

# HIGH DIMENSIONAL DATA COMPUTATION USING ZINC EXPERIMENTS

K Venkata Raju<sup>1</sup>, A Vijaya Kumar<sup>2</sup>, MD Mounika<sup>3</sup>, K Sandeep<sup>4</sup>

AssocProfessor, Department of Computer Science & Engineering, K L University, Guntur, Andhra Pradesh, India.<sup>1</sup>

AsstProfessor, Department of Computer Science & Engineering, K L University, Guntur, Andhra Pradesh, India.<sup>2</sup>

Student, Department of Computer Science & Engineering, K L University, Guntur, Andhra Pradesh, India.<sup>3</sup>

Student, Department of Computer Science & Engineering, K L University, Guntur, Andhra Pradesh, India.<sup>4</sup>

## ABSTRACT

Skyline<sup>[1]</sup> is an important operation in many applications to return a set of interesting points from a potentially huge data space. Given a table, the operation finds all tuple's that are not dominated by any other tuple's. It is found that the existing algorithms cannot process skyline on big data efficiently. This paper presents a novel skyline algorithm SSPL<sup>[1]</sup> on big data. SSPL utilizes sorted positional index lists which require low space overhead to reduce I/O cost significantly. We present a new indexing method named ZINC<sup>[2]</sup> (for Z-order indexing with Nested Code) that supports efficient skyline computation for data with both totally and partially ordered attribute domains. By combining the strengths of the Z-order indexing method with a novel nested encoding scheme to represent partial orders, ZINC<sup>[2]</sup> is able to encode partial orders of varying complexity in a concise manner while maintaining a good clustering of the PO domain values. Our experimental results have demonstrated that ZINC<sup>[2]</sup> outperforms the state-of-the-art TSS technique for various settings.

**Index Terms:** ZINC, SDC+, ZB-Tree, Skyline Computation.

## INTRODUCTION

Data mining is one of the important step in KDD process (Knowledge Discovery and Database). It's the process of extracting data from huge data set. Data mining is about processing data and identifying patterns and trends so that you can decide. Data mining principles have been around for many years, but, with the advent of *big data*, it is even more prevalent. Big data is caused the size of the information is very large. It is no longer enough to get relatively simple and straightforward statistics out of the system with large data sets. Sky line is one of the important operation in many applications to return important points from large database. Skyline has attracted extensive attention and many algorithms are proposed. A set of skyline algorithms<sup>[1]</sup>, such as Bitmap, NN, BBS, SUBSKY, and ZBtree, utilize indexes to reduce the explored data space and return skyline results. For providing skyline computation process on each data set traditionally used index-based algorithms utilize the preconstructed data-structures to avoid scanning the entire data set. It preconstructs data-structures with low space overhead. By the data-

structures, the algorithm only involves a small part of table to return the skyline<sup>[4]</sup> results. Index-based algorithms have serious limitations and the used indexes can only be built on a small and selective set of attribute combinations. Nowadays, big data is used commonly in scientific research and business application.

