# Stress-Testing Causal Claims via Cardinality Repairs

Yarden Gabbay
yardengabbay@campus.technion.ac.il
Technion
Israel

Haoquan Guan
h3guan@ucsd.edu
University of California, San Diego
USA

Shaull Almagor
shaull@technion.ac.il
Technion
Israel

El Kindi Rezig
elkindi.rezig@utah.edu
University of Utah
USA

Brit Youngmann
brity@technion.ac.il
Technion
Israel

Babak Salimi
bsalimi@ucsd.edu
University of California, San Diego
USA

## ABSTRACT

Causal analyses derived from observational data underpin high-stakes decisions in domains such as healthcare, public policy, and economics. Yet such conclusions can be surprisingly fragile: even minor data errors - duplicate records, or entry mistakes - may drastically alter causal relationships. This raises a fundamental question: *how robust is a causal claim to small, targeted modifications in the data?* Addressing this question is essential for ensuring the reliability, interpretability, and reproducibility of empirical findings.

We introduce SUBCURE, a framework for *robustness auditing via cardinality repairs.* Given a causal query and a user-specified target range for the estimated effect, SUBCURE identifies a small set of tuples or subpopulations whose removal shifts the estimate into the desired range. This process not only quantifies the sensitivity of causal conclusions but also pinpoints the specific regions of the data that drive those conclusions. We formalize this problem under both tuple- and pattern-level deletion settings and show both are NP-complete. To scale to large datasets, we develop efficient algorithms that incorporate machine unlearning techniques to incrementally update causal estimates without retraining from scratch.

We evaluate SUBCURE across four real-world datasets covering diverse application domains. In each case, it uncovers compact, high-impact subsets whose removal significantly shifts the causal conclusions—revealing vulnerabilities that traditional methods fail to detect. Our results demonstrate that cardinality repair is a powerful and general-purpose tool for stress-testing causal analyses and guarding against misleading claims rooted in ordinary data imperfections.

## 1 INTRODUCTION

Causal inference now underpins decisions in medicine, policy, and economics while powering fairness auditing and data-debiasing [75, 88], explainability [28, 59], domain-robust learning [50, 77], and core data-management tasks such as query explanation, discovery, and cleaning [27, 53, 64, 73, 85, 86]. Yet a growing body of work shows that seemingly routine data imperfections - such as misclassification, measurement noise, duplicate records, and composite corruption - can distort causal estimates and invalidate hypothesis tests [3, 12, 39, 46, 58, 60, 76]. If left unaddressed, such distortions can lead to

ineffective medical interventions, misallocated resources, or flawed algorithmic decisions, undermining the reliability of empirical findings and their downstream applications. These observations highlight a critical question: **How robust are our causal claims to common, low-level data errors?**

In this work, we ask a concrete version of the robustness question: given an observational dataset, a treatment variable $T$, an outcome variable $O$, and a causal estimation procedure (e.g., the Average Treatment Effect (ATE), defined as the expected difference in outcomes between treated and untreated groups), how many data records must be removed to shift the estimated effect into a user-specified target range? This question formalizes a data-centric notion of robustness-sensitivity not to modeling assumptions, but to minimal, targeted perturbations of the data itself. We refer to this as the problem of *cardinality repair for causal effect targeting* (abbreviated as CaRET for Cardinality Repair for causal Effect Targeting): identifying the smallest subset of tuples (or subpopulations) whose removal pushes the estimate across a specified threshold. Beyond quantifying fragility, such repairs help reveal which records or regions of the data are most responsible for driving the causal conclusion.

We present SUBCURE (SUBset CaUsal REpair), a framework that helps analysts assess how strongly their causal conclusions depend on specific parts of the data. SUBCURE supports two modes: *tuple-level* repair, which identifies individual records whose removal shifts the estimated effect, and *pattern-level* repair, which targets compact subpopulations defined by attribute–value predicates. It surfaces minimal edits that change the conclusion and highlights the regions of the data most responsible for a given estimate, enabling users to explore robustness interactively and with minimal assumptions. We now illustrate this workflow through two motivating examples, showing how small, targeted deletions can flip or fortify real-world causal conclusions.

EXAMPLE 1.1. ***Twins mortality.*** *The* Twins *dataset [47] records birth details for same-sex twins; the* treatment *is being the heavier twin, and the* outcome *is first-year mortality. Alex, a social scientist, adjusts for gestational age, birth weight, prenatal care, and maternal status and obtains an ATE of −0.016—suggesting a small protective effect for the heavier twin.* Tuple-level insight: *Running* SUBCURE *in tuple mode, Alex discovers that removing just **270 records** (1.1% of the sample)*

nudges the ATE to −0.0009—a **94.4% move toward zero**. *These records are not statistical outliers under common heuristics; they appear representative on observed covariates yet collectively exert strong influence on the estimate. A quick inspection hints that many belong to extremely low-gestational-age births, where overall mortality is uniformly high—an avenue for further clinical review.* Pattern-level insight: *Switching to pattern mode,* SUBCURE *highlights a subpopulation of* **4 401 first-born twins** *weighing between* **1786 g and 1999 g** *(18.3% of the data). Deleting this slice flips the ATE to* 0.0003—a **101.8% swing** *that reverses the sign. This weight band sits just below the 2 kg clinical low-birth-weight threshold; delivery complications concentrated in this range may offset any advantage of being the heavier twin, illustrating how a narrow but medically meaningful subgroup can dictate the overall conclusion. Together, these findings show how* SUBCURE *quantifies fragility and surfaces domain-specific hypotheses—here, extreme prematurity and a critical birth-weight band—that warrant deeper investigation before drawing policy or clinical inferences.*

EXAMPLE 1.2. ***Disability and wages.*** *The* American Community Survey *(ACS) [79] is a nationwide instrument conducted by the U.S. Census Bureau. Alex studies the causal effect of* not *having a disability on annual wages, controlling for education, public-health coverage, gender, and age. The estimated ATE is $8,774, indicating that, on average, people without disabilities earn roughly nine thousand dollars more per year.* Tuple-level insight: *Running* SUBCURE *in tuple mode, Alex learns that deleting just* **8,400 records** *(0.7% of the data) lifts the ATE to $11,512—a* **31.1% increase***. A quick look shows many of the removed records cluster in low-wage service occupations with high disability rates; trimming them disproportionately amplifies high-income, non-disabled earners, inflating the gap.* Pattern-level insight: *Shifting the estimate into the target range requires removing a subpopulation comprising 11.9% of the data. In contrast, to achieve a comparable* **downward** *shift,* SUBCURE *must remove nearly 74% of the data. This asymmetry signals that the observed wage advantage is highly sensitive to a small slice of low-income disabled workers, yet remarkably stable in the opposite direction. Such asymmetry raises fairness questions: the gap may be driven by a narrow pocket of the labor market rather than a uniform disadvantage. Together, these findings illustrate how* SUBCURE *pinpoints imbalance in representation—here, a small cohort of low-income disabled respondents—helping analysts decide whether additional weighting, stratification, or data collection is needed before drawing policy conclusions.*

Our approach builds on the long-standing tradition of *sensitivity analysis*, which asks how causal conclusions shift under model-based threats such as hidden confounding, selection bias, or measurement error [11, 22, 54, 68, 78]. Classical methods typically rely on *controlled bias parameters*, or an additive-noise model to capture measurement error—while assuming the dataset itself is clean and internally coherent. Several works have proposed robustness metrics to evaluate how stable causal conclusions remain under data perturbations or removal, to overturn an inference [26], quantifying worst-case

estimate changes under data deletions [15], and estimating the worst-case ATE across subpopulations under confounder shifts [38]. Our approach quantifies ATE sensitivity by directly examining how causal conclusions shift when actual tuples are removed. This is particularly relevant for real-world datasets, which often contain messy and heterogeneous issues, such as duplicate records, miscoded or missing values, format inconsistencies, or corrupted subsets, that previous works do not capture. Our cardinality-repair strategy addresses this gap by treating the data as mutable and asking *how many records - or which subpopulations - must be removed to move an effect into a user-chosen range?* The resulting repair size provides a model-agnostic, data-centric, and intuitive robustness metric that complements and guides classical sensitivity analyses.

Our method also relates to prior work on constraint-based data cleaning. Integrity-oriented approaches enforce functional or conditional functional dependencies via tuple-level repairs [14, 16, 29, 41, 44, 45, 65]. More recent systems repair violations of conditional independence between treatment and outcome—typically to force a *null* effect (ATE = 0) [64, 75, 80]. SUBCURE subsumes these cases yet goes further on two fronts: it supports *arbitrary* effect targets (not just zero), and it pinpoints *influential subpopulations*, rather than only individual rule-breaking tuples [64, 83]. This broader scope lets analysts reason simultaneously about constraint correctness, structural bias, and the overall robustness of causal claims.

Our main contributions are as follows.

• **Problem formulation and hardness results (Section 4).** We formalize the *cardinality repair for causal effect targeting* (CaRET) problem: Given a treatment, an outcome, and a desired effect interval, find the smallest data subset whose removal shifts the estimated effect into the range. We consider two cost models: (i) tuple-level deletions, and (ii) pattern-level deletions that remove entire subpopulations defined by attribute–value predicates. We prove that both variants are NP-complete. For tuple-level, this is shown by reduction from SUBSET-SUM. For pattern-level, we show that this problem is harder than tuple-level, regardless of the causal effect function.

• **Fast, incremental search algorithms (Sections 5–6).** We develop fast, scalable cardinality-repair algorithms using two complementary search strategies. In tuple mode, we cluster the dataset using a two-stage $k$-means procedure to sample a proxy set, estimate each sampled tuple's marginal influence, and iteratively delete the one with the highest impact. Influence scores are batch-refreshed every few steps, avoiding unnecessary rescoring and reducing runtime. In pattern mode, we perform bottom-up random walks over conjunctions of attribute–value predicates. A dynamic weighting mechanism steers the walk toward impactful predicates, and exploration stops early if the candidate subgroup grows too large. These strategies yield an efficient anytime search that surfaces high-leverage subpopulations without exhaustive enumeration.

Both strategies require re-evaluating the causal effect after each candidate deletion. To make this feasible, we use two standard causal effect estimators—linear regression [24] and

inverse-propensity weighting (IPW) [69]—and develop incremental update optimizations that eliminate the need to recompute from scratch after each change. For linear regression, we cache sufficient statistics such as the covariance matrix and response cross-product, and apply low-rank updates when rows are removed. This lets us recompute the treatment effect without rebuilding the full model, and in time that depends only on the number of confounders, not the size of the dataset. For IPW, we warm-start from the previous logistic regression fit and apply a single Fisher-scoring step that adjusts only for the removed rows. This avoids full optimization over all tuples and yields updated propensity scores in constant time per deletion. These incremental estimators ensure that each search step requires only localized computation, preserving estimator fidelity while enabling SubCure to scale to large, high-dimensional datasets with interactive latency.

• **Comprehensive empirical evaluation (Section 7).** We evaluate SubCure on four public datasets and one synthetic benchmark, comparing against eight baseline methods. Sub-Cure (i) consistently identifies smaller subsets of influential data, (ii) scales to million-row, high-dimensional datasets, and (iii) achieves up to order-of-magnitude speed-ups through incremental updates without compromising accuracy. Case studies demonstrate that SubCure complements existing sensitivity analysis methods by offering additional insights into ATE sensitivity, while its identified subsets remain influential across multiple ATE estimators, highlighting the robustness of our approach.

## 2 RELATED WORK

**Constraint-based data cleaning** Traditional data cleaning research focuses on enforcing integrity constraints [19, 35], such as functional dependencies [14, 16, 41, 44, 45, 65], conditional functional dependencies [14, 29], denial constraints [20, 21], and inclusion dependencies [13]. Common repair operations include tuple deletion (cardinality repair [2, 56]), insertion [75], and value updates [64]. Recent work [64, 75, 83] has explored enforcing statistical constraints that represent conditional dependence or independence between attributes. Capuchin [75] computes optimal repairs by adding or removing tuples to enforce conditional independence (CI) constraints. The authors of [6] use GANs to generate data that satisfies CI constraints, with a focus on learning generative models rather than repairing the data. The works most closely related to ours are OTClean [64] and SCODED [83], both of which enforce statistical constraints representing conditional dependence or independence. OTClean employs optimal transport to probabilistically adjust data values to satisfy CI constraints; since it modifies values rather than removing tuples, its results are not directly comparable to ours. SCODED focuses on identifying the most influential tuples that violate a CI constraint. Its goal is to provide explanations rather than repairs, and thus it reports the top-$k$ most influential tuples instead of minimizing the subset to remove. In contrast, SubCurenot only enforces CIs (by driving the ATE to zero), but also supports targeted

causal repair, shifting the ATE into a user-defined range. Beyond identifying influential tuples, SubCure is also capable of uncovering meaningful subpopulations whose removal has a significant effect on the ATE.

**Sensitivity analysis**: A large body of work has developed techniques to assess how causal conclusions might shift in the presence of various threats [23, 66]. For instance, [68] quantifies the influence of unmeasured confounding in observational studies. [11, 54] extended this perspective to selection bias, presenting sensitivity analysis tools for evaluating how causal estimates would change under varying assumptions about the selection process. [22] developed a framework for omitted variable bias, enabling interpretable and flexible assessments of robustness to hidden confounders. [78] addressed measurement error through the simulation extrapolation method, which corrects causal estimates by modeling the effect of added noise and extrapolating to the noise-free case.

