

1 Skyline Computation with Noisy Comparisons^{*}

2 Benoît Groz¹, Frederik Mallmann-Trenn^{2,3}, Claire Mathieu^{2,3}, and
3 Victor Verdugo³

4 ¹ Université Paris-Saclay, CNRS, LRI
5 `benoit.groz@lri.fr`

6 ² King's College London
7 `frederik.mallmann-trenn@kcl.ac.uk`

8 ³ CNRS & IRIF
9 `claire.m.mathieu@gmail.com`

10 ⁴ London School of Economics
11 `v.verdugo@lse.ac.uk`

12 **Abstract.** Given a set of n points in a d -dimensional space, we seek to
13 compute the *skyline*, i.e., those points that are not strictly dominated by
14 any other point, using few comparisons between elements. We adopt the
15 noisy comparison model ([13]) where comparisons fail with constant proba-
16 bility and confidence can be increased through independent repetitions of a
17 comparison. In this model motivated by Crowdsourcing applications, Groz
18 & Milo [16] show three bounds on the query complexity for the skyline
19 problem. We improve significantly on that state of the art and provide
20 two output-sensitive algorithms computing the skyline with respective
21 query complexity $O(nd\log(dk))$ and $O(ndk\log(k))$, where k is the size of
22 the skyline. These results are tight for low dimensions.

23 **Keywords:** Skyline · Noisy comparisons · Fault-tolerance · CrowdSourc-
24 ing.

25 1 Introduction

26 Skylines have been studied extensively, since the 1960s in statistics [5], then in algo-
27 rithms and computational geometry [21] and in databases [6, 10, 14, 20]. Depending
28 on the field of research, the *skyline* is also known as the set of *maximum vectors*,
29 the *dominance frontier*, *admissible points*, or *Pareto frontier*. The skyline of a set
30 of points consists of those points which are not strictly dominated by any other
31 point. A point p is *dominated* by another point q if $p_i \leq q_i$ for every coordinate
32 (attribute or dimension) i . It is *strictly dominated* if in addition the inequality is
33 strict for at least one coordinate; see Figure 1.

34 *Noisy comparison model, and parameters.* In many contexts, comparing attributes
35 is not straightforward. Consider the example of finding *optimal* cities from [16].

36 *To compute the skyline with the help of the crowd we can ask people ques-*
37 *tions of the form “is the education system superior in city x or city y ?” or*

* This paper is NOT eligible for the best student paper award.

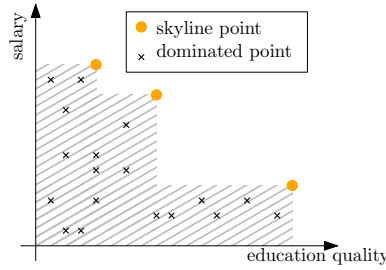


Fig. 1: Given a set of points X , the goal is to find the set of *skyline points*, i.e., points are not dominated by any other points.

38 “can I expect a better salary in city x or city y ”. Of course, people are likely
 39 to make mistakes, and so each question is typically posed to multiple people.
 40 Our objective is to minimize the number of questions that need to be issued
 41 to the crowd, while returning the correct skyline with high probability.

42 Thus, much attention has recently been given to computing the skyline when
 43 information about the underlying data is uncertain [24], and comparisons may give
 44 erroneous answers. In this paper we investigate the complexity of computing skylines
 45 in the noisy comparison model, which was considered in [16] as a simplified model for
 46 crowd behaviour: we assume queries are of the type *is the i -th coordinate of point p*
 47 *(strictly) smaller than that of point q ?*, and the outcome of each such query is indepen-
 48 dently correct with probability greater than some constant better than $1/2$ (for defini-
 49 finiteness we assume probability $2/3$). As a consequence, our confidence on the relative
 50 order between p and q can be increased by repeatedly querying the pair on the same
 51 coordinate. Our complexity measure is the number of comparison queries performed.

52 This noisy comparison model was introduced in the seminal paper [13] and has
 53 been studied in [7, 16]. There are at least 2 straightforward approaches to reduce
 54 problems in this model to the noiseless comparison setting. One approach is to take
 55 any “noiseless” algorithm and repeat each of its comparisons $\log(f(n))$ times, where
 56 n is the input size and $f(n)$ is the complexity of the algorithm. The other approach is
 57 to sort the n items in all d dimensions at a cost of $nd\log(nd)$, then run some noiseless
 58 algorithm based on the computed orders. The algorithms in [13, 16] and this paper
 59 thus strive to avoid the logarithmic overhead of these straightforward approaches.

60 Three algorithms were proposed in [16] to compute skylines with noisy compar-
 61 isons. Figure 2 summarizes their complexity and the parameters we consider. The
 62 first algorithm is the reduction through sorting discussed above. But skylines often
 63 contain only a small fraction of the input items (points), especially when there are
 64 few attributes to compare (low dimension). This leads to more efficient algorithms
 65 because smaller skylines are easier to compute. Therefore, [16] and the present
 66 paper investigate the complexity of computing skylines expressed as a function of
 67 three parameters: $n = |X|$, the number of input points; d , the number of dimensions;
 68 and $k = |\text{skyline}(X)|$, the size of the skyline (output). There is a substantial gap
 69 between the lower bounds and the upper bounds achieved by the skyline algorithms

70 in [16]. In particular, the authors raised the question whether the skyline could
 71 be computed in $o(nk)$ for any constant k when $k \ll n$. In this paper, we tighten
 72 the gap between the lower and upper bounds and settle this open question.

73 *Contributions.* We propose 2 new algorithms that compute skylines with probability
 74 at least $1 - \delta$ and establish a lower bound:

- 75 – Algorithm **SkyLowDim-Search**(X, δ) computes the skyline in $O(nd \log(dk/\delta))$
 76 query complexity and $O(nd \log(dk/\delta) + ndk)$ overall running time.
- 77 – Algorithm **SkyHighDim-Search**(X, δ) computes the skyline in $O(ndk \log(k/\delta))$
- 78 – $\Omega(nd \log k)$ queries are necessary to compute the skyline.
- 79 – Additionally, we show that Algorithm **SkyLowDim-Search** can be adapted
 80 to compute the skyline with $O(nd \log(dk))$ comparisons in the noiseless setting.

81 Our first algorithm answers positively the above question from [16]. Together with
 82 the lower bound, we thus settle the case of low dimensions, i.e., when there is a
 83 constant c such that $d \leq k^c$. Our 2 skyline algorithms both shave off a factor k from
 84 the corresponding bounds in the state of the art [16], as illustrated in Figure 2
 85 with respect to query complexity. We point out that **SkyLowDim-Search** is a
 86 randomized algorithm: it needs to sample the input. In our bounds we guarantee
 87 that the combined probability of incorrect comparisons and poor sample choice
 88 is low: this is because we tailor the sample size to the desired accuracy. But having
 89 a randomized algorithm is still a weakness in the sense that our approach cannot
 90 yield a "trust-preserving" algorithm: even in the extreme case where comparison
 91 queries all return a correct answer (noiseless setting), our algorithm still relies on
 92 sampling and therefore has some probability of failing to return the skyline within
 93 the running time bound. However, we show that for the specific case of the noiseless
 94 setting, our algorithm can be adapted to compute the skyline in $O(nd \log(dk))$.

95 As a subroutine for our algorithms, we developed a new algorithm to eval-
 96 uate disjunctions of boolean variables with noise ("OR"). We believe algorithm
 97 **NoisyFirstTrue** to be interesting in its own right: it returns the first positive
 98 variable in input order, with a running time that scales linearly with the position
 99 of that variable in the input order.

[16]	$O(nd \log(nd/\delta))^\dagger$	$O(ndk \log(dk/\delta))$	$O(ndk^2 \log(k/\delta))$	d : dimension n : # input points k : # skyline points δ : error rate tolerated
this paper	—	$O(nd \log(dk/\delta))^\dagger$	$O(ndk \log(k/\delta))$	
best when:	$k \in \Omega(n)$	$d \leq k^c \leq n$	$k \ll d$	

Fig. 2: Query complexity of skyline algorithms depending on the values of k . For † -labeled bounds, the running time is larger than the number of queries.

