# Use Web Crawler & Doc2vector
# to analysis Similarity Litigations

Bo-Hong Liu

*Department of Industrial Engineering and Engineering Management, National Tsing Hua University, Taiwan*

**Abstract.** The general search engine is to search by keyword, although you can get a very large number of search results on the page, but still need to interpret the human to confirm whether the document content is required by the user. Legal staff in court proceedings will require the precedent as the basis for the relevant judgments, it is necessary to find similar cases to clarify the case of litigation, this study to crawl Web crawler search engine to find the relevant case articles, pre-treatment in the article Afterwards, the document value of Doc2vec is given, which can be used to make the similarity analysis of case based on the vector's directionality and vector matrix. It can provide relevant case documents of interest to the user and help the legal person to search the precedent of judgment Efficiency and accuracy.

*Keywords*: *Web crawler, Doc2vec, Litigation,* Word2vec, Precedent

## 1. Introduction

In today's era of online information development, the transmission of information is no longer based on paper records alone. It is based more on electronic records and on-line information, and even directly places information on the cloud, thus making information more transparent and readable. Government agencies and non-profit organizations are also happy to share information with people, including the court litigation online.

However, a large number of databases and sources is difficult to find the information really looking for, even using Google and other search engines still need to manually interpret, facing the massive information crawling is only the tip of the iceberg and found Of the information is true or difficult to determine, so in the text and other analysis will be inefficient problems. Therefore, in the face of a large amount of information, we can crawl the data by using the way of web crawler, and then use the doc2vec analysis method to find out the most similar articles to meet the needs of users, saving time and doing the right thing.

## 2. Literature review

### 2.1 Web Crawler

Web crawlers is an automated program that methodically scans through web pages to create an index of the data it is set to look for. This process is called Web crawling.

Web crawlers can copy all the pages they visit for later processing by a search engine which indexes the downloaded pages so the users can search much more efficiently. Crawlers consume resources on the systems they visit and often visit sites without tacit approval. Mechanisms exist for public sites not wishing to be crawled to make this known to the crawling agent. For instance, including a robots.txt file can request bots to index only parts of a website, or nothing at all.

A Web crawler starts with a list of URLs to visit, called the seeds. As the crawler visits these URLs, it identifies all the hyperlinks in the page and adds them to the list of URLs to visit, called the crawl frontier. URLs from the frontier are recursively visited according to a set of policies. If the crawler is performing archiving of websites it copies and saves the information as it goes. The archives are usually stored in such a way they can be viewed, read and navigated as they were on the live web. (As showing Figure 1)
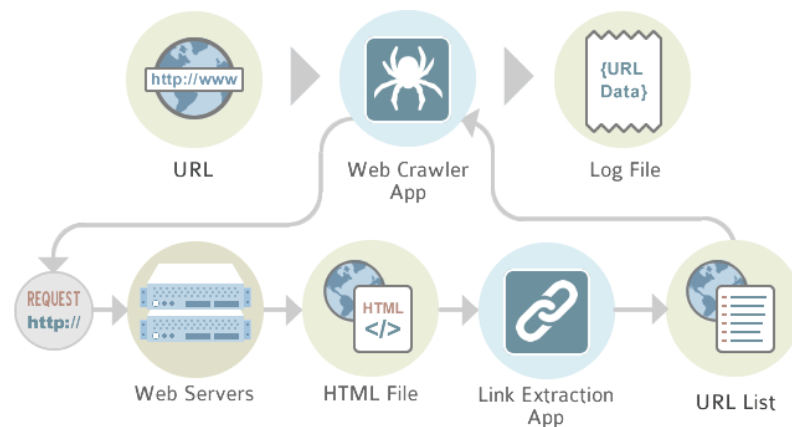


Figure 1

### 2.2 Natural Language Processing (NLP)

NLP is a method for computers to understand human language in a smart and useful way. By NLP, let the computer understand the meaning of a sentence or article, nouns, verbs, adjectives and adverbs are the basis of the literary meaning NLP from the text of the sentence will be non-structural features of the text presented as analytical data NLP is used to analyze text, allowing machines to understand how human's speak. This human-computer interaction enables real-world applications like automatic text summarization, sentiment analysis, topic extraction, named entity recognition, parts-of-speech tagging, relationship extraction, stemming, and more. NLP is commonly used for text mining, machine translation, and automated question answering.