People will expect to get results quickly and they do not want to wait for several hours. For that we propose a novel skyline algorithm on big data, skyline with sorted positional index lists (SSPL), to return skyline results efficiently<sup>[1]</sup>. The algorithm utilizes the preconstructed data structures which require low space overhead to reduce I/O cost significantly. SSPL consists of two phases: obtaining the candidate positional indexes (phase 1) and retrieving the skyline results (phase 2). In phase 1, SSPL first retrieves the sorted positional index lists  $\{L_1; L_2; \dots; L_m\}$  involved by skyline criteria  $\{A_1; A_2; \dots; A_m\}$  in a round-robin fashion. A mathematical analysis is proposed to compute scan depth  $d$  of the lists in phase 1. It is guaranteed that the candidate positional indexes corresponding to the skyline results are contained in the first  $d$  elements in  $\{L_1; L_2; \dots; L_m\}$ . In phase 1, SSPL performs pruning on any candidate positional index retrieved from  $\{L_1; L_2; \dots; L_m\}$  to discard the candidate whose corresponding tuple is not skyline result. This paper proposes general rules and mathematical analysis for pruning operation. Phase 1 ends when there is a candidate positional index seen in all lists of  $\{L_1; L_2; \dots; L_m\}$ . In phase 2, SSPL exploits the obtained candidate positional indexes to compute skyline results by a selective and sequential scan on the table. At first glance, the sorted positional index lists for SSPL are similar to the sorted column files and. However, the most significant idea for SSPL is its pruning operation. Unlike the sorted column files which are used to support sorted retrieval mainly, the sorted positional index lists are the data structures to facilitate pruning and reduce the candidate tuples significantly. Although SSPL is an approximate method to obtain skyline results<sup>[4]</sup>, its probability of correctness is extremely high. The extensive experiments are conducted on two sets of terabyte synthetic data and a set of gigabyte real data, and the experimental results show that compared to the existing algorithms; SSPL involves up to six orders of magnitude fewer tuples, and obtains up to three orders of magnitude speedup. Skyline Sorted Positional Index List algorithm have serious limitations and it fails to process the sequential execution in data sets.

Consider the procedure of SSPL we have to extract the efficient features of big data with computation and other features like data assessment we have to introduce Z-order data set comparison for efficient big data computation.

For example a set of data records  $D$  a skyline query returns the interesting subset of records of  $D$  that are not dominated (with respect to the attributes of  $D$ ) by any records in  $D$ . A data record  $r_1$  is said to dominate another record  $r_2$  if  $r_1$  is at least as good as  $r_2$  on all attributes, and there exists at least one attribute here  $r_1$  is better than  $r_2$ .

There has been a lot of research on the skyline query computation problem, most of which are focused on data attribute domains that are *totally ordered (TO)*, where the best value for a domain is either its maximum or minimum value. However, in many applications, some of the attribute domains are *partially ordered (PO)* such as interval data (e.g. temporal intervals), type hierarchies, and set-valued domains, where two domain values can be incomparable. A number of recent research work has started to address the more general skyline computation problem where the data attributes can include a combination of TO and PO domains.

The first method that proposed for the more general skyline query problem is SDC+, which is an extension of BBS index method for totally ordered domains. SDC+ works an approximate representation of each partially ordered domain by transforming it into two totally ordered domains such that each partially ordered value is presented as an interval value.

A new index method has been proposed for computing skyline queries for TO domains called ZB-tree<sup>[3]</sup>. It has better performance than BBS. It is the extension of B+ -trees, is based on interleaving the bit string representations of attribute values using the Z-order to achieve a good clustering of the data records that facilitates efficient data pruning and minimizes the number of dominance comparisons.

## RELATED WORK

### INDEX BASED ALGORITHM:

Index-based skyline algorithms utilize the preconstructed data-structures to avoid scanning the entire data set.

Tan<sup>[3]</sup> make use of bitmap to compute skyline of a table  $T(A_1; A_2; \dots; A_d)$ . Given a tuple  $x = (x_1; x_2; \dots; x_d) \in T$ ,  $x$  is encoded as a  $b$ -bit bit-vector,  $b = \sum_{i=1}^d k_i$  ( $k_i$  is the cardinality of  $A_i$ ). We assume that  $x_i$  is the  $j_i^{\text{th}}$  smallest value in  $A_i$ , the  $k_i$ -bit bit-vector representing  $x_i$  is set as follows: bit 1 to bit  $j_i - 1$  are set to 0, bit  $j_i$  to bit  $k_i$  are set to 1. The encoded table is stored as bit-transposed files, let  $BS_{ij}$  represent the bit file corresponding to the  $j^{\text{th}}$  bit in the  $i^{\text{th}}$  attribute  $A_i$ . It is given that a tuple  $x = (x_1; x_2; \dots; x_d) \in T$  and  $x_i$  is the  $(j_i)$ th smallest value in  $A_i$ . Let  $A = BS_{1j_1} \& BS_{2j_2} \& \dots \& BS_{dj_d}$  where  $\&$  represents the bitwise and operation. And let  $B = BS_{1(j_1-1)} | BS_{2(j_2-1)} | \dots | BS_{d(j_d-1)}$  where  $|$  represents the bitwise or operation. If there is more than a single one-bit in  $C = A \& B$ ,  $x$  is not a skyline tuple. Otherwise,  $x$  is a skyline tuple.