100 *Technical core of our algorithms.* The algorithm underlying the two bounds for
 101 $k \ll n$ in [16] recovers the skyline points one by one. It iteratively adds to the skyline

102 the maximum point, in lexicographic order, among those not dominated by the
 103 skyline points already found. ⁵ However, the algorithm in [16] essentially considers
 104 the whole input for each iteration. Our two algorithms, on the opposite, can identify
 105 and thus discard some dominated points early. The idea behind our algorithm
 106 **SkylineHighDim** is that it is more efficient to separate the two tasks: (i) finding a
 107 point p not dominated by the skyline points already found, on the one hand, and (ii)
 108 computing a maximum point (in lexicographic order) *among those dominating p* , on
 109 the other hand. Whenever a point is considered for step (i) but fails to satisfy that
 110 requirement, the point can be discarded definitively. The $O(ndk)$ skyline algorithm
 111 from [11] for the noiseless setting also decomposes the two tasks, although the point
 112 they add to the skyline in each of the k iteration is not the lexicographic maximum.

113 Our algorithm **SkylineLowDim** can be viewed as a 2-steps algorithm where
 114 the first step prunes a huge fraction of dominated points from the input through
 115 discretization, and the second step applies a cruder algorithm on the surviving
 116 points. We partition the input into buckets for discretization, identify “skyline
 117 buckets” and eliminate all points in dominated buckets. The bucket boundaries
 118 are defined by sampling the input points and sorting all sample points in each
 119 dimension. In the noisy comparison model, the approach of sampling the input
 120 for some kind of discretization was pioneered in [7] for selection problems, but
 121 with rather different techniques and objectives. One interesting aspect of our
 122 discretization is that a fraction of the input will be, due to the low query complexity,
 123 incorrectly discretized yet we are able to recover the correct skyline.

124 Our lower bound constructs a technical reduction from the problem of identifying
 125 null vectors among a collection of vectors, each having at most one non-zero coordi-
 126 nate. That problem can be studied using a two-phase process inspired from [13].

127 *Related work.* The noisy comparison model was considered for sorting and searching
 128 objects [13]. While any algorithm for that model can be reduced to the noiseless
 129 comparison model at the cost of a logarithmic factor (boosting each comparison so
 130 that by union bound all the comparisons required are correct), [13] shows that this
 131 additional logarithmic factor can be spared for sorting and for maxima queries,
 132 though it cannot be spared for median selection. [25], [15] and [7] investigate the
 133 trade-off between the total number of queries and the number of rounds for (variants
 134 of) top-k queries in the noisy comparison model and some other models. The noisy
 135 comparison model has been refined in [12] for top-k queries, where the probability of
 136 incorrect answers to a comparison increase with the distance between the two items.

137 Other models for uncertain data have also been considered in the literature: in
 138 some, the location of each point is determined by a probability distribution over a set
 139 of locations, whereas in other models the data is incomplete [18, 23]. Some previous
 140 work [2, 26] model uncertainty about the output by computing a ρ -skyline: points
 141 having probability at least ρ to be in the skyline. We refer to [4] for skyline com-
 142 putation using the crowd and [22] for a survey in crowdsourced data management.

143 Our paper aims to establish the worst-case number of comparisons required
 144 to compute skylines with output-sensitive algorithms, i.e., when the cost is

⁵ The difference between those two bounds is due to different subroutines to check dominance.

145 parametrized by the size of the result set. While one of our algorithm is ran-
 146 domized, we do not make any further assumption on the input (we do not assume
 147 input points are uniformly distributed, for instance). In the classic *noiseless* com-
 148 parison model, the problem of computing skylines has received a large amount of
 149 attention [6, 19, 21]. For any constant d , [19] show that skylines can be computed
 150 in $O(n \log^{d-2} k)$. When $d \in \{2, 3\}$, Barbay et al. [3] provide stronger efficiency
 151 guarantees with “instance-optimal” algorithms. [9] investigates the constant factor
 152 for the number of comparisons required to compute skyline, when $d \in \{2, 3\}$. The
 153 technique does not seem to generalize to arbitrary dimensions, and the authors ask
 154 among open problems whether arbitrary skylines can be computed with fewer than
 155 $dn \log n$ comparisons. To the best of our knowledge, our $O(nd \log(dk))$ is the first
 156 non-trivial output-sensitive upper bound that improves on the folklore $O(dnk)$
 157 for computing skylines in arbitrary dimensions. Many other algorithms have been
 158 proposed that fit particular settings (big data environment, particular distributions,
 159 etc), as evidenced in the survey [17], but those works are further from ours as
 160 they generally do not investigate the asymptotic number of comparisons. Other
 161 skyline algorithms in the literature for the noiseless setting have used bucketing. In
 162 particular, [1] computes the skyline in a massively parallel setting by partitioning
 163 the input based on quantiles along each dimension. This means they define similar
 164 buckets to ours, and they already observed that the buckets that contain skyline
 165 points are located in hyperplanes around the “bucket skyline”, and therefore those
 166 buckets only contain a small fraction of the whole input.

167 *Organization.* In Section 2, we recall standard results about the noisy comparison
 168 model and introduce some procedure at the core of our algorithms. Section 3
 169 introduces our algorithm for high dimensions (Theorem 4) and Section 4 introduces
 170 the counterpart for low dimensions (Theorem 6). Section 5 establishes our lower
 171 bound (Theorem 7).

172 2 Preliminaries

173 The complexity measured is the number of comparisons in the worst case. Whenever
 174 the running time and the number of comparisons differ, we will say so. With respect
 175 to the probability of error, our algorithms are supposed to fail with probability
 176 at most δ . Following standard practice we only care to prove that our algorithms
 177 have error in $O(\delta)$: 5δ , for instance. This is because we can run the algorithm with
 178 an adjusted value for the parameter ($\delta' = \delta/5$) while maintaining the asymptotic
 179 complexity of our algorithms.

180 Given two points, $p = (p_1, p_2, \dots, p_d)$ and $q = (q_1, q_2, \dots, q_d)$ point p is *lexicograph-*
 181 *ically* smaller than q , denoted by $p \leq_{\text{lex}} q$, if $p_i < q_i$ for the first i where p_i and q_i
 182 differ. If there is no such i , meaning that the points are identical, we use the id of
 183 the points in the input as a tie-breaker, ensuring that we obtain a total order.

184 In the noisy comparison model, we call an algorithm *trust-preserving* ([16]⁶) if for
 185 every $\delta < 1/3$ it is guaranteed to return the correct answer with probability at least

⁶ [25] calls such algorithms *fault-tolerant*

186 $1 - \delta$ whenever the input comparisons are correct with probability at least $1 - \delta$. We
 187 next describe and name algorithms that we use as subroutines to compute skylines.

188 Algorithm **NoisySearch** takes as input an element y , an ordered list (y_1, y_2, \dots, y_m) ,
 189 accessible by comparisons that each have error probability at most p , and a pa-
 190 rameter δ . The goal is to output the interval $I = (y_{i-1}, y_i]$ such that $y \in I$.

191 Algorithm **NoisySort** relies on **NoisySearch** to solve the **noisy sort prob-**
 192 **lem**. It takes as input an unordered set $Y = \{y_1, y_2, \dots, y_m\}$, and a parameter δ . The
 193 goal is to output an ordering of Y that is the correct non-decreasing sorted order.

194 Algorithm **NoisyMax** returns the maximum item in the unordered set Y
 195 whose elements can be compared, but we will rather use another variant: algorithm
 196 **MaxLex** takes as input an unordered set $Y = \{y_1, y_2, \dots, y_m\}$, a point x and a
 197 parameter δ . The goal is to output the maximum point in lexicographic order
 198 among those that dominate x . Algorithm **SetDominates** is the boolean version
 199 whose goal is to output whether there exists a point in Y that dominates x .

200 Algorithm **NoisyOr** takes as input a list (y_1, y_2, \dots, y_m) of boolean elements
 201 that can be compared to **true** with error probability at most p (typically the result
 202 of some comparison or subroutines such as **SetDominates**). The original goal
 203 was to output whether at least one of the elements is true. But we rather adopt the
 204 enhanced version discussed in [16] which solves the *first positive variable problem*.
 205 The goal is to output the index of the first element with value **true** (and $m + 1$,
 206 which we assimilate to **false**, if there are none).

