

Marrying Top-k with Skyline Queries: Operators with Relaxed Preference Input and Controllable Output Size

The two most common paradigms to identify records of preference in a multi-objective setting rely either on dominance (e.g., the skyline operator) or on a utility function defined over the records' attributes (typically, using a top- k query). Despite their proliferation, each of them has its own palpable drawbacks. Motivated by these drawbacks, we identify three hard requirements for practical decision support, namely, personalization, controllable output size, and flexibility in preference specification. With these requirements as a guide, we combine elements from both paradigms and propose two new operators, ORD and ORU. We perform a qualitative study to demonstrate how they work, and evaluate their performance against adaptations of previous work that mimic their output.

CCS Concepts: • **Information systems** → *Top-k retrieval in databases*.

Additional Key Words and Phrases: Top- k query; Skyline; Multi-dimensional datasets

ACM Reference Format:

. 2023. Marrying Top-k with Skyline Queries: Operators with Relaxed Preference Input and Controllable Output Size. *ACM Trans. Datab. Syst.* 37, 4, Article 111 (August 2023), 33 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

In the era of ubiquitous access to the Internet, users are presented with numerous alternatives to cover their everyday needs. Choosing from the available alternatives generally entails the consideration of multiple, often conflicting aspects. Indeed, multi-objective optimization has been a traditional research topic [30, 45, 65], whose practical relevance has increased in the current, fully connected reality. For a large set of alternatives (i.e., d -dimensional records), there are two main paradigms to determine those of most interest to the user, namely, based on *dominance* or *ranking by utility*.

The first paradigm considers that a record *dominates* another if all its attributes are more preferable. The *skyline* includes the records that are not dominated [15], while the *k-skyband* those dominated by a maximum of $(k - 1)$ others [59]. The dominance paradigm is intuitive to the user. On the downside, it has two major shortcomings: (i) it is not personalizable, reporting the same result for every user, and (ii) its output size (i.e., the number of reported records) is uncontrollable, and often overwhelming [13, 31].

Personalization (i.e., serving the specific preferences of an individual user) is a hard requirement for decision support, especially nowadays, when large amounts of personal information are available via smartphones, fitness trackers, online activities, etc. Regarding the output size, Hick's law, known since the 50's, suggests that controlling the number of results presented to the user is essential to the quality of the decision and to the user experience [35, 38]. That law has been used as a cornerstone in eCommerce applications, meta-search engines, etc [33, 36]. Dictating the output size

Author's address:

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Association for Computing Machinery.

0362-5915/2023/8-ART111 \$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

is crucial also because of design considerations, such as display size, device capabilities, connection speed, etc. Hence, another hard requirement is for *output-size specified* (OSS) operators.

The second paradigm, ranking by utility, associates each record with a score via a (user-specific) function over the records' attributes. Most commonly, the utility function is a weighted sum, with user preferences expressed by d per-attribute weights w_i (together comprising a preference vector \mathbf{w}). This linear type of scoring has been the most proliferate since the inception of ranking by utility [24, 43], and is shown by user studies to effectively model the way humans assess tradeoffs in real-life multi-objective decisions [61].

Ranking by utility, in the form of a top- k query, is both personalized and OSS. Its Achilles' heel, however, lies in specifying the "correct" weights, since a small change in \mathbf{w} can drastically alter the top- k result [40, 81]. Vector \mathbf{w} is assumed to be either input directly by the user or somehow mined (e.g., via online behavior and review mining [42, 72], pairwise comparisons of example records [41, 61], or some other preference learning technique [25]). In the former case, a user cannot be reasonably expected to quantify with absolute precision the relative importance of the various attributes. In the latter, preference learning methods come with an understanding that the mined \mathbf{w} is only an estimate. This shortcoming motivates a third hard requirement for practical decision support, which is *relaxed preference input*, i.e., some flexibility in the specified preferences.

In this paper, we aim to bring together the strong points of both paradigms (dominance-based and ranking by utility), while avoiding their drawbacks. In particular, we propose two operators that uphold all the three hard requirements we have established, namely, (i) being OSS, (ii) being personalized, and (iii) having a relaxed preference input. To achieve personalization, we employ linear scoring, due to its demonstrated effectiveness in modeling human decision making [61]. However, we consider the input \mathbf{w} as a best-effort estimate. We therefore relax it, by incrementally expanding it equally on all directions in the preference domain. Conceptually, at the original \mathbf{w} , this corresponds to ranking by utility (e.g., a top- k query at \mathbf{w}). As the expansion radius grows, it gradually shifts towards standard dominance, including in the output additional records that cater to alternative preferences, similar to \mathbf{w} . The stopping radius is indirectly (yet strictly) determined by the desired output size m .

Research on both standard paradigms has considered their individual weaknesses, but no existing work satisfies all three hard requirements. The skyline literature includes formulations that control the output size by loosening the definition of dominance [16, 46, 73], identifying representatives [34, 47, 50, 66], or considering subspaces [17, 70]. For example, Lin et al. [50] report the m skyline members that dominate the most non-skyline records, while Chan et al. [17] shortlist the m records that belong to the most subspace skylines. These definitions aim to produce the most competitive or the most representative skyline records in a general sense, without a specific user in mind, thus lacking personalization.

Centered more on utility, studies on regret-minimizing sets report m representative records from the dataset. Typically, they define the regret ratio as the relative difference between the utility of the top-scoring record in the selected subset and the top-scorer in the entire dataset. Their objective is to minimize the aggregate (usually, the maximum) regret ratio across every possible utility function [58, 75], i.e., the reported subset is meant to satisfy as it best can all possible users, without an intent for personalization.

Two recent studies, [21] and [55], attempt to relax the preference input in ranking by utility. That is, they assume that the preference input is a convex polytope R instead of a vector \mathbf{w} . Concordantly, they report the records that could be among the k most preferable for any $\mathbf{w} \in R$ (for $k = 1$ and $k \geq 1$ in [21] and [55], respectively). Unfortunately, these methods are not OSS and, worse yet, come without even an estimate of the output size, i.e., the user/application is in the dark on whether R is too large or too small to produce, even approximately, the required number of records. Furthermore,

Operator	Personalized	OSS	Flexible Input
Skyline/ k -Skyband	✗	✗	✓
Top- k	✓	✓	✗
OSS skylines	✗	✓	✓
Regret-minimizing sets	✗	✓	✓
Fixed-region techniques	✓	✗	✓
Proposed (ORD and ORU)	✓	✓	✓

Table 1. Multi-objective queries and their properties

these approaches may remove the need for a particular \mathbf{w} , but require specifying a polytope R in the preference domain. Deciding R is left to the user or application, a choice that becomes tougher considering that the dynamics in the preference domain are hard to gauge. In usability terms, specifying the output size m is arguably more tangible and more relatable to the user/application than specifying a polytope in the preference domain.

Our operators, ORD and ORU, satisfy all three hard requirements. They expand the preference input \mathbf{w} in a similar way, however, they retain a stronger flavor of either paradigm each. ORD employs an adaptive notion of dominance that is guided by m , while ORU sticks closer to ranking by utility. Hard requirements aside, practicality also demands responsiveness and scalability. We make geometric observations and establish propositions that lead to efficient processing. Our algorithms are orders of magnitude faster than adaptations of previous work, which can merely simulate the ORD/ORU output, and still, without OSS guarantees.

In Table 1, we summarize the properties of existing multi-objective queries, and juxtapose them with our operators. In Section 2, we offer a more comprehensive description of related work.

2 RELATED WORK

The two traditional alternatives to determine the most preferable records from a dataset D with d attributes, are based on dominance and on ranking by utility score. A record *dominates* another if it is at least as preferable in all dimensions, and strictly more preferable in at least one dimension. The records that are not dominated by any other comprise the *skyline* [15] while, more generally, those dominated by fewer than k form the *k -skyband* [59]. In contrast, in the ranking approach, the score of a record is typically defined as the weighted sum of its attributes for a vector of d user-specific weights. The top- k set includes the k records with the largest scores [39].

For large, indexed datasets, the most common processing algorithms in both cases follow the *branch-and-bound* methodology. BBS [59] visits index nodes and data records in increasing distance from the top corner of the data space (i.e., the corner with the maximum possible attribute values), using a min-heap to organize them by that distance. It maintains as skyline/ k -skyband the records dominated by none/fewer than k records encountered so far. BBR [67] computes the top- k set by visiting index nodes and data records in decreasing (upper bound of) score, using a max-heap. The first k records popped from the heap are the top- k .

OSS Skylines: The size of the skyline is uncontrollable and oftentimes very large [31]. That being a major shortcoming, there have been several approaches to limit it.

Chan et al. [16] consider that a record \mathbf{r}_i *m -dominates* another \mathbf{r}_j for an $m \leq d$ if there is a subspace of m dimensions where \mathbf{r}_i dominates \mathbf{r}_j . A smaller m implies a smaller skyline, thus controlling its size. Koltun and Papadimitriou [46] propose ϵ -dominance, where the attributes of a record \mathbf{r}_i are multiplied by $(1 + \epsilon)$ to check whether it dominates another record \mathbf{r}_j . In the same spirit, Xia et al. [73] increment the attributes of \mathbf{r}_i by an absolute δ value on all (appropriately

scaled) dimensions. Other studies aim to select the m most representative skyline records. The *dominance count* of a record, i.e., the number of records it dominates, has been used as a measure of its importance [29, 69, 78]. By that intuition, Lin et al. [50] choose the m skyline records that dominate the most other records. Lee and Hwang [47] propose a pivot-based space partitioning for that problem, while Gao et al. [26] enhance it by favoring representatives that dominate the less frequently dominated non-skyline records. The latter's performance is improved by Han et al. [34]. By a different, distance-based intuition, Tao et al. [66] choose as representatives the m skyline records that minimize the distance from the remaining skyline members.

Sarma et al. [62] pick m records from the skyline, so that the probability that a random user would click on one of them is maximized. Assuming that a record r_i is interesting if its attributes exceed a certain threshold per dimension, and that the distribution of the threshold values is known, they propose approximate and sampling methods to select the m representatives. Magnani et al. [51] consider various measures of diversity and significance, and assume a linear combination of these two factors as the objective function that the m chosen skyline records must maximize.

Another approach considers membership in *subspace* skylines [60, 68]. Chan et al. [17] report the m skyline records that appear in the most subspace skylines. Vlachou and Vazirgiannis [70] measure importance according to dominance in different subspaces, and assume propagation of importance via dominance links. Another formulation considers that some attributes are more important [44, 53]; Lee et al. [48] select representatives according to skyline membership in the induced prioritized subspaces.

Most OSS skylines do not take into account a user's personal preferences. An exception, in abstract terms at least, are Bartolini et al. [11], who consider that record attributes correspond to user-specific ratings. If a user has not provided ratings for records r_i and r_j , but at least a fraction of similar users have indicated ratings where r_i dominates r_j , the same is assumed for the user at hand too. The required fraction indirectly controls the skyline size. The focus in [11] is to infer dominance when user ratings (i.e., record attributes) are missing. In contrast, in our target applications the records' attributes are given and no information for other users is required. Another distinction between OSS skylines and our work is that they consider $k = 1$, i.e., once dominated, a record is eliminated. Instead, our operators may dig deeper, to larger k values, because they can rely on the personal preferences (roughly) specified by w .

Regret Minimization: Work on *regret-minimizing sets* (RMS) produces an m -sized subset $S \subset D$ that tries to satisfy as it best can any possible user. In the original formulation [58], the regret ratio for a user is defined as the relative difference between the maximum utility score in S and that in the entire D . The objective for S is to minimize the maximum regret ratio for any possible user. There have been many follow-up studies (e.g., [8, 76]), considering also RMS variants, most notably k -RMS [19] (where the regret ratio reflects the difference between the top-scorer in S and the top- k -th in D), minimizing the *average* regret ratio [80], defining regret based on the rank of records [9], etc. A survey is given in [75]. RMS studies are not concerned with personalization. Also, even if fed with our operators' stopping radius, they cannot reproduce our output. For example, to solve classic RMS [58], it suffices to consider only skyline records. In contrast, our output may also include records below the skyline.

Inspired by RMS, but aiming for personalization, *interactive regret minimization* (IRM) involves the user in the search process [57]. Initially oblivious of her preferences, IRM goes through multiple rounds of interaction. In each round, it presents her with a number of records and asks her to choose the best, thus learning her (latent) preference vector increasingly well. When the regret ratio is guaranteed to be small enough (or the actual top-scorer is found), the last record chosen becomes the answer for this user. The original IRM method [57] involves artificial records in its

interactions, which is resolved in [74]. The latter is enhanced in [82] by asking the user to sort the presented records (instead of just choosing the best). IRM assumes a different query processing model altogether, requiring active user involvement. Moreover, its objective is to eventually identify the one record with maximum utility, and thus considers only records on the skyline or convex hull.

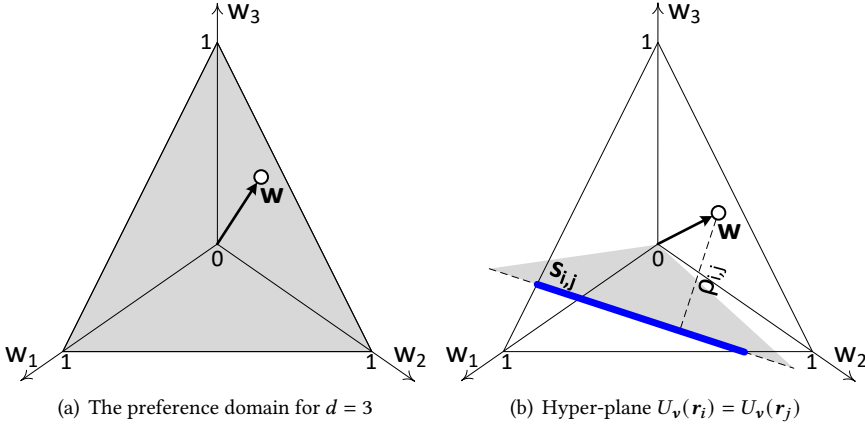
Fixed-region Techniques: The closest related studies to our work are [21] and [55]. Given a convex preference polytope R , Ciaccia and Martinenghi [21] define that \mathbf{r}_i R -dominates \mathbf{r}_j if \mathbf{r}_i scores higher than \mathbf{r}_j for any $\mathbf{w} \in R$. They propose an R -dominance test which checks one linear condition per extreme vertex of R , and compute the R -skyline (i.e., the records that are not R -dominated by any other) by integrating that test into standard skyline algorithms. They also introduce an operator that reports as *potentially optimal* every \mathbf{r}_i that is the top record for at least one $\mathbf{w} \in R$. To check a record for potential optimality, they solve a linear programming (LP) problem defined according to the extreme vertices of R . Mouratidis and Tang [55] extend potential optimality to $k \geq 1$, i.e., they identify the records that appear in the top- k result for at least one $\mathbf{w} \in R$. In a more advanced variant, they explicitly report every possible (order-insensitive) top- k set for any $\mathbf{w} \in R$. They first disqualify records R -dominated by k or more others. Among the remaining candidates, they determine the top- k -th in each partition of R and, accordingly, the (order-insensitive) prefix of the top- k set.

In terms of practicality, the operators in [21] and [55] lack the OSS property, meaning that the user/application cannot determine (or even predict) the size of the output. The techniques themselves cannot be extended to our problem, because they rely on R being fixed and given in advance. For example, their R -dominance and LP tests are defined according to the extreme vertices of R , and are contingent on these vertices being fixed and known. Furthermore, they require R to be a convex polytope. In contrast, in our case R is not specified, and the preference region (even if it were given in advance) is effectively a hyper-sphere, not a polytope. If we approximate hyper-spheres with hyper-cubes and make repetitive calls for different side-lengths of R in an exploratory manner, the approaches in [21] (for $k = 1$) or [55] (for general k) could somehow simulate our operators, but even with that slack, they would require an excessive number of trials/executions to produce an output of exactly m records. In other words, a second compromise is necessary, i.e., allow them to terminate when the output size is “almost” m (e.g., within a 10% deviation). Our framework not only produces output of the exact desired size (strictly OSS), but it also reports *order-sensitive* top- k results anywhere within its stopping radius ρ .

Related Top- k Work: On the top- k front, there are studies for unspecific or unknown preference vector \mathbf{w} that are somewhat related to our work. For example, Soliman et al. [64] compute the most probable top- k result if \mathbf{w} is a random, uniformly distributed vector. Uncertain records/attributes have also been considered, leading to probabilistic top- k outputs [7, 22, 77]. On the other hand, Zhang et al. [81] compute the preference region that corresponds to a given top- k result, in a task which, loosely speaking, is inverse from ours.

