Fair and Actionable Causal Prescription Ruleset

Benton Li Cornell University USA cl2597@cornell.edu Nativ Levy
Technion
Israel
nativlevymail@gmail.com

Brit Youngmann
Technion
Israel
brity@technion.ac.il

Sainyam Galhotra Cornell University USA sg@cs.cornell.edu Sudeepa Roy Duke University USA sudeepa@cs.duke.edu

ABSTRACT

Prescriptions, or actionable recommendations, are commonly generated across various fields to influence key outcomes such as improving public health, enhancing economic policies, or increasing business efficiency. While traditional association-based methods may identify correlations, they often fail to reveal the underlying causal factors needed for informed decision-making. On the other hand, in decision making for tasks with significant societal or economic impact, it is crucial to provide recommendations that are justifiable and equitable in terms of the outcome for both the protected and non-protected groups. Motivated by these two goals, this paper introduces a fairness-aware framework leveraging causal reasoning for generating a set of actionable prescription rules (ruleset) toward betterment of an outcome while preventing exacerbating inequalities for protected groups. By considering group and individual fairness metrics from the literature, we ensure that both protected and non-protected groups benefit from these recommendations, providing a balanced and equitable approach to decision-making. We employ efficient optimizations to explore the vast and complex search space considering both fairness and coverage of the ruleset. Empirical evaluation and case study on real-world datasets demonstrates the utility of our framework for different use cases.

1 INTRODUCTION

Prescriptions, or actionable recommendations, are commonly generated across various fields to influence key outcomes such as improving product satisfaction, enhancing economic policies, or increasing business efficiency. Policymakers in government, decision-makers in businesses, and data scientists in various fields, often rely on data-driven approaches to identify potential actions to influence an outcome of interest, such as increasing income levels or loan approval rates. While association or prediction-based methods are extensively used in practice to draw useful insights from data, they typically identify correlations among variables and may fail to reveal the underlying causal factors, i.e., which actions may result in an improved outcome, needed for informed decision-making.

Causal analysis or causal inference, therefore, is considered one of the most important requirements to generate prescriptions that are actionable and aligned with human reasoning [38]. Causal inference, and in particular observational studies for causal inference on collected data (when controlled trials are impossible due to cost or ethical reasons), have been extensively studied in the statistics and artificial intelligence (AI) literature for several decades [66, 79].

Motivated by this foundational work on causal inference, the notion of causality has also influenced the field of database research. The causal models from AI have been extended to relational databases [82], and causality has been incorporated into various data management tasks such as finding responsibilities of inputs toward query answers [57, 58, 60], explanations for query answers [77, 109], data discovery [25, 108], data cleaning [68, 83], hypothetical reasoning [26], and large system diagnostics [6, 27, 33, 55].

If-then rules are generally considered interpretable by humans [15, 35, 46, 71, 97]. We give a concrete example of the difference between association and causation in generating prescriptions or recommended actions in the form of if-then rules below:

Example 1.1. Importance of causal prescriptions: Consider the Stack Overflow (SO) annual developer survey [2], where respondents from around the world answer questions about their jobs and demographics. A sample of the dataset with a subset of the attributes (there are 20 attributes) is presented in Table 1. Alice, a researcher in the United Nations (UN) finance department, is interested in discovering ways to increase the salaries of high-tech employees worldwide. She is looking for a set of actionable recommendations to raise the overall average salary. Using association-based approaches [15, 46], she may discover that individuals residing in the US who identify as straight or heterosexual tend to earn higher salaries (see Section 7.2 for full details). However, this observation merely indicates a correlation: people living in the US, for example, generally earn more than those outside the country. Their comparatively higher salaries are primarily attributable to the country's economy and are unrelated to their sexual orientation. Thus, this observation cannot be used as a prescription rule to increase salary. Our causal analysis, on the other hand, reveals that individuals aged 25-34 with dependents would benefit from working as front-end developers. This results in a \$44,009 annual salary increase on average. Another observation is that students should pursue an undergraduate major in CS. This can boost their salary by \$22,174 per year (see details in Section 6).