207 **Theorem 1** ([13], [16]). *When the input comparisons have error probability at*
 208 *most $p = 1/3$, the table below lists the number of comparisons performed by the*
 209 *algorithms to return the correct answer with success probability $1 - \delta$:*

Algorithm	NoisyOr	NoisyMax	NoisySort	NoisySearch	SetDominates	MaxLex
Comparisons	$O(m \log \frac{1}{\delta})$	$O(m \log \frac{1}{\delta})$	$O(m \log \frac{m}{\delta})$	$O(\log \frac{m}{\delta})$	$O(m d \log \frac{1}{\delta})$	$O(m d \log \frac{1}{\delta})$

210 *Furthermore, these algorithms are trust preserving. This means that when the*
 211 *input comparisons already have error probability at most δ , we can discard from*
 212 *the complexity the dependency in δ (replacing δ by some constant).*

213 We first refine the complexity of **NoisyOr** and call **NoisyFirstTrue** the re-
 214 fined algorithm which only spends constant time per variable it processes, and
 215 which identifies correctly all processed variables with high probability.

216 **Theorem 2.** *Algorithm **NoisyFirstTrue** solves the first positive variable prob-*
 217 *lem with success probability $1 - \delta$ in $O(j \cdot \log(1/\delta))$ where j is the index output by*
 218 *the algorithm. Furthermore, the algorithm is trust-preserving.*

219 *Proof.* The proof, left for the Appendix, shows that the error (resp. the cost) of
 220 the whole algorithm is dominated by the error (resp. the cost) of the last iteration.

Algorithm NoisyFirstTrue(x_1, \dots, x_n, δ) (see Theorem 2)
input: $\{x_1, \dots, x_n\}$ set of boolean random variables, δ error probability
output: the index j of the first positive variable, or $m+1$ (=false).

```

1:  $i \leftarrow 1$ 
2:  $\delta' \leftarrow \delta/2$ 
3: while  $i \leq n$  do
4:    $j \leftarrow \text{NoisyOr}(x_1, \dots, x_i, \delta')$ 
5:   if  $\text{CheckVar}(x_j, \delta'/2^i)$  then
6:     return  $j$ 
7:   else
8:      $i \leftarrow 2 \cdot i$ 
9: return false

```

221 3 Skyline computation in high dimension

222 We first assume that an estimate \hat{k} of k is known in advance. We will show afterwards
223 how we can lift that assumption.

224 We are now ready to give the full description of our algorithm **SkylineHighDim**.

Algorithm SkylineHighDim(k, X, δ) (see Theorem 3)
input: $X = \{p_1, \dots, p_n\}$ set of points, \hat{k} upper bound on skyline size, δ error probability
output: $\min\{\hat{k}, \text{skyline}(X)\}$ skyline points w.p. $1 - \delta$

```

1: Initialize  $S \leftarrow \emptyset$ ,  $i \leftarrow 1$ 
2: while  $i \neq -1$  and  $|S| < \hat{k}$  do
3:    $i' \leftarrow$  index of the first point  $p_{i'}$  not dominated by current skyline points.7
4:   {Find a skyline point dominating  $p_i$ }
5:   Compute  $p^* \leftarrow \text{MaxLex}(p_{i'}, \{p_i, \dots, p_n\}, \delta/(2\hat{k}))$ 
6:    $S \leftarrow S \cup \{p^*\}$ 
7:    $i \leftarrow i'$ 
7: Output  $S$ 

```

225 **Theorem 3.** *Given $\delta \in (0, 1/2)$ and a set X of data items, **SkylineHighDim**(X, δ)*
226 *outputs a subset of X which, with probability at least $1 - \delta$, is the first $\min(|X|, \hat{k})$*
227 *skyline points. The running time and number of queries is $O(nd\hat{k} \log(\hat{k}/\delta))$.*

228 *Proof.* Each iteration through the loop adds a point to the skyline S with probability
229 of error at most δ/\hat{k} . The final result is therefore correct with success probability $1 - \delta$.

⁷ This point can be computed using algorithm **NoisyFirstTrue** on the boolean variables:
 $\neg \text{SetDominates}(S, p_i, \delta/(2\hat{k})), \dots, \neg \text{SetDominates}(S, p_n, \delta/(2\hat{k}))$, where we denote
by \neg the negation. This means that $\neg \text{SetDominates}(S, p_n, \delta/(2\hat{k}))$ returns true when
the procedure **SetDominates**($S, p_n, \delta/(2\hat{k})$) indicates that p_n is not dominated.

230 The complexity is $O((i' - i) * d \hat{k} \log(\hat{k}/\delta))$ to find a non-dominated point $p_{i'}$ at line
 231 3, and $O(nd \log(\hat{k}/\delta))$ to compute the maximal point above $p_{i'}$ at line 4. Summing
 232 over all iterations, the running time and number of queries is $O(nd \hat{k} \log(\hat{k}/\delta))$.

233 Algorithm **SkylineHighDim**(X, δ) can only return the skyline in $O(ndk \log(k/\delta))$
 234 if it is provided with a good estimate of the skyline cardinality $\hat{k} \in O(k)$. We next
 235 show how to guarantee the complexity by trying a sequence of successive values
 236 for \hat{k} . The successive values in the sequence grow exponentially to prevent failed
 237 attempts from penalizing the complexity.

Algorithm SkyHighDim-Search(X, δ) (see Theorem 4)

input: X set of points, δ error probability

output: skyline(X) w.p. $1 - \delta$

1: Initialize $j \leftarrow 0, \hat{k} \leftarrow 1$
 2: **repeat**
 3: $j \leftarrow j + 1; \hat{k} \leftarrow 2\hat{k}; S \leftarrow \mathbf{SkylineHighDim}(\hat{k}, X, \delta/2^j)$
 4: **until** $|S| < \hat{k}$
 5: Output S

238 **Theorem 4.** Given $\delta \in (0, 1/2)$ and a set X of data items, **SkyHighDim-Search**(X, δ)
 239 outputs a subset of X which, with probability at least $1 - \delta$, is the skyline. The running
 240 time and number of queries is $O(ndk \log(k/\delta))$.

241 *Proof.* The proof is relatively straightforward and left for the Appendix.

242 4 Skyline computation in low dimension

243 Let us first sketch our algorithm **SkylineLowDim**(k, X, δ). The algorithm works
 244 in 3 phases. The first phase partitions input points in buckets. We sort the i -th
 245 coordinate of a random sample to define $s + 1$ intervals in each dimension $i \in [d]$,
 246 hence $(s + 1)^d$ buckets, where each bucket is a product of intervals of the form $\prod_i I_i$;
 247 then we assign each point p of X to a bucket by searching in each dimension for
 248 the interval I_i containing p_i .

249 The second phase eliminates irrelevant buckets: those that are dominated by
 250 some non-empty bucket and therefore have no chance of containing a skyline
 251 point. With high probability the bucketization obtained from the first phase will
 252 be "accurate enough" so that we will be able to identify efficiently the irrelevant
 253 buckets, and will also guarantee that the points in the remaining buckets form
 254 a small fraction of the input (provided k is small). In phase 3, we thus solve the
 255 skyline problem on a much smaller dataset, calling Algorithm **SkylineHighDim**
 256 to find the skyline of the remaining points. ⁸

⁸ Alternatively, one could use an algorithm provided by Groz and Milo [16], it is only important that the size of the input set is reduced to n/k to cope with the larger runtime of the mentioned algorithms.

257 4.1 Emptiness testing, and domination relationships between buckets

258 Our bucketization does not guarantee that all points are assigned to the proper
 259 bucket. In particular, empty buckets may erroneously be assumed to contain
 260 some points. To drop the irrelevant buckets, we thus design a subroutine **First-**
 261 **Nonempty-Bucket** that processes a list of buckets, and returns the first bucket
 262 that really contains at least one point. Incidentally, we will not double-check the
 263 emptiness of every bucket using this procedure, but will only check those that may
 264 possibly belong to the skyline: those that we will define more formally as buckets
 265 of type (i), (ii) and (iv) in the proof of Theorem 5. We could not afford to "fix" the
 266 whole assignment as it may contain too many buckets.

