller lexicon in the second phase of recognition instead of the large vocabulary in order to improve recognition.

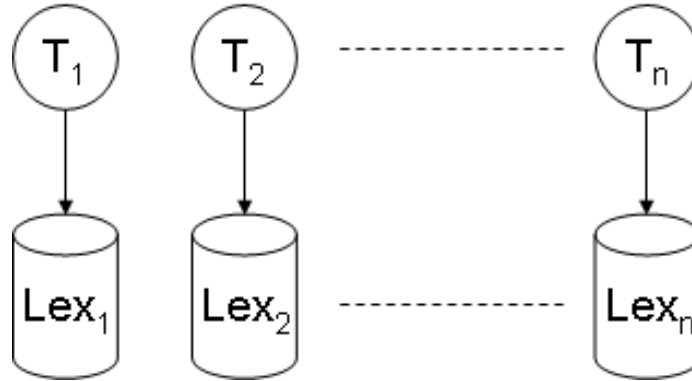


Figure 3.3.1: Generating topic specific lexicons.

Document categorization is a well researched area in the field of information retrieval when it comes to text categorization. Various topic models have been proposed that treat documents as bag of word models (SVD decomposition followed by cosine similarity measures [72], Naive Bayes [41] etc). For a given handwritten document to be categorized, the raw output of a recognizer is usually extremely noisy. Work presented in [51] demonstrated how the accuracy of document categorization reduces drastically as the accuracy of OCR reduces. This problem is further magnified in the case of handwriting recognition. For a large lexicon the word error rate for the top choice is generally very high. However, the error rates for the same lexicon when taking top N outputs of the recognizer into consideration are significantly lower than in the first case (see Table 3.3.1). Figure 3.3.3 shows a sentence from the IAM database with top-10 choices for each word from the word recognizer. We propose to utilize this fact and categorize a document based on n-best hypotheses for a word

Correct Choice in	Percentage
Top 1	32.22
Top 10	43.50
Top 20	65.89
Top 50	73.24

Table 3.3.1: Top- $N$  Recognition Accuracies of the Word Recognizer.

using a word recognizer.

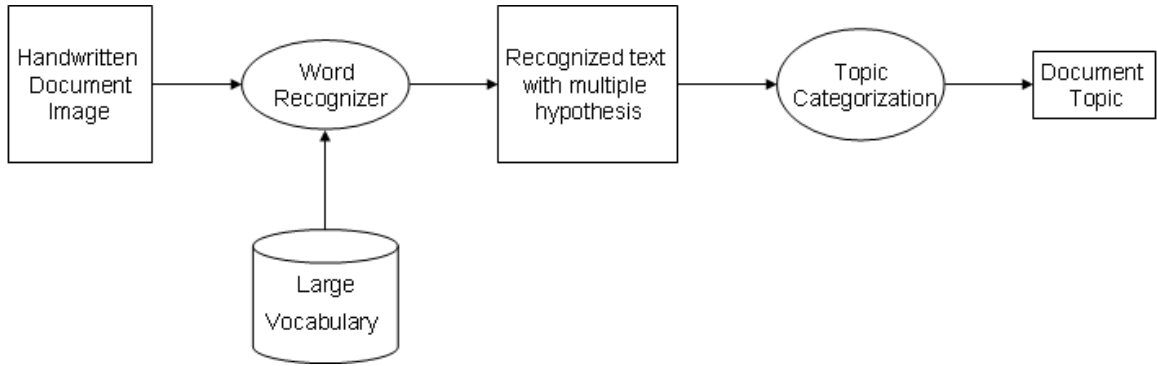


Figure 3.3.2: Topic Categorization in handwritten document images.

Assume that a document is related to a single topic. This assumption, can be relaxed later, however, it is valid for handwritten documents being considered. The document is composed (says) of a set of words  $W = \{w_1, w_2, \dots, w_n\}$ . Let us assume that we fix the documents to a set  $T = \{t_1, t_2, \dots, t_n\}$ . Then the probability of the topic given the words in the document can be estimated as  $P(T = t_i | w) = \frac{P(W|T=t_i)P(T=t_i)}{P(W)}$ . Whereas, in a system where we are dealing with text, the term  $P(W|T = t_i)$  can be estimated from the word occurrences of  $w_i$  given a topic labeled document. In our

1 MOVE to stop Mr. Gaitskell from nominating any more Labour life Peers

L	Move	Ve	Stop	Hr	Gaitskell	Arm	nominating	Dug	Nell	Labour	Use	Peers
A	Have	Po	Top	Mr	cranwell	Anna	Dominating	Dry	Note	Labora	We	Beers
x	Moke	Do	Tap	Air	farewell	Arne	Conducting	Cry	Cold	Labourer	Who	Years
b	Love	Leo	Stag	Dr	capsule	Mam	Commanding	Ous	Noel	Calais	Ass	Fears
d	Mont	No	Bag	Mrs	capable	Pyrne	Comrades	Any	Cork	Colour	Use	Piers
j	Most	As	Step	Her	tatsfield	Mono	Nomination	Buy	Nail	Latent	Us	Poets
h	Host	Co	Shop	Ris	fumbled	Mana	Amusing	Day	Cock	Lobby	Ago	Pass
K	Mole	Re	Toss	Nor	fumble	Ann	Longimanus	Airy	Well	Latitude	Af	Pests
So	Hole	Is	Stood	His	carroll	None	Reminding	Pug	Mod	Lambert	Act	Vary
he	Moves	To	Odd	No	funnel	from	Maudling	aw	Cool	Cables	Life	peas

Figure 3.3.3: Example of N-best choices from an HR for an IAM sentence.

case each  $w_i$  is a set of n-best hypotheses for each word image instead of being an atomic unit. The model can be thought of as a bag of words where each word instead of being one entry is a bunch of hypotheses with associated posterior probabilities. Thus, our model would have to consider a document as a bag of as many bags as there are words in the document. The size of each bag would be the value of ‘n’ as in the n-best hypotheses. Figure 3.3.2 illustrates the task of categorization of a handwritten document.

## 3.4 Models for Topic Categorization

Researchers have used techniques like Naïve Bayes, Maximum Entropy Models, Conditional Random Fields, Hidden Markov Models and more recently Latent Dirichlet Allocation for text classification. As described in Section 3.3, our task can be compared to classification with noisy data. We cannot directly apply these methods due to low accuracies and hence we have adapted the Naïve Bayes model and used the top-n information as a feature in the Maximum Entropy Model.

### 3.4.1 Naïve Bayes

Naïve Bayes classifier is a simple model for text classification that assumes independence amongst words in a give context. Because of the independence assumption the task of learning is greatly simplified. A document  $d_i$  is assigned to a category  $c_j$  under a model  $\theta$  using the formula

$$j_0 = \arg \max_{j=1\dots m} P(c_j|d_i, \theta) = \arg \max_{j=1\dots m} P(c_j|\theta)P(d_i|c_j, \theta)$$

The independence assumption dictates

$$P(d_i|c_j, \theta) = \prod_{k=1}^{|d_i|} P(w_k|c_j, \theta)$$

where  $|d_i|$  is the length of the document  $d_i$ .

If  $|D|$  denotes the total number of documents, then the class prior is given by

$$P(c_j|\theta) = \frac{\sum_{i=1}^{|D|} P(c_j|d_i)}{|D|}$$

McCallum and Nigam [41] describe how two variants of the Naïve Bayes could be utilized for text classification. The first variant only depends on the presence or absence of words and thus can be characterized by a distribution that is based on a *multi-variate Bernoulli*. This model however fails to capture the word counts and occurrences which is an important feature in the document topic. Also, in our case since the accuracies are low and we utilize top- $n$  ( $n = 10$  in our case) outputs from the recognizer, this model is not suited for our task. On the other hand the second variant of Naïve Bayes where the document is represented by the set of word occurrences captures this information. Here it can be characterized by a distribution that is a *multinomial*. This is similar to what can be referred to as the ‘*unigram language model*’. In this case the generative model is given by



$$P(w_t|c_j, \theta) = \frac{1 + \sum_{i=1}^{|D|} B_{it} P(c_j|d_i)}{2 + \sum_{i=1}^{|D|} P(c_j|d_i)}$$

where

$$B_{it} = \begin{cases} \sum P(im = w_t) & w_t \in d_i \\ 0 & w_t \notin d_i \end{cases}$$

Note that in the equation of  $B_{it}$ , we have replaced the count of the word by the probability  $\sum P(im = w_t)$ . This is the probability from the word recognizer for the image  $im$  being word  $w_t$ . This is what we will refer to as the *partial count* of the word. Thus, we were easily able to introduce the top-n choices from the word recognizer by replacing the count for that choice by its probability score.

