

BDA Mini Project Report

News article scraping and keyword extraction

Sheetal Thakkar 1521002

Farheen Kamal 1521008

I. Introduction

Scraping in simple terms mean to remove an unwanted covering or a top layer from something. There is a wide range of applications of Nearest Neighbour Search, one of which is document analysis. Document analysis is a form of qualitative research in which documents are interpreted by the researcher to give voice and meaning around an assessment topic. There are primarily three types of documents used for analysis viz. Public records, personal records and physical evidences. Document Analysis holds its application in Plagiarism detection, News Aggregator, Text Summarization and so on. The idea for this project was incepted considering the applications and interest in this domain.

II. Problem Definition

This project aims at creating a script to scrape news articles, provided the URL of the article. It also gives keywords and summary of the content using NLP functions provided by Newspaper library. The script also extracts keywords using MapReduce strategy and both these results are compared to get an insight of similarity.

III. Implementation details

This project uses Newspaper, a Python 3 library for extracting and curating articles. Keyword extraction is done using the library functions and also using MapReduce, results from both these approaches are then compared. Text summarization is also implemented, besides, this work also extracts the title, top images and videos from the articles.

Steps for implementation are:

1. Implement mapper and reducer to truncate the common 2 letter and 3 letter words and get the count of frequency of each word in the content. Arrange these words in descending order, and get the top ten keywords.
2. Execute the mapper and reducer using shell script.
3. In the main Python script, import all the necessary libraries and scrape the article using Newspaper.

IV. Result and Analysis

```
@sheetal-Inspiron-3558:~/BDA
sheetal@sheetal-Inspiron-3558:~/BDA$ python3 BDA_News_Aggregator.py
Enter the URL of the news article: http://fox13now.com/2013/12/30/new-year-new-laws-obamacare-pot-guns-and-drones/
-->TITLE: New Year, new laws: Obamacare, pot, guns and drones
****

-->AUTHORS/ PUBLISHER:
Cnn Wire
****

-->PUBLISH DATE: 2013-12-30 00:00:00
****

-->ARTICLE CONTENT: By Leigh Ann Caldwell
WASHINGTON (CNN) — Not everyone subscribes to a New Year's resolution, but Americans will be required to follow new laws in 2014.
Some 40,000 measures taking effect range from sweeping, national mandates under Obamacare to mar
****

-->SUMMARY: Oregon: Family leave in Oregon has been expanded to allow eligible employees two weeks of paid leave to handle the death of a fami
ly member.
Arkansas: The state becomes the latest state requiring voters show a picture ID at the voting booth.
Minimum wage and former felon employment workers in 13 states and four cities will see increases to the minimum wage.
New Jersey residents voted to raise the state's minimum wage by $1 to $8.25 per hour.
California is also raising its minimum wage to $9 per hour, but workers must wait until July to see the addition.
****

-->KEYWORDS(NLP Library):
[ leave , 3 ][ family , 3 ][ states , 2 ][ minimum , 4 ][ law , 3 ][ drones , 3 ][ laws , 4 ][ national , 2 ][ state , 4 ][ pot , 1 ]
[ latest , 3 ][ wage , 5 ]
****

-->KEYWORDS (MAP REDUCE):
[ minimum , 4 ][ state , 4 ][ family , 3 ][ leave , 3 ][ drones , 3 ][ latest , 3 ][ Oregon: , 2 ][ national , 2 ][ eligible , 2 ][ Il
linois , 2 ][ Rhode , 2 ][ passed , 2 ][ through , 2 ][ California: , 2 ]
****

-->SIMILARITY: 70.0 %
****

-->ARTICLE HTML: <!DOCTYPE html>
<!--[if lt IE 7]> <html class="no-js lt-ie9 lt-ie8 lt-ie7"> <![endif]-->
```

There is around 60% of similarity in the keywords extracted using both the approaches. We also get the summary and HTML script of the article.

V. Application

1. News Aggregator: News Aggregator, is client software or a web application which aggregates syndicated web content such as online newspapers, blogs, podcasts, and video blogs (vlogs) in one location for easy viewing. These applications combine articles from different sources to avoid redundancy of information displayed on their portal. Hence, extracting keywords from the articles and comparing these keywords to check if certain percent of similarity exists can help us in combining articles from different sources.

2. Plagiarism detection: It is the process of locating instances of plagiarism within a work or document. The widespread use of computers and the advent of the Internet has made it easier to plagiarize the work of others. This work holds application in this area as well.

VI. Conclusion

We understood the importance and application of Big Data Analytics. We implemented different approaches to extract keyword and did a comparative study.

VII. References

- [1] <https://github.com/codelucas/newspaper>
- [2] <http://newspaper.readthedocs.io/en/latest/>
- [3] <https://stackoverflow.com/>
- [4] <https://www.python.org/>