



# Efficient Query Repair for Aggregate Constraints

Shatha Algarni  
University of Southampton  
University of Jeddah  
s.s.algarni@soton.ac.uk

Boris Glavic  
University of Illinois,  
Chicago  
bglavic@uic.edu

Seokki Lee  
University of Cincinnati  
lee5sk@ucmail.uc.edu

Adriane Chapman  
University of Southampton  
adriane.chapman@soton.ac.uk

## ABSTRACT

In many real-world scenarios, query results must satisfy domain-specific constraints. For instance, a minimum percentage of interview candidates selected based on their qualifications should be female. These requirements can be expressed as constraints over an arithmetic combination of aggregates evaluated on the result of the query. In this work, we study how to repair a query to fulfill such constraints by modifying the filter predicates of the query. We introduce a novel query repair technique that leverages bounds on sets of candidate solutions and interval arithmetic to efficiently prune the search space. We demonstrate experimentally, that our technique significantly outperforms baselines that consider a single candidate at a time.

### PVLDB Reference Format:

Shatha Algarni, Boris Glavic, Seokki Lee, and Adriane Chapman. Efficient Query Repair for Aggregate Constraints. PVLDB, 19(2): 252 - 264, 2025. doi:10.14778/3773749.3773762

### PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/ShathaSaad/Query-Refinement-of-Complex-Constraints>.

## 1 INTRODUCTION

Analysts are typically well versed in writing queries that return data based on obvious conditions, e.g., only return applicants with a master’s degree. However, a query result often has to fulfill additional constraints, e.g. fairness, that do not naturally translate into conditions. While for some applications it is possible to filter the results of the query to fulfill such constraints this is not always viable, e.g., because the same selection criterion has to be used for all job applicants. Thus, the query has to be repaired such that the result set of the fixed query satisfies all constraints. Prior work in this area, including query-based explanations [15, 35] and repairs [9] for missing answers, work on answering why-not questions [5, 15] as well as query refinement / relaxation approaches [26, 29, 36] determine why specific tuples are not in the query’s result or how to fix the query to return such tuples. In this work, we study a more general problem where the *entire result set* of the query has to fulfill some constraint. The constraints we study in this work are expressive enough to guarantee that query results adhere to legal and ethical regulations, such as fairness.

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.  
Proceedings of the VLDB Endowment, Vol. 19, No. 2 ISSN 2150-8097.  
doi:10.14778/3773749.3773762

**EXAMPLE 1 (MOTIVATING EXAMPLE).** Consider a job applicant dataset  $D$  for a tech-company that contains six attributes: ID, Gender, Field, GPA, TestScore, and OfferInterview. The attribute OfferInterview was generated by an external ML model suggesting which candidates should receive an interview. The employer uses the query shown below to prescreen candidates: every candidate should be a CS graduate and should have a high GPA and test score.

**Q1:** `SELECT * FROM D WHERE Major = 'CS' AND TestScore ≥ 33 AND GPA ≥ 3.80`

**Aggregate Constraint.** The employer wants to ensure that interview decisions are not biased against a specific gender using statistical parity difference (SPD) [4, 27]. Given two groups, e.g., male and female, and a binary outcome attribute  $Y$  where  $Y = 1$  is assumed to be a positive outcome (OfferInterview=1 in our case), the SPD is the difference between the probability for individuals from the two groups to receive a positive outcome. In our example, the SPD can be computed as shown below ( $G$  is Gender and  $Y$  is OfferInterview). We use  $\text{count}(\theta)$  to denote the number of query results satisfying condition  $\theta$ . For example,  $\text{count}(G = M \wedge Y = 1)$  counts male applicants with a positive label.

$$\text{SPD} = \frac{\text{count}(G = M \wedge Y = 1)}{\text{count}(G = M)} - \frac{\text{count}(G = F \wedge Y = 1)}{\text{count}(G = F)}$$

The employer would like to ensure that the SPD between male and female is below 0.2. The model generating the OfferInterview attribute is trusted by the company, but is provided by an external service and, thus, cannot be fine-tuned to improve fairness. However, the employer is willing to change their prescreening criteria by expressing their fairness requirement as long as the same criteria are applied to judge every applicant to ensure individual fairness. This can be achieved using an aggregate constraint  $\text{SPD} \leq 0.2$ . That is, the employer desires a repair of the query whose selection conditions are used to filter applicants. We present additional motivating examples beyond fairness in [3].

In this work, we model constraints on the query result as arithmetic expressions involving aggregate queries evaluated over the output of a *user query*. When the result of the user query fails to adhere to such an *aggregate constraint (AC)*, we would like the system to fix the violation by *repairing* the query by adjusting its selection conditions, similar to [24, 29]. Specifically, we are interested in computing the top- $k$  repairs with respect to their distance to the user query. The rationale is that we would like to preserve the original semantics of the user’s query as much as possible. Moreover, instead of assuming a single best repair, we consider returning  $k$  repairs ranked by their distance to the original query to allow users to choose the one that best matches their intent. ACs significantly generalize the cardinality constraints supported in prior work on query repair for fairness [13, 14, 26] and on query relaxation &

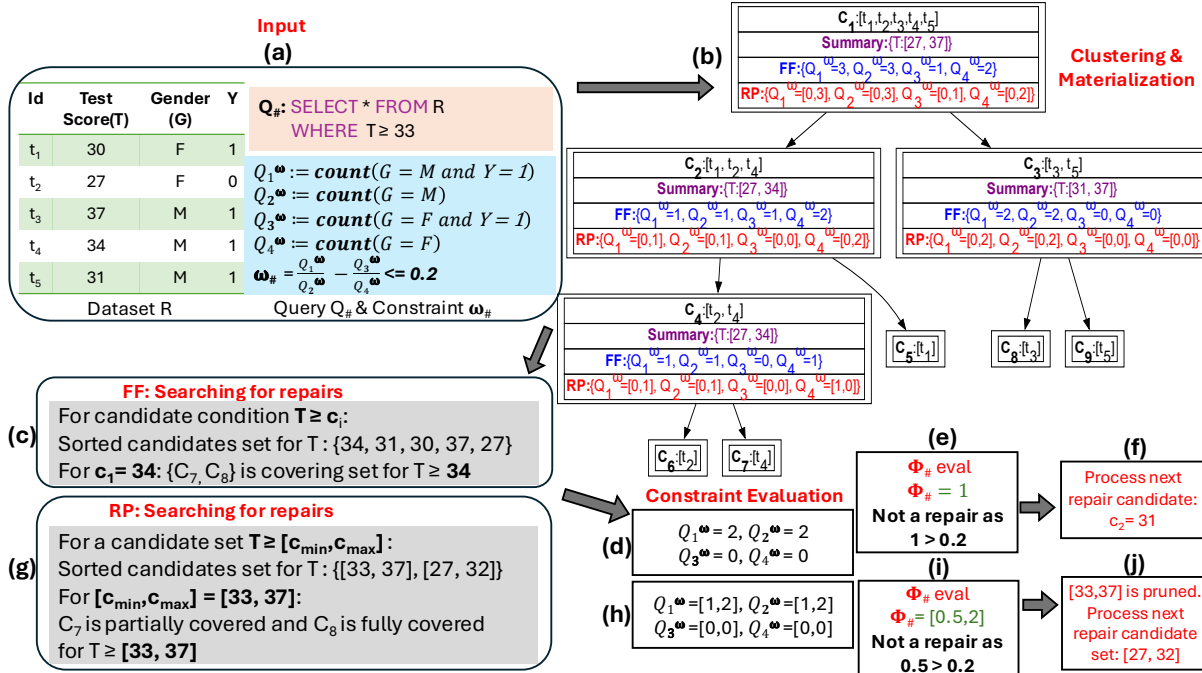


Figure 1: Overview of query repair with aggregate constraints using range-based pruning.

refinement [25, 29]. By allowing arithmetic combinations of aggregation results we support common fairness measures such as SPD that cannot be expressed as cardinality constraints. Our work is applicable to any use case where uniform criteria have to be applied to select a set of entities subject to additional constraints, e.g., a government agency has to solicit contractors, 20% of which should be local (see [3] for a detailed example). New challenges arise from the generality of ACs as ACs are typically not monotone, invalidating most optimizations proposed in related work.

A brute force approach for solving the query repair problem is to enumerate all possible candidate repairs in order of their distance to the user query. Each candidate is evaluated by running the modified query and checking whether it fulfills the aggregate constraint. The algorithm terminates once  $k$  repairs have been found. The main problem with this approach is that the number of repair candidates is exponential in the number of predicates in the user query. Furthermore, for each repair candidate we have to evaluate the modified user query and one or more aggregate queries on top of its result. Given that the repair problem is NP-hard in the number of predicates we cannot hope to avoid this cost in general.

**Reusing aggregation results.** Nonetheless, we identify two opportunities for optimizing this process. When two repair candidates are similar (in terms of the constants they use in selection conditions), then typically there will be overlap between the aggregate constraint computations for the two candidates. To exploit this observation, we use a kd-tree [6] to partition the input dataset. For each cluster (node in the kd-tree) we materialize the result of evaluating the aggregation functions needed for a constraint on the set of tuples contained in the cluster as well as store bounds for the attribute values within the cluster (as is done in, e.g., zonemaps [31, 40]). Then to calculate the result of an aggregation function for a repair

candidate, we use the bounds for each cluster to determine whether all tuples in the cluster fulfill the selection conditions of the repair candidate (in this case the materialized aggregates for the cluster will be added to the result), none of the tuples in the cluster fulfill the condition (in this case the whole cluster will be skipped), or if some of the tuples in the cluster fulfill the condition (in this case we apply the same test to the children of the cluster in the kd-tree). We refer to this approach as *Full Cluster Filtering (FF)*. In contrast to the brute force approach, FF reuses aggregation results materialized for clusters. Continuing with Example 1, consider the kd-tree in Figure 1(b) which partitions the input dataset  $R$  in Figure 1(a) into a set of clusters. Here, we simplify  $Q_1$  from Example 1 by considering only a single condition,  $TestScore(T) \geq 33$ , but use the same aggregate constraint  $\omega_{\#}$ . Consider cluster  $C_2$  in Figure 1(b), where the values of attribute  $TestScore(T)$  are bounded by [27, 34]. For a repair candidate with a condition  $T \geq 37$  the entire cluster can be skipped as no tuples in  $C_2$  can fulfill the condition. In contrast, for condition  $T \geq 30$ , all tuples in  $C_3$  satisfy the condition, since the values of attribute  $T$  are bounded by [31, 37] (see Figure 1(b)).

**Evaluating multiple candidate repairs at once.** We extend this idea to bound the aggregation constraint result for sets of repair candidates at once. We refer to this approach as *Cluster Range Pruning (RP)*. A set of repair candidates is encoded as intervals of values for the constant  $c_i$  of each predicate  $a_i$  op  $c_i$  of the user query, e.g.,  $c_1 \in [33, 37]$  as shown in Figure 1(g). We again reason about whether all / none of the tuples in a cluster fulfill the condition for every repair candidate from the set. The results are valid bounds on the aggregation constraint result for any candidate repair within the set. Using these bounds, we validate or disqualify complete candidate sets at once.

We make the following contributions in this work:

- A formal definition of query repair under constraints involving arithmetic combinations of aggregate functions in Section 2.
- We present an optimized algorithm for the aggregate constraint repair problem that reuses aggregation results when evaluating repair candidates (Section 3) and evaluates multiple repair candidates by exploiting sound bounds that hold for all repair candidates in a set (Section 4).
- A comprehensive experimental evaluation over multiple datasets, queries, and constraints in Section 5. Compared to the state of the art [26], we cover significantly more complex constraints.

## 2 PROBLEM DEFINITION

We consider a dataset  $D = R_1, \dots, R_z$  consisting of one or more relations  $R_i$ , an input query  $Q$  that should be repaired, and an aggregate constraint. The goal is to find the  $k$  queries that fulfill the constraint and minimally differ from  $Q$ .

**User Query.** A user query  $Q$  is a select-project-join (SPJ) query, i.e., a relational algebra expression of the form:

$$\pi_A(\sigma_\theta(R_1 \bowtie \dots \bowtie R_l))$$

We assume that the selection predicate  $\theta$  of such a query is a conjunction  $\theta = \theta_1 \wedge \dots \wedge \theta_m$  of comparisons of the form  $a_i \text{ op}_i c_i$ . For numerical attributes  $a_i$ , we allow  $\text{op}_i \in \{<, >, \leq, \geq, =, \neq\}$  and for categorical attributes  $a_i$  we only allow  $\text{op}_i \in \{=, \neq\}$ . We use  $Q(D)$  to denote the result of evaluating  $Q$  over  $D$ .