Several works have proposed robustness measures to assess the stability of causal inferences under data removal or perturbations. KonFound [26] is a robustness measure that quantifies the degree of hidden bias needed to overturn an inference. It assesses sensitivity to unobserved confounders through hypothetical value perturbations, identifying the "switch point" where an effect becomes null. ZamInfluence [15] is a general-purpose robustness measure that measures how much an estimate can change when a fraction of the data is removed, treating the estimator as a black box. It searches for worst-case deletions that either flip the sign of the estimate or alter its statistical significance, making it applicable across diverse statistical models. WTE [38] is a robustness measure for confounders shift that estimates the worst-case ATE across subpopulations of a given size, providing conservative guarantees of external validity under distributional changes. Our approach offers a complementary perspective. Assuming confounders are observed, we examine ATE sensitivity by removing actual tuples or subpopulations. Rather than identifying a single "switch point", our method explores how the ATE responds to data removal and how it can be shifted toward any user-specified target range, offering a data-centric robustness measure that can support and inform traditional sensitivity analyses.

**Machine unlearning**: Machine unlearning focuses on efficiently removing the influence of data points from trained models without full retraining [30, 33, 81]. This is particularly relevant in settings requiring compliance with data privacy regulations ("the right to be forgotten" [67]). Prior work has proposed exact and approximate unlearning techniques across various models, including regression models [33, 37, 82], and deep neural networks [32, 61]. In this work, we leverage unlearning techniques to accelerate causal effect estimation after data subset removal, enabling efficient exploration of data subsets that influence the ATE.

## 3 BACKGROUND ON CAUSAL INFERENCE

**Causal inference and ATE.** We use Pearl's model for *observational causal analysis* [62]. Causal inference aims to quantify

the effect of a *treatment variable* ($T$) on an *outcome variable* ($O$). For instance, to estimate the impact of disability on annual wage. The gold standard for establishing causality is through *randomized controlled trials*. In such experiments, a population is randomly divided into two groups: the **treated group**, which receives the intervention ($\mathrm{do}(T = 1)$ for a binary treatment), and the **control group**, which does not ($\mathrm{do}(T = 0)$). This random assignment ensures that, on average, all other factors are balanced between the groups, isolating the effect of the treatment. A widely used metric for estimating causal effects is the **Average Treatment Effect (ATE)**, which measures the difference between the average outcome observed in the treated and in the control group [62, 72].

$$ATE(T, O) = \mathbb{E}[Y \mid \mathrm{do}(T = 1)] - \mathbb{E}[Y \mid \mathrm{do}(T = 0)] \quad (1)$$

The do-operator $\mathrm{do}(T = 1)$ formalizes the idea of manipulating the treatment, distinguishing between mere observation where $T = 1$ and actively setting $T = 1$ to analyze its causal impact.

Randomized controlled experiments are often impractical or unethical in real-world scenarios. For instance, we cannot randomly assign disabilities to individuals to study its effect on wage. In such cases, *observational causal analysis* provides a viable alternative, allowing for sound causal inference, albeit under additional assumptions. The primary challenge in observational studies arises from *confounding factors*. These are attributes that influence both the treatment assignment and the outcome, thereby obscuring the true causal effect. Consider our example of understanding the causal effect of disability on annual wage (Example 1.2). Unlike a randomized experiment, disability status in this dataset is not randomly assigned. Instead, it is often correlated with other attributes such as gender and age. These attributes act as confounders because they can influence both the likelihood of having a disability and annual wage.

To address confounding variables, Pearl's causal model offers a framework for obtaining unbiased causal estimates by accounting for these confounding attributes, denoted as $\mathbf{Z}$. This approach relies on key assumptions, including the unconfoundedness and overlap assumptions which we do not elaborate on here. Given a set of confounding variables $\mathbf{Z}$, the ATE can be adjusted to the following form:

$$ATE(T, O) = \mathbb{E}_Z \left[\mathbb{E}[O \mid T{=}1, \mathbf{Z}{=}z] - \mathbb{E}[O \mid T{=}0, \mathbf{Z}{=}z]\right] \quad (2)$$

This adjusted ATE, as shown in Equation (2), can be estimated directly from observed datasets.

In this work, we assume that the user provides a sufficient set of confounding variables $\mathbf{Z}$. In practice, such a set can be identified using a causal discovery algorithm (e.g., [31, 87]) and applying criteria such as Pearl's backdoor criterion [62], or alternatively, through domain expertise.

**ATE Estimators.** The ATE in Eq. 2 can be estimated in high-dimensional settings using standard methods such as *linear regression* [24] and *inverse propensity weighting (IPW)* [69]. We focus on these two estimators because their algebraic structure makes them especially amenable to the incremental updates

required for efficient repair. Linear regression models the outcome as a linear function of the treatment and confounders and reads the ATE from the treatment coefficient. IPW estimates the probability of treatment (the *propensity score*, typically via logistic regression) and re-weights each unit by the inverse of this probability to correct for treatment imbalance. While these methods rely on simplifying assumptions, our incremental downdate techniques are estimator-agnostic and can be readily transferred to more sophisticated methods such as doubly robust estimation [10] or double machine learning [18]

## 4 PROBLEM FORMULATION

We consider a single-relation database over a schema $\mathbb{A}$. The schema is a vector of attributes $\mathbb{A}=(A_1, \ldots, A_m)$, where each $A_i$ is associated with a domain $\mathrm{dom}(A_i)$, which can be categorical or continuous. A database instance $D$ populates the schema with a set of tuples of the form $t=(a_1, \ldots, a_m)$ where $a_i \in \mathrm{dom}(A_i)$. We use bold letters to represent a subset of attributes $\mathbf{A} \subseteq \mathbb{A}$, and an uppercase letter to represent a subset of tuples from the database instance $\Gamma \subseteq D$.

Consider a binary[1] treatment variable $T \in \mathbb{A}$, and a numerical outcome variable $O \in \mathbb{A}$. The causal effect, determined by ATE (Equation (2)), of $T$ on $O$ when evaluated on the database instance $D$ is denoted by $ATE_D(T, O)$, where $D$ indicates the causal effect was estimated over the database instance $D$. We assume that the user provides a sufficient set of confounding variables to control for. For simplicity, we do not explicitly refer to them, although they are considered when estimating the causal effect of $T$ on $O$.

If the ATE of $T$ on $O$ is substantially different from the expected causal effect, it indicates a *causal inconsistency* in the database instance $D$. Our objective is to identify a set of tuples $\Gamma \subseteq D$ responsible for this discrepancy, as they are the candidate tuples contributing to the shift.

We assume the presence of a desired causal effect, denoted as $ATE_d$, given by, e.g., a domain expert, relevant literature, or estimated over other dataset versions. Given an error bound $\epsilon > 0$, we assume that $ATE_D(T, O) \notin [ATE_d - \epsilon, ATE_d + \epsilon]$ (i.e., the current ATE is outside the target range). Our goal is to find a minimal size set of tuples $\Gamma \subseteq D$ such that, after removing $\Gamma$ from $D$ we bring the database to a state where the causal effect of $T$ on $O$ is within the range $[ATE_d - \epsilon, ATE_d + \epsilon]$. The set $\Gamma$ is the tuples responsible for the causal inconsistency.

If $\Gamma$ is small relative to the dataset size, this indicates that the ATE can be shifted into the desired range by removing only a small fraction of the data. Examining these tuples allows the user to identify data regions that strongly influence the causal estimate and to assess the sensitivity of the causal conclusions to small changes in the data.

We define the Cardinality Repair for Causal Effect Targeting (abbreviated as CaRET) problem.

PROBLEM 1 (CaRET). *Let $D$ be a database instance over a schema $\mathbb{A}$. We are given a binary treatment variable $T \in \mathbb{A}$, and*

---

[1]For simplicity, we consider a binary treatment variable. However, our framework can be generalized to handle non-binary variables as well.

an outcome variable $O \in \mathbb{A}$. Let $ATE_D(T, O)$ denote the average treatment effect of $T$ on $O$, as estimated over the database $D$. Given a desired causal estimate $ATE_d$ and a threshold $\epsilon > 0$, find a minimal size subset of tuples $\Gamma \subseteq D$ s.t: $ATE_{D \setminus \Gamma}(T, O) \in [ATE_d - \epsilon, ATE_d + \epsilon]$.

Our framework assumes a sufficient set of observed confounding variables, in contrast to classical sensitivity analysis, which explicitly addresses unobserved confounders. Prior work typically models hidden confounding or measurement error using strong parametric bias models (e.g., additive-noise or multiplicative-bias frameworks) while assuming a clean dataset. By contrast, our approach is designed to handle messy data, including duplicates, outliers, and corrupted subsets, thereby relaxing assumptions about data quality.

We consider two data repair settings, both restricted to tuple deletions, with the goal of removing the smallest possible subset. The first setting allows removal of arbitrary tuples, following cardinality repair research for database constraints [2, 56]. The second, more constrained setting, removes an entire subpopulation defined by a pattern, aiming to minimize the size of that subpopulation. For both settings, we analyze the computational complexity of the corresponding optimization problems and propose efficient algorithms to solve them.

*Tuple Removal.* We first focus on cardinality repair [2, 56], which finds a minimal set of tuples to remove to achieve a target causal effect. We show that the CaRET problem is NP-complete via a reduction from the SUBSET-SUM problem [40].

PROPOSITION 4.1. *The CaRET problem is NP-complete.*

In Section 5.1, we introduce an efficient greedy heuristic algorithm to address this problem.

*Pattern Removal.* In this part, we consider the removal of a subpopulation. To specify a subpopulation, we use *patterns* [42, 71, 80, 85] that comprise conjunctive predicates on attribute values. Formally, consider a database $D$ with attributes $T, O,$ $A_1, \ldots, A_n$. Assume the domain of $A_i$ is $\mathrm{dom}(A_i)$. A *conjunctive pattern* $\psi$ is a formula of the form $\bigwedge_{i \in I} A_i = a_i$ where $I \subseteq \{1, \ldots, n\}$ and $a_i \in \mathrm{dom}(A_i)$ for all $i \in I$. We only consider patterns that do not refer to the outcome variable $O$.

Our patterns are restricted to conjunctions of equality predicates, in line with past work on explanations that deem such predicates understandable [4, 25, 70]. The advantage of using such patterns lies in their ability to *explain* the removed subset. We leave for future work the consideration of a richer class of patterns, including continuous variables and disjunction.

For a pattern $\psi = (A_1 = a_1 \wedge \ldots \wedge A_l = a_l)$, we denote by $\psi(D)$ the set of tuples form $D$ that satisfy $\psi$. That is, $\psi(D) = \{t \mid \bigwedge_{i=1}^{l} t[A_i] = a_i, t \in D\}$. In *pattern deletion*, we delete from $D$ the set of tuples $\psi(D)$ for some pattern $\psi$. Our goal is to identify the smallest subpopulation in terms of tuples from $D$, defined by some pattern $\psi$, such that after removing $\psi(D)$ from $D$, the estimated ATE is within the target range. More formally,

PROBLEM 2 (CaRET - PATTERN REMOVAL). *Let $D$ be a database instance over a schema $\mathbb{A}$. We are given a binary treatment*

$T \in \mathbb{A}$, *and an outcome* $O \in \mathbb{A}$. *Let* $ATE_D(T, O)$ *denote the causal effect of $T$ on $O$, as estimated over the database $D$. Given a desired causal estimate $ATE_d$ and a threshold $\epsilon > 0$, find a minimal size subpopulation defined by a pattern $\psi$ such that:*

$$ATE_{D \setminus \psi(D)}(T, O) \in [ATE_d - \epsilon, ATE_d + \epsilon]$$

In order to analyze the complexity of the problems, we implicitly refer to their decision-problem versions, where an upper bound on the size of the removed subpopulation is given. This analogy is sound in the sense that there is a polynomial-time solution to either Problem 1 or Problem 2 if and only if there is a polynomial-time solution for the decision version.

We can show that Problem 2 is NP-complete, by reducing it from the CaRET problem (Problem 1) where any tuple can be deleted from the data (Proposition 4.1).

PROPOSITION 4.2. *Problem 2 is NP-complete.*

We remark that for both problems, membership in NP is witnessed by listing the subpopulation that is removed (for Problem 1) and by the pattern removed (for Problem 2), where we can verify the solution in polynomial time by removing the subpopulation and computing the new ATE. We note that both Problem 1 and Problem 2 remain hard even in the case where there are no confounding variables. In light of this observation, we have the following in particular.

PROPOSITION 4.3. *Problem 1 and Problem 2 are NP-complete for AVG and ATE.*

## 5 THE SUBCURE FRAMEWORK

In this section, we introduce the two underlying algorithms of SUBCURE - one for tuple-based removal (SUBCURE-TUPLE), and the second for the restricted setting where only a subpopulation can be removed (SUBCURE-PATTERN). Both algorithms adopt a heuristic approach to efficiently find a solution, motivated by the fact that the underlying optimization problems are NP-complete. As we demonstrate in our experimental evaluation (Section 7), these algorithms scale well to large, high-dimensional datasets and consistently produce solutions that are smaller in size compared to existing and naive solutions.

*On ILP and SAT formulations.* A natural question is whether our problem could be encoded as an Integer Linear Program (ILP) or SAT instance, similar to the Generalized Deletion Propagation (GDP) framework [52]. However, the key constraint $ATE_D(T, O) \in [ATE_d - \epsilon, ATE_d + \epsilon]$ is inherently nonlinear. Even without confounders, the ATE is defined as the difference between two means, making it nonlinear in the decision variables. When confounders are included, the ATE is estimated, e.g., via linear regression, whose coefficients depend nonlinearly on the dataset (e.g., through matrix inversion). A possible linear relaxation assumes fixed numbers of tuples removed from each subgroup, which makes the constraint linear but requires solving a separate ILP for each combination of removal sizes, an infeasible approach as dimensionality grows. While a reduction to SAT is theoretically possible, it would require

encoding nonlinear arithmetic and is computationally impractical. Therefore, instead of relying on ILP or SAT solvers, our methods focus on a practically scalable algorithmic design. For completeness, we also compare against a naive brute-force baseline that computes the exact optimum (see Section 7).