### 3.4.2 Maximum Entropy

The Naïve Bayes however suffers from the independence assumption which may not be valid. This is often referred to as the bag-of-words model where the words are exchangeable and the topic does not depend on the mutual dependence of the words. To account for this we modified the Maximum Entropy Model. The motivating idea behind maximum entropy is the most uniform model should be preferred when no information is present and whatever information is present should constraint that

uniform model. Thus, the Maximum Entropy model prefers the most uniform model satisfying the constraint:

$$P(c|d) = \frac{\exp \sum_i \lambda_i f_i(d,c)}{\sum_c \exp \sum_i \lambda_i f_i(d,c)}$$

The  $f_i$ 's are any real-valued functions describing features of the (document,class) relationship. In order to include partial counts as in the Section 3.4.1, we used a real-valued feature function

$$f_{w,c'}(d, c) = \begin{cases} 0 & c \neq c', \\ \frac{N(d,w)}{|d|} & \text{otherwise} \end{cases}$$

where  $N(d, w)$  is the ‘partial’ count (probability from recognizer) of the word  $w$  in document  $d$ . With this representation, a word occurring often in one class will yield a higher weight for that word-class pair than for the other classes. In addition to this, it is easy to extend this notion to our partial counts by simply replacing the count of the words with the sum of the probabilities from the handwriting recognizer.

As is mentioned in Nigam et al [47], since the Maximum Likelihood training overfits, we use a Maximum Aposteriori training with Gaussian priors over feature functions. Just like in [47], we use a Gaussian prior with mean at zero and a diagonal covariance

matrix. The prior probability of the model is the product of Gaussians of each feature value  $\lambda_i$  with variance  $\sigma_i^2$ .

$$\prod_i \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{\frac{-\lambda_i^2}{2\sigma_i^2}}$$

The training has been implemented using the Improved Iterative Scaling (IIS) algorithm which is of the Quasi-Newton family of numerical algorithms. We used the Mallet (<http://mallet.cs.umass.edu>) package from University of Massachusetts at Amherst to train the model.

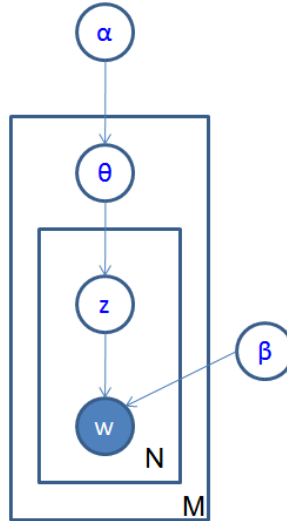


Figure 3.4.1: Plate representation of the LDA.

### 3.5 Latent Dirichlet Allocation

We now relax the condition of one topic per document to include a hierarchy of topics inside a document. *Latent Dirichlet Allocation* (LDA) [9] is a hierarchical probabilistic model that perceives a document corpus as being generated by some parameterized *random* process. The goal of this model is to capture the intra-document and inter-document statistical relationships. The formulation of the model is based on the intuitive notion that a document can be characterized as a set of topics and each word within the document belongs to one of these topics. The traditional approach has been to model the documents in a vector space wherein each document is a data point in the space and the words are the dimensions or the axes. The vectors are computed using the *tf-idf* scheme [60] where the normalized term *tf* signifies the importance of the word to the document in which it occurs based on frequency count and the normalized term *idf* signifies the importance of the word in the whole corpus based on the number of documents in which it occurs. The vector space model (VSM) has been used widely in search engines as it reduces the document to being represented by a set of words that can be easily queried. However, it suffers from the problems of *synonymy* and *polysemy*<sup>1</sup>. Moreover, it fails to capture any of the statistical features of the corpus except for the frequency count and hence cannot

---

<sup>1</sup>synonymy - different words conveying the same meaning; polysemy - same word conveying different meanings

be used for the classification task. In order to overcome these shortcomings, *latent semantic indexing* [13] was introduced. It is so called because the documents are represented not by words within the document but by the underlying *latent concepts* that aim to capture the semantic relations between the words. It involves decomposing the matrix  $A^T A$  (where  $A$  is a matrix consisting of document row vectors) using *singular value decomposition* and considering only the eigenvectors corresponding to the top few eigenvalues. Thus LSI is able to capture high correlations between the documents and the number of dimensions are considerably reduced.

### 3.5.1 LDA Model

LDA is essentially a hierarchical model similar to that of pLSI [19]. In this model, each document is modeled as a random finite mixture over a set of topics and each topic is modeled as a multinomial distribution over the words in the whole corpus. A way to compare LDA with LSI and pLSI is that the documents are represented by topic probabilities in LDA, by a reduced dimensional matrix in LSI and by the list of mixing proportions in pLSI. A *document* is a sequence of  $N$  words denoted by  $W = (w_1, w_2, \dots, w_N)$  where  $w_N$  is the  $n^{th}$  word in the sequence. A *corpus* is a collection of  $M$  documents denoted by  $\mathcal{D} = \{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_M\}$ . The words from a vocabulary are indexed by  $1 \dots V$ . The topics are denoted by  $z$  and there are  $k$  topics

within a corpus.

In LDA, each document is characterized by a random mixture of latent topics and each topic is characterized by a distribution over words. The generative process starts by choosing a document  $\mathbf{w}$  in a corpus  $\mathcal{D}$ . For this document, choose  $\boldsymbol{\theta} \sim \text{Dir}(\boldsymbol{\alpha})$ . For each of the  $N$  words, choose a topic  $z_n \sim \text{Multinomial}(\boldsymbol{\theta})$  and then choose a word  $w_n$  from  $p(w_n|z_n, \boldsymbol{\beta})$ . Here  $p(w_n|z_n, \boldsymbol{\beta})$  is a multinomial probability conditioned on the topic  $z_n$ .

The random variable  $\boldsymbol{\theta}$  has  $k$ -dimensions and has the following probability density on a  $(k-1)$  simplex :

$$(3.5.1) \quad p(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1}$$

where  $\boldsymbol{\alpha}$  is a  $k$ -vector with components  $\alpha_i > 0$  and  $\Gamma(x)$  is the Gamma function. The Dirichlet distribution is chosen as it is the conjugate prior of the multinomial distribution. Given the parameters of the model, the joint distribution of topic mixture  $\boldsymbol{\theta}$ , a set of  $N$  topics  $\mathbf{z}$  and a set of  $N$  words  $\mathbf{w}$  is as follows:

$$(3.5.2) \quad p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = p(\boldsymbol{\theta}|\boldsymbol{\alpha}) \prod_{n=1}^N p(z_n|\boldsymbol{\theta}) p(w_n|z_n, \boldsymbol{\beta})$$

where  $p(z_n|\boldsymbol{\theta})$  is  $\theta_i$  for the unique  $i$  such that  $z_n^i = 1$ . The marginal distribution of a document is:

$$(3.5.3) \quad p(\mathbf{w}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int p(\boldsymbol{\theta}|\boldsymbol{\alpha}) \left( \prod_{n=1}^N \sum_{z_n} p(z_n|\boldsymbol{\theta}) p(w_n|z_n, \boldsymbol{\beta}) \right) d\boldsymbol{\theta}$$

The probability of the corpus is obtained by taking the product of the marginal probabilities of all the documents :

$$(3.5.4) \quad p(\mathcal{D}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int p(\boldsymbol{\theta}_d|\boldsymbol{\alpha}) \left( \prod_{n=1}^{N_d} \sum_{z_{dn}} p(z_{dn}|\boldsymbol{\theta}_d) p(w_{dn}|z_{dn}, \boldsymbol{\beta}) \right) d\boldsymbol{\theta}_d$$

The parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\alpha}$  are corpus level parameters and are sampled only once in the entire process. The parameter  $\boldsymbol{\theta}_d$  is sampled once for each document. The variables  $z_{dn}$  and  $w_{dn}$  are sampled once for each word in each document.

The equation 3.5.3 can also be explained by the theorem of *exchangeability*. For a frequentist, the observations are always independent whereas a probabilist will consider the observations independent conditioned on some hidden variable and then the idea would be to find out the conditional probability distribution of this hidden quantity. Hence in LDA, the topic  $z$  is the hidden quantity and we have to compute the posterior distribution of topics.

