

# DATA ORGANIZATION AND KEYWORD SEARCH: A DETAILED REVIEW

Dharmendra Singh<sup>1</sup>, Pooja<sup>2</sup>, Dhawal Vyas<sup>3</sup>

<sup>1</sup>M.Tech Scholar, <sup>2,3</sup>Assistant Professor, Department Of Computer Science & Engineering, Chandravati Education Charitable Trust Group of Institution Bharatpur.

**Abstract:** Data is the part of every document. In other word we can say that the data is the building block of the document. This paper reviews about the characterization of the data as well as the concepts of various types of search techniques used for searching in the document.

**Keywords:** Keyword Search, Graphs, DISCOVER, Structured Data, Unstructured Data

## I. INTRODUCTION

With the advance in the field of the information innovation there are number of powerful strategies and productive procedures of keyword search which are as of now being used. Keyword search method is to a great extent utilized for searching unstructured data. With time it has brought about improvement of number of methods of rating and positioning of question comes about and to appraise the adequacy of those systems. In database group principle center is around tremendous accumulation of the structured data which brought about advancement of number of artificial procedures and techniques for handling or executing the structured inquiries on the database.

In the present time, the blend of database strategies and the information recovery procedures is extremely crucial. With the colossal development of web and expanding clients of web requested prerequisite of keyword search systems and to broaden idea of keyword search over social data. Keyword search methods are exceptionally helpful for breaking down both the structured and in addition the unstructured data which contains the extensive measure of the textual information. In our research paper we will investigate different keyword search systems and we will likewise attempt to break down the territories on which we can work to enhance execution of keyword search algorithms.

## II. STRUCTURED AND UNSTRUCTURED DATA

Structured Data is one in which data is sorted out as far as structures i.e. relations or tables and that structure will take after a strict database mapping Like in SQL. These tables additionally compose the data as far as lines and sections, where lines allude to tuples or records and segments alludes to qualities and all tables are limited together with some cardinality or connections e.g. one to many, numerous to numerous et cetera as show in fig. 1.

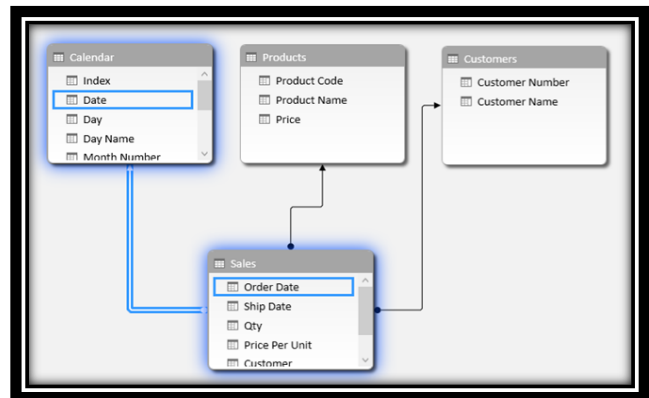


Fig 1. Structured Data

Unstructured data is absolutely inverse to structured data. It contains data that don't sorted out in any predefined Schema. It can be in any shape like Audio, Videos, JPEG Files, Pdfs, Text Files and so on and it is typically alludes to information that doesn't live in a traditional row-column database.



Fig 2. Unstructured Data

## III. KEYWORD SEARCH

Information recovery is the way toward social event information by utilizing keywords from the applicable record and that report can be unstructured or structured data. It conceals its many-sided quality from client by giving conceptual view. As client don't have any information about pattern and some other question handling dialect, he can search through dynamic interface by putting keywords. By utilizing Keyword Search client can submit keyword to search motors (Internet Search) or structured data and thusly it restores a rundown of records to client as per positioning. Positioning of reports are given in view of the keywords match and event of keyword coordinate specifically record. Positioning is given in diving request of event of keyword coordinate and the archive with most extreme event get higher need.

