

K-Dominant Skyline Join Queries: Extending the Join Paradigm to K-Dominant Skylines

Anuradha Awasthi Arnab Bhattacharya Sanchit Gupta Ujjwal Kumar Singh
Dept. of Computer Science and Engineering, Indian Institute of Technology, Kanpur, India
{aawasthi,arnabb,sanchitg,ujjkumsi}@cse.iitk.ac.in

ABSTRACT

Skyline queries enable multi-criteria optimization by filtering objects that are worse in all the attributes of interest than another object. To handle the large answer set of skyline queries in high-dimensional datasets, the concept of k -dominance was proposed where an object is said to dominate another object if it is better (or equal) in at least k attributes. This relaxes the full domination criterion of normal skyline queries and, therefore, produces lesser number of skyline objects. This is called the k -dominant skyline set. Many practical applications, however, require that the preferences are applied on a joined relation. Common examples include flights having one or multiple stops, a combination of product price and shipping costs, etc. In this paper, we extend the k -dominant skyline queries to the join paradigm by enabling such queries to be asked on joined relations. We call such queries *KSJQ* (*k-dominant skyline join queries*). The number of skyline attributes, k , that an object must dominate is from the combined set of skyline attributes of the joined relation. We show how pre-processing the base relations helps in reducing the time of answering such queries over the naïve method of joining the relations first and then running the k -dominant skyline computation. We also extend the query to handle cases where the skyline preference is on aggregated values in the joined relation (such as total cost of the multiple legs of the flight) which are available only after the join is performed. In addition to these problems, we devise efficient algorithms to choose the value of k based on the desired cardinality of the final skyline set. Experiments on both real and synthetic datasets demonstrate the efficiency, scalability and practicality of our algorithms.

Keywords

Skyline query; K-dominant skyline query; Join; Aggregation; K-dominant skyline join query; KSJQ

1. INTRODUCTION AND MOTIVATION

Skyline queries are widely used to enable multi-criteria decision making in databases [3]. Consider a scenario where a person wants to buy a good house. Her preferences are low cost, proximity to market, quiet neighborhood, etc. In a real scenario, it is almost

impossible to find a single house that is best in all her preferences. The skyline query helps her narrow down the choices by filtering out houses that are *worse* (or equal) than some other house in *all* the preferences. Assuming rationality, her choices cannot lie outside the *skyline set*.

In high-dimensional spaces, however, the skyline set becomes less useful due to its impractically large size. The size tends to increase exponentially with dimensionality [28]. This happens as it becomes harder for any object to dominate another object in *all* the attributes. The problem is especially severe in real datasets where the data is generally anti-correlated in nature. For example, a quiet neighborhood closer to a market is likely to be more costly.

There have been various works on handling the large cardinality of skyline sets in high dimensions. The most prominent is that of k -dominant skylines where, instead of being better in all the d dimensions, an object need only be better in some $k < d$ dimensions to dominate another object [4]. As a result, it becomes easier for an object to be dominated which leads to lesser number of skylines.

The k -dominant skylines are quite useful in real scenarios. In the example of housing discussed above, if there are many attributes, it may be rare that one house is better than another on all the counts. Instead, if a smaller number of attributes, say $k = 2$, is specified, there are more chances of finding a house that has a lower cost and a quieter neighborhood (but may not be closer to a market) than another house. As a result, more houses can be filtered, and the retrieved skyline set becomes more manageable and useful.

To the best of our knowledge, however, the k -dominant skyline queries have not been explored for *multiple* relations.¹ Suppose there are two relations having d_1 and d_2 skyline attributes. After joining, a bigger relation with $d_1 + d_2$ skyline attributes is formed. (We discuss the different variants and restrictions later.) A k -dominant skyline, where $k < d_1 + d_2$, is then sought on this joined relation.

A real-life example of this situation happens often in flight bookings. Suppose a person wants to fly from city A to city B. Her preferences are lower cost, lower duration, higher ratings and higher amenities. While the basic skyline works for direct flights, in many cases, a flight route from A to B includes one (or more) stopovers. Thus, a valid flight path contains the join of all flights from city A to other cities and from those cities to city B where the intermediate city is the same. The preferences a user would want now applies to the *entire* flight path and not a single leg of the journey. The skyline is, therefore, needed on the joined relation [2, 21].

Once more, it is harder for a flight combination to dominate another flight combination in all the skyline attributes over the two relations. Here, the k -dominant skyline query is a natural choice, where k is less than the total number of the skyline attributes in the

¹A preliminary version of this paper will appear as a poster [1].

joined relation.

With the increase in the number of attributes, the size of the skyline set increases further and, hence, computing the k -dominant skylines becomes even more relevant and useful.

The naïve approach first creates the joined relation and then subsequently computes the k -dominance. This strategy, while straightforward, is inefficient and impractical for large datasets.

A further practical consideration in the flight example is that a user is not really bothered about the individual legs, but rather the *total* cost and *total* duration of the journey. Thus, the skyline preferences should be applied on the *aggregated* values of attributes from the base relations. Note that the aggregated values are available only after the join and are, therefore, harder to process efficiently.

In this paper, we explore the question of finding k -dominant skylines on joined relations, where the skyline preferences can be on both aggregated and individual values. Apart from posing the problem, our main contribution is to push the skyline operator before the join as much as possible, thereby making the whole algorithm efficient and practical.

In addition to finding k -dominant skylines, an important question that often arises in practical applications is how to choose a “good” value of k ? While there is no universal answer, one of the guiding principles is the number of skyline objects finally returned [4]. The basic idea of skyline queries is to serve as a filter for poor objects and, thus, a user may find it easier to specify a value of δ objects that she is interested in examining more thoroughly rather than a value of k . Since the size of the skyline set increases with k , the “optimal” value of k may be then taken as the smallest one that returns at least δ skyline objects, or the largest one that returns at most δ skyline objects.

We address the above question in the context of k -dominant skylines over joined relations.

In sum, our contributions are:

1. We propose the problem of finding k -dominant skyline queries over joined relations. We term such queries *KSJQ*.
2. We design efficient algorithms to solve the *KSJQ* problem.
3. We devise ways to arrive at a good value of “ k ” by specifying a threshold size of the k -dominant skyline set.

The rest of the paper is organized as follows. Sec. 2 sets the background on joins of k -dominant skylines. Using this, different problem statements are defined in Sec. 3. Sec. 4 outlines the related work. Various optimizations that can improve the efficiency of the problem are explained in Sec. 5. Algorithms that use these optimizations are described in Sec. 6. Sec. 7 analyzes the experimental results before Sec. 8 concludes.

2. BACKGROUND

2.1 Skylines

Consider a dataset R of objects. For each object u , a set of d attributes $\{u_1, \dots, u_d\}$ are specified, which form the *skyline attributes*. For each of the skyline attributes, without loss of generality, the preference is assumed to be less than ($<$), i.e., a lower value is preferred over a higher one. An object u *dominates* another object v , denoted by $u \succ v$, if and only if, for all the d skyline attributes, u_i is preferred over or equal to v_i , and there exists at least one skyline attribute where u_j is strictly preferred over v_j . The *skyline set* $S \subseteq R$ contains objects that are *not* dominated by any other object [3]. In other words, every object in the non-skyline set is dominated by at least one object in the dataset.

2.2 K-Dominant Skylines

The k -dominant skyline query [4] relaxes the definition of domination between objects. An object u k -dominates another object v , denoted by $u \succ_k v$, if and only if, for at least k of the d skyline attributes, u_i is preferred over or equal to v_i , and there exists at least one skyline attribute where u_j is strictly preferred over v_j . The k -dominant skyline set contains objects that are not k -dominated by any other object. The k attributes are *not* fixed and can be *any* subset of the d skyline attributes.

The k -dominant query is particularly problematic when $k \leq \frac{d}{2}$ since then two objects can dominate each other. Even when $k > \frac{d}{2}$, the k -dominance relationship is *not transitive* and may be even *cyclic*: $u \succ_k v \succ_k w \succ_k u$.

2.3 Multi-Relational Skylines

Consider two datasets R_1 and R_2 with d_1 and d_2 skyline attributes respectively. The *join* of the two relations (using appropriate join conditions) forms the dataset $R = R_1 \bowtie R_2$ with $d = d_1 + d_2$ skyline attributes. The *multi-relational skyline* is extracted from R [21].

In a significant variant of the problem, a number of skyline attributes in each relation are marked for *aggregation* with corresponding attributes from the other relation [2]. Thus, a attributes out of d_1 in R_1 and d_2 in R_2 are aggregated in the joined relation R . As a result, R contains $(d_1 - a) + (d_2 - a) + a = d_1 + d_2 - a$ skyline attributes. The skyline query is then asked over these attributes. The aggregation function is assumed to be *monotonic*; this ensures that if the base values of two tuples $u \in R_1$ and $t \in R_2$ are preferred over the base tuples of two other tuples $v \in R_1$ and $s \in R_2$ respectively, the aggregated value of $u \bowtie t \in R$ will be necessarily preferred over the aggregated value of $v \bowtie s \in R$.

The case for more than two base relations can be handled by cascading the joins.

3. PROBLEM STATEMENTS

We propose two main variants of what we call the *K-DOMINANT SKYLINE JOIN QUERY (KSJQ)* problem.

The first variant works on relations for which the skyline attributes are strictly *local*.

The number of skyline attributes in the first and second relations, R_1 and R_2 , are d_1 and d_2 respectively.

$$R_1 = \{h_{1_1}, \dots, h_{1_m}, s_{1_1}, \dots, s_{1_{d_1}}\} \quad (1)$$

$$R_2 = \{h_{2_1}, \dots, h_{2_m}, s_{2_1}, \dots, s_{2_{d_2}}\} \quad (2)$$

$$R = \{h_1, \dots, h_m, s_{1_1}, \dots, s_{1_{d_1}}, s_{2_1}, \dots, s_{2_{d_2}}\} \quad (3)$$

where h_i captures the join of h_{1_i} with h_{2_i} in the joined relation $R = R_1 \bowtie R_2$, and s_{1_i}, s_{2_i} are the skyline attributes.

