

N-ary Relation Approach for Open Domain Question Answering System Based on Information Extraction through World Wide Web

Roma Yadav, S.R.Tandan

Abstract— In this paper, we have presented n-ary relation based open domain question answering system for Extraction Information from an oversized assortment of document against arbitrary questions. We proposed two algorithms to extract entity and relationship from string and to extract answer for queried question. Our proposed algorithm works on both online and offline mode with greater accuracy

Index Terms— N-ary Relation, Open Domain, Question Answering System, Knowledge Base, Information Extraction, NLP.

I. INTRODUCTION

Question answering system came into news when [1] IBM's question answering system, Watson, defeated the two greatest Jeopardy champions, Brad Rutter and Ken Jennings, by a significant margin. It is typically a computer program that can answer in natural language, of a natural language question.

A. Question Answering System

Question Answering (QA) is a computer science discipline of information retrieval and natural language processing (NLP). The discipline deals with building a system which can analyze a query in human's natural language and can answer automatically in the form of natural language. Question Answering System (QAS) deals with a variety of question types, such as what, how, when, where, what, hypothetical, cross-lingual, semantically constrained etc.

B. Classification of Question Answering System

Question answering system can be classified into various grounds. It can be classified as:

1) Based on Domain of Questions [2, 3]

- *Closed Domain Question Answering System*: This kind of QAS deals with specific domain of topic of interest. e.g. question based on medical queries. These kinds of QAS are easy to implement [4]. Natural Language processing (NLP) system uses domain specific knowledge for extraction of answers. Sometimes closed domains refer to only specific type of question [5]. It may be descriptive rather than procedural.
- *Open Domain Question Answering System*: This kind of QAS are made to deal with almost any kind of questions

[4, 6]. The system relies on the knowledge base which may provide desired information such as local text, web pages, other databases etc. It requires bigger knowledge base.

2) Based on Response

- *Question Answering System*: eg. yahoo answer, forum, wiki answer etc
- *Automatic Question Answering System*: eg. IBM Watson

3) Based on Interaction

- *Interactive Question Answering System*
- *Non-Interactive Question Answering System*

C. Problem Description and Solution strategy

We are making effort to overcome the following disadvantages of available solutions:

- Relying on the offline data so that updated answer can be retrieved
- Storing data locally is overhead
- Speed is low

To overcome above Problem we propose following

- Do not rely on offline data instead utilize World Wide Web.
- To increase speed retrieve filtered page from a search engine it will save searching time complexity
- Perform scoring locally and make simple. No need to depend on the heavy algorithm.
- We have modeled n-ary query and answer model.

II. RELATED WORK

Question Answering System is studied by many researchers. Since it is not much older field, we find relatively small literature on it. Some important studies are as follows:

(Afader et al., 2013) [7] Studied question answering as a machine learning problem, and induce a function that maps open-domain questions to queries over a database of web extractions. The trained the function that will map a natural language question to a query over a database D. The precision was ~77%. (Niranjan et al., 2012) [8] Introduced the Rel-grams language model, which is analogous to an n-grams model, but is computed over relations rather than over words. (Alan Ritter et al., 2012) [9] Presented a scalable and open-domain approach to extracting and categorizing events from status messages of Twitter users. They used basian technique to extract event based information from twitter. The extraction is a 4-tuple representation of events which includes a named entity, event phrase, calendar date, and event type. (Etzioni et al., 2011) [10] Described the extraction of unstructured data based on structure of English grammar. They introduced Open Information Extraction paradigm which is basis for question-answering system.

Roma Yadav, M.Tech Scholar, CSE, Dr. C.V. Raman University, Bilaspur(C.G.), India

S.R.Tandan, Assistant Professor, CSE, Dr. C.V. Raman University, Bilaspur(C.G.), India

N-ary Relation Approach for Open Domain Question Answering System Based on Information Extraction through World Wide Web

(Daya C. Wimalasuriya et al., 2010) [11] Introduced Ontology Based Information Extraction using different classification techniques support vector machines (SVM), maximum entropy models and decision trees have been used in IE.

(Thomas Lin et al., 2010) described information extraction as common sense and are denoted by $f(a, b)$ where f is relation between attribute a and b . They have employed Open Information Extraction (Open IE) for relation extraction.

(Michele Banko et al., 2007) [34] performed experiments over a 9,000,000 Web page corpus that compare TextRunner with KnowItAll, a state-of-the-art Web IE system. TextRunner achieved an error reduction of 33% on a comparable set of extractions.

