

## COVER NOTE OF THE PVLDB SUBMISSION WITH PAPER ID = 1501 TITLED WITH “ROBUST BEST POINT SELECTION UNDER UNRELIABLE USER FEEDBACK”

Dear PC Chair and the area coordinator,

First of all, we sincerely thank the reviewers for their insightful and constructive comments. We revised the paper according to their suggestions. We are including our response and our updated version. Thank you for helping to review the paper.

The following shows our paper information.

- Paper Title: Robust Best Point Selection under Unreliable User Feedback
- Paper ID: 1501

Thank you for your attention.

Regards,

Qixu Chen and Raymond Chi-Wing Wong

## I SUMMARIZATION OF REVISION

We revised the user study section thoroughly in response to the reviewers’ feedback. Specifically, in the current submission, we did the following to the user study.

- (1) We demonstrated that a linear utility function is a suitable approximation for the user’s preference in our scenario to address Reviewer 4’s concern regarding the linear utility assumption.
- (2) We added a new part (i.e., Part 3 for the current submission) into the existing user study to address Reviewer 4’s concern regarding our algorithm not performing well when only a limited rounds of interactions are allowed.
- (3) We conducted a second user study to provide evidence of persistent errors in real-world scenarios to address Reviewer 4’s concern about the presence of errors in real user input.

## II RESPONSE TO REVIEWERS’ COMMENTS

### II.1 Revisions Corresponding to Reviewer 4’s Comments

REVIEWER 4’S COMMENT 1. *“Thanks a lot for the revised paper. The quality is much improved. The authors have incorporated most of my suggestions. While I appreciate the authors’ effort a lot, I was not able to see the experimental improvement on the targeted problem as depicted in Figure 21.*

*Part 3 in Section 6.5 (user study) corresponds to the scenario with real user input and this is where I look at for performance comparison. It is very likely that real user input contains little persistent/random error or that the real scenario does not have linear utility function (see O1, D1.1/1.2), and as a result the baselines like HD-PI performs better than the proposed approach (higher hit rate when allowing the same number of rounds of user input).*

*Part 4 in the User study “artificially” ingested persistent error and it indeed demonstrated the usefulness of the proposed method. However, this synthetized persistent error does not seem to appear in the real user study.”*

**Response.** Thank the reviewer for this comment. The comment shows two major concerns. The first concern is whether linear utility function works well in the real user study. The second concern is whether the persistent/random errors in the real user input has an impact on the performance of the algorithms. The reviewer wonders if these two concerns are the two possible reasons why our algorithm’s hit rate is worse than other methods in Part 3 of our user study.

For the first concern, note that it is impossible to obtain the user’s “real” underlying preference function. Therefore, following the existing works [8, 34, 41, 42, 44] in this line of research, we model the user’s utility function to be linear. We want to argue that in the scenario of our user study, the linear utility function models the user’s preference reasonably well. Observe that in Part 1 of the user study, we did not ingest any artificial error. Thus, in this part, the utility vector is also learned based on the *real user input*, and the linear vector learnt by *SS-score* obtains over 90% hit rate. Given this high hit rate, it is sufficient to say that linear function can approximate the user’s preference reasonably well. In the current submission, we added the following sentence in Section 6.5 to clearly state this point.

*“Given that SS-score obtains a hit rate exceeding 90%, we conclude that the linear utility function approximates the user’s preference reasonably well.”*

For the second concern, as stated in the user study (Section 6.5), we computed an estimation of empirical user error rate of 5.8% using *SS-score*, which is not a negligible percentage. The effect of user errors can also be validated by the performance of *HD-PI* in Part 1. Note that *HD-PI* also assumes the linear utility function and returns the best point using the learned utility vector, but, it does not address user errors. If there is little user error involved, or the user errors have a small impact on the algorithms’ performance, even if the linear function cannot perfectly model user’s function, the performance of *HD-PI* should at least align with the performance of *SS-score* in Part 1 (i.e., a 90% hit rate). However, the hit rate of *HD-PI* is only 60%. The 30% gap of hit rate reflects the large impact of user errors.

Furthermore, we conducted the second user study to show the presence of persistent error in real user input. The details for this user study can be found in Section 6.5. The results of this user study demonstrate the existence of persistent user errors. Under its setting, at least 33% of the users are prone to persistent errors during interaction.

Next, we explain why *SS-score* does not perform well when only a limited number of rounds are allowed. As mentioned in the response in our previous revision, there is a trade-off between round efficiency and error handling capacity. Since our algorithms prioritize error handling, it is expected that the round efficiency of our algorithms are lower than algorithms that do not consider error handling (e.g., *HD-PI*) or have a lower error handling capacity (e.g., *Verify-Point*). This can also be observed from our other experiments, such as in Figure 14 (b), where our algorithms need more rounds

to finish compared with *HD-PI* and *Verify-Point*. Since in Part 3 (of the previous submission), the algorithms were forced to stop after round 10, *SS-score* was not able to reduce the candidate size to  $l$  and had to randomly select  $l$  points from a relatively large number of candidates. Consequently, its hit rate greatly decreased as expected.

To show that the hit rate drop of *SS-score* in Part 3 (of the previous submission) is attribute to the small round limit set in this part, we made the following changes to the user study. We removed the old Part 2 and Part 4 in our previous submission from the main paper and put them into technical report since the reviewer suggested that they do not correspond to a real user case. In the current submission, this user study consists of 3 parts, and the setting of Part 2 and 3 are as follows.

*“Part 2 and Part 3 follow the design of Part 1 with the only difference that in these two parts, each algorithm was required to stop and return a list of at most  $l$  options after  $t$  rounds. In Part 2,  $t$  is set to 10. In Part 3,  $t$  is set to the number of rounds required by *SS-score* to terminate (which is typically 12 to 13).”*

After these changes, Part 3 of the previous submission became Part 2 of the current submission, and we introduced a new setting as Part 3 of the current submission. It is worth mentioning that in the current submission, the setting of Part 3 is essentially the same as that of Part 2, except with a larger round limit (i.e., 12 to 13). We observe that in the current Part 3, *SS-score* outperforms all other baselines, and its hit rate is comparable to that in Part 1. In the current submission, we updated Figure 20 to show the results of these 3 parts and stated the following.

*“In Part 2, the hit rate of *SS-score* is lower than *HD-PI* and *Verify-Point*. This is explainable because there is a trade-off between round efficiency and error handling capacity. Since *SS-score* prioritizes error handling, it is expected that its round efficiency is lower than algorithms that either do not consider error handling or have a lower error handling capacity. Consequently, when forced to stop after round 10, *SS-score* has to randomly select  $l$  points from a relatively large number of candidates, resulting in a lower hit rate. However, as can be observed in Part 3, when allowed to use slightly more rounds (i.e., 2 to 3 more rounds), the hit rate of *SS-score* exceeds other baselines by a large margin, aligning with its performance in Part 1.”*

We also wish to point out that the user study setting suggested by the reviewer (i.e., interacting for a limited number of round and then terminating) is slightly different from the scenarios considered in this paper. We acknowledge that there are many situations where the user expects to see the results within only a limited number of rounds. However, as stated in Section 1, in this paper we focus on situations where finding the best point is the most important to the user, while missing the best point can lead to unforgettable and unchangeable consequences, such as crucial lifelong decisions like choosing the university and purchasing a house. In these applications, the accuracy of finding the best option is the most important aspect and any potential user error during interaction must be properly handled. Moreover, the user is typically willing to answer a few more questions to reduce the chance of missing the best option. Therefore, although our algorithm’s performance in Part 3 is not better than the baseline, it does not undermine the effectiveness of our methods in the scenarios mainly considered in this paper. To emphasize this, we added the following sentence in Section 6.5.

*“As stated in Section 1, there are situations where missing the best point can cause unforgettable and unchangeable consequences, such as choosing the university and purchasing a house. In these scenarios, the user’s primary concern is the accuracy of finding the best point, even if slightly more questions need to be answered.”*

### III RESPONSE TO META-REVIEWERS’ COMMENTS

META REVIEWER’S COMMENT 1. “Reviewer 4 (Question 3) has noted key issues that require the need for the claims made in your paper to be qualified in view of the lack of complete validation in user study. Specifically, limitations of your technique and their implication for real use cases must be noted and put in perspective.”

**Response.** Please read our response in Section II. We addressed all concerns raised by Reviewer 4 on the user study.

# Robust Best Point Selection under Unreliable User Feedback

Qixu Chen

The Hong Kong University of Science and Technology  
Kowloon, Hong Kong  
qchenax@connect.ust.hk

Raymond Chi-Wing Wong

The Hong Kong University of Science and Technology  
Kowloon, Hong Kong  
raywong@cse.ust.hk

## ABSTRACT

The task of finding a user’s utility function (representing the user’s preference) by asking them to compare pairs of points through a series of questions, each requiring him/her to compare 2 points for choosing a more preferred one, to find the best point in the database is a common problem in the database community. However, in real-world scenarios, users may provide unreliable answers due to two major types of errors, namely *persistent* errors and *random* errors. Existing interaction algorithms either simply assume that all answers provided by the user are reliable, or are capable of handling *random* errors only, which can lead to finding *undesirable* points, ignoring persistent errors. To address this challenge, we propose more generalized algorithms that are robust to both persistent and random errors made by the user. Specifically, we propose (1) an algorithm that asks an asymptotically optimal number of questions, and (2) an algorithm that asks an even smaller number of questions empirically, with provable performance guarantee. Our experiments on both real and synthetic datasets demonstrate that our algorithms outperform existing methods in terms of accuracy, even with a small number of questions asked.

## 1 INTRODUCTION

When faced with a database containing millions of points, an end user may be only interested in finding his/her favorite point in the database. For example, a user may want to buy a car with a cheap price and a high horsepower from a car database, and s/he might only investigate the cars that excel in these aspects. Here, price and horsepower are some *attributes* that a user will consider when buying a car. To help the user efficiently find the most interesting point, *multi-criteria decision-making queries* are proposed to return a representative set from the database which consists of potentially interesting points. Two popular queries are the top- $k$  query [25, 33] and the skyline query [6]. The top- $k$  query returns  $k$  points from the dataset with the highest *utility* w.r.t a utility function. However, in practice, this utility function may not be known in advance. The second query, namely the skyline query, is based on a concept called *dominance*. Specifically, a point  $p$  *dominates* another point  $q$  if  $p$  is not worst than  $q$  in any attribute and is better than  $q$  in at least one attribute. The skyline query returns all points that are not dominated by any other points in the database. Although the knowledge of the user’s utility function is not required, the skyline query does not guarantee a controllable return size, which might be as large as the entire database in the worst case [29, 32].

Recently, a novel interactive framework [31, 42, 44] was proposed to combine the advantage of both the top- $k$  query (which has a fixed return size) and the skyline query (which does not need an exact utility function). Without any required knowledge of the user in advance, it asks the user a number of *rounds* of simple questions and learns the user’s preference progressively, and recommends

points based on the learned preference. A widely applied form of questions [31, 34, 37, 41, 42, 44] is to present 2 points in each round, and asks the user to select the preferred one. Consider the car purchasing scenario. The interactive framework simulates a sales assistant that asks Alice to indicate her preference among several pairs of cars and make recommendations based on her answers.

Although with good experimental performance, one key factor that prevents the existing interactive algorithms from being more practical is that they take no consideration of the *reliability* of the user feedback and implicitly assume that the user is *always* correct. However, in practice, the user’s feedback can be *unreliable* due to various reasons. For example, during the interaction, even though Alice wants a car with good horsepower and costs as low as possible, she may indicate that she prefers a more expensive car to a cheap car due to a *careless mistake* (e.g., mis-clicking) or some *cognitive bias* (e.g., she mistakenly thinks that a car with a higher price must have better horsepower). Although a sales assistant can observe Alice’s real intent by considering her overall choices, making this error when interacting with the existing algorithms will make them believe that Alice is willing to spend more. Consequently, these algorithms prune a set of cars from further consideration, possibly including the real best car for Alice.

It is worth mentioning that making small errors in the decision-making process can cause unforgettable and unchangeable consequences. For example, a common type of error that occurs in the trading market is known as a “fat finger error”, which refers to a mistake made by a trader in the trading system by clicking or pressing the wrong key. In 2018, Samsung Security made a wrong transaction worth 100 billion dollars due to a fat finger error, which could have resulted in a loss of 428 million dollars, equivalent to 12.17% of the company’s market capitalization [2]. Another example is about one of the critical milestones of one’s life: the selection of the tertiary school. In 2020, a student in Mainland China achieved a top-tier score in the National College Entrance Examination in China (also called gaokao). However, he was admitted to a low-ranked college with a similar name to his target university since he confused the name of the two schools in the tertiary school selection system, which may cause a huge impact on his future life [43].

Motivated by the deficiency of existing algorithms, in this paper, we study the problem called the *interactive best point retrieval* problem under unreliable user input, which is more realistic. Roughly speaking, our problem is to find in a dataset  $D$  the best point (i.e., the point with the highest *utility*) for a user w.r.t. his/her personalized utility function  $f$ , which we do not know in advance and needs to be learned by asking the user to answer several rounds of questions. We focus on a type of question that is widely adopted [31, 34, 37, 41, 42, 44], which is to display two points from  $D$  and asks the user to select the preferred point. Due to the simplicity of this question type, the techniques developed for this question type

can be easily extended to other types of questions (e.g., displaying more than two points in each question and asking for the favorite one). In Section A.4 of the Appendix, we show how the proposed algorithms can be applied to other question types as well.

There are two major types of user errors that cause the unreliability in the user’s answers, namely *random errors* and *persistent errors*. Random errors [8, 20] occur due to unintentional reasons such as mis-clicking and pressing the wrong key, which means that the answer to the same question may be different when asked again due to this careless mistake. On the other hand, persistent errors [17, 20, 24] are caused by cognitive biases or other sources of interference, and the answer to the same question may be *consistently* wrong. For example, when comparing a car priced at \$5000 and another car originally priced at \$10000 but currently offered at a 50% discount, despite the former having slightly better specifications and offering higher utility, Alice may consistently choose the latter. This decision is influenced by the psychological phenomenon known as the “anchoring bias”, wherein individuals are swayed by the allure of discounts [46]. Compared to random errors, persistent errors are harder to be handled, since in the case of a random error, the user’s real preference can be revealed by repeating the same question several times. The techniques developed in this paper can handle both types of errors, or even a combination of them.

Unfortunately, most existing interactive algorithms which do not consider user errors will return *undesirable* points based on *wrongly learned* utility functions. Their accuracy in returning the best point under the setting of unreliable user feedback is unsatisfactory (e.g., smaller than 70% on a 4-d dataset with 1 million points). Moreover, direct adaptations of existing algorithms considering user errors from the field of “learning to rank” and machine learning turn out to be inefficient since they ask too many questions. In our experiment, algorithms in this line of research, such as *Active-Ranking* [20] and *Pref-Learn* [34], ask more than 60 questions when the input dataset is large (i.e., 1 million points), which is undesirable. Users will lose the plot and get frustrated if they need to answer excessive questions [26, 35].

The most closely related work is [8] which aims at finding the best point considering random user errors. However, the techniques developed in [8] can only handle random errors by asking the *same* pair of points multiple times and taking the *majority vote*, which cannot be adapted to address persistent errors because the majority vote also results in incorrect preferences in this scenario. Compared to [8], our algorithms are more generalized and practical since they effectively handle both types of errors in a unified way. This advancement is attributed to a novel and previously unexplored geometrical concept proposed in this paper called the *confidence region*. The confidence region exhibits several interesting properties when dealing with user errors, and we leverage these properties in the design of our algorithms.

**Contributions.** We summarize our contributions as follows. Firstly, we study the *interactive best point retrieval* problem considering both persistent errors and random errors. Secondly, we show a lower bound on the number of questions needed (also called *round complexity*) for our problem. Thirdly, we study a novel geometrical concept called the *confidence region*, which is instrumental in effectively handling both types of errors. Fourthly, we propose (1) an algorithm with asymptotically optimal round complexity;

and (2) an algorithm with even better empirical performance. Both algorithms return the best point with provable guarantee. Lastly, we conducted comprehensive experiments to demonstrate the superiority of our algorithms. They maintain high accuracy (e.g., nearly to 100% in most experiments) with only a small number of questions, but existing approaches either ask too many questions (e.g., twice as many as ours) or are much more inaccurate (e.g., more than 10% less accuracy than ours).

**Organizations.** The rest of the paper is organized as follows. Section 2 shows the related work. The formal definition of the studied problem is given in Section 3. In Section 4, we show the lower bound on the number of rounds required by this problem and describe the general framework of our algorithms. We present the details of proposed algorithms in Section 5 and show experimental results in Section 6. Section 7 concludes the paper.

## 2 RELATED WORK

Various queries are proposed to assist the multi-criteria decision-making. The *preference-based queries* return points based on the preference or the expected point of the user, and *interactive queries* involve user interaction to find points that may interest the user.

Besides the top- $k$  and skyline queries mentioned in Section 1, many other types of preference-based queries are proposed. The  $k$ -nearest neighbors query (kNN) [39] returns  $k$  points that are closest to an example point given by the user based on their Euclidean distances. The similarity query [38] defines a more complicated distance function and finds  $k$  closest points to a given point using this function. The problem with these two queries is that an example point must be provided by the user in advance, which may be too demanding in many situations. Some recent studies aim to combine the ideas of top- $k$  and skyline. [10, 29, 30] consider returning a fixed number of points to the user whose preference is estimated to lie in a region of the function space. However, if this region is large, the output size must also increase to guarantee the retrieval of high-quality points. In the worst case, it degenerates to the original skyline query. [1, 9, 32] propose  $k$ -regret minimizing query, which computes a set such that for any utility function, there exists a point in the set whose *regret ratio*  $\leq \epsilon$ , a user parameter. However, to guarantee a small  $\epsilon$ , the output size could be large for some datasets (e.g.,  $> 1000$  points when  $\epsilon = 0.5\%$  [1]).

The interactive queries learn the user’s preference by interacting with the user and return points based on the learned preference. [28] proposes the interactive skyline query, which reduces the skyline size by learning the relative skyline importance of the user. During the interaction, the user is asked to partition points into *superior* and *inferior* groups. However, the output size is still uncontrolled even if the preference for all attributes is obtained. The interactive similarity query [4, 5, 37] interacts with the user to learn a distance function, which is then used to find  $k$  points that are closest to a query point. However, during the interaction, the user is asked to assign *relevant scores* to hundreds of points, this quantitative question is too demanding for most users.