PROBLEM 1 (KSJQ). *Given two datasets R_1 and R_2 having d_1 and d_2 skyline attributes respectively, find k -dominant skylines from the joined relation $R = R_1 \bowtie R_2$ having $d = d_1 + d_2$ skyline attributes.*

We restrict k to be at least one more than the dimensionality in the base relations, i.e., $\max\{d_1, d_2\} < k < d$. This restriction constrains at least some skyline attributes from each relation to satisfy the preferences. In other words, the k preferred attributes must span *both* the relations. However, there is no restriction over how many skyline attributes from each relation must be satisfied. Denoting the number of skyline attributes that are chosen from the two relations by k_1 and k_2 where $k_1 + k_2 = k$, this implies that $1 \leq k_1 \leq d_1$ and $1 \leq k_2 \leq d_2$.

The second variant works on the *aggregate* version of the problem. Assuming a total of a aggregate attributes, the number of

attributes on which only local skyline preferences are applied are $d_1 - a = l_1$ and $d_2 - a = l_2$ respectively. The joined relation R , therefore, contains $l_1 + l_2$ local attributes and a aggregate attributes.

PROBLEM 2 (AGGREGATE KSJQ). *Given two datasets R_1 and R_2 having d_1 and d_2 skyline attributes respectively, of which a attributes are used for aggregation, find k -dominant skylines from the joined relation $R = R_1 \bowtie R_2$ having $d_1 + d_2 - a$ skyline attributes.*

Once more, we assume that $\max\{d_1, d_2\} < k \leq d$. Expressing in terms of local and aggregate attributes, $\max\{l_1, l_2\} + a < k \leq l_1 + l_2 + a$.

The third and fourth problems address the tuning of the value of k . The value can be tuned in two ways.

PROBLEM 3 (AT LEAST δ). *Given a KSJQ query framework and a threshold number of skyline objects δ , determine the smallest value of k that returns at least δ skyline objects.*

PROBLEM 4 (AT MOST δ). *Given a KSJQ query framework and a threshold number of skyline objects δ , determine the largest value of k that returns at most δ skyline objects.*

Prob. 4 is directly linked with Prob. 3. If k^* is the answer to Prob. 3, then the answer to Prob. 4 must be $k^* - 1$. There are two corner cases. First, if $k^* = 1$, then the answer to Prob. 4 should be trivially 1 as well. Second, if k^* -dominant skyline returns *exactly* δ skyline objects (or $k^* = d$), then the answer to Prob. 4 is k^* as well. Thus, henceforth, we focus only on Prob. 3.

4. RELATED WORK

The skyline operator was introduced in databases by [3] by adopting the maximal vector or the Pareto optimal problem [15]. Several indexed [14], [17] and non-indexed algorithms [3, 5, 22] have been since proposed to retrieve the skyline set.

Analyses of the cardinality of the skyline set [9, 28] have shown that the number of skyline objects can grow exponentially with the increase in dimensionality. Consequently, several attempts have been made to restrict the size of the skyline set. These mostly include notions of approximate skylines or representative skylines [6, 8, 10, 12, 13, 16, 23, 25, 26, 27].

A completely different approach—*k-dominant skylines*—was proposed by [4]. It uses subsets of skyline attributes to control the cardinality of the skyline set. Although the parameter k is an input to the problem, the authors also proposed a way to derive the smallest k that will guarantee a δ number of skylines.

There have been many recent works on k -dominant skylines including extensions to high-dimensional spaces [18, 19], using parallel processing [24], cardinality estimation [11] as well as in update-heavy datasets [20] and combined datasets [7].

The skyline join problem where skylines are retrieved from joined relations, both components of which contain skyline attributes, was introduced in [21]. This work, however, simply assumed that the skyline attributes remain unchanged from the base relations to the joined one. The aggregate skyline join queries [2] removed this restriction by enabling skyline preferences to be posed on aggregate values of attributes formed by combining attributes from the base relations. The aggregation function was assumed to be monotonic.

In this paper, we pose the k -dominant skyline problem in the join paradigm (including the aggregated version). To the best of our knowledge, this is the first work in this direction.

fno	destination	cost	dur	rtg	amn	category
11	C	448	3.2	40	40	SS_1
12	C	468	4.2	50	38	NN_1
13	D	456	3.8	60	34	SN_1
14	D	460	4.0	70	32	NN_1
15	E	450	3.4	30	42	SN_1
16	F	452	3.6	20	36	SS_1
17	G	472	4.6	80	46	SN_1
18	H	451	3.7	20	37	SS_1
19	E	451	3.7	40	37	NN_1

Table 1: Flights from city A (f_1).

5. OPTIMIZATIONS

In this section, we describe the various optimization schemes that can be used to speed-up the process of finding k -dominant skylines in joined relations. Sec. 6 uses these optimizations to design the algorithms.

5.1 Join Attributes

We assume an *equality* join condition. Two base tuples $u \in R_1$ and $v \in R_2$ can be joined to form $t = u \bowtie v \in R$ if and only if all their join attributes match. Referring to Eq. (3), $\forall_{j=1, \dots, m} h_{1_j} = h_{2_j}$. (We discuss the relaxation of this assumption in Sec. 6.6.)

ASSUMPTION 1 (EQUALITY JOIN). *The join conditions are all based on equality, i.e., the join is an equality join.*

All the tuples in a base relation are, thus, divided into *groups* according to the value of the join attributes. In every group, the values of the join attributes, h_{i_1}, \dots, h_{i_m} ($i = 1, 2$), are the same.

In the flight example, since the joining criterion is destination of first flight to be the source of the second flight, the first base relation is divided on the basis of destinations while the second one is divided on sources.

5.2 Grouping

Based on the k -dominance properties within each group, each base relation is partitioned into different sets as follows.

A tuple may be a k -dominant tuple for the entire base relation. However, even when it is not, it may not be k -dominated by any other tuple in its group. It is then a k -dominant tuple when only its group is concerned.

Based on the above notion of k -dominance within a group, each base relation R_i is divided into 3 *mutually exclusive* and *exhaustive* sets SS_i , SN_i , and NN_i :

$$R_i = SS_i \cup SN_i \cup NN_i \quad (4)$$

We next define the three sets.

DEFINITION 1 (SS). *A tuple u is in SS if u is a k -dominant skyline in the overall relation; consequently, it is a k -dominant skyline in its group as well.*

DEFINITION 2 (SN). *A tuple u is in SN if u is a k -dominant skyline only in its group but not in the overall relation.*

DEFINITION 3 (NN). *A tuple u is in NN if u is not a k -dominant skyline in its group; consequently, it is not a k -dominant skyline in the overall relation as well.*

Consider the examples in Table 1 and Table 2. We assume that all the attributes have lower preferences². We set $k = 3$ for both

²Although the preferences for ratings and amenities are generally the other way round, we stick to lower preferences for ease of understanding.

fno	source	cost	dur	rtg	amn	category
21	D	348	2.2	40	36	SS_2
22	D	368	3.2	50	34	NN_2
23	C	356	2.8	60	30	SN_2
24	C	360	3.0	70	28	NN_2
25	E	350	2.4	30	38	SN_2
26	F	352	2.6	20	32	SS_2
27	G	372	3.6	80	42	SN_2
28	H	350	2.4	35	37	SN_2

Table 2: Flights to city B (f_2).

the relations. The categorization of the tuples are shown in the last column.

Table 3 shows the joined relation.

The important outcome of this division into the three sets is that it allows certain tuples to be automatically designated as k -dominant skylines (or not) *without* computing the join, as explained next.

5.3 Deciding about Skylines before Join

We first analyze a simple case where out of the final k attributes in the joined relation, if a joined tuple t k -dominates another joined tuple s , t must k_1 -dominate s in the first base relation and must k_2 -dominate s in the second base relation ($k = k_1 + k_2$).

Note that, in practice this situation will not be specified by any user. We are only going to use it to explain the basic concepts and then build upon it later by removing this assumption. For the theorems and observations in this section, the first base relation is divided into the three sets SS_1 , SN_1 , NN_1 according to k_1 -domination while the second base relation is divided into SS_2 , SN_2 , NN_2 using k_2 -domination.

The first theorem shows that a tuple formed by joining the SS counterparts is *always* a k -dominant skyline.

THEOREM 1. *The tuples in the set $(SS_1 \bowtie SS_2)$ are k -dominant skylines.*

PROOF. Consider a joined tuple $t' = u' \bowtie v'$ formed by joining the tuples $u' \in SS_1$ and $v' \in SS_2$. Assume that $t = u \bowtie v$ k -dominates t' , i.e., $t \succ_k t'$. Since $u' \in SS_1$, no tuple and, in particular, u can k_1 -dominate u' . Similarly, $v \not\succeq_{k_2} v'$. Therefore, $\nexists t, t \succ_k t'$. Hence, t' is a k -dominant tuple. \square

The flight combination (16,26) in Table 3 is an example.

The second theorem shows the reverse: composite tuples formed by joining a base tuple in NN can *never* be k -dominant skylines.

