

A Survey on Efficient Exploration of Algorithms and Pulled out From Scholarly Big Data

A.R.Sanas^{1*}, P.S.Patil²P.G. Student, Dept. of Computer Science Engineering, Savitribai Phule Pune University, Pune, India^{1*}Professor, Dept. of Computer Science Engineering, Savitribai Phule Pune University, Pune, India²

ABSTRACT: Algorithms are important and essential part of computational science research work which is regularly published in scholarly articles. To solve the scaling problem of handling more documents, we examine an intelligent system designed with the goal of automatically identifying algorithms. In this paper, we propose a method to develop an algorithm search engine. The proposed system analyzes a document to discover any algorithm that may be there in the document. If any algorithm is found in the document, the document text is further analyzed to extract additional information about the algorithm. Machine learning and rule based approaches are used to discover algorithm representation. All discovered algorithms and their associated metadata are indexed and offered for searching throughout a text query interface. Along with this, particularly important and relevant textual content can be accessed the platform and highlight portions by anyone with different levels of knowledge. In support of lectures and self-learning, the highlighted documents can be shared with others. But different levels of learners can-not use the highlighted part of text at same understanding level. The problem of predicting new highlights of partly high-lighted e-learning documents can be solved by us.

KEYWORDS: Algorithms, Pseudo codes, Scholarly big data, Sentence Extractor, Steaming, TFIDF

I. INTRODUCTION

To solve problems, Algorithms are used everywhere in Computer Science and over an exact methodology. Each and every aspect of human life has been affected by algorithms. Improvement of the current algorithms and development of new algorithms for new and unsolved problems is being made by researchers. To search for algorithms, the current state of the art search engines however, are not optimized. Documents that contain an algorithm and those that do not cannot be distinguished by them. The search results contain a variety of unwanted, irrelevant documents as a result. For example, a query to search for algorithm for KNN will return lots of documents that contain KNN in them even though they do not offer any details about the algorithmic aspects of KNNs. The relevant documents might not show up in the top results due to inappropriate ranking schemes. For algorithm searchers who are inexperienced in document search, the problem is very bad. Particularly important and relevant textual content can be accessed the platform and highlight portions by anyone with different levels of knowledge. In support of oral lessons or individual learning, the highlighted documents can be shared with the learning community. However for learners with different levels of knowledge, highlights are often incomplete or unsuitable.

We desire to have a system that automatically discovers and extracts algorithms from scholarly digital documents. Such a system could prove to assist algorithm indexing, searching, and a wide range of potential knowledge discovery applications and studies of the algorithm evolution, and presumably increase the productivity of users. This becomes a challenge for identification and extraction since algorithms represented in documents do not conform to specific styles, and are written in arbitrary formats. To overcome this disadvantage we propose this system. In this system we first detect PCs and APs by using different methods then we find the textual meta data that can be instantly extracted from various documents with the use of Application Program Interface and generate classification models fitted to different

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 2, February 2018

levels of knowledge from a set of highlighted documents to forecast new highlights. We provide linking, indexing of the extracted meta data and make it search able to the user and it can increase the productivity of user with the help of this system.

II. RELATED WORK

Prasenjit Mitra et al have studied that to identify and extract algorithm representations in a heterogeneous pool of scholarly documents, a novel set of scale able techniques used by AlgorithmSeer. Especially, hybrid machine learning approaches are said to discover algorithm representations. Then, techniques to pull out textual meta data for each algorithm are used. The user searching some algorithm on the CiteSeerX data set, this site gives so many documents with his relevant search. On the search result all document index with his best ranking and extract all data related document. This document is in the form of synopsis, on the search keywords and the algorithm and gives out-put to user. Especially they suggested detecting algorithms in scholarly documents, a set of scale able machine learning based methods and finally they justified how algorithms are indexed and made searchable. All the extracted algorithms and their related textual meta-data are then cataloged using SOLR18, which then makes the algorithms searchable.

Elena Baralis et al have studied that the concept of high-lighter. They have introduced about a HIGH-LIGHTER is a new approach to inevitably generating highlights of learning documents. Issue of automatically generating document highlights is resolved by them. High-lights are graphical signs that are regularly used to mark part of the textual content. For example, the most substantial parts of the text can be underlined, colored, or circled. The significance of text highlights in learning purpose. Teachers and learners can easily share the highlighted documents through e-learning platforms. Nevertheless, the manual generation of text high-lights is time absorbing. So minimizing such problem they generates classification models fitted to different levels of knowledge from a set of highlighted documents to forecast new highlights, which are delivered to learners to enhance the quality of their learning experience. To start the process of highlighting learning documents, they use text classification techniques. It appraises the capability level of the highlighting users to drive the generation of new highlights.

