

# Similar Content Retrieval From Wikipedia Citation

Rahul Gupta  
MT20065

Deepti Gupta  
MT20085

Palak Tiwari  
MT20103

Ravi Rathee  
MT20066

Mohit Ghai  
MT20308

## 1. PROBLEM DESCRIPTION

The drastic growth in the content available on the web has been observed in the past few years. Wikipedia saw exponential growth since 2006, and the total number of articles on the site reached 55 million, out of which 6.2 million are in English. This massive collection of articles poses new kinds of difficulties to users, such as retrieving semantically similar information. As lexically similar texts may not relate semantically like food and bread are semantically equivalent however lexically different. It is time-consuming to visit various cited links and read content manually when the data is vast and diverse. Even if we try reading content manually, not all cited text is relevant, leading to wastage of time. Sometimes the user needs context to the quoted text and requires to search manually for the cited text and may end up getting some relevant data or may deviate from the context. This problem needs the system to be more imaginative and efficient and increase the user readability and experience.

We aim to propose an information retrieval system that would retrieve semantically similar data from the noted link, noted link here means Wikipedia's highlighted internal links and external citations. We propose to make a chrome extension plugin which will return the most relevant article from citations based on user's query. Our approach will save users precious time by not surfing over and over for the cited context. It helps the user to extract the vital information more engagingly. It also avoids the situation when the user visits multiple articles for the cited text and may often deviate from the context. This will enhance the overall knowledge of the context of the user as well.

## 2. LITERATURE REVIEW

[1] **Semantic Cosine Similarity** Cosine similarity is a widely implemented metric in information retrieval and related studies. Metric tries to model the document as a vector of terms and similarity between documents can be calculated by finding cosine values between the document term vectors. Even in search engines, similarity between documents and user query are stored in descending order of similarity. Paper summarises related work done in three sections mainly, A) Cosine Similarity B) WordNet C) Semantic Similarity between Words and finally discusses about disadvantages of using Cosine Similarity. Paper also proposes two enhancements over previously used cosine formula. The two en-

hancements were A) Dimension equalization of both term vectors. B) If there exists a synonym pair in both vectors choose the first one as dimension name for both vectors. For example, rodent is chosen as dimension's name from a pair of "rodent-mouse".

[2] **Automatic extraction of semantic relations from Wikipedia**, - Recently Patrick Arnold and Erhard Rahm presented Automatic Extraction of Semantic Relations patterns from voluminous unstructured Wikipedia articles using a finite state machine. The machine is capable of parsing the first few lines of the wiki article. The proposed system only parses the first 750 words of the article to get the abstract information. The paper describes the system where input is a sentence with a list of predefined anchor terms. If the sentence contains an anchor term, FSM discovers the semantic pattern starting with the anchor term and ending in the final state. The system can capture at least two semantic ways. After the extraction of object terms, subject terms, and field terms, semantic relations are found. These relations can be used to make huge updated thesaurus as a knowledge for tasks such as retrieving semantic ontology mappings. The paper also discusses that method is capable of detecting entity articles from concept articles with the use of categories. For this, regular expressions are used to cover few types to discard pieces under such categories.

As discussed in the paper, FSM faced parsing errors due to the sentences being too complex. Recall problem in the term extraction due to the removal of complex expressions with parenthesis which had relevant terms.

[3] **Hybrid Geometric Approach for Measuring Similarity Level Among Documents and Document Clustering** - Propose a method for measuring similarity levels among documents and document clustering. This involves a few steps, which in order are: 1) Carefully choosing features to reflect underlying semantics. 2) Creating a Vector Space Model. 3) Comparing algorithms used to measure similarities between documents like Non Geometric measures and Geometric measures - cosine similarity and Euclidean Distance 4) Drawbacks of above algorithms are discussed.

In this paper, the authors discuss a novel approach to overcome the drawbacks of methods discussed earlier: the Novel geometric similarity measurement method.

