

Efficient Ranked Keyword Using AML

K.Padmapriya¹, S.P.Karthik²

¹M.E-SE Final Year, ²Assistant Professor-CSE

^{1,2}Bharath Niketan Engineering College, Aundipatti, Tamil Nadu, India,

¹ padmakannan88@gmail.com, ²karthik.sp6@gmail.com

Abstract— Entity Recognition is process of identifying predefined entities such as person names, products, or locations in a given document. This is done by finding all possible substrings from a document that match any reference in the given entity dictionary. Approximate Membership Extraction (AME) method was used for finding all substrings in a given document that can approximately match any clean references but it generates many redundant matched substrings because of approximation (rough calculation), thus rendering AME is not suitable for real-world tasks based on entity extraction. We propose a web-based join framework which combines a web search along with the approximate membership localization. Our process first provides a top n number of documents fetched from the web using a general search using the given query and then approximate membership localization(AML) is applied on these documents using the clear reference table and extracts the entities form the document to form the intermediate reference table using Edit distance Vector, Score Correlation.

Keywords— Data Mining, Approximate Membership Localization Algorithm, P-Prune, Entity Recognition.

I. INTRODUCTION

The task of Entity Recognition is to identify predefined entities such as person names, products, or locations in this document. With a potentially large dictionary, this entity recognition problem transforms into a Dictionary-based Membership Checking problem, which aims at finding all possible substrings from a document that match any reference in the given dictionary. With the growing amount of documents and the deterioration of documents' quality on the web, the membership checking problem is not trivial given the large size of the dictionary and the noisy nature of documents, where the mention of the references can be approximate and there may be mentions of non-relevant references.

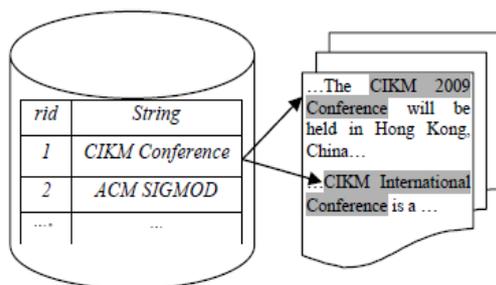


Figure: 1 using AME

The approximation is usually constrained by a similarity function (such as edit distance, jaccard, cosine similarity, etc.) and a threshold within [0, 1], such that slight mismatches are allowed between the substring and its corresponding dictionary reference. For instance, given a list of conference names like “CIKM Conference,” “ACM SIGMOD Conference,” Conference” as shown on the left part of Fig. 1, the task is to find matches from the text on the right, such as “CIKM 2009 Conference” and “CIKM International Conference,” although they do not match the string “CIKM Conference” in the dictionary exactly. The dictionary-based approximate membership checking process is now expressed by the Approximate Membership Extraction (AME), finding all substrings in a given document that can approximately match any clean references.

In this project, we propose a new type of membership checking problem: Approximate Membership Localization (AML). AML only aims at locating true mentions of clean references. AML targets at locating non-overlapped substrings in a given document that can approximately match any clean reference. We also

increase the performance of the AML process making it as a Web-Based Join Framework.

Propose the project as a web-based join framework which uses a list of elements T with an attribute $T:X$ and a clean reference list R with an attribute $R:A$, the problem is to create from the web an intermediary table RT containing valued correlations between two attributes $T:X$ and $R:A$ in order to perform a join between T and R . Given that the information available on the web can be dirty and noisy, RT shall contain the likelihood associated with its entries. Based on the hypothesis that there exist web documents containing elements of $T:X$ that also contain the elements of $R:A$, we use the elements of $T:X$ as a query for a search engine to retrieve the ranked list of documents.