S.Bhatia et al have taken into account an algorithm search engine for software developers. To solve any problem developer first develop algorithms. Algorithms can be crucial and are very important for absolute software projects. In this system they propose an algorithm search engine that keeps abreast of the latest algorithmic developments. Using a pdf to text converter all the documents in the repository are first converted into text. To find algorithms which are then indexed along with their associated meta-data the extracted text is then examined. The query processing engine approves the query from the user by using the query interface then searches the index for relevant algorithms, and finally shows a ranked list of algorithms to the user.

J.B.Baker et al have studied that comparing approaches to mathematical document analysis from pdf. Even for born digital documents in standard formats document analysis of mathematical texts is a challenging problem. In the context of PDF documents they suggest solution to addressing this problem. For character recognition together with a virtual link network for structural analysis they identified OCR approach. To direct extraction of symbol information they use other approach from the PDF file with a two stage parser to extract layout and expression structures. With respect to character identification and structural analysis of mathematical expressions they match the efficiency and correctness of the two techniques quantitatively and qualitatively with respect to layout analysis.

J.Kittler et al have studied that combining classifiers. They focus on classifier combination. They develop structure for classifier combination and to make a decision many current schemes can be taken into account as special cases of compound classification where all the representations are used collectively. To derive the commonly used classifier combination schemes such as the product rule, sum rule, min rule, max rule, median rule, and majority voting they used various assumptions and using various approximations. The different combination schemes are then equated experimentally. A surprising outcome is come out of this. That is the combination rule evolved under the most restrictive assumptions; the sum rule outperforms other classifier combinations schemes. They investigate the sensitivity of various schemes to estimation errors to explain this empirical finding. The sum rule is most flexible to

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 2, February 2018

estimation errors as shown by the sensitivity analysis. They follow two steps. In first step they give theoretical ideas of many existing classifier combination schemes for combining the decisions of multiple experts, each employing a distinct pattern representation. In second step to improve the understanding of their properties they analyze the sensitivity of these schemes to calculation errors.

C.L.Giles. et al have studied that finding algorithms in scientific articles. Algorithms are very important part of computer science. To solve any problem first required algorithm. In this system to check for the presence of an algorithm they first examine a document. The document text is further examined to identify sentences that describe the algorithm if an algorithm is detected. Algorithm specific meta data from the document is also pulled out and indexed additionally. To calculate the connection of algorithms to a given user query this algorithm specific information is then utilized and the algorithms are submitted in decreasing order of their connection to the user. In this system they deliver a vertical search engine that identifies the algorithms present in documents and pull out and forms the related meta data and textual description of the identified algorithms. In response to user queries this algorithm gives specific information utilized for algorithm ranking.

SaurabhKataria et al studied that an important source of information that is largely under-utilized are two dimensional plots in digital documents on the web. They explain how data and text can be pulled out inevitably from these 2-D plots. For extracting data and text from two-dimensional plots they advanced automated methods from digital documents and implement it to documents published on the web. This method minimizes the time absorbing manual process of retrieving this data. The algorithm pulls out axes, the ticks on the axes, the text labels associated with the ticks and the labels of the axes. To extract each data point symbol and its textual description from the legend it discovers the legend as a text-dominated box in the figure and pulls out the lines from the legend and segments the lines. To identify their shapes and records the values of the X and Y coordinates for that point they developed a tool. Overlapping data points can be addressed by them to overcome the problem of segmenting. The data and text extraction from the 2-D plots are fairly accurate as indicated by experimental results.

D.M.Blei .et al have studied that Latent Dirichlet Allocation (LDA) technique. They developed LDA, relating to probabilistic model for accumulation of distinct data. LDA is opponent to methods such as LSI and pLSI in the setting of dimension reduction for document collection and other discrete corpora and a basic model. It is also planned to be illustrative of the way in which probabilistic models can be scaled up to offer useful circumstantial machinery in domains involving multiple levels of structure. As a probabilistic module, LDA can be easily implanted in a more complex model a property that is not possessed by LSI. This permits a richer structure in the potential topic space and in specific allows a form of document clustering that is different from the clustering that is achieved via shared topics. LDA is a three level hierarchical Bayesian model, in which each item of a collection is designed as a limited mixture over an underlying set of topics. Each topic is, in turn, modeled as an in finite mixture over an underlying set of topic probabilities. In the other words of text modeling, the topic probabilities offer an definite representation of a document. They present effective inference techniques based on variations methods and an EM algorithm for empirical Bayes parameter estimation.