This algorithm proposed by authors called the "TS-SS Algorithm" computes between vectors from two diverse perspectives. It generates the similarity value between two vectors from the angle and ED between them and the difference be-

tween their magnitudes.

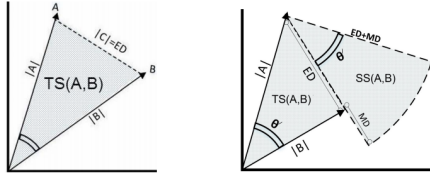


Figure 1: TS and SS Algorithm

This algorithm has two parts:-

A) TS Algorithm - used to generate different similarity values.

B) SS Algorithm (Sector's Area Similarity) TS alone is not sufficient to interpolate vector's differentiation precisely, and hence we need an SS algorithm. TS and SS complete each other, and that is why we combine them by multiplying them together.

[4] **A New Approach for Implicit Citation Extraction** - In this paper, the problem is divided into two and then conquered. Two subproblems are topic modeling and implicit citation detection (similarity detection between the text and the citation) and then proposed techniques are applied. LDA (Latent Dirichlet Allocation) is used for topic modeling to get the semantically similar context as lexical similarity might not always give appropriate contexts. For implicit citation detection, two-word embeddings (vectors) are determined. Sentence2vec and Topic2vec are used. However, both are determined using the concepts of Doc2vec. Unsupervised learning techniques are used as the Wikipedia dataset is not annotated. The paper also shows a real-world example to demonstrate the proposed method.

The limitation of this paper is, the dataset used is not multilingual, and also the accuracy of the proposed system is not determined.

[5] **Semantic Analysis of Wikipedia documents using Ontology** - With the drastic growth in web content, there is the need for better query results from the server-side. The authors proposed a system where they fetched the user query, extracted the critical term using an algorithm like KPE, and ranked it. For the next step, they applied Cosine Similarity. The Semantic Web will offer a way for solving the problem at the architecture level. Each page possesses semantic metadata that records additional detail concerning the Web page itself. For evaluation, they used precision, recall, f-measure, and accuracy. The results achieved by them were 0.73 Accuracy, 0.6 Precision, 0.87 recall, and 0.78 F-measure. However, the techniques used are unsupervised. Adding supervised methods for ontology learning may add to the results

[6] **BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding** - BERT or Bidirectional Encoder Representations from Transformers is an encoder only based Transformer network used for Natural Language Processing tasks. The major application of Bert is to extract features i.e. convert sentence into its vector representation. BERT uses attention mechanisms and con-

siders both left and right context of the given sentence. It is trained on two phases MLM (Masked Language Model) & NSP (Next Sentence Prediction). There are two types of BERT models BERT<sub>BASE</sub> (110 million parameters) & BERT<sub>LARGE</sub> (340 million parameters). We can either use pre-trained model of BERT or do fine-tuning based on our data. Other applications of BERT include Classification, Summarization, Question-Answer tasks etc.

[7] **Annotating documents by Wikipedia concepts** - This paper demonstrates what all Wikipedia pages in the target Wikipedia document are the most relevant. For this, they rank the pages as well. There are four types of Wikipedia pages connected to a single Wikipedia page denoted by the target document. Authors propose an algorithm that will retrieve the most semantically relevant pages in the link structures of pages. According to the algorithm, firstly title of the target document is cleaned. These concepts are currently ranked as per the formula stated in the paper, and the top-ranked pages are considered the most relevant.

[8] **Identifying Citing Sentences in Research Papers Using Supervised Learning** - Sugiyama et al. designed a method to predict whether a sentence requires citation or not using a supervised learning approach. They used ACL Anthology Reference Corpus. Extracting the features like Bigram, Unigram, Previous and Next sentence, Position, Proper Nouns, and Orthographies. Constructing the classifier Maximum Entropy and Support Vector Machine and classify the sentence whether it needs to be cited or not. They achieved an accuracy of 0.82 using both Maximum Entropy and SVM with  $C = 0.9$ . The results showed that for training accurate models in both ME and SVM, next sentence and proper nouns are effective features. The limitation is that the system doesn't consider any similarity between the sentences. We try to explore similarities using different techniques.

