



An Ideology for Text Scooping Method

Govindraj Chittapur¹, Chetan Konnur²Assistant Professor¹, Student²

Department of MCA

Basaveshwar Engineering College Bagalkot, India

Abstract:

Identifying the similarity between articles is the complex task, here some of the systems like Stanford CS and Crossref are used to compare similarity status between the articles and the user has to pay for it. The system has to perform very quick similarity comparison between the articles using minimum amount of time and the user can able to submit many articles with different names at a time and user can get the similarity status of submitted articles. And the system should generate unique group id for every cluster of the articles. This system should be used to identify the similarity of the articles which are submitted by the users in two ways; first one is to submit the article by pasting the description directly with different article name and second one is by pasting the URL of the article. After the successful submission of articles those will be cleaned by removing the stop words then system should provide unique token id for the each word exist in the cleaned article and then tokenized clean article will be used in the frequency computation, log likelihood computation, RV coefficient computation and clustering process and these each modules will be executed one by one sequentially and finally the processed articles will be grouped with the similarity status and unique group id.

I. INTRODUCTION

System Architecture for Proposed System

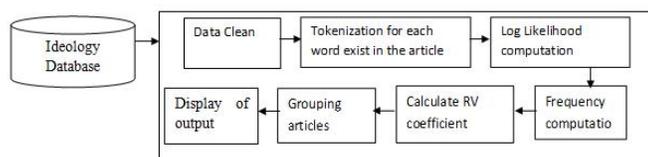


Figure.1. System Architecture for “An Ideology for Text Scooping Method”

In this proposed system the articles will be submitted in two ways as shown in the above architecture diagram, in the initial stage the articles will be submitted in the format of offline submission means here the user will submit the article by pasting the article description in the given text area and the user can also submit the article by pasting the URL of the article, in this article description will be extracted by the URL by converting the URL in to DOM tree. As shown in the above architecture diagram the articles will be stored in the format of online submission or offline submission.

Here in the above diagram the stop words database will be used to store the preserved words which are used in the data cleaning process and these stored stop words will be removed from the article description. And in the data store section the other data like login information, Likelihood computation, clustered articles, cleaned articles, RV coefficient results and other related data will be stored in this block. In the data cleaning section the data will be cleaned using the stop words which are stored in the stop words database, and the articles which are submitted by the users will be cleaned in the process, here all the stop words are erased and the submitted articles will be cleaned and doesn't have any meaning for the description. In the tokenization section the cleaned articles will be used to tokenize the words, here the

cleaned article descriptions which contain core words will be used to tokenization, and here each and every word will be tokenized and unique token id will be given for the each and every word of the article. Here in log likelihood block will be used to calculate the similarity between two articles by using the below equations

In this we need to compute

Expected freq = freq of word * freq of others / total freq

Log Likelihood (principle) = $2 * \sum_{i=1}^N \log(g_i/eg_i)$

$G_1 * \log(g_1/eg_1) + G_2 * \log(g_2/eg_2) + \dots + G_n \log(g_n/eg_n)$

G_1 is freq of first word in the article

G_2 is freq of second word in the article G_n is freq of nth word in the article

In the frequency computation module the submitted articles will be used for calculating the frequency of the repeated words in the comparing articles. RV coefficient is the module which is used to calculate the coefficient of the articles and here the words exist in the cleaned description are tend to intersection of words means here the same words exist in the comparing articles will be calculated using

$\sum (x-x^*) (y-y^*/x^{**}y^{**})$

This is used to get the intersection or common words between two articles. If $RV > 0.02$ then compared articles are similar or else they are not similar In the final module themes the compared articles will be grouped with the unique group id and shows the status on which are same and which articles are not same among the given threshold value.

II. LITERATURE SURVEY

Identifying similarity between the articles is a important research topic in research and development project management here some of the authors mentioned with their algorithms, D. Saravana Priya and Dr. M Karthikeyan [1] has proposed Clustering analysis, Fuzzy SOM, knowledge based agent, NGRA

Algorithm, Ontology, R&D and Text mining, Kuwar Aditya et al. [2] Has proposed Self Organiz-ing Method (SOM), Pravin Shinde et al. [3] has proposed Clustering, Classification, Self Organizing Map, Text Min-ing, Optimization, Latent Semantic Indexing, Mr Bhushan Medage et al. [4] proposed Data mining ,Text mining, ab-stract selection, Project Evaluation, Guide Allocation, John Butler et al. [5] proposed Simulation; Ranking and Selec-tion; Multiple Attribute Utility Theory.

