# Document Structure Recognition

Shan Zhong
University of South Carolina, USA
zhongs@email.sc.edu

## Abstract

Electronic documents are a common way to store information. Contracts and data tables from company's public disclosures information, government's orders can be easily accessed and download by internet. These data are valuable, usually people just download and clean them manually, but when we have thousands of websites, documents to download and analysis, then it becomes problematic

During this summer, I lead a small team with two undergrad students from two top universities of China. They have basic programming skills with also some knowledge in statistical science. We together worked on the document recognition problem.
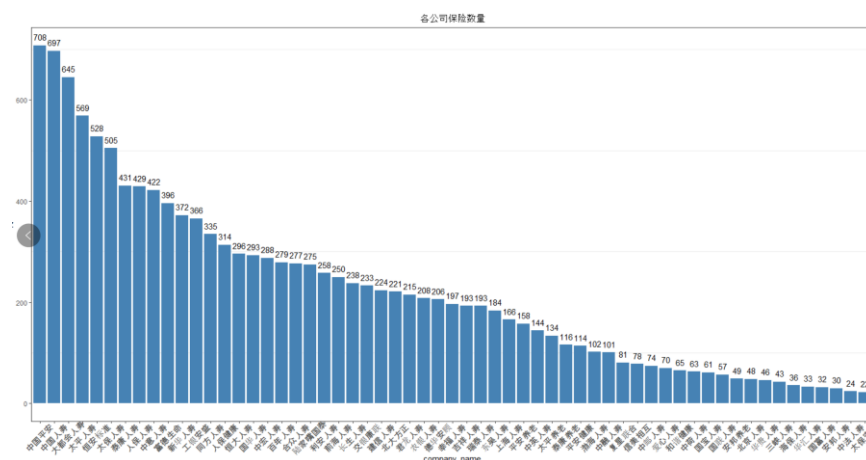
## Goal and purpose

The main problem is currently there is no good way to organize and analyze them in a universal way. As they have a variety of typesetting and formatting (docx, doc, pdf, images，directly posted online….). I also aware that now days there exist tools like Pdf Plumber that

can recognize lines and cells in excel, (https://github.com/jsvine/pdfplumber , paper :

https://dspace.cc.tut.fi/dpub/bitstream/handle/123456789/21520/Nurminen.pdf?sequence=3),

plug in tools like tesseract (https://github.com/tesseract-ocr/ ) that can do deep learning text

recognition, powerful clustering packages like Scikit-learning(https://github.com/scikit-

learn/scikit-learn) in python. So being able to utilize and incorporate existing methods that

can solve part of the problem and then integrate them into one algorithm that can organize

and clean dataset as a human do can be valuable and time-saving, so this is the main purpose

for this paper.

## Data

Figure1: companies with the number of products they have in y axis.



There are 91 insurance companies in china that sells health product, we first get a list for all

of them, then find their company website. We used python based scrapping tools Scrapy

(https://docs.scrapy.org/en/latest/) to scrape insurance contracts, public disclosure

information, insurance coverages description, and premium rate table for their products.

Finally, we found there are around 20,000 products (including those stopped selling) existed. We made a data base for all of them.

How de download the data is we first created a url filename list from the scraping tool's observation. Some company websites were not that easy to scrap, so we created a virtual machine using windows pro docker to render these websites. We then managed to get a list of all downloadable items from their website. We cleaned these datasets first, then download items according to the links on these lists. Save them locally according to different company different product types and separated them into different folders according to some rule created by our algorithm. Here we created algorithms that can process thousands of files and do operations simultaneously, for example check if successfully download, check if we failed to download and then retry. Some website has software that preventing us downloads lots of material at the same time, so we also adjust algorithm for that.

Figure 2: the url lists

Figure 3, python scrapy code frame with each company, each python file representing a company

## Method

1. **Integrating documents into one universal format**

From scrapping tools we downloaded thousands of documents to analysis and counted as our sample data. The first thing we need to do is to make them have at least the same format. First, we tried converting all of the documents into a txt format, then we found we lost the structure, front, location of the text, although we can make the text still follow some order and we can still use it to preform analysis

We then transfer all documents into a pdf format, thus information like structures can be remained. We designed several algorithms to detect their structure, read their content, and save them in a more universal way (one dimensional dataset, for example: one column indicates if it is the title or subtitles, another column keeps the whole paragraphs, one column

to keep the notes...). Why we want to do so is when just facing insurance documents, there are many similar patterns such as "what are our coverages", "when we refuse to pay" as documents subtitle. We can then use regular expressions to recognize them.