#### IV. KEYWORD SEARCH TECHNIQUES CLASSIFICATION

Keyword search techniques are classified into two main groups:

##### 4.1 Schema Based Keyword Search

##### 4.2 Graph based Keyword Search

##### 4.1 Schema Based Keyword Search:

Schema based approaches support keyword search over relational database (like SQL) utilizing execution of SQL summons [1]. These procedures are combination of vertices and edges including tuples and keys (essential and remote key). Each tuple in database utilizes as vertex and edges characterize interdependency among tuples.

On account of RDBMS, keyword search utilizing the Schema Based Approach is performed through making use SQL. Mapping Based approach working is partitioned into the two primary advances:

- Determine how to make and produce SQL questions so as to discover the structures among tuples.
- Determine how to assess the questions which are produced in step (I) effectively.

Discover:

Discover is procedure empowering their client to search into database by means of keywords with any question dialect Knowledge. As indicated by searching keywords Discover First Create Candidate organize chart of tuples and relations at that point diagram yield most limited arrangement first.

It plays out every single searching operation in two noteworthy strides as.

(i) Candidate Network Generator: It helps in creating all candidate networks of relations, which are known as join expressions one that produce the joining networks of tuples.

(ii) Plan Generator: It manufactures anticipates the effective and the best possible evaluation of the arrangement of candidate networks, by making utilization of the chances to reuse common sub expressions of the candidate networks [3]. Keeping in mind the end goal to produce the ideal execution design as for the real cost,

DISCOVER make utilization of the eager calculation. One primary part of the DISCOVER is that, it performs keyword search without utilizing the prerequisite of the client to know outline of the database.

For positioning the outcome, DISCOVER restores a monotonic score aggregation function [5]. The principle disadvantage of this calculation is that the cost of producing CNs set is high [4].

Spark:

The interest for RDBMS to help keyword search on content data is expanding as there is wide increment in the content data put away in the relational databases. The current keyword search techniques are not achievable for the content data search. The fundamental point of these systems is to center around viability and productivity of the keyword question search [7]. Spark concept devise another positioning equation by making utilization of the current information recovery procedures. The principle utilization of the Spark is that it takes a shot at vast scale genuine databases (Eg. Client Relationship Management) by taking in consideration both the RDBMS adequacy and proficiency.

It influences utilization of the Top-k to join calculation which incorporates two effective inquiry handling algorithms for positioning function.

(a) Dealing with Non-monotonic scoring function.

A non-monotonic function is a function that is expanding and diminishing on various interims of its area.

For instance, consider our underlying case  $f(x)$  measures up to  $x^2$ . We saw that this function is expanding on the interim  $x$  is more prominent than 0, and diminishing on the interim  $x$  is under 0. Since the function is expanding and diminishing on various interims of its area, the function is a non-monotonic function. Fundamentally, if a function isn't expanding on its whole space or diminishing on its whole area, at that point the function isn't monotonic, and we say that it is non-monotonic.

(b) Skyline Sweeping Algorithm.

The Skyline Point Algorithm includes:

(I) Block Nested Loop.

(ii) Divide and Conquer.

(iii) Plane-Sweep.

(iv) Nearest Neighbor Search.

(v) Branch and Bound Algorithms.

(I) Outline of Block Nested Loop

1. Look over a rundown of point and test each point for predominance criteria.

2. Rundown of potential horizon focuses seen so far are kept up by fulfilling a solitary dimension, each went by point is contrasted and all components in the rundown. The rundown is appropriately refreshed.

3. This calculation completes a ton of repetitive work. It has no provision for early termination. Add up to work done relies upon the request in which focuses were experienced..

(ii) Divide-and-Conquer

1. The algorithm recursively separates expansive datasets into smaller partitions. The algorithm continues till each smaller partition of the dataset fits in the primary memory.

2. We figure the halfway horizon for each partition utilizing any in-memory approach and later consolidate these fractional horizon points to shape the last horizon query.