III. DESIGN ISSUES FOR “AN IDEOLOGY FOR TEXT SCOOPING METHOD”

Here the articles should be submitted in two ways and each way consist its own interface, the user can submit articles using URL of the article or by directly pasting the article description. And here the URL of the article should be of html type then only the DOM tree procedure will work on div tags and the description should be limited in size.

IV. IMPLEMENTATION OF “AN IDEOLO-GY FOR TEXT SCOOPING METHOD”

An Ideology for Text Scooping Method is a pro-posed solution where the similar types of articles to be find out from the several set of articles, in the earlier solution the manual process is used to find the duplicate and similar types of articles and the manual process is very time con-suming and the result is not accurate and it uses lot of man power with other resources, this process does not consist any of the algorithms and the articles should be analyzed to process the complete task. And now in the proposed solution the several algorithms are used to differentiate the article, here text scooping means the text mining and this is done by the text mining algorithm, this proposed solution consists several stages in it like submit articles, data cleaning, tokenization, frequency computation, log likely hood, RV coeffi-cient, clustering these are the sequential stages to perform the appropriate solution to find the similarity in the article. First thing is there will be a registration and login, perform registration as admin or user, user can submit article offline or online, the admin can view all articles perform data cleaning, tokenization, freq computation, log likelihood computation, RV coefficient algorithm and grouping. Then there will be all users, they can view themes for articles submitted by n users and view grouping of all articles, Now coming to registration in this we should give first name last name user id password email id and then step number 2 is perform basic validation like first name cannot be empty then the last name cannot be empty then the user id should not be empty then the password cannot be empty email cannot be empty all these basic validation should per-form after the successful validation go to the next step oth-erwise show validation error for the user, if validations are successful retrieve list of user id’s from the application like u1,u2,u3 etc till Un we check whether the given user id is in the list if yes validation error otherwise user is created and registered successfully, this is the registration process. Then after that next basic module is login for this username and password for admin and user, basic valida-tion’s like username cannot be empty and password cannot be empty after this if basic validations are successful then retrieve list of existing user names and check whether given username is in usernames list after that if it is there then pro-ceed otherwise validation error, then in the next step retrieve the actual

password for the user password actual then com-pare actual password whether it is equal to given password from the user if both are same means continue otherwise invalid user/password error, otherwise what we do is if password validation is successful then obtain the login type if login type is double equals to 1 then provide admin func-tionality otherwise provide user functionality. In the initial stage of article submission the articles are submitted in two ways 1. Directly copy and paste the content of article which is exist in the local machine 2. Past-ing the URL link, in these two ways we can submit the ar-ticles by giving the appropriate article name. After the successful submission of the articles the articles will be cleaned on the data cleaning method by using pre reserved stop words, and in this proposed solution con-sist an option for inserting and removing stop words, accord-ing to these saved stop words the data cleaning process will be done and in this process the unwanted data which is simi-lar to the stop words only are removed from the articles and after this data cleaning the article descriptions will be mea-ningless and the cleaned article will be stored in the format of cleaned, article name and clean article.

V. RESULT AND DISCUSSION:

The below screenshots are useful to understand the interface and the procedure of the process



Figure.2. Article submission using URL

The system will fetch the article description from div tags using DOM tree technique

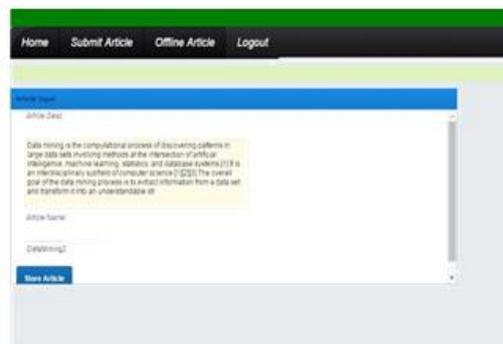


Figure.3. Article submission by pasting description of the article And here the description size should be limited

| Clean ID | Article Name | Clean Article |
|----------|--------------|---|
| 25 | Article2 | contextual learning based constructivist |
| 26 | Article1 | contextual learning based constructivist |
| 27 | DataMining | navigation search analytics information |
| 28 | DataMining2 | data mining computational process discov |
| 29 | DataMinin... | navigation search analytics information |
| 30 | DTAA1 | navigation search analytics information |
| 31 | Electroma... | principle article send electromagnetic wave |
| 32 | HELLO | navigation search analytics information |
| 33 | KohArticle1 | definition word symbol explaining meaning |
| 34 | Test Article | article responsible lecture motivation user |