3 PROBLEM FORMULATION

We consider that the available options are represented as d -dimensional records $\mathbf{r} = \langle x_1, x_2, \dots, x_d \rangle$ in a dataset D indexed by a spatial access method, e.g., an R-tree [12, 63]. We make the convention that the larger the attributes the better, yet our findings adapt easily to cases where some/all attributes are to be minimized. Given a preference vector \mathbf{v} of d non-negative weights w_i , the utility score of a record \mathbf{r} is defined as their inner product, i.e., $U_{\mathbf{v}}(\mathbf{r}) = \sum_{i=1}^d x_i \cdot w_i$. Accordingly, the top- k result comprises the k records with the highest scores. Ordering D by utility is independent from the magnitude of \mathbf{v} [37, 49], thus we assume preference vectors where $\sum_{i=1}^d w_i = 1$. In other words,

Fig. 1. Preference domain, and minidist $\rho_{i,j}$ example

the domain of the preference vectors, called *preference domain*, is the unit $(d - 1)$ -simplex in a space whose d axes correspond to the w_i values, i.e., the simplex $\Delta^{d-1} = \{\mathbf{v} \in \mathbb{R}_+^d \mid \sum_{i=1}^d w_i = 1\}$. For $d = 3$, the preference domain is an equilateral triangle, shown in gray in Figure 1(a). Effectively, any valid preference vector is represented as a vertex in that triangle. For $d = 4$, the preference domain is a tetrahedron, and so on.

Let \mathbf{w} be a best-effort estimate of the user's preference vector, henceforth called the *seed*, and consider the preference vectors \mathbf{v} within distance ρ from \mathbf{w} , i.e., where $|\mathbf{v} - \mathbf{w}| \leq \rho$. If a record \mathbf{r}_i scores at least as high as another \mathbf{r}_j for every such vector \mathbf{v} , and strictly higher for at least one of them, we say that \mathbf{r}_i ρ -dominates \mathbf{r}_j . The records that are ρ -dominated by fewer than k others form the ρ -skyband. This general notion includes the ρ -skyline as a special case for $k = 1$. Note that a larger ρ implies a larger ρ -skyband. In the extreme settings, $\rho = 0$ renders the ρ -skyband equivalent to a traditional top- k query at \mathbf{w} , while $\rho = \infty$ makes it equivalent to the standard k -skyband¹. We may now define our first operator, abbreviated as ORD to stress its QSS property, relaxed input, and stronger dominance-oriented flavor.

DEFINITION 1 (ORD). *Given the seed vector \mathbf{w} and the required output size m , ORD reports the records that are ρ -dominated by fewer than k others, for the minimum ρ that produces exactly m records in the output.*

Observe that user and application are both transparent to ρ , which relieves them from being concerned with the complex dynamics of the preference domain. The appropriate ρ is determined automatically by our framework, according to the desired output size m . Our second operator shares that trait too, but follows more closely the ranking by utility paradigm, thus its abbreviation, ORU.

DEFINITION 2 (ORU). *Given the seed vector \mathbf{w} and the required output size m , ORU reports the records that belong to the top- k result for at least one preference vector within distance ρ from \mathbf{w} , for the minimum ρ that produces exactly m records in the output.*

¹If \mathbf{r}_i scores no lower than \mathbf{r}_j for every vector in the preference domain, and higher for at least one vector therein, then (and only then) \mathbf{r}_i dominates \mathbf{r}_j [20].

While beyond the requirements of Definition 2, a byproduct of our ORU algorithm is the reporting of the specific (order-sensitive) top- k result for any vector within radius ρ from \mathbf{w} . This enables additional applications, like determining the most stable [81] or the most representative [64] top- k results in the vicinity of \mathbf{w} , according to the volume of the preference regions that produce them.

Our ORD/ORU techniques require no precomputation other than a general-purpose spatial index on D . This implies that updates in D affect only (and are readily supported by) the index. Also, it enables the integration of common predicates into our framework. For example, should the user impose arbitrary range predicates (e.g., price between \$150 and \$200, size between 400ft² and 600ft², etc), we may execute a multi-dimensional range query on D , followed by ORD/ORU in the selected part of the index/dataset.

Multi-objective querying generally loses its meaning in high dimensions. For instance, for more than a handful of dimensions almost every record tends to belong to the skyline [17, 31], while utility-wise the scores of all records tend to converge [56, 79]. We, hence, focus on low-dimensional settings. A final remark is that although we position our work within preference-based record shortlisting for a human user, our techniques apply to general multi-objective scenarios where the suitability of available options is defined by a linear function over the options' attributes.

4 OSS DOMINANCE-BASED OPERATOR

The output of ORD is a ρ -skyband, and in particular the one for the smallest ρ that includes m records. We make several observations that lead to an efficient ORD processing methodology.

4.1 Observations and Main Idea

Properties of Candidate Records: Without loss of generality, assume that no two records coincide or score the same for the seed vector \mathbf{w} . Unless a record belongs to the traditional k -skyband, it cannot belong to the top- k result for any preference vector [15]. Hence, for any ρ , the ρ -skyband is a subset of the k -skyband. Therefore, the latter includes all the candidates we may need for ORD. Consider a record \mathbf{r}_i among them. The remaining candidates fall into three categories regarding their potential to ρ -dominate \mathbf{r}_i :

- Records that score lower than \mathbf{r}_i for the seed vector \mathbf{w} cannot ρ -dominate it for any radius ρ (because any ρ includes the seed itself). By the way, a corollary of this is that the top- k records of \mathbf{w} belong to the ρ -skyband for every ρ .
- Records that dominate \mathbf{r}_i in the traditional sense, score higher for any preference vector, thus they ρ -dominate \mathbf{r}_i for any ρ .
- The remaining records (i.e., those that do not dominate \mathbf{r}_i but score higher than it for \mathbf{w}) ρ -dominate \mathbf{r}_i for a non-empty range of ρ values, as we explain next.

Consider a record \mathbf{r}_j that falls in the third category. As such, it does not dominate \mathbf{r}_i . Also, since $U_{\mathbf{w}}(\mathbf{r}_j) > U_{\mathbf{w}}(\mathbf{r}_i)$, record \mathbf{r}_j is not dominated by \mathbf{r}_i either. Every pair of records that do not dominate each other define a hyper-plane with equation $U_{\mathbf{v}}(\mathbf{r}_i) = U_{\mathbf{v}}(\mathbf{r}_j)$ that cuts through the preference domain, i.e., it divides Δ^{d-1} into two non-empty parts. Since \mathbf{r}_j scores higher than \mathbf{r}_i for the seed \mathbf{w} , it holds that $U_{\mathbf{v}}(\mathbf{r}_i) < U_{\mathbf{v}}(\mathbf{r}_j)$ in the entire part that includes \mathbf{w} , while $U_{\mathbf{v}}(\mathbf{r}_i) > U_{\mathbf{v}}(\mathbf{r}_j)$ in the other part. Assuming $d = 3$, Figure 1(b) illustrates in gray a hyper-plane with equation $U_{\mathbf{v}}(\mathbf{r}_i) = U_{\mathbf{v}}(\mathbf{r}_j)$.

Let $\mathbf{s}_{i,j}$ be the intersection of hyper-plane $U_{\mathbf{v}}(\mathbf{r}_i) = U_{\mathbf{v}}(\mathbf{r}_j)$ with the preference domain Δ^{d-1} . Geometrically, $\mathbf{s}_{i,j}$ is a $(d - 2)$ -simplex, e.g., for $d = 3$ it is a line segment, shown bold in Figure 1(b). Let also $\rho_{i,j}$ be the minimum distance between \mathbf{w} and any $\mathbf{v} \in \mathbf{s}_{i,j}$. For any preference vector in Δ^{d-1} that is within distance $\rho_{i,j}$ from \mathbf{w} , record \mathbf{r}_j scores higher than \mathbf{r}_i , i.e., \mathbf{r}_j ρ -dominates \mathbf{r}_i for every $\rho \leq \rho_{i,j}$. In contrast, it does not ρ -dominate \mathbf{r}_i for $\rho > \rho_{i,j}$. In implementation terms, we can compute the mindist $\rho_{i,j}$ using a quadratic programming solver [32, 54] with (squared) distance as

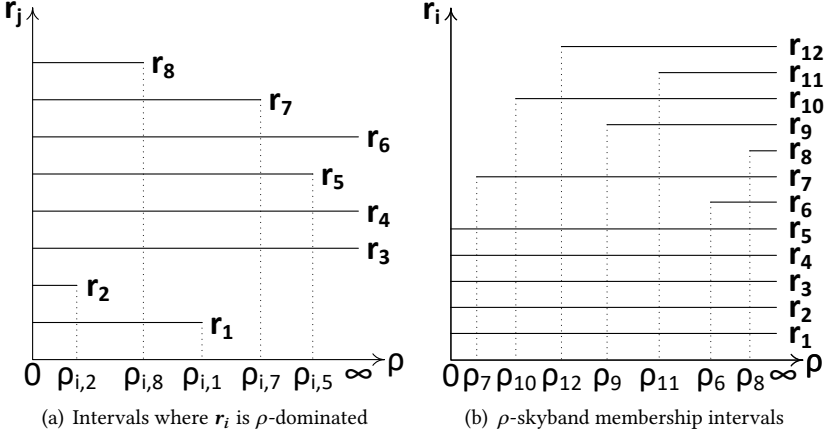


Fig. 2. Determining the ORD output

the minimization objective, subject to the linear constraints that define $s_{i,j}$, i.e., $U_v(r_i) = U_v(r_j)$ and $\sum_{i=1}^d w_i = 1$.

Inflection Radius: By considering all records r_j in the third category against r_i , they are each mapped into an interval of ρ values where they ρ -dominate it. Figure 2(a) offers an example. Assume that $k = 5$ and that the 5-skyband includes 8 records that score higher than r_i for w . Out of these, r_i is dominated in the traditional sense by 3 (i.e., r_3 , r_4 , r_6), thus their infinite intervals. The remaining 5 do not dominate r_i , hence, they fall in the third category; they each have a finite mindist $\rho_{i,j}$, mapped to intervals as illustrated. By sweeping the intervals from left to right, we can easily identify the ρ value past which r_i is dominated by fewer than k others, i.e., it becomes part of the ρ -skyband. We call that value the *inflection radius* of r_i and denote it as ρ_i . In our example, $\rho_i = \rho_{i,7}$.

A Preliminary Approach: Based on the above, a first-cut ORD solution is to compute the entire k -skyband, and for each record in it, to derive the inflection radius. Then, to output the m of them with the smallest inflection radii. To exemplify, assume that the 5-skyband includes 12 records; in Figure 2(b), we map each of them to an interval of ρ values where it belongs to the ρ -skyband, according to its inflection radius. These intervals have different meaning from Figure 2(a). In Figure 2(a), all intervals refer to r_i , helping to compute its own inflection radius ρ_i . In contrast, Figure 2(b) is a global representation of the ρ -skyband for different ρ values. Specifically, if we sweep the chart with a vertical line, the intervals that intersect the line at any position, indicate the ρ -skyband members for the ρ value that corresponds to that position. In our example, assuming that $m = 8$, the ORD output is set $\{r_1, r_2, r_3, r_4, r_5, r_7, r_{10}, r_{12}\}$, which corresponds to the ρ -skyband for $\rho = \rho_{12}$.

An interesting insight is that the ORD output may vary from standard ranking by utility (top- k) all the way to traditional dominance-based querying (k -skyband), depending on m . On the one hand, by definition, the top- k records are the only members of the ρ -skyband for $\rho = 0$, which corresponds to the smallest possible m (i.e., $m = k$). On the other hand, every k -skyband member will appear to the ρ -skyband for a sufficiently large ρ or, equivalently, for a sufficiently large m . To visualize these extremes, at the leftmost position in Figure 2(b) the sweeping line intersects only

the intervals of the top- k records (\mathbf{r}_1 to \mathbf{r}_5), while at the rightmost² the entire k -skyband. That said, although this generality is welcome, the practical strength of ORD is for m values in between these extremes.

4.2 Efficient ORD Processing

The ORD processing idea described so far is a foundation that offers an abstract-level understanding of our method. However, an efficient solution must address several performance issues. Primarily, we would want to avoid computing the entire k -skyband in the beginning of the process. Indeed, the k -skyband may include numerous records, many times more than the m required [13, 31]. Ideally, we want to limit the number of considered candidates to as tight a superset of the ORD output as possible. The algorithm we present next serves that objective.

We first invoke a progressive k -skyband retrieval that fetches its members one by one, and place them into a candidate set. Importantly, unlike standard k -skyband computation, we enforce that its members are fetched in decreasing score order for \mathbf{w} (we will explain how shortly). This retrieval order is essential, because when a new candidate \mathbf{r}_i is fetched, we can definitively compute its inflection radius ρ_i already, without having to derive the entire k -skyband. The rationale is that only k -skyband records with higher score may ρ -dominate \mathbf{r}_i , and these are guaranteed to be fetched before it.

We keep fetching new k -skyband members in that fashion, until the candidate set reaches size $(m + 1)$. At that stage, we evict the candidate with the largest inflection radius. Also, an important algorithmic shift takes place. Let $\bar{\rho}$ be the maximum inflection radius of the remaining m candidates. The $\bar{\rho}$ -skyband is guaranteed to include at least the m existing candidates, thus $\bar{\rho}$ upper bounds the eventual stopping radius of the algorithm. Therefore, from this point onwards, we switch to fetching $\bar{\rho}$ -skyband members (instead of k -skyband members), still in decreasing order of score for \mathbf{w} . The switch can be performed transparently, as we elaborate later.

To exemplify, consider Figure 2(b). Assume that $k = 5$, $m = 11$, and that we have fetched the $m + 1 = 12$ depicted records, in the order indicated by their subscripts, i.e., \mathbf{r}_i was fetched i -th. To bring the candidates down to $m = 11$, we discard \mathbf{r}_8 , as it has the largest inflection radius. Practically, that sets $\bar{\rho}$ to ρ_8 .

As new candidates are fetched, and evictions are made to keep their total number to m , the current $\bar{\rho}$ keeps shrinking. Meanwhile, as $\bar{\rho}$ shrinks, records tend to ρ -dominate more others. This implies that the $\bar{\rho}$ -skyband retrieval becomes increasingly more selective, thus filtering out more aggressively regular k -skyband members that cannot participate in the ORD result. When the $\bar{\rho}$ -skyband module cannot fetch any more records, the candidate set is finalized as the ORD result. The latter corresponds to the ρ -skyband for ρ equal to the maximum inflection radius across its members.

The ORD algorithm relies on a progressive k -skyband module, with the extra requirement to fetch records in decreasing score according to \mathbf{w} . We use an adaptation of BBS [59] where we visit index nodes and records in decreasing (upper bound of) score for \mathbf{w} , using a max-heap. Once the $(m + 1)$ -th record is fetched, we shift to $\bar{\rho}$ -skyband computation using the exact same heap as per normal, but replace the regular dominance tests of BBS with $\bar{\rho}$ -dominance, for the current $\bar{\rho}$ value. The visiting order by score and the use of $\bar{\rho}$ -dominance tests instead of regular dominance are permissible modifications to vanilla BBS, because its correctness is guaranteed as long as no record \mathbf{r}_j fetched after another \mathbf{r}_i may dominate \mathbf{r}_i [59]. That property is upheld by our visiting order (by score), both initially for regular dominance, and after the shift to $\bar{\rho}$ -dominance. Indeed, $U_{\mathbf{w}}(\mathbf{r}_i) > U_{\mathbf{w}}(\mathbf{r}_j)$

²The largest meaningful ρ (although visualized as ∞) is the distance between \mathbf{w} and its furthest point in Δ^{d-1} , since that ρ already covers the entire preference domain.