First, given that the mention attribute in web pages might not be exact, we need to locate clean references approximately mentioned in given documents, for this AML method is presented second, besides given attribute, some other attributes might also occur in the web pages. To link given attribute with them, we use score correlation. We use Edit Distance Vector which is the simple function for finding the similarity when compare with cosine similarity. Edit Distance Vector reduce the extra complexity in computation because it is secondary step of extraction there is no need of high evaluation.

II. ALGORITHMS USED

Data mining can be used for extracting previously unknown patterns like groups of data records as well as to extract word and calculate between two data items. The algorithm used to find these dependencies between the data items is AML algorithm, while the algorithm used to extract groups of records is K-means clustering.

Approximate Membership Localization (AML) Algorithm

Algorithms for the AML problem based on the two assumptions below:

Assumption 1: any approximate mention m that matched with a reference consists of consecutive

words in a document, i.e., each m is a substring.

Assumption 2: only substrings whose length is up to a length threshold L are of interest, so we may as well require that $|m| \leq L$.

For the purpose of similarity functions, we regard strings as sets of words. For any word w , we use $wt(w)$ to denote its Inverse Document Frequency (IDF) weight. Given a string $S = w_1, w_2, \dots, w_k$, the weight of s is the sum of weights of all its constituents and can be denoted as $wt(s) = \sum_{1 \leq i \leq k} wt(w_i)$. The weighted Jaccard similarity of two strings s_1 and s_2 can be calculated as: $WJS(s_1, s_2) = wt(s_1 \cap s_2) / wt(s_1 \cup s_2)$. It is worth mentioning that the algorithms presented in this paper can be applied to any similarity functions such as (weighted) jaccard similarity, edit distance [9], [22], as long as it satisfies: $sim(s_1, s_2) \leq wt(s_1 \cap s_2) / wt(s_1)$ (and also $sim(s_1, s_2) \leq wt(s_1 \cap s_2) / wt(s_2)$).

In this section, we introduce two algorithms to solve the AML problem. The first algorithm is based on AME. Since the AML results are a subset of the AME results, we can do AME first, then remove the redundant pairs from the AME results second. However, this algorithm uses a lot of extra time for generating these redundant pairs and then removing them. As an alternative, a more efficient approach is to prune the potential redundant substrings before generating them, so that less time is required for generating and verifying the remaining candidate pairs.

The project as a web-based join framework which uses a list of elements T with an attribute $T:X$ and a clean reference list R with an attribute $R:A$, the problem is to create from the web an intermediary table RT containing valued correlations between two attributes $T:X$ and $R:A$ in order to perform a join between T and R . Given that the information available on the web can be dirty and noisy, RT shall contain the likelihood associated with its entries. Based on the hypothesis that there exist web documents containing elements of $T:X$ that also contain the elements of $R:A$, we use the elements of $T:X$ as a query for a search engine to retrieve the ranked list of documents Docs

First, given that the mention attribute in web pages might not be exact, we need to locate clean references approximately mentioned in given documents, for this AML method is presented second, besides given attribute, some other attributes might also occur in the web pages. To link given attribute with them, we use score correlation. We use Edit Distance Vector which is the simple function for finding the similarity when compare with cosine similarity. Edit Distance Vector reduce the extra complexity in computation because it is secondary step of extraction there is no need of high evaluation.

Scoring Correlations:

Scoring correlations for each document depends upon three relevant parameters frequency, distance and document importance.

Frequency freq: the number of times each reference is mentioned in each document of Docs.

Distance dist: the distance between the mention of each clean reference and the position of T:x.

Document importance imp: documents retrieved on the web are of different importance w.r.t. their relevance to the query, i.e., their ranks in a web search engine results.

The following equation is used to find the score correlation

$$p(r, T, x) = \frac{\sum_{d \in Docs} imp(d) \cdot score(r, d)}{\sum_{d \in Docs} imp(d)}$$

where $imp(d)$ is the importance of document d and $score(r; d)$ is the local score of r in d .

