



Automatic Extraction of Top-K Lists from Web

Ashish N. Patil¹, Shital N. Kadam²

H.O.D, C.S.E Department, DACOE College, Karad, India¹

BE Student, C.S.E Department, DACOE College, Karad, India²

Abstract: This Now a days there is very busy schedule and all the users wants to complete their task within a less time. Whenever users fire any Top-k query as an input on World Wide Web then they got number of links as an output. The links which are received by users may contain non-useful information or garbage data. Sometimes these links may contain audios, videos, and Twitter and Facebook comments which is not useful for users. To overcome these problems we develop new proposed Top-k extraction system with better performance. By using this system, user will save its time to get proper result. By using this system, when user enter Top-k list as a query it will give user direct Top-k list as a result in tabular format. For obtain this result Top-k extraction system we use Top-k extraction algorithm. This system includes all web pages. The web pages is in the structured, unstructured and semi-structured format. Also gives results in less time. This paper is summaries with data which is useful to user that are extract from top-k web pages, which are web pages that web pages includes top k instances of a query which is user fire on web. That's why top-k lists are highly valuable, give richer data and information.

Keywords: Extraction of web information, Top-k lists, Extraction of list, Web mining, Data mining, Context, Web pages.

I. INTRODUCTION

This Now a days the largest source of information is World Wide Web. A lot of recent work has concentrated on earning knowledge from structured information on web, in particular, from web tables. However, it is problematic how much valuable knowledge we can extract from lists and web tables. It is true that the total number of web tables is bulky in the entire collection, but only a very poor percentage of them contain exact information. An even smaller percentage of them contain information understandable without context.

In this paper, on behalf of focusing on structured data (such as tables) and avoiding context we concentrate on text that we can know, and then we use the context to represent less structured or almost free-text data, and guide their extraction. Specifically, we focus on a beneficial and useful source of information on the web, which we call top-k web pages.

A top-k web page specifies k items of a exact interest. In most conditions, the information is in simple language text which is not directly machine understandable, although the information has the same format or different style. But most considerable, the title of a top-k page often clearly discovers the context, which makes the page understandable and extractable.

For obtaining our goal to find top k instances, when we enter search query we get unstructured data from web. Parser parse unstructured data and show top 10 closest URLs with titles. Algorithm find dust within URLs and removes that dust. Levenshtein algorithm sorts those URLs according to distance and display that distance.

HTML parser parse pages, extract data and display expected result. If that pages does not contain expected result, then parser display closest page according to search query.

II. LITERATURE SURVEY

An Automatic Extraction of Top-k Lists from the Web
Important source of structured information on the web is links. This paper is concerned with "top-k list" pages, which are web pages that specify a list of k instances of a particular query. Examples include "the 10 tallest building in the world" and "the top 20 best cricket players in India". We present an efficient algorithm that fetch the target lists with great accuracy even when the input pages contain other non-useful data of the same size or errors. The extraction of such lists can help develop existing knowledge bases about general consideration.

AUTOMATIC EXTRACTION OF DATA FROM DEEP WEB PAGE

There is large amount of information accessible to be mined from the World Wide Web. The information on the Web is in the form of structured and unstructured objects, which is known as data records. Such data records are necessary because essential information are available in these pages, e.g. lists of products and there detail information. It is important to extract such data records to provide proper information to user as per their concern. Manual approach, supervised learning, and automatic techniques are used to solve these problems. The manual method is not relevant for huge number of pages. It is a challenging work to retrieve appropriate and beneficial



information from Web pages. Presently, numbers of web retrieval systems called web wrappers, web crawler have been invented. In this paper, some current techniques are inspected, then our work on web information extraction is presented. Experimental analysis on large number of real input web URL address selections indicates that our algorithm properly extracts data in most cases.

Survey on web mining techniques for Extraction of top k list

Today finding proper result within less time is important need but one more problem is that very poor percentage information available on web is useful and interpretable and which consumes lot of time to extract. The method for extracting information from top k web pages which contains top k instances of interested topic needed to deals with system. In contrast with other structured data like web tables Information in top-k lists contains valuable and exact information of rich, and interesting. Therefore top-k list are of higher quality as it can help to develop open domain knowledge bases to applications such as search for truth result.

