



iJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 6 Issue: VI Month of publication: June 2018

DOI: <http://doi.org/10.22214/ijraset.2018.6004>

www.ijraset.com

Call: ☎ 08813907089

E-mail ID: ijraset@gmail.com

Ontology Based Semantic Document Clustering Using LDA Algorithm

Jheel Bagani¹, Prof. Dr. Abha Chaubey²

^{1, 2}Shri Shankaracharya group of Institution, Dept. of Computer Science and Engineering, Bhilai, Chhattisgarh, India

Abstract: Document clustering is the process of clustering or dividing the documents contents into various sub parts through which each of the parted contents can be identified individually. The clustered documents are of collection of similar semantics information grouped together. In this paper we proposes a Nobel method using LDA algorithm for ontological document clustering.

Keywords: Clustering, document clustering, semantic document clustering.

I. INTRODUCTION

Today we are facing a consistently expanding volume of text documents. The voluminous texts streaming over the Internet, gigantic collection of documents in computerized libraries and repositories and computerized data, for example, blogs and emails are growing up quickly day by day. It brings the difficulties that how to arrange text documents successfully and proficiently. The issue of clustering has been considered generally in the database and statistics writing in the context of a wide variety of data mining tasks. The clustering issue can be characterized as discovering gatherings of similar objects in the data. The closeness between two unique objects is estimated with the use of a comparability function. The issue of clustering can be exceptionally profitable in the text area, as the items to be clustered can be of various granularities, for example, documents, paragraphs, sentences, terms and so on. So as to apply the majority of the clustering algorithms, two things are required: speaking to a protest, and a likeness measure between objects. A clustering algorithm finds a parcel of an arrangement of articles that satisfies some standard in view of a comparability measure.

The text clustering is the issue of consequently grouping of free text documents. The groups are ordinarily depicted by an arrangement of keywords or expressions that portrayed the regular substance of the documents in the group. To perform a clustering procedure, the articles ought to have some sort of attributes to quantify the separation or likeness among the objects. These qualities are typically called as features of the object. The majority of the proposition in this field consider the document as an arrangement of words. In these portrayals, each include relates to a solitary word found in the document set. As a document set may contain a few thousand of words, these outcomes in a high impracticable dimensionality. To diminish the document space dimensionality some word diminishment strategies are connected in the pre-processing phase. The most widely recognized strategy to diminish the quantity of various words is to wipe out the words with low data value. These words are called stop words. The stop words are gathered into a dictionary or then again a list. Another path for diminishment depends on the factual properties of the words: the occasional and the visit words are filtered out from the original text.

II. PROPERTIES OF TEXT DOCUMENTS

Naive systems don't commonly function admirably for clustering text data because text data has various unique properties which requires specific algorithms. The recognizing qualities of text portrayal are as takes after:

- A. The dimensionality of the text portrayal is extremely thick while the basic data is sparse. In other words, the lexicon from which documents are drawn might be of the order of 100, however the document itself may contain just a couple of hundred words. This issue turns out to be significantly more genuine when a document to be clustered is short, for example, sentences or tweets.
- B. While the lexicon of a given document might be vast, the words are ordinarily related with each other. It implies that the quantity of key parts in data is substantially littler than the element space. This requires the watchful outline of algorithms which can represent word connections in the clustering procedure.
- C. The quantity of words in various documents may change broadly and subsequently it is fundamental to standardize the document portrayals properly amid clustering assignment.

The sparse and high dimensional portrayal of distinctive documents requires the plan of text-particular clustering algorithms for document representation and processing. On the off chance that the quantity of words is equivalent to N , the number of expressions containing k words is N^k . Accordingly, all together to diminish computational cost, particular algorithms are required for various periods of the text document clustering process. Numerous methods have been proposed to streamline document portrayal for enhancing the precision of coordinating a document with a query [1].

III. PROPERTIES OF TEXT CLUSTERING

The significant research in text clustering has been done in the context of following two sorts of text data: Dynamic Applications and Heterogeneous Applications. Today, a vast measure of text data is being made by unique applications, for instance, informal organizations or online talk applications. Such streaming applications must be pertinent in the event of non-filtered text. Text applications are progressively emerging in heterogeneous applications where the text is accessible in the context of hyperlinks and other heterogeneous sight and sound data.