Although the type of interaction varies, one widely adopted interaction is to display a pair (or set) of points in each round and ask the user to select the preferred one [31, 34, 37, 41, 42, 44]. [31] proposes an interactive regret minimization query, which aims

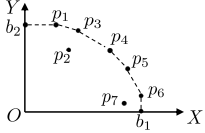


Figure 1: The upper hull

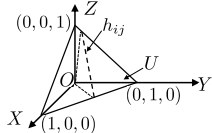


Figure 2: Utility space in 3d

at minimizing the regret ratio of the return set with a fixed size by learning the user’s utility function using artificial points. To overcome the deficiency of using unreal points, [44] introduces *UH-Simplex* and *UH-Random*, which only display points inside the database. [42] proposes *HD-PI* and *RH* that return one of the top- $k$  points of the user. However, these algorithms all assume that the user makes no error, and make decisions based on user’s answers without questioning their reliability. It is hard to directly adapt these algorithms to handle user errors, and their performance degenerates when the user makes mistakes.

The most closely related work is [8] in which they propose algorithms *Verify-Point* and *Verify-Space* that return the best point considering *random user errors*. However, their methods have several drawbacks. Firstly, they address random user errors by taking majority votes on repetitive questions. This technique, however, cannot handle *persistent errors* since the feedback to the same question will be the same. In our experiments, their accuracy decreases by more than 10% when addressing persistent errors compared with addressing random errors. Also, repeating the same question is not suitable for many applications. For example, Alice will be confused if the seller asks her to compare the same pair of cars three times. Secondly, they still display artificial points during the interaction. The algorithms proposed by us overcome all these limitations.

In the field of machine learning (ML) and information retrieval (IR), the robustness to erroneous user input is considered in some algorithms [12, 13, 20, 21, 34, 36]. However, since their focuses (e.g., finding the exact ranking) are typically different from ours, they tend to be extravagant in the number of pairwise comparisons used, and thus, are inefficient for the user to interact with. [34] proposes *Preference-Learning* to learn user’s preference and addresses user errors by introducing a slack variant in a linear SVM. It tends to ask extra questions since the major focus of this work is to approximate a preference vector. [20] aims at computing the ranking of all points and uses the majority vote to resolve conflicts in comparison results. This line of work needs more questions to find the ordering of non-best points, which is not a concern in our problem. [12] finds top- $k$  points with initial partial ordering information and handle errors using majority votes. They require more questions than ours since they do not consider the geometric relation between points.

Compared with the existing work, our work has the following advantages. Firstly, we do not require any knowledge of the user’s utility function in advance and we guarantee a small return size, which overcomes the limitation of traditional top- $k$  and skyline queries. Secondly, we reduce the user’s effort by only requiring the user to answer a small number of simple questions, while some existing algorithms ask a large number of questions (e.g., [13, 20, 34]) and others ask difficult questions (e.g., [5, 28]). Finally, our algorithms allow the user to be possibly unreliable without any significant impact on the quality of the points returned, which entails more practical value compared to some existing work that assumes perfect user feedback (e.g., [31, 42, 44]).

### 3 PROBLEM DEFINITION

In this section, we provide the formal definitions for the *interactive best point retrieval problem*, the *random errors*, and the *persistent errors*. The commonly used symbols are summarized in Table 1. The input to our problem is a  $d$ -dimensional dataset  $D$ . It may contain more than  $d$  attributes, but we assume that the user is interested in exactly  $d$  of them. Since the number of attributes considered in human decision-making is typically limited, similar to [8, 41, 42], we focus on the case where  $d$  is not too large (e.g.,  $d \leq 7$ ), although the developed techniques can be applied to any  $d$ . The  $i$ -th dimensional value of a point  $p$  is denoted by  $p[i]$  for  $i \in [1, d]$ , and the range of value for each attribute is normalized to  $[0, 1]$ . For each attribute, we assume that a higher value is preferred, but our techniques can be easily extended to the case where a small value is preferred.

Following [29, 31, 34, 42, 44, 45], we focus on the family of linear utility functions. That is, the family of functions that express the utility of a point  $p$  as  $f(p) = u \cdot p$ , where  $u$  is a  $d$ -dimensional non-negative vector called the *utility vector*. The utility vector  $u$  encodes the user’s preference, where a higher value at the  $i$ -th dimension (i.e.,  $u[i]$ ) implies that the user is more concerned about the  $i$ -th attribute. We are interested in returning a set of points with a fixed size  $l$  containing the best point, i.e., the point with the highest utility w.r.t  $u$ , which we denote by  $p_{\max} = \arg \max_{p \in D} u \cdot p$ . Note that scaling  $u$  has no influence on the rank of points as well as the best point. Therefore, we assume that  $\sum_{i=1}^d u[i] = 1$ .

The utility vector  $u$  of a user is initially unknown and will be learned by interacting with the user. The interaction continues for several rounds and the user will answer one question in each round. In the rest of the paper, we use the term “round” and “question” interchangeably. We adopt a popular type of questions [31, 34, 37, 41, 42, 44] which is to display two points, namely  $p_i$  and  $p_j$ , from  $D$ , and the user is asked to choose the preferred point. Let  $>$  and  $<$  denote the *ground truth* relation between two points’ utilities, and  $>$  and  $<$  denote the preference *indicated* by the user. That is,  $p_i > p_j$  if  $u \cdot p_i > u \cdot p_j$ , and  $p_i > p_j$  if the user indicates that s/he prefers  $p_i$  to  $p_j$ . The existing algorithms assume that the user *always* selects the point with a higher utility (i.e., if  $p_i > p_j$ , then  $p_i > p_j$  with 100% certainty). In practice, however, the user occasionally chooses the point with a lower utility (e.g., due to the reasons described in Section 1). We say that the user makes an *error* if  $p_i > p_j$  but the indicated preference is  $p_j > p_i$ , and we assume that the user makes an error when comparing each pair of points with an error rate at most  $\theta$ . An error is called a *persistent error* if the answer to the same pair of points is *consistently wrong*; otherwise, it is a *random error*. Notably, persistent errors are harder to be handled compared to random errors. This is because when a user has a random error, if the same question involving two points are asked to him/her *multiple* times, it is possible that s/he could indicate the correct preference in one of the questions, giving a chance for the system to know the inconsistency from the user for the same question implying a possible error. But, when a user has a persistent error, if the same question is asked multiple times, it is not possible that s/he could indicate the correct preference, which means that this increases the difficulty of inferring his/her preference. From now on, we assume that all the errors made are persistent since if we can handle persistent errors, the case for random errors could

$D$	the input dataset
$d$	the dimensionality of $D$
$n$	size of $uhull(D)$
$u$	the utility vector
$\theta$	the upper bound of user error rate
$p_i$	data point in $D$
$P_i$	partition of $p_i$
$h_{i,j}$	the hyperplane related to $p_i$ and $p_j$
$l$	the return size
$t$	number of rounds
$R^i (R_t^i)$	the $i$ -th confidence region (just after round $t$ )
$X^i (X_t^i)$	set of partitions belonging to $R^i (R_t^i)$

**Table 1: Commonly-used Symbols**

be handled automatically. In practice,  $\theta$  can be set to a reasonable number larger than the real error rate, and the exact value of the real error rate need not be known. According to [8, 23],  $\theta$  can be naturally assumed to be less than 0.5 for reasonable users. In their studies,  $\theta$  is at most 5%. Following [14–16, 18, 19, 22, 36], we assume that the user errors are independent across different pairs of points.

In this paper, we are interested in the following problem:

**PROBLEM 1.** *Given a dataset  $D$  with size  $n$ , an error rate upper bound  $\theta$  and a return size  $l$ , how to interact with the user to find a set with size at most  $l$  such that the probability that this set contains the best point is maximized?*

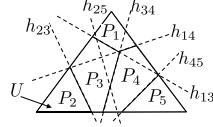
## 4 ALGORITHM FRAMEWORK

In this section, we present the general framework of our proposed algorithms. We first introduce some useful concepts in Section 4.1. Then, in Section 4.2, we present the algorithm framework and prove the lower bound on the required number of rounds for our problem.

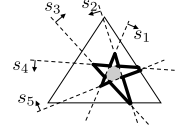
### 4.1 Preliminaries

In geometry, the *convex hull* of a dataset  $D$ , denoted by  $conv(D)$ , is the smallest convex set containing all points in  $D$  [27]. A point  $p \in D$  is a *vertex* of  $conv(D)$  if  $p \notin conv(D/\{p\})$ . Let  $b_i$  be the  $d$ -dimensional point with its  $i$ -th coordinate being 1 and all other coordinates being 0. Furthermore, let  $B = \{b_i | 1 \leq i \leq d\}$  and  $O$  be the origin. Consider the set of points that are both in  $D$  and in  $conv(D \cup B \cup \{O\})$ . We call these points the *upper hull vertices* and denote this set by  $uhull(D)$ . For example, in Figure 1, we visualize a 2-d dataset  $D = \{p_i | i \in [1, 7]\}$ . Its upper hull vertices are  $p_1, p_3, p_4, p_5$  and  $p_6$ . One important conclusion is that the best point must be in  $uhull(D)$  [29, 44]. Therefore, the set of upper hull vertices constitutes possible candidates of the best point. Let  $n$  denote the size of  $uhull(D)$ . For the ease of illustration, in the rest of this paper, when we say  $D$ , we mean  $uhull(D)$  unless otherwise specified.

Recall that all possible utility vectors form a set  $U = \{u | u[i] \geq 0 \text{ and } \sum_{i=1}^d u[i] = 1\}$ .  $U$  is called the *utility space* and is a  $(d-1)$ -dimensional convex polytope. Consider a 3-dimensional example in Figure 2. The utility space  $U$  is a planar triangle with vertices  $(1, 0, 0)$ ,  $(0, 1, 0)$  and  $(0, 0, 1)$ . For any two points  $p_i$  and  $p_j$  in  $D$  where  $i < j$ , we can construct a hyperplane  $h_{ij}$  that has its normal  $p_i - p_j$  and passes through the origin.  $h_{ij}$  intersects with  $U$  and divides the utility space into two *halfspaces* [27]. The halfspace



**Figure 3: Example on partitions**



**Figure 4: Example on confidence regions**

above (resp. below)  $h_{ij}$  is denoted by  $h_{ij}^+$  (resp.  $h_{ij}^-$ ), which contains all utility vectors  $u$  that satisfy  $u \cdot p_i > u \cdot p_j$  (resp.  $u \cdot p_i < u \cdot p_j$ ).

In each round, two points from  $D$ , denoted by  $p_i$  and  $p_j$ , are displayed to the user. When the user indicates the preference between them, we could learn that the utility vector of the user is in  $h_{ij}^+$  or  $h_{ij}^-$ . Let  $s$  denote a halfspace indicated by the user and  $S$  the set of all halfspaces indicated by the user so far. Besides, let  $s_{ij}$  denote the halfspace selected by the user when  $p_i$  and  $p_j$  are displayed (i.e.,  $s_{ij}$  could be either  $h_{ij}^+$  or  $h_{ij}^-$  based on the answer provided by the user). We say that a region (or a point) is supported by a halfspace  $s$  if this region (or this point) lies completely in  $s$ . Given a set  $S'$  of halfspaces, we say that a region (or a point) is supported by set  $S'$  if this region (or this point) is supported by each halfspace in  $S'$ . When the user makes no error, the utility vector  $u$  must lie in the intersection of all halfspaces in  $S$  (i.e.,  $u \in \cap_{s \in S} s$ ). Many existing algorithms [41, 42, 44, 45] utilize this property to approximate the utility vector, which explains why they are prone to user errors: When the wrong halfspace is chosen, the real  $u$  will not locate in the intersection of all halfspaces. Consequently, the resulting estimation of  $u$  is less accurate and the best point may be missed.

For each point  $p_i \in D$ , its corresponding *best partition*  $P_i$ , or *partition* for short, is the set of utility vectors in  $U$  that give  $p_i$  the highest utility score. This means that for any  $u \in P_i$  and for any  $p_j \in D/\{p_i\}$ ,  $u \cdot p_i > u \cdot p_j$ . Therefore,  $P_i$  is the intersection of all  $h_{ij}^+$  for  $p_j \in D/\{p_i\}$  and the utility space  $U$ , i.e.,  $P_i = (\cap_{p_j \in D/\{p_i\}} h_{ij}^+) \cap U$ , which is also a  $(d-1)$ -dimensional convex polytope. As an example, consider Figure 3 where we show the partitions corresponding to a 3-d dataset containing 5 points in its upper hull. The utility space  $U$  is the outer triangle. The partitions  $P_1$  to  $P_5$ , which are 2-d polygons, are bounded by *solid* lines and the intersection of (some)  $h_{ij}$  and  $U$  are shown as *dashed* lines.

Next, we introduce an important concept, called *confidence region*, in Definition 1.

**DEFINITION 1.** *The  $i$ -th confidence region, denoted by  $R^i$ , is the maximal area in the utility space  $U$  that is supported by at least one set of  $|S| - i$  halfspaces in  $S$ . Mathematically,  $R^i = \bigcup_{S' \in S_{|S|-i}} (\cap_{s \in S'} s)$  where  $S_{|S|-i}$  is the set of all  $(|S| - i)$ -subsets of  $S$ .*

If the user makes at most  $i$  errors during the interaction, the utility vector is not supported by at most  $i$  halfspaces in  $S$ . In this case, we immediately have  $u \in R^i$  by its definition. Consider the example shown in Figure 4. The dashed lines represent the hyperplanes asking the user, and the arrow on each hyperplane shows the halfspace indicated by the user. Assume that the user has indicated 5 halfspaces so far, i.e.,  $S = \{s_1, s_2, s_3, s_4, s_5\}$ .  $R^0$  is the area supported by all 5 halfspaces, which is the gray pentagon.  $R^1$  is the area supported by at least one set of all halfspaces except 1 halfspace (i.e., 4 out of 5 halfspaces), which corresponds to the star-shaped area bounded by bold lines. If the user makes at most 1 error when indicating  $s_1, s_2, s_3, s_4$  and  $s_5$ , his/her utility vector must lie in  $R^1$ . Although  $R^0$

must be convex, when  $i \geq 1$ ,  $R^i$  could be non-convex and is in fact made of the union of several disjoint convex polytopes. In Figure 4,  $R^1$  is the union of  $R^0$  and 5 small convex polytopes (triangles in this example) adjacent to  $R^0$ .

We have the following observation.

OBSERVATION 1. For any  $i' \geq i$ ,  $R^i \subseteq R^{i'}$ .

We say that a partition  $P_j$  overlaps with a confidence region  $R^i$  if  $P_j \cap R^i \neq \emptyset$ . Based on Observation 1, if  $P_j$  overlaps with  $R^i$ ,  $P_j$  also overlaps with all  $R^{i'}$ s where  $i' \geq i$ . On the other hand, if  $P_j$  does not overlap with  $R^i$ , it also does not overlap with all  $R^{i'}$ s where  $i' \leq i$ . Since when the user makes at most  $i$  errors, the real  $u$  falls in  $R^i$ , it follows that only points whose partitions overlap with  $R^i$  can be the best point, while those whose partitions do not overlap with  $R^i$  cannot be. For example, if the user makes no errors, then only points whose partitions overlap with  $R^0$  can be the best point. Therefore, when determining which points should be in the returned set, we adopt a strategy that first includes points whose partitions overlap with  $R^0$ , followed by points whose partitions overlap with  $R^1$ ,  $R^2$  and so on.

To make sure that the return size is at most  $l$ , we only maintain confidence regions that guarantee to overlap with at most  $l$  partitions when questions related to all pairs of points (i.e., all possible questions) are asked. The relation between  $R^i$  and the number of partitions it overlaps is shown in Lemma 1.

LEMMA 1. If questions related to all pairs of points are asked, the confidence region  $R^i$  overlaps with at most  $2i + 1$  partitions.

PROOF SKETCH. Denote the set of points whose partitions overlap with  $R^i$  by  $RS$ . Since each pair of points in  $RS$  has been compared, there are in total  $|RS|(|RS| - 1)/2$  comparisons. Each comparison yields one loser. Any points in  $RS$  cannot lose more than  $i$  times, otherwise, it will not be in  $R^i$ . Thus, we have  $|RS|(|RS| - 1)/2 \leq i|RS|$ , so  $|RS| \leq 2i + 1$ . The complete proofs of theorems and lemmas in this paper can be found in Section B of the Appendix.  $\square$

COROLLARY 1. Let  $t_l$  be the number of questions required such that  $R^k$  overlaps with at most  $l$  partitions. If  $k \leq \lfloor \frac{l-1}{2} \rfloor$ , then  $t_l$  is at most the total number of possible questions. In practice,  $t_l$  is much smaller than the total number of possible questions.

If  $k \leq \lfloor \frac{l-1}{2} \rfloor$ , Corollary 1 suggests that  $t_l$  is at most the number of all possible questions. In practice,  $t_l$  is typically much smaller than the number of all possible questions (often less than 30 in our experiments). In the rest of this paper, we set  $k = \lfloor \frac{l-1}{2} \rfloor$  unless otherwise specified.

## 4.2 The General Framework

In this section, we introduce the general framework of our proposed algorithms. We maintain  $k + 1$  confidence regions, namely  $R^0, R^1, \dots, R^k$  (Later in Section 5, we will see that we need not maintain the real  $R^i$ s. Instead, they are only maintained *conceptually*). Note that  $R^i$  in Lemma 1 (and  $R^k$  in Corollary 1) is the final  $R^i$  ( $R^k$ ) after some questions are asked. In our algorithms, since we have not asked any questions at the beginning,  $R^i$ s are initialized to the entire utility space and will be updated when more questions are

asked. We use a subscript  $t$ , i.e.,  $R_t^i$ , to denote  $R^i$  just after  $t$  halfspaces are indicated. Lemma 2 shows the rule of updating  $R_t^i$  where  $i \in [0, k]$ , when the  $(t + 1)$ -th halfspace is indicated.