Figure.4. List of all submitted articles

Each of the cleaned articles consist its own clean ID

| Article Name | Token ID | Token Name |
|--------------|----------|----------------|
| Art:le2 | 1126 | contextual |
| Art:le2 | 1127 | learning |
| Art:le2 | 1128 | based |
| Art:le2 | 1129 | constructivist |
| Art:le2 | 1130 | theory |
| Art:le2 | 1131 | teaching |
| Art:le2 | 1132 | learning |
| Art:le2 | 1133 | learning |
| Art:le2 | 1134 | takes |
| Art:le2 | 1135 | take |

Figure.5. Tokenized each word of the article and the ID will be unique

| Frequency ID | Article Name | Token Name | Frequency |
|--------------|-------------------|-----------------|-----------|
| 1 | Electromagnetic s | principle | 1 |
| 2 | Electromagnetic s | article | 1 |
| 3 | Electromagnetic s | send | 1 |
| 4 | Electromagnetic s | electromagnetic | 1 |
| 5 | Electromagnetic s | waves | 1 |
| 6 | Electromagnetic s | destination | 1 |
| 7 | Electromagnetic s | demodulate | 1 |
| 8 | Electromagnetic s | apack | 1 |
| 9 | Electromagnetic s | lpack | 1 |
| 10 | Electromagnetic s | spack | 1 |
| 11 | Test Article | article | 1 |
| 12 | Test Article | responsible | 1 |
| 13 | Test Article | testing | 1 |
| 14 | Test Article | population | 1 |
| 15 | Test Article | user | 1 |
| 16 | Test Article | it | 1 |

Figure.6. Frequency calculation over the articles according to the number of repeated token names and generates unique fre-quency ID

| Log Likelihood ID | Article Name | Token Name | Frequency | Expected Frequency | Log Likelihood |
|-------------------|-------------------|-----------------|-----------|--------------------|------------------|
| 2053 | Electromagnetic s | principle | 1 | 0.9 | 11.3789358910453 |
| 2054 | Electromagnetic s | article | 1 | 0.9 | 11.3789358910453 |
| 2055 | Electromagnetic s | send | 1 | 0.9 | 11.3789358910453 |
| 2056 | Electromagnetic s | electromagnetic | 1 | 0.9 | 11.3789358910453 |
| 2057 | Electromagnetic s | waves | 1 | 0.9 | 11.3789358910453 |
| 2058 | Electromagnetic s | destination | 1 | 0.9 | 11.3789358910453 |
| 2059 | Electromagnetic s | demodulate | 1 | 0.9 | 11.3789358910453 |
| 2060 | Electromagnetic s | apack | 1 | 0.9 | 11.3789358910453 |
| 2061 | Electromagnetic s | lpack | 1 | 0.9 | 11.3789358910453 |
| 2062 | Electromagnetic s | spack | 1 | 0.9 | 11.3789358910453 |
| 2063 | Electromagnetic s | principle | 1 | 0.9 | 11.3789358910453 |
| 2064 | Electromagnetic s | article | 1 | 0.9 | 11.3789358910453 |
| 2065 | Electromagnetic s | send | 1 | 0.9 | 11.3789358910453 |
| 2066 | Electromagnetic s | electromagnetic | 1 | 0.9 | 11.3789358910453 |
| 2067 | Electromagnetic s | waves | 1 | 0.9 | 11.3789358910453 |

Home Articles Stopword Analysis Data Cleaning Likelihood Computation Duplicate Articles Logout

Article Similarity Input

Left Article Name: DataMining

Right Article Name: DataMining2

Type: 0.02

Compare Similarity

Figure.7. Comparing articles by selecting both left and right different articles with threshold value

Home Articles Stopword Analysis Data Cleaning Likelihood Computation Duplicate Articles Logout

