# Review on Phase Based Semantic Search Using Suffix Tree Clustering

Priyanka[1], Umesh Goyal[2]

P.G. Student, Department of Computer Science & Engineering, Delhi Institute of Technology Management and Research, Faridabad, Haryana, India[1]

Assistant Professor, Department of Computer Science & Engineering, Delhi Institute of Technology Management and Research, Faridabad, Haryana, India [2]

**ABSTRACT:** In the substance mining area, noticeable techniques use the sack of-words models, which address a record as a vector. These methods ignored the word gathering information, and the incredible packing result obliged to some remarkable spaces. This paper proposes review of another closeness measure reliant on the postfix tree model of substance files. It separates the word progression information and after that figure the likeness between the substance records of the corpus by applying a postfix tree closeness that solidifies with TF-IDF weighting system. Preliminary outcomes on standard record benchmark corpus that exhibit that the new message similarity measure is convincing. Differentiating and the outcomes of the other two progressive word gathering based procedures, our proposed system achieves an improvement of about 15% on the typical of F-Measure score**.**

**KEYWORDS**:  Crawler, Term Frequency–Inverse Document Frequency, Clustering, Document Model, Similarity Measure.

## I.  INTRODUCTION

A postfix tree is a data structure that yields gainful string planning and addressing. Expansion trees have been examined and used broadly, and have been associated with key string issues, for instance, finding the longest repeated, inferred string matches, string relationships, and substance weight. Following, we delineate the expansion tree data structure – its definition, improvement figurings and rule characteristics. The going with the portrayal of expansion tree takes after Dan Gusfield's exceedingly endorsed book on strings, trees, and groupings. One critical difference is that we treat chronicles as groupings of words, not characters. A postfix tree of a string is basically a diminished trie of the considerable number of augmentations of that string. In logically definite terms. A postfix tree T for an m-word string S is a set up composed tree with decisively m leaves numbered 1 to m. Each inside center, other than the root, has at any rate two children and each edge is named with a nonempty sub-arrangement of articulations of S. No two beats of a center point can have edge imprints beginning with a comparable word. The key part of the expansion tree is that for any leaf I, the association of the edge blemishes in transit from the root to leaf I unequivocally spells out the postfix of S that starts at the position I, that is it spells out S[i..m]. In circumstances where one postfix of S arranges a prefix of another expansion of S by then no postfix tree consenting to the above definition is possible since the route for the essential expansion would not end at a leaf. To avoid this, we acknowledge the last articulation of S does not show up wherever else in the string. This shields any postfix from being a prefix to another expansion. To achieve this we can incorporate a completion character, which isn't in the language that S is taken from, beyond what many would consider possible of S. Definition: The name of a center in the tree is portrayed as the connection, all together, of the sub-strings naming the edges of the route from the root to that center.

In any case, papery reports are experiencing a change to electronic record bit by bit. This change is irreversible on the grounds that electronic archives are more secure and simpler to spare and use than papery records. Each association has a database, which contains enormous volumes of electronic archives. The Word Wide Web is such a database, and