LEMMA 2. For each  $i \in [0, k]$  and each  $t \geq 0$ , let  $s$  be the halfspace indicated in round  $t + 1$ . The relation between  $R_t^i$  and  $R_{t+1}^i$  is as follows:

$$\begin{aligned} R_{t+1}^i &= R_t^i \cap s & \text{if } i = 0 \\ R_{t+1}^i &= (R_t^i \cap s) \cup (R_t^{i-1} \cap s^-) & \text{if } i \in [1, k] \end{aligned}$$

where  $s^-$  is the complement of  $s$ .

PROOF SKETCH. We prove this lemma using induction. The special case where  $i = 0$  is trivially correct since  $R_t^0$  is the region supported by all  $t$  halfspaces. We then show that if at round  $t'$ , all  $R_{t'}^i$ s are correct, then at round  $t' + 1$ , the resulting  $R_{t'+1}^i$  is indeed the maximal region supported by at least  $t' + 1 - i$  halfspaces.  $\square$

As a running example, consider Figure 5 where  $k = 1$  and  $t = 2$ . Assume that the user indicated  $h_{34}^+$  and  $h_{14}^+$  in round 1 and round 2, respectively (Recall that  $h_{ij}^+$  is the halfspace containing  $P_i$ ). Then,  $R_2^0$  is the gray region and  $R_2^1$  is the area bounded by bold lines. If the user indicates  $h_{13}^+$  the next round, we know from Lemma 2 that  $R_3^0 = R_2^0 \cap h_{13}^+$ , and  $R_3^1 = (R_2^1 \cap h_{13}^+) \cup (R_2^0 \cap h_{13}^-)$ . The resulting  $R_3^0$  and  $R_3^1$  are shown in Figure 6, where  $R_3^0$  is the gray region and  $R_3^1$  is the area bounded by bold lines.

LEMMA 3. For each  $i \in [0, k]$  and each  $t \geq 0$ , we have (1)  $R_{t+1}^i \subseteq R_t^i$  and (2)  $R_t^i \subseteq R_{t+1}^{i+1}$ .

PROOF SKETCH. (1) can be proved using Lemma 2 and the fact that  $R_t^{i-1} \subseteq R_t^i$ . To prove (2), observe that for any point  $v \in R_t^i$ ,  $v$  is supported by a set of  $(t - i)$  halfspaces. The same set of  $(t - i)$  halfspaces makes  $v$  also in  $R_{t+1}^{i+1}$ , which implies that  $R_t^i \subseteq R_{t+1}^{i+1}$ .  $\square$

In round  $t$ , a partition  $P_j$  is said to belong to a confidence region  $R_t^i$  if either (1)  $R_t^i$  ( $i \geq 1$ ) overlaps with  $P_j$  but  $R_t^{i-1}$  does not, or (2)  $R_t^i$  ( $i = 0$ ) overlaps with  $P_j$ . That is, a partition belongs to the smallest confidence region that overlaps with it. Note that knowing the partitions overlapping with each  $R^i$  could help us easily derive the partitions belonging to each  $R^i$ , and vice versa. Similarly, we say that a point  $p_j$  belongs to  $R_t^i$  if its corresponding partition  $P_j$  belongs to  $R_t^i$ . From Lemma 3, confidence regions can only shrink by each round. Due to this shrinking behavior, a partition  $P_j$  that belongs to  $R_t^i$  in round  $t$  may no longer belong to  $R_{t+1}^i$  in round  $t + 1$ . If this happens, we say that  $P_j$  is detached from  $R_{t+1}^i$  (or simply  $R^i$  if the context is clear). Note that after being detached from  $R_{t+1}^i$ ,  $P_j$  will belong to  $R_{t+1}^{i+1}$  in round  $t + 1$  (by Lemma 3,  $P_j$  overlaps with  $R_{t+1}^{i+1}$  since  $R_t^i \subseteq R_{t+1}^{i+1}$ ). When more and more new halfspaces are indicated, a partition is gradually detached from  $R^0, R^1, R^2$  until  $R^k$ . Since we are not interested in partitions that do not overlap with  $R^k$ , if a partition  $P_i$  is detached from  $R^k$ , we say that  $P_i$  (and point  $p_i$ ) is discarded. We use  $X^i$  to denote the set of partitions belonging to  $R^i$ , and use  $X_t^i$  to denote  $X^i$  just after round  $t$ .

Consider our running example in Figure 5. Assume that  $k = 1$  and  $t = 2$ , and recall that  $R_2^0$  is the gray region and  $R_2^1$  is the area bounded by bold lines. Since both  $P_1$  and  $P_3$  overlaps with  $R_2^0$ , we know that both  $P_1$  and  $P_3$  (and thus  $p_1$  and  $p_3$ ) belong to  $R_2^0, P_2$



overlaps with  $R_2^1$  but not  $R_2^0$ , so it belongs to  $R_2^1$ . In the next round ( $t = 3$ ), consider the updated  $R_3^0$  and  $R_3^1$  in Figure 6, where  $R_3^0$  is the gray region and  $R_3^1$  is the area bounded by bold lines. Since  $P_3$  no longer overlaps with  $R_3^0$ , it is detached from  $R_3^0$  and belongs to  $R_3^1$  instead.  $P_2$  is detached from  $R_3^1$  and is thus discarded.

**Stopping Condition.** The algorithm stops when the number of partitions overlapping with  $R^k$  is at most  $l$  (i.e.,  $|\bigcup_{i=0}^k X^i| \leq l$ ). Then, it returns the points that correspond to these partitions.

We are now ready to present the lower bound on the number of rounds required to reach the stopping condition.

**THEOREM 1.** *Given an input size  $n$ , a parameter  $l$  and an error rate upper bound  $\theta$ . There exists a dataset such that any question-asking strategy must ask  $\Omega(\frac{l}{\theta} + n)$  questions to discard at least  $n - l$  points.*

**PROOF SKETCH.** We first prove that there exists a dataset  $D$  such that the following property holds for any  $p_i \in D$  and its partition  $P_i$ : For any non-empty set of halfspaces  $HS \subseteq \{h_{ab}^* | a, b \neq i, * \in \{+, -\}\}$ , if the intersection of halfspaces in  $HS$  forms a non-empty region, then  $P_i$  also overlaps with this region. This property implies that, for any  $p_i \in D$ , if in some round  $P_i$  is detached from some confidence region  $R^m$  where  $m \in [0, k]$ , the halfspace indicated in this round must satisfy one of the following: (1) This halfspace is related to  $p_i$  and some other point  $p_j$ , and  $p_i$  is less preferred to  $p_j$ ; or (2) This halfspace makes  $R^m$  an empty region. We then show that to discard at least  $n - l$  partitions, the round complexity for these two cases are  $\Omega(ln)$  and  $\Omega(\frac{l}{\theta} + n)$ , respectively. Since  $\Omega(ln) > \Omega(\frac{l}{\theta} + n)$ , the lower bound is  $\Omega(\frac{l}{\theta} + n)$ .  $\square$

## 5 THE SS AND FC ALGORITHMS

In this section, we introduce two algorithms, namely *SS* (*Shape-Sampling*) (Section 5.1), and *FC* (*FindCycles*) (Section 5.2), which are developed based on the framework described in Section 4.2. These algorithms differ primarily in the strategies they use to choose questions in each round. For *SS*, we propose two (sub-)strategies that are empirically demonstrated to require only a small number of questions. On the other hand, *FC* employs a strategy that is shown to have an asymptotically optimal round complexity.

### 5.1 The SS Algorithm

In this section, we introduce *SS* by addressing two sub-problems, namely (1) what data structures are required and how to update them in each round, and (2) how to design question selection strategies such that a small number of questions is required.

**5.1.1 Data Structure Maintenance.** While it is possible to directly compute the exact shapes of confidence regions (i.e.,  $R^i$ s), updating these non-convex regions can be computationally expensive when  $d$  increases. However, it is worth noting, as discussed in Section 4.2, that the stopping condition relies solely on the sets of partitions belonging to each confidence region (i.e.,  $X^i$ s). As long as  $X^i$ s can be obtained, maintaining the exact shapes of  $R^i$ s will not be necessary. Therefore, in algorithm *SS*,  $R^i$ s are only kept *conceptually*, and  $X^i$ s are determined using some pre-computed results on a set of randomly sampled points (and thus the name *Shape-Sampling*). By

sacrificing a small degree of accuracy in finding the best point, *SS* achieves an efficient processing time, as demonstrated in Section 6.

*SS* consists of two phases, namely the *pre-processing phase* and the *running phase*. The task of the pre-processing phase is to uniform-randomly sample a number  $Y$  of points from each partition  $P$  using techniques developed in existing studies (e.g., [7]). To distinguish between data points and sample points, we use  $q$  to denote a sampled point and  $P(q)$  to denote the partition where  $q$  is sampled from. After the pre-processing phase is finished, the sampled points can be stored so that future runs can skip this phase and start directly with the running phase. In the running phase, *SS* maintains  $k + 1$  sets, namely  $Q^0, Q^1, \dots, Q^k$ , where  $Q^i$  stores the set of sample points that falls in  $R^i$  where  $i \in [0, k]$ . We use  $Q_t^i$  to denote  $Q^i$  just after round  $t$ .

Next, we describe how to maintain  $Q_t^i$ s and  $X_t^i$ s in the running phase. Initially, when  $t = 0$ , since all partitions overlaps with  $R_0^0, R_0^1, \dots, R_0^k$ , each of  $Q_0^0, Q_0^1, \dots, Q_0^k$  stores the entire set of sample points. In round  $t$ , let  $s$  be the halfspace indicated at this round, the empty-initialized  $Q_t^i$  is constructed using  $Q_{t-1}^i$  and  $Q_{t-1}^{i-1}$  as follows: (1) For each point  $q \in Q_{t-1}^{i-1}$ ,  $q$  is inserted to  $Q_t^i$  if  $q \in s$ ; (2) If  $i > 0$ , for each point  $q \in Q_{t-1}^{i-1}$ ,  $q$  is inserted to  $Q_t^i$  if  $q \notin s$ . It is easy to verify using Lemma 3 that with the above updating rules, each  $Q_t^i$  stores the sample points in  $R_t^i$ . After  $Q_t^i$ s are obtained, with the assumption that the number of sampled points is adequately large (to be discussed next),  $X_t^i$ s can be determined as follows: (1) When  $i = 0$ ,  $X_t^i = \{P(q) | q \in Q_t^i\}$ ; (2) When  $i > 0$ ,  $X_t^i = \{P(q) | q \in Q_t^i\} / \{P(q) | q \in Q_t^{i-1}\}$ . *SS* stops when the size of  $\bigcup_{i=0}^k X_t^i$  is at most  $l$ , and it returns the points that correspond to the partitions in  $\bigcup_{i=0}^k X_t^i$ .

To illustrate *SS*, we use Figure 5 as a running example. Assume that  $l = 3$  and  $k = 1$ , and the number of sample points is sufficiently large. Initially,  $Q_0^0$  and  $Q_0^1$  each contains a set of all sample points. We have  $X_0^0 = \{P_1, P_2, P_3, P_4, P_5\}$  and  $X_0^1 = \emptyset$ . After the user indicated  $h_{34}^+$  and  $h_{14}^+$  in the first 2 rounds,  $Q_2^0$  contains all sample points in  $R_2^0$  (i.e., the gray region) and  $Q_2^1$  contains all sample points in  $R_2^1$  (i.e., the area bounded by bold lines). Then,  $X_2^0 = \{P_1, P_3\}$  and  $X_2^1 = \{P_2\}$ . Since  $|\bigcup_{i=0}^k X_2^i| = |\{P_1, P_2, P_3\}| = 3 \leq l$ , *SS* stops and returns  $\{P_1, P_2, P_3\}$ .

Since *SS* uses sample points to decide  $X^i$ s, when it terminates, even if the partition corresponding to the best point still overlaps with  $R^k$ , the best point will not be returned if there is no sample point in the intersection of this partition and  $R^k$ . To make this probability small, the sample size should be sufficiently large. Lemma 4 decides a sufficient number of samples which guarantees that this case happens with a probability at most a user parameter  $\epsilon$ .

**LEMMA 4.** *The intersection of hyperplanes in set  $\{h_{ij} | p_i, p_j \in D\}$  divides the utility space into a number of cells, where each cell corresponds to a unique ranking of points in  $D$ . Given a parameter  $\epsilon$ , if the number of cells in any partition is at most  $p$ , by sampling  $Y = \frac{p}{\epsilon\epsilon}$  points from each partition, *SS* returns the best point with probability at least  $1 - \epsilon$  if the cell where the real utility vector lies overlaps with  $R^k$ .*

**PROOF SKETCH.** When the algorithm stops, if the cell where the real utility vector lies overlaps with  $R^k$ , then the best point will not



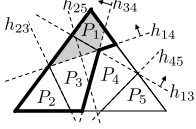


Figure 5: Confidence regions before asking  $h_{13}$

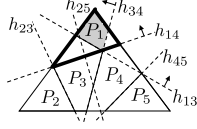


Figure 6: Confidence regions after asking  $h_{13}$

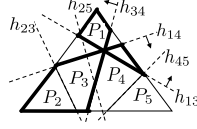


Figure 7: FC example just after round 3

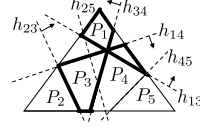


Figure 8: FC example just after round 4

be returned only if there is no sample point in this cell. We show that the upper bound of this event's probability is  $\frac{\rho}{Y}e^{-1}$ . Solving  $\frac{\rho}{Y}e^{-1} \leq \varepsilon$  yields the lemma.  $\square$

Although in the worst case  $Y = O(\frac{n^d}{\varepsilon})$  [27], in practice,  $Y$  can be set by gradually increasing it until the accuracy of returning the best point converges, and thus,  $Y$  could be set to a much smaller number. In our experiments, on a 4-d dataset with 100000 points, by sampling 1000 points from each partition, the accuracy already converges, indicating a small accuracy loss caused by sampling.

**5.1.2 Question Selection Strategies.** Next, we describe how  $SS$  chooses the question in each round so that the required number of questions can be reduced. Observe that to detach a partition  $P$  from a confidence region  $R^i$ ,  $P$  must also be detached from all  $R^{i'}$ 's where  $i' \leq i$ . Thus, intuitively, a question is more preferred if (1) it can detach partitions from  $R^i$  with a smaller  $i$  and (2) it can detach a larger number of partitions. For example, questions that can detach a large number of partitions from  $R^0$  are the most preferred. We develop two question selection strategies that comply with the above intuition, namely the *random-based selection* and the *score-based selection*. The random-based selection has a faster process time, and the score-based selection asks fewer questions empirically. In the following, we describe the two selection strategies in detail.

**The Random-based Selection.** In round  $t$ , let  $X_t^i$  be the set of partitions that belongs to  $R_t^i$ , and let  $X_t^i$  be the set of points whose partition is in  $X_t^i$ . The random-based selection runs for at most  $k+1$  iterations. In iteration  $i$  ( $i \in [1, k+1]$ ), it sets a *candidate point set*, denoted by  $CP$ , to be  $\bigcup_{j=0}^{i-1} X_t^j$ , enumerates all (order-insensitive) pairs of points consisting of points in  $CP$  and randomly permutes them, and sequentially checks each pair that has not been checked in previous iterations. If it finds a pair that has not been asked in previous questions, it selects this pair and stops. If such a pair cannot be found after depleting all pairs in this iteration, it starts the next iteration. Note that this strategy favors questions that can detach partitions belonging to  $R_t^i$  with a small value of  $i$ . It also guarantees to find an unasked pair if one exists, since otherwise, the stopping condition described in Section 4.2 is already met.

As a running example, assume that  $k=1$  and at round  $t$  we have  $X_t^0 = \{p_1, p_2, p_3\}$  and  $X_t^1 = \{p_4, p_5\}$ . In the first iteration,  $CP$  is set to  $\{p_1, p_2, p_3\}$ . We randomly permute all pairs of points consisting of points in  $CP$  (i.e.,  $(p_1, p_2)$ ,  $(p_2, p_3)$  and  $(p_1, p_3)$ ) and sequentially consider each of them. If we find a pair that is unasked before (e.g.,  $(p_1, p_2)$ ), we use this pair. If all 3 pairs are already asked, we enter the second iteration, where  $CP$  is set to  $X_t^0 \cup X_t^1$  (i.e.,  $\{p_1, p_2, p_3, p_4, p_5\}$ ).

**The Score-based Selection.** Although the random-based selection has a fast processing time, it does not fully utilize the distribution of partitions to select the optimal hyperplane. Therefore, we design the score-based selection to further reduce the number of

questions required. We first introduce a data structure that will be used for this strategy. For each partition  $P_a$  and a hyperplane  $h_{ij}$ , there are 3 possible relationships between  $P_a$  and  $h_{ij}$ : (1)  $P_a \in h_{ij}^+$ , (2)  $P_a \in h_{ij}^-$ , and (3)  $P_a$  intersects with  $h_{ij}$ . Consider the example in Figure 5. For hyperplane  $h_{13}$ ,  $P_1$  is in  $h_{13}^+$ ,  $P_2, P_3$  and  $P_5$  are in  $h_{13}^-$ , and  $P_4$  intersects with  $h_{13}$ . We maintain a table  $L$  to store these relationships, where each row corresponds to a hyperplane  $h_{ij}$ .  $L$  has 3 columns: (1)  $h_{ij}^+$ , which stores all partitions that lie in  $h_{ij}^+$ , (2)  $h_{ij}^-$ , which stores all partitions that lie in  $h_{ij}^-$ , and (3) *score* which will be explained next. The table  $L$  corresponding to Figure 5 is shown in Table 2.

Let  $Num(s, X_t^i)$  denote the number of partitions in  $X_t^i$  that lies in a halfspace  $s$ . We define the *score* of a hyperplane  $h$  at round  $t$  to be  $score_t(h) = \min(\sum_{i=0}^k \alpha^i Num(h^+, X_t^i), \sum_{i=0}^k \alpha^i Num(h^-, X_t^i))$ , where  $\alpha$  is a parameter between 0 and 1 capturing the relative priority between different  $X_t^i$ 's. In practical applications,  $\alpha$  can be determined by conducting a grid search between 0 and 1 and finding the value that minimizes the required number of rounds. We set  $\alpha = 0.2$  in this paper based on the empirical results in Section 6.2. Note that this definition gives higher scores to those hyperplanes intersecting  $R_t^i$  with a smaller  $i$  and with partitions in  $X_t^i$  more evenly distributed on each side, indicating a higher chance of detaching more partitions. In round  $t$ , the hyperplane with the highest score that has not been chosen before will be selected and its corresponding points will be displayed. Consider the example in Figure 5 where  $t=2$  and  $k=1$ . After  $h_{34}^+$  and  $h_{14}^+$  are indicated by the user in the first 2 rounds, the table  $L$  corresponding to this stage is shown in Table 2. Based on this table,  $p_1$  and  $p_3$  will be displayed to the user in the next round since  $h_{13}$  obtains the highest score.