While generating prescriptions based on causal inference may help in robust decision-making, causal prescriptions that solely consider the betterment of an outcome (like salary) are not enough in practice. It is well-known that decision-making in many high-stake applications (like hiring policy, or policy for approving loans by banks) may lead to disparate societal or economic impact on different sub-populations. As a shocking example from a recent work called CauSumX [109] that generates a set of causal explanations for an aggregated view, the explanations generated suggest that

1

male individuals do a Bachelor's degree to increase their salary while being an unmarried woman has the most adverse effect on salary (borrowed directly from Fig. 19 in [106]). We explored this further in the context of generating prescriptions and observed that prescriptions that are not fairness-aware can generate unfair outcomes to some subpopulations which we refer to as the *protected group*. Examples include women, Black, Latino, or Native Americans, individuals with a disability, countries with a weaker economy, or other protected groups specific to an application.

Example 1.2. **Importance of fair prescriptions:** Continuing Example 1.1, while those causal prescription rules are highly beneficial for the overall population, they are considerably less effective for individuals residing in countries with a low GDP (indicating a weaker economy). For this group, the average expected increase in salary is only approximately \$13,000 per year (in contrast to \$44,009 for the entire group). Consequently, implementing these rules would exacerbate the disparity between those living in countries with strong economies and those in countries with weaker economies.

The example above shows that focusing solely on maximizing utility (i.e., increasing income) can result in a scenario where only some of the population receive significant improvement, while others experience no benefit (only a small benefit for individuals from countries with weaker economies in our example). Additionally, even if a large portion of the population receives recommendations, a protected subpopulation might not share the benefits and, worse, their situation could deteriorate, exacerbating inequalities.

Examples 1.1 and 1.2 show that it is crucial to provide recommendations that are (1) *causal* for the outcome (beyond associations), and (2) also *fair or equitable* in terms of the outcome for both the protected and non-protected groups. While recent work in database research has focused on deriving *causal explanations* for individual data points, aggregated view, or entire datasets [53, 81, 107, 109], and in particular [109] has considered generating a set of causal explanations for an aggregated view that resemble a ruleset, the absence of fairness considerations in generating these causal explanations can lead to unfair outcomes for the protected group.

Our contributions. Motivated by the dual goals of generating causal and fair prescriptions for the betterment of an outcome, we introduce a fairness-aware framework leveraging causal reasoning for generating a set of actionable prescription rules (ruleset) called FAIRCAP (Fair CAusal Prescription). Following research on fairness in data management [26, 91], we assume the existence of a protected subpopulation, defined by an attribute such as gender or race for people, or GDP of a country. Motivated by the causal explanation rules for an aggregated view [109], each prescription rule in our ruleset applies to a sub-population defined by a grouping attribute, and prescribes a treatment or intervention to improve the outcome for this sub-population. Fairness constraints ensure that the expected utility of the protected population is comparable to the utility of the unprotected individuals. We borrow the notions of group and individual fairness from the fairness literature but tailor them for prescription rules. In addition to the fairness constraints, our coverage constraints ensure that a substantial fraction of the population and protected subpopulation receives at least one recommendation.

Example 1.3. Continuing Examples 1.1 and 1.2, Alice uses our proposed system, called FAIRCAP, to impose fairness and coverage constraints to discover useful and equitable recommendations for increasing salaries worldwide. In particular, Alice chooses to implement a coverage constraint to ensure that the selected rules apply to a significant portion of people worldwide, including a sufficiently large number of individuals from countries with low GDP (the protected group). She also imposes a fairness constraint to ensure that the expected gains for both protected and non-protected groups are comparable. She discovers, for example, that for individuals with 6-8 years of coding experience (a subpopulation comprising 21% of the entire dataset and 25% of the protected group), pursuing a bachelor's degree in computer science will increase the expected salary by \$14.9k for protected and by \$17.8k for non-protected. (See Section 6 for more details.) This prescription rule applies to a large portion of the population and ensures fairness by providing a similar expected gain for both protected and non-protected groups, and the allowed difference of outcomes between these two populations may be adjusted by choosing appropriate thresholds in the fairness definitions.