how to look and use this sort of content databases is a hot research point. These necessities animate us to create relating strategies to assist the clients with browsing and sort out these electronic reports increasingly compelling. A definitive objective is to assist the client in getting what the individual needs from huge data sea. Content bunching is known as a solo and natural method of collection content archives, so those comparative records will be assembled into a group, however not those unique. Subsequently, how to characterize a more precision content comparability to bunch content archives is basic to the dire necessity with the present data society. There are broadly written works on estimating the closeness between the writings [2][5]. The TF-IDF [6] (Term Frequency-Inverted Document Frequency) model is a well-known portrayal model of content records. The closeness between two content records is processed with one of a few likeness estimates dependent on two vectors, for example, cosine, Jaccard, and Euclidean separation measure. The real disadvantages of the TF-IDF strategy are as per the following: Firstly, information examination turns out to be progressively troublesome with the expanding of measurements while utilizing the customary word recurrence breaking down strategies. Also, this strategy overlooks word successions data of reports. So as to comprehend record all the more precisely, building up a similitude measure that contains highlights, which are increasingly instructive, has gotten impressive consideration as of late. The similitude among content archives has numerous applications in characteristic language preparing, data recovery, and content mining. For instance, in page recovery of the web search tool, content similitude has been demonstrated probably the best strategy for improving recovery aftereffects of a web index [1]. Moreover, web crawlers use content recovery innovation to rank query item archives as indicated by comparability among reports from high to low. The utilization of content comparability additionally can be a commitment to a content arrangement [2], copy location in website pages [3], report synopsis [4], and so forth. The inspiration of receiving postfix tree model for archive grouping can be ascribed to two angles. The first is the interest of dimensionality decrease for the content model. In the sack of-words model, it generally has a tremendous dimensionality, and definitely results from meager vectors of content archives. The subsequent one is that additions of the archive can incorporate more data than word recurrence. A successive postfix is a lot of individual words that incorporates more theoretical and relevant implications than an individual word. To address these contentions, this paper endeavors to propose a content similitude measure dependent on addition tree model, which improve viability with word grouping data in records contrasting with conventional word recurrence technique.
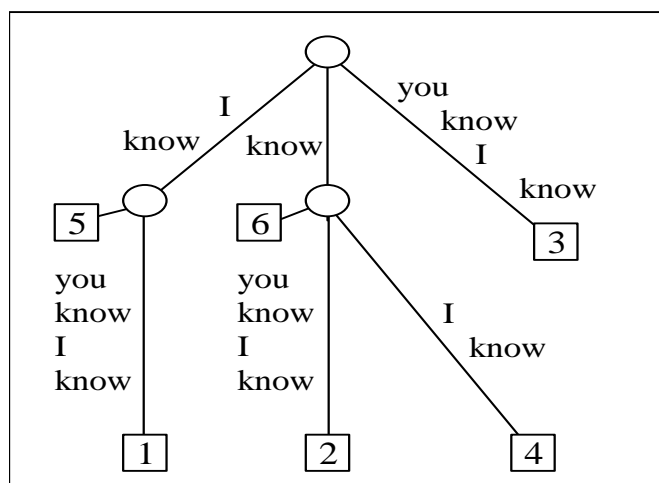


Figure 1: Example of a suffix tree, the suffix tree of the string *"I know you know I know"*. There are six leaves in this example, marked as rectangles and numbered from 1 to 6. The terminating characters are not shown in this diagram

## II. RELATED WORK

There are two general orders of substance resemblance procedures: word repeat/vector-based document model techniques and progressive word game plan based record model.

A. TF-IDF Model TF-IDF is a technique that measures the repeat of a term in a report with a factor that restrains its hugeness for its appearance in corpus. To depict, figure 1 (balanced from [7].) shows a great deal of model chronicles, with the relating term-report cross section depiction. In the system, the lines contrast with the watchwords, while the segments address the documents. The nonzero cells exhibit undeniably what terms appear in each report. At runtime, each substance report is changed into a comparative depiction, in order to use cosine or Jaccard coefficient closeness measures to enlist the substance resemblance. The TF-IDF model ignores words gathering and structure of reports. Additionally, proportion of words and chronicles are continually monstrous for a huge segment of record databases. With TF-IDF model, we should process a vector set, which has proportion of vectors and each vector has an estimation mean words number, and thusly unavoidably prompts a lower gainful enrolling. Regardless of the way that TF-IDF method has been shown to be extraordinary basically, the gigantic file data require exact portrayal rather than just term repeat. The possibility of the consistent word set relies upon the relentless thing set of the trade enlightening file. The progressive solicitation of words in a record accept a key activity of passing on the noteworthiness of the document. The distinction in the general spots of two words may change the substance of a chronicle. For example, "alliance rule" is a huge thought of data mining. In case these two words, "association" and "rule", appear in a pivot demand inside a report, like "The standard of their connection is ..." that addresses a through and through various hugeness. In any case, both of these records will be identical to each other by using TF-IDF system.