To calculate $imp(d)$, the N documents of Doc should be partitioned into B ($1 < B < N$) ranges according to their rank returned by the web search engine, so that the importance $imp(d)$ follows the normalized discount cumulative gain (NDCG) function, which is popularly used for assigning degrees of importance to web documents in a ranked list.

$$imp(d) = \frac{\log(2)}{\log(1 + \frac{[rank(d)]}{B})}$$

The document importance main depends upon the page ranking assigned to document during the web search which is not focused in the existing system. We use an inverse document based page ranking to assign rank to the documents which enhance the process of scoring correlation.

$$score(r, d) = w_a \cdot \frac{freq}{N} + (1 - w_a) \cdot \sum_{1 \leq i \leq freq} \frac{|d| - dist_i}{freq \cdot |d|}$$

P- Pruning:

Guarantee that all best match substrings can be found within. For very large reference lists, the number of these window domains, especially overlapped ones, can become an issue. In this section, we introduce how to prune domains that impossible to contain a best match substring, and how to minimize the size of the remaining domains.

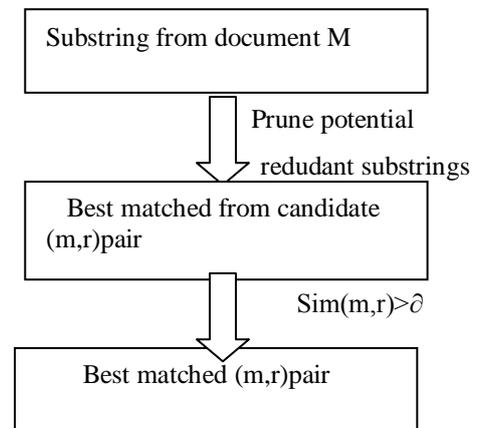


Figure: 2 pruning

Given an input document M and an entity reference list R , the task of approximate membership localization is to find all match substrings m in M for each reference r in R , such that:

1. $sim(m, r)$, where $sim()$ is a given similarity function, and tow is a given threshold between $[0, 1]$.
2. there is no substring m' that overlaps with m in any position of M which satisfies $sim(m', r') > sim(m, r)$, where r' is also a reference from R .

M- to denote the input document, R to denote the entity reference list, and

Tow-to denote the similarity threshold to the string similarity function $\text{sim}(a; b)$, where a and b are two strings

Now present the detailed algorithm of P-Prune (Potential redundancy prune). Given a reference list R, an

Inverted index is built over the words of all references in R and the strong words of each reference are also selected. For a coming document M, we generate window domains from it. Each domain is represented by $(\text{posb}; \text{pose}; r)$, where posb and pose refer to the starting and ending position of the domain in M, r is the unique reference it corresponds to. All domains are sorted according to their starting positions and stored in a set Dsort. We also use set Dact to store active domains being processed. Starting from the beginning position of M, we do:

- Step 1:** We step to the next new word position poscur , with the word $wcur$ in that position.
- Step 2:** We remove the first domain from Dsort and put it into Dact, until it can't satisfy $\text{Dsort}[1].\text{posb} = \text{poscur}$.
- Step 3:** For each domain $D = (\text{posb}, \text{pose}, r)$ in Dact, if $wcur$ presents in r, it becomes a word in D's segment, otherwise it becomes a word in D's interval. The Interval Pruning strategy is applied to all generated intervals. If $\text{poscur} = D.\text{pose}$, the Boundary Pruning and Weight Pruning are applied as well.
- Step 4:** If poscur is not included in any range of domains or it already reaches the ending position of M, then the function Locate Best Match is called to generate the best match substrings from all domains in Dact.
- Step 5:** We do step 1 to step 4 iteratively, until poscur reaches the ending position of M. We output all best match substrings we learn.

III. WORKING AND RESULTS

Over all process is the extracting document without overlap.