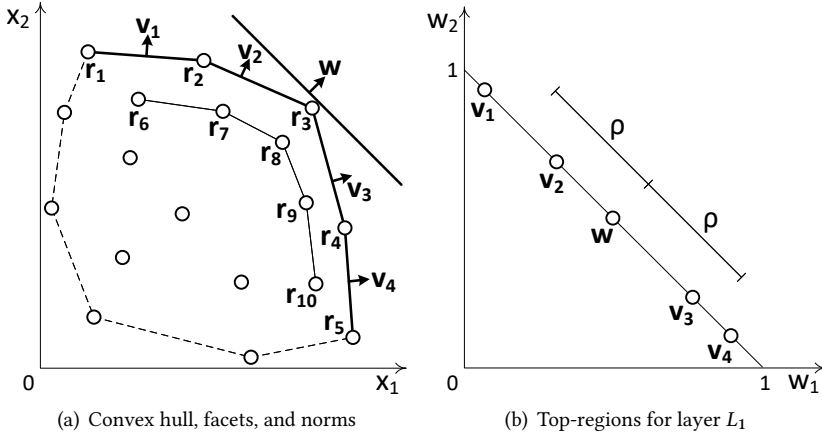


Fig. 3. Fundamental notions and principles

ensures that r_j cannot dominate nor ρ -dominate r_i for any ρ . An implementation note on the adapted BBS regards its ρ -dominance building block. That block tests whether an already-fetched $\bar{\rho}$ -skyband record r_i $\bar{\rho}$ -dominates a not-yet-fetched record r_j (or an unvisited index node whose top corner is r_j). The test is performed as explained in Section 4.1, by comparing the mindist $\rho_{i,j}$ with $\bar{\rho}$.

5 OSS UTILITY-BASED OPERATOR

Our second operator, ORU, adheres more closely to ranking by utility; it reports records that belong to the top- k for at least one preference vector within radius ρ from the seed w , for the minimum ρ that produces exactly m records. Despite the seeming similarity to ORD's definition, the top- k ranking involved in ORU renders it innately different, and its solution considerably more complex.

We present important preliminaries (in Section 5.1), a crucial theorem and algorithmic basis to process ORU (in Section 5.2), and eventually a complete implementation (in Section 5.3).

5.1 Fundamentals

The abstractions and techniques used for ORU have the notion of the *convex hull* at their core [14]. The convex hull of D is the smallest convex polytope that encloses all its records. It comprises facets, each defined by d extreme vertices (records) in general position. The outer polygon in Figure 3(a) is the convex hull of an example dataset. Facet $r_1 r_2$ is defined by extreme vertices r_1 and r_2 , etc.

A vector is *normal* to a hyper-plane when its direction is perpendicular to the hyper-plane. The *norm* of a facet on the hull is the normal vector to that facet whose sum of coordinates is 1, and is directed towards the exterior of the hull. In our example, the norm of $r_1 r_2$ is vector v_1 . Effectively, the norm of a facet can be seen as a point in the preference domain Δ^{d-1} .

The top record for a preference vector v is the one met first by a hyper-plane normal to v that sweeps the data space from the top corner to the origin [23, 28]. Hence, the top record in D is guaranteed to lie on its convex hull [18, 52]. Since in our case the weights are non-negative, the top record is among the extreme vertices of facets with non-negative norms. We call *upper hull* the part that corresponds to these facets. In Figure 3(a), the upper hull is bold (and the rest of the convex hull is dashed). For the shown w , the top record is r_3 , as it is met first by the sweeping line normal to w .

To explain the fundamentals of our methodology, **assume** that we have already computed the first k upper hull layers: the first layer, L_1 , includes the upper hull of D ; the second, L_2 , includes the upper hull of $D - L_1$; and, generally, layer L_i the upper hull of D after subtracting layers L_1 to L_{i-1} . In Figure 3(a), layer L_1 includes records r_1 to r_5 , and L_2 records r_6 to r_{10} . Note that in reality our complete ORU algorithm *does not require such precomputation*, but instead it builds *on the fly* (i.e., at query time) *only parts* of the necessary layers, thus applying to arbitrary k , avoiding precomputation costs (time and space), and extending transparently to dynamic datasets, i.e., to cases where record insertions/deletions may occur in D and invalidate precomputed information. Also, **assume** that we already know the necessary radius ρ for ORU to produce m records. Of course, this too is an assumption we will drop later (in Section 5.3).

Adjacent Set $\mathcal{A}(\mathbf{r})$: Consider a record \mathbf{r} in layer L_i . We denote by $\mathcal{F}(\mathbf{r})$ the set of L_i facets with \mathbf{r} as one of their extreme vertices, and by $\mathcal{A}(\mathbf{r})$ the records adjacent to \mathbf{r} , i.e., the L_i records (other than \mathbf{r}) that define facets in $\mathcal{F}(\mathbf{r})$. In Figure 3(a), for example, $\mathcal{F}(\mathbf{r}_3) = \{r_2r_3, r_3r_4\}$ and $\mathcal{A}(\mathbf{r}_3) = \{r_2, r_4\}$. The following Lemmas 1, 2, and 3 are crucial. Note that they refer to records within the same layer L_i .

LEMMA 1. *Given a preference vector \mathbf{v} whose top record in L_i is \mathbf{r} , if we start shifting \mathbf{v} towards any direction in the preference domain, the first record in L_i to outscore \mathbf{r} is always in $\mathcal{A}(\mathbf{r})$, i.e., among the records adjacent to \mathbf{r} . Furthermore, each of the records in $\mathcal{A}(\mathbf{r})$ is the first outscoring record for some shifting direction of \mathbf{v} .*

PROOF. Let H_r be the hyper-plane in the data space that is normal to \mathbf{v} and passes through \mathbf{r} . Record \mathbf{r} is the top for \mathbf{v} , as long as there is no record *above* H_r (i.e., in the half-space that includes the top corner of the data space). Assume that \mathbf{v} gradually shifts towards a specific direction, with H_r always passing through \mathbf{r} . As the orientation of H_r shifts together with \mathbf{v} , the first record in L_i that is met by H_r is the record \mathbf{r}_i that will outscore \mathbf{r} if we shift \mathbf{v} infinitesimally any further (in the same direction). Suppose that \mathbf{r}_i is not in $\mathcal{A}(\mathbf{r})$, i.e., it shares no common L_i facet with \mathbf{r} . At the time that H_r touches \mathbf{r}_i , according to the hypothesis, no other record in L_i should lie above H_r . This, however, is a contradiction, because the convexity of L_i implies that any hyper-plane that passes through two non-adjacent records in L_i (\mathbf{r} and \mathbf{r}_i , in this case) cuts through the interior of L_i , i.e., there is at least one other extreme vertex (record) in L_i that lies above H_r . We conclude that the first record to outscore \mathbf{r} , for any direction of shifting \mathbf{v} , must be in $\mathcal{A}(\mathbf{r})$.

It remains to show that for each record \mathbf{r}_i in $\mathcal{A}(\mathbf{r})$, there is a direction of shifting \mathbf{v} that makes \mathbf{r}_i the first outscoring record. Let f be a facet in $\mathcal{F}(\mathbf{r})$ where \mathbf{r}_i is a defining vertex. Consider the shifting of \mathbf{v} towards the norm of f , equivalently, the shifting of H_r until it falls on f . Since f is a facet of the convex hull, it leaves all L_i records towards its interior. Thus, there is no L_i record above H_r at all times until now. After H_r has fallen on f , if \mathbf{v} shifts infinitesimally towards \mathbf{r}_i , \mathbf{r}_i will become the first to outscore \mathbf{r} . \square

By Lemma 1, if \mathbf{w} shifts clockwise/anticlockwise in Figure 3(a), \mathbf{r}_4 and \mathbf{r}_2 , respectively, will be the first L_1 records to outscore \mathbf{r}_3 .

Top-region $C(\mathbf{r})$: Building on Lemma 1, our next proposition reveals an important property within L_i , and helps define the *top-region* of a record $\mathbf{r} \in L_i$, i.e., the region $C(\mathbf{r})$ in the preference domain where every vector has \mathbf{r} as its top record in L_i .

LEMMA 2. *Let \mathbf{r} be a record in L_i . \mathbf{r} is the top-scorer across all L_i records for those preference vectors \mathbf{v} that fall in the convex polytope $C(\mathbf{r})$ defined by (i.e., whose extreme vertices correspond to) the norms of the facets in $\mathcal{F}(\mathbf{r})$.*

PROOF. From Lemma 1, we infer that $C(\mathbf{r})$ is determined by records in $\mathcal{A}(\mathbf{r})$, since they are the first to outscore \mathbf{r} once \mathbf{v} leaves $C(\mathbf{r})$. In particular, each adjacent record \mathbf{r}_i corresponds to a half-space $U_{\mathbf{v}}(\mathbf{r}) \geq U_{\mathbf{v}}(\mathbf{r}_i)$ in the preference domain (simply expressing that \mathbf{r} should score no lower than \mathbf{r}_i anywhere in $C(\mathbf{r})$). $C(\mathbf{r})$ is the intersection of all these half-spaces, which (by definition [14]) is a convex polytope. Each facet of $C(\mathbf{r})$ is attributed to one of the intersected half-spaces, say, $U_{\mathbf{v}}(\mathbf{r}) \geq U_{\mathbf{v}}(\mathbf{r}_i)$ and, in effect, to one of the adjacent records, i.e., \mathbf{r}_i in this case. In general position, every extreme vertex of $C(\mathbf{r})$, say \mathbf{v}_j , corresponds to the intersection of $(d - 1)$ of its facets, i.e., to $(d - 1)$ equalities of the form $U_{\mathbf{v}}(\mathbf{r}) = U_{\mathbf{v}}(\mathbf{r}_i)$, where each \mathbf{r}_i is adjacent to \mathbf{r} . Let S be the record set composed of \mathbf{r} and these specific $(d - 1)$ adjacent records. As all records in S have the same score according to \mathbf{v}_j , by Lemma 1, such a tie is only feasible if any pair of records in S are adjacent to each other on L_i . Since, in general position, each facet on L_i is defined by d records, the records in S define a facet f in $\mathcal{F}(\mathbf{r})$, with \mathbf{v}_j as its norm. In other words, there is a direct one-to-one mapping between the facets in $\mathcal{F}(\mathbf{r})$ and the extreme vertices of $C(\mathbf{r})$. \square

By Lemma 2, $C(\mathbf{r})$ can be seen as a dual representation of $\mathcal{F}(\mathbf{r})$, where the former refers to the preference domain and the latter to the data space. Consider \mathbf{r}_3 in Figure 3(a). Facet set $\mathcal{F}(\mathbf{r}_3) = \{\mathbf{r}_2\mathbf{r}_3, \mathbf{r}_3\mathbf{r}_4\}$ translates to the top-region defined by their norms \mathbf{v}_2 and \mathbf{v}_3 , i.e., $C(\mathbf{r}_3)$ is segment $\mathbf{v}_2\mathbf{v}_3$ in the preference domain. Note that for $d = 2$, the preference domain Δ^{d-1} is a line segment.

Order Continuity: Lemmas 1 and 2, in tandem, suggest a continuity in the score order among L_i records for every \mathbf{v} . Specifically, the different top-regions for any given layer L define a partitioning of the preference domain, with adjacent records $\mathbf{r}_i, \mathbf{r}_j$ in L having neighboring top-regions $C(\mathbf{r}_i), C(\mathbf{r}_j)$ in Δ^{d-1} . Considering layer L_1 in our running example, Figure 3(b) demonstrates the partitioning of the preference domain. The top-region of \mathbf{r}_1 is the segment from point $(0, 1)$ to \mathbf{v}_1 ; of \mathbf{r}_2 from \mathbf{v}_1 to \mathbf{v}_2 ; of \mathbf{r}_3 from \mathbf{v}_2 to \mathbf{v}_3 , etc. Lemma 3 establishes a property for every vector in a top-region.

LEMMA 3. *For any preference vector $\mathbf{v} \in C(\mathbf{r})$, the top-2-nd record in L_i is always in $\mathcal{A}(\mathbf{r})$, i.e., among the records adjacent to \mathbf{r} .*

PROOF. Let $H_{\mathbf{v}}$ be the hyper-plane (in data space) that is normal to \mathbf{v} . Sweeping the data space with $H_{\mathbf{v}}$, the first encountered record in L_i is, by definition, \mathbf{r} . As sweeping continues further, the convexity of L_i ensures that $H_{\mathbf{v}}$ cuts only through the facets in $\mathcal{F}(\mathbf{r})$. Since L_i is hollow (i.e., has no records in its interior), the L_i record to be encountered next (i.e., the top-2-nd in L_i) must be an extreme vertex of a facet in $\mathcal{F}(\mathbf{r})$, i.e., a record in $\mathcal{A}(\mathbf{r})$. \square

In our example, Lemma 3 implies that the top-2-nd record in layer L_1 for any $\mathbf{v} \in C(\mathbf{r}_3)$ is either \mathbf{r}_2 or \mathbf{r}_4 . Note that all three lemmas consider a layer in isolation. For instance, although $\mathbf{w} \in C(\mathbf{r}_3)$, its top-2-nd record in the entire D is none of \mathbf{r}_2 or \mathbf{r}_4 , but \mathbf{r}_8 from L_2 .

5.2 An Algorithmic Basis for ORU

In this section, we prove an important theorem and provide an algorithmic basis to process ORU. Recall that we assumed we already know the minimum radius ρ required to produce m records, and that the first k upper hull layers are precomputed. We do not drop these assumptions yet. Here we focus on determining the top- k result for any possible preference vector within radius ρ from the seed \mathbf{w} , in order to form the ORU output.

First, we find all records in layer L_1 whose top-region has mindist to \mathbf{w} no greater than ρ . Let C be one of these regions. We already know the top record in it, say, \mathbf{r} . Considering C in isolation, our next task is to determine the top-2-nd record anywhere in it (i.e., for any possible preference vector

$\mathbf{v} \in C$), and to partition C accordingly. By Lemma 3, if we only considered L_1 , the top-2-nd record for any $\mathbf{v} \in C$ would be among those adjacent to \mathbf{r} . On the other hand, in the remaining dataset (i.e., if we ignored the records in L_1), the top-2-nd record would be in L_2 and more specifically, by Lemma 2, among the L_2 records \mathbf{r}_i whose top-region $C(\mathbf{r}_i)$ overlaps C . Theorem 1 generalizes this key observation.

THEOREM 1. *Assume that anywhere in a preference region C the (order-sensitive) top- i result is the same and it is known. Also, let L_t be the deepest layer that any of the top- i records belongs to. The top- $(i + 1)$ -th record anywhere in C must be in the union of:*

- *Set (i): The adjacent records to any member of the known top- i result in its respective layer, and*
- *Set (ii): The records in the $(t + 1)$ -th layer (i.e., L_{t+1}) whose top-region overlaps C .*

PROOF. Let S be the union of all records in the first t layers. Due to Lemma 3, when it is applied to each of the top- i records in their respective layer, if we only considered S , the top- $(i + 1)$ -th record would be in Set (i). On the other hand, if we only considered the rest of the dataset, i.e., $D - S$, the next highest-scoring record would be in L_{t+1} and specifically, by Lemma 2, in Set (ii). Hence, in the overall product set D (i.e., in the union of S and $D - S$), the top- $(i + 1)$ -th record anywhere in C must be in the union of Sets (i) and (ii). \square

Returning to our processing description for region C , the top record (the order-sensitive top- i result, in the general case) is already known and fixed anywhere in it, thus we can readily determine Set (i). We can also extract from L_2 (from L_{t+1} , in the general case) the part of the upper hull that corresponds to records in Set (ii); let us denote that part as L_{prt} . We update the upper hull L_{prt} to also cover Set (i) records, and denote its updated version as L_{upd} . Next, we apply Lemma 2 to L_{upd} to identify the top-2-nd records (the top- $(i + 1)$ -th, in the general case) for any $\mathbf{v} \in C$, and we partition C accordingly. We continue this process recursively in each produced partition until the full, order-sensitive top- k result is known anywhere in C . Repeating that process for all L_1 top-regions with mindist up to ρ , we derive all the required top- k results.

