



ISSN(Online): 2320-9801
ISSN (Print): 2320-9798

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

Content Based Document Retrieval Using Relevance Feature Discovery

Priyanka Awaje¹, R . S. Jamgekar², S. S. Joshi³,

M. E Student, Department of Computer Science and Engineering, N B Navale Sinhgad college of Engineering, Kegaon, Solapur, India

Assistant Professor, Department of Computer Science and Engineering, N B Navale Sinhgad college of Engineering, Kegaon, Solapur, India

Assistant Professor, Department of Computer Science and Engineering, N B Navale Sinhgad college of Engineering, Kegaon, Solapur, India

ABSTRACT: In data mining and knowledge engineering schema, the main focus is on the identification of text documents' feature extraction, which illustrates the user preferences' in terms of huge data patterns'. Lots of approaches are proposed earlier for data mining and text classification schemes but all are compacted with only term based methodologies. As well as all this kind of schemes are highly affected from the problem of polysemy and synonymy. Throughout the years, there has been frequently held the speculation that example based strategies ought to perform superior to anything term based ones in portraying client inclinations; yet, how to successfully utilize extensive scale designs remains a difficult issue in content mining. To make an achievement in this testing issue, this paper displays an imaginative model for pertinence highlight revelation. It finds both positive and negative designs in content reports as larger amount highlights and conveys them over low level components [terms]. It additionally orders terms into classes and upgrades term weights taking into account their specificity and their dispersions in examples. Considerable tests utilizing this model on RCV1, TREC points and Reuters-21578 demonstrate that the proposed display altogether beats both the best in class term-based strategies and the example based techniques.

KEYWORDS: Data Mining, Data Classification, Feature Extraction, Text Mining' and Classification.

I. INTRODUCTION

For identifying the relevant and irrelevant features in the text documents, a special methodology is required for mining the text sequences, which is called "Relevance Feature Discovery [RFD]". This is an especially difficult errand in present day data investigation. There are two testing issues in utilizing pattern mining procedures for discovering significance highlights in both pertinent and unimportant archives. They are:

(a) *The first is the low bolster issue, given a theme, long examples are typically more particular for the point, and however they generally show up in archives with low backing or recurrence. In the event that the base backing is diminished, a great deal of boisterous examples can be found.*

(b) *The second issue is the error issue, which implies the measures [for instance, "Backing" and "Certainty"] utilized as a part of example mining end up being not reasonable in utilizing designs for taking care of issues. For instance, a very continuous example [ordinarily a short example] might be a general example since it can be habitually utilized as a part of both important and insignificant records. Subsequently, the troublesome issue is the manner by which to utilize found examples to precisely weight helpful components.*

For a long time, we have watched that numerous terms with bigger weights are more broad since they are prone to be much of the time utilized as a part of both pertinent and superfluous reports [1]. For instance, word "LIB" might be more much of the time utilized than word "JDK"; yet "JDK" is more particular than "LIB" for portraying "Java Programming Languages", and "LIB" is more broad than "JDK" on the grounds that "LIB" is additionally as often as possible utilized as a part of other programming dialects like C or C++.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

Accordingly, in pertinence highlight revelation [RFD] we consider both terms' dispersions and specificities. Given a theme, a term's specificity portrays the degree to which the term concentrates on the subject that clients need. Be that as it may, it is extremely hard to quantify the specificity of terms on the grounds that a term's specificity relies on upon clients' points of view of their data needs.

Highlight choice is a strategy that chooses a subset of elements from information for displaying frameworks there are assortment of highlight choice techniques (e.g., Filter, Wrapper, Embedded and Hybrid methodologies, and unsupervised or semi-administered strategies) have been proposed in different fields. Highlight determination is likewise one of vital strides for content characterization and data sifting [2], which is the undertaking of doling out records to predefined classes.

II. RELATED STUDY

Highlight choice is a procedure that chooses a subset of components from information for demonstrating frameworks (see http://en.wikipedia.org/wiki/Feature_selection). Throughout the years, an assortment of highlight choice strategies (e.g., Filter, Wrapper, Embedded and Hybrid methodologies, and unsupervised or semi-administered techniques) have been proposed in different fields [6], [9], [17], [18], [19]. Highlight choice is additionally one of imperative strides for content grouping and data sifting [1], [5], [17] which is the undertaking of allocating records to predefined classes.

