

Processing Skyline Queries Based on huge and incomplete Datasets in Various Database

Mercy Bhikshavathi Kalluri*¹, Ravindra Reddy E²

*¹M.Tech Scholar, St. Mary's Group of Institutions, Guntur, Andhra Pradesh, India

²Department of CSE, Assistant Professor, St. Mary's Group of Institutions, Guntur, Andhra Pradesh, India

ABSTRACT

Personalized recommendation and therefore the process of time period data exemplify the process of huge information that within the field of Internet-of-Things (IoT) received a good extent of attention in recent literature. The integrity of huge information within the IoT is widespread. getting customized info from the incomplete information set continues to be nonplussed by looking out efficient and correct strategies at the present. Skyline query may be a wide used processing technique, particularly within the field of multi-objective call analysis and information mental image. To eliminate the negative effects on huge processing in IoT, a unique skyline preference question strategy supported huge and therefore the incomplete information set is projected during this paper. This strategy merely separates and divides huge and incomplete information set into 2 elements consistent with dimension importance and executes skyline question, severally. The strategy chiefly resolves the matter of extracting personalized info from huge and incomplete information set and improves the efficiency of skyline question on huge and incomplete information set. First, this paper presents a skyline preference question strategy supported strict bunch and implements it on dimensions that have higher importance. Second, a skyline preference query strategy supported loose bunch is enforced on dimensions that have lower importance. Finally, integrating native skyline question results, this paper calculates world skyline question results by victimization info entropy theory. The efficiency and effectiveness of Skyline Preference question (SPQ) algorithmic rule are evaluated in terms of latent period and result set size through the comparative experiments with I Skyline algorithm and sort-based incomplete information skyline algorithmic rule. an outsized variety of simulation results show that the efficiency of SPQ algorithmic rule is more than that of different common strategies.

Keywords: Preference Queries, Skylines, Skyline Queries, Algorithms, Incomplete Data, Database, information entropy.

I. INTRODUCTION

Prior to the first introduction of the skyline operator into database community by Borzsony et al. (2001), there was a problem called maximum vector problem or Pareto optimum that has been addressed by (Bentley et al., 1978). Skyline queries attempt to return a set of data items S from a database, where S contains only those data items whose dimension values are not dominated by any other data item exist

in the database. For instance, suppose there are two data items i and j comprise of a set of dimensions (attributes) belong to a database, D . Assume the data item, i belongs to S and D , while j belongs to D only. The data item i can be in the skyline set, S if and only if i has better values than j in all dimensions or at least it has value not worse than j in one dimension. The following restaurant database example in Figure 1 clarifies the concept of skyline queries in database systems. Assume the database table contains the

details of 10 restaurants, in which each restaurant consists of two dimensions. The first dimension represents the rating of each restaurant given by customers while the second dimension indicates the price of the food served in each restaurant. Figure 1(B) illustrates the representation of the restaurant database example in 2-Dimensions space. Assume a user is looking for the best restaurant in the city where the price is the cheapest and the restaurant rate is the highest. Based on skyline definition and from the database example, it can be seen that restaurants r8, r9, r3, r1 and r6 are partially dominated by r5 or r2 either based on the first dimension (price) or the second dimension (rate). While restaurants r4 and r10 are partially dominated by restaurant r7 only. It can be noticed for restaurant r8 and r5 their price values are cheapest among all other restaurants but r8 do not have the better rating than any of the other restaurants. Therefore, r8 is dominated by r5. Similarly, restaurant r7 has higher price value compare with r1 and r9; however, it has better rating value compare with any other restaurant. Hence, restaurant r7 is not dominated by any other restaurant in the restaurant database. As demonstrated in Figure 1(A), the result of applying skyline technique on restaurant database example will retrieve r5, r2 and r7 as the skylines of the restaurant database.