Kossmann propose NN algorithm to process skyline query. NN utilizes the existing methods for nearest neighbor search to split data space recursively. By a preconstructed R-tree, NN first finds the nearest neighbor to the beginning of the axes. Certainly, the nearest neighbor is a skyline tuple. Next, the data space is partitioned by the nearest neighbor to several subspaces. The

subspaces that are not dominated by the nearest neighbor are inserted into a to-do list. While the to-do list is not empty, NN removes one of the subspaces to perform the same process recursively. During the space partitioning, overlapping of the subspaces will incur duplicates, NN exploits the methods: Laisser-faire, Propagate, Merge and Fine-grained Partitioning<sup>[5]</sup>, to eliminate duplicates.

## THE SSPL ALGORITHM

This section first introduces the data-structures required by SSPL, then describes the overview of the SSPL algorithm, next shows how to perform pruning, followed that presents the implementation and analysis of SSPL, and finally introduces how to extend SSPL to cover other cases.

### Sorted Positional Index List

Given a table  $T$ , the positional index (PI) of  $t \in T$  is  $i$  if  $t$  is the  $i^{\text{th}}$  tuple in  $T$ . We denote by  $T(i)$  the tuple in  $T$  with its  $PI = i$ , and by  $T(i)(j)$  the  $j^{\text{th}}$  attribute of  $T(i)$ . The execution of SSPL requires sorted positional index lists. Given a table  $T(A_1; A_2; \dots; A_M)$ , we maintain a sorted positional index list  $L_j$  for each attribute  $A_j (1 \leq j \leq M)$ .  $L_j$  keeps the positional index information in  $T$  and is arranged in ascending order of  $A_j$ . That is  $\forall i_1, i_2 (1 \leq i_1 < i_2 \leq n)$ ;

The sorted positional index lists are constructed as follows: First, table  $T$  is kept as a set of column files  $CS = \{C_1; C_2; \dots; C_M\}$ . The schema of each column file  $C_j$  is  $C_j(PI; A_j) (1 \leq j \leq M)$ , here  $PI$  represents the positional index of the tuple in  $T$  and  $A_j$  is the corresponding attribute value of  $T(PI)$ . Then, each column file  $C_j$  is sorted in ascending order according to  $A_j$ . Because SSPL only involves  $PI$  field of column files, the  $PI$  values in column files are retained and kept as sorted positional index lists. Here we compare the sorted positional index lists with the indexes used in tree-based algorithms briefly. SSPL constructs a sorted positional index list for each attribute, only  $M$  lists are needed. SSPL reduces the space overhead of data structures from exponential to linear.

More importantly, the processing of SSPL<sup>[1]</sup> can cover all attributes, rather than limited to a small and selective set of attribute combinations in tree-based algorithms. It is noted that read/append-only is an important characteristic of big data, and update is performed in periodic and batch mode. Therefore, sorted positional index lists are worth pre-computing and will be used repeatedly until the next update. And when update operation begins, sorted positional index lists can be updated by merging the corresponding column files in big old data and relatively much smaller new data.

Skyline query processing has involved a lot of research. In this section we will assess the work that handles data with both TO and PO domains.

Importance to ZB-tree, the state-of-the-art approach for data with only TO domains was BBS. The main goal is to map each PO attribute into an approximate representation consisting of a pair of TO attributes. The transformed data is then indexed using BBS. Due to the approximate representation, this approach requires post-processing of false positive skylines. Although this limitation is alleviated with some optimization technique to allow partial progressive skyline computation, the overhead of dominance comparisons can be high.

The another state-of-the-art approach is TSS. This approach is based on BBs. Unlike the BBS approach, TSS uses a precise representation by mapping each PO domain value into an ordinal number with respect to a topological ordering of the PO domain values and a set of interval values.