(iii) Nearest Neighbor Search

1. Expect that a spatial record structure on the data points is accessible for utilize.

2. Distinguishes horizon points by rehashed application of a closest neighbor search procedure on the data points, utilizing a reasonably characterized L1 separate standard.

(iv) Branch and Bound Algorithm

A R-tree is based on the data points. Construct a need line that orchestrates protests in a MinDist requesting with respect to the root.

• Variations of the Skyline Point Query

1. Positioned horizon queries: an other inclination function is utilized rather than the base criterion.

2. Constrained horizon queries: The horizon query returns horizon points only from the data-space characterized by the constraint.

3. Identifying queries: For every horizon point in the dataset, locate the quantity of points in the dataset ruled by it.

4. K-Dominating queries recover the  $\star$  points that rule the biggest number of points in the dataset.

• Implementation of the Algorithm

Info: Given some points with two co-ordinates spoke to in a 2-D Coordinate framework.

Yield: The arrangement of Skyline points are to be distinguished.

1. Sort the points as indicated by x – organizes and allocate the lists e.g. 1,2,3...N points.
2. Make an unfilled Stack. It should entirely keep up LIFO(Last In First Out) property. The stack must have the capacity to hold the scope of files (begin point and end point of the scope of qualities). It should also store the middle point's file. Make a stack container that can hold 3 numbers.
3. Keep a worldwide esteem least.
4. Divide the points based on its middle component into left and right. We continue to do as such and now embed another component in the stack whose new range is the present left range and also store the middle of the left range. This procedure continues till the left side has only a solitary point.
5. Stamp the single point as a horizon point as it is the furthest left point.
6. Set least to be equivalent to this present point's y-facilitate.
7. Presently think about whatever is left of the points in the correct half (if any point has y-organize not as much as the base at that point stamp the point as a horizon and set the base equivalent to this present point's y-arrange).
8. Once every one of the qualities for the correct half are done, we pop the stack and then we continue 'stage 7' with the present right half until the point when all the horizon points are acquired.

4.2 Graph Based Keyword Search:

A data graph DG is the reasonable representation of Relational Database. In this graph algorithm is reached out to solve the keyword search queries. In this representation we have,  $DG(V, E)$ , where DG is the coordinated graph, v is the set of the vertices speak to data or record and E is the set of edges which also characterize relationship substance. In this graph, weights are dole out to the edges so as to speak to the vicinity of the corresponding tuples, for instance, we have two vertices u and v then the nearness of u and v is spoken to by a weight signified by  $We\{(u,v)\}$ .

Types of Graph Based Search:

(a) BANKS:-

BANKS stands for Browsing and Keyword Searching. BANKS framework speaks to relational model into data graph and as per coordinating keyword it enacts graph hubs.

• Structuring of BANKS algorithm:

In BANKS framework, the database is spoken to utilizing a coordinated graph and the record or tuple is spoken to as a hub in the coordinated graph. Outside Key or Primary keys are exhibited utilizing the edge, which corresponds to the connection between the corresponding tuples. The outcome handling of BANKS framework will returns a sub-graph which is spoken to in type of connecting hubs, one which coordinates the query keyword. This sub-graph can also additionally refined with a specific end goal to get the more precise or more suitable response for the keyword we have searched.

The focal hub in this graph is one which connects all the

keyword hubs and relationship between them. And this focal hub is known as root hub and answer can be considered as established coordinated tree which contains guided way from root to every keyword hub. The root hub is information hub and tree is connection tree.

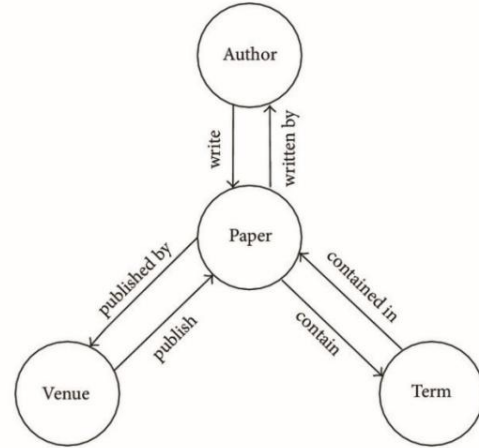


Fig. 3 DBLP Bibliography Database (A Fragment of Database.)