267 In the **First-Nonempty-Bucket** problem, the input is a sequence of pairs
 268 $[(B_1, X_1), \dots, (B_n, X_n)]$ where B_i is a bucket and X_i is a set of points. The goal is
 269 to return the first i such that $B_i \cap X_i \neq \emptyset$ with success probability $1 - \delta$. The test
 270 $B_i \cap X_i \neq \emptyset$ can be formulated as a DNF with $|X_i|$ conjunctions of $O(d)$ boolean
 271 variables each. To solve **First-Nonempty-Bucket**, we can flatten the formulas
 272 of all buckets into a large DNF with conjunctions of $O(d)$ boolean variables (one
 273 conjunction per bucket point). Using **NoisyFirstTrue** to compute the first true
 274 conjunction (while keeping tracks of which point belongs to which bucket with
 275 pointers) yields the following complexity:

276 **Lemma 1.** *Algorithm **FirstBucket** $([(B_1, X_1), \dots, (B_n, X_n)], \delta)$ solves **First-Nonempty-**
 277 **Bucket** in $O(\sum_{i \leq j} d \cdot |X_j| \log(1/\delta))$ with success probability $1 - \delta$, where j is the
 278 index returned by the algorithm (the algorithm is trust-preserving).*

279 In the second phase, Algorithm **SkylineLowDim** (k, X, δ) uses elimination. To
 280 manage ties, we need to distinguish two kinds of intervals: the trivial intervals that
 281 match a sample coordinate: $I = [x, x]$ and the non-trivial intervals $I =]a, b[$ ($a < b$)
 282 contained between samples (or above the largest sample, or below the smallest
 283 sample). To compare easily those intervals, we adopt the convention that for a non-
 284 trivial interval $I =]a, b[$, $\min I = a + \epsilon$ and $\max I = b - \epsilon$ for some infinitesimal $\epsilon > 0$:
 285 $\epsilon = (b - a)/3$ would do. We say that a bucket $B = \prod_i I_i$ is *dominated* by a different
 286 bucket $B' = \prod_i I'_i$ if in every dimension $\max I_i \leq \min I'_i$. Equivalently: we say that B'
 287 dominates B if every point (whether in the dataset or hypothetical) in B' dominates
 288 every point in B . The idea is that no skyline point belongs to a bucket dominated by
 289 a non-empty bucket. See Figure 3 for an illustration. We observe that the relative po-
 290 sition of buckets is known by construction, so deciding whether a bucket dominates
 291 another one may require time $O(d)$ but does not require any comparison query.

292 4.2 Properties satisfied by the bucket assignments

293 **Theorem 5.** *Given $\delta \in (0, 1/2)$ and a set X of data items, **SkylineLowDim** (X, δ)
 294 outputs a subset of X which, with probability at least $1 - \delta$, is the first $\min(|X|, \hat{k})$
 295 skyline points. The number of queries is $O(nd \log(d\hat{k}/\delta))$. The running time is
 296 $O(nd \log(d\hat{k}/\delta) + nd \cdot \min(\hat{k}, |\text{skyline}(X)|))$*

⁹ Note that X can contain points sharing the same coordinate meaning that the S_i are not necessarily distinct.

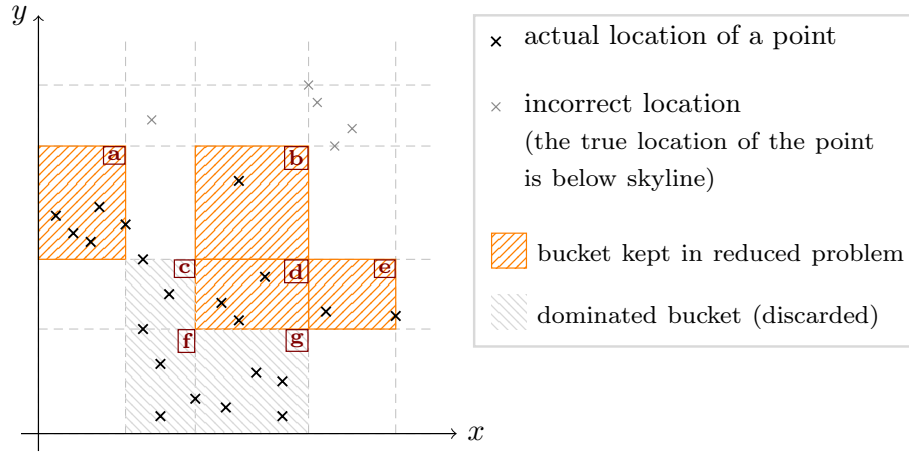


Fig. 3: An illustration of the bucket dominance and its role in **SkylineLowDim**. Here bucket b dominates c and f but not a , d , e or g . Buckets c, f, g are dominated by some non-empty bucket and therefore cannot contain a skyline point. Bucket a does not contain a skyline point, but this cannot be deduced from the bucket assignments, therefore points in bucket a are passed on to the reduced problem. In this figure we may assume to simplify that a bucket contains its upper boundary. But in our algorithm bucket a would actually contain only the 4 leftmost points, and the fifth point would belong to a distinct bucket with a trivial interval on x . . .

297 *Proof.* The proof, left for the Appendix, first shows by Chernoff bounds that the
 298 assignment satisfies with high probability some key properties: (1) few points are
 299 erroneously assigned to incorrect buckets (2) the skyline points are assigned to the
 300 correct bucket, and (3) there are at most $O(n/(d\hat{k}^2))$ points on any hyperplane
 301 (i.e., in buckets that are ties on some dimension). The proof then shows that:

- 302 – there are at most $O(n/\hat{k})$ points in the reduced problem. This is because those
 303 points belong to skyline buckets or buckets that are tied with a skyline bucket
 304 on at least one dimension (every other non-empty bucket is dominated), and
 305 property (3) of the assignment guarantees that the union of all such buckets
 306 has at most $O(n/\hat{k})$ points.
- 307 – the buckets above the skyline buckets which are erroneously assumed to contain
 308 points can quickly be identified and eliminated since they contain few points.

309 Algorithm **SkylineLowDim**(X, δ) can only return the skyline in $O(nd \log(dk/\delta))$
 310 if it is provided with a good estimate of the skyline cardinality: we must have $\hat{k} \geq k$
 311 and $\log(\hat{k}) \in O(\log(k))$. We next show how to guarantee the complexity by trying
 312 a sequence of successive values for \hat{k} . The successive values in the sequence grow
 313 super exponentially (similarly to [8, 16]) to prevent failed attempts from penalizing
 314 the complexity.

Algorithm SkylineLowDim(\hat{k}, X, δ) (see Theorem 5)
input: \hat{k} integer, X set of points, δ error probability**output:** $\min\{\hat{k}, |\text{skyline}(X)|\}$ points of skyline(X)**error probability:** δ

-
- 1: **if** $\hat{k}^5 \geq n$ or $d^5 \geq n$ or $(\log(1/\delta))^5 \geq n$ **then**
 - 2: Compute the skyline by sorting every dimension, as in [16]. Return that skyline.
 - 3: $\delta' \leftarrow \delta / (2d\hat{k})^5$ and $s \leftarrow d\hat{k}^2 \log(d^2\hat{k}^2/\delta')$
 - {Phase (i): bucketing}
 - 4: **for** each dimension $i \in \{1, 2, \dots, d\}$ **do**
 - 5: $S_i \leftarrow \text{NoisySort}$ (sample of X of size $s, i, \delta'/d$)
 - 6: Remove duplicates so that, with prob. $1 - \delta'/d$, the values in S_i are all distinct.⁹
 - 7: **for** each point $p \in X$ **do**
 - 8: Place p in set X_B associated to $B = \prod_{i=1}^d I_i$, with $I_i = \text{NoisySearch}(p_i, S_i, \delta'/(d\hat{k}))$.
 - 9: Drop all empty buckets (those that were assigned no point).
 - 10: Sort buckets into a sequence B_1, \dots, B_h so that each bucket comes before buckets it dominates.
 - {Phase (ii): eliminating irrelevant buckets}
 - 11: Initialize $X' \leftarrow \emptyset$, $i \leftarrow 1$
 - 12: **while** $i \neq -1$ **do**
 - 13: $i \leftarrow \text{FirstBucket}([(B_1, X_{B_1}), \dots, (B_h, X_{B_h})], \delta'/\hat{k})$
 - 14: $X' \leftarrow X' \cup X_{B_i}$
 - 15: **if** $|X'| > 8n/\hat{k}$ **then**
 - 16: Raise an error.
 - 17: Drop from B_1, \dots, B_h all buckets dominated by B_i , and also buckets B_1 to B_i .
 - {Phase (iii): solve reduced problem}
 - 18: Output **SkyHighDim-Search**(X', δ').
-

Algorithm SkyLowDim-Search(X, δ) (see Theorem 6)
input: X set of points, δ error probability**output:** skyline(X)**error probability:** δ

-
- 1: $\hat{k} \leftarrow (\lfloor d/\delta \rfloor)^2$
 - 2: **repeat**
 - 3: $\delta \leftarrow \delta/2$; $\hat{k} \leftarrow \hat{k}^2$; $S \leftarrow \text{SkylineLowDim}(\hat{k}, X, \delta)$
 - 4: **until** $|S| < \hat{k}$
 - 5: Output S
-

315 **Theorem 6.** Given $\delta \in (0, 1/2)$ and a set X of data items, **SkyLowDim-Search**(X, δ)
316 outputs a subset of X which, with probability at least $1 - \delta$, is the skyline. The number
317 of queries is $O(nd \log(dk/\delta))$. The running time is $O(nd \log(dk/\delta) + ndk)$.