**Aggregate Constraints (AC).** The user specifies requirements on the result of their query as an AC. An AC is a comparison between a threshold and an arithmetic expression over the result of filter-aggregation queries. Such queries are of the form  $\gamma_{f(a)}(\sigma_\theta(Q(D)))$  where  $f$  is an aggregate function – one of **count**, **sum**, **min**, **max**, **avg** – and  $\theta$  is a selection condition. We use  $Q^\omega$  to denote such a filter-aggregation query. These queries are evaluated over the user query’s result  $Q(D)$ . An aggregate constraint  $\omega$  is of the form:

$$\omega := \tau \text{ op } \Phi(Q_1^\omega, \dots, Q_n^\omega).$$

Here,  $\Phi$  is an arithmetic expression using operators  $(+, -, *, /)$  over  $\{Q_i^\omega\}$ ,  $\text{op}$  is a comparison operator, and  $\tau$  is a (constant) threshold. Aggregate constraints are non-monotone in general due to (i) non-monotone arithmetic operators like division, (ii) non-monotone aggregation functions, e.g., **sum** over the integers  $\mathbb{Z}$ , and (iii) combination of monotonically increasing and decreasing aggregation functions, e.g.,  $\max(A) + \min(B)$ .

**Query Repair.** Given a user query  $Q$ , database  $D$ , and constraint  $\omega$  that is violated on  $Q(D)$ , we want to generate a repaired version  $Q_{fix}$  of  $Q$  such that  $Q_{fix}(D)$  fulfills  $\omega$ . We restrict repairs to changes of the selection condition  $\theta$  of  $Q$ . For ease of presentation, we consider a single AC, but our algorithms can also handle a conjunction of multiple ACs, e.g., the cardinality for one group should be above a threshold  $\tau_1$  and for another group below a threshold  $\tau_2$ . Given the user query  $Q$  with a condition  $\theta := \bigwedge_{i=1}^m a_i \text{ op}_i c_i$ , a *repair candidate* is a query  $Q_{fix}$  that differs from  $Q$  only in the constants used in selection conditions, i.e.,  $Q_{fix}$  uses a condition:  $\theta' := \bigwedge_{i=1}^m a_i \text{ op}_i c'_i$ . For convenience, we will often use the vector of constants  $\vec{c} = [c'_1, \dots, c'_m]$  to denote a repair candidate and use

$\text{CAND}_Q$  to denote the set of all candidates. A candidate is a *repair* if  $Q_{fix}(D) \models \omega$ .

**Repair Distance.** Ideally, we would want to achieve a repair that minimizes the changes to the user’s original query to preserve the intent of the user’s query as much as possible. We measure similarity using a weighted linear combination of distances between the constants used in selection conditions of the user query and the repair, similar to [13, 25].<sup>1</sup> Consider a user query  $Q$  with selection condition  $\theta_1 \wedge \dots \wedge \theta_m$  and repair  $Q_{fix}$  with selection condition  $\theta'_1 \wedge \dots \wedge \theta'_m$ . Then the distance  $d(Q, Q_{fix})$  is defined as:

$$d(Q, Q_{fix}) = \sum_{i=1}^m w_i \cdot d(\theta_i, \theta'_i)$$

where  $w_i$  is a weight in  $[0, 1]$  such that  $\sum_i w_i = 1$  and the distance between two predicates  $\theta_i = a_i \text{ op}_i c_i$  and  $\theta'_i = a_i \text{ op}_i c'_i$  for numeric attributes  $a_i$  is:  $\frac{|c'_i - c_i|}{|c_i|}$ . For categorical attributes, the distance is 1 if  $c_i \neq c'_i$  and 0 otherwise. For example, for Example 1, the repair candidate with conditions **Major** = EE, **Testscore**  $\geq 33$ , and **GPA**  $\geq 3.9$  has a distance of  $1 + \frac{33-33}{33} + \frac{3.9-3.8}{3.8} = 1.026$ .

We are now ready to formulate the problem studied in this work, computing the  $k$  repairs with the smallest distance to the user query. Among these  $k$  repairs, the user can then select the repair that best aligns with their preferences. Here  $\text{top-}k_{X \in X} f(x)$  returns the  $k$  elements from set  $X$  with the smallest  $f(x)$  values.

### AGGREGATE CONSTRAINT REPAIR PROBLEM:

- **Input:** user query  $Q$ , database  $D$ , constraint  $\omega$ , threshold  $k$
- **Output:**

$$\text{top-}k_{Q_{fix} \in \text{CAND}_Q: Q_{fix}(D) \models \omega} d(Q, Q_{fix})$$

**Hardness.** To generate a repair  $Q_{fix}$  of  $Q$ , we must explore the combinatorially large search space of possible candidate repairs. For a single predicate over an attribute  $a_i$  with  $N_i$  distinct values there are  $O(N_i)$  possible repairs. Thus, the size of the candidate set  $\text{CAND}_Q$  is in  $O(\prod_{i=1}^m N_i)$ , exponential in  $m$ , the number of conditions in the user query. Unsurprisingly, the aggregate constraint repair problem is NP-hard in the schema size. In [3], we provide more details on the number of repairs for specific predicates.

## 3 THE FULL CLUSTER FILTERING ALGORITHM

We now present *Full Cluster Filtering (FF)*, our first algorithm for the aggregate constraint repair problem that materializes results of each filter-aggregation query  $Q_i^\omega$  for subsets of the input database  $D$  in a kd-tree. FF combines these aggregation results to compute the result of  $Q_i^\omega$  for a repair candidate  $Q_{fix}$  and then uses these results to evaluate the aggregate constraint (AC)  $\omega$  for  $Q_{fix}$ . Figure 1 shows an example of applying this algorithm: (b) building a kd-tree and materializing statistics, (c) searching for candidate repairs, and (d)-(e) evaluating constraints for repair candidates.

<sup>1</sup>We discuss other possible objectives used in prior work in Section 6. Other options include returning all repairs that are Pareto optimal regarding predicate-level distances or to minimize the change to the query’s result. Our algorithms can be extended to optimize for any distance metric which can be interval-bounded based on bounds for attribute values of a set of tuples.

### 3.1 Clustering and Materializing Aggregations

For ease of presentation, we consider a database consisting of a single table  $R$  from now on. However, our approach can be generalized to queries involving joins by materializing the join output and treating it as a single table. As repairs only change the selection conditions of the user query, there is no need to reevaluate joins when checking repairs. We use a kd-tree to partition  $R$  into subsets (*clusters*) based on attributes that appear in the selection condition ( $\theta$ ) of the user query. The rationale is that the selection conditions of a repair candidate filter data along these attributes.

To evaluate the AC  $\omega$  for a candidate  $Q_{fix} = [c'_1, \dots, c'_m]$ , we determine a set of clusters (nodes in the kd-tree) that cover exactly the subset of  $D$  that fulfills the selection condition of the candidate. We can then merge the materialized aggregation results for these clusters to compute the results of filter-aggregation queries  $Q_i^\omega$  used in  $\omega$  for  $Q_{fix}(D)$ . To do that, we record the following information for each cluster  $C \subseteq D$  that can be computed by a single scan over the tuples in the cluster, or by combining results from previously generated clusters if we generate clusters bottom up.

- **Selection attribute bounds:** For each attribute  $a_i$  used in the condition  $\theta$ , we store  $\text{BOUNDS}_{a_i} := [\min(\pi_{a_i}(C)), \max(\pi_{a_i}(C))]$ .
- **Count:** The total number of tuples  $\text{count}(C) := |C|$  in the cluster.
- **Aggregation results:** For each filter-aggregation query  $Q^\omega$  in constraint  $\omega$ , we store  $Q^\omega(C)$ .

An example kd-tree is shown in Figure 1(b). The user query filters on attribute *TestScore* ( $T$ ). The root of the kd-tree represents the full dataset. At each level, the clusters from the previous level are split into  $\mathcal{B}$  sub-clusters along one of the attributes in  $\theta$ .  $\mathcal{B}$ , called the branching factor, is a configuration parameter. We use  $\mathcal{B} = 2$  in the example. For instance, the root cluster  $C_1$  is split into two clusters  $C_2$  and  $C_3$  by partitioning the rows in  $C_1$  based on their values in attribute  $T$ . For cluster  $C_2$  containing three tuples  $t_1, t_2$ , and  $t_4$ , we have  $\text{BOUNDS}_T = [27, 34]$  as the lowest  $T$  value is 27 (from tuple  $t_2$ ) and the highest value is 34 (tuple  $t_4$ ). The value of  $Q_2^\omega = \text{count}(\text{Gender}(G) = M)$  for  $C_2$  is 1 as there is one male in the cluster. Consider a repair candidate with the condition  $T \geq 37$ . Based on the bounds  $\text{BOUNDS}_T = [27, 34]$ , we know that none of the tuples satisfy this condition. Thus, this cluster and the whole subtree rooted at the cluster can be ignored for computing the AC  $\omega_\#$  for the candidate.

For ease of presentation we assume that the leaf nodes of the kd-tree contain a single tuple each. As this would lead to very large trees, in our implementation we do not further divide clusters  $C$  that contain less tuples than a threshold  $S$ . i.e.,  $|C| \leq S$ . We refer to this parameter as the *bucket size*.

### 3.2 Constraint Evaluation for Candidates

The FF algorithm (Algorithm 1) takes as input the condition  $\theta'$  of a repair candidate, the root node of the kd-tree  $C_{root}$ , and returns a set of disjoint clusters  $C$  such that the union of these clusters is precisely the subset of the relation  $R$  that fulfills  $\theta'$ :

$$\bigcup_{C \in \mathcal{C}} C = \sigma_{\theta'}(R) \quad (1)$$

The statistics materialized for this cluster set  $C$  are then used to evaluate the AC for the repair candidate.

#### Algorithm 1 FULLCOVERCLUSTERSET

**Input:** kd-tree with root  $C_{root}$ , condition  $\theta' = \theta'_1 \wedge \dots \wedge \theta'_m$ , relation  $R$ .  
**Output:** Set of clusters  $C$  such that  $\bigcup_{C \in \mathcal{C}} C = \sigma_{\theta'}(R)$ .

```

1:  $stack \leftarrow [C_{root}]$ 
2:  $C \leftarrow \emptyset$   $\triangleright$  Initialize result set
3: while  $stack \neq \emptyset$  do
4:    $C_{cur} \leftarrow \text{pop}(stack)$ 
5:    $in \leftarrow \text{true}$ ,  $notin \leftarrow \text{false}$ 
6:   for all  $\theta'_i = (a_i \text{ op } c'_i) \in \theta'$  do
7:      $in \leftarrow in \wedge \text{eval}_V(\theta'_i, \text{BOUNDS}_{a_i}(C_{cur}))$   $\triangleright$  All tuples fulfill  $\theta'_i$ ?
8:      $notin \leftarrow notin \vee \text{eval}_V(\neg\theta'_i, \text{BOUNDS}_{a_i}(C_{cur}))$ 
9:   if  $in$  then  $\triangleright$  All tuple in  $C$  fulfill  $\theta'$ 
10:     $C \leftarrow C \cup \{C_{cur}\}$ 
11:   else if  $\neg notin$  then  $\triangleright$  Some tuples in  $C$  may fulfill  $\theta'$ 
12:     for all  $C \in \text{children}(C_{cur})$  do  $\triangleright$  Process children
13:        $stack \leftarrow stack \cup \{C\}$ 
14: return  $C$ 

```

**Table 1: Given the bounds  $[a, \bar{a}]$  for the attribute  $a$  of a condition  $a \text{ op } c$  or  $a \in [c_1, c_2]$ , function  $\text{eval}_V$  does return true if the condition evaluates to true for all values in  $[a, \bar{a}]$ . For RP, we consider a range  $[c, \bar{c}]$  (corresponding to a set of candidates) or two ranges  $[c_1, \bar{c}_1]$  and  $[c_2, \bar{c}_2]$  for operator  $\in$ .  $\text{reval}_V$  determines whether for every  $c \in [c, \bar{c}]$ , the condition is guaranteed to evaluate to true for every  $a \in [a, \bar{a}]$  while  $\text{reval}_\exists$  determines whether for some  $c \in [c, \bar{c}]$ , the condition may evaluate to true for  $a \in [a, \bar{a}]$ .**

