



# A Query Based Approach to Find Proportionate Entities

Mr. M. Ajith<sup>1</sup>, Mr. Ch. Ramesh Babu<sup>2</sup>

<sup>1</sup>M.Tech Student, Dept. of cse, Malla Reddy Engineering College (Autonomous), Hyderabad, Telangana India.

<sup>2</sup> Associate Professor, Dept. of cse, Malla Reddy Engineering College (Autonomous), Hyderabad, Telangana India.

<sup>1</sup>ajithreddymatta5@gmail.com

<sup>2</sup>rameshbabu.ch@mrec.ac.in

---

**Abstract**— Comparing entities by user is difficult to decide whether a query is comparable or not. To help higher intellectual process and to check entities share common functionality however has characteristics. One potential approach is a query based approach to find proportionate entities. The method is used bootstrapping approach which identifies the proportional queries and extracts the proportionate entities. This can be done by decision makers to check whether the given query is comparative or not. A successive pattern introduced that is Indicative Extraction Pattern (IEP) it may be accustomed and identify proportional queries and extract comparator pairs. These techniques identify relative question identification and equivalent entity mining. Ranking technique is used to rank the equivalent entities for user' input entity and also the results show extremely connection to uses evaluation intention.

**Keywords**— Bootstrapping algorithm, Proportionate entity finding, Information extraction, IEP, Sequential pattern mining.

---

## I. INTRODUCTION

In decision-making process, comparing various things with other one is one among the required steps that feature out. But this requires high data experience. For example throughout on-line searching a portable electronic items should have elaborated data and its specifications like Printers, Scanners, digital cameras, Graphic cards etc. In such case, it becomes troublesome for an individual with negligible data to create an honest call on that portable computer to shop for and further comparison the choice for constant. Magazine like shopper Reports, laptop Magazine and online media like CNet.com create efforts in providing editorial comparison content and survey to satisfy this would like. A judgment activity, in World Wide net era, typically involves and looks at for applicable websites insertion data relating to the targeted product, discovering difficult product, and acknowledge execs constraints. In this paper focus is on finding a collection of comparable entities provided a user's input entity. As an example, provided associate degree entity, Nokia (cell phone), wish to seek out comparable entities like Micromax, iPhone, blackberry etc. To extract comparable entities from relative matter, should always first recognize whether they are relative or not.

### Terms and concepts

#### A. Information Extraction

The process of identifying the structured information from unstructured or a semi- structured readable document is that term called as information Extraction. Methods used for information extraction

- 1) Rule-based
- 2) Pattern based
- 3) Supervised learning



### **B. Ordered Pattern mining**

Sequential Pattern mining is primarily involved with finding statistical relevant patterns between knowledge examples whenever the values are delivered in an exceptionally sequence.

### **C. Comparable Entity mining**

Proportionate entity mining achieves through the extract comparable entities from the text, questions or web corpus.

### **D. POS Tags (Part-of-speech)**

Part-of-Speech of a word could be a linguistic class outlined by its syntactical or morphological behavior. Common POS classes are: noun, verb, adverb, adjective, pronoun, preposition and conjunction. Then there are several classes that arise from totally different styles of these classes. POS tags in this paper are NN: Noun, NP: noun phrases, NNP: suitable name, NNS: Noun plural, PRP: closed-class word, VBZ: Verb, present, person particular. JJR: relative Adjective. JJS: exceptional Adjective, CC: correlative combination. Our effort on comparable entity mining is expound to the study on entity and relation removal in information taking out. In step of our definition, a comparative question has got to be a question with intention to have distinction of minimum of 2 entities. Have a tendency to exploit this approach and develop a weakly supervised bootstrapping algorithm. The main aim of this algorithm is to compare queries and extract comparable entities at a time.

### **E. Comparative Questions**

A matter whose purpose is matches two or more entities and these entities are expressly mentioned within the question.