## 5.1 SubCure for Tuple Removal

A naive brute-force approach would evaluate every possible tuple subset for removal, checking whether the resulting ATE falls within the target range. However, there can be exponentially many ways of shifting the ATE through tuple deletions (as is the case for functional dependencies [45]). To this end, we introduce SubCure-tuple, a greedy, optimized algorithm that iteratively removes tuples based on their estimated influence on the ATE, prioritizing those with the highest impact.

Inspired by "leave-one-out" influence in causal counterfactuals [9, 62], we define the influence of a tuple as follows.

DEFINITION 5.1 (TUPLE INFLUENCE). *The influence of a tuple $t \in D$ on the ATE of $T$ on $O$ is defined as:*

$$\texttt{influence}(t) = ATE_D(T, O) - ATE_{D \setminus \{t\}}(T, O)$$

*where $T$ and $O$ denote the treatment and outcome variables, respectively, and $ATE_D(T, O)$ represents the ATE of $T$ on $O$ computed over the dataset $D$, while adjusting for confounding variables (omitted from the notation for clarity).*

A positive influence score indicates that removing $t$ from $D$ would decrease the ATE, while a negative score indicates that its removal would increase the ATE.

A simple greedy approach iteratively removes the most beneficial tuple -either with the highest positive or most negative influence - based on the current ATE and the target range $[ATE_d - \epsilon, ATE_d + \epsilon]$. After each removal, influence scores are recomputed, and the process repeats until the ATE falls within the target range. To reduce the computational burden of per-tuple ATE influence calculations while preserving structural diversity, we employ a two-stage cluster-based sampling strategy: an initial clustering-based representative selection, followed by an iterative, cluster-aware removal process.

**Clustering.** Using the confounding variables **Z** taken into account for the ATE estimation, the treatment variable $T$, and the outcome $O$, we apply the $k$-means clustering algorithm [49] to form $k$ clusters. If the user does not specify $k$, we set

$$k = \max\left(5, \ \min(\lfloor \sqrt{n} \rfloor, \ \frac{n}{10})\right)$$

where $n = |D|$, ensuring $5 \leq k \leq n$, that is, the number of clusters $k$ is neither too small nor too large.

Within each cluster $C_k$, we select $s$ representative tuples (by default, we set $s = 2$). Let

$$\mathcal{R}_k = \begin{cases} C_k, & |C_k| \leq s \\ \{ r_{k,1}, r_{k,2} \}, & \text{otherwise} \end{cases}$$

where $r_{k,1}$ is the tuple closest to the cluster centroid, and $r_{k,2}$ (and additional reps if $s > 2$) are chosen at evenly spaced percentiles of the distance-to-centroid distribution (e.g. 25th, 75th). We compute the influence score for all representative

tuples, then assign each non-rep point to its nearest representative. This reduces the number of influence computations at each iteration from $n$ to at most $k \cdot s$.

**Iterative Cluster-Based Sampling.** During each iteration $i$, let $\mathcal{A}_i$ denote the set of indices of tuples still in $D$. For each cluster $C_k$, we define the available subset $\mathcal{A}_k^{(i)} = C_k \cap \mathcal{A}_i$ and set $m_k = \min(5, |\mathcal{A}_k^{(i)}|)$.

We draw $m_k$ tuples from the cluster $C_k$ $\{i_{k,1}, \ldots, i_{k,m_k}\} \subseteq \mathcal{A}_k^{(i)}$ uniformly at random. We then compute their influence score and set the cluster score as the average influence score:

$$s_k = \frac{1}{m_k} \sum_{j=1}^{m_k} \texttt{influence}(i_{k,j})$$

Let $d = \text{sign}(ATE_d - ATE_{D^i}(T, O))$, where $ATE_d$ is the target ATE value and $ATE_{D^i}(T, O)$ is the current ATE (at the $i$-th iteration). We select the cluster with the largest $s_k \cdot d$. Finlay, within this cluster, we remove the tuple with the highest influence score and continue to the next iteration.

Additional optimizations we implemented include:
**(1) Sampling**: To improve scalability on large datasets, we operate on a random sample of the data. In our experiments, for large datasets (e.g., the ACS dataset), we sampled 10% of the data. We then run the SubCure-tuple algorithm on this sample to identify tuples for removal. To amplify the effect, we also remove their neighboring tuples in the full dataset, identified using a k-nearest neighbors algorithm using *scikit-learn*'s [63] implementation. Here, $k$ was set to 100.
**(2) Periodic influence recomputation**: We update influence scores every 10 iterations instead of after each removal, significantly reducing runtime while maintaining quality, as a tuple's influence typically changes only slightly per iteration.
**(3) Offline Sampling**: Sampling is performed at preprocessing if the algorithm is expected to run on a large-scale dataset, as it depends only on the attributes used and not on the specific causal query. Thus, the same sampled-subset can be reused across causal questions involving the same attributes.

## 5.2 SubCure for Pattern Removal

Next, we present an efficient algorithm called SubCure-pattern for identifying a small subpopulation, defined by a pattern, whose removal shifts the ATE into a target range. As discussed in Section 4, this optimization problem is NP-complete, which motivates the need for a heuristic approach.

**Algorithm Overview** A naive exhaustive search would enumerate all possible subpopulations (i.e., patterns) to identify the smallest one whose removal shifts the ATE into the desired range. However, this approach is computationally inefficient due to the large number of candidate subpopulations. To address this, our algorithm performs bottom-up random walks over the subpopulation lattice, effectively exploring the space without full enumeration. During these walks, we maintain a dynamic weighting mechanism to reflect the influence of individual predicates on the ATE. This guides the search toward regions of the lattice where predicate combinations are more

likely to find impactful subpopulations. Together, these strategies enable our algorithm to identify small subpopulations with significant influence on the ATE efficiently.

The pseudocode for SubCure-pattern is shown in Algorithm 1. It begins by computing the ATE over the full dataset (line 2), then generates all non-empty subgroups defined by combinations of attribute-values over all attributes in $D$ besides the treatment and outcome (line 3). This is done in a single pass over the data, inspired by the first step of the Apriori algorithm [5]. These most specific patterns correspond to the leaves of the pattern lattice and serve as starting points for random walks. Each walk starts from a leaf node (line 5) and iteratively removes one predicate at a time (line 11). For each subgroup, the algorithm checks whether its removal moves the ATE into the target range. If such a subgroup is found, it is returned immediately; otherwise, the walk continues. To limit runtime, the algorithm performs at most $k$ random walks.

**Dynamic Weighting Mechanism**: Instead of uniformly selecting a predicate to remove from a pattern (line 11 in Algorithm1), we maintain a caching mechanism that tracks how often the removal of a predicate shifts the ATE toward the desired direction, as well as the magnitude of the shift. This information is then used to assign probabilities to each predicate, prioritizing those whose removal is more likely to guide the ATE into the target range.

Additional optimizations we implemented include:
**(1) Early termination of random walks** if the current subgroup exceeds a predefined size threshold $\tau$, prompting the algorithm to abandon that path and initiate a new random walk to find small-size subpopulations.
**(2) Caching**: Previously evaluated patterns are cached to reduce redundant computation.
**(3) Sampling:** To ensure scalability on large datasets, we uniformly sample 10% of the data and run the algorithm on this subset. The identified pattern is then applied to remove all matching tuples from the full dataset. This optimization is used only for large datasets (e.g., ACS in Section 7).

## 6 INCREMENTAL ATE UPDATES

The main computational bottleneck of SubCure-tuple and SubCure-pattern lies in the repeated ATE calculation. To this end, we propose optimizations that estimate the ATE after removing data subsets without retraining the model. We focus on two widely used ATE estimators: (1) linear regression [36] and (2) inverse propensity weighting (IPW) [36]. For linear regression, we leverage its closed-form solution to provide both exact and faster approximate ATE updates. For IPW, we build upon an existing machine unlearning technique [51], to avoid full model retraining and yield updated propensity scores in constant time per deletion.

### 6.1 ATE Update for Linear Regression

We begin with the standard linear regression model:

$$\mathbf{o} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} \tag{3}$$

---

**Algorithm 1:** The SubCure-pattern Algorithm

**input** : A dataset $D$ with a treatment variable $T$ and an outcome $O$, a target ATE range defined by $ATE_d$ and $\epsilon$, and a number $k > 0$.

**output:** A subgroup to be removed defined by a pattern $\psi$ or indication of a failure.

1 /* Initial ATE value                       */
2 $v \leftarrow$ EstimateATE $(T, O, D)$
3 $\mathcal{G} \leftarrow$ GetMostSpecificGroups $(T, O, D)$
4 **for** $i \in [1, k]$ **do**
5     $\psi_i \leftarrow$ a random group from $\mathcal{G}$
6     $v_i \leftarrow$ EstimateATE $(T, O, D \setminus (\psi_i(D)))$
7     **if** $v_i \in [ATE_d - \epsilon, ATE_d + \epsilon]$ **then**
8        **return** $\psi_i$
9     **while** $\psi_i$ *is not empty* **do**
10        /* A random step on the pattern lattice       */
11        $\psi_i \leftarrow$ RemovePredicate $(\psi_i)$
12        $v_i \leftarrow$ EstimateATE $(T, O, D \setminus (\psi_i(D)))$
13        **if** $v_i \in [ATE_d - \epsilon, ATE_d + \epsilon]$ **then**
14           **return** $\psi_i$

15 **return** No solution was found

---

where $X \in \mathbb{R}^{n \times m}$ is the design matrix (one row per data record), $\mathbf{o} \in \mathbb{R}^n$ is the vector of outcomes, $\boldsymbol{\beta} \in \mathbb{R}^m$ are the model coefficients, and $\boldsymbol{\varepsilon}$ is a noise vector. The ATE of a treatment $T$ on $O$ is defined as the coefficient of $T$, where $T \in \mathbf{X}$, and $\mathbf{X}$ also includes all the confounding variables $\mathbf{Z}$. To learn the model parameters $\boldsymbol{\beta}$, we use ordinary least squares, which minimizes the squared error between the predicted and actual values. This yields the following equation:

$$X^\top X \boldsymbol{\beta} = X^\top \mathbf{o} \tag{4}$$

This gives a closed-form expression for the optimal coefficients:

$$\boldsymbol{\beta} = (X^\top X)^{-1} X^\top \mathbf{o} \tag{5}$$

Now, suppose we remove a subset $\Gamma \subseteq D$ of $r$ tuples, namely $|\Gamma| = r$. Let $X_{\text{rmv}} \in \mathbb{R}^{r \times m}$ and $\mathbf{o}_{\text{rmv}} \in \mathbb{R}^r$ denote the corresponding removed tuples of $X$ and $\mathbf{o}$. The updated data becomes:

$$X_{\text{new}} = X \setminus X_{\text{rmv}}, \quad \mathbf{o}_{\text{new}} = \mathbf{o} \setminus \mathbf{o}_{\text{rmv}}$$

Our goal is to obtain $\beta_{\text{new}}$ defined by

$$\beta_{\text{new}} = (X_{\text{new}}^\top X_{\text{new}})^{-1} X_{\text{new}}^\top \mathbf{o}_{\text{new}} \tag{6}$$

*without* recomputing the inverse from scratch.

For convenience, define

$$A := X^\top X, \qquad \Delta := X_{\text{rmv}}^\top X_{\text{rmv}}, \qquad \text{so that} \qquad A_{\text{new}} = A - \Delta$$

**Exact Update.** To get an exact update, we use the Woodbury matrix identity [34] to get the coefficient update:

$$\beta_{\text{new}} = \left[ A^{-1} + A^{-1} U (I_r - U^\top A^{-1} U)^{-1} U^\top A^{-1} \right] (X^\top \mathbf{y} - X_{\text{rmv}}^\top \mathbf{o}_{\text{rmv}}) \tag{7}$$

***Approximate Update.*** When the number of tuples to be removed is small, specifically, when the norm $\|\Delta A^{-1}\| < 1$, we can use the *Neumann series* [84] to approximate the inverse of the matrix $A - \Delta$. By setting $Z = \Delta A^{-1}$, we approximate $(I - \Delta A^{-1})^{-1}$ by summing a finite number of terms:

$$(A - \Delta)^{-1} \approx A^{-1} \sum_{k=0}^{K} (\Delta A^{-1})^k$$

We set $K = 1$, taking the first two terms:

$$\beta_{\text{new}} \approx \left(A^{-1} + A^{-1}\Delta A^{-1}\right)\left(X^{\top}\mathbf{o} - X_{\text{rmv}}^{\top}\mathbf{o}_{\text{rmv}}\right)$$

This method lets us efficiently update the solution without recalculating the full inverse. As we show in Section 7, this method is faster than the exact method and gives accurate results when the number of tuples we remove is small.

## 6.2 ATE Update for IPW

To estimate ATE using IPW method, each tuple is weighted by the inverse probability of receiving the treatment they received, as estimated by a propensity score model. The ATE is computed as the difference in the weighted averages of outcomes between treated and control groups. Here, we build upon an unlearning method for logistic regression [51], proposing an incremental update method for the IPW ATE estimator.