Op.	$\text{eval}_V$	$\text{reval}_V$	$\text{reval}_\exists$
$>, \geq$	$\underline{a} > c, \quad \underline{a} \geq c$	$\underline{a} > \bar{c}, \quad \underline{a} \geq \bar{c}$	$\bar{a} > c, \quad \bar{a} \geq c$
$<, \leq$	$\bar{a} < c, \quad \bar{a} \leq c$	$\bar{a} < \underline{c}, \quad \bar{a} \leq \underline{c}$	$\underline{a} < \bar{c}, \quad \underline{a} \leq \bar{c}$
$=$	$\underline{a} = \bar{a} = c$	$\underline{a} = \underline{c} = \bar{a} = \bar{c}$	$[a, \bar{a}] \cap [c, \bar{c}] \neq \emptyset$
$\neq$	$c \notin [a, \bar{a}]$	$[a, \bar{a}] \cap [c, \bar{c}] = \emptyset$	$\neg(a = c = \bar{c} = \bar{a})$
$\in [c_1, c_2]$	$c_1 \leq \underline{a} \wedge \bar{a} \leq c_2$	$\bar{c}_1 \leq \underline{a} \wedge \bar{a} \leq \bar{c}_2$	$[a, \bar{a}] \cap [c_1, c_2] \neq \emptyset$

**3.2.1 Determining a Covering Set of Clusters.** The algorithm maintains a *stack* of clusters to be examined that is initialized with the root cluster  $C_{root}$  (line 1). It then processes one cluster at a time until a set of clusters  $C$  fulfilling Equation (1) has been determined (lines 3-14). For each cluster  $C$ , we distinguish 3 cases (lines 6-8): (i) we can use the bounds on the selection attributes recorded for the cluster to show that all tuples in the cluster fulfill  $\theta'$ , i.e.,  $\sigma_{\theta'}(C) = C$  (line 7). In this case, the cluster will be added to  $C$  (lines 9-10); (ii) based on the bounds, we can determine that none of the tuples in the cluster fulfill the condition (line 8). Then this cluster can be ignored; (iii) either a non-empty subset of  $C$  fulfills  $\theta'$  or based on the bounds  $\text{BOUNDS}_{a_i}(C)$  we cannot demonstrate that  $\sigma_{\theta'}(C) = \emptyset$  or  $\sigma_{\theta'}(C) = C$  hold. In this case, we add the children of  $C$  to the stack to be evaluated in future iterations (lines 11-13). The algorithm uses the function  $\text{eval}_V$  shown in Table 1 to determine whether based on the bounds of the cluster  $C$ , the comparison condition  $\theta'_i$  is guaranteed to be true for all  $t \in C$ . Additionally, it checks whether case (ii) holds by applying  $\text{eval}_V$  to the negation  $\neg\theta'_i$ . Note that to negate a comparison we simply push the negation to the comparison operator, e.g.,  $\neg(a < c) = (a \geq c)$ . As the selection condition of any repair candidate is a conjunction of comparisons  $\theta'_1 \wedge \dots \wedge \theta'_m$ , the cluster is *fully covered* (case (i)) if  $\text{eval}_V$  returns true for all  $\theta'_i$  and *not covered at all* (case (ii)) if  $\text{eval}_V$  returns true for at least one comparison  $\neg\theta'_i$ .



**3.2.2 Determining Coverage.** In Table 1, we define the function  $\text{eval}_\forall$  which takes a condition  $a \text{ op } c$  and bounds  $\text{BOUNDS}_a(C)$  for attribute  $a$  in cluster  $C$  and returns true if it is guaranteed that all tuples  $t \in C$  fulfill the condition. Ignore  $\text{reval}_\forall$  for now, this function will be used in Section 4. An inequality  $>$  (or  $\geq$ ) is true for all tuples if the lower bound  $\underline{a}$  of  $a$  is larger (larger equal) than the threshold  $c$ . The case for  $<$  and  $\leq$  is symmetric: the upper bound  $\bar{a}$  has to be smaller (smaller equals) than  $c$ . For an equality, we can only guarantee that the condition is true if  $\underline{a} = \bar{a} = c$ . For  $\neq$ , all tuples fulfill the inequality if  $c$  does not belong to the interval  $[\underline{a}, \bar{a}]$ .

For the running example in Figure 1, consider a repair candidate with the condition  $T \geq 34$ , where  $c_1 = 34$ . The algorithm maintains a stack of clusters initialized to  $[C_1]$ , the root node of the kd-tree. In each iteration it takes one cluster from the stack. The root cluster  $C_1$ , has  $\text{BOUNDS}_T(C_1) = [27, 37]$ . The algorithm evaluates whether all or none of the tuples satisfy the condition. Since neither is the case, we proceed to the children of  $C_1$ :  $\{C_2 \text{ and } C_3\}$ . The same situation occurs for  $C_2$  and  $C_3$  leading to further exploration of their children:  $C_4$  and  $C_5$  for  $C_2$  and  $C_8$  and  $C_9$  for  $C_3$ . Since the coverage for  $C_4$  cannot be determined, the algorithm proceeds to process  $C_6$  and  $C_7$ . Clusters  $C_5$ ,  $C_6$  and  $C_9$  are determined to not satisfy the condition while  $C_7$  and  $C_8$  are confirmed to meet the condition and are added to  $C$ . In this example, we had to explore all of the leaf clusters, but often we will be able to prune or confirm clusters covering multiple tuples. For instance, for  $T \geq 37$ ,  $C_2$  with bounds  $[27, 34]$  with all of its descendents can be skipped as  $T \geq 37$  is false for any  $T \in [27, 34]$ .

**3.2.3 Constraint Evaluation.** After identifying the covering set of clusters  $C$  for a repair candidate  $Q_{fix}$ , our approach evaluates the AC  $\omega$  over  $C$ . Recall that for each cluster  $C$  we materialize the result of each filter-aggregation query  $Q_i^\omega$  used in  $\omega$ . For aggregate function **avg** that is not decomposable, we apply the standard approach of storing **count** and **sum** instead. We then compute  $Q_i^\omega(Q(D))$  over the materialized aggregation results for the clusters. Concretely, for such an aggregate query  $Q^\omega := \gamma_{f(a)}(\sigma_{\theta'}(Q(D)))$  we compute its result as follows using  $C$ :  $\gamma_{f'(a)}(\bigcup_{C \in C} \{Q^\omega(C)\})$ . Here  $f'$  is the function we use to merge aggregation results for multiple subsets of the database. This function depends on  $f$ , e.g., for both **count** and **sum** we have  $f' = \text{sum}$ , for **min** we use  $f' = \text{min}$ , and for **max** we use  $f' = \text{max}$ . We then substitute these aggregation results into  $\omega$  and evaluate the resulting expression to determine whether  $Q_{fix}$  fulfills the constraints and is a repair or not.

In the example from Figure 1(c), the covering set of clusters for the repair candidate with  $c_1 = 34$  is  $C = \{C_7, C_8\}$ . Evaluating  $Q_1^\omega = \text{count}(\text{Gender}(G) = M \wedge Y = 1)$  over  $C$ , we sum the counts:  $Q_1^\omega = Q_{1C_7}^\omega + Q_{1C_8}^\omega = 1 + 1 = 2$ . Similarly,  $Q_2^\omega = Q_{2C_7}^\omega + Q_{2C_8}^\omega = 1 + 1 = 2$ ,  $Q_3^\omega = Q_{3C_7}^\omega + Q_{3C_8}^\omega = 0 + 0 = 0$ ,  $Q_4^\omega = Q_{4C_7}^\omega + Q_{4C_8}^\omega = 0 + 0 = 0$  as shown in Figure 1(d). Substituting these values into  $\omega_\#$ , we obtain  $1 \leq 0.2 = \text{false}$  as shown in Figure 1(e). Since the candidate  $T \geq 34$  does not satisfy the constraint it is not a valid repair.

### 3.3 Computing Top- $k$ Repairs

To compute the top- $k$  repairs, we enumerate all repair candidates in increasing order of their distance to the user query using the distance measure from Section 2. For each candidate  $Q_{fix}$ , we apply Algorithm FF to determine a covering clusterset, evaluate the

---

#### Algorithm 2 PARCOVERCLUSTERSET

---

**Input:** kd-tree with root  $C_{root}$ , repair candidate set  $\mathbb{Q} = [[c_1, \bar{c}_1], \dots, [c_m, \bar{c}_m]]$ , condition  $\theta$   
**Output:** Partially covering cluster set  $(C_{full}, C_{partial})$

```

1:  $stack \leftarrow [C_{root}]$ 
2:  $C_{full} \leftarrow \emptyset, C_{partial} \leftarrow \emptyset$   $\triangleright$  Initialize cluster sets
3: while  $stack \neq \emptyset$  do
4:    $C_{cur} \leftarrow \text{pop}(stack)$ 
5:    $in \leftarrow \text{true}, pin \leftarrow \text{true}$ 
6:   for all  $\theta_i = (a_i \text{ op } c_i) \in \theta$  do  $\triangleright C_{cur}$  fully / part. covered?
7:      $in \leftarrow in \wedge \text{reval}_\forall(\theta_i, [c_i, \bar{c}_i], \text{BOUNDS}_{a_i}(C_{cur}))$ 
8:      $pin \leftarrow pin \wedge \text{reval}_\exists(\theta_i, [c_i, \bar{c}_i], \text{BOUNDS}_{a_i}(C_{cur}))$ 
9:   if in then  $\triangleright$  Add fully covered cluster to the result
10:     $C_{full} \leftarrow C_{full} \cup \{C_{cur}\}$ 
11:   else if pin then
12:     if isleaf $(C_{cur})$  then  $\triangleright$  Partially covered leaf cluster
13:        $C_{partial} \leftarrow C_{partial} \cup \{C_{cur}\}$ 
14:     else  $\triangleright$  Process children of partial cluster
15:       for all  $C \in \text{children}(C_{cur})$  do
16:          $stack \leftarrow stack \cup \{C\}$ 
17: return  $(C_{full}, C_{partial})$ 
```

---

constraint  $\omega$ , and output  $Q_{fix}$  if it fulfills the constraint. Once we have found  $k$  results, the algorithm terminates.

## 4 CLUSTER RANGE PRUNING (RP)

While algorithm FF reduces the effort needed to evaluate aggregation constraints for repair candidates, it has the drawback that we still have to evaluate each repair candidate individually. We now present an enhanced approach that reasons about sets of repair candidates. For a user query condition  $\theta_1 \wedge \dots \wedge \theta_m$  where  $\theta_i := a_i \text{ op } c_i$ , we use ranges of constant values instead of constants to represent such a set of repairs  $\mathbb{Q}$ :  $[[c_1, \bar{c}_1], \dots, [c_m, \bar{c}_m]]$ . Such a list of ranges  $\mathbb{Q}$  represents a set of a repair candidates:

$$\{[c_1, \dots, c_m] \mid \forall i \in [1, m] : c_i \in [c_i, \bar{c}_i]\}$$

Consider an aggregation constraint  $\omega := \tau \text{ op } \Phi(Q_1^\omega, \dots, Q_n^\omega)$ . Our enhanced approach RP uses a modified version of the kd-tree from FF to compute conservative bounds  $\Phi$  and  $\bar{\Phi}$  on the possible values of arithmetic expression  $\Phi$  that hold for all repair candidates in  $\mathbb{Q}$ . Based on such bounds, if (i)  $\tau \text{ op } c$  holds for every  $c \in [\Phi, \bar{\Phi}]$ , then every  $Q_{fix} \in \mathbb{Q}$  is a valid repair, if (ii)  $\tau \text{ op } c$  is violated for every  $c \in [\Phi, \bar{\Phi}]$ , then no  $Q_{fix} \in \mathbb{Q}$  is a valid repair and we can skip the whole set. Otherwise, (iii) there may or may not exist some candidates in  $\mathbb{Q}$  that are repairs. In this case, our algorithm partitions  $\mathbb{Q}$  into multiple subsets and applies the same test to each partition. In the following, we first discuss our algorithm that utilizes such repair candidate sets and bounds on the aggregate constraint results and then explain how to use the kd-tree to compute such bounds.

### 4.1 Computing Top- $k$ Repairs

RP (Algorithm 3) takes as input a kd-tree with root  $C_{root}$ , a user query's condition  $\theta$ , an AC  $\omega$ , a candidate set  $\mathbb{Q} = [[c_1, \bar{c}_1], \dots, [c_m, \bar{c}_m]]$ , and a user query  $Q$  and returns the set of top- $k$  repairs  $Q_{top-k}$ .