To date, numerous classifiers, for example, Naive Bayes, Rocchio, kNN, SVM and Lasso relapse [16], [1], [2], [6] have been created, furthermore numerous trust that SVM is additionally a promising classifier [13]. The characterization issues incorporate the single class and multi-class issue. The most widely recognized arrangement [17] to the multi-class issue is to deteriorate it into some autonomy double classifiers, where a parallel one is doled out to one of two predefined classes (e.g., applicable classification or unessential class). Most customary content element determination strategies utilized the sack of words to choose an arrangement of components for the multi-class issue [13].

There are some component choice criteria for content order, including archive recurrence (DF), the worldwide IDF, data increase, common data (MI), Chi-Square (χ^2) and term quality [1], [19], [20]. In this paper we concentrate on pertinent element determination in content records. Pertinence is a major exploration issue [11], [12], [15] for Web look, which talks about a reports importance to a client or a question. Nonetheless, the customary element choice strategies are not powerful to select content components for explaining pertinence issue since importance is a solitary class issue [13]. The proficient method for highlight choice for importance depends on an element weighting capacity. An element weighting capacity shows the level of data spoke to by the component events in a record and mirrors the significance of the element. The prominent term-based positioning models incorporate $tf*idf$ based procedures, Rocchio calculation, Probabilistic models and Okapi BM25 [4], [14], [17].

Pattern Taxonomy Mining [PTM] models have been proposed in earlier researches, in which, mining shut successive examples in content passages and conveying them over a term space to weight helpful elements. Concept Based Model [CBM] has additionally been proposed to find ideas by utilizing Natural Language Processing (NLP) strategies. It proposed verb-contention structures to discover ideas in sentences.

Effective Pattern Discovery for Text Mining [6], the D-Pattern Mining Algorithm is used. It describes the training process of finding the set of d-patterns. The main focus of this paper is the deploying process, which consists of the d-pattern discovery and term support evaluation. It discovered all patterns in a positive document are composed into a d-pattern. An innovative and effective pattern discovery technique which includes the processes of pattern deploying and pattern evolving, to improve the effectiveness of using and updating discovered patterns for finding relevant and interesting information.

These example (or ideas) based methodologies have demonstrated an imperative change in the adequacy [7]. A two-phase model was likewise proposed in [8], [9], which demonstrated that the coordination of the harsh investigation (a term-based model) and example scientific classification mining is the most ideal approach to plan a two-phase model for data separating frameworks. The current model watched that numerous terms with bigger weights are more broad since they are prone to be as often as possible utilized as a part of both important and unessential archives [1]. For instance, word "LIB" might be more as often as possible utilized than word "JDK"; yet "JDK" is more particular than "LIB" for portraying "Java Programming Languages"; and "LIB" is more broad than "JDK" on the grounds that "LIB" is additionally much of the time utilized as a part of other programming dialects like C or C++.

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

Subsequently, we suggest the thought of both terms' disseminations and specificities for importance feature discovery. The current techniques for discovering relevance features can be gathered into three methodologies [1].

(a) The main methodology tries to lessen weights of terms that show up in both pertinent records and superfluous reports (e.g., Rocchio-based models [10]). This heuristic is clear in the event that we accept that terms are secluded molecules.

(b) The second one depends on how regularly includes show up or don't show up in pertinent and insignificant archives (e.g., probabilistic based models [11]).

(c) The third one depends on discovering highlights through positive examples [1], [2], [3]. The proposed display further builds up the third approach by gathering highlights into three classifications: "positive particular components", "general elements", and "negative particular elements".

III. PROPOSED SYSTEM

In proposed framework concentrates on important element choice in content reports the effective method for highlight choice for significance depends on a component weighting capacity. A component weighting capacity shows the level of data spoke to by the element events in a report and mirrors the importance of the element. The famous term-based positioning models incorporate $tf \cdot idf$ based methods is used.

