

# Efficient Nearest Keyword Set Search Using Term Frequency and Inverse Document Frequency in Datasets

C. VENKATESH<sup>1</sup>, BUSHRA TAHSEEN<sup>2</sup>

<sup>1</sup>PG Scholar, Dept of CSE, DR. K.V Subba Reddy Institute of Technology, Kurnool, Andhrapradesh, India,

<sup>2</sup>Assistant Professor, Dept of CSE, DR. K.V Subba Reddy Institute of Technology, Kurnool, Andhrapradesh, India.

## ABSTRACT

We focused on multi-dimensional dataset where each data point has set of keywords in feature space allows for the development of new tools to query and explore these multidimensional dataset. It is hard to investigate the huge dataset for a given query as well as to accomplish more accuracy on user query. Frequently query will search on dataset for accurate keyword match and it does not discover the nearest keyword set for accuracy. Vishwakarma et.al. Proposed a strategy called ProMiSH (Projection and Multi Scale Hashing) that utilizes random projection and hash based index structures to accomplish high

scalability and speedup. However, the memory usage of Random projection and multi scale hashing increases when the number of dimensions of data points increases. So, there is a require to assign some weights to the keywords to find the nearest keyword set for accuracy. The Term Frequency (TF) and Inverse Document Frequency (IDF) weighting technique are utilized to discover the efficient nearest keyword set search. Performance of those techniques provide better performance in terms of accuracy, dimension reduction rate and response time.

## I.INTRODUCTION

In today's digital world the amount of data, which is developed, is increasing

day by day. There is different multimedia in which data is saved. It's very difficult to search the large dataset for a given query as well to archive more accuracy on user query. In the same time query will search on dataset for exact keyword match and it will not find the nearest keyword for accuracy. Vishwa karma et.al. proposed ProMiSH (Projection and Multi-Scale Hashing), which is capable for preparing the Nearest watchword set queries fastly. ProMiSH ensures that it finds the best possible top-k results. ProMiSH is more effective in terms of response time and accuracy and is capable of obtaining near-best results. ProMiSH-E that discovers best subset of points and ProMiSH-A which searches near optimal results with better efficiency. Nearest keyword queries always need coordination for queries, so it is difficult to develop a capable method to solve NKS queries. This paper proposes a scoring method for position of the result sets. It includes design and implementation of a TF -

IDF algorithm that can assign weights to the keywords of point. It also scores the document using cosine similarity model. The proposed system provides a considerably better performance in terms of exactness, dimension reduction rate and response time.

## II.RELATED WORK

A selection of queries semantically changed from NKS queries, have been studied on text rich spatial datasets. By Long et.al Location-definite keyword queries on the web and in the geographic information systems were previously answered using a combination of rectangular tree and inverted list. Felipet. al. developed Information retrieval tree to place objects from spatial datasets based on a grouping of their distance to the query locations and the significance of their text descriptions to the query keywords. Cong et al. included rectangular tree and inverted list file to answer a query similar to Felipe et al. using a different position function. Zhou et al. calculated text

importance and location proximity independently, and then joined the ranking scores. Cao et al. and Long et al. projected algorithms to improve a group of latitude and longitude web objects such that the group of keywords cover the query keywords and the objects in the groups are nearby to the query location and have the lowest inner-object spaces. Other correlated queries comprise collective nearest keyword search in spatial databases, top-k privileged query, top-k sites in a spatial data based on their ability on characteristic points, and best possible location queries, are explained. Tree-based indexes, such as R-Tree and M-Tree, have been widely investigated for nearest neighbour search in high dimensional spaces. These indexes neglect to scale to dimensions more than 10 due to the curse of dimensionality. Recently vishwa karma et. al. suggests a new method called projection and multiscale hashing based on random projection and hashing. By using this index they developed ProMISH, it

finds the best separation of points exactly. and it searches near best results with better efficiency. They proved that ProMISH is quicker than tree based technique, with multiple orders of magnitude performance development. And the memory usage of this technique grows slowly when the number of dimensions in data points increase. However, their work does not support scoring schemes for ranking the result sets. This paper presents a Term Frequency and Inverse Document Frequency that needs assigning weights to the keywords of points. Then, each group of points can be scored based on distance between points and weights of keywords.