Data set has been collected from:

1. <http://www.cnts.ua.ac.be/conl12003/ner>

2. <http://tjongkimsang@ua.ac.be>
3. <http://fien.demeulder@ua.ac.be>
4. <http://projects.students3k.com/ieee-base-paper-and-document-download-on-networking-domain.html>
5. <http://ieeeprojectsworld.blogspot.in/2010>

From our data set we will extract the desire document with efficient manner which support of score correlation and AML algorithm. The stepwise process is below theory,

First, To extract the documents from the web using a content based search processed using the given document attributes as query input. Content based document retrieval is process of retrieving a document by search for the given keyword within each document in the document set. We select top K resultant documents from the search results for our next level of processing. Each document contains certain rank value assigned to it based on the content matching the given query attributes.

Secondly, to remove the AML problem from our process this is to extract the clean reference approximately. It uses an optimized P-Prune algorithm, which can prune potential redundant substrings from documents before generating clean reference approximately.

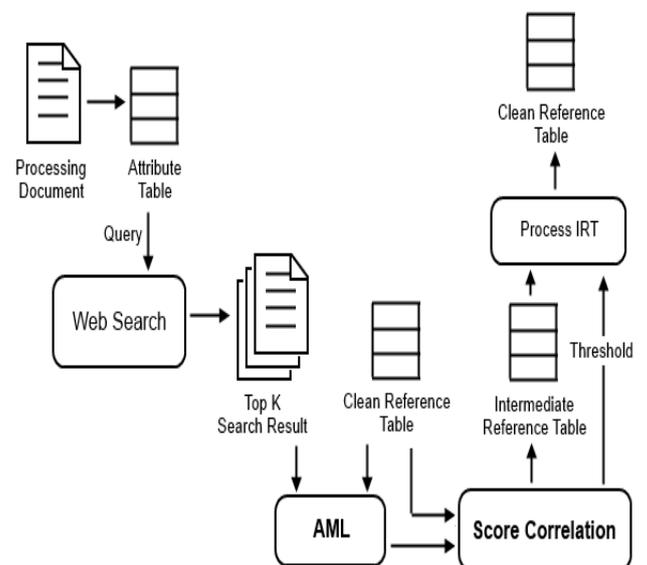


Figure 3: Framework for our method

This algorithm shows a much higher efficiency than the AME-based algorithm. After the pruning process we use a similarity function for extracting the approximate clean reference from the documents.

Thirdly, scoring correlations for each document depends upon three relevant parameters frequency, distance and document importance. The following equation is used to find the score correlation. We use an inverse document based page ranking to assign rank to the documents which enhance the process of scoring correlation.

Finally, to construct a reference table by using the results produced by the approximate membership location and score correlation modules. In this step we convert the intermediate reference table to a clean reference table considering its score correlation values. This assigned score value is used as a threshold value to filter out the approximate reference from the intermediate reference list. After that CRT will be display to user what they searching.

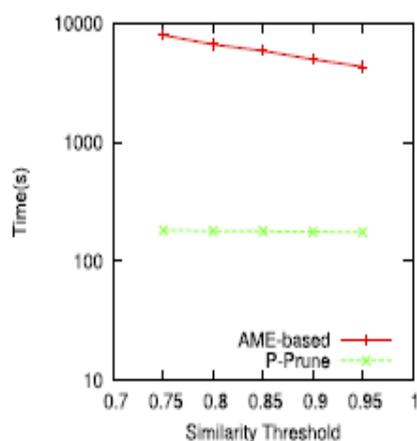


Figure 4 : AML Efficiency

IV. CONCLUSION

Formalize the AML problem and to solve it with an efficient P-Prune algorithm. The improvement of AML over AME, we apply both approaches within our proposed web-based join framework, which is a typical real-world application that greatly relies on the results of membership checking. AML-targeted solutions are more appropriate than the AME-targeted solutions for this type of real-world applications, since the matched pair results of the AML are much closer to the true matched pairs than

AME results. In Future we have to implement with large Data set and extract the documents using pruning with clear localization of the Data. also collecting keywords will be scored in cosine similarity. More efficient to AML.

