



A Genetic Programming Framework for Topic Discovery from Online Digital Library

Mrs. Vijashree Kasturba

Indus University, Ahmedabad

ABSTRACT

Various topic extraction techniques for digital libraries are projected over the past decade. usually the subject extraction system needs an oversized range of options and complex lexical analysis. whereas these options and analysis are effective to represent the applied math characteristics of the document, they did not capture the high level linguistics. during this paper, we tend to gift a replacement approach for topic extraction. Our approach combines user's click stream information with ancient lexical analysis. From our purpose of read, the user's click stream directly reflects human understanding of the high-level linguistics within the document. moreover, a simple, however effective, piecewise linear model for topic evolution is projected. we tend to apply genetic formula to estimate the model and extract topics. Experiments on the set folks congress digital library documents demonstrate that our approach achieves higher accuracy for the subject extraction than ancient strategies.

1. INTRODUCTION

Today digital libraries and net system become a fashionable and open resource for researchers. Those systems cover varied analysis topics in many alternative subjects. additional and additional analysers have found it's convenient to seek out attention-grabbing research topics, seminal analysis papers and latest progress in analysis frontiers exploitation those systems. The increasing range of individuals are contributive to those systems, either deliberately or incidentally. This has created a large set of knowledge that would doubtless be wont to improve the services of digital library and net system normally. On the opposite hand, the development of the wide adoption of on-line digital library and net system as a primary means that for information sharing has diode to an excellent booming of analysis during this field over the past fifty years. plenty of innovative tools and clever algorithms are developed to facilitate the effective and economical usage of those systems, for instances document compartmentalization [3], data retrieval [4], document classification [10] etc. On the opposite hand, with the prevailing business and subscribing services provided by the net, digital libraries and net data systems cover varied topics in many alternative analysis subjects [2]. Indeed, they need become such a fashionable and open resource out there to analysis communities. Through them, researchers will simply access an excellent quantity and a large kind of data. With those data, it's convenient to seek out attention-grabbing analysis topics, seminal analysis papers and latest progress in analysis frontiers. The development of the wide adoption of on-line digital library and net system as a primary means that for information sharing diode to an excellent booming of analysis during this field over the past fifty years. plenty of innovative tools and clever algorithms are developed to facilitate the effective and economical usage of those systems, for instances document compartmentalization [3], data retrieval [4], document classification [10] etc. Recently, Web 2.0, the new generation of net, has become more and more in style. The necessary characteristic for net two.0 is that the interactive usage of the net content and also the wealth of helpful meta-data. not like ancient net, during which most users passively receive information; in net two.0, the users are taking part in a far additional interactive role. once browsing the net, the users perform many alternative actions, like bookmarking web site, tagging documents with key words, creating recommendations to friends, and sharing documents among colleagues etc. The new thought and technology of "Web two.0" has introduced alternative innovative usages and applications for the net. Examples embody linguistics net, personal diary, social networks etc. [9]. Digital library, as a special quite net system, are influenced by the new net two.0 technology. Researchers have recently began to make the most of net two.0 and incorporate the meta information to boost the services of digital library systems. for instance, the authors [10] use the history tagging data to form recommendation and suggestion on future articles for the users.

However, the matter of "topic discovery" is incredibly difficult due to (i) the massive gap between low level feature and high level linguistics and (ii) the unceasingly ever-changing of the analysis topics over time. ancient approaches concentrate on analyzing and police work topics exploitation document content options and document structure analysis

[1]. whereas these options and analysis are capable to represent the applied math characteristics of the documents, they didn't capture the high level linguistics content of the document. and that they typically need an oversized quantity of coaching information so as to be effective. This any limits their relevancy for pursuit new analysis topics, wherever the coaching information is usually distributed. as luck would have it, the interactive usage of net two.0 generates legion user information, that partly helps to resolve the scantness of coaching information downside. one among the easy however effective information in on-line systems are the clicking Stream. Those information are recorded from the user's browsing session and keep in kind of web-logs. they will be simply accessed through meta-data attributes related to the document. they're extremely related to finish user's perception and closely follow the analysis trend development. These characteristics build them ideal candidate to boost the performance of "topic discovery" for on-line digital library. during this paper, we tend to gift a replacement approach to include the clicking Stream information for "topic discovery".

A. Background and Motivation

A photograph of the clicking stream information is shown in Table I. the clicking stream contains an inventory of five-item tuples. every tuple records five completely different items of data regarding the user's question: (1) the time stamp for query, (2) the anonymous user's positive identification, (3) the question phrase issued by the user, (4) the document link that the user clicked and (5) the outline of the clicked document (usually the abstract of the document). Before going into details on a way to method the clicking stream information, we tend to justify the motivation for exploitation the mass click stream information as a promising supply for the task of topic discovery. Firstly, user's clicks represent direct measuring for topic relevancy from human's purpose of read.