---

**Algorithm 3** Top-k Repairs w. Range-based Pruning of Candidates

**Input:** kd-tree with root  $C_{root}$ , constraint AC  $\omega$ , repair candidate set  $\mathbb{Q} = [[c_1, \bar{c}_1], \dots, [c_m, \bar{c}_m]]$ , user query condition  $\theta = \theta_1 \wedge \dots \wedge \theta_m$ , user query  $Q$

**Output:** Top-k repairs  $Q_{top-k}$

```

1:  $Q_{top-k} \leftarrow \emptyset$   $\triangleright$  Queue of repairs  $Q'$  sorted on  $d(Q, Q')$ 
2:  $rcand \leftarrow \emptyset$   $\triangleright$  Queue of repair sets  $Q'$  sorted on  $d(Q, Q')$ 
3:  $queue \leftarrow [\mathbb{Q}]$   $\triangleright$  Queue of repair candidate sets  $Q'$  sorted on  $d(Q, Q')$ 
4: while  $queue \neq \emptyset$  do
5:    $Q_{cur} \leftarrow \text{POP}(queue)$ 
6:    $Q_{next} \leftarrow \text{PEEK}(queue)$   $\triangleright$  Peek at next item in queue
7:    $(C_{full}, C_{partial}) \leftarrow \text{PARCOVERCLUSTERSET}(Q_{cur}, C_{root}, \theta)$ 
8:   if  $\text{ACEVAL}_V(\omega, C_{full}, C_{partial})$  then  $\triangleright$  All  $Q' \in Q_{cur}$  are repairs?
9:      $rcand \leftarrow \text{INSERT}(rcand, Q_{cur})$ 
10:  else if  $\text{ACEVAL}_\exists(\omega, C_{full}, C_{partial})$  then  $\triangleright$  Potential repairs?
11:    for  $Q_{new} \in \text{RANGEDIVIDE}(Q_{cur})$  do  $\triangleright$  Divide ranges
12:      if  $\text{HASCANDIDATES}(Q_{new})$  then
13:         $queue \leftarrow \text{INSERT}(queue, Q_{new})$ 
14:   $Q_{top-k} \leftarrow \text{TOPKCONCRETECAND}(rcand, k)$   $\triangleright$  Top k repairs
15:  if  $|Q_{top-k}| \geq k$  then  $\triangleright$  Have k repairs?
16:    if  $d(Q, Q_{next}) > d(Q, Q_{top-k}[k])$  then  $\triangleright$  Rest inferior?
17:      break
18: return  $Q_{top-k}$ 

```

---

The algorithm maintains three priority queues: (i)  $Q_{top-k}$  is a queue of individual repairs that eventually will store the top-k repairs. This queue is sorted on  $d(Q, Q_{fix})$  where  $Q_{fix}$  is a repair in the queue; (ii)  $rcand$  is a queue where each element is a repair candidate set  $\mathbb{Q}$  encoded as ranges as shown above. For each  $\mathbb{Q}$  we have established that for all  $Q_{fix} \in \mathbb{Q}$ ,  $Q_{fix}$  is a repair. Queue  $rcand$  is sorted on the lower bound  $d(Q, Q_{fix} \in \mathbb{Q})$  of the distance of any repair in  $\mathbb{Q}$  to the user query. Finally, (iii)  $queue$  is a queue where each element is a repair candidate set  $\mathbb{Q}$ . This queue is also sorted on  $d(Q, Q_{fix} \in \mathbb{Q})$ . In each iteration of the main loop of the algorithm, one repair candidate set from  $queue$  is processed.

The algorithm initializes  $queue$  to the input repair candidate set  $\mathbb{Q}$ . We call the algorithm with a repair candidate set that covers the whole search space (line 1-3). The algorithm's main loop processes one repair candidate  $Q_{cur}$  at a time (line 5) while keeping track of the next candidate  $Q_{next}$  (line 6) until a set of top-k repairs fulfilling AC  $\omega$  has been determined (lines 4–17). For the current repair candidate set  $Q_{cur}$ , we use function  $\text{PARCOVERCLUSTERSET}$  (Algorithm 2) to determine two sets of clusters  $C_{full}$  and  $C_{partial}$  (line 7). For every cluster  $C \in C_{full}$ , all tuples in  $C$  fulfill the condition of every repair candidate  $Q_{fix} \in Q_{cur}$  and for every cluster  $C \in C_{partial}$ , there may exist some tuples in  $C$  such that for some repair candidates  $Q_{fix} \in Q_{cur}$ , the tuples fulfill the condition of  $Q_{fix}$ . We use these two sets of clusters to determine bounds  $[\Phi, \bar{\Phi}]$  on the arithmetic expression  $\Phi$  of the AC  $\omega$ . The algorithm then uses these bounds to distinguish between three cases (line 8-13): (i)  $\omega$  is guaranteed to hold for every  $Q_{fix} \in Q_{cur}$ . We use function  $\text{ACEVAL}_V$  which takes  $C_{full}$  and  $C_{partial}$  as input to test for this case. If case (i) applies then all  $Q_{fix} \in Q_{cur}$  are repairs and we add  $Q_{cur}$  to  $rcand$  (lines 8-9); (ii) some repair candidates  $Q_{fix} \in Q_{cur}$  may fulfill the AC. In this case  $Q_{cur}$  will be split and further examined in future iterations (lines 10–13). We test for case (ii) using function  $\text{ACEVAL}_\exists$ ; (iii) no

$Q_{fix} \in Q_{cur}$  is a repair and we can discard  $Q_{cur}$ . Case (iii) applies if both  $\text{ACEVAL}_V$  and  $\text{ACEVAL}_\exists$  return false. We will discuss these functions in depth in Section 4.3.

For example, if  $\omega := 0.7 \leq \Phi$  and we compute bounds  $[\Phi, \bar{\Phi}] = [0.5, 1]$  that hold for all  $Q_{fix} \in Q_{cur}$ , then  $\text{ACEVAL}_V$  returns false as some  $Q_{fix} \in Q_{cur}$  may not fulfill the constraint. However,  $\text{ACEVAL}_\exists$  return true as some  $Q_{fix} \in Q_{cur}$  may fulfill the constraint. In this case, the algorithm partitions  $Q_{cur}$  into smaller sub-ranges  $Q_{new}$  using the function  $\text{RANGEDIVIDE}(Q_{cur})$  (line 11). Assume that  $Q_{cur} = [[c_1, \bar{c}_1], \dots, [c_m, \bar{c}_m]]$ .  $\text{RANGEDIVIDE}$  splits each range  $[c_i, \bar{c}_i]$  into a fixed number of fragments  $\{[c_{i_1}, \bar{c}_{i_1}], \dots, [c_{i_l}, \bar{c}_{i_l}]\}$  such that each  $[c_{i_j}, \bar{c}_{i_j}]$  is roughly of the same size and returns the following set of repair candidate sets:

$$\{[[c_{1j_1}, \bar{c}_{1j_1}], \dots, [c_{mj_m}, \bar{c}_{mj_m}]] \mid [j_1, \dots, j_m] \in [1, l]^m\}$$

That is, each  $Q_{new}$  has one of the fragments for each  $[c_i, \bar{c}_i]$  and the union of all repair candidates in these repair candidate sets is  $Q_{cur}$ . We use  $l = 2$  in our implementation. The function  $\text{HASCANDIDATES}$  (line 12-13) checks whether each range in  $Q_{new}$  contains at least one value that exists in the data. This restricts the search space to only include candidates with constants that actually appear in the data (or are the minimum / maximum of an attribute's domain). The rationale for this test is that for any other value  $c'$  not in this set and supported comparison operators, there exists a value  $c$  in the set (occurring in the data or a domain bound) for which the selection on  $c$  and  $c'$  returns exactly the same result. Hence, it is sufficient to only consider such repair candidates and repair sets that do not contain such a candidate can be pruned. For example, if the dataset contains only values 8 and 10 for a given attribute, then applying a filter  $a \leq 9$  would yield the same result as  $a \leq 8$ , since no data points lie between 8 and 10. If this condition is satisfied,  $Q_{new}$  is inserted into the priority queue  $queue$  to be processed in future iterations. In each iteration we use function  $\text{TOPKCONCRETECAND}$  (line 14) to determine the  $k$  repairs  $Q_i$  across all  $\mathbb{Q} \in rcand$  with the lowest distance to the user query  $Q$ . If we can find  $k$  such candidates (line 15), then we test whether no repair candidate from the next repair candidate set  $Q_{next}$  may be closer to  $Q$  than the  $k$ th candidate  $Q_{top-k}[k]$  from  $Q_{top-k}$  (line 16). This is the case if the lower bound on the distance of any candidate in  $Q_{next}$  is larger than the distance of  $Q_{top-k}[k]$ . Furthermore, the same holds for all the remaining repair candidate sets in  $rcand$ , because  $rcand$  is sorted on the lower bound of the distance to the user query. That is,  $Q_{top-k}$  contains exactly the top-k repairs and the algorithm returns this set.

## 4.2 Determining Covering Cluster Sets

Similar to FF, we can use the kd-tree to determine a covering cluster set  $C$ . However, as we now deal with a set of candidate repairs  $\mathbb{Q}$ , we would have to find a cluster set  $C$  such that for all  $Q_{fix} \in \mathbb{Q}$  we have:  $Q_{fix}(D) = \bigcup_{C \in C} C$ . Such a covering cluster set is unlikely to exist as for any two  $Q_{fix} \neq Q'_{fix} \in \mathbb{Q}$  it is likely that  $Q_{fix}(D) \neq Q'_{fix}(D)$ . Instead we relax the condition and allow clusters  $C$  that are *partially covered*, i.e., for which some tuples in  $C$  may be in the result of some candidates in  $\mathbb{Q}$ . We modify Algorithm 1 to take a repair candidate set as an input and to return two sets of clusters:  $C_{full}$  which contains clusters for which all tuples fulfill the selection condition of all  $Q_{fix} \in \mathbb{Q}$  and  $C_{partial}$  which contains clusters that are only

partially covered, i.e., may contain tuples that fulfill the condition of some  $Q_{fix} \in \mathbb{Q}$ .

Analogous to Algorithm 1, the updated algorithm (Algorithm 2) maintains a stack of clusters to be processed that is initialized with the root node of the kd-tree (line 1). In each iteration of the main loop (line 3-16), the algorithm determines whether all tuples of the current cluster  $C_{cur}$  fulfill the conditions  $\theta_i$  for all repair candidates  $Q_{fix} \in \mathbb{Q}$ . This is done using function  $\text{reval}_\forall$  (line 7). Additionally, we check whether it is possible that at least one tuple fulfills the condition of at least one repair candidate  $Q_{fix} \in \mathbb{Q}$ . This is done using a function  $\text{reval}_\exists$  (line 8). If the cluster is fully covered we add it to the result set  $C_{full}$  (line 10). If it is partially covered, then we distinguish between two cases (line 11- 16). Either the cluster is a leaf node (line 12-13) or it is an inner node (line 14-16). If the cluster is a leaf, then we cannot further divide the cluster and add it to  $C_{partial}$ . If the cluster is an inner node, then we process its children as we may be able to determine that some of its children are fully covered or not covered at all.

Table 1 shows how conditions are evaluated by  $\text{reval}_\forall$  and  $\text{reval}_\exists$ . For a condition  $a > c$ , if the lower bound of attribute  $a$  is larger than the upper bound  $\bar{c}$ , then all tuples in the cluster fulfill the condition for all  $Q_{fix} \in \mathbb{Q}$ . The cluster is partially covered if  $\bar{a} > c$  as then there exists at least one value in the range of  $a$  and constant  $c$  in  $[c, \bar{c}]$  for which the condition is true.

In the example from Figure 1, a repair candidate  $[[33, 37]]$  is evaluated. Recall that the single condition in this example is  $T \geq c$ .  $C_{root}$  has  $\text{BOUNDS}_T = [27, 37]$ . The algorithm first applies  $\text{reval}_\forall$  to check if all tuples in  $C_{root}$  satisfy the condition. Since  $27 \not\geq 37$ , the algorithm proceeds to evaluate the condition for partial coverage using  $\text{reval}_\exists$ . Since  $C_1$  is partially covered and not a leaf, the algorithm continues by processing  $C_1$ 's children,  $C_2$  and  $C_3$ . For  $C_3$ , a similar situation occurs: the lower bound of the attribute,  $a = 31$ , is not greater than the upper bound of the constant,  $\bar{c} = 37$  and we have to process additional clusters,  $C_8$  and  $C_9$ . The same holds for  $C_2$  and we process its children:  $C_4$  and  $C_5$ . Additionally,  $C_4$  fails  $\text{reval}_\forall$  but satisfies partial coverage with  $\text{reval}_\exists$ , necessitating evaluation of its children,  $C_6$  and  $C_7$ . Finally, the algorithm applies  $\text{reval}_\forall$  and  $\text{reval}_\exists$  if necessary to the clusters  $C_5$ ,  $C_6$ ,  $C_7$ ,  $C_8$ , and  $C_9$ , confirming that  $C_8 \in C_{full}$  and  $C_7 \in C_{partial}$ , as  $t_3.T = 37 \geq c$  is true for all  $c \in [33, 37]$  and  $t_4.T = 34 \geq c$  is may be true for some  $c \in [33, 37]$ .