Extracting general from web document

In this paper, author proposed a new different technique for extraction of general lists from the web. Method uses basic premises on visual rendering of list and structural arrangement of items. The aim of system was to minimize the restrictions of existing work which deals with the principle of extracted lists. Several visual and structural features were combined for obtaining goal.

III.BACKGROUND

Big Data:

Big Data is a terminology used to mean a massive volume of both structured and unstructured data that is so huge it is difficult to process using existing database and software techniques. In most enterprise schemes the volume of data is too large or it speeds too fast or it exceeds current processing capacity.

Big Data has the potential to help companies improve operations and make faster, more intelligent outcomes. This data, when caught, formatted, manipulated, reserve, and analysed can help a company to gain useful insight to increase revenues, get or retain customers, and improve operations.

Companies are increasingly looking to acquire actionable insights into their data. Number of large data projects comes from the need to answer specific business queries. An enterprise can boost sales, increase efficiency, and improve operations, customer service and risk management with the right big data analytics platforms in place.

Due to increasing efficiency and optimizing operations the business area obtains the most attention. The big data analytics used by 62 percent of respondents to improve speed and reduce complexity.

The human side of big data analytics

Basically, the value and effectiveness of big data build upon the human operators tasked with forgiving the data and formulating the proper queries to direct big data projects. Some big data tools allow less technological end users to make various predictions from everyday business data. Still, some other tools are appearing, such as Hadoop appliances. This tools are used by various businesses to implement a suitable compute infrastructure to tackle big data projects, while minimizing the need for hardware and distributed compute software know-how.

Big data infrastructure demands

The need for big data velocity requires unique demands on the underlying compute infrastructure. For quickly process huge volumes and varieties of data can overwhelm a single server or server cluster computing power is required. To achieve the desired velocity, organizations must apply adequate compute power to big data tasks. This can possibly demand hundreds or thousands of servers that can distribute the work and operate collaboratively.

Data mining:

Computer science has interdisciplinary subfield is Data mining .It is the process of analysis patterns in large data sets involving methods at the intersection of artificial intelligence, machine learning, statistics, and database systems. The data mining process has a goal that is to extract information from a data set and transform it into an understandable structure for further use. Different from the raw analysis step, it involves database and information management aspects, information pre-processing, model and inference considerations, interestingness metrics, complexity considerations, post-processing of detection structures, visualization, and online updating. "Knowledge discovery in databases" process, or KDD is the analysis step of Data mining.

One of the analytical tool for analysing data is data mining software. This software gives users to analyse data from many different views or angles, classify it, and summarize the relationships identified. Basically, data mining is a technique to finding relations or patterns among hundreds of thousands of fields in large relational databases.

The Knowledge Discovery in Databases (KDD) process has following stages:

- (1) Selection
- (2) Pre-processing
- (3) Transformation
- (4) Data Mining
- (5) Interpretation/Evaluation.

Six common classes of task are involved in Data mining:

Anomaly detection – The identification of not usable data records, that might be interesting or data errors that require further inspection.

Association rule learning – Discovers for relationships between variables. For example, a supermarket might



collected data on customer purchasing needs. Using association rule learning, the supermarket can determine which products are quickly buy by the customer and use this information for marketing purposes. This is sometimes referred to as market basket analysis to improve market business.

Clustering – it is the task of searching groups and structures in the information that are in some way or different "similar", without using known structures in the information.

Classification –it is the task of generate known structure to apply to new information. For example, an e-mail program might attempt to classify an e-mail as "spam" or not.

Regression – attempts to search a function which models the data with the small error.

Summarization – providing a more exact representation of the data set, including analysis and report generation.
Address or URL your paper, you must type out the address or URL fully in Regular font.

WebCrawler:

To create entries for a search engine index crawler visits Web sites and reads their pages and other information in order. A crawler is a program &major search engines on the Web all have such a program. Such programs are also known as a "spider" or a "bot." To visit sites that have been submitted by their owners as new or updated, for such cases Crawlers are typically programmed. Entire sites or particular page can be selectively visited and indexed. Crawlers crawl through a site a page at a time that's why we can call it as 'Crawler. Basically, WebCrawler is a meta search engine that gives the top search results from Google Search and Yahoo! Search. There are various options to search for images, audio, video, news, yellow pages and white pages which are provided by WebCrawler.