Based on the question selection strategy applied, there are two variants of  $SS$ , which are called *SS-random* and *SS-score*, respectively. The time complexities for each round of these two variants are presented in Theorem 2.

**THEOREM 2.** *Given an input size  $n$ , dimensionality  $d$  and the return size  $l$ . Let  $Y$  be the number of samples from each partition. The time complexity for the pre-processing phase is  $O((|V|+Y)dn^2)$ , where  $|V|$  is the maximum number of vertices in all partitions. In the running phase, the time complexities in each round of *SS-random* and *SS-score* are  $O(Yl dn)$  and  $O(Yl dn + n^3)$ , respectively.*

**PROOF SKETCH.** The time complexity for the pre-processing phase is  $O((|V|+Y)dn^2)$  since computing all partitions takes  $O(|V|dn^2)$  time and sampling all  $Yn$  points takes  $O(Ydn^2)$  time [7]. In each round of the running phase, the time required to update  $Q^i$ 's and compute  $X^i$ 's is  $O(Yl dn)$ , and the time required to decide the next question for *SS-random* and *SS-score* are  $O(t)$  and  $O(n^3)$ , respectively, where  $t$  is the current number of rounds. Since typically  $t \ll Ydn$ , the total time complexity then follows.  $\square$

	$h_{ij}^+$	$h_{ij}^-$	score ( $\alpha = 0.2$ )
$h_{13}$	$\{P_1\}$	$\{P_2, P_3, P_5\}$	1
$h_{14}$	$\{P_1\}$	$\{P_2, P_4, P_5\}$	0.2
$h_{23}$	$\{P_2\}$	$\{P_1, P_3, P_4, P_5\}$	0.2
$h_{25}$	$\{P_2\}$	$\{P_5\}$	0
$h_{34}$	$\{P_2, P_3\}$	$\{P_4, P_5\}$	0
$h_{45}$	$\{P_1, P_2, P_3, P_4\}$	$\{P_5\}$	0

Table 2: Table L

It is worth mentioning that  $|V|$  is not large in typical scenarios. From [40],  $|V| = O(m^{\lfloor \frac{d}{2} \rfloor})$  where  $m$  is the maximum number of halfspaces bounding a polytope. Typically,  $m \ll n$  although it can be  $n - 1$  in the worst case. For our experiment where  $d = 5$  and  $n = 10,000$ , the value of  $m$  is smaller than 100. Besides, the value of  $d$  is not large (at most 7 in most cases) due to the limited number of attributes considered by humans in decision-making [8, 30, 42, 44].

Lastly, Theorem 3 bounds the probability that the best point is returned by SS. Intuitively, when  $l = 10$ ,  $\theta = 0.05$ ,  $T = 20$  and  $\varepsilon = 0.01$ , Theorem 3 guarantees that SS find the best point with probability at least 82%. Note that this is a loose bound. In our experiments, if we set  $l = 10$ , our algorithms return the best point with probability at least 99%.

**THEOREM 3.** *Given an error rate upper bound  $\theta$ , a return size  $l$ , and a parameter  $\varepsilon$  as defined in Lemma 4. If SS terminates in  $T$  rounds, then the probability that the best point is returned by SS is at least  $1 - \varepsilon - e^{-\frac{((l-1)/2) - \theta T}{((l-1)/2) + \theta T}}$ .*

**PROOF SKETCH.** SS does not return the best point only if either (1) there is no sample point in the cell containing the real utility vector; or (2) the cell containing the real utility vector does not overlap with  $R^k$ . By Lemma 4, (1) happens with probability at most  $\varepsilon$ . By Chernoff inequality, (2) happens with probability at most  $e^{-\frac{((l-1)/2) - \theta T}{((l-1)/2) + \theta T}}$ . Summing them together completes the proof.  $\square$

## 5.2 The FC Algorithm

Although SS typically requires only a small number of questions in practice, it does not have a guarantee on the number of questions needed before termination. In this section, we introduce our second algorithm FC, which has an asymptotically optimal number of rounds. We first introduce an important concept called cycle.

**DEFINITION 2.** *A set of  $m \geq 3$  points  $\{p_1, p_2, \dots, p_m\}$  form a **cycle** if the user indicates that  $p_1 > p_2, \dots, p_{m-1} > p_m$ , and  $p_m > p_1$ . Two cycles are called **disjoint** if there are no two points,  $p_i$  and  $p_j$ , such that  $p_i > p_j$  appears in both cycles.*

Due to possible user errors, we do not assume transitivity in the indicated preferences. That is,  $p_i > p_j$  and  $p_j > p_k$  (indicated by a user) do not imply  $p_i > p_k$ . Lemma 5 reveals the relation between the occurrences of cycles and confidence regions.

**LEMMA 5.** *If there are  $i + 1$  disjoint cycles, then the  $i$ -th confidence region  $R^i = \emptyset$ .*

**PROOF SKETCH.** We first show that the intersection of halfspaces corresponding to a cycle is an empty region. Since  $R^i = \bigcup_{S' \in S_{|S|-i}} (\cap_{s \in S'} s)$ , where  $S'$  has a cardinality of  $|S| - i$ , if there are

$i + 1$  disjoint cycles, no matter how we set  $S'$ ,  $\cap_{s \in S'} s$  will be empty since there is at least one cycle in  $S'$ . Thus,  $R^i$  is also empty.  $\square$

Intuitively, FC has two stages, namely Stage 1 and Stage 2, and each stage consists of several rounds. Similar to SS, FC only keeps confidence regions conceptual and utilizes the same sampling technique described in Section 5.1.1 to maintain  $X^i$ s, i.e., the sets of partitions belonging to each confidence region. After each round, it checks if it stops by checking if  $|\bigcup_{i=0}^k X^i| \leq l$ . In Stage 1, FC adopts a question selection strategy that is different from the two strategies described in Section 5.1.2. This strategy attempts to quickly obtain disjoint cycles (if cycles appear during interaction), thereby reducing some confidence regions to empty and detaching a large number of partitions from them. When Stage 1 ends, it is shown later in Lemma 6 that at most  $k$  disjoint cycles are obtained. Let  $j$  be the number of disjoint cycles obtained just after Stage 1 (where  $j \leq k$ ). By Lemma 5, this means that  $R^0, R^1, \dots, R^{j-1}$  all become empty. Then, FC enters Stage 2. At this stage, a partition either belongs to one of the non-empty confidence regions  $R^j, R^{j+1}, \dots, R^k$ , or is already discarded. In this stage, FC chooses questions following either the random-based selection or the score-based selection in Section 5.1.2, until the stopping condition is met.

Next, we describe the details of these 2 stages. In Stage 1, FC maintains a set  $PS$  of points, initialized to an empty set. Initially, the algorithm randomly selects two points from  $D$ , asks the user's preference on them, and inserts them into  $PS$ . Then, it randomly picks a point from  $D$  that has not been picked before, which is called the *focusing point*, denoted by  $p_f$ , and asks the user's preferences between  $p_f$  and each point in  $PS$ . After all points in  $PS$  are compared with  $p_f$ ,  $p_f$  is inserted into  $PS$ . FC then selects a new point from  $D$  as the new focusing point and repeats the above process. Stage 1 stops when one of the following two conditions is met: (1) it already lasts for  $\max(\frac{k(k-1)}{2}, \frac{3k}{\theta})$  rounds just after a focusing point is inserted into  $PS$ ; or (2)  $\bigcup_{i=0}^{k-1} X^i = \emptyset$ . Lemma 6 states several important properties when Stage 1 ends.

**LEMMA 6.** *When Stage 1 ends, the following properties hold:*

1. *The number of obtained disjoint cycles is at most  $k$ .*
2. *If Stage 1 ends on Condition (1), then the expected number of obtained disjoint cycles is at least  $\max(0, k - 17)$ .*

**PROOF SKETCH.** To prove Property 1, assume that Stage 1 ends after round  $t$ . Since it does not end after round  $t - 1$ ,  $\bigcup_{i=0}^{k-1} X_{t-1}^i \neq \emptyset$ , it follows that  $R_{t-1}^{k-1} \neq \emptyset$ . By Lemma 3,  $R_{t-1}^{k-1} \subseteq R_t^k \neq \emptyset$ . Thus, the number of disjoint cycles must not exceed  $k$  (Lemma 5). To prove Property 2, we show that if Stage 1 ends on Condition (1), the user makes at least  $k$  errors with probability at least  $1 - \frac{1}{k}$ . We then show that if at least  $k$  errors are made, the expected number of disjoint cycles is at least  $k - 16$ . The expected number of obtained disjoint cycles is thus at least  $(1 - \frac{1}{k})(k - 16) \geq k - 17$ .  $\square$

FC then enters Stage 2, in which it needs to further discard the remaining partitions. To do so, FC follows either the random-based selection or the score-based selection in Section 5.1.2 to select a pair of points and display them to the user in each round. The algorithm terminates when at most  $l$  partitions remain.

To illustrate *FC*, assume that  $l = 3$  and  $k = 1$ , and in the first 3 rounds of Stage 1 the user indicated  $p_1 > p_3$ ,  $p_3 > p_4$  and  $p_4 > p_1$ , which forms a cycle. The related halfspaces  $h_{13}^+$ ,  $h_{34}^+$  and  $h_{14}^-$  and the resulting confidence regions are shown in Figure 7, where  $R_3^0$  is an empty region and  $R_3^1$  is the union of the three regions bounded by bold lines. Since  $X_3^0 = \emptyset$  and  $X_3^1 = \{P_1, P_2, P_3, P_4\}$ , Stage 1 ends because Condition (2) is satisfied. After the user indicated  $h_{23}^-$  in round 4, the resulting  $R_4^1$  is the union of the three regions bounded by bold lines shown in Figure 8. Since  $X_4^0 = \emptyset$  and  $X_4^1 = \{P_1, P_3, P_4\}$ ,  $|\bigcup_{i=0}^k X_4^i| = |\{P_1, P_3, P_4\}| = 3 \leq l$ . *FC* stops and returns  $\{p_1, p_3, p_4\}$ .

Theorem 4 presents the major conclusion of *FC*.

**THEOREM 4.** *Given an input size  $n$ , an output size  $l$  and an error rate upper bound  $\theta$ , *FC* returns a set of points with size at most  $l$  using  $O(\frac{1}{\theta} + n)$  rounds on expectation.*

**PROOF SKETCH.** Firstly, note that Stage 1 finishes within  $O(\frac{1}{\theta} + n)$  rounds as long as  $l = O(\sqrt{n})$ . When entering Stage 2, there are 2 cases: (1) Condition 1 is satisfied, and (2) Condition 1 is not satisfied, but Condition 2 is satisfied. We then show that the expected round complexity combining these 2 cases is  $O(n)$ . Therefore, the total expected round complexity is  $O(\frac{1}{\theta} + n)$ .  $\square$

**COROLLARY 2.** *The round complexity of *FC* is asymptotically optimal.*

Lastly, Theorem 5 bounds the probability of *FC* returning the best point.

**THEOREM 5.** *Given an error rate upper bound  $\theta$ , a return size  $l$ , and a parameter  $\epsilon$  as defined in Lemma 4. If *FC* terminate in  $T$  rounds, then the probability that the best point is returned by *FC* is at least  $1 - \epsilon - e^{-\frac{((l-1)/2 - \theta T)^2}{((l-1)/2 + \theta T)}}$ .*

**PROOF SKETCH.** The proof is nearly identical to the proof of Theorem 3.  $\square$

## 6 EXPERIMENT

### 6.1 Experimental Setup

Our experiments were conducted on a computer with 3.10 GHz CPU and 64GB RAM. All programs were implemented in C/C++.

**Datasets.** We conducted experiments on synthetic and real datasets. Statistics of the datasets are summarized in Section A of the Appendix. For synthetic datasets, we generated *anti-correlated* datasets using a dataset generator developed for skyline operators [6]. For real datasets, we used 3 real datasets: *AirQuality*, *Weather*, and *HTRU*. *AirQuality* has 420,478 tuples with 4 attributes, *Weather* includes 96,483 weather records with 6 attributes, and *HTRU* has 17,898 points with 7 attributes. Each dimension is normalized into the range of  $[0, 1]$ . We preprocessed all the datasets to contain only the skyline points, which are the possible best points for any utility function.

**Algorithms.** We compare our proposed algorithms, namely *FC*, *SS-score* and *SS-random*, against the competitor algorithms, including (a) algorithms that do not consider user errors, namely *HD-PI* [42], *UH-Simplex* [44] and *UtilApprox* [31], and (b) algorithms that consider user errors, namely *Verify-Point* [8], *Active-Ranking* [20]

and *Pref-Learn* [34]. Note that [8] also proposes another algorithm *Verify-Space*. Since its performance is similar to *Verify-Point* in [8], we do not include it here. To make the comparison fair, we adapt each of them to return at most  $l$  points that are most likely to be the best point. The adaptations are summarized below.

(1) Algorithms *HD-PI*, *Verify-Point* and *UH-Simplex* all maintain a set of candidate points during their interaction processes. We return all points in the candidate set when its size is no more than  $l$ . Specifically, since the candidate set maintained by *HD-PI* stores the possible top- $k$  points, we set  $k$  to 1. The set maintained in *Verify-Point* stores the possible best points, which need no further adaption. In *UH-Simplex*, the candidate set contains points with regret ratios possibly lower than a parameter  $\epsilon$ . Following [8, 42], we set  $\epsilon$  to  $1 - f(p_2)/f(p_1)$ , where  $p_1$  and  $p_2$  are the best and second best points according to the utility vector, which is equivalent to finding the best point. (2) Algorithm *Active-Ranking* aims at learning the entire ranking of all points by interacting with the user. We return the top- $l$  points after the entire ranking is obtained. (3) Algorithm *Pref-Learn* interacts with the user to learn the user's utility vector. Algorithm *UtilApprox* returns points with regret ratios smaller than a parameter  $\epsilon$  by estimating the user's utility vector. For these two algorithms, we return the top- $l$  points w.r.t to the learned utility vector after the learning processes finish. Specifically, for *Pref-Learn*, we set its error threshold to  $10^{-6}$  since according to [34], the learnt vector is very close to the theoretical optimum if the error threshold is less than  $10^{-5}$ . For *UtilApprox*, we set  $\epsilon$  in the same way as *UH-Simplex* ( $\epsilon = 1 - f(p_2)/f(p_1)$ ) to find the best point.

**Parameter Setting.** We evaluate the performance of each algorithm by varying different parameters: (1) the dataset size  $N$ , (2) the dimensionality  $d$ , (3) the user error rate upper bound  $\theta$ , (4) the return size  $l$ , (5) the parameter  $\alpha$  in the score-based selection (Section 5.1.2), and (6) the parameter  $Y$  in Theorem 2 controlling the number of samples from each partition. The default setting for each synthetic dataset is  $N = 100,000$  and  $d = 4$ . The default value of  $\theta$  is 0.05, which, according to the human reliability assessment data in [23], is a reasonable upper bound for the human error rate. According to the results in Section 6.2, we set the default value  $l = 5$ ,  $\alpha = 0.2$ , and  $Y = 1000$ .

**Performance Measurement.** The performance of each algorithm is evaluated by the following measurements: (1) *Accuracy* which is the probability that the best point is returned. Formally, accuracy is defined as  $\frac{T_{ret}}{T_{tot}}$  where  $T_{tot}$  is the total number of trails and  $T_{ret}$  is the number of times the best point is returned. (2) *Number of questions* required to return the points. (3) *Processing time* which is the average processing time to decide the next question. We report the processing time per question since compared to the total processing time, it is a more informative metric for evaluating the algorithm's responsiveness during user interaction. Each setting is repeated 100 times and the average value is reported. In each repetition, we randomly sample a vector  $u$  from the utility space as the underlying utility vector, and the point with the highest utility with respect to  $u$  in the dataset is regarded as the real best point.

The rest of the paper is organized as follows. In Section 6.2, we analyze the impact of various parameters on the performance of our algorithms. We then present the experimental results on synthetic datasets (Section 6.3) and real datasets (Section 6.4). The findings

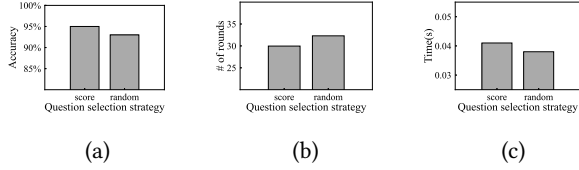


Figure 9: Effect of question selection strategies on FC

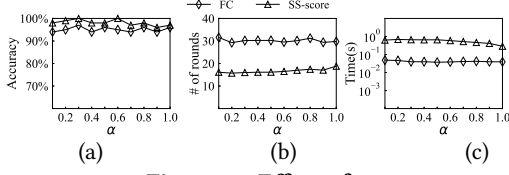


Figure 11: Effect of  $\alpha$

from a user study are discussed in Section 6.5. Finally, we provide a summary of the experiments in Section 6.6.

## 6.2 Experiments on Parameter Setting

In this section, we evaluate the impact of different parameter settings, including  $l$ ,  $\alpha$  and  $Y$ , on our algorithms.

We studied the effect of using the score-based selection and the random-based selection in Stage 2 of *FC*. Figure 9 shows the results. We observe that the score-based selection obtains a higher accuracy and requires fewer questions, while the processing time of the random-based selection is slightly faster. Since the increase in accuracy and round efficiency is considered more important than a small gain in processing time, we chose the score-based selection as the default question selection strategy for Stage 2 of *FC*.

Figure 10 analyzes the impact of varying the value of  $l$  from 1 to 9 on our algorithms. According to Figure 10 (a) and (b), when  $l$  increases, the accuracies and the number of questions of all our algorithms also increase. This is because with a larger value of  $l$ , the best point is less likely to be discarded, but more questions are required to discard other points. Based on these results, we select  $l = 5$  as our default setting since it yields a high level of accuracy while asking a small number of questions.