318 *Proof.* For iteration j , the probability of error is $\delta/2^j$, and the cost is given by
319 Theorem 5. Consequently, we obtain the complexity we claim by summing those
320 terms over all iterations.

321 *Remark 1.* In the noiseless setting, we could adopt the same sampling approach
 322 to assign points to buckets and reduce the input size. On line 18 we could use any
 323 noiseless skyline algorithm such as the $O(ndk)$ algorithm from [11], or our own
 324 similar **SkyHighDim-Search** which can clearly run in $O(dnk)$ in the noiseless
 325 case. The cost of the bucketing phase remains $O(nd\log(\hat{dk}/\delta))$. The elimination
 326 phase becomes rather trivial since all points get assigned to their proper bucket,
 327 and therefore there is no need to check buckets for emptiness as in Line 13. By
 328 setting $\delta = 1/k$ failures are scarce enough so that the higher cost of $O(ndk)$ in
 329 case of failure is covered by the cost of an execution corresponding to a satisfying
 330 sample. Consequently, the expected query complexity is $O(nd\log(dk))$, and the
 331 running time $O(nd\log(dk) + ndk)$.

332 Better yet: we can replace random sampling with quantile selection to obtain a
 333 deterministic algorithm with the same bounds. Algorithms for the *multiple selection*
 334 problem are surveyed in [9]. Actually, our algorithm can be viewed as some kind of
 335 generalization to higher dimensions of an algorithm from [9] which assigns points
 336 to buckets before recursing, the buckets being the quantiles along one coordinate.

337 5 Skyline Lower Bound

338 **Theorem 7.** *Let A be an algorithm that computes the skyline with error probability*
 339 *at most $1/2$. Then the expected number of queries of A is $\Omega(dn\log k)$.*

340 *Proof.* The proof is left for the Appendix.

341 References

- 342 1. Afrati, F.N., Koutris, P., Suciu, D., Ullman, J.D.: Parallel skyline queries. In:
 343 15th International Conference on Database Theory, ICDT '12, Berlin, Germany,
 344 March 26-29, 2012. pp. 274–284 (2012). <https://doi.org/10.1145/2274576.2274605>,
 345 <https://doi.org/10.1145/2274576.2274605>
- 346 2. Afshani, P., Agarwal, P.K., Arge, L., Larsen, K.G., Phillips, J.M.: (approximate) uncer-
 347 tain skylines. In: Proceedings of the 14th International Conference on Database Theory.
 348 pp. 186–196. ICDT '11, ACM (2011). <https://doi.org/10.1145/1938551.1938576>,
 349 <http://doi.acm.org/10.1145/1938551.1938576>
- 350 3. Afshani, P., Barbay, J., Chan, T.M.: Instance-optimal geometric algorithms. *J. ACM*
 351 **64**(1), 3:1–3:38 (2017)
- 352 4. Asudeh, A., Zhang, G., Hassan, N., Li, C., Zaruba, G.V.: Crowdsourcing pareto-
 353 optimal object finding by pairwise comparisons. In: Proceedings of the 24th ACM
 354 International on Conference on Information and Knowledge Management. pp.
 355 753–762. ACM (2015)
- 356 5. Barndorff-Nielsen, O., Sobel, M.: On the distribution of the number of admissible
 357 points in a vector random sample. *Theory of Probability and its Applications* **11**(2),
 358 249–251 (1966), <http://search.proquest.com/docview/915869827?accountid=15867>
- 359 6. Börzsönyi, S., Kossmann, D., Stocker, K.: The skyline operator. In: Proceedings
 360 of the 17th International Conference on Data Engineering. pp. 421–430. IEEE
 361 Computer Society (2001), <http://dl.acm.org/citation.cfm?id=645484.656550>

- 362 7. Braverman, M., Mao, J., Weinberg, S.M.: Parallel algorithms for select and partition
363 with noisy comparisons. In: Proceedings of the Forty-eighth Annual ACM Symposium
364 on Theory of Computing. pp. 851–862. STOC '16 (2016)
- 365 8. Chan, T.M.: Optimal output-sensitive convex hull algorithms in two and
366 three dimensions. *Discrete & Computational Geometry* **16**(4), 361–368 (1996).
367 <https://doi.org/10.1007/BF02712873>, <http://dx.doi.org/10.1007/BF02712873>
- 368 9. Chan, T.M., Lee, P.: On constant factors in comparison-based geometric
369 algorithms and data structures. *Discrete & Computational Geom-*
370 *etry* **53**(3), 489–513 (2015). <https://doi.org/10.1007/s00454-015-9677-y>,
371 <https://doi.org/10.1007/s00454-015-9677-y>
- 372 10. Chomicki, J., Ciaccia, P., Meneghetti, N.: Skyline queries, front and back. *SIG-*
373 *MOD Rec.* **42**(3), 6–18 (Oct 2013). <https://doi.org/10.1145/2536669.2536671>,
374 <http://doi.acm.org/10.1145/2536669.2536671>
- 375 11. Clarkson, K.L.: More output-sensitive geometric algorithms (extended ab-
- 376 stract). In: 35th Annual Symposium on Foundations of Computer Science,
377 Santa Fe, New Mexico, USA, 20–22 November 1994. pp. 695–702 (1994).
378 <https://doi.org/10.1109/SFCS.1994.365723>, [https://doi.org/10.1109/SFCS.1994.](https://doi.org/10.1109/SFCS.1994.365723)
379 [365723](https://doi.org/10.1109/SFCS.1994.365723)
- 380 12. Davidson, S.B., Khanna, S., Milo, T., Roy, S.: Top-k and clustering with
381 noisy comparisons. *ACM Trans. Database Syst.* **39**(4), 35:1–35:39 (2014).
382 <https://doi.org/10.1145/2684066>, <https://doi.org/10.1145/2684066>
- 383 13. Feige, U., Raghavan, P., Peleg, D., Upfal, E.: Computing with noisy
384 information. *SIAM Journal on Computing* **23**(5), 1001–1018 (1994).
385 <https://doi.org/10.1137/S0097539791195877>
- 386 14. Godfrey, P., Shipley, R., Gryz, J.: Algorithms and analyses for maximal vector
387 computation. *The VLDB Journal* **16**(1), 5–28 (2007). [https://doi.org/10.1007/s00778-](https://doi.org/10.1007/s00778-006-0029-7)
388 [006-0029-7](https://doi.org/10.1007/s00778-006-0029-7), [http://dx.doi.org/10.1007/s00778-](http://dx.doi.org/10.1007/s00778-006-0029-7)
389 [006-0029-7](http://dx.doi.org/10.1007/s00778-006-0029-7)
- 390 15. Goyal, N., Saks, M.: Rounds vs. queries tradeoff in noisy computation. *Theory of*
391 *Computing* **6**(1), 113–134 (2010)
- 392 16. Groz, B., Milo, T.: Skyline queries with noisy comparisons. In: Proceedings of the
393 34th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems.
394 pp. 185–198. PODS '15, ACM (2015). <https://doi.org/10.1145/2745754.2745775>,
395 <http://doi.acm.org/10.1145/2745754.2745775>
- 396 17. Kalyvas, C., Tzouramanis, T.: A survey of skyline query processing. *CoRR*
397 **abs/1704.01788** (2017), <http://arxiv.org/abs/1704.01788>
- 398 18. Khalefa, M.E., Mokbel, M.F., Levandoski, J.J.: Skyline query processing for
399 incomplete data. In: 2008 IEEE 24th International Conference on Data Engineering.
400 pp. 556–565. IEEE (2008)
- 401 19. Kirkpatrick, D.G., Seidel, R.: Output-size sensitive algorithms for finding maximal
402 vectors. In: Proceedings of the First Annual Symposium on Computational
403 Geometry. pp. 89–96. SCG '85, ACM (1985). <https://doi.org/10.1145/323233.323246>,
404 <http://doi.acm.org/10.1145/323233.323246>
- 405 20. Kossmann, D., Ramsak, F., Rost, S.: Shooting stars in the sky: An online algorithm
406 for skyline queries. In: Proceedings of the 28th International Conference on
407 Very Large Data Bases. pp. 275–286. VLDB '02, VLDB Endowment (2002),
408 <http://dl.acm.org/citation.cfm?id=1287369.1287394>
- 409 21. Kung, H.T., Luccio, F., Preparata, F.P.: On finding the maxima of a set of
410 vectors. *J. ACM* **22**(4), 469–476 (1975). <https://doi.org/10.1145/321906.321910>,
411 <http://doi.acm.org/10.1145/321906.321910>
- 412 22. Li, G., Wang, J., Zheng, Y., Franklin, M.: Crowdsourced data management: A survey
(2016)