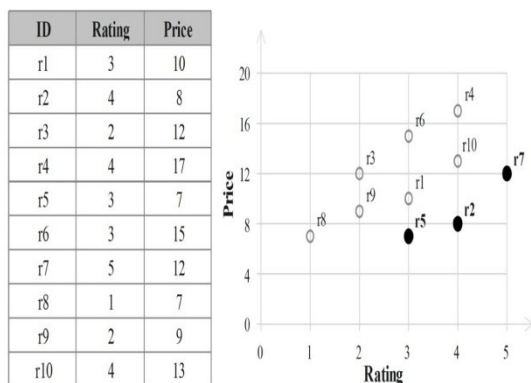


Figure 1. Example of skyline query Incomplete Dataset

A database is said to be incomplete if and only if values of some dimensions are not present (missing). Otherwise, it is complete.

With these complicated challenges, it is impractical to directly apply skyline approaches designed for the complete database on a database with partial incomplete data due to the prohibitive cost and the exhaustive pair wise comparisons that need to be conducted between data items. Most importantly, it is highly unlikely to obtain the correct set of skylines from the database. This is because a skyline process might suffer from the issue of cyclic dominance and losing the transitivity property of skyline which results into incorrect results or in the worst case might return an empty set of skylines. The following running example clarifies further these challenges. For instance, we assume that a tourist is looking to rent a car somewhere in a city with lower rent cost per hour, less consumption and a good rating given by other users (tourists). A car renting database consists of three data items with three dimensions, namely: $d_1 = rph$, $d_2 = cpl$ and $d_3 = r$, where rph represents the rent price per hour, while cpl indicates the consumption rate per liter and lastly r denotes the rating value of each car by users. We also assume that there are some data items with missing values, which means that some dimension values are not present. Thus, the data items of the database are as follows. $c_1 (5, 17, -)$ $c_2 (6, -, 5)$ $c_3(-, 20, 3)$. The symbol $(-)$ denotes that the value is missing in that dimension of the data item. Based on the common dimensions with non-missing values it can be concluded that c_1 dominates c_2 as c_1 is better than c_2 (lower is better) at first dimension rph , while c_2 dominates c_3 as c_2 had a better rating than c_3 (third dimension, r). However, comparing c_1 with c_3 demonstrate that c_3 dominates c_1 in the second dimension, cpl . Hence, car c_1 does not dominate c_3 which therefore means the dominance relation is not transitive. Furthermore, car c_3 is dominating c_1 , which means that the dominance relation is cyclic. From this example, all three cars are being

dominated by one another and therefore the skyline process fails to decide the best car for the tourist.

In this study, we are focusing on examining and discussing the previous approaches proposed for processing skyline queries with incomplete data in three different databases, namely: centralized, distributed and cloud databases. These approaches have been analyzed thoroughly describing the strengths and the weakness of these approaches. Furthermore, we also discuss the current research trends related to skyline queries in some other environments and highlight the most issues concern skyline queries in these contexts.

II. LITERATURE REVIEW

Lian et. al., 2011[1] proposed a probabilistic least influenced set query techniques over uncertain objects in databases. The high probability query objects do not affect the retrieval of uncertain object. The technique is improved with the influence of monochromatic and dichromatic applications that answers well the given queries. Minmax, Candidate pruning and Regional pruning methods effectively reduces the search space by integrating the query procedures. The method is tested over both static and moving objects over business planning datasets.

Lin and Hong, 2012[2] proposed a frequent mining datasets in transactional databases. The technique is proposed to find the availability of data in uncertain dataset using probability values 0 or 1. The processing in uncertain database is carried out with UF-tree structure that involves uncertain frequent pattern tree. The same item sets are merged in tree branch even if the transaction probabilities are different. The technique with uncertain frequent pattern tree seems better computation UF – tree.

Azeem et. al., 2015[3] designed a framework over uncertain query ranks and dataset distribution threshold tuned and ceiling-tuned distributed top-k algorithm for finding the top-k tuples. Efficient ranks and reduced communication rounds are achieved

while finding the tuples with efficient computation over distributed database. The threshold tuned algorithm proved to be effective than the other one in terms of finding top-k tuples.

Ahmed et. al., 2016[4] combined Weighted Uncertain Interesting Pattern tree structure with additional measures like wUConf and uConf to find the frequent item set. This technique helps in mining the uncertain dataset using correlated patterns. Firstly, the weights of individual measures are calculated using correlated patterns and secondly using the measures, the performance of mining is calculated. The method helps in finding the valued patterns faster than conventional approaches.