In our example, assume that ρ is as shown in Figure 3(b). The L_1 top-regions with mindist from \mathbf{w} up to ρ correspond to \mathbf{r}_2 , \mathbf{r}_3 , and \mathbf{r}_4 . Focusing on $C(\mathbf{r}_3)$ (i.e., segment $\mathbf{v}_2\mathbf{v}_3$), we know already that the top record is \mathbf{r}_3 , and seek to find the top-2-nd. Set (i) includes \mathbf{r}_2 and \mathbf{r}_4 . To determine Set (ii), we refer to L_2 . Figure 4(a) illustrates the L_2 top-regions. Among them, those that overlap $C(\mathbf{r}_3)$ are $C(\mathbf{r}_7)$, $C(\mathbf{r}_8)$, and $C(\mathbf{r}_9)$. Thus, we form L_{prt} as the part of L_2 that corresponds to \mathbf{r}_7 , \mathbf{r}_8 , \mathbf{r}_9 . Updating L_{prt} to also cover Set (i) (i.e., \mathbf{r}_2 , \mathbf{r}_4), results in the upper hull L_{upd} in Figure 4(b). L_{upd} suggests that the top-2-nd record is one of \mathbf{r}_2 , \mathbf{r}_8 , \mathbf{r}_4 . Furthermore, by Lemma 2, their L_{upd} top-regions (determined by facet norms \mathbf{v}_a and \mathbf{v}_b) help partition $C(\mathbf{r}_3)$ according to which exactly among them is the top-2-nd. Figure 4(c) presents the induced partitioning. The top-2 result is $\{\mathbf{r}_3, \mathbf{r}_2\}$ for preference vectors in $\mathbf{v}_2\mathbf{v}_a$; $\{\mathbf{r}_3, \mathbf{r}_8\}$ in $\mathbf{v}_a\mathbf{v}_b$; and $\{\mathbf{r}_3, \mathbf{r}_4\}$ in $\mathbf{v}_b\mathbf{v}_3$. The process repeats recursively in order to determine the top-3-rd record in each of the three partitions, and so on.

Observe that we may not need to reach as deep as the k -th layer, since members of Set (i) could prevent those of Set (ii) to enter the result, thus giving more “width” to the ORU search (in the data space) than “depth”. For instance, in Figure 4(b), it could be the case that none of the L_2 records belongs to L_{upd} , i.e., that the top-2-nd record comes from L_1 (i.e., \mathbf{r}_2 or \mathbf{r}_4) for any $\mathbf{v} \in C(\mathbf{r}_3)$.

5.3 Dropping Assumptions; Complete ORU

In this section, we describe our complete ORU algorithm. So far, we have made two impractical assumptions, i.e., that we have already computed the first k upper hull layers, and that we know in advance the necessary radius ρ to output m records. Here we drop these assumptions; the

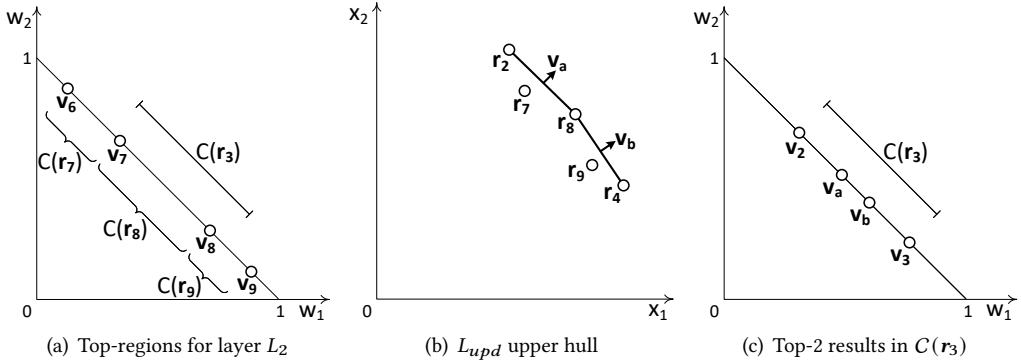


Fig. 4. Applying Theorem 1

first so that our algorithm is precomputation-free, and the second because it defies our problem formulation.

Without any precomputed layers or known ρ , our first step is to produce an overestimate of ρ , denoted as $\bar{\rho}$, which ensures an output size of at least m . That overestimate can be the radius required so that ORU's output for $k = 1$ includes m records. This radius, for any k , is guaranteed to produce at least as many records.

A straightforward approach to derive $\bar{\rho}$ (based on $k = 1$) would be to compute the upper hull of the entire dataset D , get the m top-regions with the smallest mindist to \mathbf{w} , and use the largest mindist among them as $\bar{\rho}$. That, however, may be too costly. Ideally, we would want to localize the upper hull computation to just the vicinity of \mathbf{w} . To achieve this, we exploit the fact that the ρ -skyline is a superset of ORU's output for the same ρ and $k = 1$, as follows directly from the definition of ρ -dominance.

We use an *incremental* ρ -skyline algorithm, which supports “get next” calls to extend a ρ -skyline for the immediately larger ρ around \mathbf{w} that admits exactly one new record to it. Details on that technique (for its general, ρ -skyband version) are presented in Section 5.3.2. We initialize that algorithm and prompt it until the ρ -skyline includes m records. Next, we compute their upper hull L_{tmp} . In general, not all ρ -skyline records will make it to L_{tmp} . If that is the case, we keep prompting the ρ -skyline algorithm, and updating the upper hull L_{tmp} to cover the additional records, until L_{tmp} includes m extreme vertices (records). We use as $\bar{\rho}$ the final radius reported by the ρ -skyline algorithm. On top of this, note that the final L_{tmp} is guaranteed to include all the part of layer L_1 that ORU processing could possibly need for an output of size m . In other words, the final L_{tmp} can serve already as layer L_1 in subsequent processing.

Using the obtained $\bar{\rho}$, we compute the ρ -skyband (for the actual k specified in the input). That can be done with a standard k -skyband algorithm by simply replacing regular dominance tests with $\bar{\rho}$ -dominance ones (described in the last paragraph of Section 4.2). The derived $\bar{\rho}$ -skyband is a guaranteed superset of the ORU output, thus we place its members into a *candidate set* M .

Even with $\bar{\rho}$ available, we have only an overestimate of the actual radius required to output m records. Thus, computing upper hull layers directly on M would be computationally wasteful. To circumvent this, and to also ensure an output of exactly m records, we employ an advanced, adaptive technique which: (i) gradually expands ρ around \mathbf{w} , moving from $\rho = 0$ towards $\bar{\rho}$, (ii) computes new upper hull layers on M as and when needed only, and (iii) progressively outputs confirmed candidates (i.e., records guaranteed to be in the ORU result) while the search is ongoing.

The latter property enables the tightening of $\bar{\rho}$ on the fly, and hence the shrinking of the $\bar{\rho}$ -skyband and the elimination of candidates from M , so that the layer (i.e., upper hull) computations execute on increasingly fewer records. This improved implementation is presented next.

5.3.1 Gradually expanding ρ . What we already have is the initial overestimate $\bar{\rho}$, (the necessary part of) layer L_1 , and the candidate set M . Subsequent ORU search expands concurrently (i) in terms of radius, from 0 towards $\bar{\rho}$, and (ii) in terms of layer depth, computing on demand (the necessary parts of) deeper L_i layers.

We treat the structure of all upper hull layers as an implicit tree, and apply the *best-first* approach to gradually explore that tree in increasing distance from the seed \mathbf{w} . In particular, we maintain a min-heap Q that organizes known top- i results (for $i \leq k$) and their respective preference regions C , with mindist to \mathbf{w} as their key. Let \mathbf{r} be the top record according to \mathbf{w} . We start with the top-region of \mathbf{r} (i.e., $C(\mathbf{r})$), whose mindist is by definition 0. First, we partition $C(\mathbf{r})$ according to the possible top-2-nd records, which requires computing L_2 (i.e., the upper hull of record set $M - L_1$) and applying Theorem 1, as demonstrated in Section 5.2. Each produced partition is pushed into Q with key equal to its mindist to \mathbf{w} , and associated with its (now known) top-2 result. Second, for each L_1 record \mathbf{r}_i that is adjacent to \mathbf{r} , we push its top-region $C(\mathbf{r}_i)$ into Q (with $\{\mathbf{r}_i\}$ as its top-1 result). Then, we iteratively pop the heap. For each region C popped from Q , we distinguish two cases:

Case 1: If C corresponds to a top- i result (with $i < k$), we partition it according to the different top- $(i + 1)$ -th records in C (using Theorem 1, as in Section 5.2), and push the produced partitions into the heap (associated with their, now known, top- $(i + 1)$ results). Applying Theorem 1 might require computing a new upper hull layer on the candidate set; details and an optimization are discussed later in this section. Importantly, if C corresponds to a top-1 result $\{\mathbf{r}_i\}$, we additionally push into the heap the top-regions of its adjacent records (omitting any that were pushed into Q previously, to avoid duplication). The reason is that, unlike best-first search in an actual tree, at the “root level” of our implicit structure, we initially did not push into Q the top-regions of all L_1 records, but only those neighboring $C(\mathbf{r})$. That was in order to save mindist calculations and unnecessary push operations, since many L_1 top-regions may lie too far from \mathbf{w} to affect ORU processing. Instead, we use the continuity implied by Lemmas 1 and 2 to gradually push top-regions from L_1 into Q , only when one of their neighbors is popped.

Case 2: If C corresponds to a top- k result, it is considered *finalized*. That is, the top- k result and its respective region C are appended to the ORU output. Observe that, as we explained in Section 3, our algorithm goes beyond Definition 2, to output not only records, but also specific order-sensitive top- k results, together with the preference regions that produce them.

The process terminates when the output includes m distinct records. The mindist of the last finalized region is the minimum radius ρ that appears in Definition 2. In a nutshell, our ORU methodology explores (i.e., partitions or finalizes, for $i < k$ and for $i = k$, respectively) regions C in increasing distance from \mathbf{w} , utilizing the implicit tree structure to dismiss those too distant to affect the result.

Figure 5 shows the implicit tree for our running example. Each node N_j represents a preference region (i.e., a segment, for $d = 2$) and its respective top- i result. The root corresponds to the L_1 top-region that includes \mathbf{w} , i.e., $C(\mathbf{r}_3)$. First, we partition $C(\mathbf{r}_3)$ into 3 regions, namely, $\mathbf{v}_2 \mathbf{v}_a$, $\mathbf{v}_a \mathbf{v}_b$, $\mathbf{v}_b \mathbf{v}_3$, as demonstrated in Figure 4. Associated with their top-2 results, they conceptually form nodes N_4 , N_5 , N_6 , and are pushed into Q . We also push the top-regions of L_1 records adjacent to \mathbf{r}_3 , i.e., $C(\mathbf{r}_2)$, $C(\mathbf{r}_4)$, associated with their top-1 results (nodes N_2 , N_3). Then, iterative popping commences.

The first popped node is N_5 (with mindist 0, since it actually includes \mathbf{w}), whose top-2 result is $\{\mathbf{r}_3, \mathbf{r}_8\}$. If $k = 2$, it is finalized and its top- k records are output directly (Case 2). Popping continues

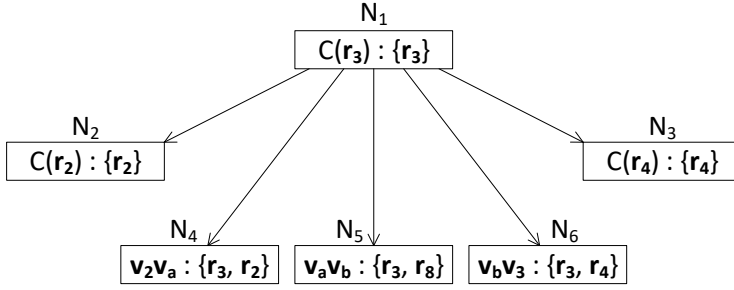


Fig. 5. Implicit tree

with N_4 and N_6 , which are finalized too; their top- k results contribute two new records to the output (i.e., r_2 and r_4). If $m = 4$, the process terminates. Otherwise, the next popped node is N_2 , for which we know the top-1 result (i.e., N_2 falls under Case 1). We will need to partition it by Theorem 1, and push into Q its resulting “children” (not illustrated). Importantly, since N_2 belongs to L_1 , we will also need to push its neighboring top-regions that were not encountered before, i.e., $C(r_1)$ (not illustrated). The implicit tree is constructed gradually, with new nodes formed for each Case 1 pop.

An important point on Case 1 is that partitioning C according to its different top- $(i + 1)$ -th records might require computing a new upper hull layer. That is, if layer L_{t+1} (referring to the value of t in Theorem 1 for C) was not previously computed, we need to compute it now. A naïve approach is to simply remove the first t layers from the candidate set M and compute the upper hull on the remaining candidates. An improvement however is possible, by shrinking M . Recall that our initial $\bar{\rho}$ estimation assumed we needed m records in the $\bar{\rho}$ -skyline. Letting τ be the number of distinct records which (i) belong to the top- k results already finalized, and (ii) are not members of the $\bar{\rho}$ -skyline, we can roll back the incremental ρ -skyline computation so that it includes only $m - \tau$ records. This backtracking effectively reduces $\bar{\rho}$, and in turn enables the shrinking of M to only keep ρ -skyband records (for the actual k input) for the reduced $\bar{\rho}$. The shrinking can be done trivially if we record the inflection radius for each record in the ρ -skyline during the original $\bar{\rho}$ estimation, and for each record in the initial candidate set (i.e., in the ρ -skyband for the original $\bar{\rho}$ estimate).

5.3.2 Incremental ρ -skyband. The OSS property and m being dictated by the user/application is a central point of our motivation, and thus a hard requirement for our operators. However, the ORU algorithm requires as a building block an incremental ρ -skyband module (IRD). While ORU requires this for the ρ -skyline only (i.e., $k = 1$), here we address the arbitrary k version, for generality.

The IRD challenge is that, unlike ORD, no ρ -dominance can ever be used to narrow down the search, because every k -skyband member may be output after a sufficient number of “get next” calls. Thus, the key question is how to serve these calls without computing the entire k -skyband. The main idea is to progressively fetch k -skyband members, but only output them when their inflection radius is no larger than a gradually growing threshold $\underline{\rho}$, introduced later.

IRD invokes a regular k -skyband algorithm to progressively fetch its members, in decreasing score order for \mathbf{w} . We use BBS [59] for that building block, but amend its default record/index node visiting order to order by score, as we did in Section 4.2. Let set T hold the k -skyband records fetched by BBS so far. As ensured by the score-based fetching order, we can compute the exact inflection radius for each record in T at the time it was fetched.

An invariant of branch-and-bound algorithms, like BBS, is that at any point during execution, their heap contents (records and index nodes) represent the not-yet-considered part of the dataset. Let S be the set of all records and nodes currently in the heap. For simplicity, we extend notation r_i to nodes too, since BBS anyway represents them by the top corner of their minimum bounding box. For each $r_i \in S$ we can compute an inflection radius ρ_i based on the current set T . However, that ρ_i is just a lower bound of the actual inflection radius, because BBS may have not yet fetched in T all the k -skyband records with score larger than $U_w(r_i)$, i.e., T may currently not include all records that ρ -dominate r_i .

Since S serves as a representation of all unexplored records, the minimum ρ_i among the members of S serves as an overall lower bound ρ for any non-fetched record. Therefore, every record in T with inflection radius no greater than ρ has a confirmed order in the output of IRD. In other words, these records are guaranteed to comprise the ρ -skyband for radius ρ .

Consider Figure 2(b), where $k = 5$. When IRD is first invoked, we execute BBS to progressively fetch the first k records, i.e., r_1 to r_5 , which by definition are the top- k . They are placed into set T and also output directly by IRD. When prompted with a “get next” call, IRD resumes BBS to progressively fetch new k -skyband members into T . Whenever BBS fetches a new record, we update ρ according to the current contents of the BBS heap (i.e., of set S). Let r_i be the not-yet-output record in T with the smallest inflection radius. If $\rho_i \leq \rho$, IRD outputs r_i and pauses. Otherwise, we keep fetching new records by BBS and updating ρ after each retrieval, until ρ becomes at least as large as the inflection radius of one record in T . That record is output by IRD as the next ρ -skyband member. Subsequent “get next” calls are served by resuming this process.

Returning to our example, consider that IRD receives a “get next” call (after its initialization, which reports the top-5 records all at once). It resumes BBS, but as it fetches r_6 , r_7 , r_8 into T , assume that ρ does not become as large as any of the inflection radii in T after any of these retrievals. However, when r_9 is fetched too, suppose that the updated ρ becomes greater than (or equal to) ρ_7 . IRD outputs r_7 and pauses (until it receives another “get next” call).

6 FASTER, EXACT ORU PROCESSING

In this section, we present an alternative, fundamentally different approach for (exact) ORU processing, namely, the *fast ORU* algorithm (FORU). It differs from basic ORU in Section 5 in three core aspects:

- Basic ORU starts with an overestimate of ρ , which it gradually shrinks until the output size drops to m . A disadvantage of this approach is that it examines a larger part of the preference domain (and therefore more candidate records) than necessary, thus wasting computations. In contrast, FORU uses the seed w as the starting point for a gradual expansion until the output includes m records.
- Basic ORU is centered on (and partitions the explored part of the preference domain according to) order-sensitive top- k results. Order sensitivity is not required by Definition 2 or by its applications, and thus leads to an unnecessary waste of computations. FORU abolishes order-sensitivity in the top- k results that lead its execution.
- Basic ORU starts with preference regions which it partitions in order to determine the possible top- k results therein. In contrast, FORU starts with possible (order-insensitive) top- k results and determines the regions that correspond to them.