NLP is characterized as a hard problem in computer science. Human language is rarely precise, or plainly spoken. To understand human language is to understand not only the words, but the concepts and how they're linked together to create meaning. Despite language being one of the easiest things for humans to learn, the ambiguity of language is what makes natural language processing a difficult

problem for computers to master.

*2.3 Word2vec*

The word vector (also known as word embedding or representation) is a technique widely used in Natural Language Processing (NLP). Using a vector to represent each word, so that convert a sentence into a vector of words to represent it and numerically. This research use Word2vec which develop by Google to convert the words into vectors.

This tool provides an efficient implementation of the continuous bag-of-words and skip-gram architectures for computing vector representations of words. These representations can be subsequently used in many natural language processing applications and for further research.

Word2vec uses distributed representations of text to capture similarities among concepts. Showing as Figure 2, it understands that Paris and France are related the same way Berlin and Germany are (capital and country), and not the same way Madrid and Italy are. This chart shows how well it can learn the concept of capital cities, just by reading lots of news articles -- with no human supervision:
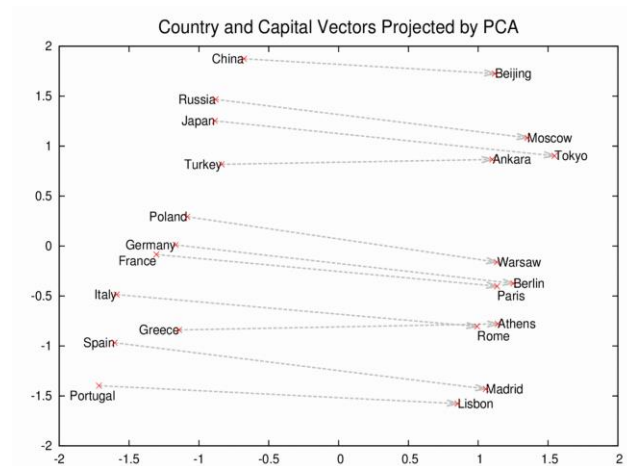


Figure 2

The model not only places similar countries next to each other, but also arranges their capital cities in parallel. The most interesting part is that we didn't provide any supervised information before or during training. Many more patterns like this arise automatically in training.

*2.4 Doc2vec*

The approach for learning paragraph vectors is inspired by the methods for learning the word vectors proposed by Quoc Le and Tomas Milolov (2014). The inspiration is that the word vectors are asked to contribute to a prediction task about the next word in the sentence. So despite the fact that the word vectors are initialized randomly, they can eventually capture semantics as an indirect result of the prediction task. Users use this idea in paragraph vectors in a similar manner. The paragraph vectors are also asked to contribute to the prediction task of the next word given many contexts sampled from the paragraph. The sentence vector is unique by combining word vectors with shareable word vectors.Doc2vec attempts to predict the probability of words in the context of sentence and sentence vectors: Learning training is done via a set of documents or sentences, given the vector so that the

3

probability of the next word appearing in the sentence vector is predicted.

In calculation, the two will be accumulated or connected, as the output layer. During the training of a sentence or paragraph, the paragraph id remains the same and contains the same paragraph vector, which is equivalent to using the semantic meaning of the sentence each time the word is predicted In the predictive phase, through a fixed word vector to infer that a paragraph id is assigned to the sentence to be predicted by the sentence vector, the vector parameters of the word vector in the training phase are invariant, and the sentence to be predicted is trained by using the gradient descent. Finally, convergence can get the sentence to be predicted paragraph vector.