In Figure 11, we varied  $\alpha$  from 0.1 to 1.0 to study its impact on *FC* and *SS-score* (note that  $\alpha$  is a parameter in the score-based selection so it does not affect *SS-random*, which uses the random-based selection). Changing  $\alpha$  does not significantly affect the accuracies of the algorithms. We chose  $\alpha = 0.2$  since it minimized the number of questions needed for both algorithms. In Figure 12, we studied the influence of increasing parameter  $Y$  from 50 to 5000 on the performance of *FC*, *SS-score* and *SS-random*. As  $Y$  increases, all algorithms exhibit improved accuracies, a higher number of questions, and a longer processing time. The time required by the pre-processing phase also increases. We chose  $Y = 1000$  as it provides a satisfactory level of accuracy while maintaining a fast processing time.

## 6.3 Experiments on Synthetic Datasets

Figure 13 summarizes our study on the algorithms' performance when handling different types of user errors: random errors, persistent errors and a combination of both. In this figure, "random" means that all errors are random errors, "persist" means that all errors are persistent errors, and "combined" means that 50% of the errors are persistent errors and the rest are random errors. Except for *Verify-Point*, all algorithms exhibit similar performance across the three error types, since they avoid asking repeated questions, making persistent and random errors equivalent. *Verify-Point*

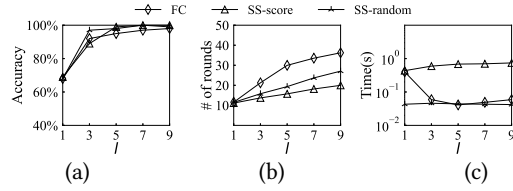


Figure 10: Effect of  $l$  (and  $k$ )

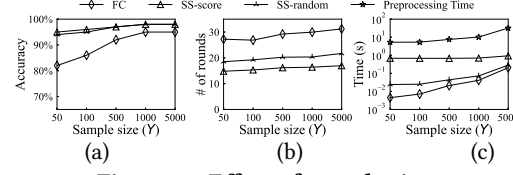


Figure 12: Effect of sample size

achieves slightly lower accuracy than our algorithms when dealing with random errors. However, since its technique (i.e., asking the same question several times) fails to address persistent errors, its accuracy decreases when the fraction of persistent errors increases. Since persistent errors are harder to be handled, we present algorithm performance assuming all user errors to be persistent for the rest of this section.

In Figure 14, we compare the performance of our algorithms (*FC*, *SS-score*, and *SS-random*) with existing methods on 4-d synthetic datasets of varying sizes (from 100 to 100 million). As shown in Figure 14 (a), our algorithms consistently outperform existing methods in terms of accuracy, with a widening performance gap when the dataset size increases. They achieve over 10% higher accuracy than the closest competitor (*UtilApprox*) for large input sizes (100 million). Among our algorithms, *SS-score* is the most round-efficient, asking at most 10 more questions compared to the most round-efficient one (*HD-PI*). *SS-random* asks slightly more questions than *SS-score*. However, it scales well on large datasets (100 million points) since it determines the next question within 0.7 seconds.

Figure 15 shows the effect of varying  $\theta$  from 0 to 0.15. According to Figure 15 (a), our algorithms achieve the highest accuracies, decrease at the slowest rates and remain above 75% even with high error rates (e.g., 0.15), but all other methods fall below 65%. Increasing  $\theta$  does not have a significant impact on the number of questions required by our algorithms, except for *FC*, whose number of questions decreases when  $\theta$  increases since cycles can be obtained more quickly when more inconsistencies are involved.

Figure 16 shows our algorithms' scalability with increasing dimensionality  $d$ . Our algorithms consistently achieve higher accuracies for all dimensional settings compared to existing methods, and the difference in accuracy grows even larger when  $d$  increases. Besides, they require only 5 to 8 additional questions per dimensionality increase. Although *SS-score* takes around 10 seconds for  $d = 5$  (since finding the best hyperplane in the large table  $L$  is time-consuming), the processing time of *SS-random* and *FC* is still within 1 second.

## 6.4 Experiments on Real Datasets

We compared the performance of our algorithms against the existing methods on 3 real datasets, namely *AirQuality*, *Weather* and *HTRU*. The results are summarized in Figure 17, 18 and 19, respectively. Our methods obtain higher accuracies than existing algorithms on all 3 datasets. In particular, the accuracies of *FC*, *SS-score*, and *SS-random* are consistently above 90%. Although existing algorithms *HD-PI* and *Verify-Point* require 5 to 7 fewer rounds than

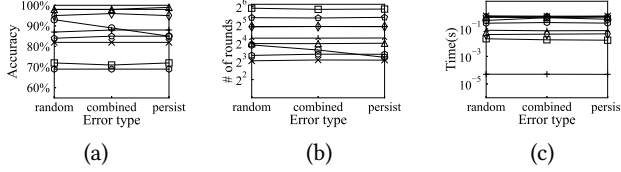


Figure 13: Effect of different types of errors

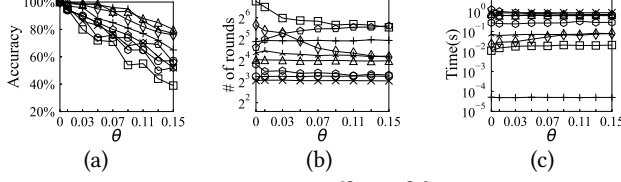


Figure 15: Effect of  $\theta$

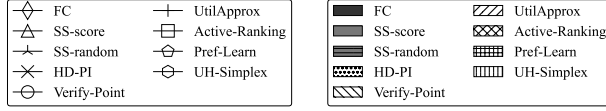


Figure 17: Results on dataset *AirQuality*

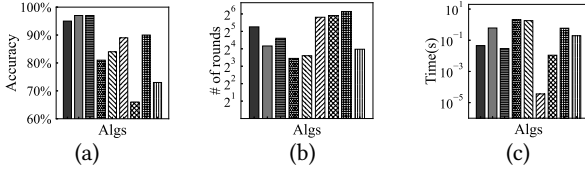


Figure 18: Results on dataset *Weather*

our most round-efficient algorithm (*SS-score*), their accuracies are 10% to 20% lower than ours, which is not satisfactory. Our algorithm runs at an interactive speed since they process each question within 0.6 seconds on all 3 datasets.

## 6.5 User Study

We conducted two user studies on a real dataset *Airbnb* [11]. *Airbnb* consists of 6809 Airbnb rentals in Amsterdam with 4 attributes: daily price, cleanliness rating, location rating and the number of reviews. Following [8, 42], we randomly sampled 1000 Airbnb rentals and recruited 25 participants.

(1) The first user study aims to study user errors' impact on algorithm performance and our algorithms' effectiveness. We compared *SS-score* against existing methods, namely *Verify-Point*, *HD-PI*, *Active-Ranking* and *Pref-Learn*. For *Pref-Learn*, since it is hard to obtain the user's real utility vector, it is re-adapted following [8, 42]: It maintains an estimated utility vector  $u$ . If 75% [34] of some randomly selected questions answered by the user can be correctly predicted using  $u$ , it stops and returns the top- $l$  points w.r.t.  $u$ . We set  $l = 5$  and derive that  $k = 2$ , following our experimental settings. This user study contains 3 parts. In Part 1, the user interacted with each algorithm for several rounds until the algorithm returns a list of at most  $l$  items. During each round, two Airbnb rental options were shown, and the user was asked to select the preferred option. Once the list was returned, the user was required to select one item from the list as the *selected favorite rental option* of the algorithm. After the selected favorite rental options of all algorithms are chosen by the user, the user was asked to select one option among all of these selected favorite rental options as his/her *tentative best point*. Then, a set of additional questions were asked to confirm whether

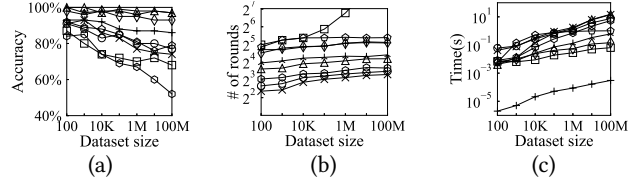


Figure 14: Effect of input size on 4d datasets

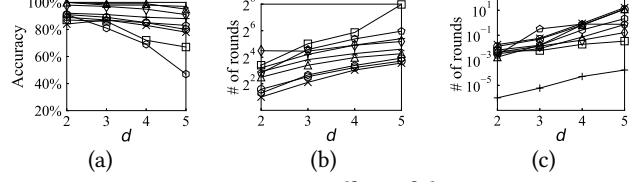


Figure 16: Effect of  $d$

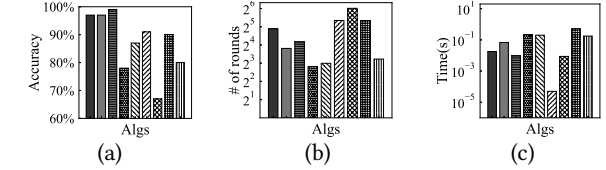


Figure 19: Results on dataset *HTRU*

this *tentative best point* was indeed preferred to each of the other selected favorite rental options. Whenever each of these selected favorite rental options, says  $p$ , was more preferred to the tentative best point  $p_{best}$ , an additional question was asked to compare these two points (i.e.,  $p$  and  $p_{best}$ ) and the point preferred by the user with more questions became the new tentative best point. After we finish the process, the current tentative best point is regarded as the best point of the user.

Part 2 and Part 3 follow the design of Part 1 with the only difference that in these two parts, each algorithm was required to stop and return a list of at most  $l$  options after  $t$  rounds. In Part 2,  $t$  is set to 10. In Part 3,  $t$  is set to the number of rounds required by *SS-score* to terminate (which is typically 12 to 13). Since some algorithms cannot guarantee the size of the returned list when forced to stop, they are re-adapted as follows: (1) For *SS-score*, we randomly return  $l$  points whose partition is in  $\bigcup_{i=0}^k X^i$ . (2) For *Verify-Point* and *HD-PI*, we randomly return  $l$  points from their candidate sets. (3) For *Active-Ranking*, we return the top- $l$  points in the topological order resulting from the  $t$  comparisons. For each part, we evaluated algorithms' performance with the following metrics: (1) the *hit rate* which is the probability that the best point is included in the returned list, (2) the number of rounds required to return the list, and (3) the average processing time to decide the next question. The average scores of all participants are reported.

Figure 20 displays the results. In Part 1, *SS-score* achieves over 90% hit rate, significantly outperforming other competitors. The processing time of our algorithm is also short since it determine the next question within 0.01 seconds. Given that *SS-score* obtains a hit rate exceeding 90%, we conclude that the linear utility function

approximates the user’s preference reasonably well. In Part 2, the hit rate of *SS-score* is lower than *HD-PI* and *Verify-Point*. This is explainable because there is a trade-off between round efficiency and error handling capacity. Since *SS-score* prioritizes error handling, it is expected that its round efficiency is lower than algorithms that either do not consider error handling or have a lower error handling capacity. Consequently, when forced to stop after round 10, *SS-score* has to randomly select  $l$  points from a relatively large number of candidates, resulting in a lower hit rate. However, as can be observed in Part 3, when allowed to use slightly more rounds (i.e., 2 to 3 more rounds), the hit rate of *SS-score* exceeds other baselines by a large margin, aligning with its performance in Part 1. As stated in Section 1, there are situations where missing the best point can cause unforgettable and unchangeable consequences, such as choosing the university and purchasing a house. In these scenarios, the user’s primary concern is the accuracy of finding the best point, even if slightly more questions need to be answered.

Based on the results obtained in Part 1, we estimated how frequently users make errors using algorithm *SS-score*, using the properties that if  $Y$  is large enough, the best point belongs to the  $i$ -th confidence region only if at least  $i$  errors are made, and the best point is not in the returned set only if at least  $k + 1$  errors are made. We regard the best point found in Part 1 by our user study as the real best point of the user (since we used additional checking questions to make sure that this point is the real best point with high probability) and record necessary information that decides the confidence region this point belongs to. The user error rate can be estimated as  $\frac{\sum_x e_x}{\sum_x t_x}$ , where  $e_x$  is the number of errors made by user  $x$  (calculated based on the above properties) and  $t_x$  is the number of rounds used by user  $x$ . Among 272 questions asked by *SS-score*, users made 16 errors, resulting in a 5.8% empirical error rate.

(2) Our second user study aims to show the presence of persistent user errors in the real-world scenarios. In this user study, each Airbnb has 5 attributes: labeled price, discount, cleanliness rating, location rating and the number of reviews. For brevity, we will refer to cleanliness rating, location rating, and the number of reviews as the “other attributes” throughout this section. The term *labeled price* represents the original price before the discount, while the *final price* refers to the price after the discount has been applied (i.e., final price = labeled price  $\times$  (100% – discount)). For example, given an option with a labeled price of \$300 and a 20% discount, its final price is  $\$300 \times (100\% - 20\%) = \$240$ . We have a hypothesis that some people in the world have the (wrong) impression that the final price *could* be low when they see a high discount rate. We could regard this impression as persistent errors in our user study. The user study consists of 3 settings. In Setting 1, we display all 5 attributes for each Airbnb, and the utility function is assumed to be linear w.r.t. all 5 attributes. In Setting 2, we still display all 5 attributes for each option, but the utility is assumed to be linear w.r.t. the final price (instead of the labeled price and the discount) and other attributes. Then, in Setting 3, only the final price and other attributes are displayed, and the utility is assumed to be linear w.r.t. the final price and other attributes. In each setting, the user interacts with algorithm *SS-score* until a list of at most  $l$  points are returned, and is then asked to select one option from the list as the *selected favorite rental option* of this setting. We refer to the selected favorite

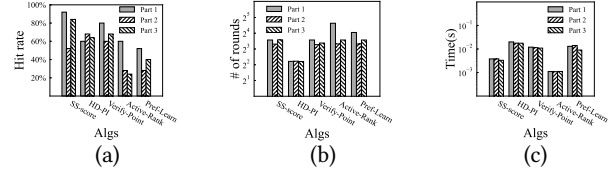


Figure 20: Results on user study

rental option of Setting 1, 2 and 3 as  $p_1$ ,  $p_2$  and  $p_3$ , respectively. After obtaining  $p_1$ ,  $p_2$  and  $p_3$ , the user is required to choose one options among them as the best point. Notably, if  $p_3$  is chosen as the best point, it indicates that the user considers the final price, rather than the labeled price and discount, when making their decision. For users who selected  $p_3$  as the best point, an additional question is asked: The user is presented with  $p_2$ , including its labeled price, discount, and computed final price, and is asked if they think the final price is higher than their expectation. Note that when the user selected  $p_2$  in Setting 2, only the labeled price and discount of  $p_2$  are presented to the user. Given that the user considers final price when making the decision, if the computed final price of  $p_2$  is higher than the user’s expectation, it suggests that the user is prone to persistent errors, possibly due to the false impression that the final price could be low because of the high discount rate. Among the 21 users who selected  $p_3$  as the best point, 33% (7) of them found that the final price of  $p_2$  was higher than expected, indicating that they have a tendency to make persistent errors. Moreover, 24% (5) users found that their best point was not in the list of options recommended by Setting 2, indicating that 24% users missed their best point due to persistent errors.

## 6.6 Summary

The experiments demonstrated the superiority of our proposed algorithms, namely *FC*, *SS-score* and *SS-random*, over existing approaches. (1) We are efficient and effective. We achieve nearly 100% accuracy in most of the experiments using a small number of rounds, which consistently outperforms existing algorithms. (2) We are scalable to the input size and dimensionality. For example, on 7-d dataset *HTRU*, *SS-score* and *SS-random* finish with around 20 questions and achieve over 98% accuracy, but algorithms *UtilApprox* and *Active-Ranking* obtain lower accuracies with even more rounds. (3) We are capable of handling many persistent errors. Even when the error rate is high (e.g., 0.15), all our algorithms still achieve more than 75% accuracy, which is at least 10% higher than other existing approaches.

## 7 CONCLUSION

In this paper, we propose interactive algorithms that return the user’s best point with high confidence even when the user makes persistent errors and random errors during the interaction, which makes them more robust than existing approaches. Specifically, we proposed algorithm *FC*, requiring an asymptotically optimal number of rounds, and algorithm *SS*, requiring a small number of questions empirically. Both algorithms return the best point with a provable guarantee. We conducted extensive experiments to demonstrate that our algorithms are efficient and effective when addressing both types of errors. In the future, we study how to extend our solutions to return top- $k$  points of the user.