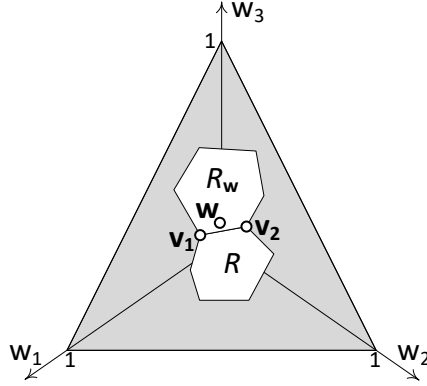


Fig. 6. Neighbor regions, neighbor sets, and Lemma 4 example

6.1 Fundamentals and Algorithmic Outline

In the context of FORU, any reference to top- k result/set refers to its order-insensitive definition. Also, in case of utility ties in the top- k -th position, the top- k set is assumed to include all tied records (i.e., its cardinality may exceed k). A central concept in FORU is the *result region*.

Result Region $\mathcal{R}(S)$: Given a top- k set S , the result region $\mathcal{R}(S)$ is the preference region that includes those and only those vectors \mathbf{v} whose top- k set is S . Following directly from its definition, $\mathcal{R}(S)$ is the intersection of half-spaces $U_{\mathbf{v}}(\mathbf{r}) \geq U_{\mathbf{v}}(\mathbf{r}')$ for each pair of $\mathbf{r} \in S$ and $\mathbf{r}' \in D - S$. Formally,

$$\mathcal{R}(S) = \bigcap_{\mathbf{r} \in S, \mathbf{r}' \in D - S} \{ \mathbf{v} \in \Delta^{d-1} \mid U_{\mathbf{v}}(\mathbf{r}) \geq U_{\mathbf{v}}(\mathbf{r}') \} \quad (1)$$

If $\mathcal{R}(S) = \emptyset$, set S is not a possible top- k result and we call it *infeasible*. On the other hand, if $\mathcal{R}(S) \neq \emptyset$, set S is a *feasible* top- k result and, geometrically, $\mathcal{R}(S)$ (as the intersection of a finite number of half-spaces) is a convex polytope [14].

Let $S_{\mathbf{w}}$ be the top- k set for the seed vector \mathbf{w} , and $\mathcal{R}_{\mathbf{w}}$ be the respective result region. FORU initially determines every possible top- k set when the preference vector moves infinitesimally outside $\mathcal{R}_{\mathbf{w}}$, and derives accordingly the result regions for these top- k sets. After computing all the neighbor regions to $\mathcal{R}_{\mathbf{w}}$, the process is repeated to determine their own neighbor regions, thus progressively tessellating the preference domain around \mathbf{w} until the union of top- k sets encountered hits the desired output size m .

Neighbor Region/Set: Consider a non-empty result region $\mathcal{R}(S)$. Formally, we define as *neighbor region* every non-empty result region $\mathcal{R}(S')$ that shares at least one extreme vertex with $\mathcal{R}(S)$. Similarly, we refer to S' as a *neighbor set* to S . Figure 6 illustrates the result region $\mathcal{R}_{\mathbf{w}}$ that (includes the seed vector \mathbf{w} and) corresponds to $S_{\mathbf{w}}$, and the result region $\mathcal{R}(S)$ of another top- k set, S . As $\mathcal{R}(S)$ has a common extreme vertex with $\mathcal{R}_{\mathbf{w}}$ (e.g., \mathbf{v}_1), it is a neighbor region to $\mathcal{R}(S)$, and set S is a neighbor set to $S_{\mathbf{w}}$.

LEMMA 4. *Given two neighboring result regions $\mathcal{R}(S)$ and $\mathcal{R}(S')$, and letting \mathbf{v} be the extreme vertex they have in common, the top- k set of \mathbf{v} is a superset of S and a superset of S' .*

PROOF. Consider a feasible set S and recall Equation 1. Any vector \mathbf{v} on the boundary of $\mathcal{R}(S)$ lies in at least one hyper-plane $U_{\mathbf{v}}(\mathbf{r}) = U_{\mathbf{v}}(\mathbf{r}')$, where $\mathbf{r} \in S$, $\mathbf{r}' \in D - S$. That is, the top- k set of \mathbf{v}

includes every record in S but, also, at least one extra record $\mathbf{r}' \notin S$. We can similarly show that the top- k set for any vector \mathbf{v} on the boundary of $\mathcal{R}(S')$ includes S' and at least one extra record. From the hypothesis that \mathbf{v} is a common extreme vertex of $\mathcal{R}(S)$ and $\mathcal{R}(S')$ it follows that \mathbf{v} lies on both their boundaries, and thus its top- k set is a superset of S and of S' . \square

Referring to Figure 6, assume that $k = 3$, that $S_{\mathbf{w}} = \{\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3\}$, and that the top- k set for \mathbf{v}_1 is $S_{\mathbf{v}_1} = \{\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3, \mathbf{r}_4, \mathbf{r}_5\}$. According to Lemma 4, for any neighbor region that has \mathbf{v}_1 as an extreme vertex (like $\mathcal{R}(S)$ in the figure), the respective top- k set S is a subset of $S_{\mathbf{v}_1}$. FORU uses that fact to discover each and every neighbor region to $\mathcal{R}_{\mathbf{w}}$ that shares \mathbf{v}_1 as an extreme vertex; recall that, although $S_{\mathbf{w}}$ (and thus $\mathcal{R}_{\mathbf{w}}$ and its extreme vertices) is clear how to derive, the neighbor regions are unknown. To compute the neighbor regions (and thus expand the tessellation around \mathbf{w}) we exploit the fact that $S \subset S_{\mathbf{v}_1}$ and consider every possible k -combination from $S_{\mathbf{v}_1}$ as a candidate top- k set, e.g., $S = \{\mathbf{r}_1, \mathbf{r}_4, \mathbf{r}_5\}$. For each of the candidate top- k sets, we apply Equation 1, using however $S_{\mathbf{v}_1}$ in place of D ; i.e., in our example where $S = \{\mathbf{r}_1, \mathbf{r}_4, \mathbf{r}_5\}$, the only non-result records \mathbf{r}' we need to consider in Equation 1 are \mathbf{r}_2 and \mathbf{r}_3 . The feasible among the candidate sets, i.e., those where $\mathcal{R}(S) \neq \emptyset$, lead to the neighbor regions that share \mathbf{v}_1 with $\mathcal{R}_{\mathbf{w}}$, like the one illustrated in Figure 6.

An important point is that, as we explained in the beginning of the section, the top- k set for a neighbor region could include more than k records. Although we consider combinations of exactly k records from $S_{\mathbf{v}_1}$, even if there are ties in the top- k -th position and some tying records are cut out from a k -combination, such records will not be missed by the overall process, because they are bound to appear in another feasible k -combination.

By repeating the above investigation for each extreme vertex of $\mathcal{R}_{\mathbf{w}}$, we can derive all neighbor regions/sets to $\mathcal{R}_{\mathbf{w}}$, completing one hop in the expansion. The tessellation could continue in a “breadth-first” manner, completing one hop after the other. This, however, is not the most efficient approach, since expansion is bounded by distance (i.e., the stopping radius ρ) instead of number of hops. For that reason, FORU employs a min-heap Q that organizes encountered extreme vertices and result regions by minimum distance to \mathbf{w} . This way, it incorporates their top- k sets in the overall output in the appropriate order, and prioritizes the derivation of new result regions in increasing distance from \mathbf{w} to save computations. In the following section, we present the technical specifics of FORU in the form of pseudo-code.

6.2 Technical Specifics

Algorithm 1 summarizes the FORU algorithm. Line 1 performs a regular top- k computation using a standard technique, like BBR [67] (the same holds for Lines 9 and 34). Line 3 initializes \mathcal{G} , an array of size m , which at the end of the process will hold the output of FORU. \mathcal{G} maintains records sorted in ascending order of a key (a distance to \mathbf{w} , in particular), initialized to ∞ for its unoccupied slots.

Line 4 initializes T_s , whose purpose is to avoid examining the same k -combination of records multiple times. Specifically, its elements are top- k sets, starting with the seed’s top- k set (i.e., $S_{\mathbf{w}}$). Each k -combination encountered by FORU (in Line 19) is book-kept as a separate element in T_s . That is essential, since the same k -combination could occur while examining different extreme vertices whose top- k sets happen to overlap. Similarly, the role of T_v is to avoid considering anew extreme vertices that have already been accounted for (in Line 30) in the context of another neighbor result region.

As explained previously, the min-heap Q (initialized in Line 6) organizes the encountered extreme vertices and result regions. If a heap element corresponds to an extreme vertex \mathbf{v} , the element includes the respective top- k set $S_{\mathbf{v}}$, with key the distance between \mathbf{v} and \mathbf{w} . If the heap element corresponds to a result region $\mathcal{R}(S)$, it includes the respective k -combination S , with key the mindist between $\mathcal{R}(S)$ and \mathbf{w} . When we pop the heap (Lines 13 and 27), the element with the smallest key is

removed from the heap, and its record set (S_v or S) and key value ($\text{dist}(v, w)$ or $\text{mindist}(\mathcal{R}(S), w)$, respectively) become available.

Let S_v be the top- k set and $\text{dist}(v, w)$ be the key of the popped heap entry. For every record r in S_v , we update \mathcal{G} (in Lines 15-18). Specifically, if r is not already in \mathcal{G} and $\text{dist}(v, w)$ is smaller than the m -th key in \mathcal{G} , we insert r into \mathcal{G} with key $\text{dist}(v, w)$ (evicting the previous m -th record from \mathcal{G} to keep its size to m). If r is already in \mathcal{G} with a key greater than $\text{dist}(v, w)$, we update its key to $\text{dist}(v, w)$. The latter situation may occur if r was previously encountered in the top- k set of another popped heap entry. As an example, consider Figure 6, where \mathcal{R} is pushed into the heap only after vertex v_1 has been popped, despite the fact that \mathcal{R} is actually nearer to w . If a record r in S_{v_1} had been inserted into \mathcal{G} , and the same record is now encountered in the top- k set of \mathcal{R} , the key of r in \mathcal{G} will be updated to $\text{mindist}(\mathcal{R}, w)$. The process to update \mathcal{G} when the popped heap entry corresponds to a result region is similar (Line 29).

FORU terminates when the m -th key in \mathcal{G} is no greater than the key at the top of the heap, at which point \mathcal{G} is finalized and output. The radius ρ in Definition 2 corresponds to the m -th key in the finalized \mathcal{G} . Another termination scenario is when the heap becomes empty, which means that the entire preference domain has been explored. The latter case may occur when the possible top- k records throughout the preference domain are fewer than m .

A low-level optimization regards Lines 23 and 24. To save computations, we may first check the feasibility of S via linear programming [14] before spending the (typically higher) cost to derive the exact geometry of $\mathcal{R}(S)$ via half-space intersection. In Line 25, we compute $\text{mindist}(\mathcal{R}(S), w)$ using quadratic programming [32, 54].

As a note on Line 23, recall that $\mathcal{R}(S)$ is defined by Equation 1, using S_v in place of D for the half-space intersection. S_v , being a top- k set, is expected to be small, i.e., to include just slightly over k records. Unlike Line 23, however, the derivation of \mathcal{R}_w in Line 2 is more complex. Specifically, \mathcal{R}_w can still be computed by Equation 1, but the non-result records we need to consider are all the k -skyband records (i.e., by using the k -skyband in place of D). Surely, using the k -skyband instead of the entire D in Equation 1 is a relief. However, the k -skyband may still include numerous records, implying a large number of half-spaces to intersect. We may reduce the number of half-spaces based on the following two observations.

First, let S_{NR} be the set of k -skyband records that do not belong to S_w (i.e., to the top- k set for w). These are the only possible candidates to enter the top- k result for a preference vector other than w . Consider that a non-result record $r_i \in S_{NR}$ dominates (in the traditional sense) another non-result record $r_j \in S_{NR}$. For any preference vector v , it holds that $U_v(r_i) \geq U_v(r_j)$. Therefore, referring to Equation 1, the condition that a result record $r \in S_w$ stays ahead of r_i utility-wise (i.e., the condition $U_v(r) \geq U_v(r_i)$) is stricter, and hence subsumes the condition that r stays ahead of r_j . Formally, the half-space $U_v(r) \geq U_v(r_j)$ is not a bounding half-space of \mathcal{R}_w and can thus be ignored. In essence, we can remove from S_{NR} those records r_j that are dominated by any other record $r_i \in S_{NR}$. The shrunk version of S_{NR} includes all non-result records that may affect \mathcal{R}_w .

Second, consider a result record $r_i \in S_w$ that dominates another result record $r_j \in S_w$. For any preference vector v if holds that $U_v(r_i) \geq U_v(r_j)$. Therefore, the condition that r_j stays ahead of a non-result record r' utility-wise (i.e., the condition $U_v(r_j) \geq U_v(r')$) is stricter, and thus subsumes the condition that r_i stays ahead of r' . In other words, the half-space $U_v(r_i) \geq U_v(r')$ is not a bounding half-space of \mathcal{R}_w and can be ignored. Letting S_R be the set of result records $r_j \in S_w$ that do not dominate any other result record $r_i \in S_w$, and using the shrunk set S_{NR} from the previous observation, \mathcal{R}_w is the intersection $\bigcap_{r \in S_R, r' \in S_{NR}} \{v \in \Delta^{d-1} | U_v(r) \geq U_v(r')\}$.

Algorithm 1 FORU(\mathbf{w}, D, m)

```

1:  $S_{\mathbf{w}} \leftarrow$  the top- $k$  set for  $\mathbf{w}$ 
2:  $\mathcal{R}_{\mathbf{w}} \leftarrow$  the result region of  $S_{\mathbf{w}}$ 
3: Initialize  $\mathcal{G}$  as a size- $m$  sorted array of records, and insert each record from  $S_{\mathbf{w}}$  into it with key 0
4: Initialize the set of encountered  $k$ -combinations  $T_s \leftarrow \{S_{\mathbf{w}}\}$   $\triangleright T_s$  is a set of  $k$ -combinations
5: Initialize the set of encountered extreme vertices  $T_v \leftarrow \emptyset$ 
6: Initialize an empty min-heap  $Q$ 
7: for each extreme vertex  $\mathbf{v}$  of  $\mathcal{R}_{\mathbf{w}}$  do
8:    $T_v \leftarrow T_v \cup \mathbf{v}$ 
9:    $S_{\mathbf{v}} \leftarrow$  the top- $k$  set for  $\mathbf{v}$ 
10:  Push set  $S_{\mathbf{v}}$  into  $Q$  with key the distance between  $\mathbf{v}$  and  $\mathbf{w}$ 
11: while  $Q$  is non-empty and the  $m$ -th key in  $\mathcal{G}$  is greater than the key at the top of  $Q$  do
12:   if the top of  $Q$  is a vertex  $\mathbf{v}$  then
13:      $\langle S_{\mathbf{v}}; \text{dist}(\mathbf{v}, \mathbf{w}) \rangle \leftarrow Q.\text{pop}()$ 
14:     for each record  $\mathbf{r} \in S_{\mathbf{v}}$  do  $\triangleright$  Update  $\mathcal{G}$  with the records in  $S_{\mathbf{v}}$ 
15:       if  $\mathbf{r}$  is not in  $\mathcal{G}$  and  $\text{dist}(\mathbf{v}, \mathbf{w})$  is smaller than the  $m$ -th key in  $\mathcal{G}$  then
16:         Insert  $\mathbf{r}$  into  $\mathcal{G}$  with key  $\text{dist}(\mathbf{v}, \mathbf{w})$ 
17:       else if  $\mathbf{r}$  is in  $\mathcal{G}$  with key greater than  $\text{dist}(\mathbf{v}, \mathbf{w})$  then
18:         Update the key of  $\mathbf{r}$  to  $\text{dist}(\mathbf{v}, \mathbf{w})$ 
19:     for every  $k$ -combination  $S$  of records from  $S_{\mathbf{v}}$  do
20:       if  $S \in T_s$  then
21:         continue  $\triangleright S$  has been encountered before; skip it
22:        $T_s \leftarrow T_s \cup \{S\}$ 
23:        $\mathcal{R}(S) \leftarrow$  the result region of  $S$ 
24:       if  $\mathcal{R}(S) \neq \emptyset$  then  $\triangleright S$  is a feasible top- $k$  set
25:         Push set  $S$  into  $Q$  with key the mindist between  $\mathcal{R}(S)$  and  $\mathbf{w}$ 
26:   else  $\triangleright$  The top of  $Q$  is a feasible result region  $\mathcal{R}(S)$ 
27:      $\langle \mathcal{R}(S); \text{mindist}(\mathcal{R}(S), \mathbf{w}) \rangle \leftarrow Q.\text{pop}()$ 
28:     for each record  $\mathbf{r} \in S$  do  $\triangleright$  Update  $\mathcal{G}$  with the records in  $S$ 
29:       Same as Lines 15-18 using key  $\text{mindist}(\mathcal{R}(S), \mathbf{w})$  for  $\mathbf{r}$ 
30:     for each extreme vertex  $\mathbf{v}$  of  $\mathcal{R}(S)$  do
31:       if  $\mathbf{v} \in T_v$  then
32:         continue  $\triangleright \mathbf{v}$  has been encountered before; skip it
33:        $T_v \leftarrow T_v \cup \mathbf{v}$ 
34:        $S_{\mathbf{v}} \leftarrow$  the top- $k$  set for  $\mathbf{v}$ 
35:       Push set  $S_{\mathbf{v}}$  into  $Q$  with key the distance between  $\mathbf{v}$  and  $\mathbf{w}$ 
36: Return  $\mathcal{G}$ 
    
```

7 APPROXIMATE ORU PROCESSING

Our experiments show that the (exact) algorithm for ORD offers sub-second response times even for multi-million record datasets. The (exact) ORU algorithm, on the other hand, can take 2 orders of magnitude longer to terminate, i.e., in the order of minutes for several settings. In this section, we describe algorithms for approximate ORU processing that offer a controllable trade-off between accuracy and running time, and come with proven approximation guarantees.