To achieve doc2vec there are two ways; namely, Distributed Bag of Words (DBOW) and Distribution Memory (DM), as shown in Figure 3 below and Figure 4, in a shorter sentence, you do not need a given word vector, the sentence vector Be trained to predict words.
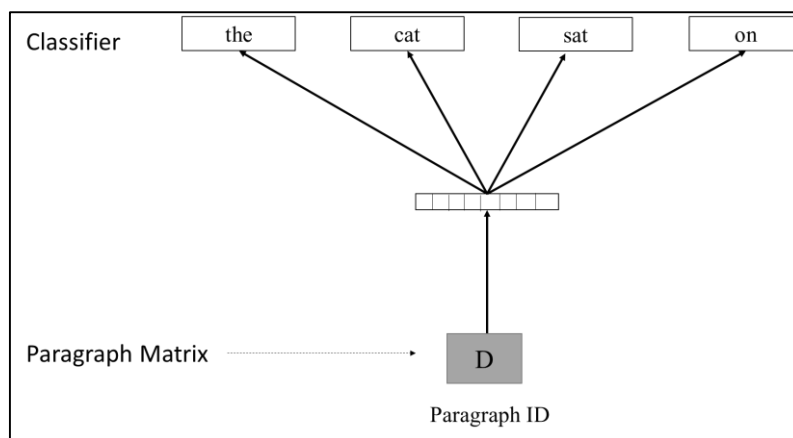


Figure 3    Distributed Bag of Words (DBOW)Model

The second is to learn a textual vector schema that predicts the fourth word ("on") in terms of three words ("the", "cat", "sat"), and the words that are typed are projected matrix w to predict words
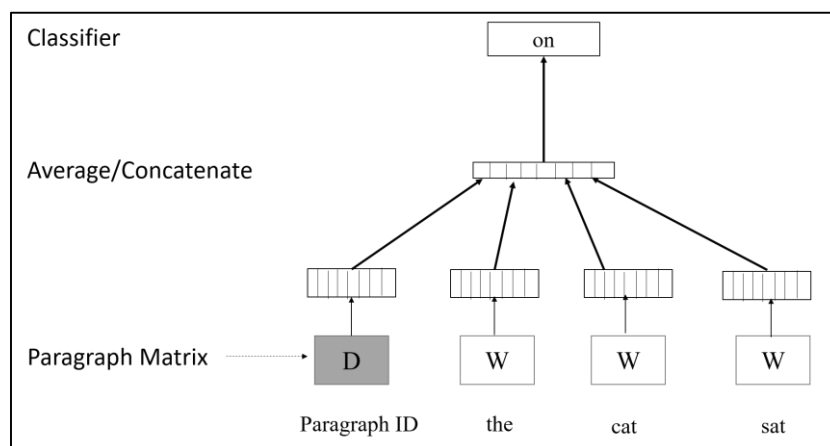


Figure 4    Distributed Memory(DM) Model

*2.5 Programing Language-Python*

Python has become a popular programming language in recent years. It is a powerful, fast, easy-to-read programming language. This report will rely heavily on python for web crawlers and dec2vec

tools.

## 3. Methodology

It is difficult to parse a large number of documents on these networks. Therefore, it must be done in stages.