### 3.5.2 Inference and Parameter Estimation

The inference problem is to evaluate the posterior distribution of the hidden variables given a document to find the topic to which the document belongs (similar to pLSI):

$$(3.5.5) \quad p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{p(\boldsymbol{\theta}, \mathbf{z}, \mathbf{w} | \boldsymbol{\alpha}, \boldsymbol{\beta})}{p(\mathbf{w} | \boldsymbol{\alpha}, \boldsymbol{\beta})}$$

Equation 3.5.3 can be re-written as:

$$(3.5.6) \quad p(\mathbf{w} | \boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \int \left( \prod_{i=1}^k \theta_i^{\alpha_i - 1} \right) \left( \prod_{n=1}^N \sum_{i=1}^k \prod_{j=1}^V (\theta_i \beta_{ij})^{w_n^j} \right) d\boldsymbol{\theta}$$

The above equation is intractable due to the coupling between  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}$ . Blei et al [8] discuss variational approximation and Gibbs sampling to compute the posterior distribution  $p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{w}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ .



## 3.6 Experiments

### 3.6.1 Dataset

For the purpose of our experiments we have used the publicly available dataset called the IAM database [40]. It has the following characteristics

- 657 writers contributed samples of their handwriting
- 1'539 pages of scanned text
- 5'685 isolated and labeled sentences
- 13'353 isolated and labeled text lines
- 115'320 isolated and labeled words
- ~15k lexicon entries

There are 13 topic categories ranging from Press to Religion as listed in Table 3.6.1. More details on the database can be found in [40]. This dataset is unconstrained (UK) English and the pages are classified according to the topics. Figure 3.6.1 shows an image from the dataset.

S. No.	Label	Topic Category
1.	A	Press:reportage
2.	B	Press:editorial
3.	C	Press:reviews
4.	D	Religion
5.	E	Skills, trades and hobbies
6.	F	Popular lore
7.	G	Biographies, essays
8.	H	Miscellaneous
9.	J	Learned and scientific writings
10.	K	General fiction
11.	L	Mystery
12.	M	Science fiction
13.	N	Adventure

Table 3.6.1: Topic Categories of the IAM database.

### 3.6.2 Results

We divided each category into a training set (60%) and trained our models on this data with the above target categories. The test set (40%) of the data was used to evaluate the classification performance. Figure 3.6.2 shows the performance of the classification algorithms. Taking top-5 topics from the classification results performs similar to classification on the actual transcript of the document.

After classification, the lexicon was reduced based on the highest mutual information of the words with the topic and then recognition was performed again with the new set of words. The immediate advantage of the method is the availability of the lexicon,

	Top(%)	Top-10(%)
Before Reduction	32.22	43.4
After Reduction	40.01	69.1

Table 3.6.2: Recognition Accuracies before and after lexicon reduction.

since it is a by-product of the classification model. On average the size of the lexicon lies between 1.5 – 2k which is an average reduction of about 85%. The recognition accuracies before and after lexicon reduction by topic classification are reported in Table 3.6.2.

### 3.7 Combination of Recognizers

It was observed that the HMM recognizer performed well on the large lexicons (in our experiment) as compared to the WMR. In addition, it is commonly observed that the WMR outperforms the HMM for smaller lexicons. Though this observation may not be a general claim, it motivated us to combine the two systems. In this experimental setup, we reduced the lexicon by topic categorization (as described earlier) using the HMM. After the topic categorization and the selection of the corresponding lexicons, we performed the second phase of recognition with the WMR (instead of the HMM). The encouraging results in Table 3.7.1 suggest that indeed the strength of the recognizers can be exploited for improving the accuracies even further.

	Top(%)	Top-10(%)
Before Reduction	32.22	43.4
After Reduction	40.01	69.1
After Reduction WMR	51.00	74.1

Table 3.7.1: Recognition Accuracies before and after lexicon reduction using combination of two HR systems.

## 3.8 Conclusions

In this chapter we have presented a novel technique of lexicon reduction for handwriting recognition using topic categorization. We improved the accuracy of the recognizer significantly. The methods are statistical and trainable. In addition the computational overhead is minimal. Another major contribution of this technique is that the lexicon is available readily after training and there is no need for manual creation of topic-wise lexicons. As noted, this method is statistical and can also be applied to multiple languages and domains. This area of research has also been addressed in text classification and our future plan is to adapt our techniques for this domain.

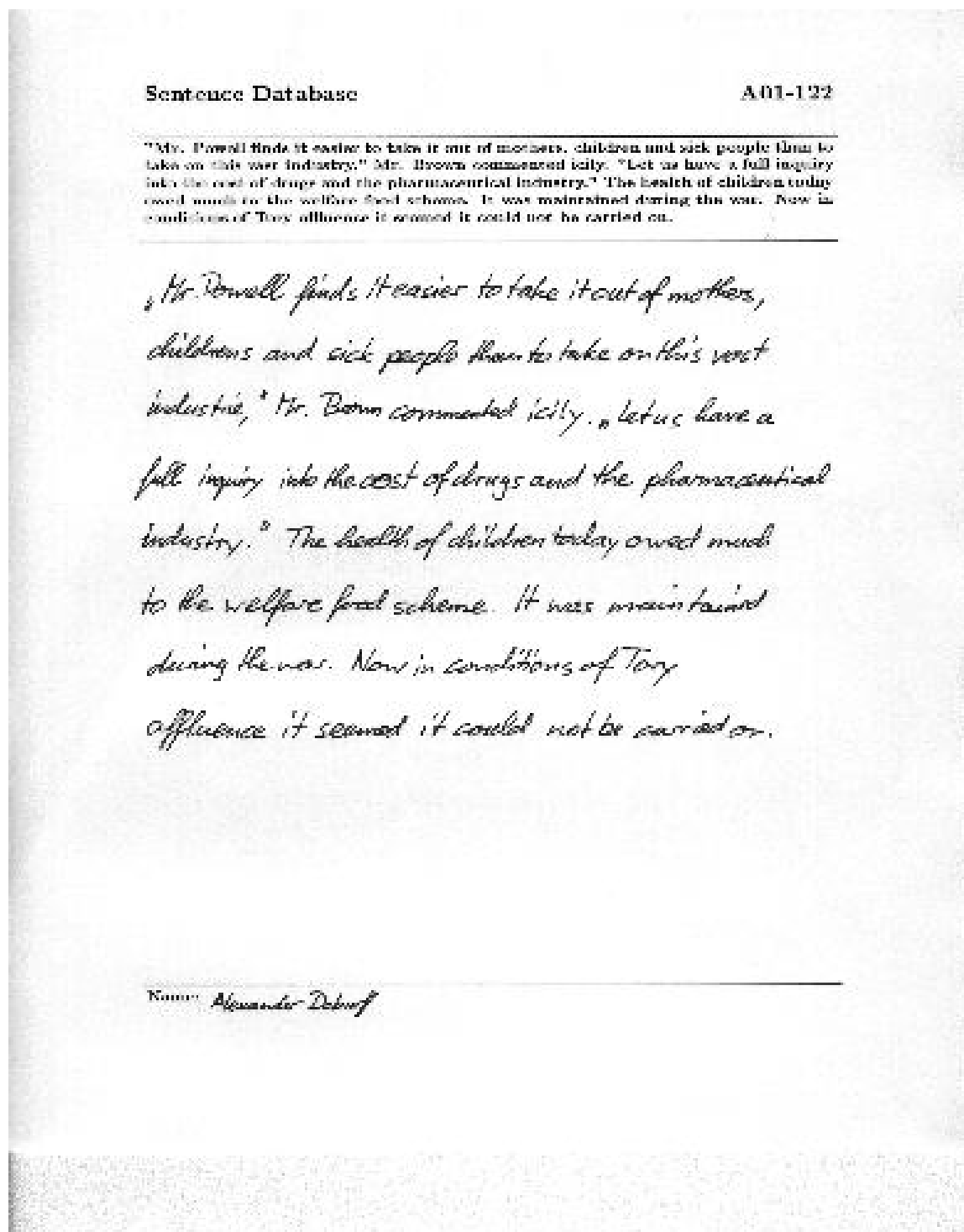


Figure 3.6.1: Example image from IAM Database.