THEOREM 2. *The tuples in the sets $(SS_1 \bowtie NN_2)$, $(SN_1 \bowtie NN_2)$, $(NN_1 \bowtie SS_2)$, $(NN_1 \bowtie SN_2)$, and $(NN_1 \bowtie NN_2)$ are not k -dominant skylines.*

PROOF. Consider a joined tuple $t' = u' \bowtie v'$ formed by joining the tuples $u' \in NN_1$ and $v' \in R_2$. Therefore, there must exist a tuple u in the same group that k_1 -dominates u' , i.e., $\exists u, u \succ_{k_1} u'$. Consider the tuple t formed by joining u with v' . Since u is in the same group as u' , this tuple necessarily exists. As t dominates t' in k_1 attributes and is equal in k_2 , overall, it dominates in $k_1 + k_2 = k$ attributes, i.e., $t \succ_k t'$. Therefore, t' is not a k -dominant skyline. This covers the cases $(NN_1 \bowtie SS_2)$, $(NN_1 \bowtie SN_2)$, and $(NN_1 \bowtie NN_2)$. The cases for $(SS_1 \bowtie NN_2)$ and $(SN_1 \bowtie NN_2)$ are symmetrical. \square

The flight combinations (11,24), (13,22), (14,21), (14,22), (12,23), and (12,24) in Table 3 exemplify the above cases. For example, the tuple (11,24) is dominated by (11,23).

However, nothing can be concluded surely about the rest of the joined tuples.

OBSERVATION 1. *The tuples in the sets $(SS_1 \bowtie SN_2)$ and $(SN_1 \bowtie SS_2)$ are most likely to be k -dominant skylines, although that is not guaranteed.*

PROOF. Consider a joined tuple $t' = u' \bowtie v'$ formed by joining the tuples $u' \in SS_1$ and $v' \in SN_2$. Consider the dominator $v \in R_2 \succ_{k_2} v'$. Since v' is a k_2 -dominate skyline in R_2 , v is in a different group than v' . Therefore, v cannot join with u' , i.e., $\nexists t = u' \bowtie v, v \succ_{k_2} v'$. Although no tuple u can k_1 -dominate u' , there may exist u whose k_1 attributes have the same value as u' . If it happens that u is join-compatible with v , then and only then, $t = u \bowtie v$ exists and $t \succ_k t'$ (since it is equal in k_1 and better in k_2 attributes). The above situation is quite unlikely, although not impossible. The case for $(SN_1 \bowtie SS_2)$ is symmetrical. \square

Thus, while (11,23) and (13,21) are k -dominant skylines, (18,28) is not (Table 3). Flight 18 has same k_1 values as 19 while 28 is dominated by 25. Therefore, (18,28) is k -dominated by (19,25).

The observation for $(SN_1 \bowtie SN_2)$ is similar.

OBSERVATION 2. *A tuple in the set $(SN_1 \bowtie SN_2)$ may or may not be a k -dominant skyline.*

PROOF. Consider a joined tuple $t' = u' \bowtie v' \in (SN_1 \bowtie SN_2)$. Consider the tuples u and v such that $u \succ_{k_1} u'$ and $v \succ_{k_2} v'$. Since $u' \in SN_1$, u is in a different group than u' . Similarly, v is in a different group than v' . If and only if u and v are join compatible, they will join to form $t = u \bowtie v$ which then k -dominates t' . Otherwise, no such t exists, and t' is a k -dominant skyline. \square

Consider (15,25) in Table 3. Since its dominators, flights 11 and 21 respectively, are not join compatible (11 reaches city C while 21 takes off from city D), the flight combination (11,21) is not valid. Consequently, no joined tuple dominates it, and (15,25) becomes a k -dominant skyline. On the other hand, for (17,27), the dominators 16 and 26 do join (the city F is common) to form the tuple (16,26). As a result, (17,27) is not a k -dominant skyline.

The overall situation is summed up in Table 4.

5.4 Skyline Attributes in Joined Relation

We next do away with the assumption that k_1 and k_2 attributes have to be satisfied separately from the first and second relations. It is simply required that the joined relation returns k -dominant skylines with no restriction on how k is broken up between the base relations. However, to ensure that the skyline preferences are respected for *at least* one attribute in every base relation, we assume that $k > \max\{d_1, d_2\}$.

A brute-force way to find the skylines is to generate all combinations of k_1 and k_2 such that $k_1 + k_2 = k$ and, then, combine the answer sets using the results obtained in the previous section. However, it can be done more efficiently as explained next.

We consider the following cases:

$$k'_1 = k - d_2 \quad k'_2 = k - d_1 \quad (5)$$

For all combinations of k_1 and k_2 such that $k_1 + k_2 = k$, the inequalities, $1 \leq k'_1 \leq k_1 \leq d_1$ and $1 \leq k'_2 \leq k_2 \leq d_2$, hold.

Consider the example in Table 3. If $k = 7$, then $k'_1 = k'_2 = 3$. Thus, categorization and k -dominant skyline sets remain the same.

We first establish the following lemma on the monotonicity of the number of skyline attributes.

LEMMA 1. *If tuple u is a j -dominant skyline, it is also a i -dominant skyline for any $i \geq j$.*

PROOF. Consider u to be a j -dominant tuple. Assume that, however, it is not i -dominant for some $i > j$. Thus, there exists

fno	stop-over	f1.cost	f1.dur	f1.rtg	f1.amn	f2.cost	f2.dur	f2.rtg	f2.amn	categorization	skyline
(11,23)	C	448	3.2	40	40	356	2.8	60	30	$SS_1 \bowtie SN_2$	yes
(11,24)	C	448	3.2	40	40	360	3.0	70	28	$SS_1 \bowtie NN_2$	no
(12,23)	C	468	4.2	50	38	356	2.8	60	30	$NN_1 \bowtie SN_2$	no
(12,24)	C	468	4.2	50	38	360	3.0	70	28	$NN_1 \bowtie NN_2$	no
(13,21)	D	456	3.8	60	34	348	2.2	40	36	$SN_1 \bowtie SS_2$	yes
(13,22)	D	456	3.8	60	34	368	3.2	50	34	$SN_1 \bowtie NN_2$	no
(14,21)	D	460	4.0	70	32	348	2.2	40	36	$NN_1 \bowtie SS_2$	no
(14,22)	D	460	4.0	70	32	368	3.2	50	34	$NN_1 \bowtie NN_2$	no
(15,25)	E	450	3.4	30	42	350	2.4	30	38	$SN_1 \bowtie SN_2$	yes
(16,26)	F	452	3.6	20	36	352	2.6	20	32	$SS_1 \bowtie SS_2$	yes
(17,27)	G	472	4.6	80	46	372	3.6	80	42	$SN_1 \bowtie SN_2$	no
(18,28)	H	451	3.7	20	37	350	2.4	35	39	$SS_1 \bowtie SN_2$	no
(19,25)	E	451	3.7	40	37	350	2.4	30	38	$NN_1 \bowtie SN_2$	no

Table 3: Joined relation ($f_1 \bowtie f_2$).

	SS_2	SN_2	NN_2
SS_1	yes (Th. 1)	likely (Obs. 1)	no (Th. 2)
SN_1	likely (Obs. 1)	may be (Obs. 2)	no (Th. 2)
NN_1	no (Th. 2)	no (Th. 2)	no (Th. 2)

Table 4: Fate of k -dominant skylines.

	SS_2	SN_2	NN_2
SS_1	yes (Th. 3)	likely (Obs. 3)	no (Th. 4)
SN_1	likely (Obs. 3)	may be (Obs. 4)	no (Th. 4)
NN_1	no (Th. 4)	no (Th. 4)	no (Th. 4)

Table 5: Joins of groups.

v that i -dominates u , i.e., v is better than u in i attributes. Then, v must be better than u in some subset j of these i attributes, which is a contradiction. Thus, u must be i -dominant for every $i \geq j$. \square

The following theorems and observations implicitly use Lemma 1.

THEOREM 3. *The tuples in the set ($SS_1 \bowtie SS_2$) are k -dominant skylines.*

PROOF. Consider $t' = u' \bowtie v' \in (SS_1 \bowtie SS_2)$. There are two cases to consider: (i) when there does not exist any tuple u that shares k'_1 attributes with u' , and (ii) when there exists a tuple u that share k'_1 or more attributes with u' .

In the first case, if $\exists t, t \succ_k t'$, then t must dominate t' in at least k'_1 attributes corresponding to u' since it can dominate t in at most d_2 attributes corresponding to v' . Since this is impossible, t' is a k -dominant skyline.

In the second case, u' can be dominated by $u \in R_1$ in at most $d_1 - 1$ attributes (otherwise $u' \notin SS_1$). Thus, to dominate t' , u must combine with $v \in R_2$ such that v dominates v' in at least $k - (d_1 - 1) = k_2 + 1$ attributes, which is impossible since $v' \in SS_2$. Hence, t' is a k -dominant skyline. \square

THEOREM 4. *The tuples in the sets ($SS_1 \bowtie NN_2$), ($SN_1 \bowtie NN_2$), ($NN_1 \bowtie SS_2$), ($NN_1 \bowtie SN_2$), and ($NN_1 \bowtie NN_2$) are not k -dominant skylines.*

PROOF. In each of the sets, at least one base tuple is in NN . Consider a joined tuple $t' = u' \bowtie v'$ where $u' \in NN_1$. Surely, $\exists u, u \succ_{k'_1} u'$ exists and u is in the same group as u' . Therefore, u is join compatible with v' . Consider $t = u \bowtie v'$. It dominates t' in $k'_1 + d_2 = k$ attributes. Thus, t' is not a k -dominant skyline. \square

For example, $(16,26) \in SS_1 \bowtie SS_2$ and $(14,22) \in NN_1 \bowtie NN_2$ (in Table 3) is a k -dominant skyline and not a k -dominant skyline respectively.

OBSERVATION 3. *A tuple in the set ($SS_1 \bowtie SN_2$) or ($SS_2 \bowtie SN_1$) is most likely to be a k -dominant skyline, although that is not guaranteed.*

PROOF. Consider $t' = u' \bowtie v' \in (SS_1 \bowtie SN_2)$. There are two cases to consider. In the first case when there does not exist any tuple u that share k'_1 attributes with u' , a joined tuple t can dominate t' in at most $k'_1 - 1 + d_2 = k - 1$ attributes. Therefore, t' is a k -dominant skyline.