### 3. BASELINE RESULTS

We developed a chrome extension which fetches the query highlighted by the user and returns the most relevant paragraph from the citations present in the given query. For finding the relevant paragraph from the citations we use the Doc2vec model for the document and Word2vec for the query as mentioned in [4] For finding the similarity between the paragraph and the given query we use different similarity measures like Cosine Similarity, Word Movers distance

The system will get a cited text from Wikipedia as an input and will return the semantically relevant Paragraph from the cited articles as an output.

- **Input from FrontEnd** - User can highlight some part of the wikipedia text, and our plugin will take that highlighted text as input and other required details of the articles along with citations and hyperlinks in the articles as input.
- **Scraping data** - Query and title name is fetched using javascript and sends a post request to the flask server. Now using page title citations are scrapped from the current page. Citation occurring in the query is extracted. Using the extracted citations from the query, documents are scrapped. Each extracted articles are

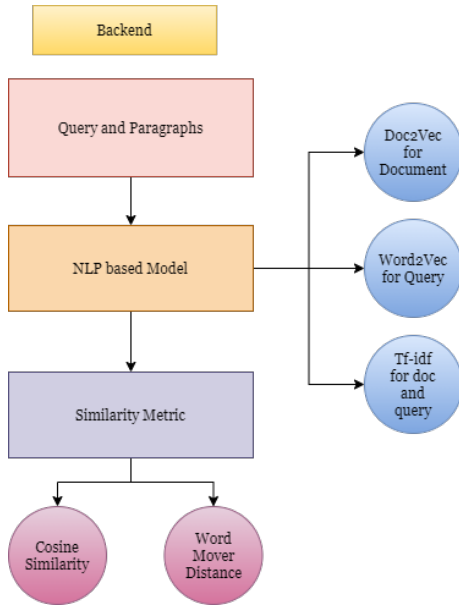


Figure 2: Backend Flow

appended together which will further be broken into paragraphs.

- **Cleaning** - Now the appended articles is broken down into paragraphs and we apply some pre-processing to clean data like, removal of junk symbols, tokenization, stop words, removing punctuation and applied lemmatization.
- **Embeddings** - We used three types of word embeddings as mentioned below:-
  - **Word2Vec** - Word2Vec is a combination of 2 techniques - CBOW (Continuous Bag of Words) and Skip-Gram (Skip -Gram model). It helps to keep semantic information and to preserve relation between different words. We get 300 dimensional embeddings for each word. Then we take mean for all the vectors and resultant query vector of 300 dimensions.
  - **Document2Vec** - It is used to create embeddings for group of words taken collectively. It represents each document as a vector. The document is splitted into words. Then tokens are generated. This token list is passed as input to doc2vec model. It gives a 300 length dimensions embeddings for each paragraph.
  - **Tf-idf** - We applied Tf-idf on both query and document after building the corpus from various wikipedia documents. After building the corpus we vectorized the input query and the document where for each paragraph present in the document different vectors are made which later on will be checked for similarity againsts the query vector.
- **Similarity Measures** - We applied similarity measuring technique like Cosine similarity, that is a metric to measure the similarity between document vector and query vector by determining the cosine angle and

Word Mover Distance to find the similarity between the document and the query.

- **Output to FrontEnd** - Most relevant results by cosine similarity are sorted in descending order and top results are returned, which potentially helps user in finding similar paragraphs. On using the Word Mover Distance also we get the most relevant result after sorting the results in descending order according to the scores.

Query: 'an American public organization that focuses on non-profit activities and creative media ventures.'