## REFERENCES

- [1] Pankaj K Agarwal, Nirman Kumar, Stavros Sintos, and Subhash Suri. 2017. Efficient algorithms for k-regret minimizing sets. *arXiv preprint arXiv:1702.01446* (2017).
- [2] Yongkil Ahn. 2019. The economic cost of a fat finger mistake: a comparative case study from Samsung Securities’s ghost stock blunder. *Journal of Operational Risk* 16, 2 (2019).
- [3] David Avis and Komei Fukuda. 1991. A pivoting algorithm for convex hulls and vertex enumeration of arrangements and polyhedra. In *Proceedings of the seventh annual symposium on Computational geometry*. 98–104.
- [4] Ilaria Bartolini, Paolo Ciaccia, and Marco Patella. 2014. Domination in the probabilistic world: Computing skylines for arbitrary correlations and ranking semantics. *ACM Transactions on Database Systems (TODS)* 39, 2 (2014), 1–45.
- [5] Ilaria Bartolini, Paolo Ciaccia, and Florian Waas. 2001. FeedbackBypass: A new approach to interactive similarity query processing. In *VLDB*. 201–210.
- [6] Stephan Borzsony, Donald Kossmann, and Konrad Stocker. 2001. The skyline operator. In *Proceedings 17th international conference on data engineering*. IEEE, 421–430.
- [7] Apostolos Chalkis and Vissarion Fisikopoulos. 2020. volesti: Volume approximation and sampling for convex polytopes in  $\mathbb{R}^d$ . *arXiv preprint arXiv:2007.01578* (2020).
- [8] Qixu Chen and Raymond Chi-Wing Wong. 2023. Finding Best Tuple via Error-prone User Interaction. In *Proceedings of the 39th IEEE International Conference on Data Engineering*.
- [9] Sean Chester, Alex Thomo, S Venkatesh, and Sue Whitesides. 2014. Computing k-regret minimizing sets. *Proceedings of the VLDB Endowment* 7, 5 (2014), 389–400.
- [10] Paolo Ciaccia and Davide Martini. 2017. Reconciling skyline and ranking queries. *Proceedings of the VLDB Endowment* 10, 11 (2017), 1454–1465.
- [11] Airbnb dataset. 2023. <http://insideairbnb.com/get-the-data/>.
- [12] Eyal Dushkin and Tova Milo. 2018. Top-k sorting under partial order information. In *Proceedings of the 2018 International Conference on Management of Data*. 1007–1019.
- [13] Brian Eriksson. 2013. Learning to top-k search using pairwise comparisons. In *Artificial Intelligence and Statistics*. PMLR, 265–273.
- [14] Moein Falahatgar, Yi Hao, Alon Orlitsky, Venkatesh Rajkumar, and Vaishakh Ravindrakumar. 2017. Maxing and ranking with few assumptions. *Advances in Neural Information Processing Systems* 30 (2017).
- [15] Moein Falahatgar, Ayush Jain, Alon Orlitsky, Venkatesh Rajkumar, and Vaishakh Ravindrakumar. 2018. The limits of maxing, ranking, and preference learning. In *International conference on machine learning*. PMLR, 1427–1436.
- [16] Moein Falahatgar, Alon Orlitsky, Venkatesh Rajkumar, and Ananda Theertha Suresh. 2017. Maximum selection and ranking under noisy comparisons. In *International Conference on Machine Learning*. PMLR, 1088–1096.
- [17] Barbara Geissmann, Stefano Leucci, Chih-Hung Liu, and Paolo Penna. 2018. Optimal sorting with persistent comparison errors. *arXiv preprint arXiv:1804.07575* (2018).
- [18] Reinhard Heckel, Nihar B Shah, Kannan Ramchandran, and Martin J Wainwright. 2019. Active ranking from pairwise comparisons and when parametric assumptions do not help. *The Annals of Statistics* 47, 6 (2019), 3099–3126.
- [19] Reinhard Heckel, Max Simchowitz, Kannan Ramchandran, and Martin Wainwright. 2018. Approximate ranking from pairwise comparisons. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 1057–1066.
- [20] Kevin G Jamieson and Robert Nowak. 2011. Active ranking using pairwise comparisons. *Advances in neural information processing systems* 24 (2011).
- [21] Yiling Jia, Huazheng Wang, Stephen Guo, and Hongning Wang. 2021. Pairrank: Online pairwise learning to rank by divide-and-conquer. In *Proceedings of the Web Conference 2021*. 146–157.
- [22] Sumeet Katariya, Lalit Jain, Nandana Sengupta, James Evans, and Robert Nowak. 2018. Adaptive sampling for coarse ranking. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 1839–1848.
- [23] Barry Kirwan. 2017. *A guide to practical human reliability assessment*. CRC press.
- [24] Rolf Klein, Rainer Penninger, Christian Sohler, and David P Woodruff. 2011. Tolerant algorithms. In *Algorithms—ESA 2011: 19th Annual European Symposium, Saarbrücken, Germany, September 5–9, 2011. Proceedings 19*. Springer, 736–747.
- [25] Jongwuk Lee, Gae-won You, and Seung-won Hwang. 2009. Personalized top-k skyline queries in high-dimensional space. *Information Systems* 34, 1 (2009), 45–61.
- [26] Alchemer LLC. 2023. <https://www.alchemer.com/resources/blog/how-many-survey-questions/>.
- [27] De Berg Mark, Cheong Otfried, van Kreveld Marc, and Overmars Mark. 2008. *Computational geometry algorithms and applications*. Springer.
- [28] Denis Mindolin and Jan Chomicki. 2009. Discovering relative importance of skyline attributes. *Proceedings of the VLDB Endowment* 2, 1 (2009), 610–621.
- [29] Kyriakos Mouratidis, Keming Li, and Bo Tang. 2021. Marrying top-k with skyline queries: Relaxing the preference input while producing output of controllable size. In *Proceedings of the 2021 International Conference on Management of Data*. 1317–1330.
- [30] Kyriakos Mouratidis and Bo Tang. 2018. Exact processing of uncertain top-k queries in multi-criteria settings. *Proceedings of the VLDB Endowment* 11, 8 (2018), 866–879.
- [31] Danupon Nanongkai, Ashwin Lall, Atish Das Sarma, and Kazuhisa Makino. 2012. Interactive regret minimization. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*. 109–120.
- [32] Danupon Nanongkai, Atish Das Sarma, Ashwin Lall, Richard J Lipton, and Jun Xu. 2010. Regret-minimizing representative databases. *Proceedings of the VLDB Endowment* 3, 1-2 (2010), 1114–1124.
- [33] Peng Peng and Raymond Chi-Wing Wong. 2015. k-hit query: Top-k query with probabilistic utility function. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*. 577–592.
- [34] Li Qian, Jinyang Gao, and HV Jagadish. 2015. Learning user preferences by adaptive pairwise comparison. *Proceedings of the VLDB Endowment* 8, 11 (2015), 1322–1333.
- [35] QuestionPro. 2023. <https://www.questionpro.com/blog/optimal-number-of-survey-questions/>.
- [36] Wenbo Ren, Jia Kevin Liu, and Ness Shroff. 2019. On sample complexity upper and lower bounds for exact ranking from noisy comparisons. *Advances in Neural Information Processing Systems* 32 (2019).
- [37] Gerard Salton. 1989. Automatic text processing: The transformation, analysis, and retrieval of. *Reading: Addison-Wesley* 169 (1989).
- [38] Thomas Seidl and Hans-Peter Kriegel. 1997. Efficient user-adaptable similarity search in large multimedia databases. In *VLDB*, Vol. 97. 506–515.
- [39] Zhixuan Song and Nick Roussopoulos. 2001. K-nearest neighbor search for moving query point. In *International Symposium on Spatial and Temporal Databases*. Springer, 79–96.
- [40] Csaba D Toth, Joseph O’Rourke, and Jacob E Goodman. 2017. *Handbook of discrete and computational geometry*. CRC press.
- [41] Weicheng Wang and Raymond Chi-Wing Wong. 2022. Interactive mining with ordered and unordered attributes. *Proceedings of the VLDB Endowment* 15, 11 (2022), 2504–2516.
- [42] Weicheng Wang, Raymond Chi-Wing Wong, and Min Xie. 2021. Interactive Search for One of the Top-k. In *Proceedings of the 2021 International Conference on Management of Data*. 1920–1932.
- [43] A Student with Top-tier Score Admitted by mediocre University (Chinese version only). 2020. [https://news.southcn.com/node\\_6854f1135c/4357641930.shtml](https://news.southcn.com/node_6854f1135c/4357641930.shtml).
- [44] Min Xie, Raymond Chi-Wing Wong, and Ashwin Lall. 2019. Strongly truthful interactive regret minimization. In *Proceedings of the 2019 International Conference on Management of Data*. 281–298.
- [45] Jiping Zheng and Chen Chen. 2020. Sorting-based interactive regret minimization. (2020), 473–490.
- [46] Yi Zong and Xiaojie Guo. 2022. An experimental study on anchoring effect of consumers’ price judgment based on consumers’ experiencing scenes. *Frontiers in Psychology* 13 (2022), 794135.



Dataset	size	dimensionality	skyline size
4d100	100	4	87
4d1k	1000	4	403
4d10k	10000	4	1484
4d100k	100000	4	3780
4d1m	1000000	4	6699
4d10m	10 million	4	12437
4d100m	100 million	4	19445
2d100k	100000	2	40
3d100k	100000	3	538
5d100k	100000	5	11352
AirQuality	420478	4	3948
Weather	96483	6	2110
HTRU	19898	7	11720

Table 3: Dataset Statistics

This appendix provides supplementary materials that are not included in the main paper due to space constraints. Section A presents the additional experiments. Section B provides detailed proofs of lemmas and theorems in the paper.

## A ADDITIONAL EXPERIMENTS

In this section, we present the experiments that are not included in the main body due to space constraints. Section A.1, A.2, A.3 and A.4 show the supplementary results on parameter settings, additional experiments on synthetic datasets, additional experiments on user study and experiments on question type flexibility, respectively. Statistics of synthetic and real datasets used in our experiments can be found in Table 3. For each synthetic dataset, we name it using its dimensionality and size. For example, “4d100” means the synthetic dataset with dimensionality 4 ( $d = 4$ ) and size 100 ( $N = 100$ ).

### A.1 Additional Experiments on Parameter Setting

**A.1.1 Experiments Showing the Advantages of Preprocessed Datasets Containing Skyline Points Only on Existing Algorithms.** To show that preprocessing the dataset to contain only the skyline points does not make the comparison among algorithms unfair, we also record the performance of our competitors which do not need preprocessing on the raw 4-d synthetic dataset with 100,000 points. In particular, since *HD-PI*, *Verify-Point* and *UH-Simplex* all require this preprocessing step [8, 42, 44], we run the remaining competitors, namely *UtilApprox*, *Active-ranking* and *Pref-Learn* on the raw input. Table 4 presents a comparison of their performance on both the processed and raw datasets, with “proc” indicating results on the processed dataset (i.e., the dataset with skyline points only) and “raw” indicating results on the raw dataset. We observe that when running on the raw input, all three algorithms tend to ask more questions and obtain lower accuracies. This is expected since with the raw input they also need to consider a large number of non-skyline points. Particularly, the processing time of *UtilApprox* increases by two orders of magnitude, and *Active-ranking* asks 180 times more questions due to its objective of ranking all points before reporting the best tuple.

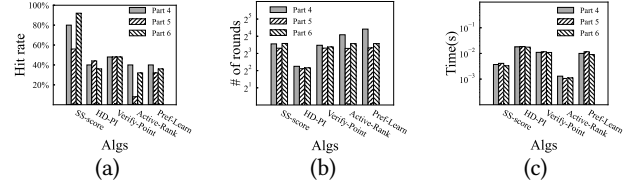


Figure 21: Results on additional user study

## A.2 Additional Experiments on Synthetic Datasets

**A.2.1 Experiments Showing P95, Median and Accumulated Processing Times.** In addition to the average processing time, we also report other time statistics in Table 5, including the processing time at the 95th percentile (P95), median processing time, and accumulated processing time of all rounds. These statistics were collected from the 4-dimensional synthetic dataset with a size of 100,000. It is worth noting that at least 95% of questions of our algorithms can be determined in 3 seconds, which proves that our proposed algorithms run at an interactive speed.

**A.2.2 Experiments on Time Required for the Pre-processing Phase.** We studied the time required for the pre-processing phase of our algorithms on synthetic datasets and real datasets. The results are summarized in Table 6. All synthetic datasets are named using the convention stated at the beginning of Section A. On all tested datasets, including the synthetic dataset with the highest dimensionality (i.e., 5d100k), the synthetic dataset with the largest size (i.e., 4d100m), and the 7-d real dataset *HTRU*, the pre-processing phase requires less than 2 minutes. We conclude that the pre-processing phase is efficient given that this phase only needs to be run once and its results can be stored and reused for future runs.

**A.2.3 Experiments on the Effect of Different User Preferences.** Since different users have distinct utility vectors, we also studied the algorithms’ performance when utility vectors are sampled from different regions in the utility space  $U$ . To sample  $u$  from different regions, in each setting, we select a non-empty set  $I \subseteq \{1, 2, \dots, d\}$  where  $d$  is the dimensionality, randomly set  $u[i]$  in the range  $[0, 0.5]$  if  $i \in I$ , and randomly set  $u[i]$  in the range  $[0.5, 1]$  if  $i \notin I$ . The utility vectors are then normalized such that  $\sum_{1 \leq i \leq d} u[i] = 1$ . The performance of our algorithms on 4-d synthetic datasets with size 100,000 are reported in Table 7, 8 and 9. For completeness, the performance of existing algorithms is also reported in Table 10, 11 and 12. In these tables, each row corresponds to a different choice of  $I$ , and the average of all choices of  $I$  is reported as well. We observe that our algorithms consistently outperform existing algorithms on almost every setting of  $I$  in terms of accuracy. Furthermore, our algorithms exhibit stability with minimal performance variations across different  $I$  values. In contrast, some existing algorithms (e.g., *Active-Ranking*) show sensitivity to changes in  $I$ , resulting in accuracy fluctuations exceeding 20%.

## A.3 Additional Experiments on User Study

**A.3.1 Experiments on capability of handling persistent errors.** To further evaluate our algorithm’s capability of handling persistent errors, we conducted an additional user study comprising three

	UtilApprox (proc)	UtilApprox (raw)	Pref-Learn (proc)	Pref-Learn (raw)	Active-ranking (proc)	Active-ranking (raw)
Accuracy	88%	87%	85%	84%	72%	70%
# of rounds	30.08	30.35	41.57	43.13	58.32	10671.50
Time (s)	$5.22 \times 10^{-5}$	$2.01 \times 10^{-3}$	$7.78 \times 10^{-1}$	1.26	$1.90 \times 10^{-2}$	$1.12 \times 10^{-3}$

**Table 4: Effect of skyline preprocessing**

Algorithm	average (s)	P95 (s)	median (s)	accumulated (s)
FC	$4.10 \times 10^{-2}$	$2.02 \times 10^{-1}$	$1.10 \times 10^{-2}$	1.25
SS-score	$6.84 \times 10^{-1}$	2.51	$2.70 \times 10^{-1}$	$1.08 \times 10^1$
SS-random	$7.20 \times 10^{-2}$	$2.02 \times 10^{-1}$	$1.88 \times 10^{-2}$	2.69
HD-PI	$8.85 \times 10^{-1}$	1.46	$2.15 \times 10^{-1}$	7.26
Verify-Point	$6.42 \times 10^{-1}$	$7.05 \times 10^{-1}$	$1.12 \times 10^{-1}$	$1.17 \times 10^1$
UtilApprox	$5.22 \times 10^{-5}$	$8.34 \times 10^{-5}$	$5.39 \times 10^{-5}$	$1.57 \times 10^{-3}$
Active-Ranking	$1.90 \times 10^{-2}$	$7.93 \times 10^{-2}$	$6.78 \times 10^{-3}$	1.11
Pref-Learn	$7.78 \times 10^{-1}$	$1.41 \times 10^{-1}$	$5.74 \times 10^{-2}$	$3.23 \times 10^1$
UH-Simplex	$2.66 \times 10^{-1}$	3.13	$6.61 \times 10^{-2}$	2.77

**Table 5: Time statistics**

Dataset	Time (s)
4d100	3.58
4d1k	4.83
4d10k	6.75
4d100k	9.24
4d1m	10.64
4d10m	56.75
4d100m	94.91
2d100k	0.18
3d100k	0.90
5d100k	46.73
AirQuality	5.35
Weather	35.02
HTRU	85.19

**Table 6: Time for pre-processing phase**

	FC	SS-score	SS-random
{1}	93%	98%	96%
{2}	92%	99%	96%
{3}	96%	100%	97%
{4}	93%	97%	100%
{1,2}	93%	96%	99%
{1,3}	96%	98%	99%
{1,4}	97%	99%	99%
{2,3}	97%	100%	100%
{2,4}	98%	99%	98%
{3,4}	97%	100%	97%
{1,2,3}	93%	99%	96%
{1,2,4}	98%	100%	97%
{1,3,4}	96%	99%	98%
{2,3,4}	96%	100%	95%
average	95%	99%	98%

**Table 7: Effect of different  $u$  on our algorithms (Accuracy)**

parts, called Part 4, Part 5 and Part 6, that follows the setting of Part 1, Part 2 and Part 3 in our first user study, respectively.

However, in these parts, we want to introduce a *known* persistent error. Therefore, we make one change in the first round for each algorithm by explicitly requiring the user to select the *less preferred* one, and the remaining rounds are kept the same by asking the user to select the *preferred* one. This (real) less preferred one in the first round is recorded as the more preferred one in the system/algorithm, which is considered as a *known* persistent error “artificially” introduced by these 3 parts. If in later rounds of the algorithm, the same question was selected again to ask this user, the recorded answer would be used automatically. The reason for this design is to introduce some persistent errors in each algorithm.

Figure 21 displays the results for Part 4, 5 and 6. Compared to their counterparts (shown in Figure 20), the hit rates of all algorithms in Part 4 (resp. Part 6) are generally lower than that in Part 1 (resp. Part 3) due to the introduced persistent error. For example,

*HDPI* and *Verify-Point* experiences a hit rate drop of over 20% and 30% in Part 4, respectively. However, the hit rate of *SS-score* remains relatively stable in both Part 4 and 6, demonstrating its ability to handle persistent errors. Notably, even in Part 5, where algorithms are forced to stop after round 10, *SS-score* still achieves a higher hit rate compared to *HDPI* and *Verify-Point*, even though the latter two algorithms have better round efficiency.