The cluster hypothesis recommends that "nearly related documents have a tendency to be important to a similar request" [2], i.e. comparable documents are accepted to be pertinent to the same queries put to a web index. This has made numerous scientists trust that a valid clustering could make seek time shorter as bunches could be recovered rather than documents. [3, 4] have demonstrated an effective method to cluster web index comes about. The internet searcher Vivisimo1 uses clustering on recovered documents.

The Scatter/Gather framework (or any clustering strategy) has additionally been proposed for perusing document collections [6]. A document collection is introduced to a user as a set of clusters. The user may check one or a few bunches for promote examination and demand that these are re-clustered giving an all the more adjusted grouping. Along these lines, the user may iteratively and intelligently investigate the collection and get a review of its substance and additionally discover specific subjects that show up in it. Clustering techniques can be used to consequently group the recovered documents into a list of important classifications, as is accomplished by Enterprise Search motors, for example, Northern Light and Vivisimo or open source programming such as Carrot2. Likewise, Google is known to use clustering strategies to coordinate certain sites with a query, since the site can be seen as a collection of subjects (multi-point document), and a query itself is a theme or a blend of a few themes.

At long last, with the ascending of informal communities as of late, for example, Facebook and Twitter, more semantic data are accessible now that pass on a lot of data. On Twitter, there are surmised 95M tweets every day, which is proportional to 1100 tweets for every second. Scientists from the Northeastern University School of Computer and Information Sciences, and Harvard Medical School has built up a creative way of following the country's state of mind utilizing tweets. All these scientists demonstrate the energy of social processing in giving precise appraisals on numerous sorts of issues, at no cost and on a huge scale. Moreover, document clustering methods can be used to bunch tweets into pertinent themes, in help of the present negligible 'Patterns' function used by Twitter. For all of these reasons, discover document clustering systems significant and in this way worth considering.

IV. LITERATURE SURVEY

Tamara G. Kolda et al. [1], the huge measure of textual information accessible today is pointless unless it can be successfully and productively sought. The objective in information retrieval is to discover documents that are pertinent to a given client query. Author can speak to and document collection by a lattice whose (i, j) section is nonzero just if the i th term shows up in the j th report; hence each document relates to a column vector. The query is additionally spoken to as a column vector whose i th term is nonzero just if the i th term shows up in the query.

M. Andrian et al. [2], an information retrieval strategy, latent semantic indexing, is utilized to naturally distinguish traceability joins from framework documentation to program source code. The consequences of two investigations to recognize interfaces in existing programming frameworks (i.e., the LEDA library, and Albergate) are displayed.

Jen-Yuan Yeh et al. [3], this paper proposes two ways to deal with address text summarization: changed corpus-based approach (MCBA) and LSA-based T.R.M. approach (LSA + T.R.M.). The first is a trainable summarizer, which considers a few features, including position, positive keyword, negative keyword, centrality, and the likeness to the title, to produce outlines.

Yoshihiko Gotoh et al. [4], in this paper, an approach for building blend dialect models (LMs) based on some thought of semantics is examined. To this end, a strategy known as latent semantic examination (LSA) is utilized. The approach epitomizes corpus determined semantic information and can show the shifting style of the text. Utilizing such information, the corpus texts are grouped in an unsupervised way and blend LMs are consequently made.

Chun-Ling Chen et al. [5], with the quick development of text documents, document clustering has turned out to be one of the fundamental systems for sorting out vast measure of reports into few significant clusters. Be that as it may, there still exist a few

difficulties for report clustering, for example, high dimensionality, scalability, exactness, and significant group names, covering clusters, and separating semantics from texts.

Jeroen De Knijff et al. [6], this paper proposes a structure to consequently develop scientific categorizations from a corpus of text reports. This system first concentrates terms from reports utilizing a grammatical form parser. These terms are then sifted utilizing area congruity, space agreement, lexical union, and auxiliary significance. The rest of the terms speak to ideas in the scientific classification.

Andreas Hotho et al. [7], Text document clustering assumes a critical part in giving natural route and perusing instruments by sorting out huge arrangements of reports into few important clusters. The sack of words portrayal utilized for these clustering strategies is regularly inadmissible as it overlooks connections between imperative terms that don't occur truly. Keeping in mind the end goal to manage the issue, we incorporate center ontologies as foundation learning into the way toward clustering text documents.