Our main contributions are as follows.

(1) We develop a framework that generates a set of prescription rules to enhance an outcome of interest (Section 4). A prescription rule consists of a grouping pattern and an intervention pattern, representing the target subpopulation and the actionable recommendation for that group, respectively. The strength of the conditional causal effect (Section 3) of this intervention on the subgroup is used to measure the expected utility of a rule. Our objective is to identify the smallest set of rules that maximizes overall expected utility. We refer to this problem as the Prescription Ruleset Selection problem. We adopt several notions of fairness (individual vs. group, statistical parity vs. bounded group loss) from the literature to define the fairness constraints for our problem. In addition, coverage constraints (for individual rules or for a group) ensure that the solution for the Prescription Ruleset Selection problem is applied to a sufficient number of individuals and to minimize inequalities. We show NP-hardness for different variants of the problems and properties (matroid) useful in our algorithms.

- (2) We develop a general three-step algorithm named FAIR-CAP to solve the optimization problem of selecting a fair prescription ruleset (Section 5). The first step involves mining frequent grouping patterns using the Apriori algorithm [5]. In the second step, we employ a lattice-based algorithm to find high utility and fair intervention patterns for grouping patterns identified in the previous step. Finally, the third step applies a greedy approach to determine a solution. FAIRCAP can be easily adapted to accommodate all variants of the Prescription Ruleset Selection problem.
- (3) We provide a detailed case study (Section 6) and experimental analysis (Section 7) to evaluate our framework and algorithms. The case study shows the qualitative difference of different variants of our problem for different choices of the fairness and coverage constraints. The experiments include two datasets, three baselines, and 18 variations of our problem with different constraints. Our evaluations suggest that fairness may come at the cost of expected utility for everyone. However, without fairness constraints, we often observe a significant disparity between the

Table 1: A subset of the Stack Overflow dataset.

ID	Gender	Ethnicity	Age	Role	Education	Country	Undergrad Major	Salary
1	Male	White	26	Data Scientist	PhD	US	Computer Science	180k
2	Non-binary	White	32	QA developer	Bachelor's degree	US	Mechanical Eng.	83k
3	Male	South Asian	29	C-suite executive	Bachelor's degree	India	Computer Science	24k
4	Female	East Asian	21	Back-end developer	Bachelor's degree	China	Computer Science	19k

Table 2: Positioning of our framework w.r.t. previous work.

Related V	Causal	Fairness	Entire dataset	
	[53]	√	Х	Х
Aggregate Query	[58, 77, 81, 107, 109]	✓	×	✓
Result Explanation	[49, 61]	Х	×	X
_	[102]	X	Х	✓
Interpretable	[15, 46]	Х	Х	√
Prediction Models				
Multi dimensional	[64, 65, 93]	Х	✓	√
data aggregation	[22, 43]	X	Х	✓
FAIRC	√	√	√	

protected and non-protected. We also observe that achieving individual fairness is harder than group fairness, as most high-utility or high-coverage rules are unfair. Lastly, we show that FAIRCAP can generate prescription rules over large datasets in a reasonable time.

We discuss related work in Section 2, review background on causal inference in Section 3, and discuss the limitations of our framework and future work in Section 8.

2 RELATED WORK

Table 2 outlines the key distinctions between FAIRCAP and prior work. The columns in bold emphasize our key contributions: we generate **causality-based** prescription rules aimed at improving outcomes for the **entire datasets**, while also **considering fairness**. In contrast, other approaches either produce non-causal rules (as shown in the Causal column), target only subsets of the data (as indicated in the Entire dataset column), or disregard fairness considerations (as highlighted in the Fairness column).

Rule mining. Association rule mining has been extensively studied [45] and is used to identify relationships between items that frequently co-occurring in datasets. These techniques are applied across various fields, such as data analysis and outcome improvement. Notable algorithms include STEM [30], FP-Growth [41], AIS [116], and the Apriori algorithm [5]. We leverage the Apriori algorithm to identify sufficiently large subpopulations for which we will generate causal interventions. Rule-based interpretable prediction models [15, 46, 104] often leverage association rule mining to generate predictive rules [46, 47], with the goal of balancing high predictive accuracy with interpretability [44, 47, 52, 80, 86].