B. Frequent Word Sequence Document Model : For the reason that report addressed by perpetual word courses of action can give indications of progress sway than normal words, [8][9] analyzed progressive word groupings, and proposed a postfix tree methodology for substance bundling. The expansion tree serves the document as a string rather than a pack of words. It inspected the sharing file part by making a postfix tree. Nevertheless, this technique did not diminish the estimation of document model and did not consider semantic information of ceaseless words in reports. [10] further inspected and discussed ordinary word progressions subject to expansion tree, and showed documents as a model of nonstop word courses of action, which diminished components of report model suitably. [11] vanquished the detriments of the above research works had that two terms can have a comparable repeat in their reports, yet one term contributes more to the significance of its sentences than the other term. It separated activity word structures of the sentence and combined standard TF-IDF development and given different thoughts with different weighting. Regardless, an issue of these counts is that none of them portrayed a utilitarian equivalence to evaluate the closeness of inquiry and reports. [12] proposed a postfix tree comparability model for gathering count, eventually, their strategy did not have any huge bearing pruning procedure to decrease the component of the expansion tree model. [13] proposed Maximal Frequent Sequences (MFS) used for substance grouping. A progressive gathering is maximal if it's definitely not a subsequence of some other unremitting plan. The standard idea of MFS is to use maximal ordinary groupings of reports as substance features in the vector space model and after that used k-plans to pack files. MFS is a system that changes over word repeat to relentless word gathering in record grouping. Its display depends upon the reasonability of using MFS for record depiction in gathering, and the ampleness of k-suggests. [14] proposed another Frequent Term-Based Clustering (FTC) procedure for record gathering. The motivation of FTC is to make record bundles that spread between gatherings as few as would be judicious. FTC starts with an unfilled set, and it continues picking one bundle delineation from the course of action of remaining progressive word groupings until all the ordinary word progressions are picked. In every movement, FTC picks one of the remaining nonstop word groupings that have the tiniest entropy spread (EO) regard. In FTC, a gathering up-and-comer is addressed by an unending word progression and each up-and-comer's EO is resolved. In this manner, FTC will by and large select gathering candidate as a file bundle, of which its number of reports is short time occasion frequencies of these records are tremendous. In any case, it will cause colossal whole bunches with only a couple of chronicles, routinely a pack has one isolated record.

| Document | Content |
|---|---|
| $d_1$ | Human machine interface for Lab ABC computer applications |
| $d_2$ | A survey of user opinion of computer system response time |
| $d_3$ | The EPS user interface management system |
| $d_4$ | System and human system engineering testing of EPS |
| $d_5$ | Relation of user-perceived response time to error measurement |
| $d_6$ | The generation of random, binary, unordered trees |
| $d_7$ | The intersection graph of paths in trees |
| $d_8$ | Graph minors IV: Widths of trees and well-quasi-ordering |
| $d_9$ | Graph minors: A survey |

| Term | Documents | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ | $d_7$ | $d_8$ | $d_9$ |
| human | .58 | 0 | 0 | .49 | 0 | 0 | 0 | 0 | 0 |
| interface | .58 | 0 | .57 | 0 | 0 | 0 | 0 | 0 | 0 |
| computer | .58 | .44 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user | 0 | .32 | .42 | 0 | .46 | 0 | 0 | 0 | 0 |
| system | 0 | .32 | .42 | .72 | 0 | 0 | 0 | 0 | 0 |
| response | 0 | .44 | 0 | 0 | .63 | 0 | 0 | 0 | 0 |
| time | 0 | .44 | 0 | 0 | .63 | 0 | 0 | 0 | 0 |
| EPS | 0 | 0 | .57 | .49 | 0 | 0 | 0 | 0 | 0 |
| survey | 0 | .44 | 0 | 0 | 0 | 0 | 0 | 0 | .63 |
| trees | 0 | 0 | 0 | 0 | 0 | 1 | .71 | .51 | 0 |
| graph | 0 | 0 | 0 | 0 | 0 | 0 | .71 | .51 | .46 |
| minors | 0 | 0 | 0 | 0 | 0 | 0 | 0 | .7 | .63 |

Figure 2. Sample corpus (adapted from [7]).

### III. PROPOSED ALGORITHM

The proposed scenario identified with stage set up together semantic looking concerning web that offer customer to more results related to that subject. Customer may pick any appropriate result related to that subject. Using this investigation paper customer will get more inquiry results that were concealed in view of nonattendance of customer data. Normally customer sends some picked inquiry to web searcher and as shown by that request, web crawler sends some picked results related to that request so around then various information are concealed. With the help of this paper customer will get progressively semantic sentences results related to that request. On reason of these request result customer may pick anyone of them. In the wake of picking that question it sends to the web search device and recoups reasonable result.