In the second case, assume that such a u with equal k'_1 attributes exists. Note, however, that u cannot be better in any other attribute as then u' would be k'_1 -dominated by u . Since $v' \in SN_2$, there exists a $v \in R_2$ that dominates v' but is in a different group. If and only if u and v are join compatible, the joined tuple $t = u \bowtie v$ exists. Comparing t against t' , we see that the k'_1 attributes corresponding to u' (or u) are same. Therefore, for t' to k -dominate t , v' must dominate v in all the $k - k'_1 = d_2$ attributes. This is unlikely, although not impossible.

The proof for ($SS_2 \bowtie SN_1$) is similar. \square

For example, consider $(18,28) \in SS_1 \bowtie SN_2$ in Table 3 against $(19,25)$. While $k'_1 = 3$ attributes for 18 and 19 are same, 25 dominates 28 in $d_2 = 4$ attributes. Thus, overall, $(19,25)$ dominates $(18,28)$ in $3+4 = 7$ attributes and, thus, $(18,28)$ is not a k -dominant skyline. On the other hand, $(11,23) \in SS_1 \bowtie SN_2$ is a k -dominant skyline since there is no tuple that shares $k'_1 = 3$ attributes with 11.

OBSERVATION 4. *A tuple in the set ($SN_1 \bowtie SN_2$) may or may not be a k -dominant skyline.*

PROOF. Consider a joined tuple $t' = u' \bowtie v' \in (SN_1 \bowtie SN_2)$. Consider the tuples u and v that dominate u' and v' respectively. Since u and v are in different groups than u' and v' , they may not be join compatible. If they are, the joined tuple $t = u \bowtie v$ dominates t' in at least $k'_1 + k'_2$ attributes. The tuple t may additionally dominate t' in some other attributes such that it overall k -dominates t' . If these two conditions are not met, t' becomes a k -dominant skyline. \square

Out of the two tuples in $SN_1 \bowtie SN_2$ in table 3, while $(15,25)$ is a k -dominant skyline because the dominators 11 and 21 are not in the same group and, therefore, cannot join, the flight combination $(17,27)$ is not a k -dominant skyline as the dominators 16 and 26 join to form $(16,26)$ which overall k -dominates $(17,27)$.

Table 5 summarizes the situation.

It is important to note that for determining the groups in the base relations, it must be the minimum number of attributes considered, i.e., k'_1 and k'_2 ; otherwise, the correctness of the above theorems may be violated.

5.5 Unique Value Property

For a tuple in $SS_1 \bowtie SN_2$ (and $SN_2 \bowtie SS_1$) to be *not* a k -dominant skyline, the number of attributes in which $u \in R_1$ is same as $u' \in SS_1$ must be at least k'_1 . If the base relations follow a *unique value property* where it is guaranteed that for any k'_i ($i = 1, 2$) number of attributes, two tuples will be unique, then the processing becomes simpler.

DEFINITION 4 (UNIQUE VALUE PROPERTY). *A relation R has the unique value property (UVP) with respect to i if for each subset of i skyline attributes of R , all the tuples are unique, i.e., no two tuple will have exactly the same values in any i -sized subset of attributes.*

If $k'_i = 1$, every tuple for any attribute must be unique. While real-valued attributes generally follow that, categorical attributes do not. However, for reasonable values of k'_1 and k'_2 , real datasets having a mix of both real-valued and categorical attributes generally follow the UVP.

The UVP is extremely useful since it ensures that the tuples in $SS_1 \bowtie SN_2$ and $SN_2 \bowtie SS_1$ become k -dominant skylines as shown next.

THEOREM 5. *If relations R_1 and R_2 follow UVP with respect to k'_1 and k'_2 attributes respectively, the tuples in the sets $(SS_1 \bowtie SN_2)$ and $(SN_1 \bowtie SS_2)$ are k -dominant skylines.*

PROOF. Consider a joined tuple t in either of the two sets. The generic situation, as shown in Obs. 3 has two cases. While the first case makes t a k -dominant skyline, the UVP precludes the second case. Thus, t is always a k -dominant skyline. \square

Although we present Th. 5 for the sake of completeness, we do not assume it in our experiments (Sec. 7).

5.6 Aggregate Attributes

We next consider the case of aggregation where the k -dominant skyline is sought over attributes that attain values aggregated from attributes in the base relation.

We assume that there are l_1 local and a aggregate skyline attributes in R_1 and l_2 local and a aggregate skyline attributes in R_2 . The total number of skyline attributes in R is, therefore, $l_1 + l_2 + a$ attributes. As earlier, the final skyline query is asked over $k < l_1 + l_2 + a$ attributes.

The groups in the base relations are partitioned based on both the local and aggregate attributes. Out of the final number of skyline attributes k , a of them can be aggregate. Thus, the minimum number of local attributes that must be dominated in each base relation is $k'_1 = k - a - l_2$ and $k'_2 = k - a - l_1$. The categorization of the base relations into the sets SS , SN and NN are done on the basis of $k'_1 = k'_1 + a$ and $k'_2 = k'_2 + a$. Since $d_1 = a + l_1$ and $d_2 = a + l_2$, these definitions are same as earlier (Sec. 5.4).

We use the following assumption about the *monotonicity* property of the aggregate attributes.

ASSUMPTION 2 (MONOTONICITY). *If the value of attribute u_1 dominates that of u_2 and the value of attribute v_1 dominates that of v_2 , the aggregated value of $u_1 \oplus v_1$ will dominate the aggregated value of $u_2 \oplus v_2$, where \oplus denotes the aggregation operator.*

Since the categorization remains the same, the fate of the joined tuples using the aggregation remains exactly the same as earlier (as summarized in Table 5).

Table 6 shows the joined relation obtained from Table 1 and Table 2 with the cost values aggregated.

Algorithm 1 KSJQ: Naïve Algorithm

Input: Relations R_1, R_2 ; Number of attributes k
Output: k -dominant skyline set T
 1: $D \leftarrow R_1 \bowtie R_2$
 2: $T \leftarrow k$ -dominant skyline(D, k)
 3: **return** T

In the example, considering $k = 6$ with $a = 1$, $k'_1 = 6 - 1 - 3 = 2$ and $k'_2 = 6 - 1 - 3 = 2$ but $k'_1 = k'_1 + 1 = 3$ and $k'_2 = k'_2 + 1 = 3$ remain as earlier.

6. ALGORITHMS

In this section, we describe the various algorithms for answering the k -dominant skyline join queries (Sec. 5.4). We consider the general case where the datasets do not follow UVP (Sec. 5.5) and, where in addition to local skyline attributes, there are aggregate ones as well (Sec. 5.6).

6.1 Naïve Algorithm

The naïve algorithm (Algo. 1) simply computes the join of the two relations first (line 1) and then computes the k -dominant skylines from the joined relation (line 2) using any of the standard k -dominant skyline computation methods [4]. Being the most basic, it suffers from two major disadvantages. First, the join can require a very large time, thereby rendering the entire algorithm extremely time-consuming and impractical. The second shortcoming is the non-progressive result generation. The user has to wait a fairly large time (at least the complete joining time) before even the first skyline result is presented to her. In online scenarios, the progressive result generation is quite an attractive and useful feature.

6.2 Target Set

To alleviate the problems of the naïve algorithms, we next propose two algorithms, *grouping* and *dominator-based*, that use the concepts of optimization from Sec. 5.

However, before we describe them, we first explain and define the concept of *target sets*. Although a tuple in a set marked by “may be” or “likely” is not guaranteed to be a k -dominant skyline, it needs to be checked against only a small set of tuples, called its *target set*.

Formally, a target set for a tuple u' in a base relation is the set of tuples $\tau(u')$ that can potentially combine with other tuples from the other base relation and k -dominate a joined tuple formed with u' . In other words, for a joined tuple $t' = u' \bowtie v'$, there may exist v such that $t = u \bowtie v$ where $u \in \tau(u')$ k -dominates t' . No tuple outside the target set $\tau(u)$ of u may combine with any other tuple and dominate t .

DEFINITION 5 (TARGET SET). *The target set for a tuple $u' \in R_i$ is the set of tuples $\tau(u') \subseteq R_i$ such that $\forall u \notin \tau(u'), \exists v, v', t = u \bowtie v \succ t' = u' \bowtie v'$.*

The definition is one-sided: a tuple in the target set may or may not join and dominate t' , but no tuple outside the target set can join and dominate t' . The utility of a target set is easy to understand. To check whether t' is a k -dominant skyline, u' needs to be checked only against its target set and nothing outside it.

The *join* of target sets for the base tuples produces the potential dominating set for the joined tuple.

The target set for a tuple $u' \in SS$ constitutes itself and the set of tuples $\{u\}$ that has at least k'_1 attributes same as u' . The augmentation is required to guarantee the correctness as explained in

fno	stop-over	cost	f1.dur	f1.rtg	f1.amn	f2.dur	f2.rtg	f2.amn	categorization	skyline
(11,23)	C	804	3.2	40	40	2.8	60	30	$SN_1 \bowtie SN_2$	yes
(11,24)	C	808	3.2	40	40	3.0	70	28	$SN_1 \bowtie NN_2$	no
(12,23)	C	824	4.2	50	38	2.8	60	30	$NN_1 \bowtie SN_2$	no
(12,24)	C	828	4.2	50	38	3.0	70	28	$NN_1 \bowtie NN_2$	no
(13,21)	D	804	3.8	60	34	2.2	40	36	$SN_1 \bowtie SN_2$	yes
(13,22)	D	824	3.8	60	34	3.2	50	34	$SN_1 \bowtie NN_2$	no
(14,21)	D	808	4.0	70	32	2.2	40	36	$NN_1 \bowtie SN_2$	no
(14,22)	D	828	4.0	70	32	3.2	50	34	$NN_1 \bowtie NN_2$	no
(15,25)	E	800	3.4	30	42	2.4	30	38	$SN_1 \bowtie SN_2$	yes
(16,26)	F	804	3.6	20	36	2.6	20	32	$SS_1 \bowtie SS_2$	yes
(17,27)	G	844	4.6	80	46	3.6	80	42	$SN_1 \bowtie SN_2$	no
(18,28)	H	801	3.7	20	37	2.4	35	39	$SS_1 \bowtie SN_2$	no
(19,25)	E	801	3.7	40	37	2.4	30	38	$NN_1 \bowtie SN_2$	no

Table 6: Joined relation ($f_1 \bowtie f_2$): aggregate.