Lin et. al., 2016[5] used potential high-utility item set mining technique to discover the frequent patterns in the dataset. The itemset has high utility than others and high existence probabilities are founded out with tuples uncertainty model. An upper bound mining algorithm uses level wise search and repeated scanning is avoided using PU list algorithm. The scalability of the system is improved with mining datasets and not generating the candidates over real life and synthetic datasets. The algorithms is tested with total patterns, runtime and memory usage with two algorithms.

Lian and Chen, 2016 [6] used probabilistic top-k dominating query technique to identify the k-uncertain objects in the database. The search speech using this technique is avoided with pruning techniques. Furthermore, arbitrary subspaces are used to compute the probabilistic top-k query and testing is done to prove its effectiveness in query retrieval.

The above techniques proved efficient in terms of uncertain database with uncertain queries but the usage of skyline queries in such database brings more robustness in algorithm computations. The integration of above techniques with novel approaches is applied over uncertain database to find the skyline queries are discussed in the following section.

Skyline Approaches

There are severe drawbacks in applying skyline queries over user preferences that leads to impractical skyline or non-dominating results from skyline queries. To avoid such congenital problems in skyline querying various researches have spotted the problems and eradicated it with special solutions. Few of which are presented here.

To handle impractical skyline size, Lee et. al., 2012[7] proposed interactive preference elicitation framework that collects user preference and updates skyline iteratively. This technique reduces the user interaction and increase the size of skyline reduction. The query formulation remains intuitive and efficient over several experimental testing.

For reverse skyline operations, Gao et. al., 2012 [8] precompiled skyline behavior using reverse skyline using skyline approximation technique that provides a node with multiple accessing. This technique is further improved with constrained reverse skyline query for constrained specific region. Full reuse and global skyline based reverse skyline algorithms are designed and tested over real and synthetic datasets.

The continuous monitoring on skyline dynamic data, Lu et. al., 2013 [9] proposed an approach with 11 algorithms for each stage processing. The main aim is to distribute the centralized load and thus reducing the bandwidth consumption. The experimental evaluation on wide area sensor networks is tested with both real and synthetic datasets.

The sliding technique with probabilistic threshold is processed over uncertain datasets by [10] Zhanget al., 2013 [10]. The technique initially computes the element type and then the skyline size is determined to identify the query composition. The probabilistic query technique is applied over continuous high-speed datasets to identify the top tuples with multiple probability thresholds.

Lian and Chan, 2013[11] proposed probabilistic group subspace skyline technique to reduce uncertainty in sensor data analysis. The technique

looks similar to [4] in terms of its faster process and pruning reduces the cost and search space [1].

Chung et. al., 2013 [12] proposed combinatorial skyline query technique to identify the perfect combinations. Past skyline query is added to it for enhancing the prediction capability of skyline queries. The brute force and decomposition algorithms are used in improving the decision-making capability, business plan, market analysis and economic research.

Lee et. al., 2013 [13] proposed Fast algorithm that processes static data on data streams. The technique is designed to handle multiple queries with filtering, discriminate analysis for discarding and determining objects, respectively.

Similar to [1], jiang et. al., 2014 [14] used monochromatic and bichromatic applications over mutual skyline query that retrieves the skyline object in dynamic environment. Reverse skyline query is used to find the query object in task allocation, market analysis and personalized matching.

To study the incomplete data in skyline environment due to privacy preservation and device failure, Gao et. al., 2014[15] proposed k-sky band query technique. Baseline and virtual point algorithm are designed to process the missing data. Concepts like shadow, expired skyline and thickness warehouse boosts the performance, and further tackles constrained and group-by skyline queries in real and synthetic data sets.

With reduced computation overhead, Yong et. al., 2015[16] identified top-k tuples using BFS and peeling algorithm with R-tree. This algorithm obtains the tuples with the concept of dependency and independency. The system is tested over synthetic datasets to measure dimensionality, cardinality, retrieval size and distribution.

To eliminate unnecessary comparison, Jongwuk Lee et. al., 2015[17] considers a cost model and developed pivot points to reduce comparisons. An algorithm is designed that combines the pivot model

with sorting and partition operation. Skytree technique is considered to include point insertion and deletion using recursive point space partitioning. Greedy algorithm is further added with sky tree to perform faster than conventional techniques.