Embedding-distance metric	Top first Result	Top second Result	Top third Result
Pretrained Doc2Vec and word2Vec with cosine similarity	Archewell press secretary Toya Holmes told Town & Country, "Founded earlier this year by The Duke and Duchess of Sussex, Archewell uplifts communities through non-profit partnerships and creative activations. It's a place where compassion matters, communities gather, and storytelling is the engine. The website has been updated to reflect the work Archewell has undertaken throughout 2020 and to create a place for people and communities around the world to share their stories. USA railway vehicle having its own means of propulsion	On the Archewell Foundation page, Harry and Meghan have listed a series of partnerships and projects the non-profit is undertaking, including those with James R. Doty's Centre for Compassion and Altruism Research and Education (CCARE), Tristan Harris' Centre for Humane Technology, Rachel Cargle's Loveland Foundation, Dr. Safiya Noble and Dr. Sarah T. Roberts' UCLA Centre for Critical Internet Inquiry, and José Andrés' World Central Kitchen.	And finally, the Archewell Productions page contains a description of the couple's new television production company, which was launched after Harry and Meghan inked a landmark deal with Netflix earlier this year. "Archewell Productions will utilize the power of storytelling to embrace our shared humanity and duty to truth through a compassionate lens," the page reads.
Pretrained Doc2Vec and word2Vec with wm distance	On the Archewell Foundation page, Harry and Meghan have listed a series of partnerships and projects the non-profit is undertaking, including those with James R. Doty's Centre for Compassion and Altruism Research and Education (CCARE), Tristan Harris' Centre for Humane Technology, Rachel Cargle's Loveland Foundation, Dr. Safiya Noble and Dr. Sarah T. Roberts' UCLA Centre for Critical Internet Inquiry, and José Andrés' World Central Kitchen.	And finally, the Archewell Productions page contains a description of the couple's new television production company, which was launched after Harry and Meghan inked a landmark deal with Netflix earlier this year. "Archewell Productions will utilize the power of storytelling to embrace our shared humanity and duty to truth through a compassionate lens," the page reads.	The Archewell Audio page features the first episode of Archewell Audio, Harry and Meghan's new podcast with Spotify, as well as a short description of the production company's aims: to "produce programming that uplifts and entertains audiences around the world," and "spotlight diverse perspectives and voices," as well as "build community through shared experiences, powerful narratives, and universal values."
Tf-idf with cosine similarity	'Just ahead of the new year, Prince Harry and Meghan Markle have updated the website for their new non-profit, Archewell. The site was launched in October with just a landing page, but now, the couple has fleshed it out, adding pages for the Archewell Foundation, Archewell Audio, and Archewell Productions.'	I am my mother's son. And I am our son's mother. Together we bring you Archewell, "the couple wrote on the homepage for their new non-profit.	Archewell press secretary Toya Holmes told Town & Country, "Founded earlier this year by The Duke and Duchess of Sussex, Archewell uplifts communities through non-profit partnerships and creative activations. It's a place where compassion matters, communities gather, and storytelling is the engine. The website has been updated to reflect the work Archewell has undertaken throughout 2020 and to create a place for people and communities around the world to share their stories."

Figure 3: Baseline results

Figure 3 shows the results from different variations on inputting the same query, the first result is from the Doc2vec on document and Word2vec on query with Cosine similarity as metric the next is with WM distance metric followed by Tf-idf Vectors on both query and document with Cosine Similarity Metric.

To validate our result input document itself contained query as one of the paragraph. When this input set was fed to the model, query text was output (most relevant paragraph).

## 4. PROPOSED METHOD

We propose a system which includes the following functionalities:

- Extracting most contextually similar paragraphs from the cited articles from the query.
- Recommending Contextually similar articles from the current article.