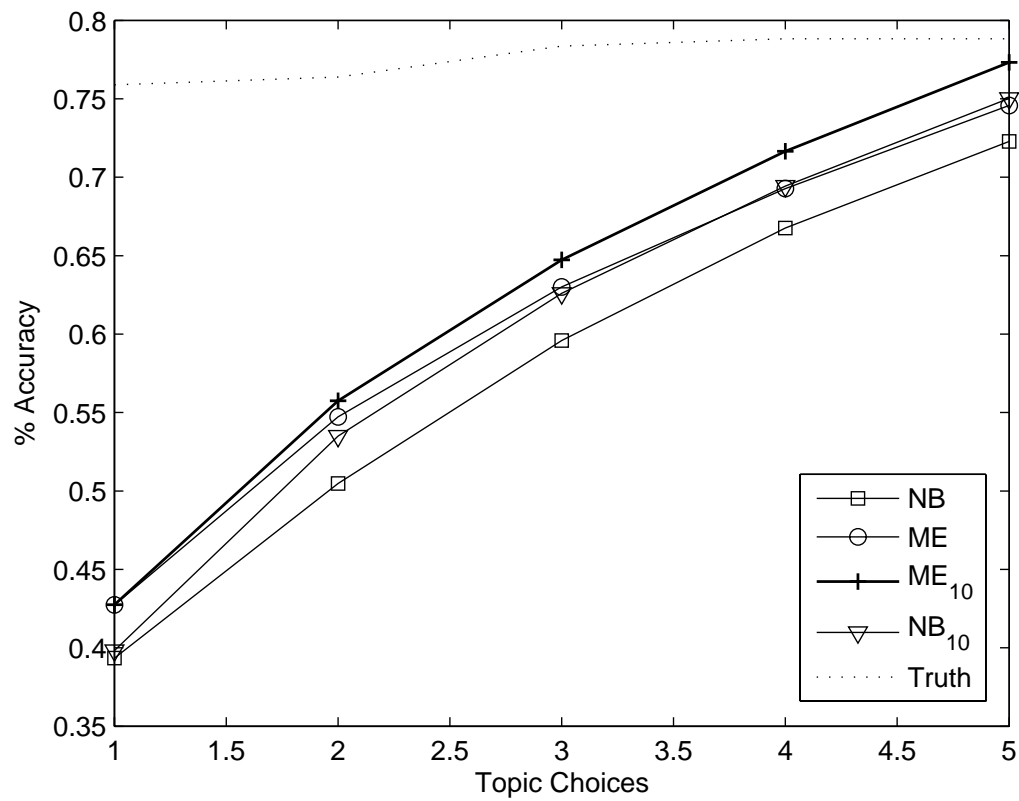


Figure 3.6.2: Topic Classification Results (NB - Naives Bayes, ME - MaxEnt, suffix 10 implies top-10 choices for each word image).

# Chapter 4

## Information Retrieval

### 4.1 Introduction

The problem of indexing and retrieving documents with low recognition rate has recently been addressed by researchers. In [44, 49, 24] different approaches to modeling typical recognition errors were proposed. In [44], 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 [49]. In [24], a language model that took common recognition

errors into account was built. This language model can 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 [39, 54, 73]. 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) [55]. They address the “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 an 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 in 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 4.2. Results of IR tests are presented and discussed in section 4.3. Finally, section 4.4 concludes the chapter.



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 health care 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) [31] is as low as around 20% as reported in [42].

## 4.2 Multiple-candidate Vector Model

### 4.2.1 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 one hand, as results in table 4.2.1 suggest, the use of *top n* ( $n > 1$ ) candidates in terms of word recognition ranks will improve the count of occurrences.

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

Table 4.2.1: Recognition rates of PCR forms for different  $n$ .

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.

### 4.2.2 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

$$\begin{aligned}
 \mathbf{tf}_{i,j}^{\mathbf{ocr}} &= E\{tf_{i,j}|O\} \\
 &= E\left\{\frac{freq_{i,j}|O}{\sum_l freq_{l,j}}\right\} \\
 &= \frac{E\{freq_{i,j}|O\}}{\sum_l freq_{l,j}}
 \end{aligned} \tag{4.2.1}$$

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  $\sum_l freq_{l,j}$  stands for the total number of terms in document  $j$ .  $\sum_l freq_{l,j}$  is a normalization factor and can be estimated by counting the total number of word image. Let the set of words in a document  $d$  be denoted by  $w(d)$ , then

$$E\{freq_{i,j}|O\} = \sum_{w \in w(d_j)} P(w = t_i) \quad (4.2.2)$$

Therefore,

$$tf_{i,j}^{ocr} = E\{tf_{i,j}|O\} = \frac{\sum_{w \in w(d_j)} P(w = t_i)}{\sum_l freq_{l,j}} \quad (4.2.3)$$

Equation 4.2.3 is an analytical expression of  $tf_{i,j}^{ocr}$ .

**Theorem 4.1** *Given probabilities  $P(w = t_i)$ ,  $\forall t_i$ , we will be able to calculate the **OCR-tf** defined as the expectation  $E\{tf_{i,j}|O\}$  of the term frequency in the ground truth of document  $d_j$ .*

However, most word recognizers do not 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 naive normalized probability distribution by fitting to a Zipfian distribution was presented in [20]. We present a novel method based on the following motivation. Given an image 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  $rank(w_i) = i (1 \leq i \leq n)$ , we assume [18],  $P(w = w_i | w \in L)$  is an exponential function of  $i$ . 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 divided by the total number of words recognized during the test.

While calculating OCR-TF using equation 4.2.3 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  $c_t$  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  $P^{(i)} \geq c_t$  are taken into account.

### 4.3 Test results and discussions

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

Physician: Kalra, Tejinder

CARE IN PROGRESS ON ARRIVAL: None Citizen PEPD/Other First Responder Other EMS PAD used

MECHANISM OF INJURY: MVA (✓) seat belt used Unarmed assault GSW Machinery Knife

CHIEF COMPLAINT: dizziness

SUBJECTIVE ASSESSMENT: at ground lying supine in bed at 1:22 AM x 3 years, brought to MRI for up to 10 min to 20 min, moved to stretcher via short cut

PRESENTING PROBLEM: Allergic Reaction Unconscious/Unresp. Shock Major Trauma OB/GYN Syncope Seizure Head Injury Trauma-Blunt Burns Stroke/CVA Behavioral Disorder Spinal Injury Trauma-Penetrating Environmental General Illness/Malaise Substance Abuse (Potential) Fracture/Dislocation Soft Tissue Injury Heat Gastro-Intestinal Distress Poisoning (Accidental) Amputation Bleeding/Hemorrhage Cold Cardiac Related (Potential) Diabetic Related (Potential) Cardiac Arrest Pain Other BLS MRI tx Hazardous Materials Obvious Death

PAST MEDICAL HISTORY: None Allergy to Codeine Hypertension Stroke Seizures Diabetes COPD Cardiac Other (List) Asthma

CURRENT MEDICATIONS (List):

VITALS	TIME	RESP	PULSE	B.P.	LEVEL OF CONSCIOUSNESS	GCS	R	PUPILS	L	SKIN	STATUS
SIGN	10:44	Rate: 18 Regular Shallow Labored	Rate: 64 Regular Irregular	130/70	Alert Voice Pain Unresp.	15		Normal Dilated Constricted Sluggish No-Reaction		Unremarkable Cool Pale Warm Cyanotic Moist Flushed Dry Jaundiced	C U P S
	2:29	Rate: 18 Regular Shallow Labored	Rate: 60 Regular Irregular	130/70	Alert Voice Pain Unresp.	15		Normal Dilated Constricted Sluggish No-Reaction		Unremarkable Cool Pale Warm Cyanotic Moist Flushed Dry Jaundiced	C U P S
	1:11	Rate: Regular Shallow Labored	Rate: Regular Irregular		Alert Voice Pain Unresp.			Normal Dilated Constricted Sluggish No-Reaction		Unremarkable Cool Pale Warm Cyanotic Moist Flushed Dry Jaundiced	C U P S

OBJECTIVE PHYSICAL ASSESSMENT: 69 y.o. F. at 10:44 AM, dizziness 3 years, HEENT normal, Pupils - PEARL, Neck - CVD, chest - clear, abd - soft, PMS - intact, ext -

COMMENTS: Main = Safety, Comfort, BLS Care

MEDICAL REASON FOR AMBULANCE TRANSPORT: MRI

REPORT GIVEN TO:

Figure 4.3.1: An example of NYS Pre-hospital Care Report form.

forms to these queries. An example of the handwritten region in a PCR form is shown in figure 4.3.1. Of all the 21 queries, 13 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 [42]. Subsequently, all the word images are extracted and fed into the word recognizer [31] with a lexicon of 4405 English words. The top 1, 2, 5, and 10 word recognition rates are shown in table 4.2.1. 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

$$\begin{cases} \frac{1}{S}, 1 \leq i \leq S \\ 0, \text{otherwise} \end{cases}$$

where  $S$  is a truncation threshold.