Similar to [11], saad et. al., 2015[18] proposed SkyQUD framework that computes skyline points using probabilistic technique in uncertain dataset. Harvesting is the initial phase performing elimination operation and strict selection is the final phase that extracts the dominating point.

Kulkarni and Momin, 2015 [19] optimized response time in static dataset using QPSkyline and theQPUpdateSkyline algorithms. This identifies the nearest frequent dataset and updates it in the profiler. The former algorithms is applicable over static datasets and latter on frequent real life datasets.

Jiang et. al., 2015[20]evolved a technique that selects the k-combination points with monotonic preference variable. This is done with top-k combinatorial metric skyline query technique with Incremental Combination Sorting and Top-k Combinatorial Metric Skyline Algorithms. The algorithms are designed to provide better skyline points in location based applications possessing real and synthetic datasets. Ying Zhang et. al., 2015[21] used random algorithm with ϵ -approximation guarantee to find the Top-k points in the skyline datasets. The p-skyline points are computed using p-exact and p-rand algorithms in real and synthetic datasets.

Li et. al., 2015 [22] used GDPS approach to improve the global knowledge in uncertain environment. This uses optimal feedback that optimizes the query using grid technology. Moreover, use of tuples-k selection, localized and server pruning over real and synthetic datasets improvised the correctness of datasets. hen and Lee, 2015[23] computed s-neighbor skyline query technique to find the non-dominating points. Rs-N skyline points are computed using BBS algorithm with new index tree procedure. This algorithm works well with s-neighbor skyline technique over synthetic datasets.

Park et. al., 2015 [24] used Map Reduce algorithm to compute continuous and discrete datasets in big data environment. This technique uses two algorithms (PS-QP-MR, PS-BR-MR and PS-BRF-MR) with initial pruning procedure or filtering approaches to clean the irrelevant data in synthetic datasets.

In [8], reverse skyline operation is performed for skyline query and Gao et. al., 2015[25] proposed similar reversek-skyband (RkSB) query technique with five various algorithms. This effectively returns the query skyline points with initial computation, reusing and pruning procedure. Further, ranked reverse skyline technique is applied over real and synthetic datasets to effectively minimize the score with pruning heuristics.

Gao et. al., 2015[26] proposed most desirable skyline object query technique returns the top dominating points using weights calculation and ranking procedure. Further, three algorithms and constrained procedure is added with the technique that improved the search results in real and synthetic datasets. The top dominating points in uncertain time series datasets are computed by He et. al., 2016 [27]. It is a dual step technique with three pruning technique that obtains skyline points on interval basis. Another two methods were designed to compute the probabilistic skyline in dataset and pruning is improved further with adjacent interval technique.

Lee et. al., 2016[28] proposed bucket skyline sort with bucket orders and point level orders for pair wise comparison in incomplete dataset. If no skyline point occurs, parameters with monotonicity properties are applied over synthetic and real life dataset to obtain accurate points.

Le et. al., 2016[29]defined returning skyline points using probabilistic skyline tuples [27] and then finding the tuples with all probable worlds. The probabilistic skyline tuples are found using best pro-skyline query technique with pruning having user defined threshold. Formulas were used to calculate the probabilistic skylines and pruning over search space in real life datasets.

III. RELATED WORK

This section provides a set of necessary definitions and notations that facilitate the process of understanding the concept of skyline queries in incomplete databases. This includes definitions related to the dominance theory of the skyline, the skyline queries, incomplete database, comparable, distributed databases and dynamic databases. These terms are further explained as follows.

A relation of the database D is denoted by $R(d_1, d_2, \dots, d_m)$ where R is the name of the relation with m -arity and $d = (d_1, d_2, \dots, d_m)$ is the set of dimensions.

Definition 1 Dominance

Given two data items p_i and $p_j \in D$ database with d dimensions, p_i dominates p_j (the greater is better) (denoted by $p_i > p_j$) if and only if the following condition holds: $\forall d_k \in d, p_i.d_k \geq p_j.d_k, \wedge \exists d_l \in d, p_i.d_l > p_j.d_l$.