Recent work has proposed generating rules based on causal relationships. In [70], a framework was introduced to address biases in the data for fair causal analysis. Other studies have explored the integration of fairness criteria into causal analysis [69, 114]. [92] proposed a method to optimally allocate treatments with uncertain costs that vary based on confounders. Related research has also focused on estimating heterogeneous treatment effects [100, 101, 103]. However, these approaches assume both treatment and outcome variables are known. In contrast, we assume only the outcome variable is provided and aim to identify treatments that influence the outcome for different subpopulations, potentially yielding different treatments for each group. Our approach ensures the rules apply broadly while maintaining fairness for minority groups.

We adapt the method proposed in [109] called CauSumX. CauSumX is designed to identify the treatment with the highest causal effect on the outcome for a given subpopulation, generating causal explanations for aggregate queries. A main difference is that CauSumX does not consider fairness. We empirically show that using CauSumX to generate prescription rules can lead to significant disparities between protected and non-protected populations (See Section 7.2). Another main difference is that CauSumX considers the aggregate view to generate explanations, whereas we consider the entire data, thus, the search space is significantly larger. Our primary contribution lies in introducing fairness constraints on the generated rules, making the necessary adjustments to the algorithm to scale, and conducting an experimental study to demonstrate the importance of fairness constraints.

Fairness in data management. Algorithmic fairness, especially in the context of predictions by ML algorithms for high-stake decision making, has been a prominent topic in ML and AI (e.g., [3, 13, 21, 28, 39, 56, 67, 72, 73, 88]). Popular notions of fairness include group and individual fairness [12, 29, 89, 113]. Group fairness (measured as statistical parity or equalized odds) ensures that the decisionmaking process is fair to the protected group but may be unfair towards any specific individual. In contrast, individual notions of fairness enforce that the decisions are fair towards every individual. We refer the reader to Section 4.6 for more details. In recent years, fairness has emerged as a key consideration in data management research [26, 39, 91, 94, 113]. This includes ensuring fairness during data acquisition [8, 62, 63], improving data cleaning processes to promote fairness [34, 83], and achieving fairness in ranking and in database queries [50, 111, 112]. Techniques to ensure fairness of allocated resources [20, 54] can also be extended to study scenarios where the available resources may be diverse and constrained. One of our contributions is the introduction of novel definitions of group and individual fairness for causal analysis. We defer the extension of our work to other definitions of fairness to future work.

Causal inference in data management. Causality, used as a generic term of cause-effect analysis, has been used in different contexts in data management research [58, 59, 75, 82]. This includes data discovery [25, 32, 37, 84, 108], data cleaning [68, 83], query result explanation [53, 77, 81, 107, 109], hypothetical reasoning [24], and system diagnostics [6, 33, 55]. We use the interventional notion of causal inference on observational data from AI and Statistics [66, 79] to define our prescription rules (more in Sections 3 and 4.3). to design interventions that improve an outcome of interest.

Aggregate query result explanation. A substantial line of research has focused on using provenance to explain aggregate SQL query results [11, 14, 16, 48, 49, 58, 59, 96]. Other explanation methods include (non-causal) interventions [18, 76, 77, 95, 102], and counterbalancing patterns [61]. Recent work [53, 81, 107, 109] has proposed methods that use causal inference to explain query results.

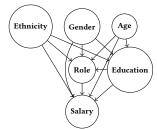


Figure 1: Partial causal DAG for the Stack Overflow dataset.

Multi dimensional data aggregation. Previous work on multidimensional data aggregation developed methods that extend the traditional drill-down and roll-up operators to find the most interesting data parts for exploration [4, 7, 40, 85, 105]. Other works have focused on assessing the similarity between data cubes [10]), or discovering intriguing data visualizations [99, 115].

Part of our goal is to identify subpopulations for which we can generate recommendations. We utilize existing solutions whenever applicable (e.g., we use the Apriori algorithm [5] to find sufficiently large subpopulations) and develop novel methods when necessary.