### III.PROBLEM STATEMENT

The amount of data which is developed is increasing day by day, thus it is very difficult to search large dataset for a given query as well to achieve more accuracy on user query.so we need a method for efficient search in multidimensional dataset. Existing system contains

ProMiSH (Projection and Multi-Scale Hashing) is done on keyword queries. By using this index they developed ProMISH, it finds the best separation of points exactly and it searches near best results with better efficiency. They proved that ProMISH is quicker than tree based technique, with multiple orders of magnitude performance development. And the memory usage of this technique grows slowly when the number of dimensions in data points increase. However, their work does not support scoring schemes for ranking the result sets.

#### IV. IMPLEMENTATION

We consider the inverted index and hashing techniques for getting the results. To improve the results efficiently we use ranking in the proposed system. Ranking gives the results in such a way that the top one will be first. Ranking provides the result based on keyword ratio. Ranking is done by tf-idf technique. Tf-idf defines term frequency and inverse document frequency. The

proposed work assigns weight to the keywords of points using Term Frequency and Inverse Document Frequency algorithm. Where each group of points can be scored based on distance between points and weights of keywords. The Term frequency and inverse document frequency weight method is constructed with two words: the first word is Term Frequency (TF) and second word is Inverse Document Frequency (IDF). We give high score to the word which is important and that appears frequently in a document. We give low score to the word means which is not an uncommon identifier, and that appears in many documents. Some Common words like "the" and "for" which are shown in many documents. We can compute the term frequency as the number of times keyword appear in a document divided by total number of words in a document. Inverse Document Frequency provides how much information the word provides that is either common or rare across all the

document. We can compute inverse document frequency as the total number of documents by number of documents that containing a term „t“. For example, document having 100 words, the term 'rat' appear 13 times, the Term frequency of the word 'rat' is  $TF(\text{rat}) = 13/100$  i.e. 0.13. IDF (Inverse document frequency) of a word is measure of how significant that term is throughout the web. For example, say the term 'rat' appears 10 million times in the whole corpus (i.e. web). Let's assume there are 0.4 million documents that contain such a huge number of 'rat', then the IDF (i.e.  $\log \{DF\}$ ) is given by the total number of documents divided by the number of documents containing the term 'rat'. IDF (Inverse document frequency) of a word is measure of how significant that term is throughout the web. For example, say the term 'rat' appears 10 million times in the whole corpus (i.e. web). Let's assume there are 0.4 million documents that contain such a huge number of 'rat', then the IDF (i.e.  $\log$

{DF}) is given by the total number of documents divided by the number of documents containing the term 'rat'.

$$IDF(\text{rat}) = \log(10,000,000/400,000)$$

i.e. 1.63

$$\therefore W_{\text{rat}} = (TF * IDF)_{\text{rat}} = 0.13 * 1.63 = 0.2119.$$

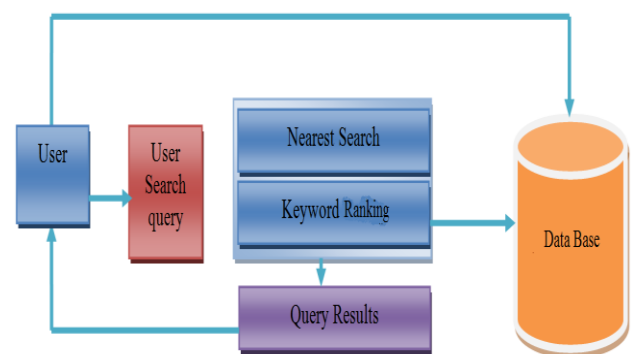


Figure: System Architecture

## V. MODULES

1. Multi-dimensional data
2. nearest Keyword
3. Indexing
4. Hashing.
5. Ranking

### 1) Multi-dimensional Data

Multi-dimensional datasets is data analysing procedure that covers data into two categories. They are data

dimensions and measurements. Keyword-based search in text-rich multi-dimensional datasets provides many peculiar applications and tools. Multi-dimensional datasets has data points and each point has a set of keywords. The presence of keywords in feature space allows for the development of new tools for querying and analysing these multi-dimensional datasets. Algorithms may take huge amount of time to terminate a multi-dimensional dataset which consists of millions of points. Therefore, there is a necessary with an efficient algorithm that scales with dataset dimensions, and yields practical query efficiency on huge datasets. In multi-dimensional spaces, it is very difficult for users to provide meaningful coordinates, and our main challenge is to provide keywords as input.