For the given input query we fetch the citations present in the query, both external and internal links, here internal means links that are from some wikipedia page itself and external means those from internet that are cited by the authors. We scrap document id's for the internal links. Then we used Doc2vec model for the documents to get the embeddings. For query we used Word2vec to get the vector embeddings. For finding the similarity between the document vectors and the input query we used Cosine Similarity

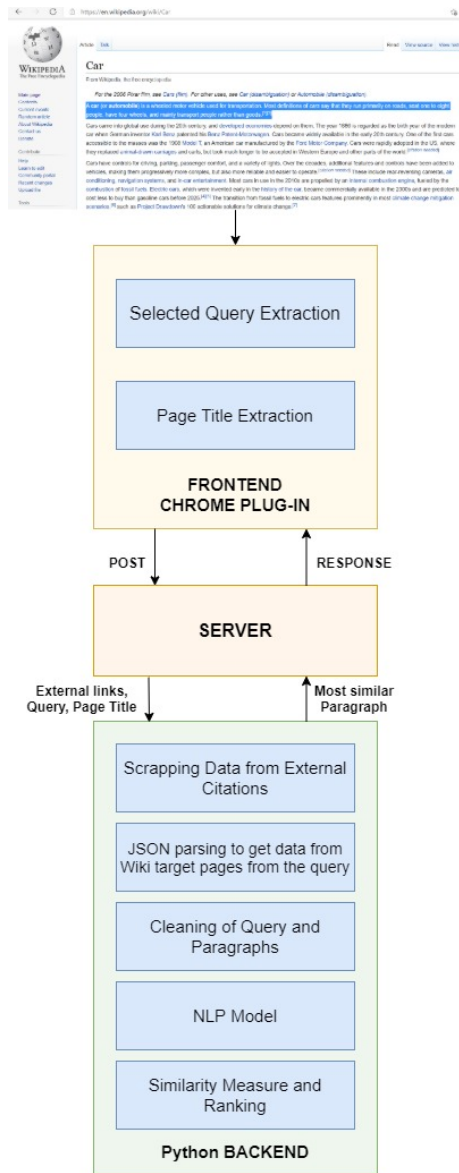


Figure 4: Pipeline Flow

measure and WMD (Word Mover's Distance).

For Recommending Contextually similar articles to the user for the query or for the current page depending on users choice, for this topic modeling can be used. This can also be done by finding the similarity of the introductions paragraphs of different articles and the given query or the introduction paragraph of the current page and thereby ranking those articles to get top five most relevant articles for users to read. For this, scrapping of the topic and extraction of different topics from the the Wikipedia dumps. Creating vectors of the current and other articles using above models that is Word2Vec and Doc2Vec. Finally finding similarity between the vectors and ranking them to get the most similar topics to recommend.

## 5. REFERENCES

- [1] F. Rahutomo, T. Kitasuka, and M. Aritsugi, "Semantic cosine similarity," 2012.
- [2] P. Arnold and E. Rahm, "Automatic extraction of semantic relations from wikipedia," *International Journal on Artificial Intelligence Tools*, vol. 24, p. 1540010, 04 2015.
- [3] A. Heidarian and M. J. Dinneen, "A hybrid geometric approach for measuring similarity level among documents and document clustering," in *2016 IEEE Second International Conference on Big Data Computing Service and Applications (BigDataService)*, 2016, pp. 142–151.
- [4] C. Jebari, M. J. Cobo, and E. Herrera-Viedma, "A new approach for implicit citation extraction," in *Intelligent Data Engineering and Automated Learning – IDEAL 2018*, H. Yin, D. Camacho, P. Novais, and A. J. Tallón-Ballesteros, Eds. Cham: Springer International Publishing, 2018, pp. 121–129.
- [5] P. Banik, S. Gaikwad, A. Awate, S. Shaikh, P. Gunjgur, and P. Padiya, "Semantic analysis of wikipedia documents using ontology," in *2018 IEEE International Conference on System, Computation, Automation and Networking (ICSCA)*, 2018, pp. 1–6.
- [6] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," *CoRR*, vol. abs/1810.04805, 2018. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [7] K. Sugiyama, T. Kumar, M. Kan, and R. C. Tripathi, "Identifying citing sentences in research papers using supervised learning," in *2010 International Conference on Information Retrieval Knowledge Management (CAMP)*, 2010, pp. 67–72.
- [8] P. Schönhofen, "Annotating documents by wikipedia concepts," in *2008 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*, vol. 1, 2008, pp. 461–467.