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

Table 4.3.1: MAP values of retrievals using different *tf-idf* schemes (with feedback and  $c_t=0.1$ ).

The results of interpolated average precisions at 11 different recall values are compared in Figure 4.3.3. The Mean Average Precisions (MAP) and R-Precisions of the results of ordinary recognition text, multi-candidate text with naive *tf-idf* scheme and

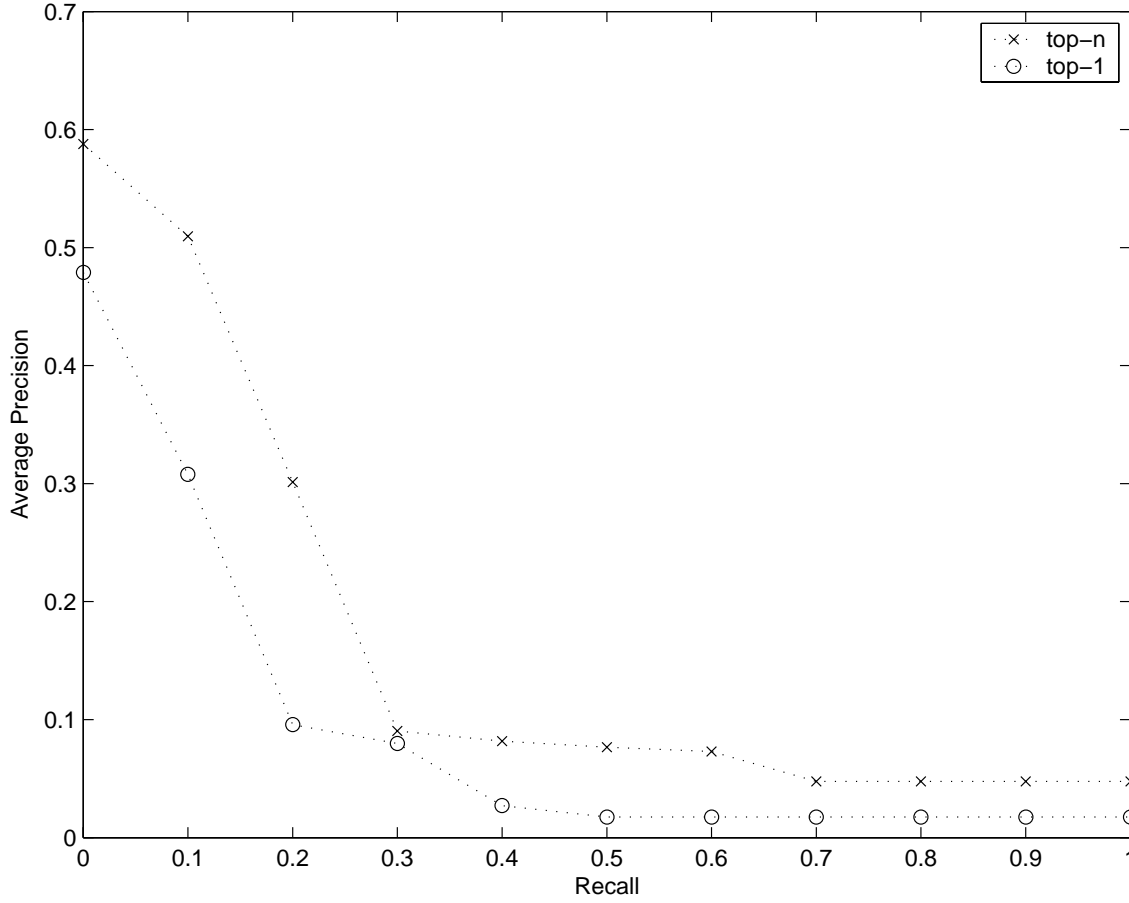


Figure 4.3.2: Performances without pseudo relevance feedback.

the proposed TF-IDF scheme are compared in table 4.3. The performances shown in figure 4.3.3 and table 4.3 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 4.3 that the scheme of exponential probability is superior the scheme of assigning equal weight to each candidate.



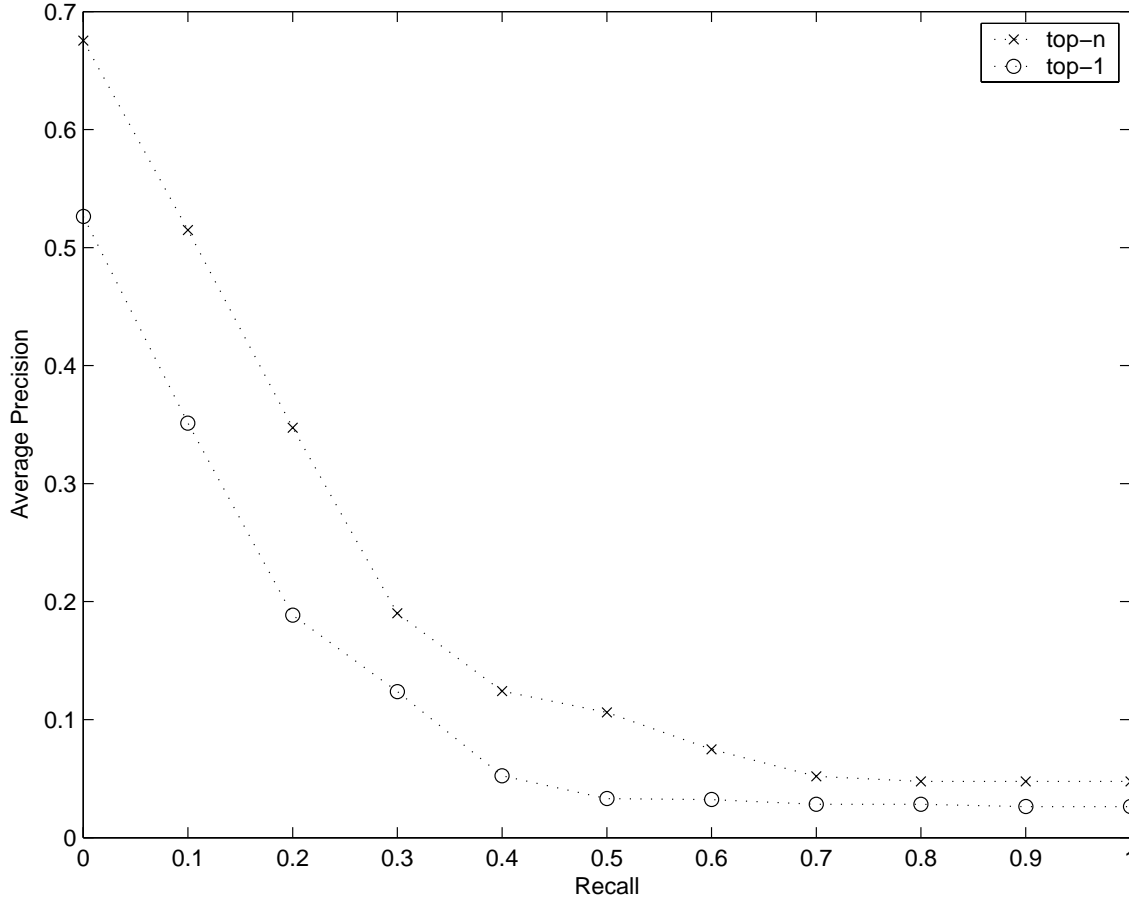


Figure 4.3.3: Performances with pseudo relevance feedback.

The effect of choosing different values is tested and the result is shown in figure 4.3.4.

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 sufficient top word candidates, whereas too small  $c_t$  value seems to involve too many irrelevant candidates which lead to poor performance.

