# Network PDF file unstructured information extraction

Pengchanghuan,Sunwei,Fuxiaohan

Shanghai Jiaotong University

598095762@qq.com

## Abstract

Modern society is paying more and more attention to the collection and analysis of data.Now we face the annual reports of many companies, we need to intelligently extract and analyze the tables in them.Here comes our project:Network PDF file unstructured information extraction.With our own intelligent analysis extraction program, we have a very high recognition extraction rate (for our test of 32 PDF reports it has a 93% recall rate).

## 1 Introduction

Modern society is paying more and more attention to the collection and analysis of data.Now we face the annual reports of many companies, we need to intelligently extract and analyze the tables in them. Once successfully acquisition of economic data for most companies, quickly access to market information, the company can make correct decision and successfully occupy a favorable position in the market, so Financial data search engine is of great significance.

The problem we have to face, however, is that the vast majority of financial statements are released in the form of PDF documents.

PDF(Portable Document Format), The purpose of designing PDF file format for Adobe company is to support the information publishing and publishing of Multimedia Integration on cross platform, especially to provide support for network information publication. To this end, PDF has an incomparable advantage over many other electronic document formats. The PDF file format encapsulates text, fonts, formats, colors, and graphics and images independent of device and resolution in a file. The format file can also contain hypertext links, voice and dynamic images and other electronic information, support special documents, integration and security, reliability is higher.

However, the PDF file is a form of documentation designed to prevent modifications (in fact, it is difficult to modify them).So our problem first is to overcome the difficult problem of PDF documents.

We will initially take the text of the PDF file operation, turn it into a CSV file, then the balance sheet to generate a description (how specific description is to design, including fuzzy matching), finally scan the text, positioning the balance sheet, and data extraction.

Our core technical issues are named entity recognition, report description, density based report scanning, positioning, and so on.

As far as our progress is concerned with our own intelligent analysis extraction program, we have a very high recognition extraction rate (for our test of 32 PDF reports it has a 93% recall rate).

## 2 Our process

### 2.1PDF2CSV

Most of the time we get the PDF format.However, there are many limitations to PDF files.For example, you can't copy text, edit data tables, and cannot extract problems like editing Word files.

Tabula[1] can help you extract data tables from PDF files and save them in CSV format so that you can easily access data and use it for the second time.It is an open source free program.

First, we import tabula-java as well as Python library xlwt in our python program.After you enter the correct parameter to these tools,you can transform the entire PDF file into a CSV file.In the CSV file,text and data are separated by commas and line breaks.

Unlike PDF format, CSV is a convenient file type, so we can search, modify and extract the contents of CSV files freely.The Figure 1 shown below is an example of extracted CSV file.



Figure 1:Extracted CSV file styles.

### 2.2Preliminary extraction

We use the table header specific known as a template for the corresponding part of the matching CSV file, extract the corresponding part of the data, and then

stored in our original database, in order to collect the training data we need for fuzzy matching that have the word form and characteristics.

These features include text content, numeric type, numeric size, literal numeric distribution, text density, and so on.

In this experiment, we take 27 known headers of 5 different reports as our initial samples, and all of them are successfully matched to the data in the CSV file.

We can have an example below in Figure 2:



Figure 2:Extracted tables and data which will be used as raw samples.

### 2.3 Autodetect module

We obtain a similar approach to intelligent extraction of tabular data features in a single article called:"Research on web core block extraction algorithm based on DOM node text density"[2].Then we develop our own autodetect module.

We find that the appearance of tables is related to many parameters.Such as text content, numeric type, numeric size, literal numeric distribution, text density, and so on.We match the model trained by the existing data in the database with the contents in the CSV file,we compare the feature sets between them.Then select the higher matching part as a table.

This is the processes of our fuzzy

matching for intelligent table extraction.

In addition,We also take the line distribution in the graph of the PDF file as the basis for the table to appear.In this process, we also use tabula as a tools ,and we modify the code to make the tool meet our experimental requirements.

## 2.4 Visualization

At the same time, we have made some visualization of the extraction of tables to facilitate the use of developers and users.

When using, we can use manual identification methods to find the forms which have not been successfully recognized. All the tables are extracted correctly under manual surveillance

## 2.5 Convert to CSV file

This is the final part of out project,we transfer the tables we extracted into CSV form,in order to make people have a good reading experience.

## 3 Evaluation

In the test, we used 32 PDF file.We have a recall rate about 93%,such a high recall rate indicates that our program is of practical value.

## 3.1Table search

In this part we search the tables using fuzzy matching.We will have a preview below in Figure3:



Figure 3:Visualization of table search.

In this part,if there are something wrong with the table search,we can manually correct the result by putting a box on the table that is missed.

## 3.2Table extraction

This is our final part we extract our tables ,transfer them into CSV files.We will

have a preview below in Figure 4:



Figure 4:Extracted tables.

This article does not have much theoretical knowledge, this project uses practical as the goal, but also hope that the teacher can give more guidance to me.

## 5 Conclusion

With our own intelligent analysis extraction program, we have a very high recognition extraction rate (for our test of 32 report draws have 93% recall rate).

## Reference

[1] http://tabula.technology/
[2]http://xueshu.baidu.com/s?wd=paperuri%3A%284cc2a28db2c77a0e3a929d54ef9ca624%29&filter=sc_long_sign&tn=SE_xueshusource_2kduw22v&sc_vurl=http%3A%2F%2Fcdmd.cnki.com.cn%2FArticle%2FCDMD-10007-1012007246.htm&ie=utf-8&sc_us=18388595119318769597

## 4 Discussion

In our project we can't extract the title of the table,because it can exist in all directions of the table.So the next thing to do is to look for the title of the table intelligently.

In this project,the senior gave me the train of thought,then, the following program editing, data collection, tool search, and, most importantly, the results are all done by myself.