S.Bhatia et al have studied that summarizing figures, tables and algorithms in scientific publications. For document-elements like tables, figures, and algorithms in digital documents to assist quick understanding by users who are searching for these document-elements identified the problem of generating synopses by them. To identify relevant sentences from the document text using a novel set of features that utilizes content and context information related to document-elements for machine-learning techniques is used by them. To identify which sentences to select in the synopses based on their similarity with the caption and the reference sentences relating to a document element and the proximity of the sentence to the reference sentences they proposed a simple model. The model attempts to balance the information content and the length of the synopsis so that the collected descriptions are both useful and accurate. To pull out useful information (synopsis) related to document elements automatically system uses the first set of methods. They use two classifiers. To identify relevant sentences from the document text based on the similarity and the proximity of the sentences with the caption and the sentences in the document text that refer to the document-element

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 2, February 2018

Naive Bayes and support vector machine classifiers are used. They study the problem of choosing the advantageous outcome synopsis size that shoots a balance between the information content and the size of the generated synopses.

III. METHODOLOGY

In our proposed system documents are processed to find out algorithm present in documents. In this system we use two text processing steps. First is stemming and second is stop-word elimination. After stemming and stop-word elimination the sentence text is transformed into a term frequency-inverse document frequency (tf-idf) matrix to study the incidents of single terms in the sentence text. One single classification model is created and used to forecast new highlights if in the training data-set there is no information about the level of knowledge of the users. Otherwise, the knowledge level of the highlighting users is calculated because it is assumed as appropriate to perform accurate highlight predictions.

AlgorithmSeer is presented as a prototype of an algorithm search engine. Plain text is extracted from the PDF file. We use PDFBox to pull out text and modify the package also to pull out object information such as font and location information from a PDF document. Then, three sub-processes operate in parallel. First document segmentation, second PC detection, and third AP detection. The document segmentation module detects sections in the document. The PC detection module detects PCs in the parsed text file. The AP detector first cleans extracted text and repairs broken sentences after that identifies APs. At last identifying PCs and APs the final step involves linking relevant algorithms. The final output would then be a set of unique algorithms as well as highlighted points.

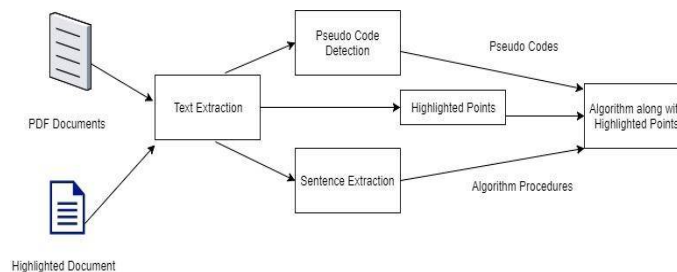


Fig. 1. Architecture of proposed system

A. Methods

TFIDF : It stands for Term Frequency-Inverse Document Frequency. The TFIDF regularly used in information retrieval and text mining. This method is used to calculate how important a word is to a document in a group. The importance increases proportionally to the number of times a word appears in the document but is 0 set by the frequency of the word in the group.

Algorithm Identification: Plain text is extracted from the PDF file. We use PDFBox to pull out text and modify the package also to pull out object information such as font and location information from a PDF document. Then, three sub-processes operate in parallel. First document segmentation, second PC detection, and third AP detection. The document segmentation module detects sections in the document. The PC detection module detects PCs in the parsed text file. The AP detector first cleans extracted text and repairs broken sentences after that identifies APs. At last identifying PCs and APs the final step involves linking relevant algorithms. The final output would then be a set of unique algorithms as well as highlighted points.

Detecting Pseudo-codes (PCs): PCs represent short and brief example of algorithms. PCs are considered as document elements. There are three methods used for detecting PCs in scholarly documents: first is Rule Based method (PC-RB), second is Machine Learning Based Method (PC-ML), and third is Combined Method (PC-CB). PCRB finds PCs by

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 2, February 2018

discovering the presence of their captions. Machine learning based (PC-ML) method directly finds the presence of PC contents. The PC-ML first finds and pulls out sparse boxes, then classifies each box whether it is a PC box or not. We suggest a collective method (PC-CB) of the PC-RB and the PC-ML using a simple heuristic as follows:

STEP1 For a given document, run both PC-RB and PC-ML

STEP2 For each PC box detected by PC-ML, if a PC caption detected by PC-RB is in proximity and then the PC box and the caption are combined.