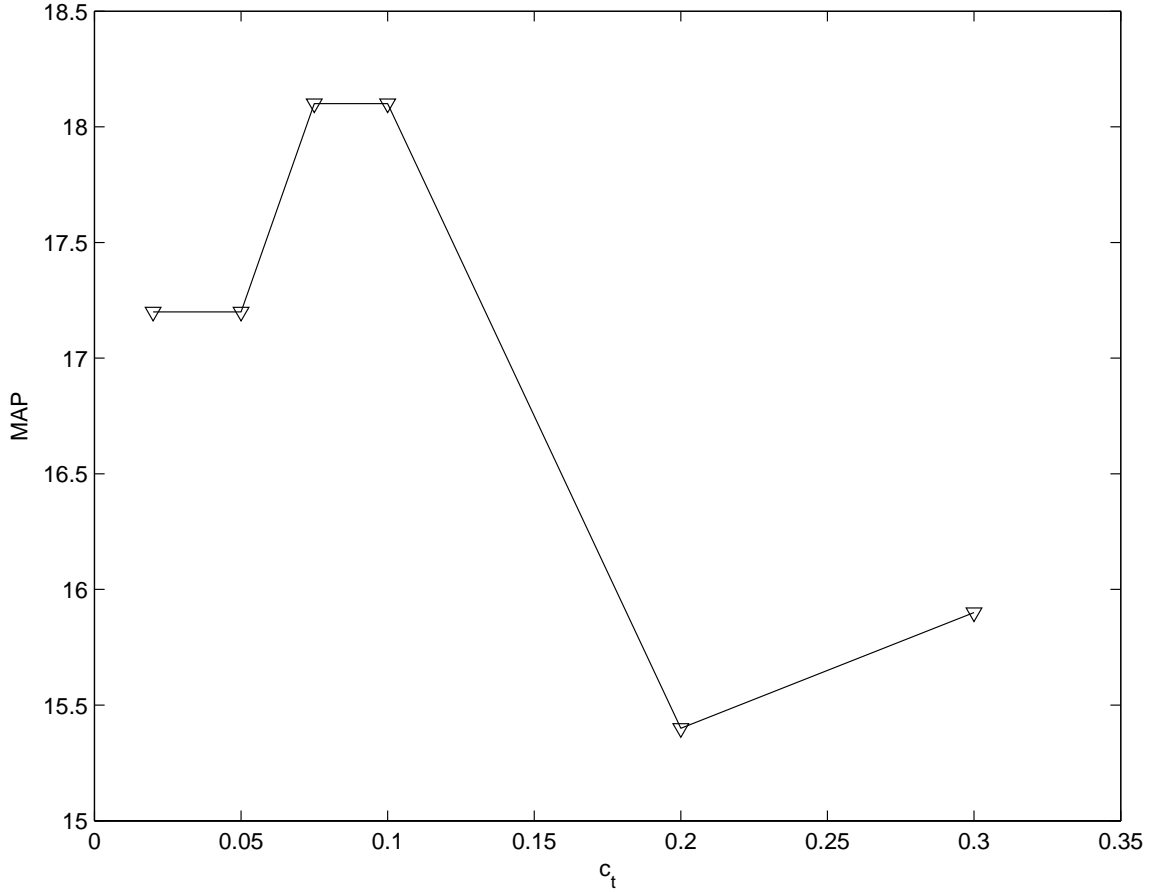


Figure 4.3.4: MAP values using different values of  $c_t \in \{0.02, 0.05, 0.075, 0.1, 0.2, 0.3\}$  are tested.

## 4.4 Overall System for Information Retrieval

We have described 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*. Our approach improved the performance of IR as compared to the method using standard Vector Model. We now include the topic

categorization results that not only improve the HR results and the retrieval performance, but can also be readily applied for query expansion. In the retrieval model we modify the convolution model as proposed by Franz et al in [15] in which they describe a probabilistic convolution model in which the probability of generating an English query stem was modeled based on the probabilities of generating that stem from an Arabic word or morpheme observed in the document. The convolution model substantially outperformed the sentence translation model (i.e. only the top choice words across a sentence), perhaps because it made use of document-wide translation probabilities. For the purpose of handwritten documents, in order to estimate generating a query  $q$  for a document  $D$ , we assume it to be

$$p(q|D) = \prod_i (\alpha \times p_0(q_i) + (1 - \alpha) \times \sum_j p(q_i|w_j)p(w_j|D))$$

where  $q_i$  are the query words,  $p_0(q_i)$  is the language model probability (possibly unigram counts),  $p(q_i|w_j)$  is the probability of recognition of image  $w_j$  as word  $q_i$  and  $p(w_j|D)$  is the probability of finding words  $w_j$  in the document  $D$ .

We used the IAM database in the same setup as in Chapter 3. The experiment is conducted as follows

- Formulate 14 queries each consisting of 3 – 6 words relevant to the complete

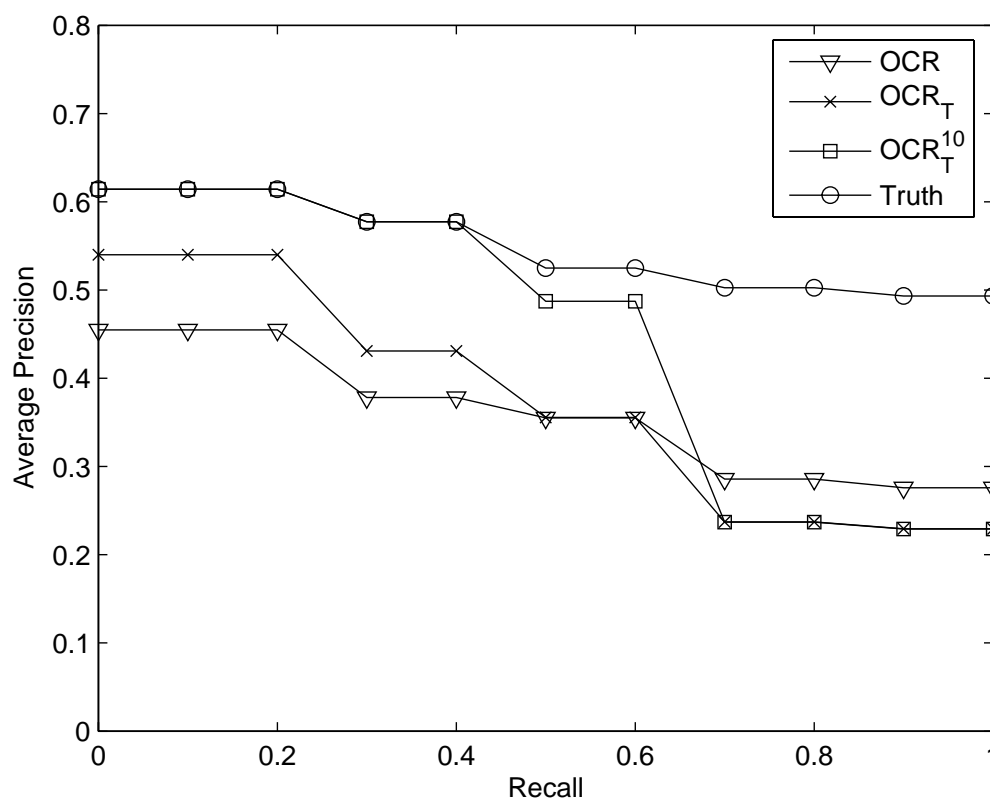


Figure 4.4.1: Precision-Recall comparison for the overall IR systems.

IAM database.

- Select documents from IAM database from the test set for the topic categorization task.
- Run the handwriting recognition system on these documents and get the output - *OCR*.

	<b>Truth</b>	$OCR$	$OCR_T$	$OCR_T^{10}$
MAP	<b>0.48</b>	0.31	0.33	0.40
R-Precision	<b>0.42</b>	0.22	0.31	0.38

Table 4.4.1: Retrieval performance for the different configurations

- Automatically categorize the documents; re-run HR and get the top-choice output -  $OCR_T$  and top-ten choice output -  $OCR_T^{10}$ .
- Generate the 10-pt. precision recall curves for the truth,  $OCR$ ,  $OCR_T$  and  $OCR_T^{10}$ .

Figure 4.4.1 shows the 10pt. precision-recall curves for the different settings. As we note that the final system that includes all contributions of our thesis outperforms the basic system ( $OCR$ ).

Table 4.4.1 lists the mean average precision and the R-precision values for the different settings. Our retrieval system performs at par with the retrieval system on the document transcripts (truth). The low MAP score, even for the truth, is because in the spirit of a fair evaluation, queries were formulated for the full IAM database where as the experiment was conducted only on the test set of the categorization experiment. Thus, not all queries returned relevant documents.