Another approach is skyline computation<sup>[1]</sup> for continuous data streaming with PO domains. The main goal is efficient skyline maintenance for streaming non-indexed data which is different from indexed based approach for static data.

Another recent approach is dynamic skyline queries which are skyline queries where the user preferences are specified at run time. Data with categorical attributes, the partial orders representing the user's value preferences for such attributes are given as part of the input skyline query.

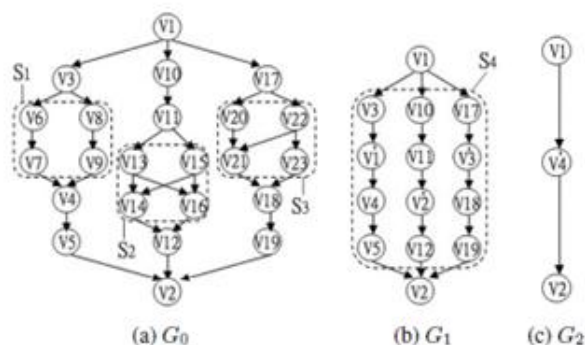
## BACKGROUND WORD

Given a table  $T(A_1; A_2; \dots; A_m), \forall t \in T$ , let us denote by  $t[j]$  the  $j$ th attribute  $A_j$  of  $t$ . Without loss of generality, let a subset of attributes  $AS_{skyline} = \{A_1; A_2; \dots; A_m\}$  be skyline criteria<sup>[1]</sup>, and the dominance relationship between tuples is defined on  $AS_{skyline}$ . For clarity, we assume that min condition only is used for skyline computation. However, the algorithm here can be extended to process any combination of condition (min or max).

**Skyline query.** Given a table  $T$ , skyline query returns a subset  $SKY(T)$  of  $T$ , in which  $\forall t_1 \in SKY(T), \nexists t_2 \in T, t_2 < t_1$ . Given tuple number  $n$  in table  $T$  and size  $m$  of skyline criteria, the expected number  $s$  of skyline results under component independence is known.  $s \approx H_{m,n}$ , here  $H_{m,n}$  is the  $m$ th order harmonic of  $n$ . For any  $n > 0, H_{0,n} = 1$ . For any  $m > 0, H_{m,0} = 0$ . For any  $n > 0$  and  $m > 0, H_{m,n}$  is inductively defined as According to the computation formula of  $H_{m,n}$ , it is found that the number of skyline results does not change significantly as the tuple number increases, while it is very sensitive to the size of skyline criteria. For example, given  $m = 3$ , when  $n$  increases from  $10^5$  to  $10^9$ ,  $s$  changes from 66 to 214. Given  $n = 10^9$ , when  $m$  increases from 2 to 5,  $s$  changes from 20 to 7,684. Although the absolute number of skyline results is large, its proportion among all tuples is rather small. For example, given  $m = 5$  and  $n = 10^9, s/n = 7.684 \times 10^{-6}$ . Given tuple number  $n$  in table  $T$  and size  $m$  of skyline criteria, the expected number  $s$  of



skyline results under component independence is known.  $s = H_{m-1}; n$ , here  $H_m; n$  is the  $m$ th order harmonic of  $n$ . For any  $n > 0, H_0; n = 1$ .



**Figure 1: Partial order reduction process generation in Z-order datasets.**

We represent a partial order by a directed graph  $G = (V; E)$ , where  $V$  and  $E$  denote, respectively, the set of vertices and edges in  $G$  such that given  $v; v_0 \in V$ ,  $v$  dominates  $v_0$  if there is a directed path in  $G$  from  $v$  to  $v_0$ . Given a node  $v \in V$ , we use  $parent(v)$  (resp.,  $child(v)$ ) to denote the set of parent (resp., child) nodes of  $v$  in  $G$ . A node  $v$  in  $G$  is classified as a minimal node if  $parent(v) = \emptyset$ ; and it is classified as a maximal node if  $child(v) = \emptyset$ . We use  $min(G)$  and  $max(G)$  to denote, respectively, the set of minimal nodes and maximal nodes of  $G$ .