Figure 4 and 5: detecting lines, and texts.



## 2. Numeric table and contract documents

How we determine the structure of text is by detecting the position of every character pieces in the document. Because of the special properties of text, we recognize each of them as a small rectangular. And we also detect horizontal and vertical lines by using python Ocr tools. We do not recognize oblique texts and line, we do not recognize pictures (which we want to focus on later). We get x, y axis position, also width and height of texts. We also get start and end position of lines.

Although there are all kinds of typesetting, but we can still categorize them into two main parts: data tables, and contract document texts. We did not encounter many pictures, and haven't get a good method to deal with right now.

**Data table:**

Figure 7 and 8: collapse every character, identify the description part, table parts, and each



cell.

Inside these document files, some of the information are saved as premium tables and

mortality tables. These tables have data in two dimensional (parameters are saved in rows and

columns relative to each value) that we want to transfer to one dimension (all parameters are

saved in columns). Sometimes it is not that easy as there may be several layers of index

columns and rows (such as firs layer of columns count gender while second layer of columns

count for insurance term periods.

Since we have used Ocr tools to detect the lines in document, we can then get the

intersections of horizontal and vertical lines. This allow us to collapse intersections within a

tolerance range into one bigger intersection as to represent for a larger cell. We can then

combine texts within each cell, align them to each layer of columns and rows to get the other

information. For example: check if a cell's x position is within the x position of a column index. In this way we transferred the data into one dimension.

The difficulty part of the above algorithm is how to implement them in an efficient way.

Also, we collected the insurance synonym list so we can save the data in a universal way, which is extremely important when creating one-dimensional dataset to save them.


**Word document:**

Figure 9 and 10, example of document position recognition



The ultimate final goal for us to analysis these documents is to find the detailed difference in the interpretation of insurance contracts, for different insurance companies and products: In the contracts there do exist similar paragraphs and regulations that allows us to cluster them together. We can then compare which product have more coverage, which product allow us to buy the product easier (insurance products have limitation on who can buy the product, for

example those who had cancer before cannot buy cancer insurance, while some product

allows in some cases)

For achieving such goals, we first tried a tolerance method: compress all the texts in one

dimension, label their x, y position as rectangles, then set the x tolerance and y tolerance of

these rectangles. After that, collapse all rectangles within the tolerance range to one bigger

rectangle. But this method will usually cause some mistakes and it also not the way that

human read and detect text structures. For example, it does not do well on the vertical axis. It

also required people to set a tolerance range.

Figure 11, an example of one dimensional data table, with text indicate the original text.

```
> pdf_extract_word(pathnames[20])
           x        y height width     doc_y page text index
1     94.681 623.561   9.432 4.260   623.561    1    4     1
2     94.681 678.163   9.432 4.260   678.163    1    5     2
3    309.601 129.552   8.369 3.780   921.552    2    7     3
4    313.306 129.552   8.369 1.890   921.552    2    .     4
5    315.271 129.552   8.369 3.780   921.552    2    1     5
6    257.404 496.397   8.369 3.780  1288.397    2    7     6
7    261.109 496.397   8.369 1.890  1288.397    2    .     7
8    263.074 496.397   8.369 3.780  1288.397    2    2     8
9    287.044 744.041   9.431 4.260  1536.041    2    1     9
10   319.079 744.041   9.431 4.260  1536.041    2    1    10
11   322.828 744.041   9.431 4.260  1536.041    2    0    11
12   286.562 245.236   8.369 3.780  1829.236    3    7    12
13   290.266 245.236   8.369 1.890  1829.236    3    .    13
```

Then we use a simple k mean clustering method on the x-axis to separate the subtitles and the

detailed information besides it. We found just by implementing the easiest clustering

algorithm, we can bring a huge increase in the accuracy

We then compare the time the algorithms needed to run in both python and R, and modified computing method to increase algorithm efficiency. We also get a review at the accuracy these algorithms gained.

