

Dissertation Proposal
Department of Computer Science and Engineering
University at Buffalo

Use of Language Models and Automatic Topic Categorization for Indexing and Retrieval of Handwritten Document Images

Faisal Farooq

Thesis Committee: Dr. Venu Govindaraju (Chair)
Dr. Peter Scott
Dr. Matthew Beal

Time: 10:00 – 11:00 am; May 12, 2006
Place: CEDAR Conference Room

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.

The objective of this research is to advance unconstrained handwriting recognition and retrieval technology so that handwritten documents can be accessed with the same ease as printed or ASCII documents. We will use statistical language models to advance the state-of-the-art in handwriting recognition. It is generally believed in the Information Retrieval community that OCR accuracy of about 70%-80% is needed for acceptable user experience. However, current handwritten word recognizers have an accuracy of about 30% for a 10,000 word lexicon. We propose a two-fold approach to improve the accuracy. First, we will develop automatic topic categorization methods to reduce the effective size of the lexicon. Secondly, we will process multiple recognition hypotheses (top n) generated by the word recognizer as it has been observed that the probability that the correct choice is within the top (say) 50 choices can be as high as 70%. We will develop a beam-search methodology to prune and re-rank the multiple hypotheses using the confidences returned by the recognizer and the topic categories. Our conjecture is that we can model the behavior of the word recognizer by analyzing the top choices and exploit this knowledge to facilitate efficient retrieval. The retrieval phase will involve a generative framework based on convolution and relevance models that will process the pruned recognition results. We will extend probabilistic convolution models that have been successfully used in the machine translation domain to overcome recognition errors.

The proposed concepts will be validated on a large publicly available dataset of over 1500 pages written by more than 600 writers covering a wide variety of topics.