## 2) Nearest Keyword

Let us consider multi-dimensional datasets where each and every data point has a set of keywords. The existence of keywords in vector area

allows for the development of new mechanism for querying and analysing these multidimensional datasets. An NKS query is a query in which user will provide set of keywords, and the outcome of that query may includes k sets of data points from which it contains all the query keywords and forms one of the top-k thick clusters in the multi-dimensional space. Location-specific keyword queries on the internet and in the GIS systems were earlier answered using R-Tree and inverted index. In nearest keyword to rank objects from spatial datasets, it is done based on a combination of their distances to the query locations and the relevance of their text descriptions based on the query keywords is done by IR-tree.

## 3) Indexing

Indexing time is the metrics to evaluate the index size for Projection And Multi Scale Hashing With Ranking. Indexing indicates the amount of time used to build



Projection And Multi Scale Hashing With Ranking variants. The memory usage and indexing time of Projection And Multi Scale Hashing With Ranking is good. Memory handling increases slowly in Projection And Multi Scale Hashing With Ranking when the number of scope in data points increases. Projection And Multi Scale Hashing With Ranking is more efficient and it takes 80% less memory and 90% less time, and is able to obtain near-optimal results.

#### 4) Hashing

The hashing technique is influenced by Locality Sensitive Hashing (LSH), which is a state-of-the-art method for nearest neighbor search in high-dimensional spaces. Hashing is done based on the records and methods used for hashing. Based on the index value obtained, we use to place the record at the particular bucket. Unlike LSH-based methods that allows only rough search with probabilistic guarantee, the index structure in

Projection And Multi Scale Hashing With Ranking supports accurate search. Random projection by hashing has become the state-of-the-art method for nearest neighbor search in the highdimensional datasets.

#### 5)Ranking

Ranking is done based on the keywords. Number of times the keyword appears in a document and number of times the keyword appeared throughout the whole dataset. In this ranking is done by using Tf-Idf technique.

#### VI.CONCLUSION

we propose Term Frequency (TF) and Inverse Document Frequency (IDF) weighting technique are utilized to discover the efficient nearest keyword set search. Performance of those techniques provide better performance in terms of accuracy, dimension reduction rate and response time.

#### VII.REFERENCES

[1]VishwakarmaSingh,B.Zong,Ambuj K.Singh,"Nearest Keyword Set Search in Multi-Dimensional Datasets",vol.28,no.3,March 2016.

- [2] C. Long, R. C.-W. Wong, K. Wang, A. W.-C. Fu, "Collective spatial keyword queries: A distance owner-driven approach," in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2013, pp. 689–700.
- [3] Z.Li, H.Xu, Y.Lu and A.Qian,"aggregate nearest keyword search in spatial databases, "in Proc .12th Int.Asia Pacific Web Conf., 2010, pp.15-21.
- [4] A.Guttman,"R-Trees:A dynamic index structure for spatial searching, "in Proc.ACM SIGMOD Int.Conf. Manage. Data Eng.,1984,pp.47-57.
- [5]X.Cao,G.Cong,C.S.JensenandB.C.Ooi,"Collective spatial keyword querying, "in proc.ACM SIGMOD Int.Conf.Manage.Data, 2011, pp.373-384.
- [6]M.Datar,N.Immorilica,P.Indyk,and V.S.Mikkorkni,"Locality sensitive hashing scheme based on p-stabledistributions,"inProc.20thAnnu. Symp.Comput.Geometry, 2004,pp.253-262.
- [7] Y.Zhou, X.Xie, C.Wang, Y.Gong, and W.-Y.ma,"Hybrid Index structure for location-based websearch,"inproc.14thACMInt.Conf. Inf.Knowl.Manage. 2005, pp.155-162.
- [8]R.Hariharan,B.Hore,G.Li,andS.Me hrotra,"processing spatial keyword queries in geographic information retrieval(GIR)system,"inproc.19thint. Conf.Sci.StatisticalDatabaseManage., 2007,p.16.
- [9]S.Vaid,C.B.Jones,H.Joho,andM.Sa nderson,"Spatio-textual indexing for geographical search on the web," in Proc.9th Int.Conf.Adv.Spatial Temporal Databases, 2005, pp.218-235.
- [10]I.DeFelipe,V.Hristidis,andN.Rish e,"Keyword search on spatial databases, "in Proc.IEEE 24th Int.Conf.Data Eng., 2008, pp.656-665.