### 4.3 Computing Bounds on Constraints

Given the cluster sets  $(C_{full}, C_{partial})$  computed by Algorithm 2, we next (i) compute bounds on the results of the aggregation queries  $Q_i^\omega$  used in the constraint, then (ii) use these bounds to compute bounds  $[\Phi, \bar{\Phi}]$  on the result of the arithmetic expression  $\Phi$  of the AC  $\omega$  over repair candidates in  $\mathbb{Q}$ . These bounds are conservative in the sense that all possible results are guaranteed to be included in these bounds. Then, finally, (iii) function  $\text{ACEVAL}_\forall$  uses the computed bounds to determine whether all candidates in  $\mathbb{Q}$  fulfill the constraint by applying  $\text{reval}_\forall$  from Table 1. For a constraint  $\omega := \tau \text{ op } \Phi$ ,  $\text{ACEVAL}_\forall$  calls  $\text{reval}_\forall$  with  $[\Phi, \bar{\Phi}]$  and  $\tau$ .  $\text{ACEVAL}_\exists$  uses  $\text{reval}_\exists$  instead to determine whether some candidates in  $\mathbb{Q}$  may fulfill the constraint. This requires techniques for computing bounds on the possible results of arithmetic expressions and aggregation functions

when the values of each input of the computation are known to be bounded by some interval.

**4.3.1 Bounding Aggregation Results.** We now discuss how to compute bounds for the results of the filter-aggregation queries  $Q_i^\omega$  of an aggregate constraint  $\omega$  based on the cluster sets  $(C_{full}, C_{partial})$  returned by Algorithm 2. As every cluster  $C$  in  $C_{full}$  is fully covered for all repair candidates in  $\mathbb{Q}$ , i.e., all tuples in the cluster fulfill the conditions of each  $Q_{fix} \in \mathbb{Q}$ , the materialized aggregation results  $Q_i^\omega(C)$  of  $C$  contribute to both the lower bound  $\underline{Q}_i^\omega$  and upper bound  $\bar{Q}_i^\omega$  as for FF. For partially covered clusters ( $C_{partial}$ ), we have to make worst case assumptions to derive valid lower and upper bounds. For the lower bound, we have to consider the minimum across two options: (i) no tuples from the cluster can lower the aggregation value if included in the computation. In this case, the cluster is ignored for computing the lower bound, e.g., this is always the case for **max**; (ii) based on the bounds of the input attribute for the aggregation within the cluster, there are values in the cluster that if added to the current aggregation result further lowers the result. For example, for **min**( $a$ ) we have to include  $\underline{a}$  if  $\underline{a}$  is smaller than the current smallest value of **min**. For **sum** we have the two cases: (i) the attribute for the aggregation has negative numbers. In this case we multiply  $\underline{a}$  with the **count** for the cluster (a lower bound on summing up all negative numbers in the cluster) and add this to the lower bound;<sup>2</sup> (ii) otherwise we should ignore this cluster for computing lower bounds. Computing the upper bound is symmetric: if there are no tuples in the cluster that would result in a larger aggregation result, e.g., for **sum** when all values of attribute  $a$  in the cluster are negative, then including any tuple from the cluster would lower the aggregation result and the cluster should be ignored. Otherwise, including some values in the cluster for the aggregation input attribute may increase the aggregation result, then we include  $\bar{a}$  (aggregation function is **min** or **max**) or  $\bar{a} \cdot |C|$  (aggregation function is **sum**) in the computation for the upper bound.

**4.3.2 Bounding Results of Arithmetic Expressions.** Given the bounds on filter-aggregation queries, we use *interval arithmetic* [17, 33] which computes sound bounds for the result of arithmetic operations when the inputs are bound by intervals. In our case, the bounds on the results of aggregate queries  $Q_i^\omega$  are the input and bounds  $[\Phi, \bar{\Phi}]$  on  $\Phi$  are the result. The notation we use is similar to [39]. Table 2 shows the definitions for arithmetic operators we support in aggregate constraints. Here,  $\underline{E}$  and  $\bar{E}$  denote the lower and upper bound on the values of expression  $E$ , respectively. For example, for addition the lower bound for the result of addition  $E_1 + E_2$  of two expressions  $E_1$  and  $E_2$  is  $\underline{E}_1 + \underline{E}_2$ .

**4.3.3 Bounding Aggregate Constraint Results.** Consider a constraint  $\omega := \tau \text{ op } \Phi$ . There are three possible outcomes for a repair candidate set: (i)  $\tau \text{ op } \Phi$  is true for all  $[\Phi, \bar{\Phi}]$  which  $\text{ACEVAL}_\forall$  determines using  $\text{reval}_\forall$  and bounds  $[\tau, \tau]$ ; (ii) some of the candidates in  $\mathbb{Q}$  may fulfill the condition, which  $\text{ACEVAL}_\exists$  determines using  $\text{reval}_\exists$ ; (iii)

<sup>2</sup>These bounds can be improved if for each filter-aggregation query with aggregation function **sum** we additionally materialize counts and separate sums for negative and positive values in a cluster for tuples that fulfill the condition of the filter-aggregation query.

none of the candidates in  $\mathbb{Q}$  fulfill the condition (both (i) and (ii) are false).

In the running example from Figure 1(g), the covering set of clusters for repair candidate set  $\mathbb{Q} := [[33, 37]]$  are  $C_{full} = \{C_8\}$  and  $C_{partial} = \{C_7\}$ . To evaluate  $Q_1^\omega = \text{count}(G = M \wedge Y = 1)$  over these clusters, the algorithm includes the materialized aggregation results for  $C_8$  for both the lower bound  $\underline{Q}_i^\omega$  and upper bound  $\overline{Q}_i^\omega$ . For the partially covered  $C_7$ , the lower bound of  $Q_{1C_7}^\omega$  is 0 for this cluster (the lowest count is achieved by excluding all tuples from the cluster), while the upper bound is 1, as there exists a male in the cluster satisfying  $Y = 1$ . Thus, we get the following bounds for  $Q_{1C_7}^\omega = [0, 1]$ . Similarly, we compute the remaining aggregation bounds:  $Q_{1C_8}^\omega = [1, 1]$ ,  $Q_{2C_7}^\omega = [0, 1]$ ,  $Q_{2C_8}^\omega = [1, 1]$ ,  $Q_{3C_7}^\omega = [0, 0]$ ,  $Q_{3C_8}^\omega = [0, 0]$ ,  $Q_{4C_7}^\omega = [0, 0]$ ,  $Q_{4C_8}^\omega = [0, 0]$ .

Next, in Figure 1(h) we sum the lower and upper bounds for each aggregation  $Q_i^\omega$  across all clusters in  $\mathbb{C}$ :  $Q_1^\omega = Q_{1C_7}^\omega + Q_{1C_8}^\omega = [1, 2]$ ,  $Q_2^\omega = Q_{2C_7}^\omega + Q_{2C_8}^\omega = [1, 2]$ ,  $Q_3^\omega = Q_{3C_7}^\omega + Q_{3C_8}^\omega = [0, 0]$ ,  $Q_4^\omega = Q_{4C_7}^\omega + Q_{4C_8}^\omega = [0, 0]$ . We then substitute the computed values  $\{Q_1^\omega, Q_2^\omega, Q_3^\omega, Q_4^\omega\}$  into  $\omega_\#$  and evaluate the resulting expression using interval arithmetic (Table 2). Given:  $\omega_\# = \underline{Q}_1^\omega / \underline{Q}_2^\omega - \underline{Q}_3^\omega / \underline{Q}_4^\omega$  the lower and upper bounds for the first term  $\underline{Q}_1^\omega / \underline{Q}_2^\omega$  are computed as:  $[\underline{E}_1 / \underline{E}_2, \overline{E}_1 / \overline{E}_2] = [1/2, 2]$ . Similarly, for the second term:  $\underline{Q}_3^\omega / \underline{Q}_4^\omega = [0, 0]$ . Applying interval arithmetic to compute the subtraction we get:  $\underline{E}_1 - \underline{E}_2, \overline{E}_1 - \overline{E}_2$ . Thus, we obtain bounds  $[\Phi_\#, \overline{\Phi}_\#] = [1/2, 2]$  (Figure 1(i)). Since  $\Phi_\# = 1/2 > 0.2$ , none of the candidates in  $\mathbb{Q} = [[33, 37]]$  can be repairs and we can prune  $\mathbb{Q}$ .

In practice, RP performs best when kd-tree clusters are homogeneous with respect to the predicate attributes  $a_i$  in  $\theta$  i.e., when most cluster bounds  $\text{BOUNDS}_{a_i}$  lie entirely above or below the repair candidate set's intervals  $[c_i, \overline{c}_i]$ . This enables efficient pruning of infeasible or fully satisfying candidate sets. The effect is especially strong when large regions of the search space can be ruled out or accepted entirely based on the aggregation constraint  $\omega$ . If predicate attributes are strongly correlated with those in the arithmetic expression  $\Phi$ , cluster inclusion often predicts the outcome of  $\Phi$ , allowing entire repair sets to be evaluated at once. Conversely, when many cluster bounds partially overlap with the predicate ranges, the algorithm must recursively partition  $\mathbb{Q}$  and evaluate more finer-grained clusters, eventually matching the cost of brute-force in the worst case.

**THEOREM 4.1 (CORRECTNESS OF FF AND RP).** *Given an instance  $(Q, D, \omega, k)$  of the aggregate constraint repair problem, Algorithm 3 computes the solution for this problem instance.*

**PROOF.** We present the proof in [3].  $\square$

## 5 EXPERIMENTS

We start by comparing the brute force approach and the baseline FF technique (Section 3.2) against our RP algorithm (Section 4) in Section 5.2. We then investigate the impact of several factors on performance in Section 5.3, including dataset size and similarity of the top-k repairs to the user query. Finally, in Section 5.4, we compare with *Erica* [26] which targets group cardinality constraints.

**Table 2: Bounds on applying an operator to the result of expressions  $E_1$  and  $E_2$  with interval bounds [39].**

op	Bounds for the expression $(E_1 \text{ op } E_2)$
+	$\underline{E}_1 + \underline{E}_2 = \underline{E}_1 + \underline{E}_2$ $\overline{E}_1 + \overline{E}_2 = \overline{E}_1 + \overline{E}_2$
-	$\underline{E}_1 - \underline{E}_2 = \underline{E}_1 - \underline{E}_2$ $\overline{E}_1 - \overline{E}_2 = \overline{E}_1 - \overline{E}_2$
$\times$	$\underline{E}_1 \times \underline{E}_2 = \min(\underline{E}_1 \times \underline{E}_2, \underline{E}_1 \times \overline{E}_2, \overline{E}_1 \times \underline{E}_2, \overline{E}_1 \times \overline{E}_2)$ $\overline{E}_1 \times \overline{E}_2 = \max(\underline{E}_1 \times \underline{E}_2, \underline{E}_1 \times \overline{E}_2, \overline{E}_1 \times \underline{E}_2, \overline{E}_1 \times \overline{E}_2)$
/	$\underline{E}_1 / \underline{E}_2 = \min(\underline{E}_1 / \underline{E}_2, \underline{E}_1 / \overline{E}_2, \overline{E}_1 / \underline{E}_2, \overline{E}_1 / \overline{E}_2)$ $\overline{E}_1 / \overline{E}_2 = \max(\underline{E}_1 / \underline{E}_2, \underline{E}_1 / \overline{E}_2, \overline{E}_1 / \underline{E}_2, \overline{E}_1 / \overline{E}_2)$

### 5.1 Experimental Setup

**Datasets.** We choose two real-world datasets of size 500K, *Adult Census Income (Income)* [21] and *Healthcare* [22], that are commonly used to evaluate fairness. We also utilize the *TPC-H* [34] benchmark, varying dataset size from 25K to 500K.