In Section 7.1, we present the fundamental *approximate ORU* algorithm (AORU), which follows closely the (exact) ORU processing paradigm, but simplifies computations when the inaccuracy is guaranteed to be within a tolerable limit. In Section 7.2, we analyze the accuracy of AORU and prove its error bound. In Section 7.3, we describe DORU, an alternative algorithm for approximate

ORU processing; DORU follows a different search paradigm from AORU that relies on lighter and easier-to-implement building blocks, while it still offers the same accuracy guarantee.

7.1 The AORU Operator

Algorithmically, AORU follows the same principles as the exact ORU algorithm in Section 5. However, the main rationale behind it is that, when an encountered region C is sufficiently small, a single top- k result may offer a good representation for any preference vector in C . By “good”, we imply both in practical terms (i.e., accuracy and efficiency), but also in terms of formal bounds, regarding any missed records (for the unaccounted vectors within C). Stated this intuition, the main challenge lies in the offer of approximation guarantees. We first describe the AORU algorithm, followed by its accuracy analysis in Section 7.2.

AORU receives an extra input, δ , which controls the accuracy of the approach and determines the bounds of the approximation. The process is identical to the regular ORU, **FIXME: (be it with the original or the improved $\bar{\rho}$ estimator from Chapter 1, whichever is faster)** but with a modification in Case 1 handling. Specifically, when the region C popped from the min-heap Q corresponds to a top- i result (for $i < k$), we compute the minimum bounding box (MBB) of C . If the main diagonal of the MBB according to the L^1 -norm (i.e., its Manhattan length) is no longer than δ , we drill-finalize C . That is, we do not partition C further. Instead, we perform a *drill* (i.e., a top- k query) at the centroid of the MBB, **FIXME: remove the clarification *and the emphasis on “drill”* if drill has been defined earlier** and insert the retrieved top- k result in the AORU output. Otherwise (i.e., if the L^1 length of the MBB diagonal is greater than δ), C is not exempted from partitioning and AORU proceeds with it just like the exact ORU algorithm. As per normal, we terminate AORU when its result size reaches m . In the event that the last drill operation retrieves more new records than necessary to fill the AORU result, we only admit (to the AORU result) the higher-ranking among them to bring its size to exactly m .

7.2 Accuracy Analysis of AORU

<THEOREM & PROOF TO BE PLACED HERE>

7.3 The DORU Operator

In this section, we present DORU, an alternative approximate operator. DORU abandons the processing principles of the exact ORU algorithm, and thus follows a different philosophy from AORU. It offers the same accuracy guarantee (as AORU), but it relies on simpler techniques, making it generally faster and also easier to implement. DORU’s core mechanism is the incremental ρ -skyband algorithm in Section 5.3.2 (for the problem’s actual value of k , not for $k = 1$) and repetitive drills at carefully selected preference vectors, called drill locations. It owes its acronym to the said drill operations.

The general idea is to perform repetitive drills around \mathbf{w} until the union of their top- k results includes at least m distinct records. In the interest of accuracy, the drill locations must be as close to \mathbf{w} as possible. On the other hand, drills whose top- k result does not include any new record (compared to those already fetched by previous drills) waste computations without contributing towards the termination of the process. We propose an approach where the drill locations are chosen in an adaptive fashion so that (i) they are gradually moving farther from \mathbf{w} , and (ii) each drill can be reasonably expected (albeit not guaranteed) to fetch previously unseen records. DORU executes in three stages.

Stage 1: In the first stage, we initialize a list L to store the drill location and the respective top- k result for every drill performed. We perform the first drill at \mathbf{w} itself. Then, we invoke the incremental

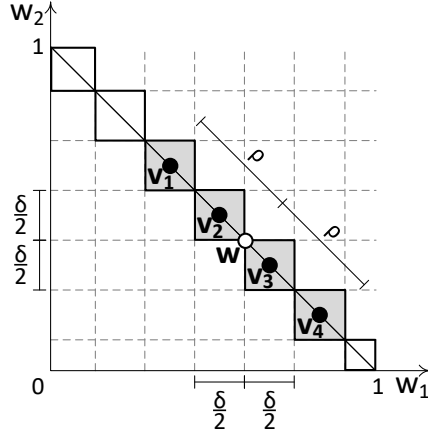


Fig. 7. Stage 2 drills

ρ -skyband algorithm (IRD) from Section 5.3.2. For each new ρ -skyband record \mathbf{r}_i that is retrieved, we perform a drill operation at a vector \mathbf{v}_i , determined as follows.

With reference to the notation introduced in Figure 1(b), letting $\rho_{i,j}$ be the inflection radius of \mathbf{r}_i (i.e., $\rho_i = \rho_{i,j}$), \mathbf{v}_i is the vector where the switch between \mathbf{r}_i and \mathbf{r}_j happens. For example, in Figure 1(b), \mathbf{v}_i is the vector in $\mathbf{s}_{i,j}$ at minimum distance from \mathbf{w} . The intuition is that \mathbf{v}_i is the closest vector to \mathbf{w} where \mathbf{r}_i first enters the ρ -skyband and therefore, very likely, where \mathbf{r}_i first appears in the top- k result. **FIXME:[Keming, if multiple records are tied for the k -th position of the top- k result, please add all of them to L – such ties are very likely and we shouldn't let them go to waste].** Computation-wise, note that vector \mathbf{v}_i is available already, since it was computed by IRD during the derivation of \mathbf{r}_i 's inflection radius.

The drill location \mathbf{v}_i , together with the retrieved top- k result, are inserted into L , and IRD is prompted to produce the next ρ -skyband record, which will determine the next drill location, and so on. Note that IRD retrieves ρ -skyband records in increasing inflection radius order, and therefore the drills are executed at increasingly farther locations from \mathbf{w} , as we originally intended. Once the number of distinct records in L hits m , the first stage of DORU terminates. The distance of the last drill location from \mathbf{w} is also recorded, as an interim radius ρ .³

Stage 2: List L could already be used to produce an (inexact) ORU output, however, that output would carry no accuracy guarantees. The second stage of DORU incorporates parameter δ and ensures that Theorem ?? is upheld. In particular, it performs another series of drill operations (irrelevant to those of the first stage), inserting new drill locations and top- k results into L .

We utilize an (implicit) regular partitioning of the preference domain according to d -dimensional hyper-cubes of equal size. The side-length of the hyper-cubes is $\frac{\delta}{d}$, so that the main diagonal of each hyper-cube has an L^1 -norm length of δ . The partitioning is centered around \mathbf{w} . Note that since the preference domain is a $(d-1)$ -simplex and the hyper-cubes are d -dimensional, we only consider hyper-cubes that overlap with the preference domain. In the 2-dimensional example of Figure 7, the relevant hyper-cubes (squares, in $d=2$) are drawn with solid boundary. Their side-length is $\frac{\delta}{d} = \frac{\delta}{2}$, except for the border ones that may be smaller (since the partitioning is centered around \mathbf{w}).

³Note that, due to the possible omission of result records by any inexact algorithm, this ρ cannot be smaller than the stopping radius of the exact ORU algorithm.

For every hyper-cube whose minimum distance to \mathbf{w} is ρ or smaller, we perform a drill at its centroid (recall that ρ is the interim radius we inherited from Stage 1). For the ρ radius illustrated in Figure 7, the necessary hyper-cubes are shown in gray. Stage 2 performs drills at their centroids (i.e., \mathbf{v}_1 to \mathbf{v}_4) and updates list L accordingly. Observe that the drill locations of Stage 2 alone, are spaced such that they already meet the conditions of Theorem ??.

A technical detail here is that for $d > 2$ the centroids of the hyper-cubes lie higher than the preference domain (i.e., $\sum_{i=1}^d w_i > 1$). The magnitude of the centroids does not affect the retrieved top- k results (hence, neither the validity of Theorem ??). However, we still normalize them according to $\sum_{i=1}^d w_i = 1$ before adding them to L . The reason, as will become clear in Stage 3, is that every top- k result in L must be accompanied by a valid drill location, i.e., one that falls inside the preference domain.

Stage 3: At this point, list L includes *the location and top- k result for every drill performed in the previous two stages*. The termination condition of Stage 1 guarantees that there are at least m distinct records in L . At the same time, the spacing of drill locations in Stage 2 ensures that the top- k results in L suffice to offer the accuracy guarantees of Theorem ?? up to distance ρ from \mathbf{w} . In Stage 3, we effectively shrink ρ by “truncating” L to include exactly m distinct records.

In particular, we sort the top- k results in L in increasing distance between their drill location and \mathbf{w} . For each top- k result, in the aforementioned order, we add its records to the DORU output (eliminating duplicates) until the latter includes m records. If the last considered top- k result introduces to the output more new records than necessary, we only include the higher-ranking ones to an output size of exactly m .

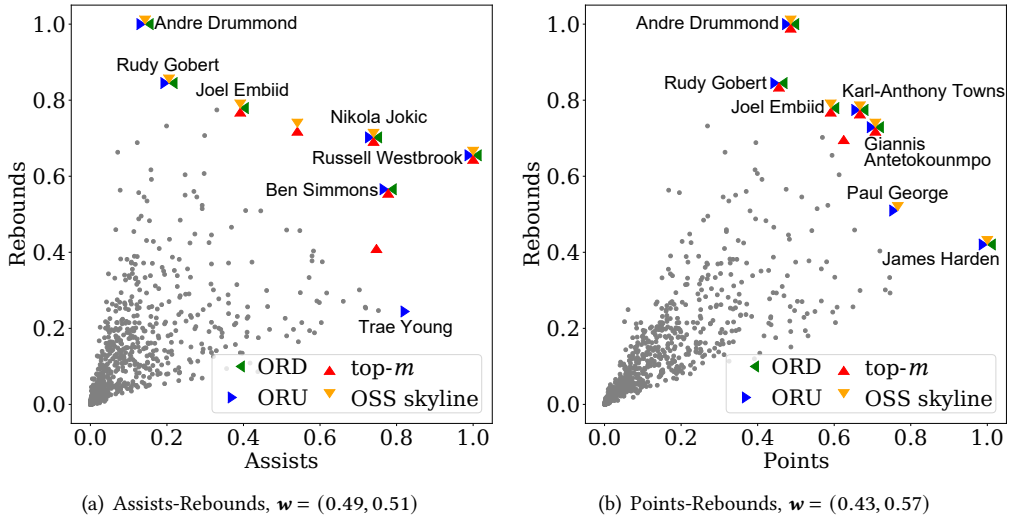
8 EXPERIMENTS

In this section, we present qualitative and performance experiments. We use both real and synthetic datasets, i.e., *HOTEL*, *HOUSE*, *NBA*, and *ANTI*, *COR*, *IND*, respectively. *HOTEL* contains 418,843 hotel records with $d = 4$ attributes [1]. *HOUSE* includes 315,265 records of $d = 6$ types of household expenses [2]. *NBA* holds $d = 8$ statistics for 21,960 NBA players [4]. The synthetic datasets represent typical distributions in multi-objective decisions [15]. Table 2 lists the problem parameters with their tested and default values (in bold). In each experiment, we vary one parameter and fix the others to their defaults. Every measurement is the average for 50 random seeds \mathbf{w} . The datasets are indexed by R-trees and kept in memory. That said, our methods are applicable to disk-based storage too. All algorithms were implemented in C++ and run on a machine with Intel i7-7700 CPU at 3.60Ghz and 32Gb RAM. We used `lp_solve` [3] as the linear programming solver, QuadProg++ [27] as the quadratic programming solver, and Qhull [10] for computational geometric primitives. **Bo: FIXME: Anonymity: Our implementation is available at [5].**

Parameter	Tested and default values
Dataset cardinality $ D $	100K, 400K , 1.6M, 6.4M, 25.6M
Dimensionality d	2, 3, 4 , 5, 6, 7
Parameter k	1, 5 , 10, 15, 20
Output size m	10, 30, 50 , 70, 90

Table 2. Parameters, tested values, and defaults

We start with qualitative results to draw a distinction between our operators and representative previous ones from Table 1.


 Fig. 8. Case study on NBA 2018-19 statistics ($k = 2, m = 6$)

8.1 Qualitative Study

First, we perform a case study to visualize how the results of ORD and ORU differ from (i) an OSS skyline, and (ii) a top- m query for the same vector \mathbf{w} . The former reports the m skyline members that dominate the most non-skyline records, following the most cited full-dimensionality OSS skyline definition [50]. We use the NBA 2018-19 season statistics for the total of its 708 players on Assists and Rebounds (in Figure 8(a)), and Points and Rebounds (in Figure 8(b)), normalized in the $[0, 1]$ range. We set $k = 2, m = 6$, and use $(0.49, 0.51)$ and $(0.43, 0.57)$ as the seed \mathbf{w} , respectively. The results of the methods are illustrated as differently oriented/colored triangles.

A first observation is that our operators report distinct results from past approaches (and from each other). For example, in Figure 8(a) only ORU reports *Trae Young* (rising stars challenge player), while half of its output records are not in the top- m result, and one third of them are not in the OSS skyline. The comparison of our operators with top- m reveals an even more interesting fact. In Figure 8(a), top- m misses *Andre Drummond* (the rebound leader in 2018-19 season), whom it would report if we only slightly revised \mathbf{w} from $(0.49, 0.51)$ to $(0.48, 0.52)$. Similarly, in Figure 8(b), top- m misses *James Harden* (the season's scoring leader), whom it would report if we revised \mathbf{w} from $(0.43, 0.57)$ to $(0.44, 0.56)$. The inclusion of these players in the result of both our operators (in the respective Figure 8(a) and 8(b) settings) confirms that they successfully employ some “width” in their search, by reporting records that are particularly strong for alternative, very similar preferences to the seed \mathbf{w} . Investigations in our full-scale experimental settings, using the Jaccard coefficient as the similarity measure (charts omitted in the interest of space), demonstrate that (i) OSS skyline is more dissimilar to our operators than top- m , as it is not guided by \mathbf{w} , and (ii) top- m becomes increasingly more dissimilar to ORD/ORU as m grows. As an indication for the omitted charts, for *IND* data and the default parameters in Table 2, the Jaccard similarity of OSS skyline to ORD is 0.25, and to ORU it is 0.24. In the same setting, the Jaccard similarity of top- m to ORD is 0.44, and to ORU it is 0.32.