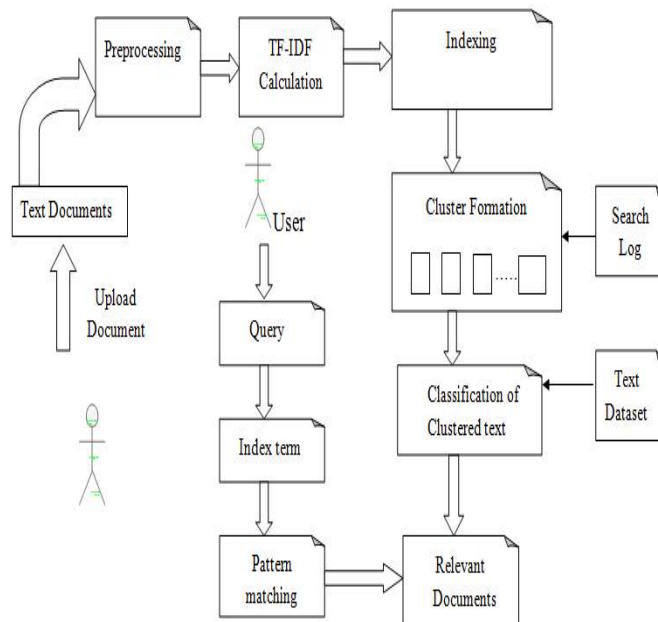


Fig.1. System Architectural Design

The framework can be comprehensively grouped into 4 modules. They are as per the following:

(A) PREPROCESSING:

In this stage the engineer stacks the dataset record/s. These records should be preprocessed which incorporates stop words evacuation here the fundamental conjunctions and joining words are expelled. Next is stemming where appropriate morphological root is recognized. Case in point words like streams, spilled and gushing have the morphological root as stream.

(B) WEIGHT COMPUTATION:

After the archives are preprocessed the following stride comprises of weight figuring. Before that the successive terms are distinguished by applying the weight



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

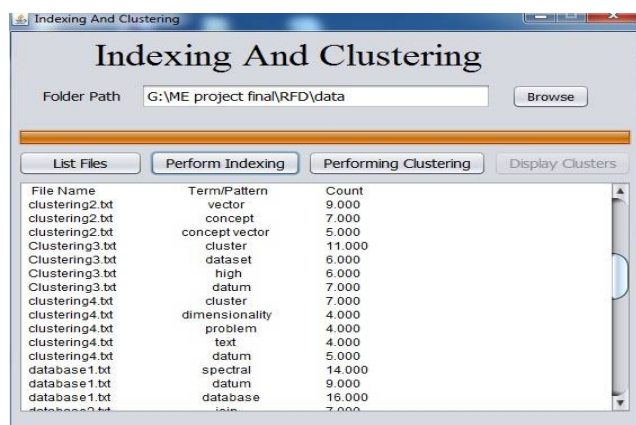
Count Algorithm. For each set the positive and negative scope sets are gotten and the specificity recipe is connected. Specificity(s) = |set of significant documents| / |set of immaterial documents| / N
Term's specificity it depicts the degree on which it concentrates on the subject.

(C) TERM CHARACTERIZATION:

Once the regular terms are recognized and the weights are figured these terms are gathered into groups. For this the grouping calculation is connected.

(D) RANKING:

In this last stage the characterized set of records are shown which are in agreement to the figured weights.