Showing as Figure 5, first, web crawlers should be implemented on the website to crawler pages, and the simulation package can be used to automatically crawler and put the contents of the documents into the database. After the establishment of the database, it can set the base as the dictionary. And then through word processing to eliminate redundant words and unnecessary noise characters, through Doc2vec to create a training module, you can give each document vector, and you can use the module for a critical comparison that Document the relationship.
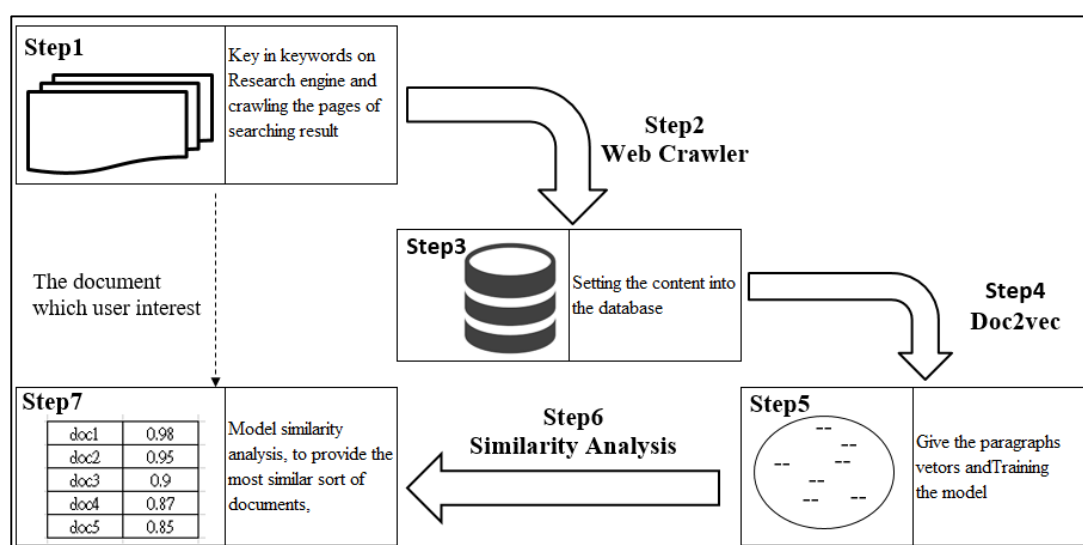


Figure 5

By using this method instead of using manual identification, not only increase the horizon and depth of the data, but also save a lot of time. For professional analysts, they can view documents more systematically.

## 4. Case Study

### 4.1Backround (As-is)

When handling legal proceedings, general law and order officials will make similar decisions in the light of the so-called precedent precedents and the application of previous norms in the light of similar precedent and circumstances as evidence in the face of similar legal cases or facts. Following the precedent as the basis of the precedents, therefore, it is necessary to find similar precedents as an important basis for legal practitioners in hearing or assisting litigation to clarify the offensive and defensive.

Through the court litigations online, official and unofficial network provide search engine for people and law politicians more conveniently search the relevant judgments of the document. Although network offer the massive data, tradition manual way to find relevant case is difficult and waste time.(Showing the Figure 6) Most of the general search engine for keyword research, but do not have

the actual semantic analysis. Users still have to find the page from the keyword in accordance with the title to read, and then enter the page to confirm whether the content is consistent, it will take a lot of time to do artificial confirmation. Search by one page can also cause the following page can't be seen. Keyword search is word search, but there is no way to ensure the importance and understand the relationship between document and document.
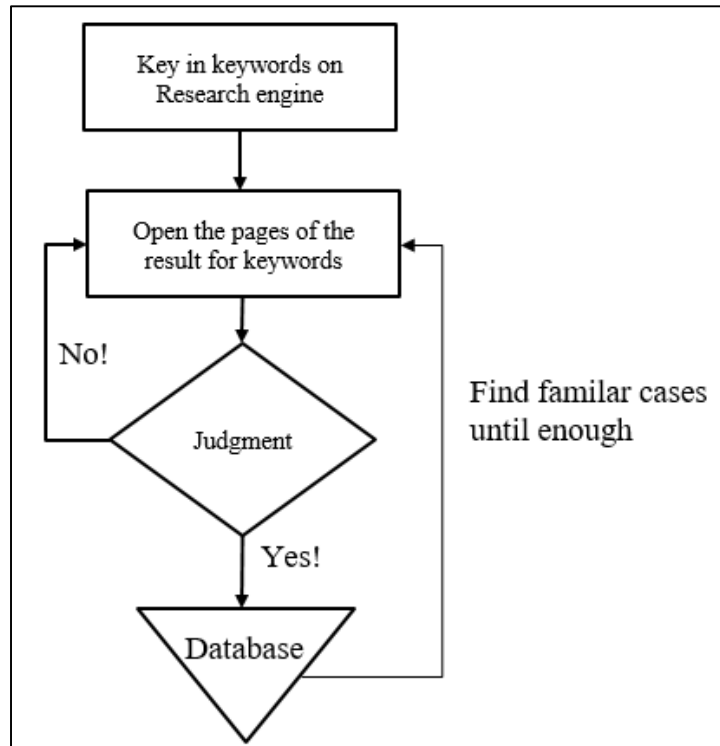


Figure 6

There are some major reasons that searching problems.