Comparison of Articles is Successful

| Right Article Sum | Left Article Sum | Left Article Mean | Right Article Mean | Left Article SD | Right Article SD | RV Coefficient | Similarity Status |
|-------------------|------------------|-------------------|--------------------|------------------|------------------|-----------------|-------------------|
| 756 | 344 | 1.72 | 2.34782608895622 | 11.6756156154611 | 28.5489663256109 | -0.000467699952 | false |

Intersection Words:

process
funding
agency
distributed
cfo

Message Information List

Message

Figure.8. Similarity status of the compared articles

| Iteration Number | Time Algo1 | Time Algo2 |
|------------------|------------|------------|
| 1 | 1 | 0.01 |
| 2 | 16 | 14 |
| 3 | 2529 | 2504 |
| 4 | 1920 | 1919 |
| 5 | 25868 | 25862 |
| 6 | 2483 | 2482 |
| 7 | 0.001 | 0.01 |
| 8 | 2 | 0.01 |
| 9 | 1 | 0.01 |
| 10 | 25 | 5 |
| 11 | 5 | 0.01 |
| 12 | 5559 | 5545 |

Figure.9. Time consumption comparison with proposed and previous system

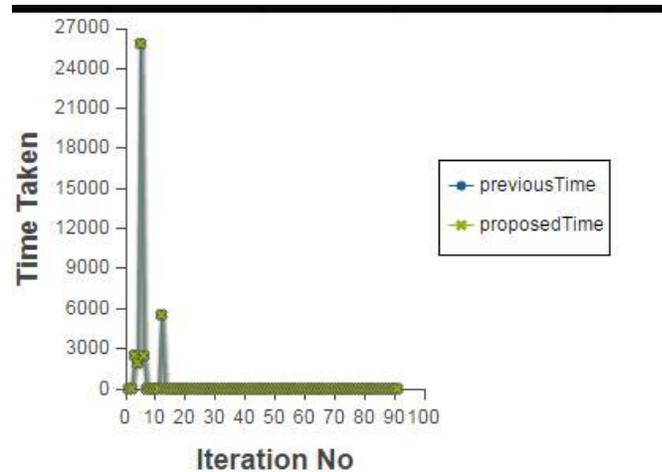


Figure.10. Graphical view of the previous and proposed system time

VI. ACKNOWLEDGMENTS

I would like to acknowledge my sincere thanks to Basa-veshwara Engineering College Department of MCA for giving support, resources for doing this research paper.

VII. REFERENCES

- [1].Q. Tian, J. Ma, and O. Liu, "A hybrid knowledge and model system for R&D project selection," Expert Syst. Appl., vol. 23, no. 3, pp. 265–271, Oct. 2002.
- [2].K. Chen and N. Gorla, "Information system project selection using fuzzy logic," IEEE Trans. Syst., Man, Cybern. A, Syst., Humans, vol. 28, no. 6, pp. 849–855, Nov. 1998.
- [3].A. D. Henriksen and A. J. Traynor, "A practical R&D project-selection scoring tool," IEEE Trans. Eng. Manag., vol. 46, no. 2, pp. 158–170, May 1999.
- [4].F. Ghasemzadeh and N. P. Archer, "Project portfolio selection through decision support," Decis. Support Syst., vol. 29, no. 1, pp. 73–88, Jul. 2000.
- [5].J. Butler, D. J. Morrice, and P. W. Mullarkey, "A multiple attribute utility theory approach to ranking and selection," Manage. Sci., vol. 47, no. 6, pp. 800–816, Jun. 2001.

VII. BIOGRAPHIES

Prof. GOVINDRAJ CHITTAPUR received the B.Sc. degree in Computer Science from the University of KUD, Dharwad, Karnataka, in 2002, the M.C.A degree from VTU, Belgaum, Karnataka, in 2005 and awarded M.sc Technology by Research from Mysore University in 2011. Has published 17 international journal and national journal has also served as Reviewer and Editorial Board Member of various international journal and conferences. His major research area includes Image and Video forensics, Machine learning and Data Mining.

Mr. CHETAN M. KONNUR received the B.Sc. degree in computer science from the Rani Channamma University, Belgaum, Karnataka, in 2014, and pursuing MCA degree under VTU, Belgaum, Karnataka.