We define $\mathbf{t}_i$ and $\mathbf{o}_i$ as the treatment and outcome values of the $i$-th unit (tuple), and $z_i$ as its assignment for the confounding variables $\mathbf{Z}$. Let $p_i = \Pr[\mathbf{t}_i = 1 \mid \mathbf{z}_i; \boldsymbol{\theta}] = \sigma(\mathbf{z}_i^{\top}\boldsymbol{\theta})$ be the propensity score of the $i$-th unit obtained from a $\ell_2$-regularised logistic model with parameters $\boldsymbol{\theta}$, where $\mathbf{z}_i$ are the confounder assignments of the $i$-th unit. For a dataset $D$ where $D| = n$, the IPW estimator estimate the ATE of $T$ on $Y$ as follows:

$$\widehat{\text{ATE}}_{\text{IPW}}(T, O) = \frac{\sum_{i=1}^{n} \frac{\mathbf{t}_i\mathbf{o}_i}{p_i}}{\sum_{i=1}^{n} \frac{\mathbf{t}_i}{p_i}} - \frac{\sum_{i=1}^{n} \frac{(1 - \mathbf{t}_i)\mathbf{o}_i}{1 - p_i}}{\sum_{i=1}^{n} \frac{(1 - \mathbf{t}_i)}{1 - p_i}} \qquad (8)$$

Suppose a subset $D_{\text{rmv}} \subseteq D$ is removed, where $|D_{\text{rmv}}| = r$. We need to update $\boldsymbol{\theta}$ without retraining the model from scratch and re-evaluate Eq. (8) efficiently. To achieve this, we use the Fisher mini-batch incremental update algorithm from [33]. The pseudocode is provided in the Appendix. This algorithm takes as input the dataset, the current model parameters, and the subset to be removed, and returns an estimate of the updated parameters. Rather than performing multiple iterations of gradient descent, the algorithm relies on a single-step update. Finally, we use the updated parameters to recompute propensity scores and re-evaluate the ATE (Eq. (8)).

## 7 EXPERIMENTAL STUDY

We present an experimental evaluation that evaluates our SUB-CURE effectiveness. We aim to address the following questions: **Q1:** How does the quality of our proposed solutions compare to naive baselines and existing methods? **Q2:** How well do our algorithms scale when applied to large, high-dimensional datasets? **Q3:** To what extent does our proposed ATE incremental update optimizations improve performance? **Q4:** Is

SUBCURE effective even when combined with alternative methods for computing ATE? **Q5:** How does SUBCURE compare to existing sensitivity analysis methods?

### 7.1 Experimental setup

All experiments were conducted on a server with an Intel(R) Xeon(R) E5-2680 v3 CPU (2.50GHz, 12 cores) and 128GB of RAM. Our algorithms were implemented in Python3, and our code and the used datasets are publicly available in [7].

By default, we focus on the linear regression estimator for ATE computation, as described in Section 6.1.

*Datasets & examined causal questions.* We examine four commonly used datasets. The datasets' statistics and corresponding causal questions are summarized in Table 1. **German Credit:** This data [8] contains details of bank account holders, including demographic and financial information, and credit risk scores. The causal query measures the impact of house ownership on the risk score. The goal is to shift this effect to zero (within a margin of ±0.01), indicating no causal relationship. **Twins**: The Twins dataset [47] is a widely used benchmark for causal inference, containing data on same-sex twins. The treatment is being the heavier twin, and the outcome is mortality within the first year of life. The initial estimated ATE is -0.016, and our objective is to shift this effect to zero (±0.001). **Stack Overflow:** The Stack Overflow Developer Survey [1] data contains responses from developers worldwide, covering topics such as professional experience, education, and salary-related information. The examined causal question is to quantify the effect of higher education on annual salary. The initial ATE is $13,236. The goal is to reduce this effect by $5,000 (± $100), resulting in a target ATE of $8,236. **ACS:** The American Community Survey (ACS) [79] dataset is a nationwide survey conducted by the U.S. Census Bureau, providing detailed demographic, social, and economic data. The causal question examines the effect of not having a disability on annual wages. The initial ATE is $8,774, meaning people without disabilities earn that much more on average compared to people with disability. The target ATE is set to $12,000 (±$500).

*Competing Baselines.* We consider the following baselines:
**OPT**: We compare our algorithms against the optimal solution to analyze their quality and runtime trade-offs. We implemented OPT-tuple and OPT-pattern using exhaustive search to obtain the optimal solution.
**SINGLE-UPDATE-SUBCURE-TUPLE**: This baseline evaluates the impact of each tuple in the data on the ATE and iteratively removes those with the highest influence until the target ATE is achieved. Unlike SUBCURE-TUPLE, the influence of each tuple is calculated only once at the start of the process.
**SCODED** [83]: SCODED is designed to identify the top-$k$ tuples that contribute most to the violation of a (in)dependence relationship between sets of variables. We use SCODED to identify tuples for removal in order to achieve a desired ATE. We evaluate both versions of SCODED, exact, and approximate.

**Table 1: Details of the datasets and corresponding causal questions we use for experiments and case studies.**

| Dataset | #Tuples | #Atts | Treatment | Outcome | Confounding Variables | Org. ATE | Tar. ATE |
|---|---|---|---|---|---|---|---|
| German | 1000 | 17 | Owning a house | Credit risk | Personal status, age | 0.13 | 0 (±0.01) |
| Twins | 23,968 | 53 | Heavier twin | Mortality | Gestational age, birth weight, prenatal care, abnormal amniotic fluid, induced labor, gender, maternal marital status, year of birth, and total previous deliveries. | -0.016 | 0 (±0.001) |
| SO | 47,702 | 21 | High Education | Annual Salary | Continent, gender, ethnicity | 13236 | 8236 (±100) |
| ACS | 1,188,308 | 17 | Not having a disability | Annual wage | Education, public health coverage, private health coverage, medicare for people 65 and older, insurance through employer, gender, age | 8774 | 12000 (±500) |

The value of $k$ is set to match the number of tuples removed by Single-Update-SubCure-tuple.

**ZamInfluence** [15]: ZamInfluence measures how much an estimate changes when a data fraction is removed, identifying deletions that flip its sign or significance. We compare against it in scenarios aiming to push the ATE toward 0.

**Outlier Detection Baselines**: We compare SubCure with two widely used outlier detection methods: **Isolation Forest** (IF) [43] and **Local Outlier Factor** (LOF) [17]. We evaluate these algorithms in two ways: first, as standalone baselines by considering the ATE after their application; second, as preprocessing steps before SubCure-tuple, reporting the total number of removals needed to reach the target ATE. We use scikit-learn's implementation [63], setting the removal budget to match the scale of Single-Update-SubCure-tuple.

For our SubCure-tuple and SubCure-pattern algorithms, as well as for the Single-Update-SubCure-tuple baseline, we evaluate both their exact and approximate variants based on the ATE incremental update methods described in Section 6.1.

The time cutoff for all algorithms was set to 10 hours. For the SubCure-pattern algorithm, we limited the size of the group to be removed to 20% of the data, and the number of random walks considered was set to 1000.

### 7.2 Quality Evaluation

For each examined causal question, we ran all algorithms. For SubCure-pattern, which is a randomized algorithm, we ran it three times and report the average runtime and the best result achieved. In all cases, the outlier detection methods (IF and LOF) failed to identify tuples that significantly affect the ATE. As a result, their standalone results are omitted from presentation; we report only their performance when used as a preprocessing step for our SubCure-tuple algorithm. The results are shown in Table 2. In the Appendix, we show heatmaps indicating the overlap between different solutions.

**Results Summary**: When only a few hundred tuples need to be removed, Single-Update-SubCure-tuple and SubCure-tuple yield nearly identical solutions, with Single-Update-SubCure-tuple being faster. As the required removals increase, Single-Update-SubCure-tuple degrades, while SubCure-tuple finds smaller solutions, highlighting the importance of recomputing influence scores. ZamInfluence identified small subsets to shift the ATE to zero, but cannot target other ATE values. SubCure-pattern consistently identified compact influential subpopulations but removed more tuples overall. The approximate update method was faster and comparable to the

exact one for small removals but less effective for larger ones. SCODED failed to reach the target ATE in all cases, and IF and LOF removed tuples with minimal ATE impact, offering no benefit even as preprocessing for SubCure-tuple.

German Credit: Single-Update-SubCure-tuple and SubCure-tuple identified the same 42 tuples (4% of the data) whose removal reduced the ATE from 0.13 to 0.007 - a 94.6% drop, with no difference between exact and approximate update methods. SubCure-pattern failed to find any subgroup smaller than 20% of the data that achieved the target ATE, though it discovered that removing a specific group (foreign workers with low checking balances and long residence) reduced the ATE by over 50%. Outlier detection using IF and LOF did not help reduce the number of tuples SubCure-tuple needed to remove, as the same 42 deletions were still required. Tuples identified by SCODED and ZamInfluence did affect the ATE, but their removal alone did not bring it into the target range.

Twins: Both Single-Update-SubCure-tuple and SubCure-tuple achieved the target ATE, but SubCure-tuple did so more efficiently, removing only 270 tuples compared to 367, underscoring the importance of recalculating influence scores. Approximate and exact ATE update methods yielded similar results. Here, ZamInfluence found the smallest solution that achieves the target ATE, but it is only applicable when the target ATE is zero. SubCure-pattern found relatively small subpopulations that achieved the target ATE, with the approximate variant identifying a group – first-born twins with birth weights between 1,786g and 1,999g – about 18% of the data whose removal reversed the ATE sign from negative to positive, increasing it by 102%. SCODED again failed to identify a solution, and outlier detection methods did not reduce the number of required deletions, indicating that the tuples selected by our method are not merely outliers.

Stack Overflow: Here again Single-Update-SubCure-tuple and SubCure-tuple identified the same 67 tuples (just 0.14% of the data), whose removal decreases the ATE by 37.5%. SubCure-pattern successfully identified impactful subpopulations; notably, the exact variant found that removing 3,744 individuals (7.8% of the data) – male respondents from the US aged 25-34, not students, with no dependents – lowered the ATE by 38%. SCODED again failed to find tuples that shift the ATE into the target range, with its approximate variant even increasing the ATE. Here again using either IF or LOF as a preprocessing step for SubCure-tuple had no effect on reducing its workload.

**Table 2: Comparison of baseline methods. Results include the achieved ATE, whether the target range was met, the number of tuples removed, and runtime. Highlighted results correspond to the fewest tuples removed.**

| Dataset | Baseline | Achieved ATE | Hit Range | #Removals | Time (s) |
|---|---|---|---|---|---|
| German Credit (org. ATE: 0.13, tar. ATE: 0 ±0.01) | SubCure-tuple (exact) | 0.006 (↓ 95.4%) | ✓ | **42** (4.2%) | 2.8 |
| | SubCure-tuple (approx) | 0.007 (↓ 94.6%) | ✓ | **42** (4.2%) | 2.78 |
| | Single-Update-SubCure-tuple (exact) | 0.006 (↓ 95.4%) | ✓ | **42** (4.2%) | 0.36 |
| | Single-Update-SubCure-tuple (approx) | 0.007 (↓ 94.6%) | ✓ | **42** (4.2%) | 0.36 |
| | SCODED (exact) | 0.02 (↓ 84.6%) | ✗ | **42** (4.2%) | 0.08 |
| | SCODED (approx) | 0.08 (↓ 38.5%) | ✗ | **42** (4.2%) | 0.03 |
| | ZamInfluence | -0.025 (↓ 292%) | ✗ | 52 (5.2%) | 0.02 |
| | IF + SubCure-tuple | 0.008 (↓ 94.1%) | ✓ | 88 (8.8%) | 2.8 |
| | LOF + SubCure-tuple | 0.008 (↓ 94.1%) | ✓ | 89 (8.9%) | 3.2 |
| | SubCure-pattern (exact) | 0.06 (↓ 53.8%) | ✗ | 126 (12.6%) | 50 |
| | SubCure-pattern (approx) | 0.07 (↓ 46.1%) | ✗ | 196 (19.6%) | 28 |
| Twins (org. ATE: -0.016, tar. ATE: 0 ±0.001) | SubCure-tuple (exact) | -0.0009 (↑ 94.4%) | ✓ | 270 (1.1%) | 92.7 |
| | SubCure-tuple (approx) | -0.0009 (↑ 94.4%) | ✓ | 270 (1.1%) | 99.9 |
| | Single-Update-SubCure-tuple (exact) | -0.0009 (↑ 94.4%) | ✓ | 367 (1.5%) | 18.5 |
| | Single-Update-SubCure-tuple (approx) | -0.0009 (↑ 94.4%) | ✓ | 367 (1.5%) | 18.6 |
| | SCODED (exact) | -0.005 (↑ 68.7%) | ✗ | 367 (1.5%) | 4.4 |
| | SCODED (approx) | -0.015 (↑ 6.2%) | ✗ | 367 (1.5%) | 4.02 |
| | ZamInfluence | -0.0005 (↑ 99.7%) | ✓ | **180** (0.75%) | 0.06 |
| | IF + SubCure-tuple | -0.0009 (↑ 94.4%) | ✓ | 714 (2.9%) | 93 |
| | LOF + SubCure-tuple | -0.0009 (↑ 94.4%) | ✓ | 714 (2.9%) | 88 |
| | SubCure-pattern (exact) | -0.0005 (↑ 96.8%) | ✓ | 4692 (19.5%) | 220 |
| | SubCure-pattern (approx) | 0.0003 (↑ 101.8%) | ✓ | 4401 (18.3%) | 441 |
| Stack Overflow (org. ATE: $13.2k, tar. ATE: $8.2k ±$100) | SubCure-tuple (exact) | 8269 (↓ 37.5%) | ✓ | **67** (0.14%) | 45.1 |
| | SubCure-tuple (approx) | 8269 (↓ 37.5%) | ✓ | **67** (0.14%) | 43.7 |
| | Single-Update-SubCure-tuple (exact) | 8269 (↓ 37.5%) | ✓ | **67** (0.14%) | 19.6 |
| | Single-Update-SubCure-tuple (approx) | 8269 (↓ 37.5%) | ✓ | **67** (0.14%) | 19.8 |
| | SCODED (exact) | 12625 (↓ 4.6%) | ✗ | **67** (0.14%) | 17.7 |
| | SCODED (approx) | 16255 (↑ 22.6%) | ✗ | **67** (0.14%) | 10.8 |
| | IF + SubCure-tuple | 8274 (↓ 37.4%) | ✓ | 465 (0.95%) | 58 |
| | LOF + SubCure-tuple | 8290 (↓ 37.3%) | ✓ | 440 (0.94%) | 44 |
| | SubCure-pattern (exact) | 8250 (↓ 37.6%) | ✓ | 3744 (7.8%) | 115 |
| | SubCure-pattern (approx) | 8140 (↓ 38.5%) | ✓ | 5673 (11.9%) | 141 |
| ACS (org. ATE: $8.7k, tar. ATE: $12k ±$500) | SubCure-tuple (exact) | 11512 (↑ 31.1%) | ✓ | **8400** (0.7%) | 1149 |
| | SubCure-tuple (approx) | 11510 (↑ 31.1%) | ✓ | **8400** (0.7%) | 1102 |
| | Single-Update-SubCure-tuple (exact) | 11503 (↑ 31.1%) | ✓ | 9600 (0.8%) | 649 |
| | Single-Update-SubCure-tuple (approx) | 11501 (↑ 31.1%) | ✓ | 9600 (0.8%) | 647 |
| | IF + SubCure-tuple | 11583 (↑ 31.8%) | ✓ | 15582 (1.42%) | 867 |
| | LOF + SubCure-tuple | 11561 (↑31.7%) | ✓ | 19582 (1.64%) | 2072 |
| | SubCure-pattern (exact) | 11869 (↑ 35.3%) | ✓ | 141168 (11.9%) | 229 |
| | SubCure-pattern (approx) | 13629 (↑ 55.3%) | ✗ | 211497 (17.8%) | 3063 |

