

Multi-Dimensional Datasets used for Keyword Search

¹EAGA DEEPIKA ²Dr. PULYALA MAHIPAL REDDY

¹PG Scholar, Department of CSE, Vaagdevi College of Engineering, Bollikunta Warangal, Telangana, Mail id: eaga.deepika@gmail.com.

²Assistant Professor, Department of CSE, Vaagdevi College of Engineering, Bollikunta Warangal, Telangana, Mail id: mahipalreddy.pulyala@gmail.com.

Abstract—Keyword-based search in text-rich multi-dimensional datasets facilitates many novel applications and tools. In this paper, we consider objects that are tagged with keywords and are embedded in a vector space. For these datasets, we study queries that ask for the tightest groups of points satisfying a given set of keywords. We propose a novel method called ProMiSH (Projection and Multi Scale Hashing) that uses random projection and hash-based index structures, and achieves high scalability and speedup. We present an exact and an approximate version of the algorithm. Our experimental results on real and synthetic datasets show that ProMiSH has up to 60 times of speedup over state-of-the-art tree-based techniques.

1 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 achieve 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. Ex: Flickr. 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 have implemented a method of efficient search in multidimensional dataset. This is associated with images as an input. Images are often characterized by a collection of relevant features, and are commonly represented as points in a multidimensional feature space. For example, images are represented using colour feature vectors, and usually have descriptive text information (e.g., tags or keywords) associated with them. We consider multidimensional datasets where each data point has a set of keywords. The presence of keywords in feature space

allows for the development of new tools to query and explore these multidimensional datasets. Our main contributions are summarized as follows.

(1) We propose a novel multiscale index for exact and approximate NKS query processing.

(2) We develop efficient search algorithms that work with the multiscale indexes for fast query processing.

(3) We conduct extensive experimental studies to demonstrate the performance of the proposed techniques.

1. Filename: It is based on image filename.

2. CBIR (Content based image search): Content based image retrieval (CBIR), also known as query by image content (QBIC) and content based visual information retrieval (CBVIR) is the application of computer vision techniques to the image retrieval problem, that is, the problem of searching for digital images in large databases. Content based image retrieval is opposed to traditional concept based approaches (see Concept based image indexing).

3. TBIR (Text based image search): Concept based image indexing, also variably named as “description based” or “text based” image indexing/retrieval, refers to retrieval from

text based indexing of images that may employ keywords, subject headings, captions, or natural language

text. It is opposed to Content based image retrieval. Indexing is a technique used in CBIR.

2. LITERATURE SURVEY

We study nearest keyword set (referred to as NKS) queries on text rich multidimensional datasets. An NKS query is a set of user provided keywords, and the result of the query may include k sets of data points each of which contains all the query keywords and forms one of the

top- k tightest clusters in the multidimensional space. Illustrates an NKS query over a set of two

dimensional data points. Each point is tagged with a set of keywords. For a query the set of points contains all the query keywords and forms the tightest cluster compared with any other set of points covering all the query keywords. Therefore, the set is the International Research result for the query Q . NKS queries are useful for many applications, such as photosharing in social networks, graph pattern search, geolocation search in GIS systems and so on. We present an exact and an approximate version of the algorithm. Our experimental results on real and synthetic datasets show

that the method has more speed up over state of the art tree based techniques. Other related queries include aggregate nearest keyword search in spatial databases, topk preferential

query, topk sites in a spatial data based on their influence on feature points, and optimal location queries. Our work is different from these techniques. First, existing works mainly focus on the type of queries where the coordinates of query points are known. Even though it is possible to make their cost functions same to the cost function in NKS queries, such tuning does not change their techniques. The proposed techniques use location information as an integral part to perform a best first search on the IRTree, and query coordinates play a fundamental role in almost every step of the algorithms to prune the search space. Moreover, these techniques do not provide concrete guidelines on how to enable efficient processing for the type of queries where query coordinates are missing. Second, in multidimensional spaces, it is difficult for users to provide meaningful coordinates, and our work deals with another type of queries where users can only provide keywords as input. Without query coordinates, it is difficult to adapt existing techniques

to our problem. Finding nearest neighbors in large multidimensional data has always been one of the research interests in data mining field. In this paper, we present our continuous research on similarity search problems. Previous work on exploring the meaning of K nearest neighbors from a new perspective in Pan KNN. It redefines the distances between data points and a given query point Q, efficiently and effectively selecting data points which are closest to Q. It can be applied in various data mining fields. A large amount of real data sets have irrelevant or obstacle information which greatly affects the effectiveness and efficiency of finding nearest neighbors for a given query data point. In this paper, we present our approach to solving the similarity search problem in the presence of obstacles. We apply the concept of obstacle points and process the similarity search problems in a different way. This approach can assist to improve the performance of existing data analysis approaches.