IV. IMPLEMENTATION DETAILS

1. Levenshtein distance algorithm

In computer science, the Levenshtein distance algorithm is used to calculate distance between two sequences. Basically, the Levenshtein distance between two words is the less number of single-character edits needed to change one word into the other word. This algorithm can be published after Vladimir Levenshtein in 1965.

Levenshtein distance also be called as edit distance. It also denote a larger family of distance metrics and it is closely related to pairwise string alignments.

• Background:

The Levenshtein algorithm is used for Microsoft Outlook or Microsoft Word and Google. So we thankful of Vladimir Levenshtein for creating this algorithm.

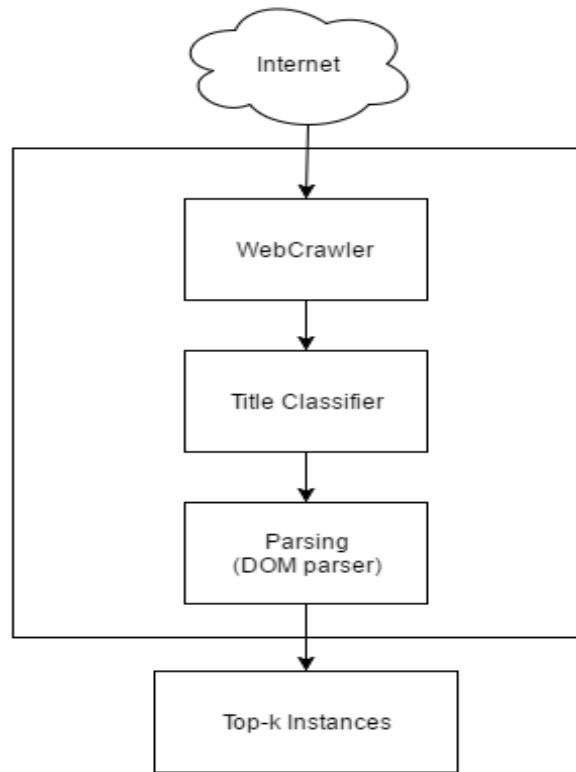


Fig. Architecture

• Example:

In example, the Levenshtein distance between "kitten" and "writing" is 3, since the following three edits change .there is no way to do this algorithm with fewer than three edits:

1. kitten→ written (substitution of "wr" for "k")
2. written→writtin (substitution of "i" for "e")
3. writtin→writing (insertion of "g" at the end).

The Levenshtein distance algorithm can have upper and lower bounds. That includes following:

- It is difference of the sizes between at least two strings.
- When the strings are equal then it is zero.
- The sum of Levenshtein distances from a third string is less than the Levenshtein distance between two strings.

2. Dust Removing Algorithm:

After the crawling, which links we get with the corresponding pages, contains the duplicate information. So, these duplicate URLs known as DUST (Duplicate URLs with Similar Text). It can be effectively explained using following example, the URLs <http://google.com/page1> and <http://page1.google.com> return the similar information. This DUST information can be created for number of reasons. For removing this DUST we created one dust list that list contains Facebook comments, twitter comments, audios, videos, YouTube, urls. By using this Dust algorithm, when user searches query then they get exact result removing the dust urls.to use this algorithm user can get its proper result within less time.



3. Flow Chart:

A web crawler also known as ants, automatic indexers, bots, web spiders, web robots. There are number of uses of web crawlers, but substantially a web crawler is used to collect or to mine data from the web. The various search engines uses WebCrawler as a means of providing up-to-date information and also used to know what's new on internet. To determine customer and market trends in a given geography, WebCrawler's are used by companies and market researchers.

After crawling, we done segmentation on that particular query. For example,



In this example,
Cricket awards 2009 is a quatifier,
Top is a criteria,
5 is a k,
Cricketers is a concept,
Score is a splitter.

Here, we use segmentation, this example is segmented into 3 parts:
Segment 1-cricket awards 2009
Segment 2-top 5 cricketers score
Segment 3-from India
Here, Segment 1 is a main part.