Detecting Algorithmic Procedures (AP): AP detection approach focus to find indication of sentences and use them to discover the accompanied APs. Two methods are advanced for detecting AP indication sentences: a rule based method (AP-RB) and a machine learning based method (AP-ML). LingPipe sentence extractor is used for extracting sentences.

Stemming: This process minimizes the words to their base or root form. This step, which can be enabled or disabled according to the users preferences, remaps the textual content to a reduced number of word roots. For example, nouns, verbs in general form, and past tenses are re-conducted to a common root form. This step is specifically useful for minimizing the bias in classification when statistics based text analyses are performed.

Stop word elimination: This process is useful to filter weakly informative words i.e. the stop words. Examples of stop words are articles, prepositions, and conjunctions. In text analyses, these words are almost uninformative for predicting highlights and thus should be ignored.

IV. EXPERIMENTAL RESULTS

Our system should provide the user the PDF page of the related document in which the algorithm is present along with this user level of thinking is also considered.

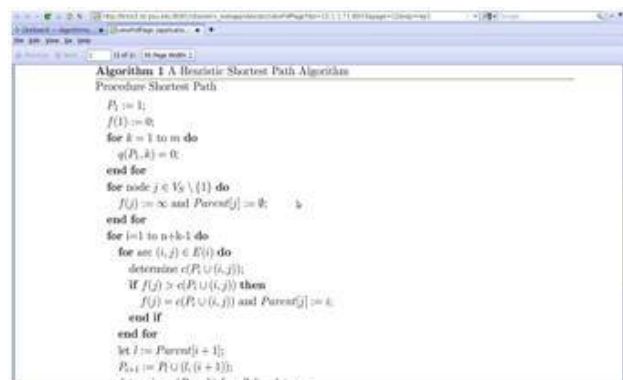


Fig. 3. Algorithm Extracted From PDF

V. CONCLUSION

Day to day researchers developed large number of an algorithm. We have discussed prototype of AlgorithmSeer. It is a search system for algorithms in large scale scholarly documents .It also provides an illustration of the actual demo system. A set of scale able machine learning based methods are proposed to detect algorithms in scholarly documents. This paper aims to detect the algorithms from collection of document then extract this algorithm from different

International Journal of Innovative Research in Science, Engineering and Technology

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Visit: www.ijirset.com

Vol. 7, Issue 2, February 2018

documents and generate unique algorithms. The result of this system provides unique algorithm to the users. Along with this it generate classification models fitted to different levels of knowledge from a set of highlighted documents to predict new highlights, which are provided to learners to improve the quality of their learning experience.

REFERENCES

- [1] Sumit Bhatia, PrasenjitMitra Senior Member and C. Lee Giles Fellow AlgorithmSeer: A System for Extracting and Searching for Algorithms in Scholarly Big Data 2016
- [2] Elena Baralis, Memberand Luca Cagliero Member Highlighter: automatic highlighting of electronic learning documents 2017
- [3] S. Kataria, W. Browuer, P. Mitra, and C. L. Giles. Automatic extraction of data points and text blocks from 2-dimensional plots in digital documents. In Pro-ceedings of the 23rd national conference on Artificial intelligence - Volume 2, AAAI08, pages 11691174. AAAI Press, 2008.
- [4] J. B. Baker, A. P. Sexton, V. Sorge, and M. Suzuki. Comparing approaches to mathematical document anal-ysis from pdf. ICDAR 11, pages 463467, 2011.
- [5] S. Bhatia, P. Mitra, and C. L. Giles. Finding algorithms in scientific articles. 2010.
- [6] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. J. Mach. Learn. Res., 3:9931022, Mar. 2003.
- [7] J.Kittler, M. Hatef, R. P. W. Duin, and J. Matas. On combining classifiers. IEEE Trans. Pattern Anal. Mach. Intell., 20(3):226239, Mar. 1998.
- [8] S. Bhatia, S. Tuarob, P. Mitra, and C. L. Giles. An Algorithm Search Engine for Software Developers. 2011
- [9] T.A. Asuncion, M. Welling, P. Smyth, and Y. W. Teh. On smoothing and inference for topic models. In Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI
- [10] G. W. Klau, I. Ljubi, P. Mutzel, U. Pfersch, and R. Weiskircher. The fractional prize-collecting Steiner tree problem on trees. Springer.