Definition 2 Skyline Queries

Select a data item p_i from the set of D database if and only if p_i is as good as p_j (where $i \neq j$) in all dimensions (attributes) and strictly better than p_j in at least one dimension (attribute). We use S_{skyline} to denote the set of skyline data items, $S_{\text{skyline}} = \{p_i \mid \forall p_j \in D, p_i > p_j\}$.

Definition 3 Incomplete Database

Given a database $D(R_1, R_2, \dots, R_n)$, where R_i is a relation denoted by $R_i(d_1, d_2, \dots, d_m)$, D is said to be incomplete if and only if it contains at least a data item p_j with missing values in one or more dimensions d_k (attributes); otherwise, it is complete.

Definition 4 Comparable

Let the data items a_i and $a_j \in R$, a_i and a_j are comparable (denoted by $a_i \varepsilon a_j$) if and only if they have no missing values in at least one identical dimension; otherwise a_i is incomparable to a_j (denoted by $a_i \not\varepsilon a_j$).

Definition 5 Dynamic Database

A database D is said to be dynamic if its contents are persistently change through many update operations such as DELETE, UPDATE or INSERT.

In dynamic databases processing skyline is quite challenging due to the continuous update of the data in the database. This frequent change might impact negatively on the validity of the skyline and force users to re-compute the skyline process whenever the data is updated. Re-evaluating the skylines by accessing the whole database is an exhaustive process that incurs longer processing time and higher overhead for the skyline process. Therefore, this approach should be avoided and provide a solution that partially scans the database when there is an update on the contents of the database.

IV. PROPOSAL WORK

Skyline Preference Query Strategy

A. Skyline Preference Query Strategy Based On Strict Clustering

The skyline preference query strategy based on strict clustering performs in two steps. Firstly, this strategy implements strict clustering on tuples from IS' and then executes skyline preference query algorithm based on attribute value ordering on tuples that aren't dominated by others.

The processing of skyline preference query algorithm based on attribute value ordering can be divided into four phases.

Phases1: Sorting attributes in non-ascending order makes tuples which more likely dominate others to be processed prior. After sorting, each dimension will generate an array $D_i, I \in [1;]$ and $D_i[j] = D_i[j \in C 1]; j \in [one; IS)$, IS_0 is the number of tuples in IS_0 . If a tuple has missing value at i th attribute, it will not be added to the array D_i . To save storage space, array D_i doesn't store complete tuple but its id. Setting a pointer ptr_i points to array D_i and all the tuples that not be dominated will be added to $Candidate_Set$. This method selects an array D_i randomly and then process the tuple pointed by pointer ptr_i . Each tuple

in Candidate_Set has two values, one is the number of being processed that is denoted as processedCount and the other is the number of character '1' in a tuple encoding the number of non-missing attributes is denoted as dimCount.

Phases 2: For the currently selected tuple, there are three cases. Case 1: If tuple p^0 has never been processed and still be in the Candidate_Set, it can be compared with the tuple p^0_j that there is no comparison between them even if p^0_j has been dominated by tuple processed before. This is because p^0_j maybe dominates p^0 I If there exists a tuple can dominate p^0 , tuple p^0 should be removed from Candidate_Set. Case 2: If tuple p^0 has never been processed but is dominated by tuple processed before that p^0 will not be in Candidate_Set anymore. Therefore, tuple p^0 only compares with the tuples which still be in Candidate_Set and there is no comparison between them. In case 1 and case 2, any tuple is dominated by tuple p^0 will be removed from Candidate_Set. Case 3: If a tuple has been processed before, it will not compare to any tuple.

Special case: If there are two tuples, p^0_i and p^0_j , p^0_i dominates p^0_j only in few dimensions (less than the half number of dimensions), then p^0_i will dominate p^0_j regarding traditional skyline rule. However, considering user preference, the miss-ing dimensions in tuples might have higher weights than those available dimensions. We should take those missing dimen-sions into consideration to ensure the accuracy of dominance relation between incomplete tuples. Judging the dominance relation of two tuples only through comparable dimensions is inaccurate. Especially, this paper de nes the weak dominance relation among tuples with missing dimensions.

Phase 3: When the number of comparisons of tuple p equals to p dimCount, which is the number of non-missing dimensions of p . Then p will be removed from the Candi-date_set to the strict skyline result set SSRS.