# Chapter 5

## Conclusions

### 5.1 Summary

In this thesis we have described how existing handwriting recognizers can be utilized for an information retrieval system for handwritten documents. Our main objective is to advance handwriting recognition and retrieval technology so that handwritten documents can be accessed with the same ease as printed or ASCII documents. Our approach has been to use the existing recognition systems as black-boxes while exploiting their capabilities to the fullest. This is particularly advantageous in scenarios

where the components of a system are developed by different parties and their internals are not accessible. This also allowed us to focus on the Information Retrieval (IR) aspects while not getting bogged down by the other necessary modules of a complete system.

The important aspects of IR we have explored in this thesis are: topic categorization, query pre- and post-processing, combining black boxes outputs, and merging retrieval results. These areas have been thus far been neglected by the HR research community for the. Our motivation stems from the encouraging results obtained when dealing primarily with machine printed documents.

We have identified the crux of the problem as being the task of improving accuracies of a handwriting recognition system. In order to address the improvement of HR, we have developed an innovative statistical phrase-based noisy channel model to correct the errors from a handwriting recognition system. It is common knowledge that the accuracy of HR is inversely proportional to the size of the lexicon. Therefore we developed a method for pruning the lexicon based on the topic categories of the handwritten documents and subsequently using a topic specific lexicon. Topic categorization provides an explicit context and reduces the search space as certain words are more likely to be used in a given topic. We have been able establish this as a promising line of research that must be further explored.

Another key idea we have developed is the notion that concise statistical information regarding a word image is available in the top  $n$  choices. For example, current recognizers perform at the top  $n$  recognition rate of about 30% for  $n = 1$  and 70% for  $n = 50$  when presented with a large lexicon. Building upon this idea, we developed a framework for convolution and relevance based retrieval models that operate on the raw recognition results captured in the top- $n$  choices. We modified the classical Vector Space Model to accommodate these top- $n$  choices and achieved improved retrieval performance as vindicated by the increase in the mean average precision. The models are robust to machine transcription errors, allowing retrieval of documents that would be typically missed by traditional transcription-based retrieval systems.

## 5.2 Contributions

The major contributions of this dissertation are in the field of retrieval of handwritten documents. We have advanced the state-of-the-art in the following areas:

- Improving Handwriting Recognition performance using an innovative noisy channel model
- Developing statistical, trainable, and scalable techniques for lexicon reduction by Topic Categorization



- Developing a methodology for utilizing top-n outputs and partial counts using convolution-based methods
- Modifying the Vector Space Model to allow Query Expansion by topic

### 5.3 Future Work

The task of retrieval from handwritten documents is still an open area of research. Further advances can be made in improving handwriting recognizers by modeling the behavior of different writers. For example, a certain writer has a particular style of writing characters e.g. loopy 'f', elongated 'p' etc. Figure 5.3.1 shows different word images from Washington's manuscripts. As an example note the style of writing the letter 'p'. These facts can be modeled using the methodology described in sections 2.2.1 and recognition accuracies improved.

As we note in section 3.6.2, the categorization performance on the OCR output reduces drastically when compared to the truth. Whereas, in this dissertation, we made advancements by exploiting top-n hypothesis and features that are suited for the HR task, topic categorization on noisy text still remains an open area. Due to lack of sufficient labeled data, the lexicon pruning experiment on the medical forms could not be performed. However, our method is very general and can be readily

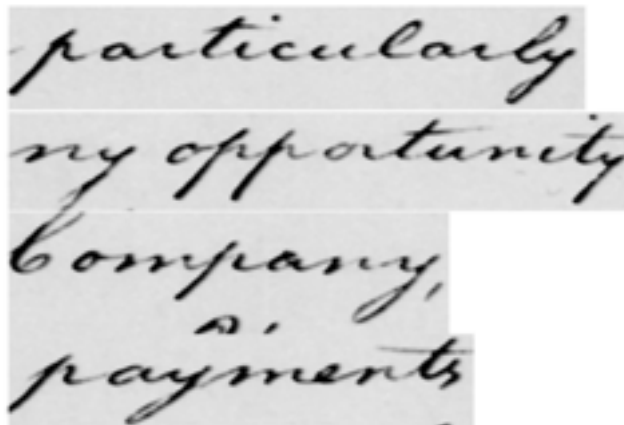


Figure 5.3.1: Written examples from George Washington. Worth noticing is the style of writing ‘*p*’

applied to this dataset as well [43].

# Bibliography

- [1] Arabic gigaword corpus. In *Linguistic Data Consortium (LDC)*, page Catalog LDC2004T17, Philadelphia, USA, 2003.
- [2] Mt'03 arabic to english mt evaluation set. In *Linguistic Data Consortium (LDC)*, page Catalog LDC2003T12, Philadelphia, USA, 2003.
- [3] Arabic news translation text. In *Linguistic Data Consortium (LDC)*, page Catalog LDC2003T12, Philadelphia, USA, 2005.
- [4] K. AL, S. R, and S. CY. A time-length constrained level building algorithm for large vocabulary handwritten word recognition. In *Proceedings of the 2nd International Conference on Advances in Pattern Recognition*, pages 127–136, 2001.
- [5] N. G. B. Zhao and K. Papineni. Inner-outer bracket models for word alignments using hidden blocks. In *Proceedings of the Human Language Technology*

- Conference (HLT-NAACL)*, Vancouver, Canada, 2005.
- [6] I. Bazzi, R. Schwartz, and J. Makhoul. An omnifont open-vocabulary ocr system for english and arabic. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 21(6):495–504, Jun 1999.
- [7] S. M. Beitzel, E. C. Jensen, and D. 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.
- [8] D. Blei, A. Ng, and M. Jordan. Latent dirichlet allocation. *Journal of Machine Learning Research*, (3):993–1022, 2003.
- [9] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent Dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022, 2003.
- [10] H. Blum. A transformation for extracting new descriptors of shape. In *Models for the perception of Speech and Visual Form*, MIT Press, pages 362–380, Cambridge, MA, 1967.
- [11] D. Bouchaffra and V. Govindaraju. A methodology for mapping scores to probabilities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(9):923–927, 1999.

- [12] 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.
- [13] S. C. Deerwester, S. T. Dumais, T. K. Landauer, G. W. Furnas, and R. A. Harshman. Indexing by latent semantic analysis. *Journal of the American Society of Information Science*, 41(6):391–407, 1990.
- [14] D. Doermann. The Indexing and Retrieval of Document Images: A Survey. Technical Report LAMP-TR-013,CFAR-TR-878,CS-TR-3876, University of Maryland, College Park, February 1998.
- [15] M. Franz, J. S. McCarley, and R. Ward. Adhoc, cross language and spoken document information retrieval. In *Proceedings of the eighth text retrieval conference (TREC8)*, page 391, 2000.
- [16] V. Govindaraju, V. Ramanaprasad, D. Lee, and S. Srihari. Reading handwritten us census forms. In *Proceedings of Third International Conference on Document Analysis and Recognition*, pages 82–85, Montreal, Canada, 1997.
- [17] N. D. Guillevis D and Y. K. Word lexicon reduction by character spotting. *Frontiers in Handwriting Recognition, 2000., Proceedings of the Seventh International Workshop on*, pages 373–382, 2000.

- [18] F. F. H. Cao and V. Govindaraju. Indexing and retrieval of degraded handwritten document. In *IJCAI Workshop on Multimodal Information Retrieval and Applications*, 2007.
- [19] T. Hofmann. Probabilistic Latent Semantic Indexing. In *Proceedings of the 22nd Annual ACM Conference on Research and Development in Information Retrieval*, pages 50–57, Berkeley, California, August 1999.
- [20] 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.
- [21] S. Impedovo, P. Wang, and H. Bunke. Automatic bankcheck processing. *Machine Perception and Artificial Intelligence*, 28, 1997.
- [22] A. Ittycheriah and S. Roukos. A maximum entropy word aligner for arabic-english machine translation. In *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, Vancouver, Canada, 2005.
- [23] J. A. J.C. Perez-Cortes, J.C. Amerngual and R. Llobet. Stochastic error-correcting parsing for ocr post-processing. In *International Conference on Pattern Recognition*, pages 4405–4408, Barcelona, Spain, 2000.
- [24] H. Jing. Using hidden markov modeling to decompose human-written summaries. *Computational Linguistics*, 28(4):527–543, 2002.