Classifier is used to sort the lists. Here, we use title classifier to classify the titles which is in the URLs. By using this classifier we can pick exact matching title.

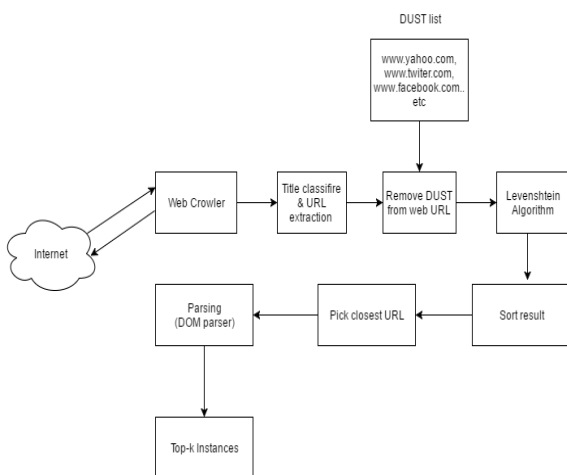


Fig. Flow Chart

After the creeping, which joins we get with the comparing pages, contains the copy data. In this way, these copy URLs known as DUST (Duplicate URLs with Similar Text). It can be adequately clarified utilizing taking after illustration, the URLs <http://google.com/page1> and <http://page1.google.com> give back the comparable data.

This DUST data can be made for number of reasons. For evacuating this DUST we made one tidy rundown that rundown contains Facebook remarks, twitter remarks, sounds, recordings, YouTube, URL's. By utilizing this Dust calculation, when client seeks question then they get correct outcome expelling the tidy urls.to utilize this calculation client can get its appropriate outcome inside less time. Now, here we use Levenshtein distance algorithm to calculate distance between two arrangements. An official suggestion of the World Wide Web Consortium (W3C) is the Document Object Model (DOM). It characterizes an interface that empowers projects to get to and upgrade the style, structure, and substance of XML records. XML parsers that bolster the DOM execute that interface.

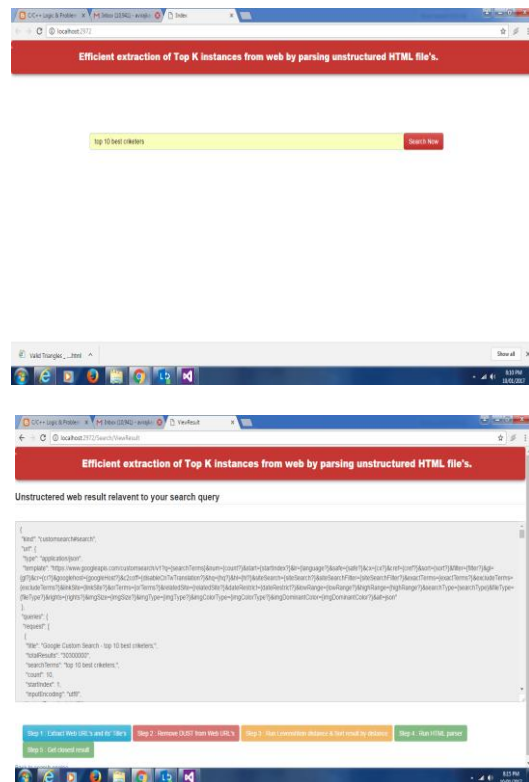
You ought to utilize a DOM parser when:

- DOM parser is utilized when you have to know a considerable measure about the structure of an archive
- You need to move parts of the report around (you might need to sort certain components, for instance)
- You need to utilize the data in the report more than once

Points of interest:

- The DOM is a typical interface for controlling archive structures. One of its plan objectives is that Java code composed for one DOM-consistent parser ought to keep running on whatever other DOM-agreeable parser without changes.

4. Screen Shots:





Efficient extraction of Top K instances from web by parsing unstructured HTML file's.