1. Man-made judgments time-consuming and laborious

2. Limited to man-made interpretation of the restrictions that only can find a smaller scope

3. To judge the searched cases of relevance need to try &error again and again until find the similar ones.

*4.2 To-be model*

The paper provide the semantic analysis to figure out the similarity of the specific litigations. By programming language, Python, to do web crawler and setting database, then convert the documents into paragraph vectors by doc2vec, after that training the model which can comparison the similarity

Step1: User typing the keywords which interest in the litigation search engine

(As showing Figure 7)    (Web provide: Find law http://lp.findlaw.com/)

Figure 7

Step2: Use Python to crawl all the text of the litigation to climb out into        database. Obtain network source code, are generally html format files, as long as    the study of html tag (tag) structure, and then parse, you can get the required    information. With the package "BeautifulSoup" which base on the properties of      html attributes and retrieve the text to find the web page source code. (Showing Figure 8)
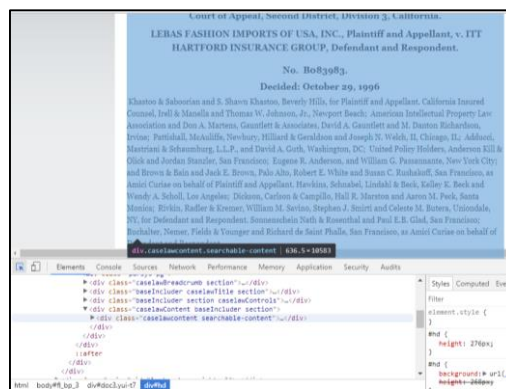


Figure 8

Through the package "selenium", it can simulate the user's behavior to            automatic web crawler and switch pages. (As showing Figure 9)
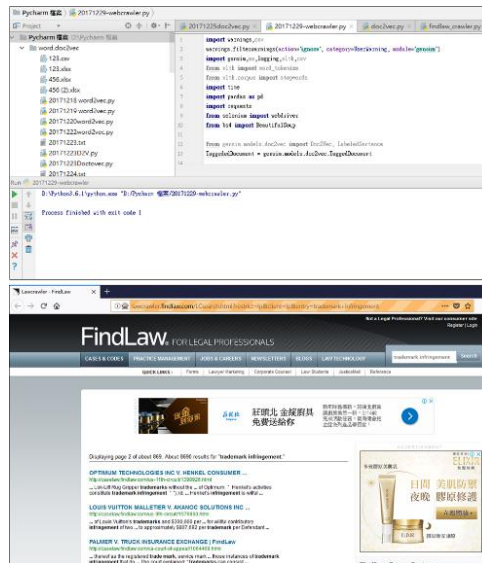
Figure 9

Step3: Capture the text into the database: Put the data into the database, and turning the data frame as Csv type. (As showing Figure 10) With the dataframe, it is more easily to manage the whole data to analysis in next step.
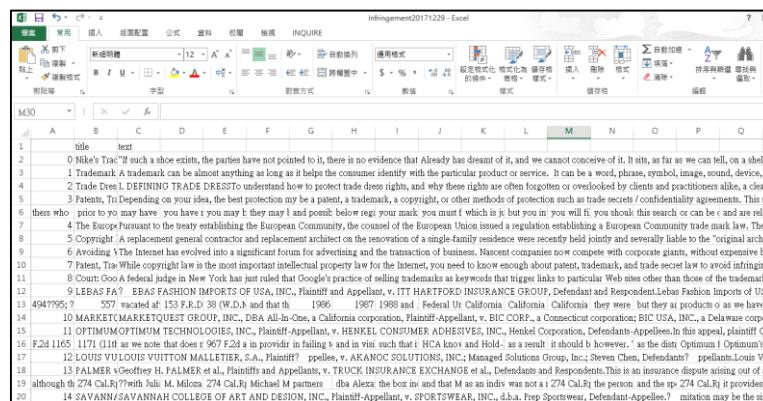


Figure 10

Step4: The use of Python for semantic, word meaning analysis, the establishment of the database model. Before that, it is important that doing text pre-processing which include wordiness, punctuation mark and specific words or signal. After that, setting each document a unique tagged id, and training as a model.