2. FEATURE EXTRACTION AND ILLUSTRATION

In this section, we tend to gift feature process for on-line click stream information.

A. Preprocessing the raw click stream accumulated from the on-line digital libraries. within the table, every row records one spherical of question and retrieval.

Each question record corresponds to one spherical of question and retrieval. However, rather than considering individual question records, we tend to preprocess the information by aggregation. In preprocessing, individual question records are sorted into larger question teams. every question cluster consists of multiple question records issued by identical user among a brief quantity. This preprocessing relies on the observation that a user is usually inquisitive about one topic among a brief fundamental quantity. Therefore, multiple queries submitted consecutively by identical user among a brief time is probably going to get on identical topics. The question cluster provides a group of extremely connected query-document pairs and exhibits higher linguistics coherence.

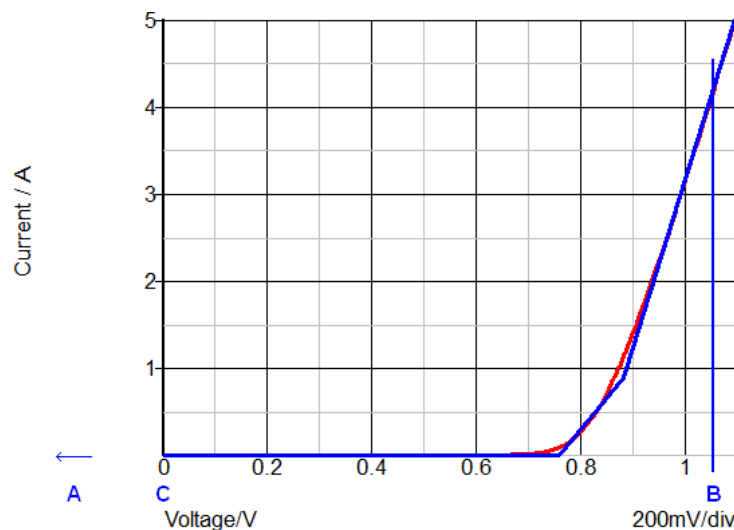


Figure 1. Convex Solver for the NMF problem of matrix size 30×100 .

B. Feature Extraction

In our approach, we discover associate degree embedding feature vector for every question cluster. This method involves 2 steps: (1) cipher the pair-wise similarity matrix between any 2 question teams and (2) realize embedding feature vector for every question cluster exploitation Non-negative Matrix factorisation.

3. TOPIC EVOLUTION MODELING AND EXTRACTION

In previous section, we've got talked regarding a way to extract options from users' click stream information and use NMF to project the options into lower dimension area. once these steps, during this section, we tend to gift our model for topic evolution.

A. Piecewise linear model In our approach, the subject evolution is shapely as piecewise linear functions. each topic is shapely as a linear perform on the evolving curve. Topic changes happen at the joint of 2 distinct linear functions. Fig. two illustrates the thought of exploitation this model for on-line topic evolution. In Fig. 2, coordinate axis represents time and coordinate axis represents the options. 3 distinct topics exist during this fundamental quantity. The red/green/blue dots correspond to individual feature embedding for question cluster from 3 topics. The dotted lines in Fig. two represent the piecewise linear model for three topics. the subject changes happen at the joint between completely different linear segments. Equation (10) shows the subject evolution model.

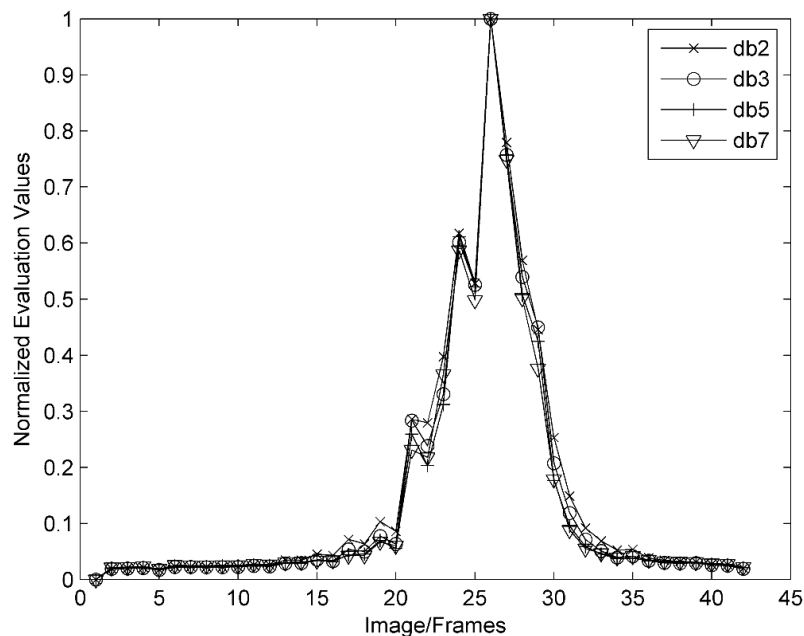


Figure 2. Illustration of piece-wise linear model for topic extraction.