As web list have million of site pages related to any request anyway more results are waste in light of the fact that customer does not send that question while they are proportionate expression of related words. To giving best result this paper give better way pick the result.
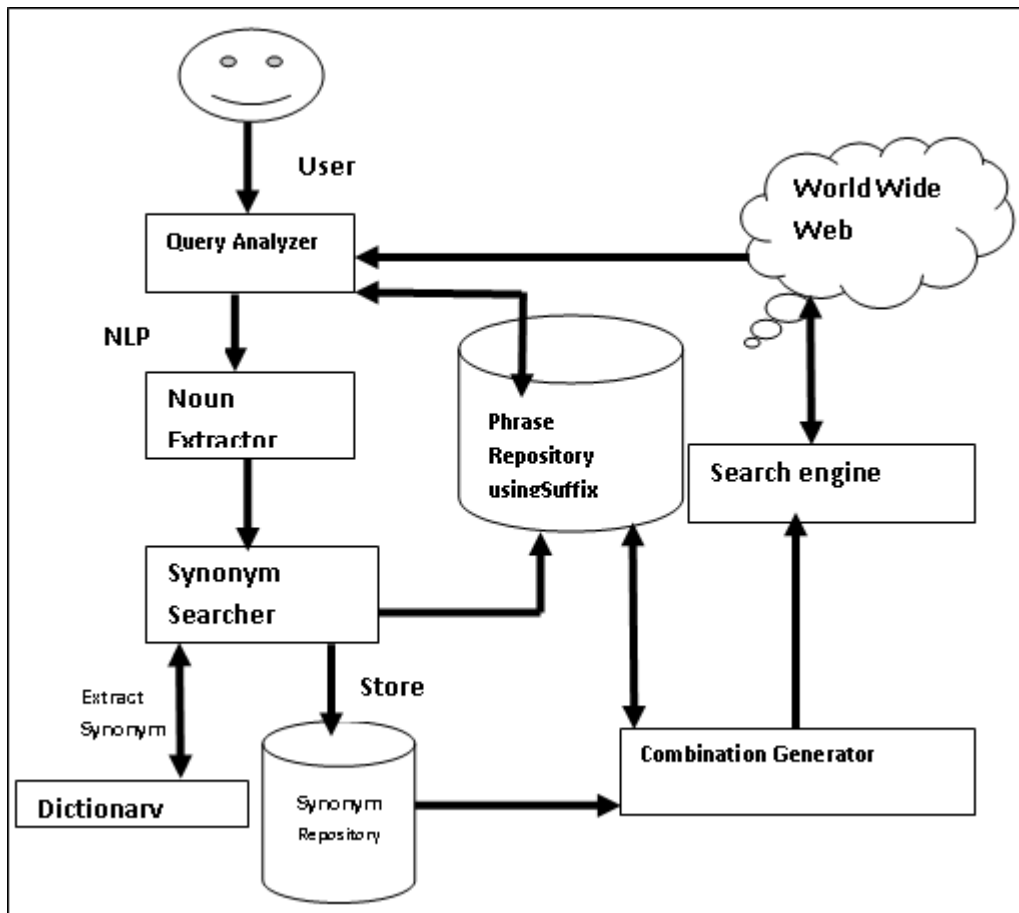
.

The proposed workflow is as under :-



Figure 3: Proposed Workflow comprising of Query Analyzer, Phase Repository based on Suffix Extraction, Noun Extractor, Synonym Searcher, Dictionary for Phase based Semantic Search Using Suffix Tree Clustering.

## IV. CONCLUSION

This proposition displays the keyword-based semantic search system search along-with term related to that keyword but user In catchphrase based semantic chase structure search simply terms related to that watchword yet customer generally search there result in sort of stage. Along these lines, word-based semantic structure shell that word phrase. This assessment paper urges the architecture to look in kind of stage so the user may glance through anything related to that term. As this paper proposes and relies upon semantic interest so the user may get many related to that arrange. Basically, this paper used the proportional word related to that organization. On occasion, the user has some specific learning related to that point, around then this proposed work or scenario might give more result related to that subject. Right, when the user gets various results then a couple of words may be there that was not visited by the user with the objective that site pages are concealed by the user which may have some extraordinary data related to that subject. With the help of this paper, the customer will get all the related stage related to customer question and on the reason, that customer may pick anyone according to his choice, therefore, using natural language processing we will achieve the desired results.