Tingting Wei et al. [8], customary clustering algorithms don't consider the semantic connections among words so that can't precisely speak to the significance of reports. To defeat this issue, presenting semantic information from cosmology, for example, WordNet has been broadly used to enhance the nature of text clustering. Notwithstanding, there still exist a few difficulties, for example, equivalent word and polysemy, high dimensionality, separating center semantics from texts, and relegating suitable depiction for the produced clusters.

Guoyu Tang et al. [9], Cross-lingual document clustering is the undertaking of consequently arranging a substantial collection of multi-lingual documents into a couple of clusters, contingent upon their substance or theme. It is notable that dialect boundary and interpretation vagueness are two testing issues for cross-lingual document portrayal.

Shibamouli Lahiri et al. [10], Keyword and key phrase extraction is an imperative issue in regular dialect preparing, with applications extending from summarization to semantic hunt to document clustering. Diagram based ways to deal with keyword and key phrase extraction evade the issue of gaining a substantial in-area preparing corpus by applying variations of PageRank algorithm on a system of words.

V. METHODOLOGY

In this segment, the proposed workflow is presented well ordered. The Fig. 1. Demonstrates the workflow of every module.

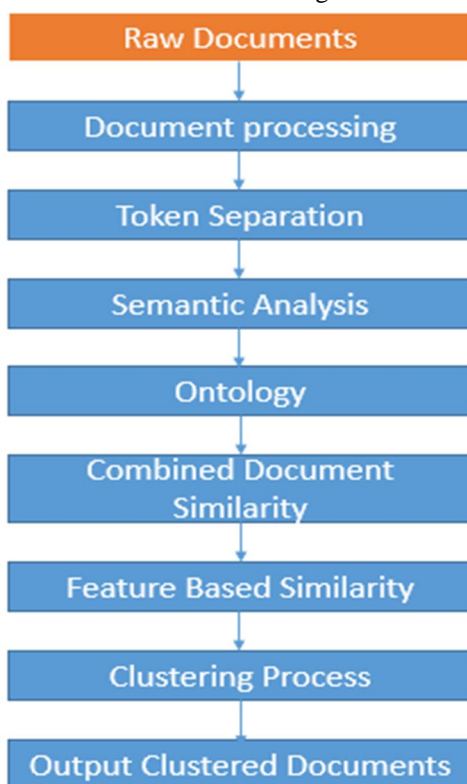


Fig. 1. Porposed System Architecture

A. Raw Documents

The twitter dataset are taken as the raw document. The twitter is taken in light of the fact that it contains short text content and can be valuable for classification of subjects.

B. Document Processing

This steps usually process and remove unnecessary content from the tweets. It cleans the tweets in short.

C. Token Separation

After the cleanup step, the tweets are further separated into tokens. The tokens represent the smallest individual unit of tweet.

D. Semantic Analysis

After that the semantic analysis are done with a specific end goal to maintain the semantic meaning of tweets.

E. Similarity Calculation

After the separation and semantic analysis of tokens, the similarity of every word in the tweets are compared with the other tweets. The same records has smallest or almost equal similarity score.

F. Clustering

The clustering of records depend on weight of the node tree. The tree represent weights consequently setting up weight depend on the below formula. The clustering formula experimentally expressed by:

$$w_{k,i} = \frac{\text{height } M_i - j + 1}{\text{height } M_i + 1}$$

Where, $w_{k,i}$ is the weight of an attribute k on the model i , $\text{height } M_i$ Is height of the model i , and j is the level of attribute node k .

VI. RESULT AND DISCUSSION

In this section, the outcome are displayed and talked about it in detail. The dataset are taken from the twitter site. The dataset contains the raw data of tweets. The depiction of dataset is shown in fig.2.

```
1 When it comes to swimming in poop water , I do n't
  think there 's such a thing as overthinking it .
2 There 's that distinction I was looking for ... LOL !
3 Spent summers in Michigan , I know how that goes ...
  but a lake vs. a gator kiddie pool are NOT the same !
4 Ewww . You can keep your gator bathwater all to
  yourself !
5 Nope , they need all that water , for one ... I bet
  cleaning them is a delightful chore !
```

Fig. 2. Shows the snapshot of sample tweets

Different modules output are displayed below.