Given the model outlined in equation (10), topic extraction becomes a drag of estimating distinct linear functions from the subject curve. The estimation of the piecewise linear perform in equation (10) involves 2 steps: (i) realize a partition of the linear curve and (ii) estimate a linear perform for every partition. This downside is nonlinear improvement downside and that we use a genetic formula to resolve it.

Fitness perform. For a given body, the fitness perform is outlined as equation (11), wherever the primary term is that the add of the fitting error from all the linear segments, and also the second term represents the penalty from incorrect estimation for the quantity of linear segments. Note that while not this penalty term, the GA formula tends to come up with too several little topic segments. In equation (11), the constant n_0 represents the previous information regarding the quantity of desired topics. it's set by the user. The constant α controls the pliability for the quantity of the topics. Larger α means that less versatile selection on the quantity of topics. In apply, we elect α adequate one hundred and twenty fifth of the overall range samples.

Selection operators. choice operators ar wont to choose individual chromosomes to that the crossover operators are going to be applied. In literature, many choice operators are projected, i.e. "Select Random", "Select Best" and "Tournament Selection". In Section IV, we tend to gift experimental results to match the performance of various selector operators. Crossover operators.

Crossover is a very important operators in GA formula, that choose people from the parental generation and interchange their genes to come up with new people (descendants). The aim here is to get descendants of higher quality that may be propagated to the long run generation and change the search to explore new regions of resolution area not



explored however. many sorts of crossover operators are explored within the organic process computing literature, that depends on the body illustration. For our downside, the crossover operator ar outlined for the binary string illustration.

4. EXPERIMENTAL RESULTS

We have enforced GA primarily based topic extraction formula exploitation the Python Genetic Evolution formula Framework -pyevolve [6]. The fitness perform, choose operators, crossover operators, mutation operators and replacement operators ar outlined as plug-in request functions. The benchmark information ar obtained from the U.S. Congress document library info (1989-2006) and also the question sessions were recorded once the users browsed the library. In total, there ar thirty three completely different topics within the info. The user's click stream ar sorted into 3432 question sessions and preprocessed into feature vectors as delineated in Section II. Table II lists the overall thirty three topics. Those topics make up four categories: (1) Economic, (2) Education, (3) Military and (4) Energy. The performance for various decisions of GA operators on the benchmark information ar given in Fig. 4. Fig. 4(a) and (b) show the results for 3 completely different GA selector operators -Random choice, Best choice, and Tournament choice (N=10), and 2 completely different crossover operators - uniform and one crossover, severally. In each Fig. 4(a) and (b), coordinate axis is that the range of generations and coordinate axis is that the average fitness perform over the populations. In Fig. 4(a), the Tournament selector and best selector performs higher.

5. CONCLUSIONS

User's click stream have recently known as a possible supply for topic extraction. during this paper, we tend to gift a GA primarily based approach to extract document topics for on-line digital library. Our approach combines user's click stream and document abstracts. There ar 3 main contributions for this paper. First, a unique approach to mix click-stream with document abstract to boost the subject discovery method is developed. Secondly, a piece-wise linear perform is projected to model the subject evolution. Thirdly, a Genetic formula is developed for topic extraction. The projected approach mechanically determines the quantity of topics, and teams question instances and library documents into completely different (a) (b) topic classes.

REFERENCES

- [1]. J. Allan, R. Papka, and V. Lavrenko. On-line new event detection and pursuit. In SIGIR, 1998.
- [2]. R. Baeza-Yates and B. Ribeiro-Neto. fashionable data retrieval. Addison-Wesley, 1999.
- [3]. C. Barry and L. Schamber. Users criteria for relevancy evaluation: A cross-situational comparison. informatics and Management, 34(2-3):219-236, 1998.
- [4]. N. J. Belkin. Intelligent data retrieval: Whose intelligence? In Proceedings of the Fifth International conference for IP, pages 25-31, 1996.
- [5]. S. Boyd and L. Vandenberghe. gibbose improvement. Cambridge University Press, 2004.
- [6]. <http://pyevolve.sourceforge.net>.
- [7]. D. Lee and H. S. Seung. Learning the components of objects by non-negative matrix factorisation. Nature, 401:788-791, 1999.
- [8]. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorisation. In NIPS, volume 13, page 556C562, 2001.
- [9]. T. Rattenbury, N. Good, and M. Naaman. Towards automatic extraction of event and place linguistics from flickr tags. In SIGIR, 2007.