**Queries.** Table 3 shows the queries used in our experiments. For Healthcare dataset (Healthcare), we use queries  $Q_1$  and  $Q_2$  from [26] and a new query  $Q_3$ . For Income, we use  $Q_4$  from [26] and new queries  $Q_5$  and  $Q_6$ .  $Q_7$  is a query with 3 predicates inspired by TPC-H's  $Q_2$ .

**Constraints.** For Healthcare and Income, we enforce the SPD between two demographic groups to be within a certain bound. Table 4 shows the details of the constraints used. In some experiments, we vary the bounds  $B_l$  and  $B_u$ . For a constraint  $\omega_i$  we use  $\omega_i^{d=p}$  to denote a variant of  $\omega_i$  where the bounds have been set such that the top-k repairs are within the first p% of the repair candidates ordered by distance. An algorithm that explores the individual repair candidates in this order has to explore the first p% of the search space. For the detailed settings see [3]. For Income (Healthcare), we determine the groups for SPD based on gender and race (race and age). For TPC-H, we enforce the constraint  $\omega_5$  which minimizes the impact of supply change disruption, where the company wants only a certain amount of expected revenue to be from certain countries. We use  $\Omega$  to denote a set of ACs.  $\Omega_6$  through  $\Omega_8$  are sets of cardinality constraints for comparison with Erica. As mentioned in Section 2, we present our repair methods for a single AC. However, the methods can be trivially extended to find repairs for a set / conjunction of constraints, i.e., the repair fulfills  $\bigwedge_{\omega \in \Omega} \omega$ . For RP (Algorithm 3), it is sufficient to replace the condition in Line 8 with  $\bigwedge_{\omega \in \Omega} \text{ACEVAL}_\forall(\omega, C_{full}, C_{partial})$  and in Line 10 with  $\bigwedge_{\omega \in \Omega} \text{ACEVAL}_\exists(\omega, C_{full}, C_{partial})$ .

**Parameters.** Recall that we use a kd-tree to perform the clustering described in Section 3.1. We consider two tuning parameters for the tree: **branching factor** - each node has  $\mathcal{B}$  children; **bucket size** - parameter  $\mathcal{S}$  determines the minimum number of tuples in a cluster. We do not split nodes with  $\leq \mathcal{S}$  tuples. When one of our algorithms reach such a leaf node we just evaluate computations over all tuples in the cluster, e.g., to determine which tuples fulfill a condition. We also control  $k$ , the number of repairs returned by our methods. The default settings are as follows:  $\mathcal{B} = 5$ ,  $k = 7$ , and  $\mathcal{S} = 15$ . The default dataset size is 50K tuples.

All algorithms were implemented in Python. Experiments were conducted on a machine with 2 x 3.3Ghz AMD Opteron CPUs (12



**Table 3: Queries for Experimentation**

SELECT * FROM Healthcare	
Q <sub>1</sub>	WHERE income >= 200K AND num-children >= 3 AND county <= 3
Q <sub>2</sub>	WHERE income <= 100K AND complications >= 5 AND num-children >= 4
Q <sub>3</sub>	WHERE income >= 300K AND complications >= 5 AND county == 1
SELECT * FROM ACSIncome	
Q <sub>4</sub>	WHERE working_hours >= 40 AND Educational_attainment >= 19 AND Class_of_worker >= 3
Q <sub>5</sub>	WHERE working_hours <= 40 AND Educational_attainment <= 19 AND Class_of_worker <= 4
Q <sub>6</sub>	WHERE Age >= 35 AND Class_of_worker >= 2 AND Educational_attainment <= 15
Q <sub>7</sub>	SELECT * FROM part, supplier, partsupp, nation, region WHERE p_partkey = ps_partkey AND s_suppkey = ps_suppkey AND s_nationkey = n_nationkey AND n_regionkey=r_regionkey AND p_size >= 10 AND p_type in ('LARGE_BRUSHED') AND r_name in ('EUROPE')

**Table 4: Constraints for Experimentation**

ID	Constraint
$\omega_1$	$\frac{\text{count}(\text{race}=1 \wedge \text{label}=1)}{\text{count}(\text{race}=1)} - \frac{\text{count}(\text{race}=2 \wedge \text{label}=1)}{\text{count}(\text{race}=2)} \in [B_l, B_u]$
$\omega_2$	$\frac{\text{count}(\text{ageGroup}=1 \wedge \text{label}=1)}{\text{count}(\text{ageGroup}=1)} - \frac{\text{count}(\text{ageGroup}=2 \wedge \text{label}=1)}{\text{count}(\text{ageGroup}=2)} \in [B_l, B_u]$
$\omega_3$	$\frac{\text{count}(\text{sex}=1 \wedge \text{PINCP} \geq 20k)}{\text{count}(\text{sex}=1)} - \frac{\text{count}(\text{sex}=2 \wedge \text{PINCP} \geq 20k)}{\text{count}(\text{sex}=2)} \in [B_l, B_u]$
$\omega_4$	$\frac{\text{count}(\text{race}=1 \wedge \text{PINCP} \geq 15k)}{\text{count}(\text{race}=1)} - \frac{\text{count}(\text{race}=2 \wedge \text{PINCP} \geq 15k)}{\text{count}(\text{race}=2)} \in [B_l, B_u]$
$\omega_5$	$\frac{\sum \text{Revenue} \text{ productsSelectedFromUK}}{\sum \text{Revenue} \text{ Selected Products}} \in [B_l, B_u]$
$\Omega_6$	$\omega_{61} := \text{count}(\text{race} = \text{race1}) \leq B_{u1}$ $\omega_{62} := \text{count}(\text{age} = \text{group1}) \leq B_{u2}$
$\Omega_7$	$\omega_{71} := \text{count}(\text{race} = \text{race1}) \leq B_{u1}$ $\omega_{72} := \text{count}(\text{age} = \text{group1}) \leq B_{u2}$ $\omega_{73} := \text{count}(\text{age} = \text{group3}) \leq B_{u3}$
$\Omega_8$	$\omega_{81} := \text{count}(\text{Sex} = \text{Female}) \leq B_{u1}$ $\omega_{82} := \text{count}(\text{Race} = \text{Black}) \leq B_{u2}$ $\omega_{83} := \text{count}(\text{Marital} = \text{Divorced}) \leq B_{u3}$

cores) and 128GB RAM. Each experiment was repeated five times and we report median runtimes as the variance is low ( $\sim 3\%$ ).

## 5.2 Performance of FF and RP

We compare FF and RP using datasets Healthcare (ACs  $\omega_1$  and  $\omega_2$  from Table 4) and Income (ACs  $\omega_3$  and  $\omega_4$ ) with the default parameter settings and queries  $Q_1$ ,  $Q_2$ , and  $Q_3$  (Table 3). In addition to runtime, we also measure number of candidates evaluated (NCE) which is the total of number of repair candidates for which we evaluate the AC and number of clusters accessed (NCA) which is the total number of clusters accessed by an algorithm.

**Runtime.** Figures 2a and 2b show the runtime of the FF and RP algorithms for Healthcare and Income, respectively. For a subset of the experiments we also report results for the Brute Force (BF) method. For a given constraint  $\omega_i$  we vary the bounds  $[B_l, B_u]$  to control what percentage of repair candidates have to be processed by the algorithms to determine the top- $k$  repairs as explained above. For example,  $\omega_1^{d=38}$  in Figure 2a for  $Q_1$  is the constraint  $\omega_1$  from Table 4 with the bounds set such that 38% of the candidate solutions have to be explored. We refer to this as the exploration distance

(ED). As expected, both FF and RP outperform BF by at least one order of magnitude in terms of runtime as shown in Figure 2a. The RP algorithm significantly reduces both NCE and NCA (e.g., Figures 2c and 2e), while the FF method maintains the same NCE as BF but decreases the NCA compared to BF (as BF does not use clusters we count tuple accesses). RP (pink bars) generally outperforms FF (blue bars) for most settings, demonstrating an additional improvement of up to an order of magnitude due to its capability of pruning and confirming sets of candidates at once. The only exception is settings where the top- $k$  repairs are found by exploring a very small portion of the search space, e.g.,  $Q_4$  with  $\omega_3$ .

**Total number of candidates evaluated (NCE).** We further analyze how NCE affects the performance of our methods (Figures 2c and 2d). RP consistently checks fewer candidates compared to FF. As observed in the runtime evaluation, the difference between the two algorithms is more pronounced when larger parts of the search space have to be explored.

**Total Number of Cluster Accessed (NCA).** The results for the number of clusters accessed are shown in Figures 2e and 2f for Healthcare and Income, respectively. Similar to the result for NCE, RP accesses significantly fewer clusters than FF.

## 5.3 Performance-Impacting Factors

To gain deeper insights into the behavior observed in Section 5.2, we investigate the relationship between the exploration distance (ED) and performance. We also evaluate the performance of FF and RP in terms of the parameters from Section 5.1.

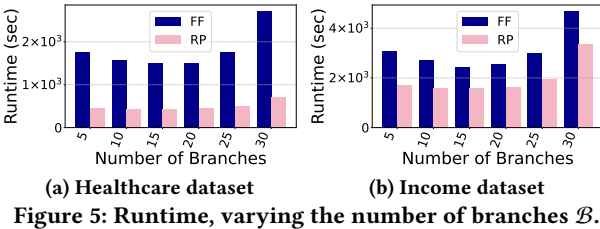
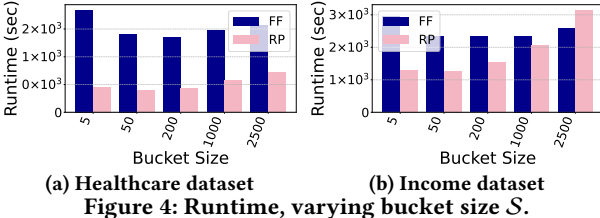
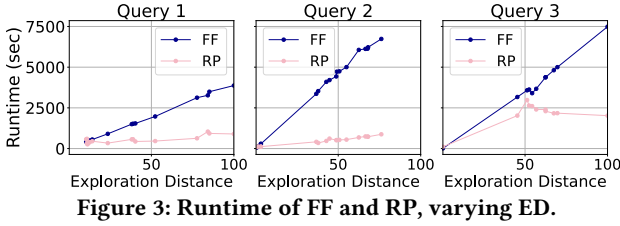
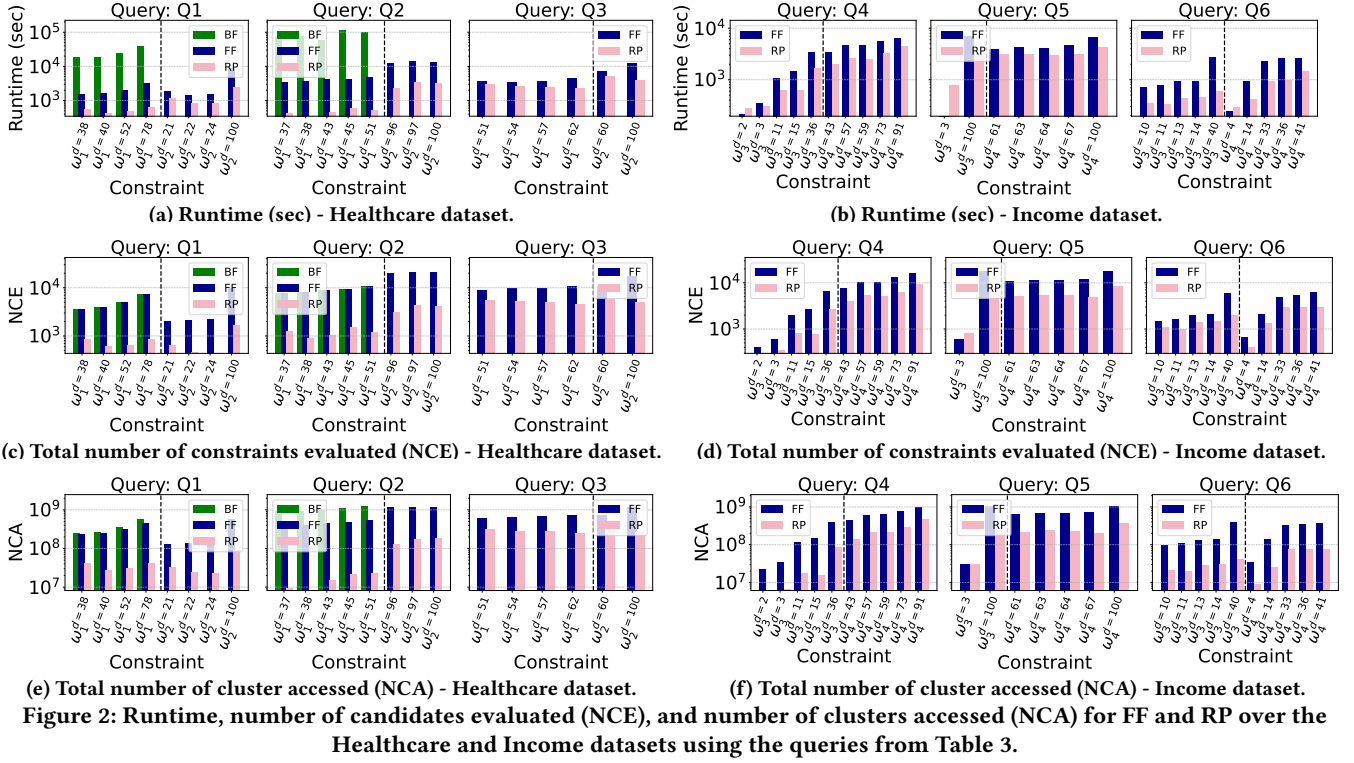
**Effect of Exploration distance.** We use queries  $Q_1$ – $Q_3$  and the constraint  $\omega_1$  on Healthcare and vary the bounds to control for ED. The result is shown in Figure 3. We present results for Income in [3]. For  $Q_1$  and  $Q_2$ , when ED is 10% or less, FF and RP exhibit comparable performance. A similar pattern is seen for  $Q_3$ , where FF performs better than RP for very low ED. The reason behind this trend is that when solutions are close to the user query (smaller ED), then there is less opportunity for pruning for RP.