(Bill Dolan et al., 2004) [13] Described unsupervised techniques for acquiring monolingual sentence-level paraphrases from a corpus of temporally and topically clustered news articles collected from thousands of web-based news sources. They employed two techniques: (1) simple string edit distance, and (2) a heuristic strategy that pairs initial (presumably summary) sentences.

III. ARCHITECTURE OF QUESTION ANSWERING SYSTEM

D. Component of a Question Answering System

Typically a QAS has following component:

- Question Classifier
- Document retrieval Component
- Filter
- Answer extraction Component
- Knowledge Base

(1) Question Classifier

Question classifier module take input a question in natural language. It determines class and types of question and class and type of answer. After the analysis of question different NLP techniques are applied over the input.

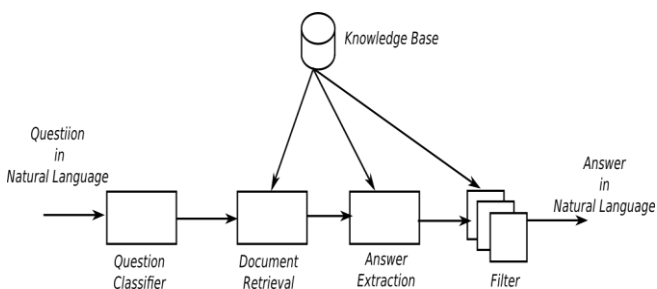


Figure 1.1: Architecture of a typical QAS.

(2) Document retrieval Component

Document retrieval module or component exploits search engines to find out relevant documents that may contain possible answer.

(3) Filter

Filter module select highly relevant document and leave trivial and non relevant documents. Then it passes output to the Answer extraction component of the question answering system.

(4) Answer extraction Component

Answer extraction component looks for answer into text depending on different context, ambiguity removing

techniques and scheme of the answer. Thus the system is able to answer most of the questions

(5) Knowledge Base

Knowledge Base is the database of the fact, document etc. It may be web document, text document or text inside database. The accuracy of the system heavily depends on the correctness and completeness of the knowledge base.

IV. METHODOLOGY

E. Stemming

Stemming is that the term utilized in linguistic morphology and knowledge retrieval to explain the method for reducing inflected (or typically derived) words to their word stem, base or root type typically a word type. several search engines treat words with constant stem as synonyms as a sort of question enlargement, a method referred to as conflation. For stemming purpose we have used tools these are Pling stemmer and snowball tools. Pling Stemmer stems associate degree English noun (plural or singular) to its singular. [15]

F. POS tagging

In corpus etymology, grammatical kind labeling (POS labeling or POST), likewise known as grammar labeling or word-classification elucidation, is that the procedure of checking up a word in an exceedingly content (corpus) as with reference to a selected grammatical feature, taking under consideration each its definition, and conjointly its setting i.e. association with near And connected words in an expression, sentence, or section.[16]

G. NER

Named-element recognition (NER) (called substance ID, element piecing and element extraction) is a subtask of data extraction that looks to find and arrange components in content into predefined classifications, for example, the names of persons, associations, areas, articulations of times, amounts, fiscal qualities, rates, and so forth we have utilized Stanford NER(Named Entity Recognizer) tools for this Purpose.[14]

H. WordNet

"WordNet is a semantic vocabulary for the English dialect. It aggregates English words into sets of equivalent words called synsets, gives short, general definitions, and records the different semantic relations between these equivalent word sets. The reason for existing is twofold: to deliver a mix of word reference and thesaurus that is all the more naturally usable, and to bolster programmed content investigation and computerized reasoning applications. WordNet recognizes things, verbs, modifiers and qualifiers in light of the fact that they take after diverse syntactic standards. [18]

I. NLP

Natural language processing (NLP) is a field of software engineering, man made brainpower(Artificial Intelligence), and computational etymology concerned with the collaborations in the middle of PCs and human (common) languages. All things considered, NLP is identified with the territory of human PC association. Numerous difficulties in NLP include common language understanding, that is, empowering PCs to get importance from human or

characteristic language data, and others include regular language era.[17]

V. PROPOSED WORK

(Afader et al., 2013) [7] The researcher has given all focuses on binary relation only. The complex questions have generally n-ary relation with more than two entity participating. If we consider more entities in the question the scoring becomes more precise. Our goal is to give a model for n-ary relation and faster algorithm to find better answer.

