

# Automatic Extraction of Reliable Web Page Lists in Data Mining

<sup>1</sup> DR.K.RAMESHWARAI AH, B TECH (CSE), MS (CSE), PHD (SE), PROFESSOR & HEAD CSE,

<sup>2</sup> BOREDA KAUSHIK VARMA, PG Scholar in C S E,

<sup>1</sup> ramhyd20@gmail.com, <sup>2</sup> Kaushikv4u@gmail.com

<sup>1,2</sup> Sreyas Institute of Engineering & Technology, Hyderabad, R.R Dist, Telangana –India

**Abstract:** In today's busy schedule finding proper information within less time is important need. When user fires top-K list or any other query, user get multiple links as output. User has to visit sites and have to search the proper result manually. But one more problem is there all data available on web is not in same format. These two problems are solved by using proposed enhanced top-k list extraction system. It will give user direct top-k list as result when user fire top-k list as query. This top-k list extraction system depends upon top-k list extraction algorithm. This system handles all web pages that may be structured, unstructured or semi- structured. Also it consumes user time to find proper information.

**Keywords:** Candidate picker, parser, ranker, top-k list, title classifier

## I. Introduction

World Wide Web contains very huge amount of information. Extracting useful information from web is called as web mining. But it takes two forms: (1) Extracting structured form information (2) Extracting natural language text. The structured form information is mainly in HTML or XML language. But information tags like <li>, <table> and <ul>. Again the question arises, „is this tabular data is valuable?“. Many times the answer is NO. User may get huge tables on web but inside those tables only small amount of information is valuable.

Understanding of context is very important because in list extraction user must know the relation between listed items. For example, table which contains two columns first is Book name and second is prize of the book. This table contains five entries. But can't get why those books are listed together, e.g. are they most famous books, are they from same regional language, are they written by same author, are they published by same publisher etc. In short we

don't know under which parameters information is collected together and what is the use of this table. To avoid this situation user must know the context. Generally the context is written in

## Natural language. But machine can't interpret natural language.

As time passes technology is going to become advanced and faster. On web when user fire any query, user will not get direct result. Result is nothing but the number of links provided by search engine; user has to visit every link and has to find proper or correct data manually. It means search engine is providing most preferable or related links not the direct result. For result user visit first link if user gets result or particular information then search is stopped. Otherwise user has to visit next link if that link contain proper information the search is stopped and if not same procedure is repeated until user will get result. This normal process takes user's lot of time.

In 2008 work performed related to extracting tabular form information from HTML pages. But disadvantage of this system is it retrieves structured tabular data only. Above we discuss about tabular data available on web. In 2009, research continued with topic mining contagious and non-contagious data records. But this method is having less accuracy. Next research was performed on hybrid approach for discovering general list and extracting it from web. It gives efficient and fast result but it provides general list only not a ranked list. In 2013 evolution is done in retrieving top-k list data from web pages. The main disadvantage of this system is, it is applicable to HTML pages only.

Due to above issue this system is focusing on rich and valuable data from web that we get with the help of top-k list. This list describes k number of items of particular type. All items from top-k list are bounded with each

other through particular parameters. Following are the reasons due to which we target top-k pages for extraction:

1. Availability of top-k data on web is large and it is rich.
2. Top-k data is of high quality.
3. Top-k list data is already ranked
4. Top-k list data has interesting semantics

The paper is organized as follows: section II introduces top-k list concept, gives top-k list examples and also multiple segments of top-k title. Section III discusses detail about proposed system which includes work flow of system and proposed algorithm. Section IV presents the mathematical model. Section V discusses applications based on proposed system. The experimental results are reported in section VI. Section VII summarizes our contributions and concludes the paper.