Obs. 3. The tuple u' must be included in the target set of u' for the same reason.

For each tuple in SS , the number of such tuples sharing at least k'_1 attributes is typically low. Hence, maintenance of target sets is quite feasible and practical.

However, the target set for a tuple in SN can be any tuple *outside* its group that dominates it. For simplicity, we consider it as the entire dataset R_i .

Similarly, the target set for NN is R_i .

6.3 Grouping Algorithm

Our first algorithm, called the *grouping* algorithm (Algo. 2), first computes the groups SS , SN and NN in the base relations. Next, the summarization from Table 5 is used. Tuples from the sets marked by “yes” are immediately output as k -dominant skyline tuples. Tuples marked by “no” are pruned and not even joined.

A tuple in a set marked by “may be” or “likely” is checked against the join of the *target* sets of the base tuples.

For a base tuple in the SS group, the target set is first augmented with tuples that share at least k'_i attributes (the *Augment* subroutine in lines 6 and 7). Then, each group of tuples is checked only against its target set. For example, in line 8, the set $SS_1 \bowtie SN_2$ is checked only against $A_1 \bowtie R_2$ for a domination in k attributes. Note that the target set for $u \in SS_1$ is only A_1 . Similarly, lines 9 and 10 handle the sets $SN_1 \bowtie SS_2$ and $SN_2 \bowtie SN_1$ respectively.

The efficiency of the grouping algorithm stems from the fact that only the tuples in $SN_1 \bowtie SN_2$ need to be compared against the entire $R_1 \bowtie R_2$. For other tuples, the decision can be taken without even joining (the “yes” and “no” cases), or the comparison set is small (for the tuples in $SS_1 \bowtie SN_2$ and $SS_2 \bowtie SN_1$).

6.4 Dominator-Based Algorithm

The grouping algorithm has the problem that for tuples in SN_i , the target set is the entire relation R_i . The next algorithm, *dominator-based* algorithm (Algo. 3) rectifies this by *explicitly* storing the set of dominators. Thus, for a tuple $\in SS_i, SN_i$, the dominator set of tuples is first obtained (lines 7 and 11). When the tuple is in SS_i , this set is empty. The dominator sets are augmented by the tuples themselves and those tuples having the same values in the required number of skyline attributes (lines 8 and 12). In general, the size of the dominator sets is quite low as compared to the entire dataset, i.e., generally $|dom(\cdot)| \ll R_i$.

The target sets for the joined tuples are composed of the joins of the dominating tuples. Each tuple in the sets marked by “likely” and “may be” is checked against these joins of the corresponding dominator sets and is added to the answer only if no dominator

Algorithm 2 KSJQ: Grouping Algorithm

Input: Relations R_1, R_2 ; Number of attributes k

Output: k -dominant skyline set T

- 1: $k'_1 \leftarrow k - d_2$
- 2: $k'_2 \leftarrow k - d_1$
- 3: $SS_1, SN_1, NN_1 \leftarrow \text{Group}(R_1, k'_1)$
- 4: $SS_2, SN_2, NN_2 \leftarrow \text{Group}(R_2, k'_2)$
- 5: $T \leftarrow (SS_1 \bowtie SS_2)$ ▷ “yes” tuples
- 6: $A_1 \leftarrow \text{Augment}(SS_1, k'_1)$ ▷ augment $u \in SS_1$ with $\{u' : u'_{k'_1} = u_{k'_1}\}$
- 7: $A_2 \leftarrow \text{Augment}(SS_2, k'_2)$ ▷ augment $v \in SS_2$ with $\{v' : v'_{k'_2} = v_{k'_2}\}$
- 8: $T_1 \leftarrow \text{CheckTarget}(SS_1 \bowtie SN_2, A_1 \bowtie R_2, k)$
- 9: $T_2 \leftarrow \text{CheckTarget}(SN_1 \bowtie SS_2, R_1 \bowtie A_2, k)$
- 10: $T_3 \leftarrow \text{CheckTarget}(SN_1 \bowtie SN_2, R_1 \bowtie R_2, k)$
- 11: **return** $T \leftarrow T_1 \cup T_2 \cup T_3$

exists (line 16).

The joins of dominating sets are substantially less in size than the target sets for the grouping algorithm (which is the join of the entire target sets). The saving is largest for tuples in $SN_1 \bowtie SN_2$.

The saving, however, comes at a cost. For each tuple in SN , *all* the dominators need to be found out. In addition, the entire dominator set needs to be stored explicitly.

The advantages of the dominator-based algorithm may not be enough to offset this overhead of time and storage, especially when there are many such tuples. Sec. 7 compares the different algorithms empirically.

6.5 Cartesian Product

When the final relation is a *Cartesian product* of the two base relations, the algorithms become considerably easier. The Cartesian product can be considered as a special case of join with every tuple having the *same* value of the join attribute. In other words, all the tuples are in the same join group. As a result, there is no SN set. A tuple is either in SS (when it is a skyline in its local relation) or in NN (when it is not a skyline). Consequently, the tables become much simpler, and the fate of all the joined (i.e., final) tuples can be concluded without the need to explicitly compute them. The tuples in $SS_1 \bowtie SS_2$ are skylines while none of the other tuples are.

6.6 Non-Equality Join Condition

In certain cases, the join condition may *not* be an equality. For example, in a flight combination, the arrival time of the first leg needs to be earlier than the departure time of the second, i.e., $f_1.\text{arrival} <$

Algorithm 3 KSJQ: Dominator-Based Algorithm

Input: Relations R_1, R_2 ; Number of attributes k
Output: k -dominant skyline set T

- 1: $k'_1 \leftarrow k - d_2$
- 2: $k'_2 \leftarrow k - d_1$
- 3: $SS_1, SN_1, NN_1 \leftarrow \text{Group}(R_1, k'_1)$
- 4: $SS_2, SN_2, NN_2 \leftarrow \text{Group}(R_2, k'_2)$
- 5: $T \leftarrow (SS_1 \bowtie SS_2)$ ▷ “yes” tuples
- 6: **for each** $u \in SS_1, SN_1$ **do**
- 7: $\text{dom}(u) \leftarrow k'_1\text{-dominators}(u)$ ▷ dominators of u with k'_1 attributes
- 8: $\text{dom}(u) \leftarrow \text{dom}(u) \cup \text{Augment}(u, k'_1)$ ▷
 $\{u' : u'_{k'_1} = u_{k'_1}\}$
- 9: **end for**
- 10: **for each** $v \in SS_2, SN_2$ **do**
- 11: $\text{dom}(v) \leftarrow k'_2\text{-dominators}(v)$ ▷ dominators of v with k'_2 attributes
- 12: $\text{dom}(v) \leftarrow \text{dom}(v) \cup \text{Augment}(v, k'_2)$ ▷
 $\{v' : v'_{k'_2} = v_{k'_2}\}$
- 13: **end for**
- 14: $T \leftarrow \emptyset$
- 15: **for each** $u \bowtie v \in (SS_1 \bowtie SN_2) \cup (SN_1 \bowtie SS_2) \cup (SN_1 \bowtie SN_2)$ **do**
- 16: $T \leftarrow T \cup \text{CheckDominators}(\text{dom}(u) \bowtie \text{dom}(v), k)$
- 17: **end for**
- 18: **return** T

f_2 .departure. In this section, we discuss how to handle such join cases when the condition is one of $<, \leq, >, \geq$.

Since the main purpose of dividing a base relation into the three sets SS, SN and NN is to ensure that certain decisions can be taken about the tuples in these sets *without* joining, all the optimizations discussed in Sec. 5 work with the following modifications.

A tuple in $SS_1 \bowtie SS_2$ can never be k -dominated and, thus, it does not matter how such a tuple is composed from the base relations. In other words, the semantics of the join condition, equality or otherwise, does not matter for $SS_1 \bowtie SS_2$ tuples.

Next consider a tuple $t' = u' \bowtie v' \in (SN_1 \bowtie SS_2)$. A tuple in the SN set, u' , is originally defined as one that is not dominated by any other tuple in the *same* group. This ensures that if u' joins with v' from the other relation, then *no* other tuple u can join with v' to dominate t' . This is the crucial property that needs to be maintained even when the join condition is non-equality.

Thus, if the join condition is $u'.arr < v'.dep$, then the set $\{u : u.arr < u'.arr\}$ is considered to be in the *same* group of u' since it can also join with v' (and can potentially dominate u). In other words, this ensures that all such u is join compatible with v' . The set SN is thus expanded to take care of the non-equality condition. Note that there may exist other tuples $\{u'' : u''.arr < v'.dep\}$ that may also join with v' ; these, however, cannot be determined locally without the knowledge of v' from the other relation and, therefore, cannot be considered. The target set of a tuple in SN is anyway the entire dataset (or its dominators). Thus, the algorithms will work correctly with the above modification.

Similarly, for the converse set, $SS_1 \bowtie SN_2$, the SN set for v' consists of $\{v : v.dep > v'.dep\}$.