REFERENCES

- [1] S.Agrawal, K. Chakrabarti, S. Chaudhuri, V. Ganti, A. Konig, and D. Xin, "Exploiting Web Search Engines to Search Structured Databases," *Proc. 18th WWW Int'l Conf. World Wide Web*, pp. 501-510, 2009
- [2] A. Arasu, V. Ganti, and R. Kaushik, "Efficient Exact Set-Similarity Joins," *Proc. 32nd VLDB Int'l Conf. Very Large Data Bases*, pp. 918-929, 2006.
- [3] R. Bayardo, Y. Ma, and R. Srikant, "Scaling Up All Pairs Similarity Search," *Proc. 16th WWW Int'l Conf. World Wide Web*, pp. 131-140, 2007.
- [4] B. Bocek, E. Hunt, and B. Stiller, "Fast Similarity Search In Large Dictionaries," *Technical Report ifi-2007.02, Dept. of Informatics Univ. of Zurich*, 2007.
- [5] H. Chan, T. Lam, W. Sung, S. Tam, and S. Wong, "A Linear Size Index For Approximate Pattern Matching," *Proc. 17th Ann. Symp. Combinatorial Pattern Matching*, pp. 49-59, 2006.
- [6] K. Jarvelin and J. Kekalainen, "Cumulated Gain-Based Evaluation Of Ir Techniques," *ACM Trans. Information Systems*, vol. 20, no. 4, pp. 422-446, 2002
- [7] N. Koudas, S. Sarawagi, and D. Srivastava, "Record linkage: Similarity Measures and Algorithms," *Proc. ACM SIGMOD Int'l Conf. Management of Data*, pp. 802-803, 2006.
- [8] L. Tanabe and W. Wilbur, "Generation Of A Large Gene/Protein Lexicon By Morphological Pattern Analysis," *J. Bioinformatics and Computational Biology*, vol. 1, no. 4, pp. 611-626, 2004.
- [9] G. Navarro and M. Raffinot, *Flexible Pattern Matching in Strings: Practical On-line Search Algorithms for Texts and Biological Sequences*. Cambridge Univ. Press, 2002.
- [10] S. Needleman and C. Wunsch, "A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins," *J. Molecular Biology*, vol. 48, no. 3, pp. 443-453.
- [11] H. Pasula, B. Marthi, B. Milch, S. Russell, and I. Shpitser, "Identity Uncertainty and Citation Matching," *Proc. Advances in Neural Information Processing Systems*, pp. 1425-1432, 2003.
- [12] S. Sarawagi and A. Bhamidipaty, "Interactive Deduplication Using Active Learning," *Proc. Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and Data Mining*, pp. 269-278, 2002.
- [13] A. Singhal, "Modern Information Retrieval: A Brief Overview," *IEEE Data Eng. Bull.*, vol. 24, no. 4, pp. 35-43, 2001.
- [14] L. Tanabe and W. Wilbur, "Generation of a Large Gene/Protein Lexicon by Morphological Pattern Analysis," *J. Bioinformatics and Computational Biology*, vol. 1, no. 4, pp. 611-626, 2004.



Mr.SP.Karthik received his B.E degree in Computer science engineering from Adiyamaan college of engineering in 2007 and Post graduate degree in computer and applications from National engineering college. His area of interest are Testing, Data mining. He is doing research (PhD) in testing from Anna University. He is the member of ISTE, IOSR.



Ms.K.Padma priya received her B.E. degree in Computer science Engg in 2010 from Ganathipathi thuls's jain engineering college, vellore and post graduate degree in Software Engineering from Anna University Chennai, India. Her areas of interest are hacking, Mobile computing, data mining. She has presented many papers in national and international conferences in various fields. As part of this paper, she is a member of ISTE