Given a partial order  $G_0$ , the key idea behind nested encodings is to view  $G_0$  as being organized into nested layers of partial orders, denoted by  $G_0 \rightarrow G_1 \rightarrow \dots \rightarrow G_{n-1} \rightarrow G_n$ , where each  $G_i$  is nested within a simpler partial order  $G_{i+1}$ , with the last partial order  $G_n$  being a total order. As an example, consider the partial order  $G_0$  shown in Fig. 2, where  $G_0$  can be viewed as being nested within the partial order  $G_1$  which is derived from  $G_0$  by replacing three subsets of nodes  $S_1 = \{v_6; v_7; v_8; v_9\}$ ,  $S_2 = \{v_{13}; v_{14}; v_{15}; v_{16}\}$  and  $S_3 = \{v_{20}; v_{21}; v_{22}; v_{23}\}$  in  $G_0$  by three new nodes  $v_{01}$ ,  $v_{02}$  and  $v_{03}$ , respectively, in  $G_1$ .  $G_1$  in turn can be viewed as being nested within the total order  $G_2$  which is derived from  $G_1$  by replacing the subset of nodes  $S_4 = \{v_3; v_{01}; v_4; v_5; v_{11}; v_{02}; v_{12}; v_{17}; v_{03}; v_{18}; v_{19}\}$  by one new node  $v_{04}$  in  $G_2$ . We refer to the new nodes  $v_{01}$ ,  $v_{02}$ ,  $v_{03}$  and  $v_{04}$  as *virtual nodes*; and each virtual node  $v_{0j}$  in  $G_{i+1}$  is said to *contain* each of the nodes in  $S_j$  that  $v_{0j}$  replaces. By viewing  $G_0$  in this way, each node in  $G_0$  can be encoded as a sequence of encodings based on the nested node containments within virtual nodes.

## PARTIAL ORDER REDUCTION ALGORITHM

Given an input partial order  $G_i$ , algorithm PO-Reduce operates as follows: Let  $S = \{S_1; \dots; S_k\}$  be the collection of regular regions in  $G_i$ ; (2) If  $S$  is empty, then let  $S = \{S_1\}$ , where  $S_1$  is an

irregular region in  $G_i$  that has the smallest size (in terms of the number of nodes) among all the irregular regions in  $G_i$ . (3) Create a new partial order  $G_{i+1}$  from  $G_i$  as follows. First, initialize  $G_{i+1}$  to be  $G_i$ . For each region  $S_j$  in  $S$ , replace  $S_j$  in  $G_{i+1}$  with a virtual node  $v_j$  such that  $parent(v_j) = parent(v)$  with  $v \in \min(S_j)$  and  $child(v_j) = child(v)$  with  $v \in \max(S_j)$ . (4) If  $G_{i+1}$  is a total order, then the algorithm terminates; otherwise, invoke the PO-Reduce algorithm with  $G_{i+1}$  as input.

When a node  $v$  in a region  $R$  is being replaced by a virtual node  $v_0$ , we say that  $v$  is *contained in*  $v_0$  (or  $v_0$  *contains*  $v$ ), denoted by  $v \in R \rightarrow v_0$ . Clearly, the node containment can be nested; for example, if  $v$  is contained in  $v_0$ , and  $v_0$  is in turn contained in  $v_{00}$ , then  $v$  is also contained in  $v_{00}$ . Given an input partial order  $G_0$ , we define the *depth* of a node  $v$  in  $G_0$  to be the number of virtual nodes that contain  $v$  in the reduction sequence computed by algorithm PO-Reduce. As an example, consider the value  $v_6$  in Fig. 2 and let  $R_0 = \{v_6; v_7; v_8; v_9\}$  and  $R_1 = \{v_3; v_1; v_4; v_5; v_{10}; v_{11}; v_{12}; v_{17}; v_{18}; v_{19}\}$ .