Next, we use real TripAdvisor data and reviews (TA) [6]. TA includes ratings for 1,850 hotels on $d = 7$ aspects (value, room, location, etc), forming dataset D . It also includes actual user reviews

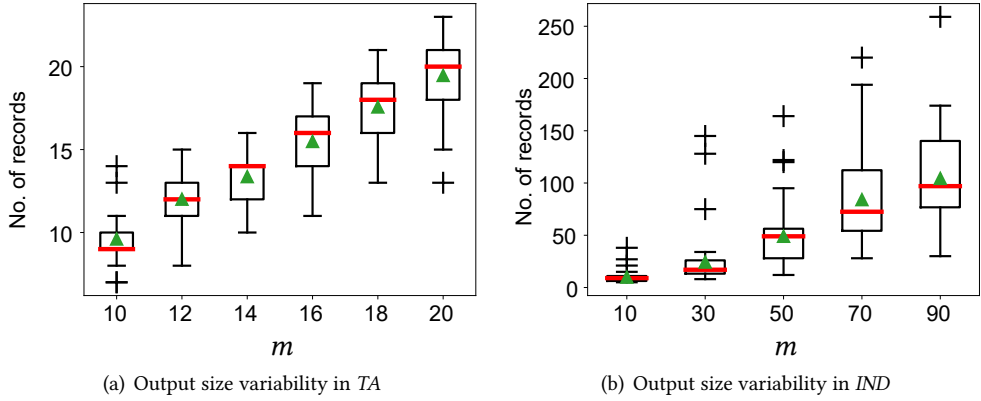


Fig. 9. Output size analysis of fixed-region techniques

for these hotels, each comprising a comment and an overall score. The reviews offer a practical example of how preference vectors could be extracted for real users. Specifically, we employ [71], an established preference mining method that estimates a user's weight vector based on her reviews. Applying [71], we get vectors for 137,563 users. The dataset and vectors from *TA* are used in the next experiment on fixed-region methods [21, 55], but at the same time they offer an end-to-end application example for our techniques.

The fixed-region methods require a preference polytope R to be specified as part of their input. We demonstrate that it is not feasible to estimate the size of R required to produce, even approximately, m records. Worse yet, even the same polytope R , when positioned at different parts of the preference domain, can produce vastly different output sizes. Specifically, we use $k = 5$ and vary m from 10 to 20 (since the dataset in *TA* is highly correlated and its 5-skyband includes only 61 hotels). We first execute ORU for 50 randomly selected *TA* users (preference vectors) and record the average stopping radius, denoted as ρ^* . We then produce a hyper-cube R with volume equal to the hyper-sphere with radius ρ^* . Next, we count the number of distinct records (*TA* hotels) output by the fixed-region top- k operator in [55] when R is positioned around each of the 50 user vectors (we use [55] since [21] works only for $k = 1$). Figure 9(a) presents in box-plot the observed variation in output size. To confirm, in Figure 9(b), we repeat this process for our full-scale setting, using random preference vectors, *IND* data, the default parameters, and standard range for m from Table 2. The box-plots indicate that even with knowledge of ρ^* , fixed-region techniques can produce dramatically different output sizes, which may hugely under- or over-shoot the target m . In contrast, our operators relieve the user/application from the need to meddle with the preference domain's complex dynamics, and abide strictly by the requested output size m . Note that the counterpart of this experiment, using ORD to compute ρ^* and producing the fixed-region R -skyband, demonstrates an even greater variability than Figure 9 for both *TA* and *IND* (charts omitted in the interest of space). For example, for target $m = 50$ in *IND* data, the fixed-region R -skyband outputs from 12 all the way to 269 records.

8.2 ORD Performance

In this section, we evaluate our ORD algorithm. For comparison, we adapted a fixed-region R -skyband technique (RSB), as described at the end of Section 2. In the adaptation, the initial hyper-cube R is sized so that its volume ratio to the preference domain equals the ratio of m to the

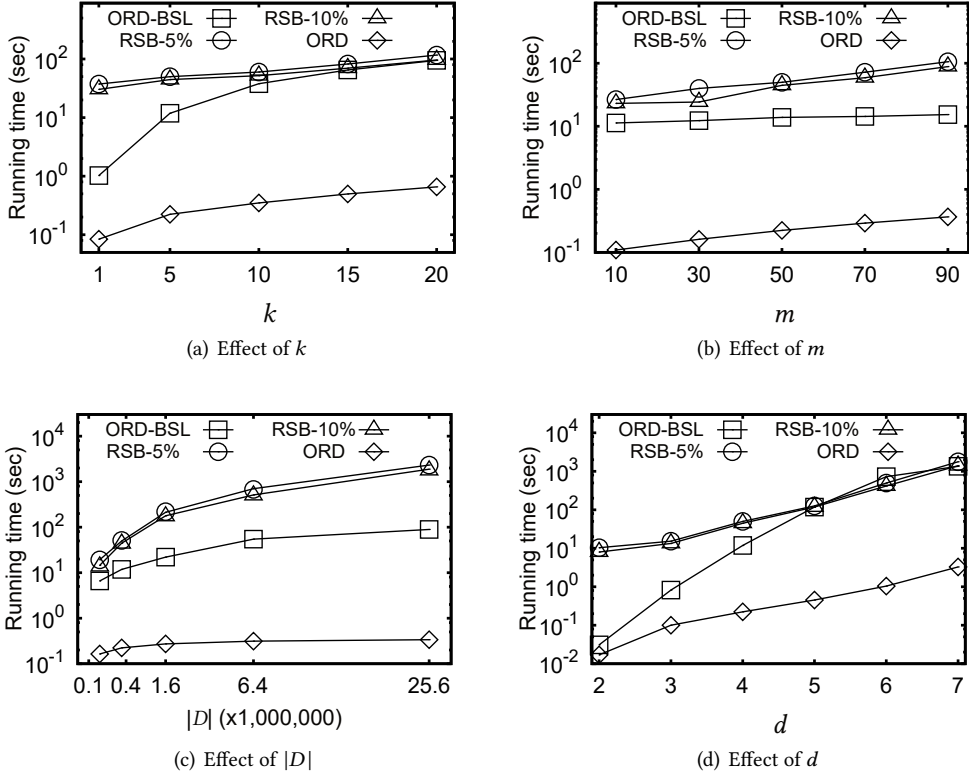


Fig. 10. ORD versus adaptations and baseline

expected k -skyband cardinality (i.e., $\frac{k \ln^{d-1} n}{d!}$, according to [31]). Based on the size of the R -skyband computed for the initial R , its volume is re-estimated proportionally to the desired m . The trials (R -skyband computation and R re-estimation) are repeated until the output is within 5% or 10% from m , resulting in two RSB versions. For the implementation of RSB, we use the index-based R -skyband module of [55], as it is considerably faster than the no-index approach of [21]. We also include baseline ORD-BSL that computes the entire k -skyband according to Section 4.1, without the enhancements in Section 4.2.

In Figure 10, we present (in logarithmic scale) the running time for all competitors versus each problem parameter, using *IND* data. ORD is 2 to 4 orders of magnitude faster than both versions of RSB, indicating the impracticality of fixed-region approaches to mimic its output, despite the ample slack given. The main reason is the numerous trials required for RSB to “converge” to an acceptable deviation from m . The runner-up is ORD-BSL, which is 1 to 2 orders of magnitude slower than the fully-enhanced ORD. The results also demonstrate ORD’s ability to scale. Its running time grows almost linearly to k and m , and sub-linearly to $|D|$. It increases more sharply with d , due to the growing complexity of its geometric building blocks. Still, even for $d = 7$ ORD takes only 3.2s.

Having established the superiority of our ORD algorithm, in Figure 11(a) we plot its running time for different synthetic distributions. Processing in *ANTI* is the slowest. The reason is that dominance (and, by deduction, ρ -dominance too) is less frequent among *ANTI* records, thus many

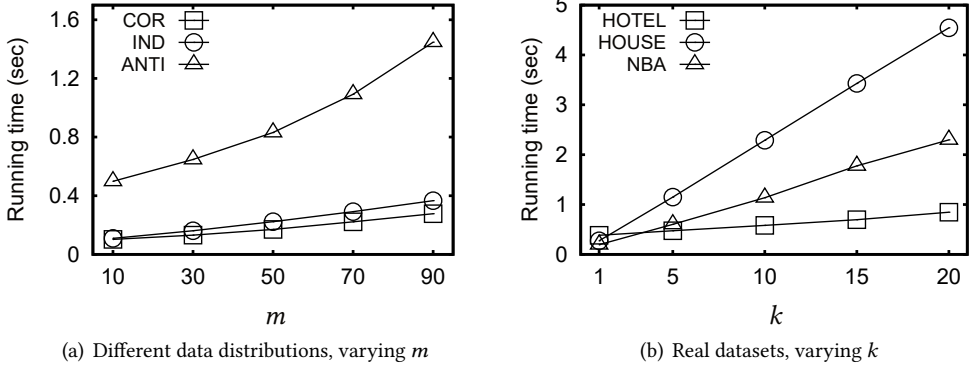


Fig. 11. ORD on different distributions and real datasets

candidates need to be considered before ORD can terminate. *COR* exhibits the inverse effect. In Figure 11(b), we use real data and vary k . *HOTEL* and *HOUSE* have comparable size, yet ORD is faster on *HOTEL*, due to its smaller dimensionality. *NBA* has the smallest cardinality, but the largest dimensionality, which explains why its line is in between the other two.

8.3 ORU Performance

Turning to ORU, we use for comparison an adaptation of the most efficient fixed-region top- k algorithm in [55], termed JAA. We employ a similar R estimation and trial approach as for RSB, allowing a deviation of up to 5% or 10% from the desired output size m . We also include ORU-BSL, a baseline that utilizes the initial overestimate $\bar{\rho}$ in Section 5.3, computes upper hull layers on the entire candidate set M , partitions the L_1 top-regions by Theorem 1, and reports the m -sized union of top- k records for the closest produced regions to \mathbf{w} .

In Figure 12, we plot the running time of all competitors versus each problem parameter, in logarithmic scale. ORU-BSL, although identical to ORU for $k = 1$, generally performs very poorly, failing to terminate within reasonable time in most large settings. This demonstrates the vital role of our gradual expansion in terms of both ρ and layer depth, presented in Section 5.3.1. Indeed, the full-fledged ORU is 2 to 4 orders of magnitude faster than ORU-BSL. ORU is also 12 to 134 times faster than JAA-10%, and even more compared to JAA-5%, confirming the general unsuitability of fixed-region approaches to our problems. Regarding ORU itself, it scales well with all parameters. Its running time increases faster than ORD because, despite the similarities in their definition, the nature of ORU is significantly more complex. To intuitively see this, determining whether a record \mathbf{r}_i belongs to the ρ -skyband of a dataset, we need a ρ -dominance test per competitor at worst. In contrast, to determine whether a record \mathbf{r}_i could appear in the top- k result for any vector \mathbf{v} within radius ρ from \mathbf{w} , we need to consider multiple combinations of competitors that may outscore \mathbf{r}_i , and different \mathbf{v} possibilities within distance ρ from \mathbf{w} .

Focusing on our ORU algorithm, in Figure 13(a), we try it on different synthetic data distributions. As expected, the problem is the toughest in *ANTI*, and the easiest in *COR*, for the reason explained in Figure 11(a). In Figure 13(b), we execute ORU on real datasets. Unlike Figure 11(b), processing in *NBA* is slower than *HOUSE*. The reason is that *NBA*'s higher dimensionality affects ORU's performance more than ORD's. ORU has a stronger geometric nature and a greater reliance on

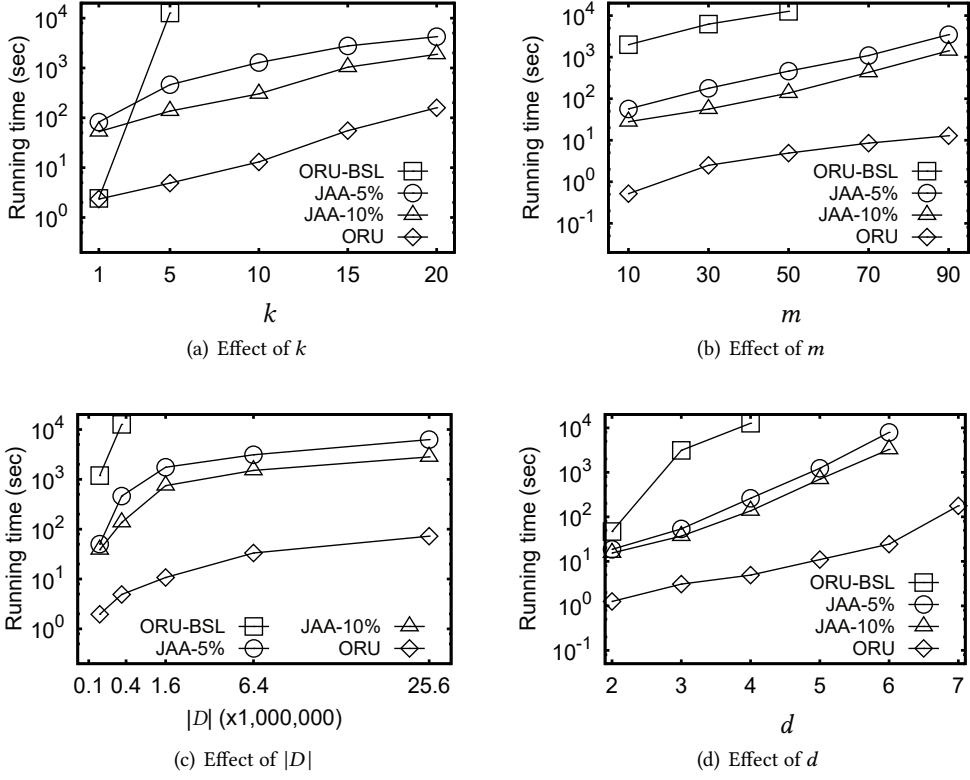


Fig. 12. ORU versus adaptations and baseline

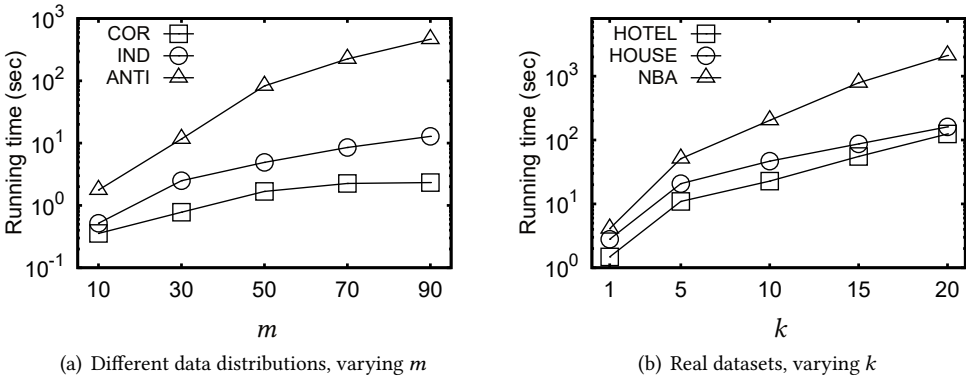


Fig. 13. ORU on different distributions and real datasets

computational geometric primitives (such as upper hull computation), whose complexity grows with d .

8.4 Discussion

In our default setting, for 400K *IND* records ORU requires 4.9s, while for 25.6M records it takes 72s. On the other hand, ORD runs in less than 1s in both cases (0.22s and 0.34s, respectively). That is, ORD is ready even for applications that require sub-second responses, while ORU is not quite there. To employ ORU in such applications, a viable approach is parallelization. At the core of its execution, ORU partitions numerous regions C under Case 1 in Section 5.3.1. With little synchronization effort, multiple regions could be (popped from ORU's min-heap and) partitioned in parallel, since determining the different top- $(i + 1)$ -th records in each of these regions is independent of the others. An orthogonal approach is to materialize the first upper hull layers of the dataset, similar to the *Onion* technique [18], or to cache and reuse the partial upper hulls from past ORU queries, together with the respective top-regions. Note that this will benefit performance even if ORU needs to consider (and thus compute on the fly parts of) upper hulls deeper than the available ones. A more general direction would be to let go of order-sensitivity in the reported top- k results (since it is anyway not required by Definition 2) or to switch to ORD processing if a region under Case 1 is deemed too small. For the former direction, algorithmic redesign will be necessary; for the latter, formal guarantees should be given to bound the induced deviation from the exact ORU answer. The last two directions could be promising for future work.