A. Tokenization of Tweets

The tweets of the users are scanned and are tokenized to form group of different words.

```
angiefeimao
animus128
annareihl
annisarizki
aqjq
arjunms
ArtiztikVizion
AshleyBenlove
ashleykristen
AwkwardAI88
ayoo_chinadol1
AyYo_Jay
badhex
```

Fig. 3. Shows the snapshot of Tokenized Tweets

B. Document Clustering

After the preparing the records are clustered. The clustering algorithm makes 12 diverse cluster for different 12 subjects.

```
askin/29 slang/29
racetrac/58 aww/false
hey/false hey/58 hey/58 meet/false hoe/58
dykin/58 boosie/58 playing/false head/false doe/58
chick/75 fil/75 clutch/75
man/false fuck/58 pj/58
daniel/58 sleep/false n't/false spell/58 baps/58
fetti/58 gettin/false stacks/58 top/false stacks/58
vae/88 lol/88
```

Fig. 4. Shows the cluster output

VII. CONCLUSION

This paper help to examine the noteworthy part of document for powerful clustering and mining. There are variety of content mining applications, which contains data inside them, this data might be of various types, for example, source of data of the documents, web logs, the connections in the documents which contains client access behavior. A considerable measure of work has been done in present days on the issue of clustering in content collections in the database and data recovery. All things considered, this work is normally intended for issue of pure content clustering in the absence of different sort of characteristics. These attributes may likewise having a great deal of information for clustering goals.

In this paper, we have proposed a novel system for twitter data clustering. The outcomes are presented as clustered tweets having similar subjects.

REFERENCES

- [1] Tamara G. Kolda, and Dianne P. O'leary. "A semidiscrete matrix decomposition for latent semantic indexing information retrieval." ACM Transactions on Information Systems (TOIS) 16, no. 4 (1998): 322-346.
- [2] Andrian Marcus, and Jonathan Maletic. "Recovering documentation-to-source-code traceability links using latent semantic indexing." In Software Engineering, 2003. Proceedings. 25th International Conference on, pp. 125-135. IEEE, 2003.
- [3] Jen-Yuan Yeh, Hao-Ren Ke, Wei-Pang Yang, and I-Heng Meng. "Text summarization using a trainable summarizer and latent semantic analysis." Information processing & management 41, no. 1, Elsevier (2005): 75-9
- [4] Yoshihiko Gotoh, and Steve Renals. "Document space models using latent semantic analysis." (1997).
- [5] Chun-Ling Chen, Frank SC Tseng, and Tyne Liang. "An integration of WordNet and fuzzy association rule mining for multi-label document clustering." Data & Knowledge Engineering 69, no. 11, Elsevier (2010): 1208-1226.
- [6] Jeroen De Knijff, Flavius Frasincar, and Frederik Hogenboom. "Domain taxonomy learning from text: The subsumption method versus hierarchical clustering." Data & Knowledge Engineering 83, Elsevier (2013): 54-69.
- [7] Andreas Hotho, Steffen Staab, and Gerd Stumme. "Ontologies improve text document clustering." In Data Mining, 2003. ICDM 2003. Third IEEE International Conference on, pp. 541-544. IEEE, 2003.
- [8] Tingting Wei, Yonghe Lu, Huiyou Chang, Qiang Zhou, and Xianyu Bao. "A semantic approach for text clustering using WordNet and Lexical chains." Expert Systems with Applications 42, no. 4, Elsevier (2015): 2264-2275.
- [9] Guoyu Tang, Yunqing Xia, Erik Cambria, Peng Jin, and Thomas Fang Zheng. "Document representation with statistical word senses in crosslingual document clustering." International Journal of Pattern Recognition and Artificial Intelligence 29, no. 02, World Scientific (2015): 1559003.
- [10] Shibamouli Lahiri, Sagnik Ray Choudhury, and Cornelia Caragea. "Keyword and keyphrase extraction using centrality measures on collocation networks." arXiv preprint arXiv: 1401.6571 (2014).



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)