Phase 4: When the candidate set is empty or all tuples are processed at least once, the rest of tuples in Candidate_Set are push back to the strict skyline clustering result set SSRS.

Algorithm 1 Dominate (p,q)

1. Input: tuple p,q
2. Output: Integer value
3. **If** $p > q$
4. **Return** 1;
5. **Else if** $p > q$
6. **If** p weight q weight
7. **Return** 1;
8. **Else if** $q > p$
9. **If** q weight p weight
10. **Return** 1;
11. **Return** 0

this is the algorithm to judge the dominance relation between two tuples.

The pseudo code of skyline preference query algorithm based on attribute value ordering is as follows:

We also describe the skyline preference query algorithm based on attribute value ordering by taking an example of the incomplete data set from Table 1. The result of IS⁰ that is sorted by attribute value is shown in Table 2. p_4 dominates p_1 , p_2 weak dominates p_4 . However, p_2 weight p_4 weight, so p_2 could not dominate p_4 . There is no dominance relation between p_4 and p_5 . Similarly, there is also no dominance relation between p_4 , p_6 and p_7 . After the rst iteration, p_4 processedCount D 1. p_1 and p_5 are removed from Candidate_Set. Then, ptr₁ also pointers to the rst tuple p_4

Table 1. The result of sorted dimensions.

D1	D2	D3	D4	D5
P4	P4	P7	P6	P3
P7	P1	P2	P2	P6
P1	P7	P4	P3	P1
P6	P3	P5	P1	P4

of array D2. Since p_4 has been processed before, it need not compare with any tuple. In the end of the second iteration, p_4 processedCount is 2. In the third iteration, the tuple of array D3, p_7 is processed and it should compare with p_1 , p_3 , p_5 and p_6 . Actually, p_7 is not dominated by p_1 and p_5 which have been removed from the Candidate_Set at the first iteration. p_7 weak dominates p_2 and p_7 weight $>$ p_2 weight, so p_7 dominates p_2 . Similarly, p_7 dominates p_3 and p_6 . By the end of this iteration, p_2 , p_3 and p_6 are removed from Candidate_Set. When the iteration is carried out to the fourth time, p_6 should be processed but it has not been at Candidate_Set hence it fulfils the condition of phase 2 of SAVO algorithm. It needs to compare with tuples still be in Candidate_Set and they should have never been compared with each other before. Currently, Candidate_Set only have two tuples, p_4 and p_7 , which have been compared with p_6 . In the end, p_6 processedCount is 1. Until the twelfth iteration, p_7 processedCount \leq p_7 dimCount \leq 3, so it can be removed from Candidate_Set to SSRS. By the eighteenth iteration, p_5 is processed and there are no remaining tuples has not been processed. Therefore, SAVO algorithm reaches to its end. We attain two tuples, p_4 and p_7 , which will become candidates' skyline points via strict clustering strategy and skyline preference query algorithm based on attribute value ordering.

B. Skyline Preference Query Strategy Based On Loose Clustering

According to the loose clustering rule described in Section III, we can attain the loose clusters on IS^{00} . But tuples in each cluster might have no dominance relation. These tuples will also be the skyline query result return to users. Thus, users need to select what they want because tuples in the result set may not meet the needs of users.

~~Algorithm 2 Skyline Preference Query Algorithm Based on Attribute Value Ordering SAVO~~