A joined tuple with NN as a component is rejected as a k -dominant skyline since the NN tuple can be always dominated by another tuple in the *same* group. Hence, similar to SN , the group of a tuple needs to be defined by taking into account the semantics of the join condition.

Considering the earlier example of $u'.arr < v'.dep$, if $u' \in NN$, then to say that $u' \in NN_1$, there must exist $u : u.arr < u'.arr$ and $u \succ_{k'_1} u'$. (The definition is suitably modified for NN_2

Algorithm 4 Finding k : Naïve Algorithm

Input: Number of skylines δ
Output: Number of attributes k

- 1: $k \leftarrow \max\{d_1, d_2\} + 1$ ▷ minimum k
- 2: **while** $k < d$ **do**
- 3: **if** $|\text{skyline}(k)| \geq \delta$ **then** ▷ actual number
- 4: **return** k ▷ return and terminate
- 5: **end if**
- 6: $k \leftarrow k + 1$
- 7: **end while**
- 8: **return** d ▷ maximum possible k

in a suitable manner.) This ensures that the joined tuple $t' = u' \bowtie v'$ will be dominated by $u \bowtie v'$. Once more, tuples of the form $\{u'' : u''.arr < v'.dep\}$ are left out due to lack of knowledge about v' .

It may happen that there does not exist any such u but there exists such an u'' . In that case, u' is classified as an SN tuple instead of an NN . This, however, only leads to extra processing of tuples joined with u' . The *correctness* is not violated as such joined tuples of the form $t' = u' \bowtie v'$ will be finally caught by $u'' \bowtie v'$ and rejected. Thus, the above modification only affects the efficiency of the algorithms, not the final result.

6.7 Algorithms for Aggregation

The algorithms that consider aggregation are essentially the same as the plain KSJQ. The only difference is that when the join is performed, the aggregation of the attributes are done as an additional step. Therefore, they are not discussed separately.

We next discuss the algorithms for finding k (Problem 3).

6.8 Finding k : Naïve Algorithm

The naïve algorithm (Algo. 4) to search for the lowest k that produces at least δ skylines starts from the least possible value of k , i.e., $\max\{d_1, d_2\} + 1$, and keeps incrementing it till the number of k -dominant skylines is at least δ . The largest possible value of k , i.e., d , is otherwise returned by default, even if it does not satisfy the δ criterion.

The algorithm is very inefficient as for each case, it computes the actual k -dominant skyline set. Further, it traverses the possibilities of k in a linear manner. The correctness is based on the fact that the number of k -dominant skylines is a monotonically non-decreasing function in k (Lemma 1).

We next design algorithms that use Table 5.

6.9 Finding k : Range-Based Algorithm

For a particular value of k , the actual number of k -dominant skylines, Δ_k , is *at least* the size of the “yes” sets, denoted by $\Delta_{k,lb}$, and *at most* the sum of sizes of the “yes”, “likely” and “may be” sets, denoted by $\Delta_{k,ub}$. These, thus, denote the lower and upper bounds respectively: $\Delta_{k,lb} < \Delta_k < \Delta_{k,ub}$.

The *range-based* algorithm (Algo. 5) uses these bounds to speed-up the process. Starting from the minimum possible $k = \max\{d_1, d_2\} + 1$, the algorithm finds $\Delta_{k,lb}$ and $\Delta_{k,ub}$ (lines 3 and 4 respectively). If $\Delta_{k,lb} \geq \delta$, then the current k is the answer (line 5). If $\Delta_{k,ub} < \delta$, then the current k cannot be the answer and k is incremented (line 7). Otherwise, i.e., if $\Delta_{k,lb} < \delta \leq \Delta_{k,ub}$, then the current k may be an answer. In this case, the *actual* k -dominant skyline set is computed. If its size is δ or greater, it is returned as the answer (line 9). Else, k is incremented (line 11), and the steps are repeated.

While this algorithm is definitely more efficient than the naïve one, it still suffers from the shortcoming that it examines a large number of k 's by incrementing it one by one. If the required k lies

Algorithm 5 Finding k : Range-Based Algorithm

Input: Number of skylines δ
Output: Number of attributes k

```
1:  $k \leftarrow \max\{d_1, d_2\} + 1$  ▷ minimum  $k$ 
2: while  $k < d$  do
3:    $\Delta_{k,lb} \leftarrow |\text{“yes” sets}|$ 
4:    $\Delta_{k,ub} \leftarrow |\text{“yes” sets}| + |\text{“likely” sets}| + |\text{“may be” sets}|$ 
5:   if  $\Delta_{k,lb} \geq \delta$  then ▷ lower bound
6:     return  $k$ 
7:   else if  $\Delta_{k,ub} < \delta$  then ▷ upper bound
8:      $k \leftarrow k + 1$ 
9:   else if  $|\text{skyline}(k)| \geq \delta$  then ▷ actual number
10:    return  $k$ 
11:   else
12:      $k \leftarrow k + 1$ 
13:   end if
14: end while
15: return  $d$  ▷ maximum possible  $k$ 
```

towards the higher end of the range, it unnecessarily examines too many lower values of k . The next algorithm does a *binary search* to reduce this overhead.

6.10 Finding k : Binary Search Algorithm

Algo. 6 shows how the binary search proceeds. It starts from the middle of the possible values of k (line 5). The values of $\Delta_{k,lb}$ and $\Delta_{k,ub}$ are computed in the same manner as earlier (lines 6 and 7 respectively).

If $\Delta_{k,lb} \geq \delta$ (line 8), then k is a potential answer. No value in the higher range can be the answer as the current k already satisfies the condition. However, there may be a lower k that satisfies the δ condition. Hence, the search is continued in the lower range to try and find a better (i.e., lesser) k .

If, on the other hand, $\Delta_{k,ub} < \delta$ (line 11), then the current estimate of k is too low. The search is, therefore, continued in the higher range.

If none of these bounds help, the actual number of k -dominant skylines, Δ_k , is found. If $\Delta_k \geq \delta$ (line 13), then the current k is a potential answer. However, a lesser k may be found and, so, the search proceeds to the lower range.

Otherwise, i.e., when $\Delta_k < \delta$ (line 16), the required k is searched in the higher range.

The algorithm stops when the lower range of the search becomes larger than or equal to the current answer (line 19), which then is the lowest k that satisfies the δ condition. It may also stop when the range is exhausted, in which case, the current value of k is returned (line 23).

The binary search based algorithm, thus, speeds up the searching through the possible range of values of k . Sec. 7 compares the three algorithms empirically.

7. EXPERIMENTAL RESULTS

We experimented with data synthetically generated using <http://randdataset.projects.pgfoundry.org/> on an Intel i7-4770 @3.40 GHz Octacore machine with 16 GB RAM using code written in Java.

We also experimented with a real dataset of two-legged flights from New Delhi to Mumbai. The details are in Sec. 7.4.

We measured the effects of various parameters on the different algorithms proposed in Sec. 6. The parameters and their default values are listed in Table 7. Note that the size of the joined relation is a derived parameter. It is equal to n^2/g for two base relations with n tuples and g groups. When the effect of a particular set of parameters are measured, the rest are held to their default values,

Algorithm 6 Finding k : Binary Search Algorithm

Input: Number of skylines δ
Output: Number of attributes k

```
1:  $l \leftarrow \max\{d_1, d_2\} + 1$  ▷ minimum  $k$ 
2:  $h \leftarrow d$  ▷ maximum  $k$ 
3:  $cur \leftarrow d$  ▷ current estimate of  $k$ 
4: while  $l < h$  do
5:    $k \leftarrow \lfloor (l + h) / 2 \rfloor$ 
6:    $\Delta_{k,lb} \leftarrow |\text{“yes” sets}|$ 
7:    $\Delta_{k,ub} \leftarrow |\text{“yes” sets}| + |\text{“likely” sets}| + |\text{“may be” sets}|$ 
8:   if  $\Delta_{k,lb} \geq \delta$  then
9:      $cur \leftarrow k$  ▷ update current estimate
10:     $h \leftarrow k - 1$  ▷ search for a lower  $k$ 
11:   else if  $\Delta_{k,ub} < \delta$  then ▷ search for a higher  $k$ 
12:      $l \leftarrow k + 1$ 
13:   else if  $|\text{skyline}(k)| \geq \delta$  then
14:      $cur \leftarrow k$  ▷ update current estimate
15:      $h \leftarrow k - 1$  ▷ search for a lower  $k$ 
16:   else if  $|\text{skyline}(k)| < \delta$  then
17:      $l \leftarrow k + 1$  ▷ search for a higher  $k$ 
18:   end if
19:   if  $l \geq cur$  then ▷ lowest  $k$  already found
20:     return  $cur$ 
21:   end if
22: end while
23: return  $cur$ 
```

Symbol	Parameter	Default value
n	Dataset size for base relation	3, 300
d	Dimensionality of base relation	7
k	Number of skyline attributes	11
a	Number of aggregate attributes	2
g	Number of join groups	10
T	Dataset type	Independent
δ	Threshold of skyline size	10, 000
N	Size of joined relation	1, 089, 000

Table 7: Parameters for experiments.

unless explicitly stated otherwise.

In the figures, the three main algorithms for KSQL are denoted as: G for grouping, D for dominator-based, and N for naïve. The overall running time for each algorithm is divided into various components: (i) time taken for computing the groups in the base relations, i.e., SS , SN , and NN , (ii) time taken for actually joining the tuples from the two base relations that cannot be pruned, (iii) time taken for finding the dominators of the tuples, and (iv) the rest of the processing. These are marked separately in the figures.

Not all the components are applicable to every algorithm, e.g., the naïve algorithm does not find groups. The components that are not applicable to an algorithm are shown to have zero costs.

The three algorithms for determining the value of k are depicted as: B for binary search, R for range-based, and N for naïve.