Thus, given an input partial order  $G_0$ , algorithm PO-Reduce outputs<sup>[10]</sup> the following:

(1) the partial order reduction sequence,  $G_0 \rightarrow G_1 \rightarrow \dots \rightarrow G_n$ , where  $G_n$  is a total order; and (2) the node containment sequence for each node in  $G_0$ . If a node  $v_0$  in  $G_0$  has a depth of  $k$ , we can represent the node containment sequence for  $v_0$  by  $v_0 \in R_0 \rightarrow v_1 \in R_1 \rightarrow \dots \rightarrow v_k \in R_k$ , where each  $v_i$  is contained in the region  $R_i$ ,  $i \in [0; k]$ .

## PERFORMANCE ANALYSIS

To evaluate the performance of our proposed ZINC, we conducted an extensive set of experiments to compare ZINC<sup>[2]</sup> against three competing methods: TSS and the two basic extensions of ZB-tree, namely, TSS+ZB and CHE+ZB. Our experimental results show that ZINC<sup>[2]</sup> outperforms the other three competing methods. Given that both TSS+ZB and CHE+ZB are also based on ZB-tree, the superior performance of ZINC demonstrates the effectiveness of our proposed NE encoding for PO domains.

**Algorithms:** We consider two variants of the main competing method, TSS: an unoptimized variant of TSS (denoted by TSS) and an optimized variant of TSS (denoted by TSS-opt). In TSS, the set of intervals associated with each data / index entry's PO value are stored explicitly with the entry, while in TSS-opt, the intervals associated with an entry are retrieved from a separate precomputed structure.

To compare the effectiveness of our proposed nested encoding scheme, we also introduced two variants of ZB-tree that are based on using different schemes to encode PO domains. The first variant, TSS+ZB, combines the TSS encoding scheme with the ZB-tree method. Each PO domain value  $v$  of a data point is encoded into a bitstring based on its ordinal value  $v$  in a topological sorting of the PO domain values. The inclusion of  $v$  in the derivation of the data

point's Z-address is important to ensure ZB-tree's monotonicity property. Each leaf node entry in TSS+ZB stores a data point  $p$  together with the interval set representation of each of  $p$ 's PO attribute values. In each internal node entry of TSS+ZB, besides storing the  $minpt$  and  $maxpt$  of the corresponding RZ-region (similar to what is done in ZB-tree), for each PO attribute  $A$ , a merged interval set for  $A$  is also stored which is the union of the interval sets for attribute  $A$  of the covered data points. In TSS+ZB, region-based dominance test is applied as follows: if (1) the Z-address of an intermediate skyline point  $pi$  dominates  $minpt$  of an internal node entry  $ej$ , and (2) the interval set of  $pi$  subsumes the interval set of  $ej$  w.r.t. every PO dimension, then the region represented by  $ej$  is dominated by  $pi$  and is pruned from consideration.

**Synthetic datasets:** We generated three types of synthetic data-sets according to the methodology in. For TO domains, we used the same data generator as [8] to generate synthetic datasets with different distributions. For PO domains, we generated DAGs by varying three parameters to control their size and complexity: *height* ( $h$ ), *node density* ( $nd$ ), and *edge density* ( $ed$ ), where  $h \in \mathbb{Z}^+$ ,  $nd, ed \in [0; 1]$ . Each value of a PO domain corresponds to a node in DAG and the dominating relationship between two values is determined by the existence of a directed path between them. Given  $h$ ,  $nd$ , and  $ed$ , a DAG is generated as follows. First, a DAG is constructed to represent a poset for the powerset of a set of  $h$  elements ordered by subset containment; thus, the DAG has  $2^h$  nodes. Next,  $(1 - nd) * 100\%$  of the nodes (along with incident edges) are randomly removed from the DAG, followed by randomly removing  $(1 - ed) * 100\%$  of the remaining edges such that the resultant DAG is a single connected component with a height of  $h$ . Following the approach in [8], all the PO domains for a dataset are based on the same DAG. Table 2 shows the parameters and their values used for generating the synthetic datasets, where the first value shown for each parameter is its default value. In this section, default parameter values are used unless stated otherwise.