## II. Top-K List

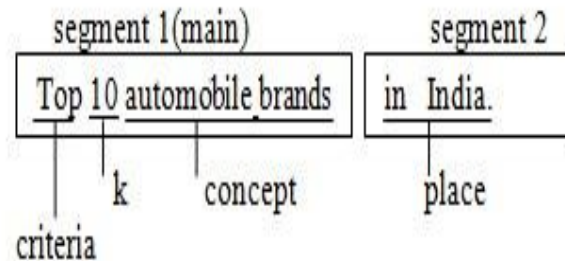
Top-k list is a list which contains k number of ranked elements. Where, k is any integer value. As compared to web tables top-k list contain rich and valuable data. Following are some examples of top-k titles:

1. Top 10 automobile brands in India.
2. Top15 android apps of 2013.
3. Five most popular social networking sites.
4. Top 10 banks in world
5. Twelve most interesting books in philosophy.

Every top-k page title contains minimum three piece of information:

- i. **k**: for example, 10, 15, five, twelve in above example. That tells how many items are in top-k page.
- ii. **concept or topic**: it tells which kind of item is retrieved. For example, automobile brands, android apps, social networking sites, banks, books etc.

- iii. **criteria for ranking**: it decides on which basis ranking is provided to provide.



Above fig shows example of top-k title. Top-k title may have many segments. The above example shows only 2 segments, First segment is main segment and second segment contain other modifiers. Segment 1 contains criteria- top, k-10 and concept-automobile brands. Second section gives information about place i.e. India. In many top-k title additional information about place or time is provide.

## III. Problem Statement

Let a web page be a pair (t, d) where t is the page title, and d is the HTML body of the page.

- A page (t, d) is a top-k page if:
- From title t we can extract a 5-tuple (k, c, m, t, l)
- Where k is a natural number, c is a noun-phrase concept defined in a knowledge base, m is a ranking criterion, t is temporal information, l is location information. Note that k and c are mandatory, while m, t, and l are optional [2].
- From the page body d we can extract k and only k items Such that:

a) Each item represents an entity that is an instance of the concept c in an is-a taxonomy;

b) The pair wise syntactic similarity of the k items is greater than a threshold.

- The top-k extraction problem can then be defined as three sub-problems (in terms of three functions):
- Title recognition F1 : (t, d) → (k, c, m, t, l)
- List extractor F2 : (k, c, d) → I

Where I is the set of terms which are instances of c and |I| = k



- Content extractor F3 : (c, d, I) →(T, S)  
Where T is a table of attribute values for the elements in I and S is its schema.

#### IV. Proposed System

The system consists of the following components:

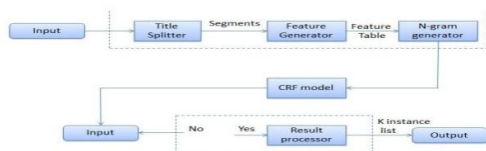
- Title Classifiers, which attempts to recognize the page title of the input web page.
- Candidate Picker, which extracts all potential top-k lists from the page body as candidate lists.

Top-K Ranker, which scores each candidate list and picks the best one;

- Content Processor, which post processes the extracted list to further produce attribute values, etc.

##### 1. Title Classifier

The title of a web page (string enclosed in <title> tag) helps us identify a top-k page. There are several reasons for us to utilize the page title to recognize a top-k page. First, for most cases, page titles serve to introduce the topic of the main body. Second, while the page body may have varied and complex formats, top-k page titles have relatively similar structure. Also, title analysis is lightweight and efficient. If title analysis indicates that a page is not a top-k page, we chose to skip this page. This is important if the system has to scale to billions of web pages.



##### 2. Candidate Picker

This step extracts one or more list structures which appear to be top-k lists from a given page. A top-k candidate should first and for most be a list of k items, Visually, it should be rendered as k vertically or horizontally aligned regular patterns. While structurally, it is presented as a list of HTML nodes with identical tag path. A tag path is the path from the root node to a certain tag node, which can be presented as a sequence of tag names.

- i) K items: A candidate list must contain exactly k items.

- ii) Identical tag path: The tag path of each item node in a candidate list must be the same.