7.1 Aggregate

We first show the results where aggregate values have been used. The aggregation function used is *sum*.

7.1.1 Effect of Dimensionality

The first experiment measures the effect of varying k . Fig. 1 shows that the running time increases sharply with k . The two different settings of dimensionality, d , and number of aggregate attributes, a , (Fig. 1a and Fig. 1b) show the robustness of this behavior. As k increases, it becomes increasingly hard to dominate a tuple in k dimensions. As a result, the number of k -dominant skyline increases heavily, thereby resulting in increased running times.

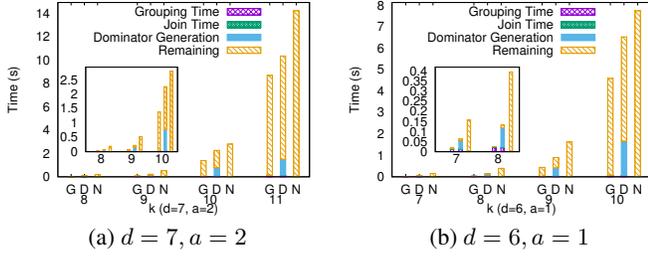


Figure 1: Effect of k .

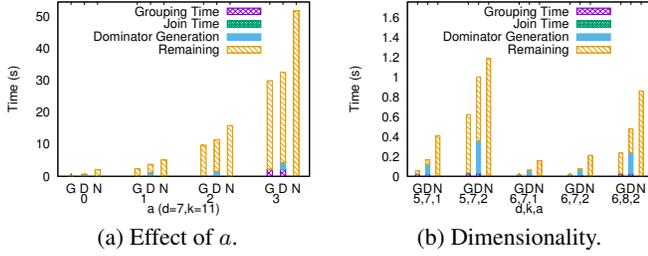


Figure 2: Effect of dimensionality.

Overall, the grouping algorithm is the fastest. The dominator-based algorithm spends a substantial amount of time in finding the dominators for each tuple. This overhead is not compensated enough in the final checking stage. This is due to the fact that the average dominator set sizes are quite large and, hence, joining large dominator sets requires a substantial amount of time. With increasing dimensionality, the time required to find the dominator sets increases as well.

As expected, the naïve algorithm performs the worst as it does not attempt any optimization at all. It is slower than the grouping algorithm by about 1.5-2 times.

The next set of experiments hold d and k constant but increases the number of aggregate attributes, a . Fig. 2a depict the results. Note that $a = 0$ signifies that no aggregate attributes are used. The trend of the results remain the same with the running time increasing with a . Once more, grouping is the best algorithm followed by dominator-based and naïve.

Fig. 2b shows a medley of results across different d , k and a . When a or k increase, the running time increases as well.

However, it seems that the reverse happens with d . Comparing the case of $d = 5, k = 7, a = 1$ with $d = 6, k = 7, a = 1$, we note that in the first case, $k'_1 = k'_2 = 3$ while in the second case, $k'_1 = k'_2 = 2$. Thus, in the second case, it is easier to find the groups and perform the joins. Consequently, it runs faster. The same reasoning holds true for $d = 5, k = 7, a = 2$ against $d = 6, k = 7, a = 2$.

We next compare $d = 5, k = 7, a = 2$ against $d = 6, k = 8, a = 2$. The values of $k'_1 = k'_2 = 4$ are same in both the cases. However, the size of dominator sets is larger in the first case. The time required to divide the base relations into the three sets is also higher. This leads to an overall higher running time.

7.1.2 Effect of Number of Join Groups

Fig. 3a shows the effect of number of join groups. Note that when $g = 1$, the join reduces to a Cartesian product. This case can be handled specially as explained in Sec. 6.5. When g is low, the number of k -dominant skylines is low since there are more chances of a tuple being k -dominated by another one in the same group. Thus, there are less number of SN tuples (none at $g = 1$) and their dominators. On the other hand, when g is high, the size of

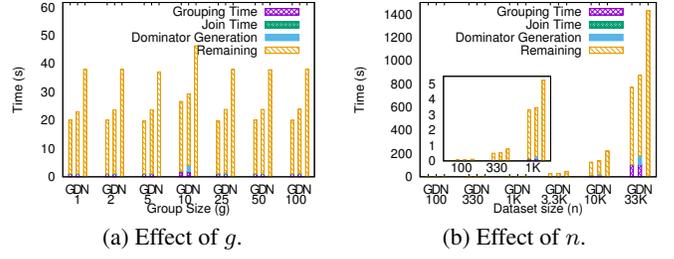


Figure 3: Scalability.

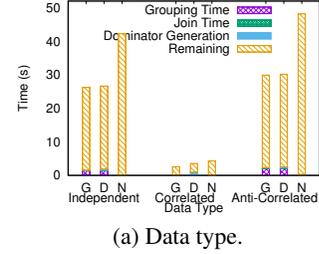


Figure 4: Type of data distribution.

the joined relation (which is n^2/g) decreases. Therefore, there are two opposing effects on the running time. Empirically, the running times are the highest at medium values.

7.1.3 Effect of Dataset Size

The next experiment varies the size of the base relations, n . Note that with increase in n , the size of the joined relation increases quadratically ($O(n^2)$). Consequently, as visible in Fig. 3b, the running time increases drastically. The scalability of the grouping algorithm, in particular, as well as the dominator-based algorithm, is sub-linear in the size of the joined relation, though.

7.1.4 Effect of Type of Data Distribution

The final set of experiments on the aggregated attributes measures the effect of type of data distribution. Fig. 4 shows that correlated datasets are the easiest to process due to higher chances of domination of a tuple by another tuple, thereby resulting in lesser number of skylines. The anti-correlated datasets are the most time consuming due to the opposing effect. The independent datasets are mid-way.

7.2 No Aggregation

The next set of experiments target the scenarios where no aggregation over the skyline attributes is done.

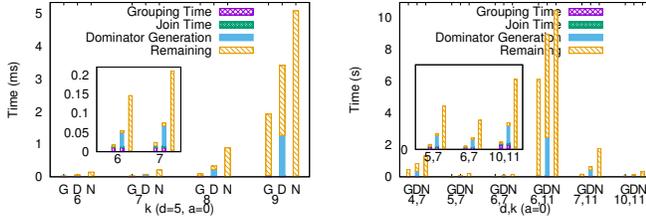
7.2.1 Effect of Dimensionality

Fig. 5a shows the effect of k when $d = 5$. Note that since $a = 0$, the possible values of k range from $d + 1 = 6$ to $2d - 1 = 9$.

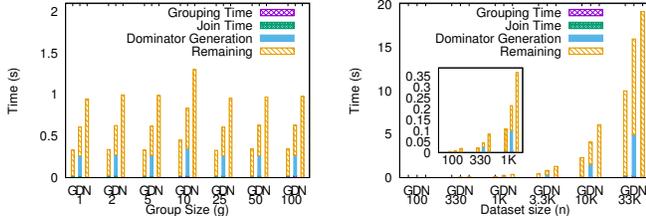
Similar to the case with aggregate attributes, the running time increases sharply with k . The grouping algorithm performs the best while the naïve is the worst. The dominator-based algorithm suffers since the time spent in finding the dominators is too large and is not sufficiently compensated later. Interestingly, since the join time is constant for the naïve algorithm irrespective of the value of k , the proportion of time spent in joining is much higher for lower k .

Fig. 5b holds k constant and varies d over two settings. When k is fixed and d increases, the values of k'_1 and k'_2 decrease. This results in faster grouping time. The dominator sets are also computed faster. Thus, the overall time decreases.

7.2.2 Effect of Number of Join Groups



(a) Effect of k . (b) Effect of d .
Figure 5: Effect of dimensionality (no aggregation).



(a) Effect of k . (b) Effect of d .
Figure 6: Scalability (no aggregation).

As explained earlier in Sec. 7.1.2, the number of join groups has two opposing effects on the running time. Therefore, as shown in Fig. 6a, the results are similar. Note that the values of d and k here are $d = 4$ and $k = 7$.

7.2.3 Effect of Dataset Size

Fig. 6b shows that the running time increases drastically with n , although the scalability is sub-linear in the size of the joined relation, i.e., n^2 . The largest dataset size we tested was for $n = 33,000$, leading to a massive size of over 10^8 for the joined relation. The fact that the grouping algorithm produces the result in less than 20s for this case establishes the practicality of the algorithms.

7.2.4 Effect of Type of Data Distribution

The effect of type of data distribution (Fig. 7) is similar to that in Sec. 7.1.4 with the anti-correlated requiring the largest amount of time and correlated the least.

7.3 Finding k

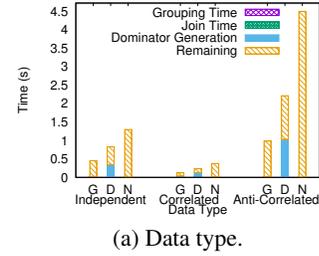
The third and final set of experiments deals with finding the value of k given a threshold δ of requisite number of k -dominant skylines. Aggregate attributes are not used.

7.3.1 Effect of Threshold δ

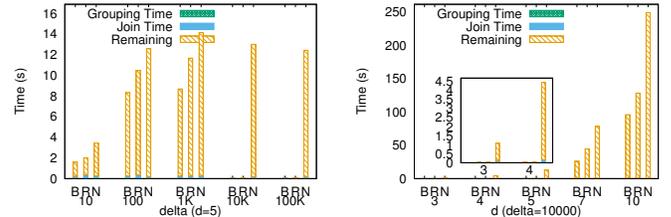
Fig. 8a shows the effect of varying δ values. The dimensionality is held fixed at $d = 5$ with $a = 0$. The possible values of k , therefore, range from $d + 1 = 6$ to $2d = 10$. The size of each base relation is $n = 3,300$ with $g = 10$, thereby resulting in more than 10^6 joined tuples (as listed in Table 7).