**Table 3: Quality results for using the incremental update optimization for IPW.**

| Dataset | Algorithm | Achieved ATE | # Removals | Time |
|---|---|---|---|---|
| German Credit | SubCure (no incremental update) | 0.009 (↓ 95.2%) | 35 (3.5%) | 35.8 |
| | SubCure | 0.009 (↓ 95.2%) | 35 (3.5%) | 25.5 |
| Twins | SubCure (no incremental update) | 0 (↑ 100%) | 300 (1.2%) | 832 |
| | SubCure | 0 (↑ 100%) | 300 (1.2%) | 712 |

ACS: Here, both Single-Update-SubCure-tuple and SubCure-tuple achieved similar target ATE values, but SubCure-tuple required fewer deletions: 8.4k tuples vs 9.6k, demonstrating its advantage when large removal is needed, as it re-evaluates tuple influence during the process while Single-Update-SubCure-tuple does not. The exact SubCure-pattern algorithm identified a meaningful subpopulation (11.9% of the data) whose removal raised the ATE by 35%, while the approximate version overshot the target range due to accumulated error from removing a large number of tuples (221k). SCODED results are omitted, as it can be applied only when the goal is to push the ATE towards zero, indicating an independence relationship.

As discussed in Example 1.2, achieving a comparable downward shift of the ATE with SubCure-tuple requires removing only 7,200 tuples (0.6% of the data). In contrast, SubCure-pattern can only find a solution that removes 889,626 tuples (74% of the data), resulting in an ATE of 5,900 (a 32.1% decrease). This subpopulation consists of individuals with no disability who hold Medicare. This asymmetry reveals an important insight: while the ATE is highly sensitive to upward shifts, it appears significantly more robust to downward ones when the intervention involves removing entire subpopulations.

*7.2.1 Estimating ATE with IPW.* We revisit our case study using IPW to estimate the ATE and apply the corresponding incremental update optimization. Since ATE update with IPW is significantly slower than with the linear model, we focus on the smaller German Credit and Twins datasets, as others exceed the time limit. Results are shown only for SubCure-tuple, as SubCure-pattern and Single-Update-SubCure-tuple show similar trends. Table 3 shows that SubCure-tuple yields very similar results across the two ATE estimators and removes almost the same tuples. For example, in German Credit, it removes 35 tuples (compared to 42 with the linear model), achieving nearly identical target ATEs. This suggests that tuples influential under one estimator are typically influential under the other - a trend further supported in Section 7.5. We also observe that using the optimization reduces runtime by approximately 20%, compared to not using the IPW incremental update. However, the linear regression-based optimization is consistently faster and produces similar outcomes.

*7.2.2 Comparison to the optimal solution.* Running an exhaustive brute-force search to obtain the optimal solution proved infeasible even for the small German Credit dataset. To nevertheless assess how closely SubCure approximates the optimal solution, we designed controlled experiments using synthetic data where an upper bound on the number of deletions is known, thus bounding the search space OPT. We generated synthetic datasets of increasing sizes (up to 50 tuples for OPT-tuples, as larger instances exceeded our time cutoff), each containing a treatment variable, an outcome,

and three confounders. The target ATE was defined as the ATE computed on each dataset with $\epsilon = 0$. We then injected 20% noisy records into the data, making the number of noisy tuples an upper bound on the solution size. We compared SubCure-tuple against OPT-tuple, and SubCure-pattern against OPT-pattern. The results are depicted in Figure 1. We observe that SubCure-tuple achieves solutions comparable in quality to those of the OPT-tuple, while being significantly faster. For SubCure-pattern, the quality of the produced solutions is also comparable to that of OPT-pattern, but the runtime differences are smaller. This behavior stems from the experimental setup: in this synthetic data, there was a single dominant optimal solution to identify, rather than multiple similar solutions, as in the other experiments. Consequently, SubCure-pattern required more time to find this solution, and in some cases, it failed to find it. In real-world scenarios, however, where numerous subpopulations may yield comparable effects, the heuristic nature of SubCure-pattern allows it to find a good solution within a limited time, albeit without formal guarantees on optimality.

We also compared SubCure-tuple and Single-Update-SubCure-tuple on larger synthetic datasets. Both algorithms produced solutions smaller than the upper bound for up to 1.3k noisy tuples. However, as the number of noisy tuples increased - requiring more deletions to reach the target ATE - the performance of Single-Update-SubCure-tuple degraded, since its single influence-score computation became outdated. Full details appear in the Appendix.

*7.2.3 Robustness to data quality issues.* To assess the robustness of SubCure to data quality issues, we conducted experiments that inject controlled noise into the datasets, following a procedure similar to prior work [65]. We consider three types of common data quality issues: outliers, duplicates, and missing values (with zeros used as placeholders), and varied the noise level from 5% to 55%. We evaluated the impact of these perturbations on two datasets, German Credit and Twins, measuring how the performance of SubCure degrades as the level of noise increases. This setup allows us to identify the point at which the system's behavior begins to deteriorate, thereby quantifying its robustness to common data quality issues. The results are depicted in Figure 2. Our results show that overall, both algorithms remain stable under varying levels and types of noise, demonstrating strong robustness. For duplicates and missing values, SubCure-tuple maintains relatively consistent performance as noise increases, indicating resilience to these perturbations. For outliers, it is somewhat less stable, since this type of noise directly affects the outcome values (by design) and, consequently, the ATE, requiring more tuples to be removed to reach the target ATE. SubCure-pattern is generally less stable across noise levels, and in some cases, failed to find a solution altogether. However, when a solution was found, its size was comparable to that obtained in the noise-free setting. This behavior is expected: because SubCure-pattern operates at the pattern level rather than

the tuple level, it cannot selectively remove corrupted tuples that disproportionately influence the ATE.

## 7.3 Efficiency Evaluation

Next, we evaluate how parameters affect algorithm runtime. Each experiment is repeated three times to account for randomization, and we report the average. We vary one parameter at a time (data size, confounders, or target ATE) while keeping others fixed. Results for the Twins dataset are omitted due to space, but similar trends were observed.
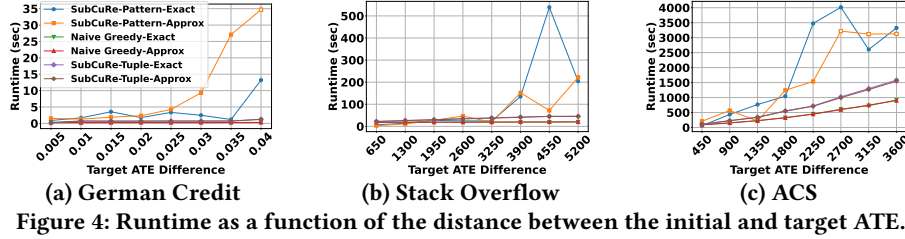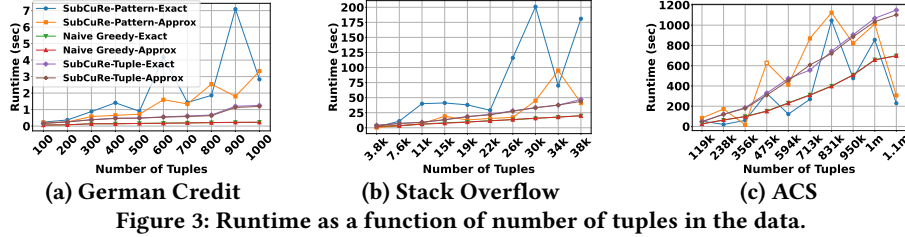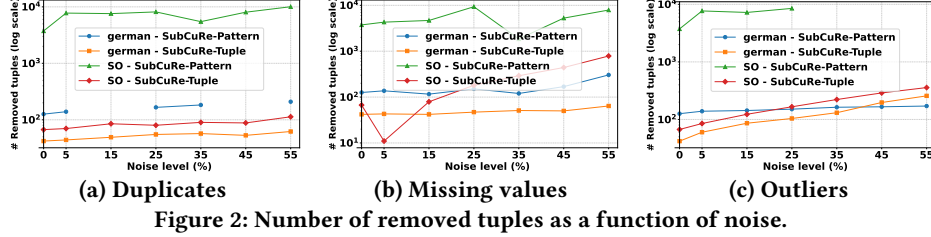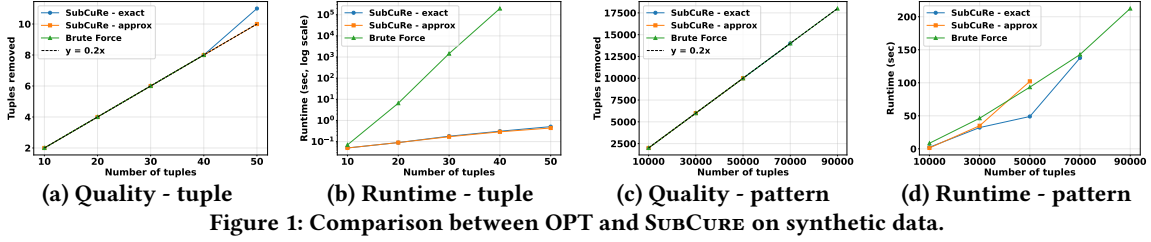
Runtime vs. number of tuples: We analyze how the number of tuples in the dataset impacts runtime. To do so, we randomly sample data subsets of increasing size. The results are shown in Figure 3. The runtime of the tuple-based algorithms (SubCure-tuple and Single-Update-SubCure-tuple) increases approximately linearly with the number of tuples, as more influence scores need to be computed. For SubCure-pattern, the runtime is less sensitive to data size. After processing around 40% of the data, most patterns have already been considered. The spikes are due to the randomized nature of the algorithm: in some runs, it quickly finds a solution, while in others, it takes longer. Also, as discussed in Section 7.2, when the size of the subset to be removed is relatively big (as in ACS and Stack Overflow), the approximate version is less accurate and thus it takes longer to find a subgroup satisfying the condition.

Runtime vs. target ATE: We examine how the distance between the initial ATE and the target ATE affects runtime. To this end, we gradually increase the ATE shift. The results are presented in Figure 4. Empty markers for SubCure-pattern indicate that no solution was found. As the gap between the initial and target ATE increases, the runtime of all algorithms grows accordingly. This is because a larger number of tuples must be removed to reach the target, leading to more iterations for SubCure-tuple and Single-Update-SubCure-tuple, and forcing SubCure-pattern to explore larger subpopulations.

Runtime vs. number of confounders: We examine how the number of confounders affects runtime by randomly selecting confounder subsets. As more variables are included in the ATE estimation, runtimes increase across all algorithms. Full details appear in the Appendix.

## 7.4 Ablation Study

We evaluate the effect of our ATE incremental update optimizations (Section 6.1) by comparing against versions of our algorithms that recompute the ATE from scratch. We focus on the smaller German Credit and Twins datasets, as the unoptimized algorithms exceeded the time limit on larger datasets. The unoptimized algorithms produced results nearly identical to their optimized counterparts. For example, on German Credit, the unoptimized SubCure-tuple (exact) selected the same 42 tuples but ran 18× slower (51s vs. 2.8s). Similarly, SubCure-pattern (exact) removed a slightly different subpopulation (132 vs. 126 tuples), reaching a nearly identical ATE (0.064 vs. 0.062), but with 2.3× longer runtime (118s vs.