### **F. Comparator**

An entity is a particularly comparative question that is to be compared. In step of the definitions, Q1 & Q2 below aren't comparative queries where as Q3 is. "Chennai" and "Hyderabad" are comparable key words.

Q1. "Which one is better?"

Q2. "Is Hyderabad the most effective city?"

Q3. "Which city is healthier Chennai or Hyderabad?"

The results are very helpful serving to user's exploration alternative of different decisions by suggesting the comparable entities supported other previous user's requests

## **II. INFORMATION EXTRACTION**

In terms of discovering connected things to associate a tending entity, the work is similar to the analysis on recommended systems that suggest things to a user. Recommended systems primarily consider similarities between things and their applied math correlations in user log information. Take an example of Amazon recommend artifact to its customers supported their own history. However recommending associate in tending item isn't similar to find equivalent item. In case of Amazon, the aim of this paper is to get their customers to feature more things to their looking carts by suggesting similar or connected things. Bootstrapping process has been shown to be a very effective in previous data mining analysis [9]. Work is related to them in terms of method discrimination bootstrapping technique is to extract entities with a selected relation. However, our task is totally different from their needs. Only extracted entities (comparator extraction) guarantee that the entities are extracted from proportionate queries (proportional query identification), which is usually not needed in the task.

## **III. WEAKLY SUPERVISED METHOD FOR COMPARATOR MINING**

Weakly supervised algorithm is pattern based approach similar to J&L method, but it is different in many aspects as a replacement for using separate CSR (Class Sequential Rule) and LSR (Label Sequential Rule) [10] This method aims to learn sequential patterns are used to identify proportional question and mine comparators at the same time.



### **Mining Indicative Extraction Patterns**

The indicative extraction pattern mining approach is based on two assumptions

- 1) Many reliable pairs are extracted by make use of sequential pattern comparator pairs, it is very likely to be an IEP.
- 2) An IEP is used for extract the comparator pairs.

The reliable pairs are mainly based on above two assumptions. Bootstrapping algorithm starts with a single IEP. It extracts a set of initial comparator Pairs. For each comparison pair, all questions containing the keywords are retrieved from a question collection and regard as relative entities from the proportional questions and comparator pairs, based on the reliability scores sequential pattern is generated and evaluated. Patterns evaluated through reliable Indicative extraction pattern and added into an IEP storage area. Then new evaluation pairs are extracting from the related group of questions using the new IEP. The new comparators are added to a consistent comparator repository and used as new seeds for pattern learning in the next iteration. The overview of bootstrapping algorithm is shown below, where the database store seed of pairs and questions archive from relevant data. Entities are extracted from reliable Questions. Comparators are also extracted from query collection to allow finding new patterns efficiently in further iterations. The process iterates until no more new patterns can be found from the related queries collection. There are two key steps in this method

- 1) Pattern Generation
- 2) Pattern Evaluation

#### **1) Pattern Generation**

To generate sequential pattern, we used surface text pattern mining method introduced in [9]. For any given proportionate query and its comparator pairs, equivalent key words in the query Replaced with symbol \$C. #start and #end, special words with symbols are attached to the starting and the ending of a sentence in the question. Then the following three kinds of sequential patterns are generated from sequences of queries

#### **Lexical Pattern**

Lexical pattern indicate sequential patterns consisting of only words and special symbols like (\$C, #start, and #end). They are generated with make use of suffix tree algorithm [5] with two constraints: every pattern must contain at least two special symbols like \$C, and its frequency collection should be more than an analytically determined by number. These patterns mainly concentrate on general related information and extract the related information.

#### **Generalized Pattern**

Generalized pattern can be formed from lexical pattern this pattern contain N words consequently, generalize lexical patterns by replacing the words with their POS tags and one or more including special symbols \$Cs. POS consist of what bases to compare the entities, these rules are specified.