9 CONCLUSION

This paper draws motivation from the known weaknesses of standard skyline and top- k queries. Based on these shortcomings, we identify three hard requirements for practical decision support in multi-objective settings: personalization; controllable output size; and flexibility in preference specification. We argue that no previous study has effectively satisfied all three requirements, and propose two new operators (ORD and ORU) to bridge that gap. Our qualitative analysis indicates that they offer a novel type of support, distinct from past practices. Also, our experiments demonstrate that our algorithms deliver practical and scalable performance. Future work could explore the directions listed in Section 8.4 and/or consider ORD/ORU in highly skewed or sparse datasets, where multi-objective querying may be meaningful in higher dimensions too.

REFERENCES

- [1] [n.d.]. Hotel dataset. <https://www.hotels-base.com>.
- [2] [n.d.]. House dataset. <https://www.ipums.org>.
- [3] [n.d.]. lp_solve. <http://lpsolve.sourceforge.net/5.5/>.
- [4] [n.d.]. NBA dataset. <https://www.basketball-reference.com>.
- [5] [n.d.]. ORD and ORU implementation. <https://github.com/ghlkm/ordu>.
- [6] [n.d.]. TripAdvisor dataset. <https://www.cs.virginia.edu/~hw5x/dataset.html>.
- [7] Lyublena Antova, Thomas Jansen, Christoph Koch, and Dan Olteanu. 2008. Fast and Simple Relational Processing of Uncertain Data. In *ICDE*. 983–992.
- [8] Abolfazl Asudeh, Azade Nazi, Nan Zhang, and Gautam Das. 2017. Efficient Computation of Regret-ratio Minimizing Set: A Compact Maxima Representative. In *SIGMOD Conference*. 821–834.
- [9] Abolfazl Asudeh, Azade Nazi, Nan Zhang, Gautam Das, and H. V. Jagadish. 2019. RRR: Rank-Regret Representative. In *SIGMOD Conference*. 263–280.
- [10] C. Bradford Barber and Hannu Huhdanpaa. [n.d.]. qhull. <http://www.qhull.org>.
- [11] Ilaria Bartolini, Zhenjie Zhang, and Dimitris Papadias. 2011. Collaborative Filtering with Personalized Skylines. *IEEE Trans. Knowl. Data Eng.* 23, 2 (2011), 190–203.
- [12] Norbert Beckmann, Hans-Peter Kriegel, Ralf Schneider, and Bernhard Seeger. 1990. The R^* -Tree: An Efficient and Robust Access Method for Points and Rectangles. In *SIGMOD Conference*. 322–331.
- [13] Jon Louis Bentley, H. T. Kung, Mario Schkolnick, and Clark D. Thompson. 1978. On the Average Number of Maxima in a Set of Vectors and Applications. *J. ACM* 25, 4 (1978), 536–543.
- [14] Mark de Berg, Otfried Cheong, Marc van Kreveld, and Mark Overmars. 2008. *Computational geometry: algorithms and applications*. Springer-Verlag TELOS.

- [15] Stephan Börzsönyi, Donald Kossmann, and Konrad Stocker. 2001. The Skyline Operator. In *ICDE*. 421–430.
- [16] Chee Yong Chan, H. V. Jagadish, Kian-Lee Tan, Anthony K. H. Tung, and Zhenjie Zhang. 2006. Finding k-dominant skylines in high dimensional space. In *SIGMOD Conference*. 503–514.
- [17] Chee Yong Chan, H. V. Jagadish, Kian-Lee Tan, Anthony K. H. Tung, and Zhenjie Zhang. 2006. On High Dimensional Skylines. In *EDBT*. 478–495.
- [18] Yuan-Chi Chang, Lawrence Bergman, Vittorio Castelli, Chung-Sheng Li, Ming-Ling Lo, and John R. Smith. 2000. The Onion Technique: Indexing for Linear Optimization Queries. In *SIGMOD Conference*. 391–402.
- [19] Sean Chester, Alex Thomo, S. Venkatesh, and Sue Whitesides. 2014. Computing k-Regret Minimizing Sets. *PVLDB* 7, 5 (2014), 389–400.
- [20] Jan Chomicki, Paolo Ciaccia, and Niccolò Meneghetti. 2013. Skyline queries, front and back. *SIGMOD Rec.* 42, 3 (2013), 6–18.
- [21] Paolo Ciaccia and Davide Martinenghi. 2017. Reconciling Skyline and Ranking Queries. *PVLDB* 10, 11 (2017), 1454–1465.
- [22] Graham Cormode, Feifei Li, and Ke Yi. 2009. Semantics of Ranking Queries for Probabilistic Data and Expected Ranks. In *ICDE*. 305–316.
- [23] Gautam Das, Dimitrios Gunopulos, Nick Koudas, and Nikos Sarkas. 2007. Ad-hoc top-k query answering for data streams. In *VLDB*. 183–194.
- [24] James S. Dyer and Rakesh K. Sarin. 1979. Measurable Multiattribute Value Functions. *Oper. Res.* 27, 4 (1979), 810–822.
- [25] Johannes Fürnkranz and Eyke Hüllermeier (Eds.). 2010. *Preference Learning*. Springer.
- [26] Yunjun Gao, Qing Liu, Lu Chen, Gang Chen, and Qing Li. 2015. Efficient algorithms for finding the most desirable skyline objects. *Knowl. Based Syst.* 89 (2015), 250–264.
- [27] Luca Di Gaspero. [n.d.]. QuadProgpp. <https://github.com/liuq/QuadProgpp>.
- [28] Shen Ge, Leong Hou U, Nikos Mamoulis, and David W. Cheung. 2013. Efficient All Top-k Computation - A Unified Solution for All Top-k, Reverse Top-k and Top-m Influential Queries. *IEEE Trans. Knowl. Data Eng.* 25, 5 (2013), 1015–1027.
- [29] Shen Ge, Leong Hou U, Nikos Mamoulis, and David Wai-Lok Cheung. 2015. Dominance relationship analysis with budget constraints. *Knowl. Inf. Syst.* 42, 2 (2015), 409–440.
- [30] A. M. Geoffrion, J. S. Dyer, and A. Feinberg. 1972. An Interactive Approach for Multi-Criterion Optimization, with an Application to the Operation of an Academic Department. *Management Science* 19, 4-part-1 (1972), 357–368.
- [31] Parke Godfrey, Ryan Shipley, and Jarek Gryz. 2007. Algorithms and analyses for maximal vector computation. *VLDB J.* 16, 1 (2007), 5–28.
- [32] Donald Goldfarb and Ashok U. Idnani. 1983. A numerically stable dual method for solving strictly convex quadratic programs. *Math. Program.* 27, 1 (1983), 1–33.
- [33] Jason Gross. 2012. Redefining Hick’s Law. *SMASHING Magazine*. <https://www.smashingmagazine.com/2012/02/redefining-hicks-law/>.
- [34] Xixian Han, Bailing Wang, Jianzhong Li, and Hong Gao. 2019. Ranking the big sky: efficient top-k skyline computation on massive data. *Knowl. Inf. Syst.* 60, 1 (2019), 415–446.
- [35] William Edmund Hick. 1952. On the Rate of Gain of Information. *Quarterly Journal of Experimental Psychology* 4, 1 (1952), 11–26.
- [36] Christian Holst. 2016. Infinite Scrolling, Pagination Or “Load More” Buttons? Usability Findings In eCommerce. *SMASHING Magazine*. <https://www.smashingmagazine.com/2016/03/pagination-infinite-scrolling-load-more-buttons/>.
- [37] Vagelis Hristidis, Nick Koudas, and Yannis Papakonstantinou. 2001. PREFER: A System for the Efficient Execution of Multi-parametric Ranked Queries. In *SIGMOD Conference*. 259–270.
- [38] Ray Hyman. 1953. Stimulus information as a determinant of reaction time. *Journal of Experimental Psychology* 45, 3 (1953), 188–196.
- [39] Ihab F Ilyas, George Beskales, and Mohamed A Soliman. 2008. A Survey of Top-k Query Processing Techniques in Relational Database Systems. *ACM Comp. Surveys* 40, 4 (2008), 11:1–11:58.
- [40] David Rios Insua and Simon French. 1991. A framework for sensitivity analysis in discrete multi-objective decision-making. *Eur. J. Oper. Res.* 54, 2 (1991), 176–190.
- [41] Kevin G. Jamieson and Robert D. Nowak. 2011. Active Ranking using Pairwise Comparisons. In *NIPS*. 2240–2248.
- [42] Thorsten Joachims. 2002. Optimizing search engines using clickthrough data. In *KDD*. 133–142.
- [43] Ralph L. Keeney and Howard Raiffa. 1976. *Decisions with Multiple Objectives: Preferences and Value Trade-Offs*. Wiley.
- [44] Werner Kießling. 2002. Foundations of Preferences in Database Systems. In *VLDB*. 311–322.
- [45] Murat Köksalan, Jyrki Wallenius, and Stanley Zionts. 2011. *Multiple Criteria Decision Making: From Early History to the 21st Century*. World Scientific Publishing Co. Pte. Ltd.
- [46] Vladlen Koltun and Christos H. Papadimitriou. 2007. Approximately dominating representatives. *Theor. Comput. Sci.* 371, 3 (2007), 148–154.

- [47] Jongwuk Lee and Seung-won Hwang. 2014. Scalable skyline computation using a balanced pivot selection technique. *Inf. Syst.* 39 (2014), 1–21.
- [48] Jongwuk Lee, Gae-won You, and Seung-won Hwang. 2009. Personalized top-k skyline queries in high-dimensional space. *Inf. Syst.* 34, 1 (2009), 45–61.
- [49] Hui Li, Tsz Nam Chan, Man Lung Yiu, and Nikos Mamoulis. 2017. FEXIPRO: Fast and Exact Inner Product Retrieval in Recommender Systems. In *SIGMOD Conference*. 835–850.
- [50] Xuemin Lin, Yidong Yuan, Qing Zhang, and Ying Zhang. 2007. Selecting Stars: The k Most Representative Skyline Operator. In *ICDE*. 86–95.
- [51] Matteo Magnani, Ira Assent, and Michael L. Mortensen. 2014. Taking the Big Picture: representative skylines based on significance and diversity. *VLDB J.* 23, 5 (2014), 795–815.
- [52] Jiri Matousek and Otfried Schwarzkopf. 1992. Linear Optimization Queries. In *ACM Symposium on Computational Geometry*. 16–25.
- [53] Denis Mindolin and Jan Chomicki. 2009. Discovering Relative Importance of Skyline Attributes. *PVLDB* 2, 1 (2009), 610–621.
- [54] Renato D. C. Monteiro and Ilan Adler. 1989. Interior path following primal-dual algorithms. part II: Convex quadratic programming. *Math. Program.* 44, 1-3 (1989), 43–66.
- [55] Kyriakos Mouratidis and Bo Tang. 2018. Exact Processing of Uncertain Top-k Queries in Multi-criteria Settings. *PVLDB* 11, 8 (2018), 866–879.
- [56] Kyriakos Mouratidis, Jilian Zhang, and HweeHwa Pang. 2015. Maximum Rank Query. *PVLDB* 8, 12 (2015), 1554–1565.
- [57] Danupon Nanongkai, Ashwin Lall, Atish Das Sarma, and Kazuhisa Makino. 2012. Interactive regret minimization. In *SIGMOD Conference*. 109–120.
- [58] Danupon Nanongkai, Atish Das Sarma, Ashwin Lall, Richard J. Lipton, and Jun (Jim) Xu. 2010. Regret-Minimizing Representative Databases. *PVLDB* 3, 1 (2010), 1114–1124.
- [59] Dimitris Papadias, Yufei Tao, Greg Fu, and Bernhard Seeger. 2005. Progressive skyline computation in database systems. *ACM Trans. Database Syst.* 30, 1 (2005), 41–82.
- [60] Jian Pei, Yidong Yuan, Xuemin Lin, Wen Jin, Martin Ester, Qing Liu, Wei Wang, Yufei Tao, Jeffrey Xu Yu, and Qing Zhang. 2006. Towards multidimensional subspace skyline analysis. *ACM Trans. Database Syst.* 31, 4 (2006), 1335–1381.
- [61] Li Qian, Jinyang Gao, and H. V. Jagadish. 2015. Learning User Preferences By Adaptive Pairwise Comparison. *PVLDB* 8, 11 (2015), 1322–1333.
- [62] Atish Das Sarma, Ashwin Lall, Danupon Nanongkai, Richard J. Lipton, and Jun (Jim) Xu. 2011. Representative skylines using threshold-based preference distributions. In *ICDE*. 387–398.
- [63] Timos K. Sellis, Nick Roussopoulos, and Christos Faloutsos. 1987. The R+-Tree: A Dynamic Index for Multi-Dimensional Objects. In *VLDB*. 507–518.
- [64] Mohamed A. Soliman, Ihab F. Ilyas, Davide Martinenghi, and Marco Tagliasacchi. 2011. Ranking with uncertain scoring functions: semantics and sensitivity measures. In *SIGMOD Conference*. 805–816.
- [65] Ralph E. Steuer. 1986. *Multiple Criteria Optimization: Theory, Computation and Application*. Wiley.
- [66] Yufei Tao, Ling Ding, Xuemin Lin, and Jian Pei. 2009. Distance-Based Representative Skyline. In *ICDE*. 892–903.
- [67] Yufei Tao, Vagelis Hristidis, Dimitris Papadias, and Yannis Papakonstantinou. 2007. Branch-and-bound processing of ranked queries. *Inf. Syst.* 32, 3 (2007), 424–445.
- [68] Yufei Tao, Xiaokui Xiao, and Jian Pei. 2007. Efficient Skyline and Top-k Retrieval in Subspaces. *IEEE Trans. Knowl. Data Eng.* 19, 8 (2007), 1072–1088.
- [69] Eleftherios Tiakas, Apostolos N. Papadopoulos, and Yannis Manolopoulos. 2011. Progressive processing of subspace dominating queries. *VLDB J.* 20, 6 (2011), 921–948.
- [70] Akrivi Vlachou and Michalis Vazirgiannis. 2010. Ranking the sky: Discovering the importance of skyline points through subspace dominance relationships. *Data Knowl. Eng.* 69, 9 (2010), 943–964.
- [71] Hongning Wang, Yue Lu, and Chengxiang Zhai. 2010. Latent aspect rating analysis on review text data: a rating regression approach. In *KDD*. 783–792.
- [72] Hongning Wang, Yue Lu, and Chengxiang Zhai. 2011. Latent aspect rating analysis without aspect keyword supervision. In *KDD*. 618–626.
- [73] Tian Xia, Donghui Zhang, and Yufei Tao. 2008. On Skylining with Flexible Dominance Relation. In *ICDE*. 1397–1399.
- [74] Min Xie, Raymond Chi-Wing Wong, and Ashwin Lall. 2019. Strongly Truthful Interactive Regret Minimization. In *SIGMOD Conference*. 281–298.
- [75] Min Xie, Raymond Chi-Wing Wong, and Ashwin Lall. 2020. An experimental survey of regret minimization query and variants: bridging the best worlds between top-k query and skyline query. *VLDB J.* 29, 1 (2020), 147–175.
- [76] Min Xie, Raymond Chi-Wing Wong, Jian Li, Cheng Long, and Ashwin Lall. 2018. Efficient k-Regret Query Algorithm with Restriction-free Bound for any Dimensionality. In *SIGMOD Conference*. 959–974.

- [77] Ke Yi, Feifei Li, George Kollios, and Divesh Srivastava. 2008. Efficient Processing of Top-k Queries in Uncertain Databases with x-Relations. *IEEE Trans. Knowl. Data Eng.* 20, 12 (2008), 1669–1682.
- [78] Man Lung Yiu and Nikos Mamoulis. 2009. Multi-dimensional top-*k* dominating queries. *VLDB J.* 18, 3 (2009), 695–718.
- [79] Albert Yu, Pankaj K. Agarwal, and Jun Yang. 2016. Top-k Preferences in High Dimensions. *IEEE Trans. Knowl. Data Eng.* 28, 2 (2016), 311–325.
- [80] Sepanta Zeighami and Raymond Chi-Wing Wong. 2019. Finding Average Regret Ratio Minimizing Set in Database. In *ICDE*. 1722–1725.
- [81] Jilian Zhang, Kyriakos Mouratidis, and HweeHwa Pang. 2014. Global immutable region computation. In *SIGMOD Conference*. 1151–1162.
- [82] Jiping Zheng and Chen Chen. 2020. Sorting-Based Interactive Regret Minimization. In *APWeb-WAIM*. 473–490.