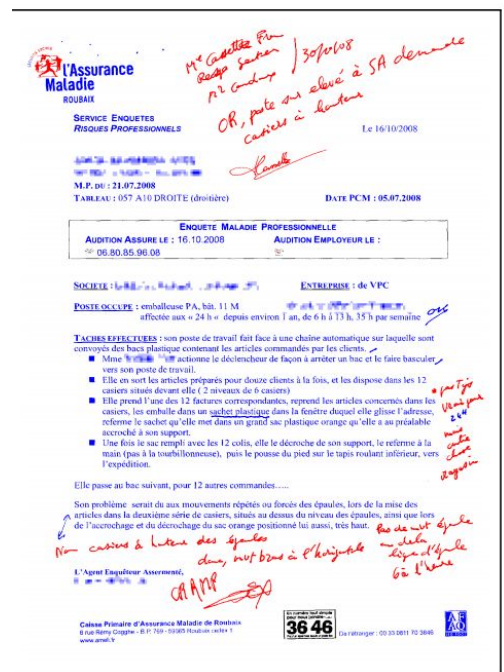


Handwritten and Printed Text Separation in Historical Documents

Objective of this work:

With the increase of digitized documents, automatic document analysis has become extremely important. The presentation of administrative documents for public introduces varieties of document types, content, quality and structure. Fundamentally speaking, documents can be skewed, noisy, overlapped with graphics, i.e., lines, unconstrained annotations, stamps.

In this thesis, existing technologies for visual semantic analysis of unstructured data will be investigated. The dataset of the work consists of 5595 images of tax forms from 1988. In order to provide a textual representation of the documents automatic text and pattern recognition processes have to be applied. An optical character recognition (OCR) system recognizes either printed or handwritten text. Hence, the task of the thesis is to separate machine printed text from handwritten text in scanned documents before feeding it to an OCR system. As a first step, the images will be preprocessed (e.g. cropping, noise filtering). Then the images will be segmented (e.g. line/word segmentation). As the last contribution the segments will be fed into a deep learning binary classifier to be separated into handwritten and printed text segments.



The project work will be supervised by **Prof. Dr. Harald Sack, Tabea Tietz and Oleksandra Vsesviatska, Information Service Engineering at Institute AIFB, KIT, in collaboration with FIZ Karlsruhe.**

Keywords:

Knowledge Graphs, Cultural Heritage, NLP

Pre-requisites:

Knowledge of Programming with Python.

Contact persons:

Tabea Tietz

tabea.tietz@fiz-karlsruhe.de

Oleksandra Vsesviatska

oleksandra.vsesviatska@fiz-karlsruhe.de