#### **Specialized Pattern**

In a few cases, a pattern can also be too general. For example, although a query have two entities like “Nokia or Samsung?” equivalent question, pattern “<\$C or \$C>” is general. And there can be many non equivalent queries matching. Pattern for instance. “True or false?” For this reason, perform pattern specialization by adding POS tags to all comparator slots. For example, from the previous pattern “<\$C or \$C>” and the question “Nokia or Samsung?”, “<\$C/NN or \$C/NN?>” will be produced as a specialized pattern.

#### **2) Pattern Evaluation**

According to first assumption, a reliability score  $R^k(p_i)$  for a candidate pattern follows:  $\forall$  at iteration k can be defined as

$$R^k(P) = \frac{N_Q(P_i \rightarrow CP_j) \forall CP_j = CP^{k-1}}{N_Q(P_i \rightarrow *)} \quad (1)$$



Where  $p_i$  is extract known reliable comparator pair's.  $cp_i^{k-1}$  indicates the reliable comparator pair repositories accumulate until the (k-1) iteration (x) means the number of questions satisfying a condition x the  $p_i \rightarrow cp_i$  denotes that  $cp_i$  can be extract from a appropriate question by applying pattern  $p_i$ . The condition  $p_i \rightarrow *$  denotes every question contain pattern. However, Equation (1) can suffer from imperfect information about reliable comparator pairs. Example Very few predictable pairs are generally discovered in early stage of bootstrapping. In this case, the value of Equation (1) might be underestimated which could affect the performance of equation (2) an individual Indicative extraction pattern from non-reliable patterns. Reduce this problem by applying above procedure. Let us denote the set of candidate patterns at the iteration k by  $R^k$ . define the support for comparator pair  $cp_i$  which can be extracted by  $R^k$  and does not exist in the current reliable set

$$S(cp_i) = N(p_k \rightarrow cp_i) \tag{2}$$

Where  $p_k \rightarrow cp_i$  means that one of the patterns can extract  $cp_i$  in certain questions. Naturally, if  $cp_i$  can be extracted by several candidate patterns in, it is likely to be extracted as a reliable one in the next iteration. Based on this perception a pair  $cp_i$  whose support is more than a threshold  $\alpha$  is regarded as a likely reliable pair. Look ahead reliability score  $R(p_i)$  is defined:

$$R^k(p_i) = \frac{N_q(p_i \rightarrow cp_i) \vee Rel^k}{N_q(p_i \rightarrow *)} \tag{3}$$

Where Rel k indicates a set of acceptable strong pairs based on RK. Include in the Equation (1) and (3), the final reliability score  $R(p_i)$  final k for a pattern is defined as follows:

$$(p_i) \text{ final } k = \lambda \cdot R^k + (1-\lambda) \cdot R(p_i) \tag{4}$$

Evaluate all candidate patterns and select pattern whose score is more than the threshold has IEPs. All necessary parameter values are empirically determined and are determine based on values.

#### IV. MINING INDICATIVE EXTRACTION PATTERN

Proportional query analysis is a grammatical expression used to make a request for information extraction or the request made using finding proportionate entities. Questions assess proportional judgment are often phrased as directed comparison, that is, one entity is compared with other one. The equivalent queries have keywords. Admin analyze the weather queries are comparative or not. IEP implementation this defines the sequential pattern identifies the starting and ending of the sentences.

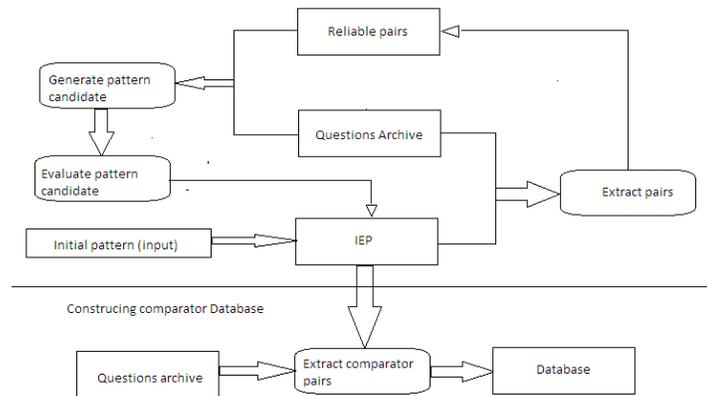


Figure: Overview of the bootstrapping Algorithm (curtesy [1])