**Effect of Bucket Size.** We now evaluate the runtime of FF and RP varying the bucket size  $\mathcal{S}$  using  $Q_1$  with  $\omega_1$  with bounds  $[0.44, 0.5]$  for the Healthcare dataset and  $Q_4$  with  $\omega_3$  using bounds  $[0.34, 0.39]$  for the Income dataset. We vary bucket size  $\mathcal{S}$  from 5 to 2500. We use the default branching factor  $\mathcal{B} = 5$ . Our algorithm chooses the number of levels to ensure that the size of leaf clusters is  $\leq \mathcal{S}$ . For example, for  $\mathcal{S} = 200$ , the tree will have 4 levels. The results of this experiment are shown in Figure 4. The advantage of smaller bucket sizes is that it is more likely that we can find a cluster that is fully covered / not covered at all. However, this comes at the cost of having to explore more clusters. In preliminary experiments, we have identified  $\mathcal{S} = 15$  to yield robust performance for a wide variety of settings and use this as the default.

**Table 5: Branching Configuration and Data Distribution**

# of Branches	# of Leaves	# of Branches	# of Leaves
5	15625	20	8000
10	10000	25	15625
15	3375	30	27000

**Effect of the Branching Factor.** We now examine the relationship between the branching factor  $\mathcal{B}$  and the runtime of FF and



RP. We use the same queries, constraints, bounds, and datasets as in the previous evaluation and vary branching factor  $B$  from 5 to 30. The corresponding number of leaf nodes in the kd-tree is shown in Table 5. As we use the default bucket size  $S = 15$ , the branching factor determines the depth of the tree. The result

shown in Figure 5 confirms that, as expected, the performance of FF and RP correlates with the number of clusters at the leaf level. For FF, branching factors of 5 and 25 yield nearly identical runtime because both have the same number of leaf nodes (15,625). A similar pattern can be observed for  $B = 10$  and  $B = 20$ . At  $B = 15$ , FF achieves the lowest runtime, as it involves the smallest number of leaves (3,375). For  $B = 30$ , the number of leaf clusters significantly increases, leading to a substantial rise in the runtime of FF. For RP, overall performance trends align with those of FF. However, RP is less influenced by the branching factor as for smaller clusters it may be possible to prune / confirm larger candidate sets at once. Both bucket size  $S$  and branching factor  $B$  impact performance, and the optimal values depend on the characteristics of the dataset and queries. The intuition is as follows: when  $S$  is too small, the resulting tree becomes too deep, leading to an excessive number of leaf clusters; when  $S$  is too large, the ability to prune effectively diminishes because clusters encompass too many data points. Likewise, if  $B$  is too large, data is distributed across many child nodes, making it harder to prune entire sub-trees; if  $B$  is too small, the tree again becomes too deep with many leaf clusters. For new datasets, we suggest starting with moderate values for both  $S$  and  $B$ , then adjusting based on the number of leaf clusters observed: if the tree has too many leaf clusters, consider increasing  $S$  or decreasing  $B$ ; if pruning is insufficient, consider decreasing  $S$  or increasing  $B$ . There are several additional factors that can affect optimal choices for these parameters: (i) strong correlations between attributes used in conditions lead to more homogeneous clusters which in turn means that larger clusters can be tolerated without significantly

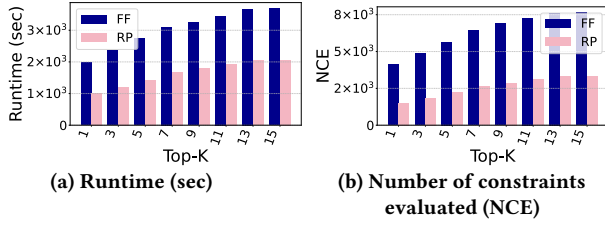


Figure 6: Performance of FF and RP varying  $k$

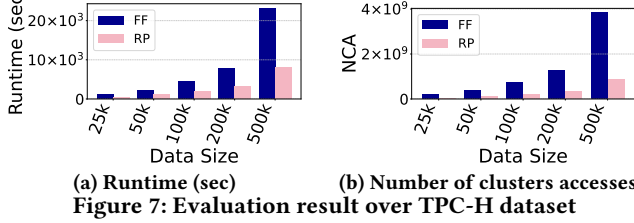


Figure 7: Evaluation result over TPC-H dataset

impacting pruning power, (ii) if attributes in user query conditions are correlated with attributes of filter-aggregation queries, then aggregation results vary widely for clusters potentially leading to a stronger separation between repairs and more pruning potential even with larger clusters. We leave automatic parameter tuning, e.g., based on measuring correlations between attributes over a sample, to future work.

**Effect of  $k$ .** In this experiment, we vary the parameter  $k$  from 1 to 15. For both FF and RP, as  $k$  increases, the runtime also increases, as shown in Figure 6a. When  $k$  is larger, the algorithms have to explore a larger fraction of the search space to find additional repairs. Similarly, the NCE as shown in Figure 6b exhibits the same increasing trend. RP consistently outperforms FF.

**Effect of Dataset Size.** Next, we vary the dataset size and measure the runtime and NCA for the TPC-H dataset using  $Q_7$  with  $\omega_5$ . Dataset size impacts both the size of the search space and the size of the kd-tree. Nonetheless, as shown in Figure 7a, our algorithms scale roughly linearly in dataset size demonstrating the effectiveness of reusing aggregation results for clusters and range-based pruning. This is further supported by the NCA measurements shown in Figure 7b, which exhibit the same trend as the runtime.

## 5.4 Comparison with Related Work

We compare our approach with Erica [26], which solves the related problem of finding all minimal refinements of a given query that satisfy a set of cardinality constraints for groups within the result set. Such constraints are special cases on the ACs we support. Erica returns all repairs that are not *dominated* (the skyline [10]) by any other repair where a repair dominates another repair if it is at least as close to the user query for every condition  $\theta_i$  and strictly closer in at least one condition. Thus, different from our approach, the number of returned repairs is not an input parameter in Erica. For a fair comparison, we compute the minimal repairs and then set  $k$  such that our methods returns a superset of the repairs returned by Erica. Our algorithms, like Erica, operate by modifying constants in predicates on attributes already present in the query and do not introduce new predicates. A key difference is that Erica supports adding constants to set membership predicates for categorical attributes, e.g., replacing  $A \in \{c_1\}$  with  $A \in \{c_1, c_2, c_3\}$ , while our

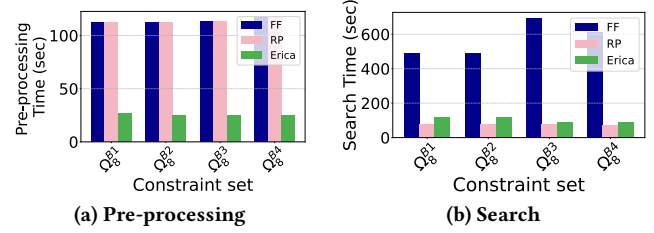


Figure 8: Runtimes of FF, RP, and Erica.

approach maps categorical values to numeric codes and adjusts thresholds. As we will discuss in Section 6, both our approach and Erica can model addition of new predicates by refining dummy predicates that evaluate to true on all inputs. To conduct the evaluation for Erica, we used the available Python implementation [1]<sup>3</sup>.

We adopt the queries, constraints, and the dataset from [26]. We compare the generated refinements and runtime of our techniques with Erica using  $Q_1$  and  $Q_2$  (Table 3) on the Healthcare dataset (50K tuples) with constraints  $\Omega_6$  and  $\Omega_7$  (Table 4), respectively.

**Generated Repairs Comparison.** We first compare the generated repairs by our approach and Erica. As mentioned above, we adjust  $k$  per query and constraint set to ensure that our approach returns a superset of the repairs returned by Erica. For  $Q_1$  with  $\Omega_6$  ( $Q_2$  with  $\Omega_7$ ), Erica generates 7 (9) minimal repairs whereas our technique generates 356 (1035), including those produced by Erica. The top-1 repair returned by our approach is guaranteed to be minimal. However, the remaining minimal repairs returned by Erica may have a significantly higher distance to the user query than the remaining top-k answers returned by our approach. For example, in  $Q_2$ , given the condition  $\text{num-children} \geq 4$  of the user query, our solution includes a refined condition  $\text{num-children} \geq 3$  whereas Erica provides a refinement  $\text{num-children} \geq 1$  which is dissimilar to the user query.

**Runtime Comparison.** For this experiment we use  $Q_4$  with  $\Omega_8$  for the 50K Income dataset, which is derived from the dataset, query, and constraint used to evaluate Erica in [26]. We use the same bounds in the constraints for both Erica and our algorithms:  $B1 := (B_{u1} = 30, B_{u2} = 150, B_{u3} = 10)$  and  $B2 := (B_{u1} = 30, B_{u2} = 300, B_{u3} = 25)$ ,  $B3 := (B_{u1} = 10, B_{u2} = 650, B_{u3} = 50)$ , and  $B4 := (B_{u1} = 15, B_{u2} = 200, B_{u3} = 15)$ . To ensure a fair comparison of execution time, we fix the number of generated repairs (i.e., top- $k$ ) in our approach to equal to the number of repairs produced by Erica. We set  $k=17$  for constraint sets  $\Omega_8^{B1}$  and  $\Omega_8^{B2}$ ,  $k=11$  for  $\Omega_8^{B3}$ , and  $k=13$  for  $\Omega_8^{B4}$ . Due to the different optimization goals, variations in the generated repairs between our approach and Erica are expected. The results shown in Figure 8b demonstrate that the total runtime of RP and Erica are comparable with RP exhibiting a significantly lower time for exploring candidate solutions while Erica has a significantly lower preprocessing time. The higher preprocessing time is expected as Erica only has to generate provenance expressions and build lists of candidate constants sorted by their distance to the constant in the user's query. In contrast, our approach has to cluster

<sup>3</sup>We replaced Erica's use of pandas for filtering data and constraint evaluation (which are implemented in C) with equivalent pure-Python loops over lists just as in our own code, so that both implementations are using the same programming language. This change ensures that our comparison highlights algorithmic differences rather than language speed.

**Table 6: Comparison of query repair techniques**

Approach	Supports sum, min, max, avg	Distance Metric	Constrains Result Sub- sets?	Repairs Joins?	Arithmetic Expressions Supported
HC [11]	✗	✗	✗	✗	✗
TQGen [30]	✗	✗	✗	✗	✗
SnS [29]	✗	✗	✗	✗	✗
EAGER [2]	✓	<i>lin. comb.</i>	✗	✗	✗
SAUNA [23]	✗	<i>L2 (result)</i>	✗	✓	✗
ConQueR [35]	✗	<i>edit-distance</i>	✗	✓	✗
FixTed [9]	✗	<i>skyline</i>	✗	✓	✗
FARQ [32]	✗	<i>Jacc. (result)</i>	✓	✗	✗
Erica [26]	✗	<i>skyline</i>	✓	✗	✗
<b>RP (ours)</b>	✓	<i>lin. comb.</i>	✓	✗	✓

the data, index the clusters in a kd-tree, and materialize summaries for each cluster, which is more computationally intensive. However, this extra work enables us to reason about complex, non-monotone constraints, which Erica’s list-based approach cannot. Furthermore, we argue that it can be beneficial to decrease search time at the cost of higher preprocessing time as some of the preprocessing results could be shared across user requests. Overall, the total runtime of RP and Erica is comparable, even though our approach does not apply any specialized optimizations that exploit monotonicity as in Erica. These results also highlight the need for our range-based optimizations, as FF is significantly slower than Erica.

