

Dissertation Proposal Defense
Enhancement and Retrieval of Low Quality Handwritten Documents

Huaigu Cao

Committee: Dr. Venu Govindaraju (Chair), Dr. Peter Scott, and Dr. Vipin Chaudhary

10:00 – 11:00 AM May 9, 2007
CEDAR Conference Room

Abstract

Recent developments in Optical Character Recognition (OCR) techniques have made it possible to convert large volumes of documents into text, PDF, XML and store them in the computer. Concurrently, advances in Information Retrieval (IR) have led to powerful tools for indexing and searching large on-line databases of text documents. However, the IR performance on OCR'ed text continues to be poor due to challenges presented by document analysis tasks, especially when the documents are of poor quality. Such documents are common in applications pertaining to historical manuscripts, contemporary forms filled in non-ideal conditions, carbon copies, etc. There is a tremendous demand for converting non-digital handwritten documents into electronic form and making them searchable.

The state-of-the-art handwriting recognition techniques can only achieve about 30% accuracy on poor quality handwritten documents with large lexicon. We propose to improve the recognition accuracy and develop probabilistic models for indexing and retrieval to enable a satisfactory search experience for users. We will develop a Markov random field (MRF) based preprocessing algorithm for document enhancement which will use local spatial constraints learnt from a training set (instead of previously used hard coded heuristic rules). The MRF will help in filtering out the form template and extracting the handwritten data.

We also propose to modify the traditional IR model to accept the top-n OCR returned choices instead of just the top choice. Each choice will be weighted in proportion to the probability of the rank returned by the recognizer and based on a probabilistic function learnt from the training data. We will modify the standard vector model (VM) for handwriting to estimate the number of occurrences of terms from word segmentation and recognition probabilities. Our experimental test bed will be a set of 8,000 hand-filled medical (carbon copy) forms. We will show the advantage of our model over VM in terms of mean average precision (MAP), R-Precision, and interpolated 11-point precisions.