Step 1: Extract Web URL's and its Title's

#	Page Title	Web URL
1	Top 10 Greatest Cricketers of All Time - Sportology	http://sportology.com/top-10-greatest-cricketers-time/
2	Top 10 Batmen of All Time (Crickets) - TheTopTen8	http://www.thetopten8.com/cricket-batmen/
3	CIVIC 15 IN REVIEW: TOP 10 BATTING PERFORMANCES - ICC	http://www.icc-cricket.com/news/2015/features-and-specials/6797/cac-15-in-review-top-10-batting-performances
4	ICC Player Rankings - Wikipedia	https://en.wikipedia.org/wiki/ICC_Player_Rankings
5	Top 10 Cricket players of 2015-2016 - YouTube	https://www.youtube.com/watch?v=mg0d7T1kMD
6	Crickets Greatest Player of All Time	http://www.topsports.com/worlds-top-greatest-all-time-crickets.htm
7	Top 10 Greatest Cricketers of All Time - YouTube	https://www.youtube.com/watch?v=RV5A6b6kvc
8	Who are the top 10 greatest cricketers of all time? Why? - Quora	https://www.quora.com/Who-are-the-top-10-greatest-cricketers-of-all-time-Why
9	Top 10 Best Cricketers 2016 - YouTube	https://www.youtube.com/watch?v=Alqz0LqYtBU
10	Top 10 Greatest Cricket-Keepers in Cricket History	http://www.totalsport.com/cricket/greatest-cricket-keepers-of-all-time/

Efficient extraction of Top K instances from web by parsing unstructured HTML file's.

Step 2: Remove DUST from Web URL's

#	Page Title	Web URL	Is DUST Present?
1	Top 10 Greatest Cricketers of All Time - Sportology	http://sportology.com/top-10-greatest-cricketers-time/	False
2	Top 10 Batmen of All Time (Crickets) - TheTopTen8	http://www.thetopten8.com/cricket-batmen/	False
3	CIVIC 15 IN REVIEW: TOP 10 BATTING PERFORMANCES - ICC	http://www.icc-cricket.com/news/2015/features-and-specials/6797/cac-15-in-review-top-10-batting-performances	False
4	ICC Player Rankings - Wikipedia	https://en.wikipedia.org/wiki/ICC_Player_Rankings	False
5	Top 10 Cricket players of 2015-2016 - YouTube	https://www.youtube.com/watch?v=mg0d7T1kMD	True
6	Crickets Greatest Player of All Time	http://www.topsports.com/worlds-top-greatest-all-time-crickets.htm	False
7	Top 10 Greatest Cricketers of All Time - YouTube	https://www.youtube.com/watch?v=RV5A6b6kvc	True
8	Who are the top 10 greatest cricketers of all time? Why? - Quora	https://www.quora.com/Who-are-the-top-10-greatest-cricketers-of-all-time-Why	False
9	Top 10 Best Cricketers 2016 - YouTube	https://www.youtube.com/watch?v=Alqz0LqYtBU	True
10	Top 10 Greatest Cricket-Keepers in Cricket History	http://www.totalsport.com/cricket/greatest-cricket-keepers-of-all-time/	False

Efficient extraction of Top K instances from web by parsing unstructured HTML file's.

Step 3: Run Levenshtein distance & Sort result by distance

#	Page Title	Web URL	Is DUST Present?	Distance
1	Top 10 Greatest Cricketers of All Time - Sportology	http://sportology.com/top-10-greatest-cricketers-time/	False	0.428571428571429
2	Who are the top 10 greatest cricketers of all time? Why? - Quora	https://www.quora.com/Who-are-the-top-10-greatest-cricketers-of-all-time-Why	False	0.4
3	Top 10 Greatest Cricket-Keepers in Cricket History	http://www.totalsport.com/cricket/greatest-cricket-keepers-of-all-time/	False	0.375
4	Crickets Greatest Player of All Time	http://www.topsports.com/worlds-top-greatest-all-time-crickets.htm	False	0.355555555555556
5	Top 10 Batmen of All Time (Crickets) - TheTopTen8	http://www.thetopten8.com/cricket-batmen/	False	0.29220981967213
6	CIVIC 15 IN REVIEW: TOP 10 BATTING PERFORMANCES - ICC	http://www.icc-cricket.com/news/2015/features-and-specials/6797/cac-15-in-review-top-10-batting-performances	False	0.141818181818182
7	ICC Player Rankings - Wikipedia	https://en.wikipedia.org/wiki/ICC_Player_Rankings	False	0.111111111111111

Efficient extraction of Top K instances from web by parsing unstructured HTML file's.