The similarity between two data points used to be based on a similarity function such as Euclidean distance which aggregates the difference between each dimension of the two data points in traditional nearest neighbor problems. In those applications, the nearest neighbor problems are solved based

on the distance between the data point and the query point over a fixed set of dimensions (features). However, such approaches only focus on full similarities, i.e.,

the similarity in full data space of the data set. Also early methods suffer from the “curse of dimensionality”. In a high dimensional space the data are usually sparse, and widely used distance metric such as Euclidean distance may not work well as dimensionality goes higher. Recent research [8] shows that in high dimensions nearest neighbor queries become unstable: the difference of the distances of farthest and nearest points to some query point does not increase as fast as the minimum of the two, thus the distance between two data points in high dimensionality is less meaningful. Some approaches are proposed targeting partial similarities. However, they have limitations such as the requirement of the fixed subset of dimensions, or fixed number of dimensions as the input parameter(s) for the algorithms.

3 INDEX STRUCTURE FOR EXACT PROMISH

We start with the index for exact ProMiSH (ProMiSH-E). This index consists of two main components. Inverted Index I_{kp} . The first component is an inverted index referred

to as I_{kp} . In I_{kp} , we treat keywords as keys, and each keyword points to a set of data points that are associated with the keyword. Let D be a set of data points and V be a dictionary that contains all the keywords appearing in D . We build I_{kp} for D as follows. (1) For each $v \in V$, we create a key entry in I_{kp} , and this key entry points to a set of data points $D_v = \{d_j \mid d_j \in D \text{ and } v \in \text{keywords}(d_j)\}$ (i.e., a set includes all data points in D that contain keyword v). (2) We repeat (1) until all the keywords in V are processed. In Fig. 2, an example for I_{kp} is shown in the dashed rectangle at the bottom. Hashtable-Inverted Index Pairs HI. The second component consists of multiple hashtables and inverted indexes referred to as HI. HI is controlled by three parameters: (1) (Index level) L , (2) (Number of random unit vectors) m , and (3) (hashtable size) B . All the three parameters are non-negative integers. Next, we describe how these three parameters control the construction of HI.

In general, HI contains L hashtable-inverted index pairs, characterized by $\{H_s, I_s\}_{s=1}^L$, where H_s and I_s are the s -th hashtable and inverted index, respectively.

Algorithm . SearchInSubset

In: F0: subset of points; Q: query keywords;
q: query size
In: PQ: priority queue of top-k results
1: rk PQ $\frac{1}{2}$ k $\frac{1}{2}$ r /* kth smallest diameter */
2: SL $\frac{1}{2}$ δ v; $\frac{1}{2}$ $_P$: list of lists to store groups per querykeyword
3: for all v \in Q do
4: SL $\frac{1}{2}$ v $_f$ o \in F0 : o is tagged with vg /* form groups */
5: end for
6: /* Pairwise inner joins of the groups*/
7: AL: adjacency list to store distances between points
8: M 0: adjacency list to store count of pairs between groups
9: for all δ vi; vj \in Q such that i q; j q; i < j do
10: for all o \in SL $\frac{1}{2}$ vi $_do$
11: for all o0 \in SL $\frac{1}{2}$ vj $_do$
12: if jjo $_o$ 0jj2 rk then
13: AL $\frac{1}{2}$ o; o0 $_j$ jo $_o$ 0jj2
14: M $\frac{1}{2}$ vi; vj $_M$ $\frac{1}{2}$ vi; vj $_p$ 1
15: end if
16: end for
17: end for
18: end for
19: /* Order groups by a greedy approach */
20: curOrder $\frac{1}{2}$ $_$
21: while Q $\frac{6}{4}$; do
22: δ vi; vj \in removeSmallestEdge(M)
23: if vi \in curOrder then

24: curOrder.append(vi); Q Q n vi
25: end if
26: if vj \in curOrder then
27: curOrder.append(vj); Q Q n vj
28: end if
29: end while
30: sort(SL, curOrder) /* order groups */
31: findCandidates(q, AL, PQ, Idx, SL, curSet, curSetr, rk)

4. PROPOSED SYSTEM

In this paper, we consider multi-dimensional datasets where each data point has a set of keywords. The presence of keywords in feature space allows for the development of new tools to query and explore these multi-dimensional datasets. In this paper, we study nearest keyword set (referred to as NKS) queries on text-rich multi-dimensional datasets. An NKS query is a set of user-provided keywords, and the result of the query may include k sets of data points each of which contains all the query keywords and forms one of the top-k tightest cluster in the multi-dimensional space. We propose ProMiSH (short for Projection and Multi-Scale Hashing) to enable fast processing for NKS queries. In particular, we develop an exact ProMiSH (referred to as ProMiSH-E) that always retrieves the optimal top-k results, and an approximate ProMiSH

(referred to as ProMiSH-A) that is more efficient in terms of time and space, and is able to obtain near-optimal results in practice. ProMiSH-E uses a set of hashtables and inverted indexes to perform a localized search.