3 BACKGROUND ON CAUSAL INFERENCE

In this section, we review the basic concepts and key assumptions for inferring the effects of an intervention on the outcome on collected datasets without performing randomized controlled experiments. We use *Pearl's graphical causal model* for *observational causal analysis* [66] to define these concepts.

Causal Inference and Causal DAGs. The primary goal of causal inference is to model causal dependencies between attributes and evaluate how changing one variable (referred to as intervention) would affect the other. Pearl's Probabilistic Graphical Causal Model [66] can be written as a tuple (U, V, Pr_U, ψ) , where U is a set of *exogenous* variables, Pr_U is the joint distribution of U, and V is a set of observed *endogenous variables*. Here ψ is a set of structural equations that encode dependencies among variables. The equation for $A \in V$ takes the following form:

$$\psi_A : \operatorname{dom}(Pa_{\mathbf{U}}(A)) \times \operatorname{dom}(Pa_{\mathbf{V}}(A)) \to \operatorname{dom}(A)$$

Here $Pa_{\mathbf{U}}(A) \subseteq \mathbf{U}$ and $Pa_{\mathbf{V}}(A) \subseteq \mathbf{V} \setminus \{A\}$ respectively denote the exogenous and endogenous parents of A. A causal relational model is associated with a directed acyclic graph (causal DAG) G, whose nodes are the endogenous variables \mathbf{V} and there is a directed edge from X to O if $X \in Pa_{\mathbf{V}}(O)$. The causal DAG obfuscates exogenous variables as they are unobserved. The probability distribution $Pr_{\mathbf{U}}$ on exogenous variables \mathbf{U} induces a probability distribution on the endogenous variables \mathbf{V} by the structural equations ψ . A causal DAG can be constructed by a domain expert as in the above example, or using existing *causal discovery* algorithms [31].

EXAMPLE 3.1. Figure 1 depicts a partial causal DAG for the SO dataset over the attributes in Table 1 as endogenous variables (we use a larger causal DAG with all 20 attributes in our experiments). Given this causal DAG, we can observe that the role that a coder has in their company depends on their education, age gender and ethnicity.

Intervention. In Pearl's model, a treatment T = t (on one or more variables) is considered as an *intervention* to a causal DAG by mechanically changing the DAG such that the values of node(s) of

T in G are set to the value(s) in t, which is denoted by do(T = t). Following this operation, the probability distribution of the nodes in the graph changes as the treatment nodes no longer depend on the values of their parents. Pearl's model gives an approach to estimate the new probability distribution by identifying the confounding factors Z described earlier using conditions such as d-separation and backdoor criteria [66], which we do not discuss in this paper. Average Treatment Effect. The effects of an intervention are often measured by evaluating Conditional Average treatment effect (CATE), measuring the effect of an intervention on a subset of records [36, 78] by calculating the difference in average outcomes between the group that receives the treatment and the group that does not (called the *control* group), providing an estimate of how the intervention by T influences an outcome O for a given subpopulation. Given a subset of the records defined by (a vector of) attributes B and their values b, we can compute $CATE(T, O \mid B = b)$ as:

$$\mathbb{E}[O \mid do(T=1), B=b] - \mathbb{E}[O \mid do(T=0), B=b]$$
 (1)

Setting $B = \phi$ is equivalent to the ATE estimate. The above definitions assumes that the treatment assigned to one unit does not affect the outcome of another unit (called the Stable Unit Treatment Value Assumption (SUTVA)) [79]¹.