A sequent pattern is named indicative extraction pattern if it will won't to determine proportional queries and extract comparators with high responsibilities. Then formally outline the unreliability score of a pattern. Once a matter matches an IEP, It's classified as a proportional query and therefore the token sequences admire the comparator slots within the IEP square measure extract as comparators.



### **1) Comparator Extraction**

The bootstrapping approach is prototype based mostly this approach to analyze the seed pairs. It's would like to extract the options supported the seed pairs. Review the patterns with many options. New comparator pair's mine from the queries collection exploitation the most recent IEP. The new comparators square measure value-added to a consistent comparator storage area and used as new seed for pattern learning within the next iteration. The method iterates till no additional new patterns will be found from the question assortment.

### **2) Decision Making**

Given a matter, choose the longest one in all patterns which may be applied to the question. Offer the choice creating supported pattern analysis. If a comparator is compared to several alternative significant comparators which may be again compared to the input entity, thought a valuable comparator in ranking.

### **3) Performance Evaluation**

Performance evaluation assess smart and the great proportionate query identification pattern and take out the great comparators and a decent comparator combine ought to occur in good comparative investigation to bootstrap the extraction and identification of methods and Then calculate Frequency pattern rate to see whether the choice creating is correct or not.

## **V. CONCLUSION**

This paper planned for analysis and enables several patterns to be enclosed and compared to other one. The equivalent entities are extracted from the proportionate queries that are obtained from the user. The simplest comparator is detected and comparators ranking methodologies is enforced for analyzing all the products and seek out their competitors. It's won't to dissolve the anomaly in entities. The results of comparator mining are often used for exchange search or product suggestion method.

Advantages of this paper

Improve recall and accuracy values Remove inconsistency in entities by exploitation comparator analysis with patterns. In future may be chance to improve pattern applications and mine the exceptional proportionate entities. And identify the comparator pair's aliases.

## **REFERENCES**

- [1] Shasha li,Chin-yew lin and Young-in Song 2013. "Comparable Entity Mining from Comparative Questions".
- [2] Raymond J. Mooney and Razvan Bunescu. 2005. Mining knowledge from text using information extraction.ACM SIGKDD Exploration Newsletter, 7(1):3-10.
- [3] Dragomir Radev, Weiguo Fan, Hong Qi, and Harris Wu and Amardeep Grewal. 2002. Probabilistic question answering on the web. Journal of the American Society for Information Science and Technology, pages 408–419.
- [4] Deepak Ravichandran and Eduard Hovy. 2002. Learning surface text patterns for a question answering system.
- [5] Gusfield .D, Algorithms on Strings, Trees, and Sequences: Computer Science and Computational Biology. Cambridge Univ. Press, 1997.
- [6] Haveliwala.T.H, "Topic-Sensitive Page rank," Proc. 11th Int'lConf. 2002.
- [7] Jeh.G and Widom.J, "Scaling Personalized Web Search," Proc. 12th Int'l Conf. World Wide Web (WWW '02), pp. 271-279, 2003.
- [8] Stephen Sunderland. 1999. Learning information extraction rules for semi-structured and free text. Machine Learning, 34(1-3):233–272.
- [9] D. Ravichandran and E. Hovy, "Learning Surface Text Patterns for a Question Answering System," Proc. 40th Ann. Meeting on Assoc. For Computational Linguistics (ACL '02), pp. 41-47, 2002.
- [10] N. Jindal and B. Liu, "Mining Comparative Sentences and Relations," Proc. 21st Nat'l Conf. Artificial Intelligence (AAAI '06),2006.