A. Models

The component of the question and answer are defined as follows:

Entity Set $E = \{ e_1, e_2, e_3, \dots e_n \}$

Relation Set $R = \{ r_1, r_2, r_3, \dots r_m \}$

Knowledge Base $K = E \times R \times E$

The query model:

Query $Q = (e_{q1}, e_{q2}, e_{q3} \dots e_{qi}, r)$ where

$e_{qi} \in E$ for $i = 1, 2, 3, \dots 1$ and

$r \in R$

The answer model Answer $A = \{ e_{a1}, e_{a2}, e_{a3}, \dots e_{aj} \}$

Where $e_{qi} \in E$ and

$e_{ai} \in A \wedge \forall e_{qk} \in Q \wedge r \in Q / (e_{ai}, e_{qk}, r) \in K$

B. Proposed Algorithms

We proposing two algorithms

- To extract entities and relation from string
- To extract answer for queried question

1) Entity relationship extractor ERE(S)

Input:

A string $S = \{ w_1, w_2, w_3 \dots w_n \}$

Where

w_i is word in the string and $i = 1, 2, 3, \dots n$

Output:

Knowledge Base Tuple $\varepsilon = (e_{s1}, e_{s2}, e_{s3} \dots e_{n-1}, r)$

Where

e_{si} is either entity corresponding extracted from string S and r is Relation among entities extracted from string S

Algorithm:

1. $n \leftarrow$ number of words in input string
2. St[n] is set of stemmer
3. **for** $i \leftarrow 1$ to n
4. St[i] \leftarrow StemOf(S[i])
5. **end for**
6. Pt[n] is set of POS(Parts of Speech) Tag
7. **for** $i \leftarrow 1$ to n
8. St[i] \leftarrow POSTOf(St[i])
9. **end for**
10. ε [n] is Knowledge base tuple
11. $j \leftarrow 1$
12. $k \leftarrow n$

13. **for** $i \leftarrow 1$ to n
14. **if** St[i] \leftarrow is noun then
15. ε [j] = S[i]
16. $j \leftarrow j + 1$
17. **else**
18. ε [n] = S[i]
19. $k \leftarrow k - 1$
20. **end if**
21. **end for**
22. **return** ε

2) AnswerExtraction (AE)

Input: A question String $Q_s = \{ w_1, w_2, w_3 \dots w_n \}$

Where

w_i is word in the string and $i = 1, 2, 3, \dots n$

Output: Answer Entity A_e

Algorithm:

1. Query tuple $\varepsilon_Q \leftarrow$ ERE(Q_s)
2. P[] \leftarrow is set of text paragraph extracted by search engine
3. Plength \leftarrow length of P
4. HighestScore $\leftarrow 0$
5. **for** page in P
6. **for** String in page
7. Score $\leftarrow 0$
8. Entity Relation tuple $\varepsilon_s \leftarrow$ ERE(String)
9. Score \leftarrow number of matched entity in ε_Q and ε_s
10. **if** HighestScore < Score then
11. HighestScore \leftarrow Score
12. $A_e \leftarrow \varepsilon_s - \varepsilon_Q$
13. **end if**
14. **end for**
15. **end for**
16. **return** A_e

VI. EXPERIMENT

We have simulated experiment using Bing API free subscription of 5000 by default data set for online mode and for offline mode we are utilizing freely available database of English Wikipedia[19]. We processed it so that search engine can index it. We are locally installing open source freeware search engine INDRI to index the Wikipedia data.

For evaluation and testing we have used TREC data for both online and offline mode. TREC is good source of question and answer list.

N-ary Relation Approach for Open Domain Question Answering System Based on Information Extraction through World Wide Web

TABLE I
THE CONFIGURATION OF THE SYSTEM

Configuration	Offline Mode	Online Mode
Data source	English Wikipedia dump	Bing API
Test data	TREC list of question and answer	TREC list of question and answer
Search Engine	INDRI	Bing Search engine
Stemmer	Pling Stemmer	Pling Stemmer
POSTagger	Stanford POST	Stanford POST
NLP Toolkit	Stanford NLP	Stanford NLP
Name entity recognizer	Stanford NER	Stanford NER
Wordnet Library	JWNL Wordnet Library	JWNL Wordnet Library
Processor	Intel Core i3	Intel Core i3
Clock rate	2.40GHz	2.40GHz
RAM	4GB	4GB

VII. RESULT AND DISCUSSION

The evaluation parameter for the system is accuracy of the answer.