- 413 23. Lofi, C., El Maarry, K., Balke, W.T.: Skyline queries in crowd-enabled databases. In:
414 Proceedings of the 16th International Conference on Extending Database Technology.
415 pp. 465–476. ACM (2013)
- 416 24. Marcus, A., Wu, E., Karger, D., Madden, S., Miller, R.: Human-powered sorts and joins.
417 Proc. VLDB Endow. **5**(1), 13–24 (2011). <https://doi.org/10.14778/2047485.2047487>,
418 <http://dx.doi.org/10.14778/2047485.2047487>
- 419 25. Newman, I.: Computing in fault tolerant broadcast networks and noisy decision
420 trees. Random Struct. Algorithms **34**(4), 478–501 (2009)
- 421 26. Pei, J., Jiang, B., Lin, X., Yuan, Y.: Probabilistic skylines on uncertain
422 data. In: Proceedings of the 33rd International Conference on Very
423 Large Data Bases. pp. 15–26. VLDB '07, VLDB Endowment (2007),
424 <http://dl.acm.org/citation.cfm?id=1325851.1325858>

425 **Appendix for the upper bounds**

426 **Proof of Theorem 2**

427 *Proof.* We denote by $\text{CheckVar}(x, \delta)$ the procedure that checks if $x = \mathbf{true}$ with
 428 er. pr. δ by majority vote, and returns the corresponding boolean. We finally denote
 429 by j' the true index of the first positive variable in x_1, \dots, x_m . We assume the input
 430 comparison oracles with error probability δ .

431 The probability that **NoisyOr** fails to identify j' for $i = 2^{\lceil \log k \rceil}$ (i.e., the first
 432 time it faces variable $x_{j'}$) is at most δ' . The probability that an incorrect index is
 433 returned (before $i \geq j'$) is at most $\sum_i \delta'/2^i$. The algorithm thus returns an incorrect
 434 index with probability at most $\delta' + \sum_i \delta'/2^i \leq \delta$.

435 **NoisyOr** requires $O(i)$ comparisons at line 4, whereas CheckVar requires $O(i)$
 436 comparisons at line 5. Replacing i with 2^h , the total cost on a successful execution
 437 is therefore $\sum_{h=1}^{\lceil \log j' \rceil} 2^h = O(j')$.

438 **Proof of Theorem 4**

439 *Proof.* For iteration j , the probability of error is $\delta/2^j$, and the cost is $O(nd\hat{k}\log(\hat{k}/\delta))$.
 440 Consequently, the probability that the algorithm fails to return the correct answer
 441 is at most: $\sum_j \delta/2^j \leq \delta$, and the running time is $O(\sum_{j=1}^{\lceil \log k \rceil + 1} nd2^j \log(2^j \times 2^j/\delta)) \in$
 442 $O(ndk\log(k/\delta))$. The complexity is $O((i' - i) * d\hat{k}\log(\hat{k}/\delta))$ to find a non-dominated
 443 point $p_{i'}$ at line 4, and $O(nd\log(\hat{k}/\delta))$ to compute the maximal point above $p_{i'}$
 444 at line 6. Summing over all iterations, the running time and number of queries is
 445 $O(nd\hat{k}\log(\hat{k}/\delta))$.

446 **Proof of Theorem 5**

447 The following Lemma lists properties that our bucketing assignment satisfies with
 448 high probability. We will show in Theorem 5 that our algorithm can compute the
 449 skyline efficiently for any assignment satisfying those properties.

450 **Lemma 2.** *Assume that the samples have been correctly ordered at line 5. With*
 451 *error probability δ/\hat{k} , the assignment performed at line 8 satisfies the following two*
 452 *properties:*

- 453 – *if I is a non-trivial interval (i.e., unless it matches the coordinate of a sample*
 454 *point), $|\{p: I = \mathbf{NoisySearch}(p_j, S_j, \delta'/(d\hat{k}))\}| \leq 4n/(d\hat{k}^2)$*
- 455 – *less than $2n/(d\hat{k}^2)$ points are (erroneously) assigned to buckets above the real*
 456 *skyline buckets.*
- 457 – *the skyline points are assigned to their correct bucket.*

458 *Proof.* Recall that $\delta' = \delta/(2d\hat{k})^5$, and that p_j denotes the j^{th} coordinate of point
 459 j . Assume the points of X are ordered w.r.t. to their j^{th} coordinate, breaking ties
 460 arbitrarily. Consider these ordered points to be divided into blocks, each one having
 461 $\ell = n/(d\hat{k}^2)$ consecutive points, except the last which may have less. In particular,
 462 the number of blocks is $d\hat{k}^2$.

463 Consider now the samples after line 5. Each block (but the last) contains at
 464 least one sample with probability at least $1 - (1 - \ell/n)^s \geq 1 - \delta'/d^2\hat{k}^2$. If one sample
 465 is indeed taken from every block (except maybe the last), the distance between any
 466 two samples is at most 2ℓ . As a consequence, the number of points p that should
 467 be assigned to any given bucket is bounded by 2ℓ , except for buckets with a trivial
 468 interval because several such buckets can be merged when removing duplicates
 469 at line 6. By Chernoff bounds, the number of points assigned to wrong buckets is
 470 at most $2n/(d\hat{k}^2)$ w.p. at least $1 - \delta'$. By union bound over all d dimensions and
 471 over all s intervals, we therefore have probability at least $1 - 3\delta'$ that one sample
 472 is taken from each block and that the total number w of points assigned to wrong
 473 buckets (over all dimensions and blocks) is less than $2n/(d\hat{k}^2)$. Consequently, with
 474 probability at least $1 - 3\delta'$ the assignment satisfies the first property. Indeed, for each
 475 dimension j and interval I , the number of points in I is bounded by 2ℓ (maximum
 476 distance between two samples) plus $2n/(d\hat{k}^2)$ (incorrect assignments into buckets):

$$|\{p: p \text{ was sorted into } I \text{ in line 8}\}| \leq 2\ell + \frac{2n}{d\hat{k}^2} = \frac{4n}{d\hat{k}^2}.$$

477 As for the number of buckets erroneously assumed to be non-empty, it is bounded by
 478 the number of points assigned to wrong buckets and is therefore at most $2n/(d\hat{k}^2)$.

479 This concludes the proof of the Lemma. We next turn to the proof of Theorem 5:

480 *Proof.* When $\hat{k}^5 \geq n$, $d^5 \geq n$ or $(\log(1/\delta))^5 \geq n$, the bounds can clearly be achieved
 481 by the other algorithms discussed previously, so we assume w.l.o.g. that $\hat{k}^5 < n$
 482 and $d^5 < n$ and $(\log(1/\delta))^5 < n$. We evaluate the cost of the algorithm assuming
 483 that (a) the samples are correctly sorted at Line 5, (b) the assignment satisfies the
 484 properties in Lemma 2, and (c) no mistakes are made at lines 13 and 18. In other
 485 words, we only accept a few mistakes at Line 8.

486 **Phase (i) Bucketing.** Line 5: by Theorem 1 (noisy sorting) the sample is
 487 sorted in $d \cdot O(\text{slog}(sd/\delta')) = O(nd \log(d\hat{k}/\delta))$. Line 8: by Theorem 1 (noisy search)
 488 the points are assigned to their bucket in $nd \cdot O(\log(sd\hat{k}/\delta')) = O(nd \log(d\hat{k}/\delta))$.

489 We will distinguish 4 kinds of (presumably) non-empty buckets (all other buck-
 490 ets are dropped at line 9): (i) those above the skyline that have been erroneously
 491 assigned some points, (ii) the buckets containing skyline points, (iii) the buckets
 492 that are dominated by buckets of type (ii), and (iv) the other (non-empty) buckets:
 493 they are not above the skyline but we do not have sufficient information to realize
 494 that they have no skyline points, because they are not dominated by any non-empty
 495 bucket. The algorithm is obviously not able to distinguish buckets of type (ii) and
 496 (iv), hence both are passed on to **SkylineHighDim** at line 18.