Step 4: HTML parsing

Please wait system started the HTML parsing. This will take

Efficient extraction of Top K instances from web by parsing unstructured HTML file's.

Ranking	Player	Played	Matches	Runs	Strike Rate	Ave	100/50	Outs	
1	Sachin Tendulkar(Batting: India)	1989 +	177	290	14892	248*	56.94	51	59
2	Ricky Ponting(Batting: AUS)	1995-2010	152	239	12363	257	53.51	39	56
3	Rahul Dravid(Batting: India)	1996 +	213	265	12324	270	52.4	32	60
4	Bruce Lantshuij(Batting: NZ)	1990-2006	131	232	11913	400*	52.89	34	48
5	Jacques Kallis(Batting: SA)	1995 +	145	246	11947	204*	57.43	40	54
6	Alan Border(Batting: AUS)	1976-1994	156	285	11174	205	50.56	27	63
7	Suresh Warngad(Batting: AUS)	1985-2004	168	289	10927	200	51.86	32	50
8	Saun Gavaskar(Batting: India)	1971-1987	125	214	10822	246*	51.12	34	45
9	Udayan Madhava Jaywardene(Batting: SRI)	1997 +	119	196	9630	374	52.62	28	38

V. CONCLUSION

This paper overcomes problem of extracting top-k lists from the web. Top- k lists are cleaner, easier to understand and more interesting for human consumption as compared to other structured data, so that it becomes an important source for data mining and knowledge discovery. The Top-k extraction system will provide the highly accurate, rich and higher quality result from TOP -k web pages which contain huge number of links. The result should contains images, names, genders, locations and Address fields.

ACKNOWLEDGMENT

We would like to acknowledge our gratitude to **Prof. G. V. Deshpande** for valuable suggestions in carrying our research work. We also take opportunity to thank my friends for supporting me.

REFERENCES

- [1] Zhixian Zhang, Kenny Q. Zhu, Haixun Wang, Hongsong Li "Automatic Extraction of Top-k Lists from the Web"
- [2] Nripendra Narayan, Das Ela Kumar "AUTOMATIC EXTRACTION OF DATA FROM DEEP WEB PAGE" IJCMS ISSN 2347 – 8527 Volume 3, Issue 1 April 2014
- [3] Priyanka Deshmane, Dr.PramodPatil, Prof Abha Pathak "Survey on web mining techniques for Extraction of top k list"
- [4] B. Naresh, PG Scholar in CS, D.Sri Ram Reddy, Associate Professor, CH. Poornima, Associate Professor and HOD (C S E). "Dynamic Data Extraction of Top-K List from the Web"
- [5] K. Lerman, C. Knoblock, and S. Minton, "Automatic Data Extraction from Lists and Tables in Web Sources," Proc. IJCAI 2001 Workshop Adaptive Text Extraction and Mining, 2001; www.isi.edu/ info-agents/ papers/ lerman01-atem.pdf.
- [6] V. Crescenzi, G. Mecca, and P. Merialdo, "RoadRunner: Towards Automatic Data Extraction from Large Web Sites," Proc. 27th Int'l Conf. Very Large Data Bases (VLDB 01), Morgan Kaufmann, 2001, pp. 109–118.
- [7] C. Knoblock and A. Levy, eds., Proc. 1998 Workshop AI and Information Integration, AAAI Press, 1998; www.isi.edu/ariadne/ aiii98-wkshp/proceedings.html.
- [8] Baeza-Yates, R. "Algorithms for string matching: A survey." ACM SIGIR Forum, 23(3-4):34–58, 1989
- [9] http://www.w3.org/TR/xhtml1/ XHTML, W3C Recommendation
- [10] http://www.zvon.org/xxl/XSLTutorial/Output/ contents.html XSLT Tutorial
- [11] http://www.w3.org/TR/xslt.html XSL Transformations, W3C Recommendation
- [12] Extracting general from web document F. Fumarola,T. Weninger, R.Barber, D.Maleba and J.Han