The ideal way of estimating the ATE and CATE is through randomized controlled experiments, where the population is randomly divided into two groups (treated and control, for binary treatments): denoted by do(T = 1) and do(T = 0) resp.) [66]. However, randomized experiments cannot always be performed due to ethical or feasibility issues. In these scenarios, observational data is used to estimate the treatment effect, which requires the following additional assumptions. The first assumption is called unconfoundedness or strong ignorability [74] says that the independence of outcome O and treatment T conditioning on a set of confounder variables (covariates) Z, i.e., $O \perp T|Z=z$. The second assumption called *overlap* or positivity says that there is a chance of observing individuals in both the treatment and control groups for every combination of covariate values, i.e., $0 < Pr(T=1 \mid Z=z) < 1$. The unconfoundedness assumption requires that the treatment T and the outcome O be independent when conditioned on a set of variables Z. In SO, assuming that only $Z = \{Gender, Age, Country\}$ affects T = Education, if we condition on a fixed set of values of Z, i.e., consider people of a given gender, from a given country, and at a given age, then T =Education and O = Salary are independent. For such confounding factors Z, Eq. (1) reduces to the following form (positivity gives the feasibility of the expectation difference):

$$CATE(T, O|B=b) = \mathbb{E}_Z \left[\mathbb{E}[O|T=1, B=b, Z=z] - \mathbb{E}[O|T=0, B=b, Z=z] \right]$$

This equation contains conditional probabilities and not do(T = b), which can be estimated from an observed data. Pearl's model gives a systematic way to find such a Z when a causal DAG is available.

4 PROBLEM FORMULATION

We consider a single-relation database over a schema \mathbb{A} . The schema is a vector of attribute names, i.e., $\mathbb{A}=(A_1,\ldots,A_s)$, where each A_i is associated with a domain $\mathsf{dom}(A_i)$, which can be categorical or continuous. A database instance D, populates the schema with a set of tuples $t=(a_1,\ldots,a_s)$ where $a_i\in\mathsf{dom}(A_i)$. We use $t[A_i]$ to denote the value of attribute A_i of tuple t.

¹This assumption does not hold for causal inference on multiple tables and even on a single table where tuples depend on each other.

Our high-level goal is to return a set of *prescription rules* (ruleset) with certain desired properties including fairness. In this section, first we define patterns on attributes, protected groups, prescription rules, and discuss the desired properties, finally defining our optimization problem in Section 4.7.

4.1 Pattern and Protected Group

To define the notion of *prescription rules*, we build upon the commonly used concept of *patterns* [22, 51, 77, 102]. Patterns are commonly used in query result explanation research [22, 51, 77, 102]. Generally, patterns are equivalent to the WHERE clause in SQL queries. However, in this work, we focus on a specific type of pattern: a conjunction of predicates, as has been done in prior query result explanation research [22, 51, 77, 102, 109].

DEFINITION 4.1 (PATTERN). Given a database instance D with schema \mathbb{A} , a predicate is an expression of the form $\varphi = A_i$ op a_i , where $A_i \in \mathbb{A}$, $a_i \in \text{dom}(A_i)$, and $\text{op} \in \{=, \neq, <, >, \leq, \geq\}$. A pattern is a conjunction of predicates $\mathcal{P} = \varphi_1 \wedge \ldots \wedge \varphi_k$.

EXAMPLE 4.1. An example pattern over the Stack Overflow dataset (Table 1) is \mathcal{P} ={Role = Designer \land Country = US}. It defines a subset of the dataset comprised of designers from the US.

Next we define *coverage* of a pattern \mathcal{P} defined by the number of tuples from D that it captures.

Definition 4.2 (Coverage of a pattern). Given a database instance D a pattern \mathcal{P} , and a tuple $t \in D$, \mathcal{P} is said to cover t if t satisfies the predicates in \mathcal{P} . The subset of tuples in D covered by \mathcal{P} is denoted by $Coverage(\mathcal{P})$.

As is common in fairness research [13, 91, 110], we assume the presence of a protected group, defined by the pattern \mathcal{P}_p . The remainder of the population (i.e., $D \setminus \mathcal{P}_p(D)$) is referred to as the non-protected group. A protected group in the Stack Overflow dataset may be defined based on sensitive attributes such as age or ethnicity. For instance, it could be defined as $\mathcal{P}_p = \{ \text{ Ethnicity} \neq \text{ White} \}$ to refer to non-white individuals.

4.2 Prescription Rules

A prescription rule outlines an intervention (treatment) designed to improve a target variable within a particular subpopulation. For example, a prescription rule might recommend that people under 25 pursue a Ph.D. to increase their salary. Before defining perception rules, we first discuss how attributes participate in such rules.