Lets the number of query is q

The number of query correctly answered query is c

The number of query not correctly answer is e

Where $q = c + e$

Then Accuracy,

The accuracy = c/q

Or

The accuracy = $c / c + e$

Or

The accuracy = $1 - \text{error rate}$

Error Rate

The error rate = e/q

Or

The error rate = $e/c + e$

Or

The error rate = $1 - \text{accuracy}$

Our system is evaluated against 527 questions from the TREC question data in offline mode and 460 against the same question answer data set from the TREC in online mode. And the result is summarized in the table as follows.

Table 2:
Evaluation of the question answer system

Mode	Total Question	Correct Answer	Error Answer	Accuracy	Error rate
Offline	527	194	333	0.37	0.63
Online	460	202	259	0.44	0.56

the performance of our system is satisfactory as is higher than performance of base paper the bench mark performance of binary relation based on paraphrased. We can achieve higher performance by increasing size of Knowledge base.

VIII. CONCLUSION

The proposed n-ary model for open domain Question answering system is novel approach for information extraction. Our proposed algorithm produced more accurate result for fired query term. We replaced answer searching and matching algorithm of Ephyra by our proposed algorithm. Implementation of algorithm is easy and reducing the time complexity.

IX. REFERENCES

- [1] Jeopardy watson, ibm. 2011. URL http://www.nytimes.com/2011/02/17/science/17jeopardy-watson.html?_r=0.
- [2] Alfonso Valencia Roser Morante, Martin Krallinger and Walter Daelemans. Machine reading of biomedical texts about alzheimer's disease. CLEF 2012 Evaluation Labs and Workshop, September 2012.
- [3] L. Hirschman and R Gaizauskas. Natural language question answering. the view from here. Natural Language Engineering (2001), pages 275 {300, 2001.
- [4] Pampapathi R Galitsky B. Can many agents answer questions better than one. First Monday, 2005.
- [5] Boris Galitsky. Natural language question answering system: Technique of semantic headers. International Series on Advanced Intelligence, 2, 2003.
- [6] J Lin. The web as a resource for question answering: Perspectives and challenges. Proceedings of the Third International Conference on Language Resources and Evaluation (LREC 2002), 2002.
- [7] AFader, L. Zettlemoyer and O. Etzioni, "Paraphrase-driven learning for open question answering." Sofia, Bulgaria, pp. 16081618, 2013
- [8] Alan Ritter, Mausam, Oren Etzioni and Sam Clark, "Open Domain Event Extraction from Twitter", Knowledge Discovery and Data Mining, 2012.
- [9] Niranjan Balasubramanian, Stephen Soderland and Mausam, "Rel-grams: A Probabilistic Model of Relations in Text", Joint Workshop on Automatic Knowledge Base Construction and Web-scale Knowledge Extraction, 2012.
- [10] Oren Etzioni, Anthony Fader, Janara Christensen, Stephen Soderland and Mausam, "Open Information Extraction: the Second Generation", International Joint Conference on Artificial Intelligence, 2011
- [11] Daya C. Wimalasuriya and Dejing Dou., "Ontology-based information extraction: An introduction and a survey of current approaches.", J. Inf. Sci. 36, 3 (June 2010), 306-323.
- [12] Michele Banko, Oren Etzioni, Stephen Soderland, and Daniel S. Weld. "Open information extraction from the web. Commun.", ACM 51, 12 (December 2008), 68-74.
- [13] Bill Dolan, Chris Quirk, and Chris Brockett. "Unsupervised construction of large paraphrase corpora: exploiting massively parallel news sources." In Proceedings of the 20th international conference on Computational Linguistics (COLING '04). Association for Computational Linguistics, Stroudsburg, PA, USA, . Article 350, 2004
- [14] Named Entity Recognizer http://en.wikipedia.org/wiki/Named-entity_recognition
- [15] Stemmer. <http://en.wikipedia.org/wiki/Stemming>
- [16] POS Tagger. http://en.wikipedia.org/wiki/Part-of-speech_tagging
- [17] Natural language processing. https://en.wikipedia.org/wiki/Natural_language_processing
- [18] WordNet. <http://en.wikipedia.org/wiki/WordNet>
Wikipedia:Database Download https://en.wikipedia.org/wiki/Wikipedia:Database_download.