- ii) V -Score:

V -Score calculates the visual area occupied by a list, since the main list of the page tends to be larger and more prominent than other minor lists. The V -Score of a list is the sum of the visual area of each node and is computed by:

$$Area(L) = \sum_{i \in L} (TextLength(n) \times FontSize(n) \wedge 2)$$

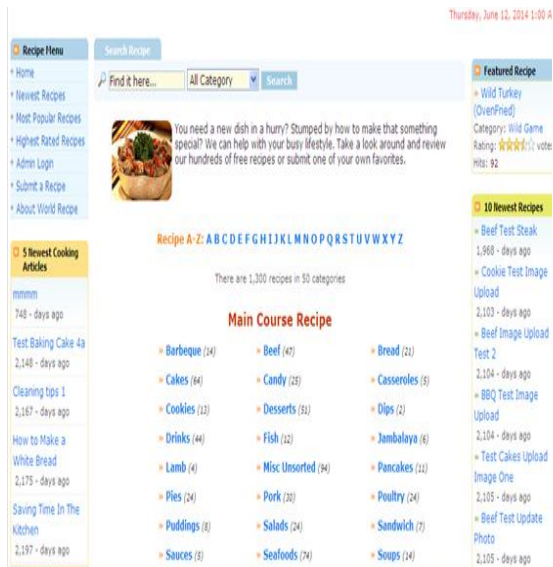
#### V. Applications

We are building a Q/A system using the top-k data to answer queries such as “tallest persons in the world”, or “What are best-selling books in 2010” directly. For this there must be a valid query provided as an input to the

product. If the query is valid then an instant result is provided otherwise there could be an unsatisfied result.

#### VI. Results

We test our system on the various online webpage and on the different domains we found that our implementation approach improves the performance of the system. N gram model improve the accuracy of result. We used html dom tree with the pruning techniques which help to minimized time complexity of our project. Following table show the performance of rule based, learning based and our proposed method.



[6] W. Wu, H. Li, H. Wang, and K. Q. Zhu, "Probase: A probabilistic taxonomy for text understanding," in SIGMOD, 2012.

[7] F. Fumarola, T. Wening, R. Barber, D. Malerba, and J. Han, "Extracting general lists from web documents: A hybrid approach," in IEA/AIE (1), 2011, pp. 285–294.

[8] G. Miao, J. Tatemura, W.-P. Hsiung, A. Sawires, and L. E. Moser, "Extracting data records from the web using tag path clustering," in WWW, 2009, pp. 981–990.

[9] M. J. Cafarella, E. Wu, A. Halevy, Y. Zhang, and D. Z. Wang, "Webtables: Exploring the power of tables on the web," in VLDB Auckland, New Zealand, 2008.

[10] Smith Tsang, Ben Kao, Kevin Y. Yip, Wai-Shing Ho and Sau Dan Lee, "Decision Trees for Uncertain Data", IEEE conference, 2011.

[11] J. Wang, H. Wang, Z. Wang, and K. Q. Zhu, "Understanding tables on the web," in ER, 2012, pp. 141–155.

[12] W. Gatterbauer, P. Bohunsky, M. Herzog, B. Krüpl, and B. Pollak, "Towards domain independent information extraction from web tables", In WWW, pages 71–80. ACM Press, 2007.

## VII. Conclusion and Future Work

Finally we would like to conclude that we have implemented the extraction of top-k list from the web. For finding out the top-k list. Also we would like to conclude that compared to other structure data top-k list are cleaner, easier to understand and more interesting for human consumption and therefore are an important source for data mining and knowledge discovery. Our project can be extended to find information from links present inside other links and try to reduce computational work.

## References

[1] Soumen Chakrabarti Mining The Web: iscovering Knowledge From Hypertext Data".

[2] Zhixian Zhang, Kenny Q. Zhu, Haixun Wang, Hongsong Li, "Automatic Extraction of Top-k Lists from the Web", IEEE, ICDE Conference, 2013, 978-1-4673-4910-9.

[3] G. Miao, J. Tatemura, W.-P. Hsiung, A. Sawires, and L. E. Moser, "Extracting data records from the web using tag path clustering".

[4] Z. Zhang, K. Q. Zhu, and H. Wang, "A system for extracting top-k lists from the web," in KDD, 2012

[5] C.-H. Chang and S.-C. Lui, "Iepad: information extraction based on pattern discovery," in WWW, 2001, pp. 681–688.