With increasing δ , the naïve algorithm requires larger running times since it keeps iterating over k . The range-based search also iterates over the values of k , although it avoids the costly full k -dominant skyline computation in the intermediate steps, if possible. When δ is very large, it falls outside the upper bounds for most values of k and, hence, the algorithm runs very fast. In this experiment, when $\delta \geq 10,000$, the largest possible $k = 10$ is returned as the answer.

For low values of δ , the iterations of both naïve and range-based



(a) Data type.
Figure 7: Type of data distribution (no aggregation).



(a) Effect of δ . (b) Effect of d .
Figure 8: Effect of dimensionality (finding k).

algorithms stop early. The answer for $\delta = 10$ is $k = 8$ while that for $\delta = 100$ and $\delta = 1,000$ are both $k = 9$.

The binary search method is faster even for medium values of δ since it avoids the iterative procedure and quickly finds the desired value of k . Overall, it is always the fastest algorithm.

7.3.2 Effect of Dimensionality d

Fig. 8b shows the complementary effect, that of varying dimensionality d while keeping δ constant at 10,000. When d is low, the chosen δ value is too large and the algorithms terminate fast. When d is increased, the algorithms need to search through a larger range and, therefore, takes a much longer time. The binary search method is consistently the fastest algorithm. It outperforms the range-based algorithm by about 1.2-1.5 times while the naïve algorithm is slower by a factor of 2-2.5.

7.3.3 Effect of Number of Join Groups

There is no appreciable effect of the number of join groups on the algorithms for finding k (Fig. 9a).

7.3.4 Effect of Dataset Size

Fig. 9b shows the effect of increasing dataset size, n . The dimensionality and threshold values are kept fixed at $d = 5$ and $\delta = 1,000$ with $g = 10$.

For very low values of n (up to 1,000), the threshold is too high, and the maximum possible $k = 10$ is required to satisfy the threshold δ .

With increasing n , the running time increases due to increasing k -dominant skyline computation times. Even for larger n , the values of k are towards the higher end. Therefore, the binary search algorithm is the most suitable one for finding k .

7.3.5 Effect of Data Type

The effect of type of data distribution is as expected with correlated being the fastest and anti-correlated the slowest (Fig. 10a).

7.4 Real Dataset

To test our algorithms on real data, we collected information about various attributes on domestic flights in India from www.makemytrip.com. The first base table contained information about

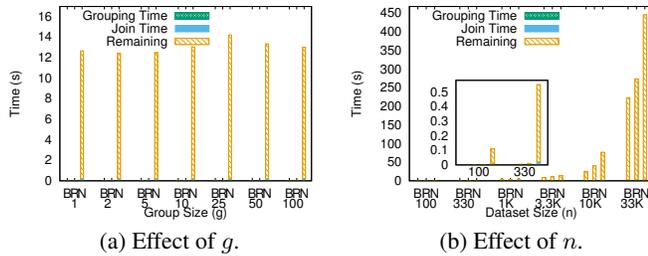
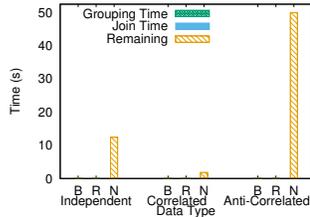


Figure 9: Scalability (finding k).



(a) Data type.

Figure 10: Type of data distribution (finding k).

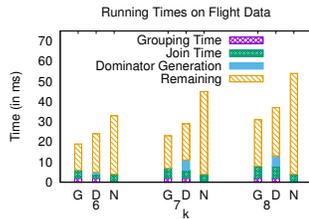


Figure 11: Real data.

192 flights from New Delhi to 13 important cities of India while the second base table contained information about 155 flights from those 13 cities to Mumbai. For each base table, 5 attributes were considered: cost, flying time, date change fee, popularity, and amenities. The first 2 attributes were aggregated while the other 3 were used as local attributes. The join was based on the equality of the intermediate city. Thus, each tuple in the joined relation contained $3 + 3 + 2 = 8$ attributes. The size of the joined relation was 2649.

We ran experiments on $k = 6, 7, 8$. The results are summarized in Fig. 11. The grouping-based algorithm performed the best followed by the dominator-based method and the naïve algorithm. All the results were produced in milliseconds highlighting the practicality of the methods.

7.5 Summary of Experiments

The experiments show that the grouping algorithm consistently outperforms the other methods in solving the KSJQ queries. For finding the right value of k , the binary search algorithm turns out to be the best method always.

8. CONCLUSIONS

In this paper, we proposed a novel query, k -dominant skyline join query (KSJQ), that incorporates finding k -dominant skylines over joined relations where the attributes may be aggregated as well. We analyzed certain optimizations for the query and used them to design efficient algorithms. In addition, given the number of final skylines sought, we also proposed efficient algorithms to find the right value of k .

In future, we would like to extend the algorithms to work in parallel, distributed and probabilistic settings.

9. REFERENCES

- [1] A. Awasthi, A. Bhattacharya, S. Gupta, and U. K. Singh. k -dominant skyline join queries: Extending the join paradigm to k -dominant skylines. In *ICDE*, page to appear, 2017.
- [2] A. Bhattacharya and B. P. Teja. Aggregate skyline join queries: Skylines with aggregate operations over multiple relations. In *COMAD*, pages 15–26, 2010.
- [3] S. Börzsönyi, D. Kossmann, and K. Stocker. The skyline operator. In *ICDE*, pages 421–430, 2001.
- [4] C.-Y. Chan, H. V. Jagadish, K.-L. Tan, A. K. H. Tung, and Z. Zhang. Finding k -dominant skylines in high dimensional space. In *SIGMOD*, pages 503–514, 2006.
- [5] J. Chomicki, P. Godfrey, J. Gryz, and D. Liang. Skyline with presorting. In *ICDE*, pages 717–719, 2003.
- [6] A. Das Sarma, A. Lall, D. Nanongkai, R. J. Lipton, and J. Xu. Representative skylines using threshold-based preference distributions. In *ICDE*, pages 387–398, 2011.
- [7] L. G. Dong and X. W. Cui. Finding k -dominant skyline for combined dataset. *Applied Mechanics and Materials*, 568:1534–1538, 2014.
- [8] Y. Gao, J. Hu, G. Chen, and C. Chen. Finding the most desirable skyline objects. In *DASFAA*, pages 116–122, 2010.
- [9] P. Godfrey. Skyline cardinality for relational processing. In *FOIKS*, pages 78–97, 2004.
- [10] Z. Huang, Y. Xiang, and Z. Lin. l -SkyDiv query: Effectively improve the usefulness of skylines. *Sc. China Information Sciences*, 53(9):1785–1799, 2010.
- [11] H.-K. Hwang, T.-S. Tsai, and W.-M. Chen. Threshold phenomena in k -dominant skylines of random samples. *SIAM J. Computing*, 42(2):405–441, 2013.
- [12] W. Jin, J. Han, and M. Ester. Mining thick skylines over large databases. In *PKDD*, pages 255–266, 2004.
- [13] V. Koltun and C. H. Papadimitriou. Approximately dominating representatives. *Theor. Comput. Sci.*, 371(3):148–154, 2007.
- [14] D. Kossmann, F. Ramsk, and S. Rost. Shooting stars in the sky: an online algorithm for skyline queries. In *VLDB*, pages 275–286, 2002.
- [15] H. T. Kung, F. Luccio, and F. P. Preparata. On finding the maxima of a set of vectors. *J. ACM*, 22(4):469–476, 1975.
- [16] X. Lin, Y. Yuan, Q. Zhang, and Y. Zhang. Selecting stars: The k most representative skyline operator. In *ICDE*, pages 86–95, 2007.
- [17] D. Papadias, Y. Tao, G. Fu, and B. Seeger. An optimal and progressive algorithm for skyline queries. In *SIGMOD*, pages 467–478, 2003.
- [18] M. Siddique and Y. Morimoto. Extended k -dominant skyline in high dimensional space. In *ICISA*, pages 1–8, 2010.
- [19] M. A. Siddique and Y. Morimoto. Efficient computation for k -dominant skyline queries with domination power index. In *ICIT*, pages 382–387, 2010.
- [20] M. A. Siddique and Y. Morimoto. Efficient maintenance of k -dominant skyline for frequently updated database. In *DBKDA*, pages 107–110, 2010.
- [21] D. Sun, S. Wu, J. Li, and A. K. H. Tung. Skyline-join in distributed databases. In *ICDE Workshop*, pages 176–181, 2008.
- [22] K.-L. Tan, P.-K. Eng, and B. C. Ooi. Efficient progressive skyline computation. In *VLDB*, pages 301–310, 2001.
- [23] Y. Tao, L. Ding, X. Lin, and J. Pei. Distance-based representative skyline. In *ICDE*, pages 892–903, 2009.
- [24] H. Tian, M. A. Siddique, and Y. Morimoto. An efficient processing of k -dominant skyline query in mapreduce. In *Data4U*, pages 29–34, 2014.
- [25] G. Valkanas, A. N. Papadopoulos, and D. Gunopulos. SkyDiver: a framework for skyline diversification. In *EDBT*, pages 406–417, 2013.
- [26] T. Xia, D. Zhang, and Y. Tao. On skylining with flexible dominance relation. In *ICDE*, pages 1397–1399, 2008.
- [27] Z. Zhang, X. Guo, H. Lu, A. K. H. Tung, and N. Wang. Discovering strong skyline points in high dimensional spaces. In *CIKM*, pages 247–248, 2005.
- [28] Z. Zhang, Y. Yang, R. Cai, D. Papadias, and A. Tung. Kernel-based skyline cardinality estimation. In *SIGMOD*, pages 509–521, 2009.