**Figure 1: Comparison between OPT and SUBCURE on synthetic data.**



**Figure 2: Number of removed tuples as a function of noise.**



**Figure 3: Runtime as a function of number of tuples in the data.**



**Figure 4: Runtime as a function of the distance between the initial and target ATE.**

50s). As shown in Figure 5, runtime differences grow with data size. For SUBCURE-TUPLE, the optimizations yield up to a 100× speedup. For SUBCURE-PATTERN, the gain is around 10× on German Credit, where small subpopulations are removed (up to 200 tuples). On Twins, where larger groups are considered (of size up to thousands of tuples), the optimizations are less effective and can even slow performance. This shows that incremental update optimizations preserve solution quality while improving runtime, especially when only a small number of tuples are removed. When larger portions of the data are removed, the overhead can outweigh the benefits.

## 7.5 Generalizability of SUBCURE

Next, we show that SUBCURE-TUPLE consistently identifies influential tuples that affect the ATE, even when alternative estimation methods are used. Specifically, we evaluate its effectiveness with two widely used estimators: Doubly Robust
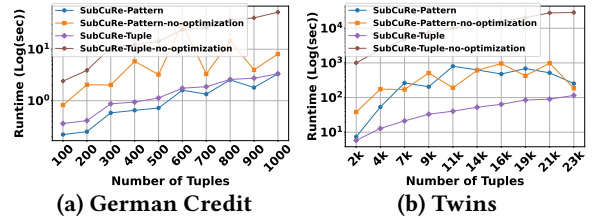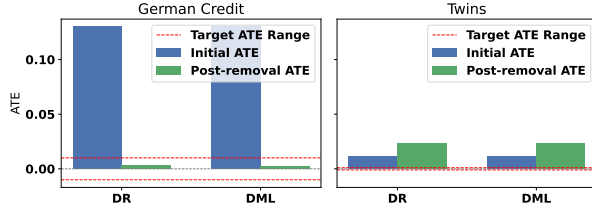


**Figure 5: Runtime as a function of number of tuples - with and without ATE incremental update optimizations**

(DR) [10] and Double Machine Learning (DML) [18]. We conduct experiments on the German Credit and Twins datasets. For each scenario, we first compute the initial ATE using the DR and DML estimators[2]. Then, we apply SUBCURE-TUPLE (with linear regression or IPW estimators) to identify influential tuples. After removing the selected tuples, we recompute the ATE using the same estimator (DR or DML). Our results

---

[2]using the implementation of EconML [57].

**Figure 6: ATE estimates before and after removing tuples identified by SUBCURE-TUPLE (using linear regression update) from the German Credit and Twins datasets under two ATE estimators: DR and DML.**

demonstrate that SUBCURE-TUPLE reliably finds tuples whose removal shifts the ATE in the desired direction, across both DR and DML. The results while using linear regression with the exact update optimization are shown in Figure 6 In German Credit, the initial ATE is 0.13 using all estimators. After removing tuples selected by SUBCURE-TUPLE- using either linear regression or IPW - the ATE, when recomputed with DR and DML, moves into the target range of $(0 \pm 0.01)$. In Twins, the initial ATE estimated by DR and DML is positive (0.011), while the ATE estimated by linear regression and IPW is negative (-0.016). Consequently, SUBCURE-TUPLE selects tuples whose removal causes an increase in ATE. Although the post-removal ATE computed with DR and DML does not reach the target range $(0 \pm 0.001)$, the shift in ATE is in the correct direction, consistent with the effect predicted by SUBCURE-TUPLE. *These results highlight the generalizability of SUBCURE-TUPLE: its selected tuples influence ATE estimates consistently, even when evaluated using different estimation strategies.*

### 7.6 Relation to Sensitivity Analysis

We demonstrate that SUBCURE provides a complementary perspective on ATE sensitivity compared to existing methods. We focus on two closely related approaches (discussed in detail in Section 2): (i) KonFound [26], which quantifies the degree of hidden bias required to overturn an inference, and (ii) WTE [38], which measures robustness under confounder shift by estimating the worst-case ATE across subpopulations of a given size. As a case study, we evaluate both methods on relevant datasets, report their results, and highlight their differences from SUBCURE-TUPLE. We exclude SUBCURE-PATTERN, as both compared methods operate at the individual tuple level.

Konfound: KonFound evaluates sensitivity to unobserved confounders by identifying the "switch point" at which an effect becomes null, akin to pushing the ATE toward 0 (effectively nullifying the causal effect). Thus, we only evaluate this method on the German Credit and Twins datasets, where the goal is to push the ATE towards 0. In the German dataset, Konfound estimated that over 50% of the data would need to be replaced to nullify the effect, while in Twins the estimate was 48%. In contrast, SUBCURE-TUPLE identified just 4.2% of the data in German and 1.1% in Twins for removal to achieve no effect. This indicates that *while results are generally robust to unobserved confounders, they can be highly sensitive to*

*specific observed data points*, highlighting the complementary perspective SUBCURE-TUPLE provides for sensitivity analysis.

WTE For WTE, we set the subpopulation size threshold to match the number of tuples removed by SUBCURE-TUPLE in each dataset. Across all datasets, WTE shows minimal change in ATE, indicating strong robustness to confounder distribution shifts. In the German Credit dataset, SUBCURE-TUPLE reduces the ATE by 95%, whereas WTE changes it by only 3%. In Twins, SUBCURE-TUPLE reduces the ATE magnitude by 94%, while WTE shifts it by less than 1%. On Stack Overflow, SUBCURE-TUPLE decreases the ATE by 38%, whereas WTE changes it by only 7%. Similarly, on the ACS dataset, SUBCURE-TUPLE increases the ATE by 31%, while WTE changes it by only 2%. These results show that *while the ATE may be largely stable under shifts in confounder distribution, it can be highly sensitive to the removal of specific tuples*, demonstrating a complementary notion of robustness captured by SUBCURE-TUPLE.

## 8 CONCLUSION & FUTURE WORK

We presented SUBCURE, a framework for auditing the robustness of causal conclusions via cardinality repair. By identifying minimal subpopulations whose removal brings the ATE within a target range, SUBCURE offers a quantitative sensitivity measure and interpretable insights into influential data regions.

SUBCURE assumes that a sufficient set of confounders is provided by the user. In practice, however, errors in confounder selection are a major source of bias in causal analysis [23, 55, 86]. Additionally, the current pattern repair algorithm is limited to conjunctions of equality predicates, which may be too coarse for datasets with continuous or ordinal attributes, potentially missing more expressive or meaningful subpopulations. Finally, our implementation currently supports a single-relation schema, and does not yet extend to multi-table relational databases. Previous work, such as [74], has extended causal models to handle multi-table data. Extending our approach to multi-table relational data is not straightforward, as it requires careful handling of how treatments, outcomes, and confounders are distributed across tables, how interventions propagate through joins, and how relational dependencies affect influence measures. Developing such an extension is non-trivial and represents an important direction for future work. Notably, prior work leveraging causal inference in data management [48, 73, 85] has also focused on single-table settings. We leave all these extensions to future work.

Several other promising directions remain for future research. A natural extension is to handle multiple causal effect estimations simultaneously. Additionally, expanding the framework to support a wider variety of intervention models beyond tuple deletion, such as tuple updates and insertions, would increase its applicability to real-world data scenarios where data can be modified in diverse ways. We also plan to explore directions that could deepen the theoretical understanding of our problem. One interesting direction is to adopt a probabilistic perspective, treating the database as a sample from an underlying distribution (similar to [38]). This

could enable reasoning about the effect of removing tuples without explicitly recomputing the ATE, for instance by using probabilistic bounds to estimate the impact of tuple removals. Such an approach might yield theoretical guarantees on the expected change in ATE. Another promising avenue is to explore approximate methods that provide (weaker) guarantees, such as using over- and under-approximations of the ATE obtained by sorting tuples based on outcome values and applying greedy selection strategies. Finally, studying the parameterized complexity of our problem, for example, with respect to output size, is another promising direction that could reveal fundamental limitations of our approach. We leave these investigations for future work.

## REFERENCES

[1] 2021. 2021 Stackoverflow Developer Survey. https://insights.stackoverflow.com/survey/2021.
[2] Foto N Afrati and Phokion G Kolaitis. 2009. Repair checking in inconsistent databases: algorithms and complexity. In *Proceedings of the 12th International Conference on Database Theory*. 31–41.
[3] Anish Agarwal and Rahul Singh. 2021. Causal inference with corrupted data: Measurement error, missing values, discretization, and differential privacy. *arXiv preprint arXiv:2107.02780* (2021).
[4] Shunit Agmon, Amir Gilad, Brit Youngmann, Shahar Zoarets, and Benny Kimelfeld. 2024. Finding Convincing Views to Endorse a Claim. *Proc. VLDB Endow.* 18, 2 (2024), 439–452. https://www.vldb.org/pvldb/vol18/p439-agmon.pdf
[5] Rakesh Agrawal, Ramakrishnan Srikant, et al. 1994. Fast algorithms for mining association rules. In *Proc. 20th int. conf. very large data bases, VLDB*, Vol. 1215. Santiago, Chile, 487–499.
[6] Kartik Ahuja, Prasanna Sattigeri, Karthikeyan Shanmugam, Dennis Wei, Karthikeyan Natesan Ramamurthy, and Murat Kocaoglu. 2021. Conditionally independent data generation. In *Uncertainty in Artificial Intelligence*. PMLR, 2050–2060.
[7] Anonymous Author(s). 2015. Git Repository. https://anonymous.4open.science/r/SubCuRe-662F/README.md.
[8] Arthur Asuncion and David Newman. 2007. UCI machine learning repository.
[9] Juhan Bae, Nathan Ng, Alston Lo, Marzyeh Ghassemi, and Roger B Grosse. 2022. If influence functions are the answer, then what is the question? *Advances in Neural Information Processing Systems* 35 (2022), 17953–17967.
[10] Heejung Bang and James M Robins. 2005. Doubly robust estimation in missing data and causal inference models. *Biometrics* 61, 4 (2005), 962–973.
[11] Matthew Blackwell. 2014. A selection bias approach to sensitivity analysis for causal effects. *Political Analysis* 22, 2 (2014), 169–182.
[12] Jasper Bogaert, Wen Wei Loh, Florian Schuberth, and Yves Rosseel. 2025. The Effect of Measurement Error on Hypothesis Testing in Small Sample Structural Equation Modeling: A Comparison of Various Estimation Approaches. *Structural Equation Modeling: A Multidisciplinary Journal* 32, 2 (2025), 215–236.
[13] Philip Bohannon, Wenfei Fan, Michael Flaster, and Rajeev Rastogi. 2005. A cost-based model and effective heuristic for repairing constraints by value modification. In *Proceedings of the 2005 ACM SIGMOD international conference on Management of data*. 143–154.
[14] Philip Bohannon, Wenfei Fan, Floris Geerts, Xibei Jia, and Anastasios Kementsietsidis. 2006. Conditional functional dependencies for data cleaning. In *2007 IEEE 23rd international conference on data engineering*. IEEE, 746–755.
[15] Tamara Broderick, Ryan Giordano, and Rachael Meager. 2020. An automatic finite-sample robustness metric: when can dropping a little data make a big difference? *arXiv preprint arXiv:2011.14999* (2020).
[16] Nofar Carmeli, Martin Grohe, Benny Kimelfeld, Ester Livshits, and Muhammad Tibi. 2021. Database Repairing with Soft Functional Dependencies. In *24th International Conference on Database Theory, ICDT 2021, March 23-26, 2021, Nicosia, Cyprus (LIPIcs, Vol. 186)*, Ke Yi and Zhewei Wei (Eds.). Schloss Dagstuhl - Leibniz-Zentrum für Informatik, 16:1–16:17.
[17] Zhangyu Cheng, Chengming Zou, and Jianwei Dong. 2019. Outlier detection using isolation forest and local outlier factor. In *Proceedings of the conference on research in adaptive and convergent systems*. 161–168.
[18] Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, Whitney Newey, and James Robins. 2016. Double/debiased machine learning for treatment and causal parameters. *arXiv preprint arXiv:1608.00060* (2016).
[19] Fei Chiang and Renée J. Miller. 2011. A unified model for data and constraint repair. In *2011 IEEE 27th International Conference on Data Engineering*. 446–457. doi:10.1109/ICDE.2011.5767833
[20] Jan Chomicki and Jerzy Marcinkowski. 2005. Minimal-change integrity maintenance using tuple deletions. *Information and Computation* 197, 1-2 (2005), 90–121.
[21] Xu Chu, Ihab F Ilyas, and Paolo Papotti. 2013. Holistic data cleaning: Putting violations into context. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*. IEEE, 458–469.
[22] Carlos Cinelli and Chad Hazlett. 2020. Making sense of sensitivity: Extending omitted variable bias. *Journal of the Royal Statistical Society Series B: Statistical Methodology* 82, 1 (2020), 39–67.
[23] Iván Díaz and Mark J van der Laan. 2013. Sensitivity analysis for causal inference under unmeasured confounding and measurement error problems. *The international journal of biostatistics* 9, 2 (2013), 149–160.
[24] Peng Ding and Fan Li. 2018. Causal inference. *Statist. Sci.* 33, 2 (2018), 214–237.
[25] Kareem El Gebaly, Parag Agrawal, Lukasz Golab, Flip Korn, and Divesh Srivastava. 2014. Interpretable and informative explanations of outcomes. *Proceedings of the VLDB Endowment* 8, 1 (2014), 61–72.
[26] Kenneth A Frank, Spiro J Maroulis, Minh Q Duong, and Benjamin M Kelcey. 2013. What would it take to change an inference? Using Rubin's causal model to interpret the robustness of causal inferences. *Educational Evaluation and Policy Analysis* 35, 4 (2013), 437–460.
[27] Sainyam Galhotra, Yue Gong, and Raul Castro Fernandez. 2023. Metam: Goal-oriented data discovery. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 2780–2793.
[28] Sainyam Galhotra, Romila Pradhan, and Babak Salimi. 2021. Explaining black-box algorithms using probabilistic contrastive counterfactuals. In *Proceedings of the 2021 International Conference on Management of Data*. 577–590.
[29] Floris Geerts, Giansalvatore Mecca, Paolo Papotti, and Donatello Santoro. 2013. The LLUNATIC data-cleaning framework. *Proceedings of the VLDB Endowment* 6, 9 (2013), 625–636.
[30] Antonio Ginart, Melody Guan, Gregory Valiant, and James Y Zou. 2019. Making ai forget you: Data deletion in machine learning. *Advances in neural information processing systems* 32 (2019).
[31] Clark Glymour, Kun Zhang, and Peter Spirtes. 2019. Review of causal discovery methods based on graphical models. *Frontiers in genetics* 10 (2019), 524.
[32] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 2020. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 9304–9312.
[33] Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. 2020. Certified data removal from machine learning models. In *Proceedings of the 37th International Conference on Machine Learning*. 3832–3842.
[34] William W Hager. 1989. Updating the inverse of a matrix. *SIAM review* 31, 2 (1989), 221–239.
[35] Ihab F Ilyas and Xu Chu. 2019. *Data cleaning*. Morgan & Claypool.
[36] Guido W. Imbens and Donald B. Rubin. 2015. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press.
[37] Zachary Izzo, Mary Anne Smart, Kamalika Chaudhuri, and James Zou. 2021. Approximate data deletion from machine learning models. In *International Conference on Artificial Intelligence and Statistics*. PMLR, 2008–2016.
[38] Sookyo Jeong and Hongseok Namkoong. 2020. Robust causal inference under covariate shift via worst-case subpopulation treatment effects. In *Conference on Learning Theory*. PMLR, 2079–2084.
[39] Daniel Kadlec, Kristin L Sainani, and Sophia Nimphius. 2023. With great power comes great responsibility: common errors in meta-analyses and meta-regressions in strength & conditioning research. *Sports Medicine* 53, 2 (2023), 313–325.
[40] Jon Kleinberg and Eva Tardos. 2006. *Algorithm design*. Pearson Education India.
[41] Solmaz Kolahi and Laks VS Lakshmanan. 2009. On approximating optimum repairs for functional dependency violations. In *Proceedings of the 12th International Conference on Database Theory*. 53–62.
[42] Yin Lin, Brit Youngmann, Yuval Moskovitch, HV Jagadish, and Tova Milo. 2021. On detecting cherry-picked generalizations. *Proceedings of the VLDB Endowment* 15, 1 (2021), 59–71.