5 CONCLUSIONS

In this paper, we proposed solutions to the problem of top-k nearest keyword set search in multi-dimensional datasets. We proposed a novel index called ProMiSH based on random projections and hashing. Based on this index, we developed ProMiSH-E that finds an optimal subset of points and ProMiSH-A that searches near-optimal results with better efficiency. Our empirical results show that ProMiSH is faster than state-of-the-art tree-based techniques, with multiple orders of magnitude performance improvement. Moreover, our techniques scale well with both real and synthetic datasets. Ranking functions. In the future, we plan to explore other scoring schemes for ranking the result sets. In one scheme, we may assign weights to the keywords of a point by using techniques like tf-idf. Then, each group of points can be scored based on distance between points and weights of keywords. Furthermore, the criteria of a result containing all the

keywords can be relaxed to generate results having only a subset of the query keywords

REFERENCES

- [1] W. Li and C. X. Chen, "Efficient data modeling and querying system for multi-dimensional spatial data," in Proc. 16th ACM SIGSPATIAL Int. Conf. Adv. Geographic Inf. Syst., 2008, pp. 58:1–58:4.
- [2] D. Zhang, B. C. Ooi, and A. K. H. Tung, "Locating mapped resources in web 2.0," in Proc. IEEE 26th Int. Conf. Data Eng., 2010, pp. 521–532.
- [3] V. Singh, S. Venkatesha, and A. K. Singh, "Geo-clustering of images with missing geotags," in Proc. IEEE Int. Conf. Granular Comput., 2010, pp. 420–425.
- [4] V. Singh, A. Bhattacharya, and A. K. Singh, "Querying spatial patterns," in Proc. 13th Int. Conf. Extending Database Technol.: Adv. Database Technol., 2010, pp. 418–429.
- [5] J. Bourgain, "On Lipschitz embedding of finite metric spaces in Hilbert space," *Israel J. Math.*, vol. 52, pp. 46–52, 1985.
- [6] H. He and A. K. Singh, "GraphRank: Statistical modeling and mining of significant subgraphs in the feature space," in Proc. 6th Int. Conf. Data Mining, 2006, pp. 885–890.
- [7] X. Cao, G. Cong, C. S. Jensen, and B. C. Ooi, "Collective spatial keyword querying,"

in Proc. ACM SIGMOD Int. Conf. Manage.Data, 2011, pp. 373–384.

[8] C. Long, R. C.-W. Wong, K. Wang, and A. W.-C. Fu, “Collectivespatial keyword queries: A distance owner-driven approach,” in Proc. ACM SIGMOD Int. Conf. Manage. Data, 2013, pp. 689–700.

[9] D. Zhang, Y. M. Chee, A. Mondal, A. K. H. Tung, and M. Kitsuregawa, “Keyword search in spatial databases: Towards searching by document,” in Proc. IEEE 25th Int. Conf. Data Eng., 2009, pp. 688–699.

[10] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni, “Locality sensitive hashing scheme based on p-stable distributions,” in Proc. 20th Annu. Symp. Comput. Geometry, 2004, pp. 253–262.

[11] Y. Zhou, X. Xie, C. Wang, Y. Gong, and W.-Y. Ma, “Hybrid index structures for location-based web search,” in Proc. 14th ACM Int. Conf. Inf. Knowl. Manage., 2005, pp. 155–162.

[12] R. Hariharan, B. Hore, C. Li, and S. Mehrotra, “Processing spatial keyword (SK) queries in geographic information retrieval (GIR) systems,” in Proc. 19th Int. Conf. Sci. Statistical Database Manage., 2007, p. 16.

[13] S. Vaid, C. B. Jones, H. Joho, and M. Sanderson, “Spatio-textual indexing for geographical search on the web,” in Proc.

9th Int. Conf. Adv. Spatial Temporal Databases, 2005, pp. 218–235.

[14] A. Khodaei, C. Shahabi, and C. Li, “Hybrid indexing and seamless ranking of spatial and textual features of web documents,” in Proc. 21st Int. Conf. Database Expert Syst. Appl., 2010, pp. 450–466.

[15] A. Guttman, “R-trees: A dynamic index structure for spatial searching,” in Proc. ACM SIGMOD Int. Conf. Manage. Data, 1984, pp. 47–57.

[16] I. De Felipe, V. Hristidis, and N. Rishe, “Keyword search on spatial databases,” in Proc. IEEE 24th Int. Conf. Data Eng., 2008, pp. 656–665.

[17] G. Cong, C. S. Jensen, and D. Wu, “Efficient retrieval of the top-k most relevant spatial web objects,” Proc. VLDB Endowment, vol. 2, pp. 337–348, 2009.