#### A.4 Experiments on Question Type Flexibility

To test the flexibility of our algorithms on different types of questions, in this section, instead of asking questions in the form of pairwise comparison, we adapted our algorithms *FC*, *SS-score* and *SS-random* to ask two new types of questions, namely in each round (a) displaying  $s > 2$  points and asking the user to select the favorite one, and (b) displaying  $s > 2$  points and asking the user to divide

	FC	SS-score	SS-random
{1}	24.21	16.38	18.28
{2}	29.87	16.43	18.01
{3}	31.19	16.21	18.11
{4}	29.13	16.28	18.24
{1,2}	26.49	16.34	18.13
{1,3}	27.14	16.74	18.79
{1,4}	33.76	16.79	18.42
{2,3}	31.88	16.14	18.53
{2,4}	24.70	16.18	19.22
{3,4}	27.93	16.07	18.72
{1,2,3}	32.76	16.45	18.18
{1,2,4}	33.11	16.35	19.03
{1,3,4}	32.16	16.64	18.77
{2,3,4}	34.46	16.13	19.23
average	29.91	16.37	18.55

**Table 8: Effect of different  $u$  on our algorithms (# of rounds)**

	FC	SS-score	SS-random
{1}	$5.15 \times 10^{-2}$	$8.05 \times 10^{-1}$	$7.91 \times 10^{-2}$
{2}	$4.11 \times 10^{-2}$	$6.77 \times 10^{-1}$	$8.20 \times 10^{-2}$
{3}	$3.84 \times 10^{-2}$	$7.84 \times 10^{-1}$	$7.94 \times 10^{-2}$
{4}	$4.10 \times 10^{-2}$	$7.88 \times 10^{-1}$	$8.18 \times 10^{-2}$
{1, 2}	$4.49 \times 10^{-2}$	$7.59 \times 10^{-1}$	$7.97 \times 10^{-2}$
{1, 3}	$4.37 \times 10^{-2}$	$7.41 \times 10^{-1}$	$7.76 \times 10^{-2}$
{1, 4}	$3.68 \times 10^{-2}$	$6.86 \times 10^{-1}$	$8.04 \times 10^{-2}$
{2, 3}	$3.88 \times 10^{-2}$	$8.20 \times 10^{-1}$	$7.88 \times 10^{-2}$
{2, 4}	$4.76 \times 10^{-2}$	$7.59 \times 10^{-1}$	$7.86 \times 10^{-2}$
{3, 4}	$4.19 \times 10^{-2}$	$8.00 \times 10^{-1}$	$8.11 \times 10^{-2}$
{1, 2, 3}	$3.72 \times 10^{-2}$	$7.54 \times 10^{-1}$	$7.80 \times 10^{-2}$
{1, 2, 4}	$3.70 \times 10^{-2}$	$7.18 \times 10^{-1}$	$7.69 \times 10^{-2}$
{1, 3, 4}	$3.78 \times 10^{-2}$	$7.35 \times 10^{-1}$	$7.83 \times 10^{-2}$
{2, 3, 4}	$3.62 \times 10^{-2}$	$6.78 \times 10^{-1}$	$7.65 \times 10^{-2}$
average	$4.10 \times 10^{-2}$	$7.50 \times 10^{-1}$	$7.92 \times 10^{-2}$

**Table 9: Effect of different  $u$  on our algorithms (Time (s))**

them into two groups called the *superior group* and the *inferior group*. We describe the results of these variants in Section A.4.1 and A.4.2, respectively.

Since the algorithms need to display  $s$  (instead of 2) points in each round, we adapted the random-based selection and the score-based selection as follows: (1) for the random-based selection, instead of returning the first pair of points whose preference is not known, when we find a pair whose preference is unknown, we record this pair and continue the selection process until the number of distinct points in the recorded pairs is at least  $s$ , then these points are displayed to the user (in case the number of distinct points in the recorded pairs is  $s + 1$ , we simply display  $s$  of them); (2) for the score-based selection, instead of returning the pair with the highest score whose preference is unknown, we scan through all pairs with unknown preferences and find  $\frac{(s-1)(s-2)}{2} + 1$  pairs with the highest score, sort them in descending order of scores and select distinct

points in these pairs until  $s$  points are selected, which are displayed to the user (the reason why  $\frac{(s-1)(s-2)}{2} + 1$  pairs are needed is to ensure that there are at least  $s$  distinct points in these pairs).

As a special treatment, since Stage 1 of *FC* has its own point selection strategy, we adapted it by first selecting two points using its own question selection strategy for Stage 1 described in Section 5.2, then selecting the remaining  $s - 2$  points using the adapted version of the score-based selection.

**A.4.1 Selecting favorite point among  $s$  displayed points.** In this set of experiments, we adapted the algorithms to display more than two points in each round. Specifically,  $s > 2$  points are displayed and the user is asked to select their favorite one among them. We simulate how the user decides which point to select among the  $s$  displayed points with a persistent error rate  $\theta$  as follows: For each pair of displayed points, we compare them with a persistent error rate  $\theta$ . The point that wins in the most number of comparisons is selected. If multiple points share the highest number of wins, we arbitrarily select one of them. For each pair of selected/non-selected points, we can use it to construct a halfspace using the technique in Section 4.1, and update the related data structures.

We summarize the results in Figure 22. From Figure 22 (a), when  $s$  is increased from 2 to 10, the accuracies of our algorithms slightly increases. This is because the selected point is likely to be ranked highly in the displayed points, thus when  $s$  grows, the percentage of incorrect halfspaces becomes smaller. In other words, the “effective” error rate is smaller than  $\theta$ . When  $s$  increases, the number of questions required by *SS-score* and *FC* decreases since more hyperplanes can be acquired in each round, making the confidence regions shrink faster. Besides, the processing time of these two algorithms increases along with  $s$  because of the increasing amount of updates on data structures in each round. Notably, for *SS-random*, its number of questions first decreases and then increases when  $s$  grows. One possible explanation is that, since *SS-random* adopts the random-based selection, when  $s$  grows, the hyperplane generated by each pair of the selected point (typically the point with the highest utility) and the non-selected point is less likely to evenly divide confidence regions, making the sizes of confidence regions shrink slower. The average processing time of *SS-random* first increases and then decreases along with  $s$ , which is explainable given the growing trend of the number of rounds required by *SS-random* (i.e., first decrease then increase), and the fact that the processing time of later rounds is typically faster than previous rounds.

**A.4.2 Dividing  $s$  points into two groups.** We adapted our algorithms to display  $s > 2$  points in each round and ask the user to divide them into two groups called the *superior group* and the *inferior group*, where the superior (resp. inferior) group contains points with higher (resp. lower) utility scores. To simulate the dividing process with a persistent error rate  $\theta$ , we first make pairwise comparisons on all pairs of  $s$  displayed points, each with a persistent error rate  $\theta$ , and sort the points in decrease order of the number of comparisons in which they win. Then, we uniform-randomly pick a number  $i$  in range  $[1, s - 1]$ , put the first  $i$  points into the superior group, and the rest  $s - i$  points into the inferior group. Note that we assume each group has at least one point, otherwise, all points will belong to the same group and there is little information from this round.

	HD-PI	Verify-Point	UtilApprox	Active-Ranking	Pref-Learn	UH-Simplex
{1}	79%	76%	92%	76%	88%	66%
{2}	84%	85%	91%	78%	92%	75%
{3}	80%	80%	84%	63%	87%	72%
{4}	73%	77%	87%	75%	91%	66%
{1, 2}	83%	77%	92%	76%	83%	85%
{1, 3}	80%	87%	93%	77%	92%	60%
{1, 4}	84%	84%	89%	53%	75%	57%
{2, 3}	72%	79%	93%	72%	82%	81%
{2, 4}	77%	79%	90%	78%	86%	55%
{3, 4}	73%	77%	89%	73%	90%	64%
{1, 2, 3}	75%	76%	87%	66%	82%	57%
{1, 2, 3}	81%	80%	92%	62%	83%	59%
{1, 2, 3}	76%	85%	88%	61%	87%	52%
{1, 2, 3}	79%	84%	85%	79%	89%	61%
average	78%	80%	89%	71%	86%	65%

Table 10: Effect of different  $u$  on existing algorithms (Accuracy)

	HD-PI	Verify-Point	UtilApprox	Active-Ranking	Pref-Learn	UH-Simplex
{1}	7.93	9.65	29.09	60.92	44.16	12.15
{2}	8.10	9.49	29.78	60.94	42.91	10.34
{3}	8.43	9.27	29.17	97.21	41.80	10.62
{4}	8.35	9.57	27.85	63.14	44.63	12.26
{1, 2}	8.11	9.48	29.53	59.48	43.55	8.69
{1, 3}	8.24	9.29	29.39	56.00	40.96	11.89
{1, 4}	8.18	9.45	30.56	59.84	44.45	12.46
{2, 3}	8.02	9.64	30.67	60.01	42.47	10.25
{2, 4}	7.99	9.13	28.76	58.33	42.93	12.29
{3, 4}	8.46	9.40	29.75	56.91	41.25	12.44
{1, 2, 3}	8.17	9.66	33.14	57.93	42.80	10.78
{1, 2, 4}	8.09	9.21	31.15	56.40	41.48	10.23
{1, 3, 4}	8.10	9.36	30.34	55.26	42.42	11.58
{2, 3, 4}	8.28	9.42	29.40	60.12	41.53	13.30
average	8.18	9.43	29.90	61.61	42.67	11.38

Table 11: Effect of different  $u$  on existing algorithms (# of rounds)

For each distinct pair of superior/inferior points, we can create a halfspace using the technique described in Section 4.1.

Figure 23 shows the algorithm performance when  $s$  varies from 2 to 10. When  $s$  grows, since more halfspaces can be obtained in each round, the number of rounds required by all algorithms tends to decrease. Meanwhile, obtaining more halfspaces means more data structures need to be updated in each round, and thus the processing time increases. Notably, the accuracies of all algorithms increases along with  $s$ , which is reasonable since with more points to be considered in each round, our simulation makes it harder to obtain wrong halfspaces.

## B RELATED PROOFS

This section details the proofs of lemmas and theorems presented in the main body of the paper.

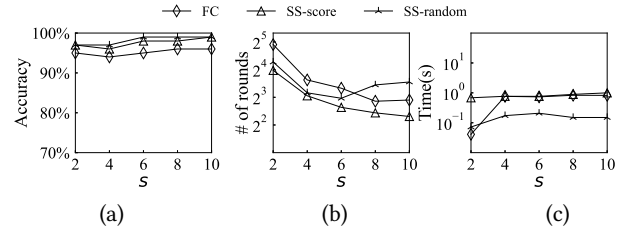


Figure 22: Varying question type: selecting favorite among  $s$  displayed points

### B.1 Proof of Lemma 1

When the results of all pairwise comparisons are acquired, denote the set of points whose partitions overlap with  $R^I$  by  $RS$ . Since each pair of points in  $RS$  has been compared, there are in total

	HD-PI	Verify-Point	UtilApprox	Active-Ranking	Pref-Learn	UH-Simplex
{1}	$3.82 \times 10^{-1}$	$6.24 \times 10^{-1}$	$5.07 \times 10^{-5}$	$1.82 \times 10^{-2}$	$7.49 \times 10^{-1}$	$2.27 \times 10^{-1}$
{2}	$3.04 \times 10^{-1}$	$7.09 \times 10^{-1}$	$5.38 \times 10^{-5}$	$1.91 \times 10^{-2}$	$7.54 \times 10^{-1}$	$2.63 \times 10^{-1}$
{3}	$2.85 \times 10^{-1}$	$6.49 \times 10^{-1}$	$5.05 \times 10^{-5}$	$1.18 \times 10^{-2}$	$7.65 \times 10^{-1}$	$4.19 \times 10^{-1}$
{4}	$3.03 \times 10^{-1}$	$6.26 \times 10^{-1}$	$5.12 \times 10^{-5}$	$1.84 \times 10^{-2}$	$7.39 \times 10^{-1}$	$4.33 \times 10^{-1}$
{1, 2}	$3.32 \times 10^{-1}$	$6.34 \times 10^{-1}$	$5.13 \times 10^{-5}$	$1.94 \times 10^{-2}$	$7.3 \times 10^{-1}$	$1.94 \times 10^{-1}$
{1, 3}	$3.24 \times 10^{-1}$	$6.48 \times 10^{-1}$	$5.12 \times 10^{-5}$	$2.06 \times 10^{-2}$	$7.63 \times 10^{-1}$	$1.76 \times 10^{-1}$
{1, 4}	$2.73 \times 10^{-1}$	$6.14 \times 10^{-1}$	$5.26 \times 10^{-5}$	$1.92 \times 10^{-2}$	$7.06 \times 10^{-1}$	$3.77 \times 10^{-1}$
{2, 3}	$2.87 \times 10^{-1}$	$6.21 \times 10^{-1}$	$5.21 \times 10^{-5}$	$1.94 \times 10^{-2}$	$7.42 \times 10^{-1}$	$2.7 \times 10^{-1}$
{2, 4}	$3.52 \times 10^{-1}$	$6.47 \times 10^{-1}$	$5.14 \times 10^{-5}$	$1.98 \times 10^{-2}$	$7.25 \times 10^{-1}$	$3.98 \times 10^{-1}$
{3, 4}	$3.11 \times 10^{-1}$	$6.37 \times 10^{-1}$	$5.23 \times 10^{-5}$	$1.97 \times 10^{-2}$	$7.44 \times 10^{-1}$	$4.12 \times 10^{-1}$
{1, 2, 3}	$2.80 \times 10^{-1}$	$6.29 \times 10^{-1}$	$5.18 \times 10^{-5}$	$2.10 \times 10^{-2}$	$7.38 \times 10^{-1}$	$1.79 \times 10^{-1}$
{1, 2, 4}	$2.74 \times 10^{-1}$	$6.42 \times 10^{-1}$	$5.96 \times 10^{-5}$	$2.02 \times 10^{-2}$	$7.36 \times 10^{-1}$	$3.48 \times 10^{-1}$
{1, 3, 4}	$2.80 \times 10^{-1}$	$6.42 \times 10^{-1}$	$5.20 \times 10^{-5}$	$2.05 \times 10^{-2}$	$7.44 \times 10^{-1}$	$3.85 \times 10^{-1}$
{2, 3, 4}	$2.69 \times 10^{-1}$	$6.25 \times 10^{-1}$	$5.03 \times 10^{-5}$	$1.90 \times 10^{-2}$	$7.86 \times 10^{-1}$	$3.86 \times 10^{-1}$
average	$3.04 \times 10^{-1}$	$6.39 \times 10^{-1}$	$5.22 \times 10^{-5}$	$1.90 \times 10^{-2}$	$7.44 \times 10^{-1}$	$3.19 \times 10^{-1}$

Table 12: Effect of different  $u$  on existing algorithms (Time (s))

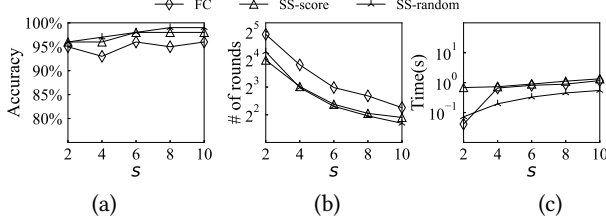


Figure 23: Varying question type: dividing  $s$  points into two groups

$|RS|(|RS| - 1)/2$  comparisons. Each comparison yields one loser, so there are  $|RS|(|RS| - 1)/2$  losers. Further, each point in  $RS$  cannot lose in more than  $i$  comparisons, otherwise, it is not in  $R^i$  since its partition is not supported by at least  $i$  halfspaces. Therefore, we have  $|RS|(|RS| - 1)/2 \leq i|RS|$ , so  $|RS| \leq 2i + 1$ .

## B.2 Proof of Lemma 2

In this section, we prove the recurrence relation in Lemma 2 by induction. Firstly, notice that when  $t = 0$ , confidence regions  $R_0^0, R_0^1, \dots, R_0^k$  all corresponds to the entire utility space, which is trivially correct. Assume that at round  $t'$ , all confidence regions are correct. Then, at round  $t' + 1$ , when a new halfspace  $s$  arrives,  $R_{t'+1}^0 = R_{t'}^0 \cap s$ , since  $R^0$  is the region supported by all halfspaces.

As for  $R_{t'+1}^i$  where  $i > 0$ , firstly, for any points in  $R_{t'+1}^i \cap s$ , it is supported by  $t' - i$  halfspaces in the first  $t'$  halfspaces and also supported by  $s$ . Therefore, it is supported by  $t' + 1 - i$  halfspaces and should lie in  $R_{t'+1}^i$ . With a similar argument, we can also prove that any points in  $R_{t'+1}^{i-1} \cap s^-$  should also lie in  $R_{t'+1}^i$ . Expect for the above cases, no other points lie in  $R_{t'+1}^i$ . Therefore, all confidence regions are correct at round  $t' + 1$ , which proves the correctness of the recurrence.

## B.3 Proof of Lemma 3

We now show the correctness of the two relations stated in Lemma 3. We first prove that  $R_{t+1}^i \subseteq R_t^i$ . Firstly, note that by the definition of confidence region,  $R_t^{i-1} \subseteq R_t^i$ . By Lemma 2,  $R_{t+1}^i = (R_t^i \cap s) \cup (R_t^{i-1} \cap s^-)$ . Since  $R_t^i \cap s \subseteq R_t^i$  and  $R_t^{i-1} \cap s^- \subseteq R_t^{i-1} \subseteq R_t^i$ , their union  $R_{t+1}^i \subseteq R_t^i$  as desired.

Next, we prove that  $R_t^i \subseteq R_{t+1}^{i+1}$ . By the definition of confidence region, for any point  $v \in R_t^i$ , there is a set  $S'$  of at least  $t - i$  halfspaces indicated by the user such that  $v \in \bigcap_{s \in S'} s$ . Notice that  $\bigcap_{s \in S'} s \subseteq R_{t+1}^{i+1}$  by the definition of  $R_{t+1}^{i+1}$ . Therefore, we have that for any point  $v \in R_t^i$ ,  $v \in R_{t+1}^{i+1}$ , thus,  $R_t^i \subseteq R_{t+1}^{i+1}$ .

## B.4 Proof of Theorem 1

In this section, We show the proof of Theorem 1. We first show that there exists a dataset  $D$  with the following property:

**Property 1.** For any  $p_i \in D$  and its partition  $P_i$ , for any non-empty set of halfspaces  $HS \subseteq \{h_{ab}^* | a, b \neq i, * \in \{+, -\}\}$ , if the intersection of halfspaces in  $HS$  forms a non-empty area, then  $P_i$  also overlaps with this area.

Consider a  $n$ -dimensional dataset  $D_n = \{p_1, p_2, \dots, p_n\}$  such that for each  $p_i, i \in [1, n]$ , its  $i$ -th dimensional value is 1 and all other dimensions have value 0. Recall that for any two points  $p_a$  and  $p_b$ , their associated hyperplane is denoted by  $h_{ab}$  with its normal  $v_{ab} = p_a - p_b$ . Further, the halfspace indicated by the user, denoted by  $s_{ab}$ , could be either  $h_{ab}^+$  or  $h_{ab}^-$ . For ease of illustration, from now on, we assume without loss of generality that  $s_{ab} = h_{ab}^+$ . Mathematically,  $s_{ab} = \{u \in U | v_{ab} \cdot u > 0\}$  where  $U$  is the utility space.