1. Input: IS^0
2. Output: SSRS

3. initialize Candidate_Set D IS^0 , SSRS D ;
iterD1
4. **Foreach** dimension s_i 2 IS^0
5. This step sorts the attribute value at s_i by non-ascending order and stores corresponding tuple id in array D_i . If tuples have the same attribute value, the smaller tuple id, the tuple has more priority to be stored.
6. The pointer ptr_i in array D_i is initialized by 1.
7. **Endfor**
8. **Foreach** p^0 2 IS^0
9. Initialize p^0 processedCount D 01
10. p^0 dimCount D the number of dimensions that values are non-missing;
11. **Endfor**
12. **While** Candidate_Set $6D$; and at least one tuple has not been processed
13. nextDim D iter% ;
14. P D D nextDim $[ptr$ nextDim]
15. P isDominated D false;
16. **If** P processedCount D 0
17. **If** P 2 Candidate_Set
18. **For** each q 2 (IS^0 -Candidate_Set)
19. **If** (P has never compared with q)
20. **If** (**dominate**(P, q) DD 0)
21. P isDominated D true;
22. **Endif**
23. **Endif**
24. **Endfor**
25. **Endif**
26. **If** P isDominated DD false
27. **Foreach** (c in Candidate_Set)
28. **If** (there has no comparison between p and c before)
29. **If** (**dominate**(P, c) DD 0)
30. P isDominated D true;
31. **Else**
32. Remove c from Candidate_Set;
33. **Endif**
34. **Endfor**
35. **Endif**
36. **If** P isDominated DD true and P in Candidate_Set
37. Remove P from Candidate_Set;


```

38. Endif
39. Endif
40. P processedCount C C;
41. If (P in Candidate_Set and P processedCount
DD P dimCount)
42. Remove P from Candidate_Set to SSRS;
43. Endif
44. IterCC;
45. ptrnextDim C CI
46. Endwhile
47. Return SSRS;

```

Considering this problem, we take the user preference as an important factor in the decision of dominance. A calculation method of dominance degree is proposed to reduce the redundancy of skyline query result set, which differs from most traditional methods. The dominance relation of tuples with user preference is determined by the degree of dominance between them. The following is the definition of the degree of dominance:

Definition 5: The Degree of Dominance

The difference of attribute value of any two tuples can be regarded as their dominance distance. The degree of dominance between any two tuples at the same cluster can be recorded as the sum of the product of dominance distance and weight on comparable dimensions. $w_{c1}, w_{c2}, \dots, w_d$ are the weights of dimensions. The missing values of each dimension of tuple p_i are denoted by $v_{ij}, j \in [C \setminus d]$. The degree of dominance between any two tuples p_i and p_j is $\text{dom}_{i;j} = \sum_{k \in [C \setminus d]} w_k (v_{ik} - v_{jk})$, wherein k represents the k th comparable dimension of a tuple. If $\text{dom}_{i;j} > 0$, then $p_j < p_i$. If $\text{dom}_{i;j} < 0$, $p_j > p_i$. If $\text{dom}_{i;j} = 0$, there are no dominance relation between p_i and p_j . In each cluster, we calculate the degree of dominance of any two tuples. After the eliminating tuples that are dominated by others, remaining tuples at clusters will be added into loose clustering skyline query result set, RSRS

V. CONCLUSIONS

Regarding the incomplete and missing information within the Field of Internet-of-Things, this paper studies the customized recommendation which mixes the large knowledge of the Internet-of-Things with user preference. We have a tendency to propose a sky-line preference skyline strategy supported huge and incomplete data set, as well as encryption and clustering strategy, 2 methods of skyline preference question supported dimension importance clump and a variety strategy supported data entropy theory. As well as sensible implementations, we have a tendency to style a comparative experiment with Skyline algorithmic program and SIDS algorithmic program regarding 3 aspects: knowledge size, knowledge dimension, and knowledge missing rate. Experimental results supported synthetic knowledge sets manifest that SPQ algorithmic program is competitively a lot of efficient and positive in results accuracy than alternative algorithms.

VI. REFERENCES

- [1]. D. Gil, A. Ferrandez, H. Mora-Mora, and J. Peral, 'Internet of Things: A review of surveys based on context aware intelligent services,' *Sensors*, vol. 16, no. 7, p. E1069, Jul. 2016.
- [2]. H. Song, D. Rawat, S. Jeschke, and C. Brecher, *Cyber-Physical Systems: Foundations, Principles and Applications*. Boston, MA, USA: Academic, 2016, p. 514. G. G. Zhang, Y. Bi, C. Li, Y. Zhang, and C. Zeng, 'Security processing model research based on massive IoT data,' *J. Chin. Comput. Syst.*, no. 9, 2090-2094, Sep. 2013.
- [3]. L. Zhang, Y. Wang, B. Y. Song, X. Li, and X. Hao, 'Geometry-based spatial skyline query in wireless sensor network,' in *Proc. 11th Web Inf. Syst. Appl. Conf.*, Tianjin, China, Sep. 2014, pp. 27-32.
- [4]. Y. Wang, B. Y. Song, J. L. Wang, L. Zhang, and L. Wang, 'Geometry-based distributed spatial skyline queries in wireless sensor networks,' *Sensors*, vol. 16, no. 4, p. 454, Apr. 2016.