Finally, we managed to design a new algorithm of identify based on incorporating several existing methods together, and with the input feed in the file path, the data will be transferred and stored in mysql data base.

Figure 12: sample mysql table



We plan to use information like front size and style to determine whether it is foot notes, title, or descriptive stuffs in the future.

**What to do with the data we get**

We transfer these word documents into a row of data, indicating each attribute we extract from the contract document. Such as coverage periods, coverage amount…etc. We built up a insurance term dictionary to look up for these words.

Figure 13, example of the word lookup package we built

| flags: | status1 | status2 | |
|---|---|---|---|
| | 1 flags = re.IGNORECASE : ignore cases | flags = re.VERBOSE : allow comments in regular expression | |
| | | | |
| index | functions | Usage and effects | Arguments |
| | 1 str_detect(pattern,dat,flags=0,negate = False) | Detect the presence or absence of a pattern in a string | negate: If TRUE, return non-matching elements. |
| | 2 str_which(pattern,dat,flags=0,negate=False) | Find positions | negate: If TRUE, return non-matching elements. |
| | 3 str_subset(pattern,dat,flags=0,negate=False) | Keep strings matching a pattern | negate: If TRUE, return non-matching elements. |
| | 4 str_count(pattern="",dat,flags=0) | Count the number of matches in a string | |
| | 5 str_starts(pattern,dat,flags=0,negate = False) | Detect the presence or absence of a pattern at the beginning of a string | negate: If TRUE, return non-matching elements. |
| | 6 str_ends(pattern,dat,flags=0,negate = False) | Detect the presence or absence of a pattern at the end of a string | negate: If TRUE, return non-matching elements. |
| | 7 str_replace(pattern,replace,dat,flags=0) | Replace matched pattern(first one) in a string | |
| | 8 str_replace_all(pattern,replace,dat,count=0,flags=0) | Replace all matched patterns in a string | count: maximum number of replacement for a string |
| | 9 str_remove(pattern,dat,flags=0) | Remove the first matched patterns in a string | |
| | 10 str_remove_all(pattern,dat,flags=0) | Remove all matched patterns in a string | |
| | 11 str_extract(pattern,dat,flags=0) | Extract the frist matching patterns from a string | |
| | 12 str_extract_all(pattern,dat,target,flags=0) | Extract all matching patterns from a string | target: choose to operate on a string or vector |
| | 13 str_locate(pattern,dat,flags=0) | Locate the position of the first patterns in a string | |
| | 14 str_locate_all(pattern,dat,flags=0) | Locate the positions of all patterns in a string | target: choose to operate on a string or vector |
| | 15 str_sub(dat,start=1,end="len") | Extract substrings from a character vector by position | choose from start position to end position |
| | 16 str_c(*pattern,sep="") | Join multiple strings into a single string | sep: add words between to pattern |
| | 17 str_split(pattern,dat,split_point="remove",max_split=0,flags=0) | Split up a string into pieces and remove split patterns | max_split: control the maximum number of pieces |
| | str_split(pattern,dat,split_point="before",max_split=0,flags=0) | Split up a string into pieces and put split patterns at the beginning of strings | max_split: control the maximum number of pieces |
| | str_split(pattern,dat,split_point="after",max_split=0,flags=0) | Split up a string into pieces and put split patterns at the end of strings | max_split: control the maximum number of pieces |
| | 18 str_trim(dat,side="both") | Trim whitespace from a string | side = "both" "right" or "left" |

Also, for conveniently look for combination of words in python. We designed our customized regular expression lookup functions. We mimic the R Stringr package and added several additional features.

We matched the illness definition by the insurance contract, to check if there are overlaps, and check the percentage covered by each different insurance product. We can then put these attributes into the column for our later analysis.

Figure 14: we have a lot of variables, there are around 100.



Our dependent variable will of course be the price, we built a contract valuation regression model to predict the actuarial value of contracts sold by each different company. As oppose to the pure actuarial model which was based on the modification of mortality table. We compared the result for both model, and present to our customers.

We used the 2010 china Critical illness life table and the mortality table in building the pure actuarial model. We found using regression model, it is much earlier to explain them to customers. We can directly figure out why one product is more expensive when doing comparison, and find out the most valuable products.

**Conclusion**

This paper is mostly about the implementation part, I still need to work on the algorithm to really discover something totally new.