Step5: According to user's needs to find similar documents

Users can interpret their own to find out the case of interest to the module, and do similarity analysis of the whole module. Finally, module provides the highest degree of similarity to the user.

In this case, assume that users are interested in case No. 9 and want to find out the documents related to them, and then use similar analysis to find similar documents. While the similarity is compared with the vector angle. The equal $\cos\theta$ is to 1, the higher the similarity is, so the relationship can be obtained according to the arrangement. If there is a low similarity among this document and other documents, the relevance of the document is low or can be changed Search again for new results.

(As Showing Figure 11)



Figure 11

For the similarity document of the 9th litigation, the similarity of the litigations is (123, 0.8105), (201, 0.792), (15, 0.784), (37, 0.763), (64, 0.752), (215, 0.751), (192, 0.7432620525360107), (143, 0.732), and so on. The ninth issue concerns the acquisition of a company and the confusion of consumers with the trademark. As a result of a manual review of similarities in litigation, it would be found that the first five articles are all litigation with trademark infringement and corporate acquisitions, and therefore can satisfy users' needs. (As showing Table1)

| Litigation | Similarity | Type | Cause |
|---|---|---|---|
| 9(test) | **1** | Trademark confusion | The Acquisition contract of acquisition causes the reverse confusion |
| 123 | 85.8% | Trademark confusion | Advertise of reverse confusion. |
| 201 | 78.6% | Trademark confusion | Similar trademarks of the two wineries cause trademark infringement. |
| 15 | 78.3% | Trademark confusion | Trademark dilution caused by internet website. |
| 192 | 77.9% | Trademark confusion | Similar trademarks of the two wineries cause trademark infringement. |
| 143 | 76.0% | Trademark confusion | Trademark license dispute. |

Table1

Compare the As-is and To-be model with several search efficiency indicators，(Showing as Table 2)

| Index | As-Is | To-Be |
| --- | --- | --- |
| Search Method | Manual | Computer assist |
| Search Scope | less than 5% | 100% |
| Search Efficiency | Low | High |
| Search Accuracy | uncertain | Systematic |
| Time Cost | high | lower |

Table 2

Table 2 shows that As-is Model traditionally performs manual interpretation, which is less efficient and accurate than To-be model, and takes more time. In the other hand, program analysis is also more systematic, which can use the similarity relationship to search for similar cases.

## 5. Conclusion

With regard to the speed of receiving information and the convenience of searching, people's lives have been surrounded by a great deal of data and information. The data analysis and quantification of words can predict the behavior of human beings.

However, in the face of a large amount of information, Find what user really want, Big data is not equal to Right data, when access to large amounts of data, but also how to make the right use of the choice of kinds of methods and tools. In this case study, though there is a wealth of documentation on the web for reference, it is better than the dark days of paper space and confined space, but the attendant is how to view the information users want. Web crawler can quickly crawl more than 300 documents, and by doc2vec to analyze the similarity between the various document vectors to sort out the document recommended to the user, it can greatly reduce the user to view one by one, but also increase the breadth of the scope of the document, for the control of litigation can be more perfect. Do the right thing is better than do the thing right.

**References**

1. http://pushwindersingh.com/what-is-web-crawlers/

2. https://code.google.com/archive/p/word2vec/

3. Le, Q., & Mikolov, T. (2014). Distributed representations of sentences and documents.

In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)* (pp. 1188-1196).

4. http://cs224d.stanford.edu/lectures/CS224d-Lecture2.pdf

5. https://rare-technologies.com/word2vec-tutorial/

6. 玩轉社群：文字大數據實作 謝邦昌、鄭宇庭 五南出版社

7. Python 初學特訓班 文淵閣工作室 碁峰出版