BANKS framework makes utilization of two kinds of datasets, DBLP and little theory database. DBLP database is appeared in Fig. 3, which also demonstrates how DBLP is converted into structured relational configuration. There were 124,612 hubs and 319232 edges on the graph.

Representation of the BANKS framework consists of vertices, edges, edge weights, hub weights. Each tuple T in database is spoken to by a hub T in the graph. Assume that there are two tuples T1& T2 with the end goal that they have outside key relationship, at that point in the graph this is spoken to by an edge from T1 to T2 and also a back edge from T2 to T1.

The weight for forward connection along outside key relationship reflects nearness relationship between the tuples and it is set to 1 as a matter of course. In the present implementation of BANKS the forward edge weight is set to  $s(R(u),R(v))$  and the invert edge weight is set to  $[s(R(v),R(u))*INv(u)]$  and the genuine weight is the base of the two as takes after:

$$b(u,v)=\min(s(R(u),R(v)), s(R(v),R(u))*INv(u))$$

BANKS is fundamentally used as a piece of demand to circulate organizational data, bibliographic data and electronic lists. BANKS causes us in expelling the critical data without the learning of blueprint of the database [8]. A customer can fundamentally remove the required data by creating a couple of keywords, by then after hyperlinks and working together with controls on the demonstrated comes to fruition.

Keyword searching in BANKS is done by influencing utilization of the region to construct positioning in light of remote key connections. BANKS essential ideal position it that it reduces the undertakings associated with disseminating social data on the web and additionally it makes it more available..

(b) Data Spot

Data Spot is a database conveying instrument and it lets the

end customer to research the substantial database without making utilization of any request vernacular. DataSpot make utilization of plan less semi-composed chart which is known as hyperbase. As showed by the possibility of the DATASPOT, the Search Server performs looks for inside the hyper base and consequently as needs be it returns either HTML pages or question API[9]. The Data Spot used as a piece of electronic inventory, business list, described promotions, help work regions and back.

•Model Description:

A graph structure which is used for the data depiction in Data Spot is known as hyper base. Hyper base includes hubs, edges, and center point names. Hubs are associated through composed edges. These organized edges are additionally designated fundamental edge and identification edge. The direct edge is an edge amidst the parent center point and the tyke center point. These game plans of kids hubs are asked for while the parent hubs are not asked. The direct edges can't make a cycle. And the second kind is identification edge, it is used to indicate reference and subject center point relationship. The reference center phenomenally perceives subject center point and subject center point can have most extraordinary single reference center. At first hubs in hyperbase address data questions and edges address relationship between them. DataSpot distributor makes an interpretation of the source data into hyperbase, which accordingly capably questioned by DataSpot server.

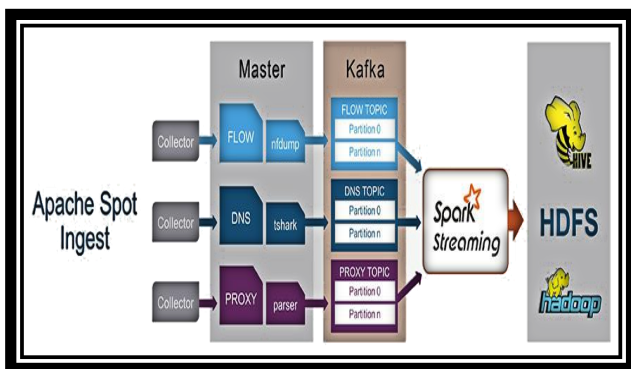


Fig. 4: The Data Spot Architecture

(c) Proximity search

Proximity search tackles general connections among objects every together question answers, these strategies obliging for the smart request sessions. In content Processing Proximity Search searches for Documents or content records where no less than two term occasions according to organize are inside decided separations. Where isolate is number of widely appealing words. Web Movie Site Database (IMDB) site impacts utilization of the proximity to search remembering the true objective to answer its database request. IMDB destinations include 140,000 motion pictures and information about in excess of 500,000 film industry workers. The thought driving is that the database can be viewed as set of connected articles, where objects address films, entertainers, boss and so on. And the separation work in light of joins isolating items [10].