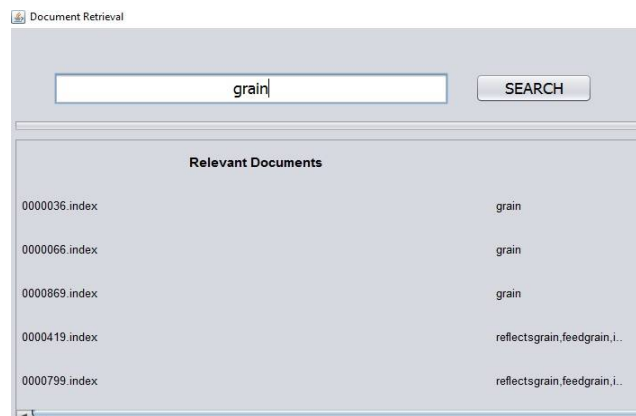
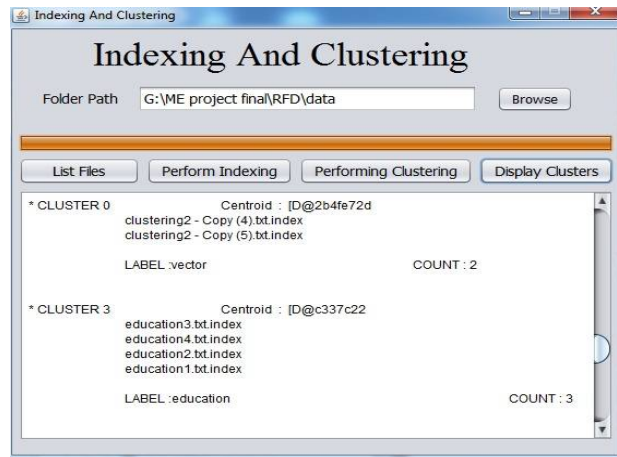




International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016



IV. RESULT TABLE

The following table illustrates the result on number of documents from dataset Reuters 21578.

Table 1: Statistical Information for RFD with different values of K

K=No. of Documents in Training Set	No of Positive Document	No of Negative Document
k=50	40	10
k=25	24	1
k=20	15	5
k=10	8	2

International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

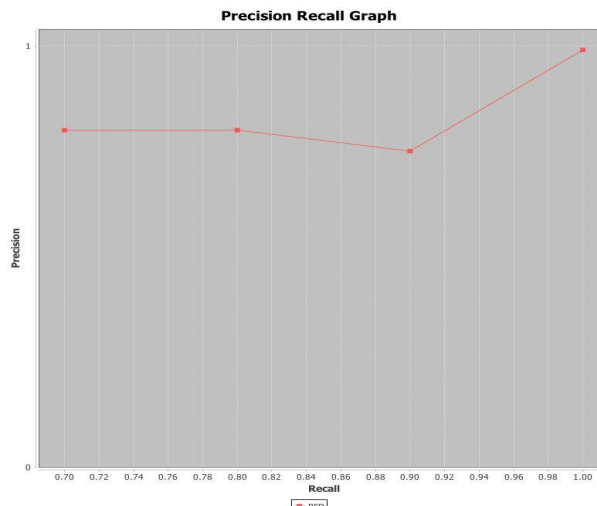


Fig 2: Precision Recall for RFD

Table 2: Statistical Information for word Similarities

Document Id	Terms /Patterns	Term Frequency
1	Price	5
1	Level	5
1	Reserve	3
1	Wheat	3
2	Prev	8
2	Nil	7
2	Feb	9
2	Total	13
2	prev feb	8
11	Bonus	3
11	Yemen	2
11	Wheat	4
11	north Yemen	2
11	wheat flour	3
36	Corn	3

V. CONCLUSION

The RFD model, introduces a strategy to discover and order low level components in light of both their appearances in the larger amount designs and their specificity. It likewise acquaints a strategy with select insignificant archives for weighting features. In this system, we build up the RFD model which demonstrates that the proposed specificity capacity is sensible and the term grouping can be successfully approximated by a component bunching strategy. The past RFD model uses two exact parameters to set the limit between the classifications. It accomplishes the expected execution, yet it requires the physically testing of an extensive number of various estimations of parameters.



International Journal of Innovative Research in Computer and Communication Engineering

(An ISO 3297: 2007 Certified Organization)

Vol. 4, Issue 11, November 2016

The proposed model uses a component grouping method to automatically group terms into the three classifications. Analyzed with the first model, the new model is a great deal more effective and achieved the palatable execution also.

The proposed framework demonstrates that the utilization of unimportance input is noteworthy for improving the execution of pertinence feature disclosure models. It provides a promising technique for creating effective text digging models for importance feature relevance based on both positive and negative criticism.