- [25] T. W. K. Papineni, S. Roukos and W. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of Association of Computational Linguists (ACL)*, pages 311–318, Philadelphia, USA, 2002.
- [26] A. Kaltenmeier, T. Caesar, J. Gloger, and E. Mandler. Sophisticated topology of hidden markov models for cursive script recognition. *Document Analysis and Recognition, 1993., Proceedings of the Second International Conference on*, pages 139–142, 20-22 Oct 1993.
- [27] S. M. Katz. Estimation of probabilities from sparse data for the language model component of a speech recogniser. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 35(3):400–401, 1987.
- [28] G. Kaufmann, H. Bunke, and M. Hadorn. Lexicon reduction in an hmm-framework based on quantized feature vectors. In *ICDAR '97: Proceedings of the 4th International Conference on Document Analysis and Recognition*, pages 1097–1101, Washington, DC, USA, 1997. IEEE Computer Society.
- [29] H. G. P. Keaton and R. Goodman. Keyword spotting for cursive document retrieval. In *Proceedings of the IEEE Workshop on Document Image Analysis*, San Juan, Puerto Rico, 1997.

- [30] S. Khoubyari and J. J.Hull. Keyword location in noisy document image. In *2nd Annual Symposium on Document Analysis and Information Retrieval*, pages 217–231, UNLV, Las Vegas, 1993.
- [31] 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(4):366–379, 1997.
- [32] G. Kim, V. Govindaraju, and S. Srihari. Architecture for handwriting recognition systems. *International Journal of Document Analysis and Recognition*, 2(1):37–44, 1999.
- [33] F. Kimura, M. Shridhar, and Z. Chen. Improvements of a lexicon directed algorithm for recognition of unconstrained handwritten words. *Document Analysis and Recognition, 1993., Proceedings of the Second International Conference on*, pages 18–22, 20-22 Oct 1993.
- [34] A. Koerich, R. Sabourin, and C. Suen. Large vocabulary off-line handwriting recognition using a constrained level building algorithm. *Pattern Analysis and Applications*, 6(2):97–121, 2003.
- [35] K. Kukich. Techniques for automatically correcting words in text. *ACM Computing Surveys*, 24(4):377–439, 1992.



- [36] D. Lopresti and J. Zhou. Robust retrieval of noisy text. In *ADL '96: Proceedings of the 3rd International Forum on Research and Technology Advances in Digital Libraries*, page 76, Washington, DC, USA, 1996. IEEE Computer Society.
- [37] R. M. L.R. Bahl, F. Jelinek. A maximum likelihood approach to continuous speech recognition. *IEEE Transactions on PAMI*, 5(2):179–190, 1983.
- [38] G. S. M.A. Jones and B. Ballard. Integrating multiple knowledge sources in a bayesian ocr post-processor. In *International Conference on Document Analysis and Recognition*, pages 925–933, St. Malo, France, 1991.
- [39] R. Manmatha and W. Croft. Word spotting: Indexing handwritten manuscripts. In *Intelligent Multi-media Information Retrieval Collection*. AAAI press, May 1997.
- [40] U. Marti and H. Bunke. The iam-database: an english sentence database for off-line handwriting recognition. *Int. Journal on Document Analysis and Recognition*, 5:39–46, 2002.
- [41] A. McCallum and K. Nigam. A comparison of event models for naive bayes text classification, 1998.
- [42] R. Milewski and V. Govindaraju. Extraction of handwritten text from carbon copy medical form images. In *Document Analysis Systems*, pages 106–116, 2006.

- [43] R. Milewski, S. Setlur, and V. Govindaraju. A lexicon reduction strategy in the context of handwritten medical forms. In *Proceedings of Eighth International Conference on Document Analysis and Recognition*, pages 1146–1150, Seoul, Korea, 2005.
- [44] 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.
- [45] G. Nagy. Twenty years of document image analysis. *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 22(1):38–62, 2000.
- [46] J. Nielsen, V. Philips, and S. Dumais. Information retrieval of imperfectly recognized handwriting. 1993.
- [47] K. Nigam, J. Lafferty, and A. McCallum. Using maximum entropy for text classification. In *In IJCAI-99 Workshop on Machine Learning for Information Filtering*.
- [48] W. B. O. Kolak and P. Resnik. A generative probabilistic ocr model for nlp applications. In *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, pages 55–62, Edmonton, Canada, 2003.

- [49] 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.
- [50] F. O. P. Koehn and D. Marcu. Statistical phrase-based translation. In *Proceedings of the Human Language Technology Conference (HLT-NAACL)*, Edmonton, Canada, 2003.
- [51] R. Price and A. Zukas. Accurate document categorization of ocr generated text. pages 97–102, College Park, MD, 2005.
- [52] N. S. R. K. Powalka and R. J. Whitrow. Word shape analysis for a hybrid recognition system. *Pattern Recognition*, 30(3):421–445, March 1997.
- [53] T. Rath, R. Manmatha, and V. Lavrenko. A search engine for historical manuscript images. In *Proceedings of SIGIR '04*, pages 369–376, Sheffield, England, 2004.
- [54] T. M. Rath and R. Manmatha. Word spotting for historical documents. *International Journal on Document Analysis and Recognition*, 2006.
- [55] 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.

- [56] R. Rosenfeld. The cmu statistical language modeling toolkit and its use in the 1994 arpa csr evaluation. In *Proceedings of the ARPA Spoken Language Systems Technology Workshop*, San Mateo, California, 1995.
- [57] M. S and G. V. Holistic lexicon reduction for handwritten word recognition. In *Proceedings of the SPIE - Document Recognition III, San Jose, CA, 1996*;, pages 224–234, 1996.
- [58] V. K. S. Madhvanath and V. Govindaraju. Syntatic methodology of pruning large lexicons in cursive script recognition. *Pattern Recognition*, 34(1):37–46, January 2001.
- [59] H. N. S. Vogel and C. Tillmann. Hmm based word alignment in statistical translation. In *Proceedings of the 16th International Conference on Computational Linguistics (COLING)*, pages 836–841, Copenhagen, Denmark, 1996.
- [60] G. Salton and M. McGill. *Introduction to Modern Information Retrieval*. McGraw-Hill Book Company, 1984.
- [61] G. Seni, R. Srihari, and N. Nasrabadi. Large vocabulary recognition of on-line handwritten cursive words. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 18(7):757–762, Jul 1996.

- [62] A. Senior and A. Robinson. An off-line cursive handwriting recognition system. *IEEE transactions on Pattern analysis and Machine intelligence*, 20(3):309–321, 1998.
- [63] M. Shridhar, G. Houle, and F. Kimura. Handwritten word recognition using lexicon free and lexicon directed word recognition algorithms. *Document Analysis and Recognition, 1997., Proceedings of the Fourth International Conference on*, 2:861–865 vol.2, 18-20 Aug 1997.
- [64] S. Srihari and E. Keubert. Integration of hand-written address interpretation technology into the united states postal service remote computer reader system. In *Proceedings of Fourth International Conference on Document Analysis and Recognition*, pages 892–896, Ulm, Germany, 1997.
- [65] A. Stolcke. Srilmm - an extensible language modeling toolkit. In *Proc. Intl. Conf. Spoken Language Processing*, Denver, Colorado, 2002.
- [66] A. Takasu. An approximate multi-word matching algorithm for robust document retrieval. In *CIKM '06: Proceedings of the 15th ACM international conference on Information and knowledge management*, pages 34–42, New York, NY, USA, 2006. ACM.
- [67] A. O. Thomas, A. Rusu, S. Mukund, and V. Govindaraju. Non-writer specific handwriting generation for the captcha application. In *Proceedings of the IEEE*

- Western New York Image Processing Workshop*, Rochester, NY, 2007.
- [68] C. Tillmann and H. Ney. Word re-ordering and a dp beam search algorithm for statistical machine translation. *Computational Linguistics*, 29(1):97–133, 2003.
- [69] P. K. U. Pal and B. Chaudhuri. Ocr error correction of an inflectional indian language using morphological parsing. *Journal of Information Science and Engineering*, 16(6):903–922.
- [70] A. Vinciarelli, S. Bengio, and H. Bunke. Offline recognition of unconstrained handwritten texts using hmms and statistical language models. *IEEE transactions on Pattern analysis and Machine intelligence*, 26(6):709–720, 2004.
- [71] H. Xue and V. Govindaraju. Hidden markov models combining discrete symbols and continuous attributes in handwriting recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 28(3):458–462, 2006.
- [72] Y. Yang and C. Chute. An example-based mapping method for text categorization and retrieval. *ACM Transactions on Information Systems*, 12(3):252–277, 1994.
- [73] 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.

- [74] M. Zimmermann and J. Mao. Lexicon reduction using key characters in cursive handwritten words. *Pattern Recogn. Lett.*, 20(11-13):1297–1304, 1999.