497 The number h of non-empty buckets is not necessarily much smaller than n as
 498 h may grow exponentially with d . Line 10 does not contribute to query complexity,
 499 but contributes $O(hd) \in O(nd)$ to the running time, using radix sort. Everything
 500 considered, the query complexity and running time of the bucketing phase are
 501 $O(nd \log(d\hat{k}/\delta))$.

502 **Phase (ii) Eliminating irrelevant buckets.** The buckets that are tested for
 503 emptiness are those of type (i), (ii) and (iv) because buckets of type (iii) are dropped

504 at line 17. The number of buckets of type (ii) is at most k . Furthermore, a bucket
 505 can be of type (iv) iff there is one dimension i such that they share the same
 506 coordinates I_i as a skyline bucket on dimension i , and the interval I_i is not trivial.
 507 Consequently, by Lemma 2, there are at most $(dk) \times 4n/(d\hat{k}^2)$ points that belong
 508 to buckets of type (iv). The number of points in buckets of type (ii) is even smaller:
 509 when such a bucket is trivial it contains only skyline points, and when it is not
 510 trivial, there is a dimension on which it is a non-trivial interval and therefore by
 511 Lemma 2 it has at most $4n/(d\hat{k}^2)$ points, hence a total of at most $k \times 4n/(d\hat{k}^2)$
 512 points in buckets of type (ii). When the estimate \hat{k} is large enough ($k \in O(\hat{k})$),
 513 the number of points in buckets of type (ii) or (iv) is therefore $O(n/\hat{k})$. The case
 514 when this is not $O(n/\hat{k})$ because the estimate is not large enough is handled on
 515 line 15. Similarly, Lemma 2 guarantees that $O(n/\hat{k})$ points have been assigned to
 516 buckets of the first kind. Therefore, the total number of points ever considered
 517 on line 13 is $O(n/\hat{k})$. The contribution of line 13 to the complexity is therefore
 518 $O(d(n/\hat{k})\log(\hat{k}/\delta'))$ by Lemma 1. Line 17 does not contribute to query complexity,
 519 but contributes $hd\hat{k} \in O(nd \cdot \min(\hat{k}, |\text{skyline}(X)|))$ to the running time.

520 Actually, we need to optimize a bit the algorithm to achieve that running time.
 521 There can be much more than \hat{k} iterations, but there are only $\min(\hat{k}, |\text{skyline}(X)|)$
 522 "relevant" iterations in which we need to drop buckets. So we first strengthen
 523 the requirement on the order at line 10, so that a bucket comes before buckets
 524 it *weakly* dominates, where B' weakly dominates B (using the notation above) if
 525 in every dimension $\max_i \leq \max I'_i$. At line 17, if B_i has already been marked as
 526 weakly dominated, we move on to the next iteration (any bucket that B_i would
 527 dominate has already been dropped). Otherwise, we iterate through the list of
 528 remaining buckets, and we perform the following operations at a cost of $O(d)$ per
 529 bucket: we drop the buckets that B_i dominates, and mark the other buckets that
 530 B_i weakly dominates. There are only $\min(\hat{k}, |\text{skyline}(X)|)$ buckets that are not
 531 weakly dominated, hence the running time.

532 **Phase (iii) Solving the reduced problem.** Finally, at line 18 the size of
 533 X' is $O(n/\hat{k})$, so its skyline can be computed in $O(nd\log(\hat{k}/\delta'))$ by **SkyHighDim-Search**.

534 We next show that the correct answer is returned with high probability. First,
 535 the probability that the algorithm fails to satisfy our requirements (a) to (c) above
 536 are respectively $d \cdot \delta'/d$, δ/\hat{k} and $\hat{k} \cdot \delta'/\hat{k} + \delta'$. So the conditions are met — hence
 537 the algorithm returns the correct output — with probability at least $1 - 4\delta$.

538 Appendix for the lower bounds

539 In this section, we exhibit an $\Omega(dn\log k)$ lower bound on the query complexity in
 540 the noisy skyline problem, denoted **Skyline**. To that end, we define a noisy vector
 541 problem, in which one is given k vectors each of length ℓ and needs to decide for each
 542 vector whether it is the all-zero vector. We prove a lower bound for this problem
 543 and reduce it to **Skyline** yielding the desired result.

544 **5.1 (k, ℓ) -Null-Vectors: Definition and Lower Bound**

545 In the (k, ℓ) -Null-Vectors the input S is a collection $\{\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^k\} \subseteq \{0, 2\}^\ell$ of vectors
 546 such that for each $i \in [k]$, $\sum_{j=1}^\ell \mathbf{v}_j^i \leq 2$, and the output is a vector $(w_1, w_2, \dots, w_k) \in$
 547 $\{0, 2\}^k$ such that for each $i \in [k]$, $w_i = \sum_{j=1}^\ell \mathbf{v}_j^i$. We define the distribution μ
 548 over vectors of $\{0, 2\}^\ell$ as follows. For each $j \in [\ell]$, $\mu(2e_j) = 1/(2\ell)$, where e_j is the
 549 canonical vector with a 1 in the j -th entry and zero elsewhere; $\mu(0, \dots, 0) = 1/2$. For
 550 inputs to (k, ℓ) -Null-Vectors, we will consider the product distribution μ^k .

551 **Lemma 3.** *For (k, ℓ) -Null-Vectors under the product distribution μ^k , if A is a*
 552 *deterministic algorithm with success probability at least $3/4$, then the worst case*
 553 *number of queries of A is $\Omega(\ell k \log k)$.*

554 *Proof.* The proof is by contradiction. Assume that A is an algorithm with success
 555 probability at least $3/4$ and worst case number of queries $T \leq (\ell k \log_3 k)/1000$. We
 556 assume that the adversary is *generous*, i.e. the adversary tells the truth for every
 557 entry (i, j) such that $v_j^i = 0$, and that lies with probability $1/3$ otherwise.

558 Generalizing the 2-phase computational model by Feige, Peleg, Raghavan and
 559 Upfal [13], we will give the algorithm more leeway and study a 4-phase computation
 560 model, defined as follows. In the first phase, the algorithm queries every entry v_j^i
 561 $(\log_3 k)/100$ times. In the second phase, the adversary reveals to the algorithm all
 562 remaining hidden entries (i, j) such that $v_j^i = 2$, except for a single random one.
 563 In the third phase, the algorithm can strategically and adaptively choose $\ell k/10$
 564 entries, and the adversary reveals their true value at no additional cost. Finally,
 565 in phase 4, the algorithm outputs $w_i = 2$ for every vector where it found an entry
 566 equal to 2, and $w_i = 0$ for the rest of the vectors.

567 To see how the two models are related, observe that since $T \leq (\ell k \log_3 k)/20$,
 568 by Markov's inequality at most a set S of $\ell k/10$ entries are queried by algorithm
 569 A more than $(\log_3 k)/2$ times, so at the end of the first phase we have queried
 570 every entry at least as many times as A , except for those $\ell k/10$ entries, and in the
 571 beginning of the third phase there is all the necessary information to simulate the
 572 execution of A , adaptively finding S (and getting those values correctly), hence
 573 the success probability of the three-phase algorithm is greater than or equal to the
 574 success probability of A . Also observe that, thanks to the definition of μ and to
 575 the generosity of the adversary, any execution where all queries to a vector lead
 576 to 0 answers must lead to an output where $w_i = 0$ —else the algorithm would be
 577 incorrect when μ selects the null vector.

578 We now sketch the analysis of the success probability of the three-phase algo-
 579 rithm. Due to the definition of μ , with probability at least $9/10$ the ground-truth
 580 input drawn from μ^k has $k/2 \pm O(\sqrt{k})$ vectors that contain an entry equal to 2.
 581 At the end of the first phase, and due the fact that the adversary is generous, we
 582 have that at most of them have been identified. There remain $k/2 \pm O(\sqrt{k})$ vectors
 583 that appear to be all zeroes, and about $(k/2)(1/3)^{(\log_3 k)/2} = (1/2)\sqrt{k}$ of those
 584 vectors contain a still-hidden entry whose true value is 2. During the third phase,
 585 all of those hidden 2's are revealed except for one. At that point, there still remain
 586 $k/2 \pm O(\sqrt{k})$ vectors whose entries appear to be all zeroes, there is a 2 hidden

587 somewhere uniformly at random, but all entries have been queried an equal number
 588 of times, all in vain. To find that remaining hidden entry (and therefore decide
 589 which w_i is equal to 2), the algorithm has no information to distinguish between
 590 the $\ell(k/2 \pm O(\sqrt{k}))$ remaining entries. Since, the algorithm may only select $\ell k/10$
 591 elements to query further, the algorithm's success probability after the fourth
 592 phase cannot be better than $(k\ell/10)/(\ell(k/2 \pm O(\sqrt{k}))) < 1/4$, a contradiction.