[43] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2012. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 6, 1 (2012), 1–39.

[44] Ester Livshits and Benny Kimelfeld. 2022. The Shapley value of inconsistency measures for functional dependencies. *Logical Methods in Computer Science* 18 (2022).

[45] Ester Livshits, Benny Kimelfeld, and Sudeepa Roy. 2020. Computing optimal repairs for functional dependencies. *ACM Transactions on Database Systems (TODS)* 45, 1 (2020), 1–46.

[46] Merilyn Lock and Walid El Ansari. 2025. New world of big data—new challenges for evidence synthesis: impact of data duplication on estimates generated by meta-analyses and the development of a framework for its identification and management. *Journal of Clinical Epidemiology* 179 (2025), 111641.

[47] Christos Louizos, Uri Shalit, Joris M Mooij, David Sontag, Richard Zemel, and Max Welling. 2017. Causal effect inference with deep latent-variable models. *Advances in neural information processing systems* 30 (2017).

[48] Pingchuan Ma, Rui Ding, Shuai Wang, Shi Han, and Dongmei Zhang. 2023. XInsight: EXplainable Data Analysis Through The Lens of Causality. *Proc. ACM Manag. Data*, Article 156 (jun 2023), 27 pages.

[49] James MacQueen. 1967. Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, Vol. 5. University of California press, 281–298.

[50] Sara Magliacane, Thijs Van Ommen, Tom Claassen, Stephan Bongers, Philip Versteeg, and Joris M Mooij. 2018. Domain adaptation by using causal inference to predict invariant conditional distributions. *Advances in neural information processing systems* 31 (2018).

[51] Ananth Mahadevan and Michael Mathioudakis. 2021. Certifiable machine unlearning for linear models. *arXiv preprint arXiv:2106.15093* (2021).

[52] Neha Makhija and Wolfgang Gatterbauer. 2025. Is Integer Linear Programming All You Need for Deletion Propagation? A Unified and Practical Approach for Generalized Deletion Propagation. *Proceedings of the VLDB Endowment* 18, 8 (2025), 2667–2680.

[53] Markos Markakis, Brit Youngmann, Trinity Gao, Ziyu Zhang, Rana Shahout, Peter Baile Chen, Chunwei Liu, Ibrahim Sabek, and Michael Cafarella. 2024. From Logs to Causal Inference: Diagnosing Large Systems. *Proceedings of the VLDB Endowment* 18, 2 (2024), 158–172.

[54] Maya B Mathur and Tyler J VanderWeele. 2020. Sensitivity analysis for unmeasured confounding in meta-analyses. *J. Amer. Statist. Assoc.* (2020).

[55] Roseanne McNamee. 2003. Confounding and confounders. *Occupational and environmental medicine* 60, 3 (2003), 227–234.

[56] Dongjing Miao, Zhipeng Cai, Jianzhong Li, Xiangyu Gao, and Xianmin Liu. 2020. The computation of optimal subset repairs. *Proceedings of the VLDB Endowment* 13, 12 (2020), 2061–2074.

[57] Microsoft Research. 2019. EconML: A Python Package for ML-Based Heterogeneous Treatment Effects Estimation. https://github.com/microsoft/EconML.

[58] Caleb H Miles, Linda Valeri, and Brent Coull. 2024. Measurement error-robust causal inference via constructed instrumental variables. *arXiv preprint arXiv:2406.00940* (2024).

[59] Tim Miller. 2019. Explanation in artificial intelligence: Insights from the social sciences. *Artificial intelligence* 267 (2019), 1–38.

[60] Linda Nab, Rolf HH Groenwold, Maarten van Smeden, and Ruth H Keogh. 2020. Quantitative bias analysis for a misclassified confounder: a comparison between marginal structural models and conditional models for point treatments. *Epidemiology* 31, 6 (2020), 796–805.

[61] Thanh Tam Nguyen, Thanh Trung Huynh, Zhao Ren, Phi Le Nguyen, Alan Wee-Chung Liew, Hongzhi Yin, and Quoc Viet Hung Nguyen. 2022. A survey of machine unlearning. *arXiv preprint arXiv:2209.02299* (2022).

[62] Judea Pearl. 2009. Causal inference in statistics: An overview. (2009).

[63] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. 2011. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* 12 (2011), 2825–2830.

[64] Alireza Pirhadi, Mohammad Hossein Moslemi, Alexander Cloninger, Mostafa Milani, and Babak Salimi. 2024. OTClean: Data Cleaning for Conditional Independence Violations using Optimal Transport. arXiv:2403.02372 [cs.LG]

[65] El Kindi Rezig, Mourad Ouzzani, Walid G. Aref, Ahmed K. Elmagarmid, Ahmed R. Mahmood, and Michael Stonebraker. 2021. Horizon: scalable dependency-driven data cleaning. *Proc. VLDB Endow.* 14, 11 (jul 2021), 2546–2554. doi:10.14778/3476249.3476301

[66] James M Robins, Andrea Rotnitzky, and Daniel O Scharfstein. 2000. Sensitivity analysis for selection bias and unmeasured confounding in missing data and causal inference models. In *Statistical models in epidemiology, the environment, and clinical trials*. Springer, 1–94.

[67] Jeffrey Rosen. 2011. The right to be forgotten. *Stan. L. Rev. Online* 64 (2011), 88.

[68] Paul R Rosenbaum. 2005. Observational study. *Encyclopedia of statistics in behavioral science* 3 (2005), 1451–1462.

[69] Paul R Rosenbaum and Donald B Rubin. 1983. The central role of the propensity score in observational studies for causal effects. *Biometrika* 70, 1 (1983), 41–55.

[70] Sudeepa Roy, Laurel Orr, and Dan Suciu. 2015. Explaining query answers with explanation-ready databases. *Proceedings of the VLDB Endowment* 9, 4 (2015), 348–359.

[71] Sudeepa Roy and Dan Suciu. 2014. A formal approach to finding explanations for database queries. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. 1579–1590.

[72] Donald B Rubin. 2005. Causal inference using potential outcomes: Design, modeling, decisions. *J. Amer. Statist. Assoc.* 100, 469 (2005), 322–331.

[73] Babak Salimi, Johannes Gehrke, and Dan Suciu. 2018. Bias in olap queries: Detection, explanation, and removal. In *Proceedings of the 2018 International Conference on Management of Data*. 1021–1035.

[74] Babak Salimi, Harsh Parikh, Moe Kayali, Lise Getoor, Sudeepa Roy, and Dan Suciu. 2020. Causal Relational Learning. In *Proceedings of the 2020 International Conference on Management of Data, SIGMOD Conference 2020, online conference [Portland, OR, USA], June 14-19, 2020*, David Maier, Rachel Pottinger, AnHai Doan, Wang-Chiew Tan, Abdussalam Alawini, and Hung Q. Ngo (Eds.). ACM, 241–256. doi:10.1145/3318464.3389759

[75] Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. 2019. Interventional fairness: Causal database repair for algorithmic fairness. In *Proceedings of the 2019 International Conference on Management of Data*. 793–810.

[76] Francesco Sarracino and Malgorzata Mikucka. 2017. Bias and efficiency loss in regression estimates due to duplicated observations: a Monte Carlo simulation. In *Survey Research Methods*, Vol. 11. 17–44.

[77] Bernhard Schölkopf, Francesco Locatello, Stefan Bauer, Nan Rosemary Ke, Nal Kalchbrenner, Anirudh Goyal, and Yoshua Bengio. 2021. Toward causal representation learning. *Proc. IEEE* 109, 5 (2021), 612–634.

[78] Yi Shang. 2012. Measurement error adjustment using the SIMEX method: An application to student growth percentiles. *Journal of Educational Measurement* 49, 4 (2012), 446–465.

[79] U.S. Census Bureau. 2025. American Community Survey (ACS) - Data. https://www.census.gov/programs-surveys/acs/data.html Accessed: 2025-02-11.

[80] Eugene Wu and Samuel Madden. 2013. Scorpion: Explaining away outliers in aggregate queries. (2013).

[81] Yinjun Wu, Edgar Dobriban, and Susan Davidson. 2020. Deltagrad: Rapid retraining of machine learning models. In *International Conference on Machine Learning*. PMLR, 10355–10366.

[82] Yinjun Wu, Val Tannen, and Susan B Davidson. 2020. Priu: A provenance-based approach for incrementally updating regression models. In *Proceedings of the 2020 ACM SIGMOD international conference on management of data*. 447–462.

[83] Jing Nathan Yan, Oliver Schulte, MoHan Zhang, Jiannan Wang, and Reynold Cheng. 2020. Scoded: Statistical constraint oriented data error detection. In *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 845–860.

[84] Kôsaku Yosida. 2012. *Functional analysis*. Vol. 123. Springer Science & Business Media.

[85] Brit Youngmann, Michael Cafarella, Amir Gilad, and Sudeepa Roy. 2024. Summarized Causal Explanations For Aggregate Views. *Proceedings of the ACM on Management of Data* 2, 1 (2024), 1–27.

[86] Brit Youngmann, Michael Cafarella, Yuval Moskovitch, and Babak Salimi. 2023. On Explaining Confounding Bias. In *ICDE*. IEEE, 1846–1859.

[87] Alessio Zanga, Elif Ozkirimli, and Fabio Stella. 2022. A survey on causal discovery: Theory and practice. *International Journal of Approximate Reasoning* 151 (2022), 101–129.

[88] Jiongli Zhu, Sainyam Galhotra, Nazanin Sabri, and Babak Salimi. 2023. Consistent Range Approximation for Fair Predictive Modeling. *Proceedings of the VLDB Endowment* 16, 11 (2023), 2925–2938.

# APPENDIX

We provide missing proofs, details, and additional experiments.

# A PROOFS

In this part, we provide formal proofs of the propositions stated in Section 4.

### A.0.1 Proof of Proposition 4.1.

PROOF OF PROPOSITION 4.1. We show a reduction from SUBSET-SUM to the decision version of CaRET. Recall that in SUBSET-SUM we are given a set of numbers $S = \{x_1, \ldots, x_n\} \subseteq \mathbb{N}$ and a target $k$, and the problem is to decide whether there exists a subset $S' \subseteq S$ such that $\sum_{x \in S'} x = k$.