- [5]. Y. Gu, G. Yu, X. J. Li, and Y. Wang, 'RFID data interpolation algorithm based on dynamic probabilistic path-event model,' *J. Softw.*, no. 3, 438-451, Mar. 2010.
- [6]. Y. Wang, B. Yin, G. H. Liu, B. Y. Song, and J. L. Wang, 'Skyline query of massive incomplete data based on combinational dimensions,' *Frontiers Comput. Sci. Technol.*, no. 4, pp. 495-503, Sep. 2016.
- [7]. S. Jeschke, C. Brecher, H. Song, and D. Rawat, *Industrial Internet of Things*. Switzerland: Springer, 2017, p. 715.
- [8]. L. L. Ding, J. C. Xin, G. R. Wang, and S. Huang, 'Efficient skyline query processing of massive data based on Map-Reduce,' *Chin. J. Comput.*, vol. 34, pp. 1786-1796, Oct. 2011.
- [9]. S. Borzsonyi, D. Kossmann, and K. Stocker, 'The skyline operator,' in *Proc. 17th ICDE*, 2001, pp. 421-430.
- [10]. J. J. Cao, X. C. Diao, S. Chen, and Y. Z. Shao, 'Data cleaning and its general system framework,' *Comput. Sci.*, vol. 39, pp. 207-211, Nov. 2011.
- [11]. Y. H. Lin, C. H. Zhang, and J. Liu, 'Realization of data cleaning based on editing rules and master data,' *Comput. Sci.*, vol. 39, pp. 174-176, Nov. 2012.
- [12]. W. Chen and Q. L. Ding, 'Cleaning method for incomplete data in data cleaning,' *Microcomput. Appl.*, no. 2, pp. 44-45, Feb. 2005.
- [13]. X. J. Wei, J. Yang, C. P. Li, and H. Chen, 'Skyline query processing,' *Softw.*, vol. 19, pp. 1386-1400, Jun. 2008.
- [14]. L. Zhu, J. H. Guan, and S. G. Zhou, 'Review of skyline computing research,' *Comput. Eng. Appl.*, no. 6, pp. 160-165, Feb. 2008.
- [15]. Z. Zhang, H. Lu, B. C. Ooi, and A. K. Tung, 'Understanding the meaning of a shifted sky: A general framework on extending skyline query,' *VLDB J.*, vol. 19, no. 2, pp. 181-201, Feb. 2010.
- [16]. Y. Y. Li, Z. Y. Li, M. X. Dong, W. Y. Qu, C. Q. Ji, and J. F. Wu, 'Efficient subspace skyline query based on user preference using mapreduce,' *Ad Hoc Netw.*, vol. 35, pp. 105-115, Nov. 2015.
- [17]. M. E. Khalefa, M. F. Mokbel, and J. J. Levandoski, 'Skyline query processing for incomplete data,' in *Proc. 24th ICDE*, 2008, pp. 556-565.
- [18]. R. Bharuka and P. S. Kumar, 'Finding skylines for incomplete data,' in *Proc. 24th ADC*, 2013, pp. 109-117.
- [19]. L. Zhang, P. Zou, Y. Jia, and L. Tian, 'Continuous dynamic skyline queries over data stream,' *J. Comput. Res. Develop.*, vol. 48, pp. 77-85, Jan. 2011.
- [20]. K. Ahmed, N. S. Naik, and M. A. Gregory, 'Enhanced distributed dynamic skyline query for wireless sensor networks,' *J. Sens. Actuators Netw.*, vol. 5, no. 1, p. 2, Mar. 2016.
- [21]. Y. Park, J. K. Min, and K. Shim, 'Parallel computation of skyline and reverse skyline queries using mapreduce,' *Proc. VLDB Endowment*, vol. 6, no. 14, pp. 2002-2013, Sep. 2013.
- [22]. Databasesports Website. [Online]. Available: <http://databasebasketball.com>