Consider any non-empty set of halfspaces  $HS$ . Assume that the intersection of halfspaces in  $HS$  forms a non-empty region, which we denote by  $\mathcal{R}$  (i.e.,  $\mathcal{R} = \bigcap_{s \in HS} s$ ). We show that there is a non-empty intersection between  $P_i$  and  $\mathcal{R}$ . Specifically, let  $e_i$  be the vector with  $e[i] = 1$  and all other entries  $e[j], j \neq i$  being 0. Since the utility of  $p_i$  w.r.t.  $e_i$  is 1 and the utilities of all other points  $p_j, j \neq i$  w.r.t.  $e_j$  is 0,  $e_i \in P_i$  by the definition of  $P_i$ . Further, for any

$s_{ab} \in HS$ , we have  $e_j \cdot v_{ab} = e_j \cdot (p_a - p_b) = 0$ . Now consider a vector  $r \in \mathcal{R}$  and a small positive number  $\lambda$ . Construct  $e' = e_i + \lambda r$ . Given a small enough  $\lambda$ , we have  $e' \in P_i$  since  $e_i \in P_i$ . Further,  $e' \in \mathcal{R}$  since for any  $s_{ab} \in HS$ , we have  $e' \cdot v_{ab} = (e_i + \lambda r) \cdot v_{ab} = \lambda r \cdot v_{ab} > 0$ . Since  $e'$  is both in  $P_i$  and  $\mathcal{R}$ , we conclude that  $P_i$  and  $\mathcal{R}$  have a non-empty intersection.

Given a dataset  $D$  with Property 1, and given a partition  $P_i$  that currently belongs to a confidence region  $R^j$ , there are only two ways to detach  $P_i$  from  $R^j$ , which we call them the *direct way* and the *indirect way*. One can verify these are indeed the only way to detach  $P_i$  from  $R^j$  given Property 1.

**The direct way.**  $P_i$  is detached from  $R^j$  if the user prefers another point  $p_h$  to  $p_i$ .

**The indirect way.**  $P_i$  is detached from  $R^j$  if  $R^j$  becomes empty.

By Lemma 3, if a partition is detached from  $R_t^j$  in round  $t$ , it will belong to  $R_{t+1}^{j+1}$  in round  $t+1$ . Since each partition initially belongs to  $R^0$  and needs to be detached from  $R^k$  (where  $k = \lfloor \frac{l-1}{2} \rfloor$ ) before it is discarded, we conclude that to discard a partition, it needs to be detached for  $k+1$  times in total. The algorithm stops only when at most  $l$  points are not discarded. Since the direct way detaches one partition from some confidence region at each time, if partitions are only detached using the direct way, then  $\Omega(kn)$  rounds are needed. It remains to bound the number of rounds required if the indirect way is used. To bound this number, we first present the following lemma.

**LEMMA 7.** *There exists a dataset such that for any question-asking strategy, the  $i$ -th confidence region  $R^i = \emptyset$  if and only if there are  $i+1$  disjoint cycles.*

**PROOF.** The “if” part is already proved in Lemma 5. For the “only if” part, we prove the following equivalent statement: There exists a dataset such that if there are no  $i+1$  disjoint cycles, then the  $i$ -th confidence region  $R^i \neq \emptyset$ . To see this, we use the same dataset  $D_n$  constructed above. Assume that the user has made  $t$  comparisons. Since there are at most  $i$  disjoint cycles in these comparisons, by selectively removing  $i$  of them, the remaining  $t-i$  comparisons will have no cycle. In other words, the remaining  $t-i$  comparisons conform with some ranking on the points involved in these comparisons. Assume without generality that they conform with the ranking  $p_1 > p_2 > p_3 > \dots > p_m$ , where  $p_1, \dots, p_m$  are the points involved in these  $t-i$  comparisons. We can then construct a utility vector  $u$  that yields this ranking. Specifically, let  $u$  be constructed such that  $u[i] > u[i+1]$  for any  $i \in [1, m-1]$ . Then,  $\forall i \in [1, m-1]$ ,  $u \cdot p_i > u \cdot p_{i+1}$  since  $u \cdot p_i = u[i] > u[i+1] = u \cdot p_{i+1}$ . This means the  $t-i$  halfspaces corresponding to the  $t-i$  preferences support a non-empty region, and thus by definition,  $R^i$  is also non-empty.  $\square$

If we  $n-l$  partitions need to be discarded using the indirect way, Lemma 7 tells us that we need  $k$  disjoint cycles to make  $R^0, R^1, \dots, R^{k-1}$  empty (Note that  $R^k$  cannot be empty, otherwise, there is no point to return.).  $k$  user errors are required to form  $k$  disjoint cycles. Since the expected number of rounds for the user to make one error is  $\frac{1}{\theta}$ , we need at least  $\Omega(\frac{k}{\theta})$  rounds to accumulate  $k$  errors. After obtaining  $k$  disjoint cycles, all partitions are either discarded or belong to  $R^k$ . Since each question detaches at

most 1 partition from  $R^k$ , We still need an extra number of  $\Omega(n)$  rounds to detach partitions from  $R^k$  until only  $l$  of them are left. Therefore, the total expected number of rounds using the indirect way is  $\Omega(\frac{k}{\theta} + n)$ . Since the number used by the direct way is  $\Omega(kn)$  and  $\Omega(\frac{k}{\theta} + n) < \Omega(kn)$ , and  $k = \lfloor \frac{l-1}{2} \rfloor$ , the lower bound is thus  $\Omega(\frac{l}{\theta} + n)$ .

## B.5 Proof of Lemma 4

In this section, we give proof to Lemma 4. When the algorithm ends, if the cell where the real utility vector lies overlaps with  $R^k$ , then the best point will not be returned only if there is no sample point inside this cell. Assume that there are  $\rho$  cells in the partition of the best point, and the total number of points sampled from this partition is  $m$ , further, assume that these  $\rho$  cells in this partition are denoted by  $c_1, c_2, \dots, c_\rho$  and have sizes  $f_1, f_2, \dots, f_\rho$  proportional to the total size of this partition, where  $\sum_{i=1}^\rho f_i = 1$ . With no prior knowledge of the distribution of the user’s utility vector, the utility vector is assumed to be uniformly distributed. The probability that a cell  $c_i$  contains the user’s utility vector and has no sample point in it is then  $f_i(1-f_i)^m$ . Thus, the probability that when  $SS$  ends, the cell containing the utility vector has no sample point in it is  $\sum_{i=1}^\rho f_i(1-f_i)^m$ , which is further upper bounded by  $\sum_{i=1}^\rho f_i e^{-f_i m}$ .  $f_i e^{-f_i m}$  takes the maximum when  $f_i = \frac{1}{m}$ . Thus,  $\sum_{i=1}^\rho f_i e^{-f_i m} \leq \frac{\rho}{m} e^{-1}$ . Solving  $\frac{\rho}{m} e^{-1} = \varepsilon$  yields  $m \geq \frac{\rho}{\varepsilon e}$ .

## B.6 Proof of Theorem 2

In this section, we show the proof of Theorem 2. Let  $|V|$  denote the maximum number of vertices in all partitions. Since a partition is bounded by at most  $n$  halfspaces, computing all partitions takes  $O(|V|dn)$  time [3]. By [7], sampling one point from a  $d$ -dimensional polytope bounded by  $m$  halfspaces can be done in  $O(md)$  time. Since each partition is bounded by at most  $n$  halfspaces, sampling all  $Yn$  points takes  $O(Ydn^2)$  time. The time complexity for the pre-processing phase is thus  $O((|V| + Y)dn^2)$ .

Next, we bound the time complexity in each round of *SS-random* and *SS-score*. For *SS-random*, since the time for updating each sample point is  $O(d)$ , the time required to update  $Q^0, \dots, Q^k$  is  $O(Yldn)$ .  $X_t^0, X_t^1, \dots, X_t^k$  can be determined by scanning through all  $Q_t^0, Q_t^1, \dots, Q_t^k$  and can be done in  $O(Yln)$  time. The random-based selection takes  $O(t)$  time, where  $t$  is the current number of rounds, since it always finds an unasked pair after checking at most  $t+1$  candidates. Since typically  $t \ll Ydn$ , the time complexity for each round of *SS-random* is  $O(Yldn)$ .

For *SS-score*, the time required to update  $Q_t^i$ s and compute  $X_t^i$ s is the same as above (i.e.,  $O(Yldn)$ ). The time for determining the next question is  $O(n^3)$ , since computing the score for all  $O(n^2)$  hyperplane candidates takes  $O(n)$  time each. Therefore, the time complexity for each round of *SS-score* is  $O(Yldn + n^3)$ .

## B.7 Proof of Theorem 3

We present the proof of Theorem 3 in this section. Firstly, *SS* does not return the best point only if (1) there is no sample point in the cell containing the real utility vector; or (2) the cell containing the real utility vector does not overlap with  $R^k$ . By Lemma 4, (1)

happens with probability at most  $\varepsilon$ . (2) happens only if the user makes more than  $k$  errors out of  $T$  questions, where  $T$  is the number of rounds used by SS before termination. Since the expected number of error made by the user is  $\theta T$ , using Chernoff inequality, this probability is upper bounded by  $e^{-\frac{(k-\theta T)^2}{k+\theta T}} = e^{-\frac{((l-1)/2-\theta T)^2}{[(l-1)/2]+\theta T}}$ . Thus, the probability that SS returns the best point is at least  $1 - \varepsilon - e^{-\frac{((l-1)/2-\theta T)^2}{[(l-1)/2]+\theta T}}$ , which completes the proof.

## B.8 Proof of Lemma 5

In this section, we give proof to Lemma 5. First, we show that the intersection of halfspaces corresponding to a cycle is an empty region. Consider a cycle consists of points  $p_1, p_2, \dots, p_m$ , where the user indicated  $p_1 > p_2, p_2 > p_3, \dots, p_{m-1} > p_m$  and  $p_m > p_1$ . A utility vector  $u$  that satisfies all these preferences does not exist, since we must have  $u \cdot p_1 > u \cdot p_m$  and  $u \cdot p_m > u \cdot p_1$ . Therefore, the intersection of halfspaces corresponding to a cycle is an empty region.

Since  $R^i = \bigcup_{\bar{S} \in \binom{S}{i}} (\cap_{s \in \bar{S}} s)$ , where  $\binom{S}{i}$  is all  $i$ -subset of  $S$ , if there are  $i+1$  disjoint cycles, no matter how we set  $\bar{S}$ , there is at least one cycle in  $S/\bar{S}$ , and  $\cap_{s \in S/\bar{S}} s$  is an empty region. The union of a set of empty regions is also empty.

## B.9 Proof of Lemma 6

The proof comprises two parts. In the first part, we show that when Stage 1 ends, the number of obtained disjoint cycles is at most  $k$ . In the second part, we show when Stage 1 ends on Condition (1), the expected number of disjoint cycles obtained is at least  $\max(0, k-17)$ .

We first prove that when Stage 1 ends, the number of obtained disjoint cycles is at most  $k$ . Assume that Stage 1 ends just after round  $t$ . We know that  $\bigcup_{i=0}^{k-1} X_{t-1}^i \neq \emptyset$ , since otherwise, Stage 1 should ends after round  $t-1$ . Thus,  $R_{t-1}^{k-1} \neq \emptyset$ . Then, by Lemma 3,  $R_{t-1}^{k-1} \subseteq R_t^k \neq \emptyset$ . Given that  $R_t^k \neq \emptyset$ , by applying Lemma 5, the number of obtained disjoint cycles cannot exceed  $k$ , since otherwise  $R_t^k$  will become empty.

We then prove that when Stage 1 ends on Condition (1), the expected number of disjoint cycles obtained is at least  $\max(0, k-17)$ . To see this, we first show that when Stage 1 ends on Condition (1), the user makes at least  $k$  errors with probability at least  $1 - \frac{1}{k}$ . We then show that if at least  $k$  errors are made, the expected number of disjoint cycles is at least  $k-16$ . Combining them together, the expected number of disjoint cycles is at least  $(1 - \frac{1}{k}) \cdot (k-16) \geq k-17$ .

To see that the user makes at least  $k$  errors with probability at least  $1 - \frac{1}{k}$ , note that when Condition 1 is satisfied, Stage 1 lasts for at least  $\max(\frac{k(k-1)}{2}, \frac{3k}{\theta})$  rounds. Let  $X$  be the random variable denoting the number of errors the user made after at least  $\frac{3k}{\theta}$  rounds. The expected number of errors is  $\mu_X = 3k$ , and the variance is  $\sigma_X^2 = 3k(1-\theta)$ . Using Chebyshev inequality, we obtain that  $P(X < k) \leq \frac{1}{k}$ .

Next, we show that if the user already makes  $k$  errors, the expected number of disjoint cycles is at least  $k-16$ . Note that when Stage 1 ends on Condition (1), there are at least  $m \geq k$  points in  $PS$  (recall that  $PS$  is the set of points maintained in Stage 1 of

FC), and every pair of these points has been compared. We constrain ourselves on finding one specific type of cycles: it is formed by 3 points, namely  $p_i, p_j$ , and a point  $p_a$  where  $i < a < j$  (i.e.,  $u \cdot p_i > u \cdot p_a > u \cdot p_j$ ), and the user indicates  $p_i > p_a, p_a > p_j$ , and  $p_j > p_i$ , where the error is  $p_j > p_i$ . To distinguish the different roles these 3 points play in a cycle, We call  $p_i$  and  $p_j$  the *basic points* (i.e., the points that are involved in the error) and  $p_a$  the *non-basic point* (i.e., the point that is not involved in the error).

We will show that by simply considering this one type of cycles, the expected number of disjoint cycles is already at least  $k-16$ . With this construction, two cycles  $(p_i, p_j, p_a)$  and  $(p_{i'}, p_{j'}, p_{a'})$  may not be disjoint only if they share at least one basic point (i.e.,  $p_i = p_{i'}$  or  $p_j = p_{j'}$ ). To see this, observe that if  $p_i \neq p_{i'}$  and  $p_j \neq p_{j'}$ , no matter what values  $a$  and  $a'$  take, the two cycles  $(p_i, p_j, p_a)$  and  $(p_{i'}, p_{j'}, p_{a'})$  cannot share a common preference, and must be disjoint. Consider a comparison of pair  $(p_i, p_j)$  where  $i < j$ , on which the user makes an error. We want to bound the probability that this error does not create a disjoint cycle. Among the remaining  $k-1$  errors, the number of errors that involve  $p_i$  or  $p_j$  is at most 4 (since  $m \geq k$  and among the remaining  $\frac{m(m-1)}{2} - 1$  pairs,  $2m-4$  of them contain  $p_i$  or  $p_j$ ). This means that on expectation, at most 4 other cycles contain  $p_i$  or  $p_j$  as a basic point. Consequently, among the cycles that contain  $p_i$  or  $p_j$  as a basic point, at most 8 points on expectation (i.e., one basic point and one non-basic point from each cycle) are from the  $(j-i-1)$  points ranked between  $p_i$  and  $p_j$ . As long as  $p_a$  where  $i < a < j$  is not chosen from these points, cycle  $(p_i, p_j, p_a)$  cannot share a common preference with any other cycles and must be disjoint.

Since each pair of points is equally probable to be pair  $(p_i, p_j)$ , we consider different values of  $i$  and  $j$ :

- Case  $j = i+1$ . Since there are  $\frac{m(m-1)}{2}$  comparisons in total and only  $m-1$  of them satisfies  $j = i+1$ , this case happens with probability  $\frac{2}{m}$ .
- Case  $j = i+2$ . This case happens with probability  $\frac{2(m-2)}{m(m-1)}$ .
- Case  $j = i+3$ . This case happens with probability  $\frac{2(m-3)}{m(m-1)}$ .
- ...

We can show that if  $y$  points between  $p_i$  and  $p_j$  are involved in some other cycles that contain  $p_i$  or  $p_j$  as a basic point, the probability that a user error when comparing  $p_i$  and  $p_j$  does not cause a disjoint cycle is upper bounded by  $\sum_{x=1}^y \frac{2(m-x)}{m(m-1)} \leq \frac{2y}{m}$ . Given that  $E[y] \leq 8$ ,  $E[\frac{2y}{m}] \leq \frac{16}{m}$ . Thus, the expected number of disjoint cycles obtained is at least  $k - \frac{16}{m}k \geq k-16$ .

Since when Stage 1 ends on Condition (1), the user makes at least  $k$  errors with probability at least  $1 - \frac{1}{k}$ , and if at least  $k$  errors are made, the expected number of disjoint cycles is at least  $k-16$ . Combining them together, the expected number of disjoint cycles is at least  $(1 - \frac{1}{k}) \cdot (k-16) \geq k-17$ .

## B.10 Proof of Theorem 4

In this section, we show the proof of Theorem 4. Firstly, since  $k = \left\lfloor \frac{l-1}{2} \right\rfloor$  and  $l$  can be regarded as a constant  $< O(\sqrt{n})$ , Stage 1 of FC runs for at most  $O(\max(\frac{k(k-1)}{2}, \frac{3k}{\theta})) = O(\frac{l}{\theta} + n)$  rounds. When FC enters stage 2, there are 2 cases: (1) Condition 1 is satisfied, and



less than  $k$  errors are made in Stage 1, and (2) Condition 1 is not satisfied, but Condition 2 is satisfied. Next, we bound the number of rounds required for these 2 cases.

For case (1), by Lemma 6, since Stage 1 ends on Condition (1), the expected number of disjoint cycles is at least  $k - 17$ . Let  $k'$  be the number of disjoint cycles that appear. We have  $E[k - k'] \leq 17$ . Then by Lemma 5, confidence regions  $R^0, R^1, \dots, R^{k'-1}$  are all empty, and all remaining partitions are either discarded or belong to  $R^{k'}, \dots, R^k$ . Since the expected value of  $k - k'$  is a constant, we only need to detach partitions from confidence regions for  $O(n)$  times in total. Since each round detaches at least one partition from one confidence region, the expected number of questions required

to discard all except  $l$  remaining partitions is  $O(n)$  (see the proof of Theorem 1 in Section B.4).

For case (2), since Condition (2) is satisfied,  $R^{k-1}$  is empty. All partitions are either discarded, or belong to  $R^k$ . Since each round detaches at least one partition from one confidence region, Stage 2 ends in  $O(n)$  rounds.

Because both cases take  $O(n)$  rounds, the expected number of rounds for both cases is  $O(n)$ . Since Stage 1 takes  $O(\frac{l}{\theta} + n)$  rounds and Stage 2 takes  $O(n)$  rounds on expectation, the total expected number of rounds for  $FC$  is  $O(\frac{l}{\theta} + n)$ .

### B.11 Proof of Theorem 5

The proof is nearly identical to the proof of Theorem 3.