REFERENCES

- [1] M. Aghdam, N. Ghasem-Aghaee, and M. Basiri, "Text feature selection using ant colony optimization," in *Expert Syst. Appl.*, vol. 36, pp. 6843–6853, 2009.
- [2] A. Algarni and Y. Li, "Mining specific features for acquiring user information needs," in *Proc. Pacific Asia Knowl. Discovery Data Mining*, 2013, pp. 532–543.
- [3] A. Algarni, Y. Li, and Y. Xu, "Selected new training documents to update user profile," in *Proc. Int. Conf. Inf. Knowl. Manage.*, 2010, pp. 799–808.
- [4] N. Azam and J. Yao, "Comparison of term frequency and document frequency based feature selection metrics in text categorization," *Expert Syst. Appl.*, vol. 39, no. 5, pp. 4760–4768, 2012.
- [5] R. Bekkerman and M. Gavish, "High-precision phrase-based document classification on a modern scale," in *Proc. 11th ACM SIGKDD Knowl. Discovery Data Mining*, 2011, pp. 231–239.
- [6] N. Zhong, Y. Li, and S.-T. Wu, "Effective pattern discovery for text mining," in *IEEE Trans. Knowl. Data Eng.*, vol. 24, no. 1, pp. 30–44, Jan. 2012.
- [7] C. Buckley, G. Salton, and J. Allan, "The effect of adding relevance information in a relevance feedback environment," in *Proc. Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 1994, pp. 292–300.
- [8] G. Cao, J.-Y. Nie, J. Gao, and S. Robertson, "Selecting good expansion terms for pseudo-relevance feedback," in *Proc. Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2008, pp. 243–250.
- [9] G. Chandrashekar and F. Sahin, "A survey on feature selection methods," in *Comput. Electr. Eng.*, vol. 40, pp. 16–28, 2014.
- [10] B. Croft, D. Metzler, and T. Strohman, *Search Engines: Information Retrieval in Practice*. Reading, MA, USA: Addison-Wesley, 2009.
- [11] F. Debole and F. Sebastiani, "An analysis of the relative hardness of Reuters-21578 subsets," *J. Amer. Soc. Inf. Sci. Technol.*, vol. 56, no. 6, pp. 584–596, 2005.
- [12] J. Eisenstein, A. Ahmed, and E. P. Xing, "Sparse additive generative models of text," in *Proc. Annu. Int. Conf. Mach. Learn.*, 2011, pp. 274–281.
- [13] G. Forman, "An extensive empirical study of feature selection metrics for text classification," in *J. Mach. Learn. Res.*, vol. 3, pp. 1289–1305, 2003.
- [14] Y. Gao, Y. Xu, and Y. Li, "Topical pattern based document modelling and relevance ranking," in *Proc. 15th Int. Conf. Web Inf. Syst. Eng.*, 2014, pp. 186–201.
- [15] X. Geng, T.-Y. Liu, T. Qin, A. Arnold, H. Li, and H.-Y. Shum, "Query dependent ranking using k-nearest neighbor," in *Proc. Annu. Int. ACM SIGIR Conf. Res. Develop. Inf. Retrieval*, 2008, pp. 115–122.
- [16] A. Genkin, D. D. Lewis, and D. Madigan, "Large-scale Bayesian logistic regression for text categorization," *Technometrics*, vol. 49, no. 3, pp. 291–304, 2007.
- [17] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," in *J. Mach. Learn. Res.*, vol. 3, no. 1, pp. 1157–1182, 2003.
- [18] J. Han, J. Pei, and Y. Yin, "Mining frequent patterns without candidate generation," in *Proc. ACM SIGMOD Int. Conf. Manage. Data*, 2000, pp. 1–12.
- [19] Y.-F. Huang and S.-Y. Lin, "Mining sequential patterns using graph search techniques," in *Proc. Annu. Int. Conf. Comput. Softw. Appl.*, 2003, pp. 4–9.
- [20] G. Ifrim, G. Bakir, and G. Weikum, "Fast logistic regression for text categorization with variable-length n-grams," in *Proc. ACM SIGKDD Knowl. Discovery Data Mining*, 2008, pp. 354–362.