Model Description:

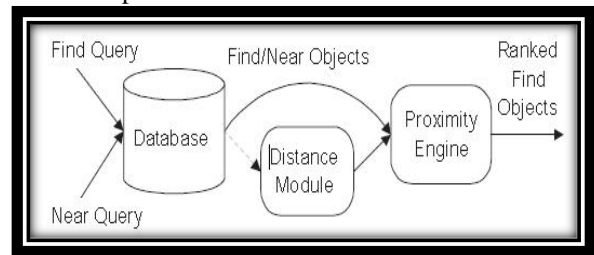


Fig. 5 Proximity Search Model

The database is addressed through graph in which data (objects) addressed as vertices and relationship. Proximity wears down the most short separation between objects. For proximity searching database is viewed as accumulation of items where objects are connected by independent capacity. The Fig. 5 demonstrates Proximity Search Model.

V. CONCLUSION

This paper reviews the categorization and gives us the detailed information about the various types of search available for searching the documents and gives the detailed information regarding the search also.

REFERENCES

- [1] Lu, Yue & Tan, Chew Lim.,” Keyword searching in compressed document images”. DCC,2003..
- [2] S. S. Pawar, A. Manapatil, A. Kadam and P. Jagtap, "Keyword search in information retrieval and relational database system: Two class view," International Conference on Electrical, Electronics, and Optimization Techniques (ICEEOT), 2016.
- [3] Q. Dong, Z. Guan and Z. Chen, "Attribute-Based Keyword Search Efficiency Enhancement via an Online/Offline Approach,” IEEE 21st International Conference on Parallel and Distributed Systems (ICPADS), 2015.
- [4] Kehinde K. Agbele, Kehinde Daniel Aruleba, Eniafe F. Ayetiran," Efficient schema based keyword search in relational databases." University of Computer Studies, Mandalay, Myanmar, International Journal of Computer Science, Engineering and Information Technology (IJCEIT) 2.6 (2012).
- [5] Sanjay Agrawal, Surajit Chaudhuri, Gautam Das,"DBXplorer: enabling keyword search over relational databases",SIGMOD,2002.
- [6] A. Karapakula, M. Puramchand and G. M. Rafi, "Coordinate matching for effective capturing the similarity between query keywords and outsourced documents," IET Chennai 3rd International on Sustainable Energy and Intelligent Systems (SEISCON 2012), Tiruchengode, 2012.
- [7] W. Tang, L. Yan, Z. Yang and Q. H. Wu, "Improved document ranking in ontology-based document search engine using evidential reasoning," in IET Software, vol. 8, no. 1, pp. 33-41, February 2014.
- [8] Shengli Wu, Jieyu Li, "Merging Results from Overlapping Databases in Distributed Information

- Retrieval",PDP,2013.
- [9] A. Lakhani, A. Gupta and K. Chandrasekaran, "IntelliSearch: A search engine based on Big Data analytics integrated with crowdsourcing and category-based search", International Conference on Circuits, Power and Computing Technologies , 2015.
- [10] Roy Goldman, Narayanan Shivakumar, Suresh Venkatasubramanian, Hector Gercia Molina "Proximity Search In Database" In Proceedings of the 24th VLDB Conference, New York, USA, 1998.
- [11] Gary Pan, SeowPoh Sun, Calvin Chan and Lim Chu Yeong,"Analytics and Cybersecurity: The shape of things to come", CPA,2015