Recall that in CaRET we are given a database instance $D$ over a schema $\mathbb{A}$, a binary treatment variable $T \in \mathbb{A}$, an outcome variable $O \in \mathbb{A}$, a desired value $ATE_d$ and $\epsilon > 0$. In the decision version of this problem we are also given an additional bound $n$, and we need to decide whether there is a subset of tuples $\Gamma \subseteq D$ s.t. $|\Gamma| \leq n$ and it holds that $ATE_{D \setminus \Gamma}(T, O) \in [ATE_d - \epsilon, ATE_d + \epsilon]$.

We turn to describe the reduction. See Example A.1 for an example. Given a SUBSET-SUM instance $(S, k)$, we define an instance of CaRET as follows: The schema $\mathbb{A}$ is defined as $\mathbb{A} = \{T, O\}$ (i.e., there are no confounding variables). For each element $x \in S$, we add to $D$ two corresponding tuples $t_x^1$ and $t_x^2$ with the following values:
- $t_x^1[T] = 1, t_x^1[O] = x$
- $t_x^2[T] = 0, t_x^2[O] = 0$

In addition, we add to $D$ the tuple $t$ where $t[T] = 1$ and $t[O] = -k$. We set the target ATE to $ATE_d = 0$ and set $\epsilon = 0$, i.e., we must reach exactly the target ATE of 0. Finally, we bound on the allowed number of tuples to delete to be at most $|S|$. Clearly the reduction can be implemented in polynomial time. We turn to show correctness.

($\Rightarrow$) Assume the SUBSET-SUM instance has a solution $S' \subseteq S$. We remove from $D_d$ all tuples $t_x^1$ for $x \notin S'$ (intuitively, we retain only the $t_x^1$ tuples that are in $S'$). Note that since $|S'| \leq |S|$, we are within the allowed number of deleted tuples. We claim that the ATE after this removal is 0. Indeed, we have $AVG(O|T = 0) = 0$ (since all the tuples with $T = 0$ have $O = 0$), and since $\sum_{x \in S'} x = k$, then $AVG(O|T = 1) = (-k + \sum_{x \in S'} x)/(|S'| + 1) = (-k + k)/(|S'| + 1) = 0$ (accounting for the remaining tuples with $T = 1$). It follows that the resulting ATE is 0.

($\Leftarrow$) Assume that the instance of our problem is solvable, with a deleted set of tuples $\Gamma$. Intuitively, we do not care about deleted tuples where $t[T] = 0$, as these tuples also have $t[O] = 0$ and do no affect the ATE. Denote by $\Gamma'$ the set of deleted tuples with $t[T] = 1$. Since the new ATE is 0, and since the contribution of $t[T] = 0$ tuples is 0, it follows that after deleting the $\Gamma'$, the average of the tuples with $t[T] = 1$ is 0, and therefore also the sum.

Since $S$ contains only non-negative numbers, it follows that the tuple $t$ where $t[O] = -k$ is not deleted (otherwise the sum would be strictly positive). Thus, all tuples in $\Gamma'$ are of the form $t[T] = 1$ and $t[O] > 0$. Denote by $S' = \{x \in S \mid t_x \notin \Gamma'\}$ the remaining tuples, it follows that $-k + \sum_{x \in S'} x = 0$,

so $\sum_{x \in S'} x = k$, and therefore the SUBSET-SUM instance is solvable.

$\square$

**Table 4: The database $D$ created by the reduction in Example A.1. The tuple naming convention is as per the reduction in Theorem 4.1.**

|         | T | O  |
|---------|---|----|
| $t_1^1$ | 1 | 1  |
| $t_3^1$ | 1 | 3  |
| $t_5^1$ | 1 | 5  |
| $t$     | 1 | -4 |
| $t_1^2$ | 0 | 0  |
| $t_3^2$ | 0 | 0  |
| $t_5^2$ | 0 | 0  |

EXAMPLE A.1. *We illustrate the reduction in Theorem 4.1. Consider the following instance of SUBSET-SUM: $S = \{1, 3, 5\}$ and $k = 4$. We construct an instance of the cardinality repair for causal effect targeting problem as depicted in Table 4, and set $\epsilon = 0$. Observe that $ATE(T, O) = \frac{1}{4}(1 + 3 + 5 - 4) - \frac{1}{3}0 = \frac{5}{4}$.*

*A solution to the SUBSET-SUM instance is the set $S' = \{1, 3\}$ since their sum is equal to 4. Therefore, by deleting from $D$ the tuples $t_5^1$ (namely keeping $t_1^1$ and $t_3^1$), the new ATE is $\frac{1}{3}(1 + 3 - 4) - \frac{1}{3}0 = 0$, as required.*

### A.0.2 Proof of Theorem 4.2.
We now turn to prove Theorem 4.2, showing that Problem 2 is NP-hard. Again, the decision version of the problem includes a parameter $n$ that bounds the size of the removed population. Our reduction is from the setting of tuple deletion, namely CaRET.

Consider therefore an instance of the cardinality repair for causal effect targeting problem, i.e., a database $D$ over schema $\mathbb{A}$, variables $T, O$, a desired value $ATE_d$, $\epsilon > 0$ and a bound $n$. We can assume without loss of generality that there are no confounding variables, as the problem remains NP-hard in this case, as we show in the proof of Theorem 4.1 (this assumption only simplifies the writing, but has no meaningful impact on the construction).

Intuitively, the reduction introduces new attributes, in such a way that every subset of tuples can be defined by a pattern, thus showing an equivalence between tuple deletion and pattern deletion.

We start with the explicit construction.

Denote by $m$ the number of tuples in $D$. We obtain a new database $D_2$ by adding to $D$ fresh attributes $S_1, \ldots, S_m$ (where $S_i$ stands for "Select $i$"), and by modifying the tuples such that for tuple $t_i$ the values of the new attributes are $t_i(S_j) = 0$ if $i \neq j$ and $t_i(S_i) = 1$. We refer to $D_2$ as the *identifier augmentation* of $D$. We demonstrate the construction in Table 5.

LEMMA A.1. *Consider a database $D$ and its identifier augmentation $D_2$. For every set of tuples $I \subseteq \{1, \ldots, n\}$, there exists a pattern $\psi_I$ such that $SAT(\psi_I) = I$.*

**Table 5: A database $D$ and its identifier augmentation $D_2$.**

|       | T | O  |
|-------|---|----|
| $t_1$ | 1 | 10 |
| $t_2$ | 1 | 8  |
| $t_3$ | 1 | 6  |
| $t_4$ | 1 | 3  |

**Table 6:** $D$

|       | $S_1$ | $S_2$ | $S_3$ | $S_4$ | T | O  |
|-------|-------|-------|-------|-------|---|----|
| $t_1$ | 1     | 0     | 0     | 0     | 0 | 12 |
| $t_2$ | 0     | 1     | 0     | 0     | 0 | 9  |
| $t_3$ | 0     | 0     | 1     | 0     | 0 | 1  |
| $t_4$ | 0     | 0     | 0     | 1     | 0 | 1  |

**Table 7:** $D_2$

PROOF. Define the pattern $\psi_I = \bigwedge_{i \notin I} S_i = 0$. We claim that $\text{SAT}(\psi_I) = I$. Indeed, for $i \in I$ we have that $t_i(S_j) = 0$ for all $j \neq i$, and in particular for all $j \notin I$. Therefore $i \in \text{SAT}(\psi_I)$.

Conversely, let $i \in \text{SAT}(\psi_I)$, then $t_i(S_j) = 0$ for all $j \notin I$. Since $t_i(S_i) = 1$, it follows that $i \in I$. □

As an example of Theorem A.1, consider the set $\{t_2, t_3\}$ in Table 5. The pattern $S_1 = 0 \wedge S_4 = 0$ exactly captures it.

REMARK A.1. *A similar construction can be obtained with disjunctive patterns, by setting the pattern to capture all entries inside the set $I$.*

We can now show the correctness of the reduction.

($\Rightarrow$). If there exists a subset $\Gamma \subseteq D$ such that $ATE_{D \setminus \Gamma}(T, O) \in [ATE_d - \epsilon, ATE + \epsilon]$ and $|\Gamma| \leq n$, consider the pattern $\psi_\Gamma$ for $D_2$ as per Theorem A.1, then removing $\psi_\Gamma$ results in the same database $D \setminus \Gamma$, and therefore yields an ATE also within $\epsilon$ from $ATE_d$.

($\Leftarrow$). If there is a pattern in $D_2$ whose removal yields an ATE within $\epsilon$ from $ATE_d$, clearly we can also select the tuples in this pattern directly, without specifying them as a tuple, in the original database $D$.

REMARK A.2. *The reduction above does not use any property of ATE. Indeed, the reduction of tuple removal to pattern removal remains correct regardless of the optimization function.*

In light of Remark A.2, we have the following in particular.

THEOREM A.2. *The pattern-deletion problem is NP-complete for AVG and ATE.*

## B  ATE UPDATE FOR IPW

Algorithm 2 depicts the Fisher mini-batch incremental update algorithm for logistic regression that we use to update the model parameters.

## C  ADDITIONAL EXPERIMENTS

*C.0.1  Experiments with Synthetic Data under Known Deletion Budgets.* Running an exhaustive brute-force search algorithm to find the optimal solution was infeasible even on large datasets. To nonetheless evaluate how close SUBCURE-TUPLE comes to the optimal solution in large scale data, we designed experiments using synthetic data where an upper bound on the number of tuples to remove is known. Specifically, we generated a synthetic dataset with 10k tuples, including a binary treatment variable, a continuous outcome, and three confounders. We set the target ATE to be the one computed over this data with $\epsilon = 0$. We then introduced an increasing

---

**Algorithm 2:** Fisher Update for Logistic Regression

**Input** : Current model parameters $\theta \in \mathbb{R}^d$, original dataset $D$, subset to be removed $D_{\text{rmv}} \subseteq D$, mini-batch size $m'$, loss function $L(\theta; D)$

**Output:** Unlearned parameters $\theta_{\text{new}}$
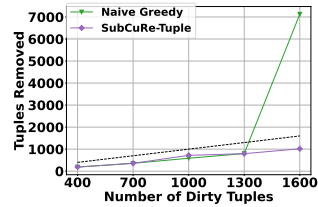
1  /* Generating Mini-Batches                              */
2  $s \leftarrow \lceil |D_{\text{rmv}}| / m' \rceil$
3  Split $D_{\text{rmv}}$ into $\{D_{\text{rmv}}^1, \ldots, D_{\text{rmv}}^s\}$
4  $D_{\text{new}} \leftarrow D \setminus D_{\text{rmv}}$
5  $\theta_{\text{new}} \leftarrow \theta$
6  **for** $i \leftarrow 1$ **to** $s$ **do**
7      /* Compute gradient and empirical Fisher (Hessian) */
8      $\Delta \leftarrow \nabla L(\theta_{\text{new}}, D_{\text{new}}^i), F \leftarrow \nabla^2 L(\theta_{\text{new}}, D_{\text{new}}^i)$
9      $\theta_{\text{new}} \leftarrow \theta_{\text{new}} - F^{-1} \Delta$
10     /* add calibrated noise for $\sigma$-certified unlearning */
11     **if** $\sigma > 0$ **then**
12         $b \leftarrow \text{SampleNoise}(d)$ /* $b \sim \mathcal{N}(0, I_d)$       */
13         $\theta_{\text{new}} \leftarrow \theta_{\text{new}} + \sigma F^{-1/4} b$
14  **return** $\theta_{\text{new}}$

number of noisy records into the data. In this setup, the number of inserted noisy tuples serves as an upper bound on the number of deletions required to reach the target ATE.

The results are shown in Figure 7. The approximate versions of SUBCURE-TUPLE and SINGLE-UPDATE-SUBCURE-TUPLE are omitted as we observed similar trends. Up to 1300 noisy tuples, all algorithms found solutions smaller than the upper bound, indicating high accuracy. However, as the number of noisy tuples increases, implying that many tuples must be removed to reach the target ATE, the performance of SINGLE-UPDATE-SUBCURE-TUPLE degrades. This inefficiency arises because SINGLE-UPDATE-SUBCURE-TUPLE computes influence scores only once, which eventually becomes outdated. A similar trend was observed in the ACS and Twins use cases.



**Figure 7: Number of removed tuples vs. number of noisy tuples. An upper bound on the solution size is marked by a dashed black line.**

*C.0.2  Factors Influencing the Performance SUBCURE-PATTERN.* We evaluate how the number of random walks affects both runtime and the algorithm's ability to identify impactful, small subpopulations. Too few walks risk missing valid solutions, while too many increase runtime. Experiments on real-world datasets show that capping random walks at 1000 provides a good trade-off between quality and efficiency. A similar balance is achieved by terminating a walk when the size of

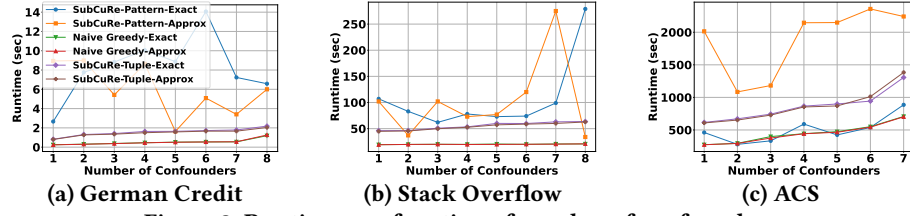(a) German Credit     (b) Stack Overflow     (c) ACS

Figure 8: Runtime as a function of number of confounders.

the current group exceeds 20% of the data. These parameters are user-configurable and can be adjusted to control runtime.

*Runtime vs. number of confounders.* The results are shown in Figure 8 illustrates the effect of the number of confounders on runtime. With more confounders, each ATE computation takes longer, and thus, as expected, it affects the runtime of all methods.