**Real dataset:** We used a real dataset on movie ratings that is derived from two data sources, *Netflix* and *MovieLens*. Netflix contains more than 100 million movie ratings submitted by more than 480 thousand users on 17770 movies during the period from 1999 to 2005. MovieLens contains more than 1 million ratings submitted by more than 6040 users on 3900 movies. Both these data sources use the same rating scale from 0 to 5 with a higher rating value indicating a more preferred movie. Our dataset consists of the ratings for 3098 of the movies that are common to both data sources.

We derived a PO attribute, named *movie preference*, for the 3098 movies as follows: a movie  $mi$  dominates another movie  $mj$  iff the average rating of  $mi$  in one data source is higher than that of  $mj$ , and the average rating of  $mi$  in the other data source is at least as high as that of  $mj$ . We also derived two TO attributes for each movie, named *average rating* and *number of ratings*, which



represent, respectively, the movie's average rating (each value is between 0.00 and 5.00) and the total number of ratings that it has received over the two data sources. The number of distinct values for these two TO domains are 501 and 219800, respectively. For each of the TO domains, a higher attribute value is preferred.

## CONCLUSION

This paper presents a novel skyline algorithm SSPL on big data. SSPL utilizes sorted positional index lists which require low space overhead to reduce I/O cost significantly. We present a new indexing method named ZINC (for Z-order indexing with Nested Code) that supports efficient skyline computation for data with both totally and partially ordered attribute domains. By combining the strengths of the Z-order indexing method with a novel nested encoding scheme to represent partial orders, ZINC is able to encode partial orders of varying complexity in a concise manner while maintaining a good clustering of the PO domain values. Our experimental results have demonstrated that ZINC<sup>[2]</sup> outperforms the state-of-the-art TSS technique for various settings.

## REFERENCES

- [1] Xixian Han, Jianzhong Li, "Efficient Skyline Computation on Big Data", IEEE Transactions On Knowledge And Data Engineering, Vol. 25, No. 11, November 2013.
- [2] Bin Liu Chee Yong Chan, "ZINC: Efficient Indexing for Skyline Computation", The 37th International Conference on Very Large Data Bases, August 29th September 3rd 2011, Seattle, Washington. *Proceedings of the VLDB Endowment*, Vol. 4, No. 3 Copyright 2010 VLDB Endowment 21508097/10/12... \$ 10.00.
- [3] C.-Y. Chan, H.V. Jagadish, K.-L. Tan, A.K.H. Tung, and Z. Zhang, "Finding K-Dominant Skylines in High Dimensional Space," Proc. ACM SIGMOD Int'l Conf. Management of Data (SIGMOD '06), pp. 503-514, 2006.
- [4] L. Chen and X. Lian, "Efficient Processing of Metric Skyline Queries," IEEE Trans. Knowledge Data Eng., vol. 21, no. 3, pp. 351-365, Mar. 2009.
- [5] M. Gibas, G. Canahuat, and H. Ferhatosmanoglu, "Online Index Recommendations for High-Dimensional Databases Using Query Workloads," IEEE Trans. Knowledge and Data Eng., vol. 20, no. 2, pp. 246-260, Feb. 2008.
- [6] P. Godfrey, "Skyline Cardinality for Relational Processing," Foundations of Information and Knowledge Systems, vol. 2942, pp. 78-97, Springer Berlin/Heidelberg, 2004.
- [7] P. Godfrey, R. Shipley, and J. Gryz, "Algorithms and Analyses for Maximal Vector Computation," The VLDB J., vol. 16, no. 1, pp. 5-28, 2007.
- [8] J. Gray and P.J. Shenoy, "Rules of Thumb in Data Engineering," Proc. 16th Int'l Conf. Data Eng. (ICDE '00), pp. 3-12, 2000.

- [9] K. Hose and A. Vlachou, "A Survey of Skyline Processing in Highly Distributed Environments," *The VLDB J.*, vol. 21, no. 3, pp. 359-384, 2012.
- [10] Y. Fang and C. Y. Chan. Efficient skyline maintenance for streaming data with partially-ordered domains. In *DASFAA*, pages 322–336, 2010.

IJAER