## REFERENCES

[1]     Meadow, C. T., Boyce, B. R., Kraft, D. H. (2000), Text Information Retrieval Systems (second edition). Academic Press.

[2]     Ko, Y., Park, J., Seo, J. (2004), 'Improving Text Categorization Using the Importance of Sentences', Information Processing & Management, vol. 40, pp. 65-79.

[3]     Theobald, M., Siddharth, J., Paepcke, A.: SpotSigs. (2008), 'Robust and Efficient Near Duplicate Detection in Large Web Collections', Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval, ACM Press, Singapore, pp.563-570.

[4]     Wang, D., Li, T., Zhu, S. (2008), 'Multi-document summarization via sentence-level semantic analysis and symmetric matrix factorization', Proceeding of the 31st Annual International ACM SIGIR Conference, ACM Press, Singapore, pp. 307-314.

[5]     Maguitman, A., Menczer, F., Roinestad, H., Vespignani, A. (2005) 'Algorithmic Detection of Semantic Similarity'. Proceeding of the 14th International World Wide Web Conference, ACM Press, Chiba, Japan, pp.107-116.

[6]     Salton, G., Wong, A., Yang, C. S. (1975), 'A vector space model for automatic indexing', Communications of the ACM, vol. 18, pp. 613-620.

[7]     Deerwester, S., Dumais, S., Furnas, T. (1990), 'Indexing by latent semantic analysis', Journal of American Society of Information Science, Vol. 41, 391-407.

[8]     Zamir, O., Etzioni, O., Madani, O., Karp, R. M. (1997), 'Fast and intuitive clustering of web documents', Proceeding of the 3rd ACM SIGKDD Conference on Knowledge Discovery and Data Mining, ACM Press, Newport Beach, California, USA, pp. 287-290.

[9]     Zamir, O., Etzioni, O. (1998), 'Web text clustering: a feasibility demonstration', Proceeding of the 28th Annual ACM SIGIR Conference on Research and Development in Information Retrieval. ACM Press, Melbourne, Australia, pp. 46-54.

[10]   Li, Y. J., Soon, M. C., John, D. H. (2008), 'Tex text clustering based on frequent word meaning sequences', Data & Knowledge Engineering, Vol. 64, pp. 381-404.

[11]   Shehata, S., Karray, F., Kamel, M. (2007), 'A Conceptbased Model for Enhancing Text Categorization', Proceedings of the 13rd ACM SIGKDD Conference on Knowledge Discovery and Data Mining, ACM Press, San Jose, California, USA, pp.629-637,

[12]   Chim, H., Deng, X. (2007), 'A new suffix tree similarity measure for document clustering' Proceeding of the 16th International Conference on World Wide Web (2007). ACM Press, Banff, Alberta, Canada, pp.121-130.

[13]   Edith, H., Rene, A.G., Carrasco-Ochoa, J.A., MartinezTrinidad, J.F. (2006), 'Document clustering based on maximal frequent sequences', Proceedings of the FinTAL2006, LNAI, vol. 4139, pp. 257-267.

[14]   Beil, F., Ester, M., Xu, X.W. (2002), 'Frequent term-based text clustering', Proceedings of the 8th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2002), pp. 436-442.

[15]   Reuters-21578 (1997), text categorization test collection, Available at: http://www.daviddlewis.com/resources/testcollections/reuters21578/, Assessed on 17 December 2010.

[16]   BBC Dataset, (2010), Machine Learning group, Available at: http://mlg.ucd.ie, Assessed on 17 December 2010.

[17]   LingPipe, (2010), Alias-i, Inc, Available at: http://www.alias-i.com, Assessed on 17 December 2010. JOURNAL OF COMPUTERS, VOL. 6, NO. 10, OCTOBER 2011 2185 © 2011 ACADEMY PUBLISHER [18] Karypis, G., (2010), CLUTO–A Clustering Toolkit, Department of Computer Science, University of Minnesota, Available at: http://www.cs.umn.edu/~karypis/cluoto/, Assessed on 17 December 2010.