## 6 RELATED WORK

**Query refinement & relaxation.** In Table 6 we compare the capabilities of several query refinement techniques for aggregate constraints in terms of supported aggregates (only **count** or also other aggregates), distance metric used to compare repairs to the original query based on distances between predicates (e.g., *lin. comb.*: linear combination of predicate-level distances, *skyline*: skyline over predicate-level distances), whether the method allows constraints that apply only to a subset of the result (some methods only constrain the whole query result), whether join conditions can be repaired, and whether they support arithmetic expressions. Li et al. [26] determine all minimal refinements of a conjunctive query by changing constants in selection conditions such that the refined query fulfills a conjunction of cardinality constraints, e.g., the query should return at least 5 answers where gender = female. A refinement is minimal if it fulfills the constraints and there does not exist any refinement that is closer to the original query in terms of similarity of constants used in predicates (*skyline*). However, [26] only supports cardinality constraints (**count**) and does not allow for arithmetic combinations of the results of such queries as shown in Table 6. Mishra et al. [29] refine a query to return a given number  $k$  of results with interactive user feedback. Koudas et al. [25] refine a query that returns an empty result to produce at least one answer. In [9, 35], a query is repaired to return missing results of interest provided by the user. Campbell et al. [13] repair top- $k$  queries, supporting non-monotone constraints through the use of constraint solvers. [11, 30] refine queries for database testing such that subqueries of the repaired query approximately fulfill cardinality constraints. [11] demonstrated that the problem is NP-hard in the number of predicates. Both approaches do not optimize for similarity to the user query. [23] relaxes a query to return approximately  $N$  results preferring repairs based on the difference between

the result of the user query and repair. Most work on query refinement has limited the scope to constraints that are monotone in the size of the query answer. Monotonicity is then exploited to prune the search space [12, 23, 29, 30, 38]. To the best of our knowledge, our approach is the only one that supports arithmetic constraints which is necessary to express complex real world constraints, e.g., standard fairness measures, but requires novel pruning techniques that can handle such non-monotone constraints. While some approaches explicitly support addition and deletion of predicates, any approach that can both relax or refine predicates is capable of supporting adding / deleting predicates: deletion by relaxing a predicate until it evaluates to true on all inputs and addition by adding dummy predicates that evaluate to true on all inputs and then either refine them (adding a new predicate) or not (decide to not add this predicate).

**How-to queries.** Like in query repair [26], the goal of how-to queries [28] is to achieve a desired change to a query’s result. However, how-to queries change the database to achieve this result instead of repairing the query. Wang et al. [37] study the problem of deleting operations from an update history to fulfill a constraint over the current database. However, this approach does not consider query repair (changing predicates) nor aggregate constraints.

**Explanations for Missing Answers.** Query-based explanations for missing answers [15, 18, 19] identify which operators of a query are responsible for the failure of the query to return a result of interest. However, this line of work does not generate query repairs.

**Bounds with Interval Arithmetic.** Prior work has highlighted the effectiveness of interval arithmetic across various database applications [17, 20, 33, 39]. For instance, [20] determines bounds on query results over uncertain databases. The work in [39] introduced a bounding technique for iceberg cubes, establishing an early foundation for leveraging interval arithmetic to constrain aggregates. Interval arithmetic has been used extensively in *abstract interpretation* [16, 17, 33] to bound the result of computations.

## 7 CONCLUSIONS AND FUTURE WORK

We introduce a novel approach for repairing a query to satisfy a constraint on the query’s result. We support a significantly larger class of constraints than prior work, including common fairness metrics like SPD. We avoid redundant work by reusing aggregate results when evaluating repair candidates and present techniques for evaluating multiple repair candidates at once by bounding their results. Our approach works best if there is homogeneity among similar repair candidates that can be exploited. Interesting directions for future work include (i) the study of more general types of repairs, e.g., repairs that add or remove joins or change the structure of the query, (ii) considering other optimization criteria, e.g., computing a skyline as in some work on query refinement, (iii) employing more expressive domains than intervals for computing tighter bounds, e.g., zonotopes [17], and (iv) supporting dynamic settings where the table, predicates, constraints, or distance metrics may change. In this regard, we may exploit efficient incremental maintenance kd-trees and aggregate summaries [7, 8]. However, our setting is more challenging as small changes to aggregation results can affect the validity of large sets of repair candidates.



## REFERENCES

- [1] 2024. Erica implementation. [https://github.com/JinyangLi01/Query\\_refinement](https://github.com/JinyangLi01/Query_refinement). Accessed on 2024-03-20.
- [2] Abdullah M. Albarak and Mohamed A. Sharaf. 2017. Efficient schemes for similarity-aware refinement of aggregation queries. *WWW* 20, 6 (2017), 1237–1267.
- [3] Shatha Algarni, Boris Glavic, Seokki Lee, and Adriane Chapman. 2025. *Efficient Query Repair for Arithmetic Expressions with Aggregation Constraints (Extended Version)*. Technical Report. <https://arxiv.org/abs/2511.00826>.
- [4] Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde, Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic, Seema Nagar, Karthikeyan Natesan Ramamurthy, John T. Richards, Diptikalyan Saha, Prasanna Sattigeri, Moninder Singh, Kush R. Varshney, and Yunfeng Zhang. 2019. AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias. *IBM J. Res. Dev.* 63, 4/5 (2019), 4:1–4:15.
- [5] Omar Benjelloun, Anish Das Sarma, Alon Y. Halevy, and Jennifer Widom. 2006. ULDBs: databases with uncertainty and lineage. In *VLDB*.
- [6] Jon Louis Bentley. 1975. Multidimensional Binary Search Trees Used for Associative Searching. *CACM* 18, 9 (1975), 509–517.
- [7] Jon Louis Bentley. 1975. Multidimensional binary search trees used for associative searching. *CACM* 18, 9 (1975), 509–517.
- [8] Jon Louis Bentley and James B Saxe. 1980. Decomposable searching problems I. Static-to-dynamic transformation. *Journal of Algorithms* 1, 4 (1980), 301–358.
- [9] Nicole Bidoit, Melanie Herschel, and Katerina Tzompanaki. 2016. Refining SQL Queries based on Why-Not Polynomials. In *TaPP*.
- [10] S. Borzsony, D. Kossmann, and K. Stocker. 2001. The skyline operator. In *ICDE*. 421–430.
- [11] N. Bruno, S. Chaudhuri, and D. Thomas. 2006. Generating Queries with Cardinality Constraints for DBMS Testing. *TKDE* 18, 12 (2006), 1721–1725.
- [12] Nicolas Bruno, Surajit Chaudhuri, and Dilys Thomas. 2006. Generating Queries With Cardinality Constraints for DBMS Testing. *TKDE* 18, 12 (2006), 1721–1725.
- [13] Felix S. Campbell, Alon Silberstein, Julia Stoyanovich, and Yuval Moskovitch. 2024. Query Refinement for Diverse Top-k Selection. *SIGMOD* 2, 3 (2024), 166.
- [14] Felix S. Campbell, Julia Stoyanovich, and Yuval Moskovitch. 2024. Rodeo: Making Refinements for Diverse Top-K Queries. *PVLDB* 17, 12 (2024), 4341–4344.
- [15] Adriane Chapman and H. V. Jagadish. 2009. Why not?. In *SIGMOD*. 523–534.
- [16] Patrick Cousot and Radhia Cousot. 1977. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL*. 238–252.
- [17] Luiz Henrique De Figueiredo and Jorge Stolfi. 2004. Affine arithmetic: concepts and applications. *Numerical algorithms* 37 (2004), 147–158.
- [18] Daniel Deutch, Nave Frost, Amir Gilad, and Tomer Haimovich. 2020. Explaining Missing Query Results in Natural Language. In *EDBT*. 427–430.
- [19] Ralf Diestelkämper, Seokki Lee, Melanie Herschel, and Boris Glavic. 2021. To Not Miss the Forest for the Trees - A Holistic Approach for Explaining Missing Answers over Nested Data. In *SIGMOD*. 405–417.
- [20] Su Feng, Boris Glavic, Aaron Huber, and Oliver A Kennedy. 2021. Efficient uncertainty tracking for complex queries with attribute-level bounds. In *SIGMOD*. 528–540.
- [21] Sorelle A. Friedler, Carlos Scheidegger, Suresh Venkatasubramanian, Sonam Choudhary, Evan P. Hamilton, and Derek Roth. 2019. A comparative study of fairness-enhancing interventions in machine learning. In *FAT*. 329–338.
- [22] Stefan Grafberger, Shubha Guha, Julia Stoyanovich, and Sebastian Schelter. 2021. Mlinspect: A data distribution debugger for machine learning pipelines. In *SIGMOD*. 2736–2739.
- [23] Abhijit Kadlag, Amol V. Wanjari, Juliana Freire, and Jayant R. Haritsa. 2004. Supporting Exploratory Queries in Databases. In *DASFAA*, Vol. 2973. 594–605.
- [24] Dmitri V Kalashnikov, Laks VS Lakshmanan, and Divesh Srivastava. 2018. Fastqre: Fast query reverse engineering. In *SIGMOD*. 337–350.
- [25] Nick Koudas, Chen Li, Anthony K. H. Tung, and Rares Vernica. 2006. Relaxing Join and Selection Queries. In *VLDB*. 199–210.
- [26] Jinyang Li, Yuval Moskovitch, Julia Stoyanovich, and HV Jagadish. 2023. Query Refinement for Diversity Constraint Satisfaction. *PVLDB* 17, 2 (2023), 106–118.
- [27] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2022. A Survey on Bias and Fairness in Machine Learning. *ACM Comput. Surv.* 54, 6 (2022), 115:1–115:35.
- [28] Alexandra Meliou and Dan Suciu. 2012. Tiresias: the database oracle for how-to queries. In *SIGMOD*. 337–348.
- [29] Chaitanya Mishra and Nick Koudas. 2009. Interactive query refinement. In *EDBT*. 862–873.
- [30] Chaitanya Mishra, Nick Koudas, and Calisto Zuzarte. 2008. Generating targeted queries for database testing. In *SIGMOD*. 499–510.
- [31] Guido Moerkotte. 1998. Small Materialized Aggregates: A Light Weight Index Structure for Data Warehousing. In *VLDB*. 476–487.
- [32] Suraj Shetiya, Ian P Swift, Abolfazl Asudeh, and Gautam Das. 2022. Fairness-aware range queries for selecting unbiased data. In *ICDE*. 1423–1436.
- [33] Jorge Stolfi and Luiz Henrique de Figueiredo. 2003. An introduction to affine arithmetic. *Trends in Computational and Applied Mathematics* 4, 3 (2003), 297–312.
- [34] TPC. 2024. TPC BENCHMARK H (Decision Support) Standard Specification Revision 3.0.1. [https://www.tpc.org/tpc\\_documents\\_current\\_versions/current\\_specifications5.asp](https://www.tpc.org/tpc_documents_current_versions/current_specifications5.asp). Accessed on 2024-03-20.
- [35] Quoc Trung Tran and Chee-Yong Chan. 2010. How to conquer why-not questions. In *SIGMOD*. 15–26.
- [36] Bienvenido Vélaz, Ron Weiss, Mark A. Sheldon, and David K. Gifford. 1997. Fast and Effective Query Refinement. In *SIGIR*. 6–15.
- [37] Xiaolan Wang, Alexandra Meliou, and Eugene Wu. 2017. QFix: Diagnosing Errors through Query Histories. In *SIGMOD*. 1369–1384.
- [38] Eugene Wu and Samuel Madden. 2013. Scorpion: Explaining Away Outliers in Aggregate Queries. *PVLDB* 6, 8 (2013), 553–564.
- [39] Xiuzhen Zhang, Pauline Lienhua Chou, and Guozhu Dong. 2007. Efficient computation of iceberg cubes by bounding aggregate functions. *TKDE* 19, 7 (2007), 903–918.
- [40] Mohamed Ziauddin, Andrew Witkowski, You Jung Kim, Dmitry Potapov, Janaki Lahorani, and Murali Krishna. 2017. Dimensions based data clustering and zone maps. *PVLDB* 10, 12 (2017), 1622–1633.