593 5.2 Reduction: Proof of Theorem 7

594 *Step 1.* Assume, for simplicity, that $d-2$ divides k . From an input $S = \{\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^k\}$
 595 to the (k, ℓ) -Null-Vectors, we first show how to construct an input \mathcal{I}_S for Skyline
 596 with n points in d dimensions and a skyline that is likely to be of size k , where
 597 $n = (\ell + d - 2)k / (d - 2)$. We first randomly permute the entries of each \mathbf{u}^i , by
 598 using k independent permutations, resulting in $S_\pi = \{\mathbf{v}^1, \mathbf{v}^2, \dots, \mathbf{v}^k\}$. Partition S_π
 599 into $k/(d-2)$ blocks of $d-2$ vectors, where for $j \in \{0, 1, \dots, k/(d-2) - 1\}$, block
 600 $S_\pi^j = \{\mathbf{v}^{j(d-2)+i} : i \in [d-2]\}$. For each block, define $\ell + d - 2$ points, as displayed
 601 (one point per row) on Figure 4, and the union over all blocks is the input \mathcal{I}_S to
 602 the Skyline. Formally, we define point $\mathbf{p}^{(t)}$ with $t = j\ell + i$ as follows.

$$\mathbf{p}_r^{(t)} := \begin{cases} j & \text{if } r = d-1, \\ n-j & \text{if } r = d, \\ 1 & \text{if } r = i-\ell \text{ and } \ell \leq i, \\ v_i^{j(d-2)+i} & \text{if } r = i \text{ and } i \in [d-2], \\ 0 & \text{otherwise.} \end{cases}$$

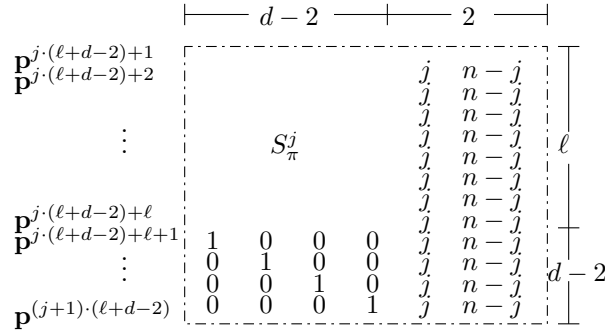


Fig. 4: Block $(j, n-j)$ of the reduction. The vectors of S_π^j placed in this block are $\mathbf{v}^{j(d-2)+1}, \mathbf{v}^{j(d-2)+2}, \dots, \mathbf{v}^{j(d-2)+(d-2)}$.

603 *Step 2.* Because of the non-domination implied by the last two coordinates of any
 604 point, the skyline of the set of points is the sum over all blocks of the skyline of

605 each block. Fix an arbitrary block and focus on the first $d-2$ dimensions. For each
 606 dimension, the corresponding column (whose first ℓ coordinates are those of some
 607 vector \mathbf{v}^i) contains exactly one 1 (on the row of some point \mathbf{p}) and possibly one
 608 2, the remaining entries being all 0. Thus it is easy to verify that \mathbf{p} is part of the
 609 skyline if and only if $\mathbf{v}^i = \mathbf{0}$.

610 From the output skyline(I) it is now easy to construct the output of the (k, ℓ) -
 611 Null-Vectors: For all blocks, for all dimensions $\leq d-2$, if $\mathbf{p} \in \text{skyline}(I)$ then $w_i \leftarrow 0$
 612 else $w_i \leftarrow 2$. This yields the correct output $\mathbf{w} = (w_1, w_2, \dots, w_k)$. Thus we derive the
 613 following observation.

614 **Observation 1** *Given the set of points skyline(\mathcal{I}_S), one can recover the solution*
 615 *to the (k, ℓ) -Null-Vectors without further queries.*

616 Furthermore, in the following we prove that the construction is likely to have
 617 k skyline points.

618 **Lemma 4.** *Let \mathcal{E} be the event that the input \mathcal{I}_S has exactly k skyline points. Then,*
 619 $\mathbb{P}(\mathcal{E}) \geq 1 - 1/k$ *as long as $k^5 \leq n$.*

620 *Proof.* First observe that, by construction, regardless of whether \mathcal{E} holds, every
 621 block contains at most $d-2$ skyline points: Consider an arbitrary block. The last
 622 two dimensions are identical for each point belonging to that block and we focus
 623 thus on the first $d-2$ dimensions. There are exactly $d-2$ points with one coordinate
 624 being 1 and all of these points are potential skyline points. In particular, take any
 625 such point \mathbf{p} and assume that the i 'th coordinate of \mathbf{p} is 1. Then \mathbf{p} is part of the
 626 skyline if and only if the vector \mathbf{v}^i is the null vector. Moreover, every block can
 627 have at most $d-2$ entries with value 2 and each such 2 eliminating one potential
 628 skyline point. Thus, there are at most $d-2$ skyline points per block.

629 Consider the vertices $\mathbf{v}^{i_1}, \mathbf{v}^{i_1+1}, \dots, \mathbf{v}^{i_2}$ of any block. We say they are *collision*
 630 *free* if the following holds: if $\mathbf{v}_{j^*}^j = 2$ for $j \in [i_1, i_2]$, then $\mathbf{v}_{j^*}^{j'} = 0$ for all $j' \in [i_1, i_2] \setminus \{j\}$.
 631 Observe that if the vertices of any block are collision free, then each of the first $d-2$
 632 dimensions is dominated by a distinct skyline point and thus there $d-2$ skyline
 633 points in that block. Thus, if the vectors of every block are collision free, then there
 634 $d-2$ skyline points per block and summing up over all $k/(d-2)$ blocks, we get
 635 that there are thus k skyline points in total.

636 Thus, in order to bound $\mathbb{P}(\mathcal{E})$ it suffices to bound the probability that all blocks
 637 are collision free.. Recall that the random permutations $\pi_1, \pi_2, \dots, \pi_k$ permute each
 638 vector \mathbf{v}^i independently. Since in a block at most k^2 pairs may collide, and each
 639 collision happens with probability $1/\ell$, the expected number of *collisions* per block
 640 is at most k^2/ℓ . The expected number of collisions over all blocks is thus, by the
 641 union bound, at most $(k/(d-2)) \cdot k^2/\ell \leq 1/k$, by assumption on k . Thus, the claim
 642 follows by applying Markov inequality.

643 *Proof (Proof of Theorem 7).* Suppose for the sake of contradiction that there exists
 644 an algorithm \mathcal{A} recovering the skyline for any input with exactly k skyline points,
 645 with error probability at most $1/10$, and using $o(nd \log k)$ queries in expectation.
 646 By Markov inequality, the probability that the number of queries exceeds 5 times

647 the expectation is at most $1/5$, so truncating the execution at that point adds $1/5$
 648 to the error probability, transforming \mathcal{A} into an algorithm B that recovers the
 649 skyline for any input with exactly k skyline points, with error probability at most
 650 $1/5 + 1/10 < 1/3$, and using $o(5nd \log k)$ queries in the worst case. We claim that
 651 this implies that one can solve the (k, ℓ) -Null-Vectors with $o(nd \log k)$ w.p. at least
 652 $1/3$ contradicting Lemma 3.

653 Let S be the input of the (k, ℓ) -Null-Vectors. We cast S as an input \mathcal{I}_S of B as
 654 described in Section 5.2. By Lemma 4, the event \mathcal{E} holds w.p. at least $1 - 1/k$ and
 655 thus there are k skyline points.

656 By assumption, B can thus compute the skyline w.p. at least $1/2 - 1/k \geq 1/3$,
 657 where we used the Union bound. Thus, by Observation 1, one can obtain w.p. at
 658 least $1/3$ the solution to (k, ℓ) -Null-Vectors using $o(nd \log k)$ queries, a contradiction.