Mutable and immutable attributes. We assume the attributes \mathbb{A} are partitioned into two disjoint sets. The first set contains the interventional attributes (e.g., programming language, education), which are the attributes that can be changed to improve the outcome. The second set contains immutable attributes (e.g., age, gender), which cannot be changed through prescription. Formally, let $\mathbb{I} \subseteq \mathbb{A}$, denote the set of immutable attributes and $\mathbb{M} \subseteq \mathbb{A}$ denote the set of mutable (interventional) attributes, where $\mathbb{M} \cap \mathbb{I} = \emptyset$ and the outcome $O \notin \mathbb{M} \cup \mathbb{I}$. This categorization intends to prohibit the infeasible or impractical recommendations to increase the outcome (e.g., changing one's age or ethnicity to improve one's income).

Our prescription rules defined below as a combination of grouping and intervention patterns are motivated by the causal explanations defined in [109]. However, the focus of this paper is to study the interplay between utility and fairness of a ruleset that was not considered in [109]. As a result, the specific objectives and optimization problem are defined differently as discussed next.

Definition 4.3 (Prescription Rule and Ruleset, Grouping and intervention patterns, and Coverage). Given a database D with mutable attributes \mathbf{M} and immutable attribute \mathbf{I} , a prescription rule r is a pair of patterns $r = (\mathcal{P}_{grp}, \mathcal{P}_{int})$, where (1) \mathcal{P}_{grp} is called the grouping pattern and consists exclusively of the immutable attributes in \mathbf{I} , and (2) \mathcal{P}_{int} is called the intervention pattern and consists exclusively of the mutable attributes in \mathbf{M} .

By overloading notations, $COVERAGE(r) = COVERAGE(\mathcal{P}_{grp})$ is called the coverage of rule r since it captures the subset of tuples in D on which the rule r applies, i.e., each tuple $t \in D$ is either covered or not by \mathcal{P}_{grp} of rule r. \mathcal{P}_{int} defines the recommended intervention in the prescription rule aimed at betterment of the outcome O for the subgroup that COVERAGE(r) defines.

Given a set of prescription rules R (called a ruleset), $COVERAGE(R) = \bigcup_{r \in R} COVERAGE(r)$, i.e., coverage of a ruleset corresponds to the subset of tuples in D that are covered by at least one of the rules in R.

For a prescription rule $r = (\mathcal{P}_{\texttt{grp}}, \mathcal{P}_{\texttt{int}})$, the intervention pattern $\mathcal{P}_{\texttt{int}}$ partitions the tuples defined by $\mathcal{P}_{\texttt{grp}}$ into treated (T = 1 if $\mathcal{P}_{\texttt{int}}$ evaluates to true for a tuple) and control groups (T = 0 if $\mathcal{P}_{\texttt{int}}$ evaluates to false). This partition is then used to assess the causal effects of the intervention $\mathcal{P}_{\texttt{int}}$ on the outcome O within the subpopulation Coverage(r) which the rule r applies to.

Example 4.2. An example prescription rule suggests that individuals aged 25-34 with dependents, should work as front-end developers. (\mathcal{P}_{grp} : age = 25-34 \(\triangle \text{dependents} = yes \)), the intervention is working as front-end developers (\mathcal{P}_{int} : role = frontend developer). The expected CATE value is \$44,009, namely, the expected salary increase for a 25-34-year-old individual with dependents working as a frontend developer is \$44,009 per year (compared to a 25-34-year-old individual with dependents working in a different role).

4.3 Utility of a Prescription Ruleset

Utility of a single rule: To evaluate the effectiveness of a prescription rule $r = (\mathcal{P}_{\mathsf{grp}}, \mathcal{P}_{\mathsf{int}})$ toward improving the outcome O, we define its utility. The utility measures the expected impact (as CATE) of the recommended intervention on the outcome O within the subpopulation $\mathsf{COVERAGE}(r)$ where r applies to. We define the overall utility, and utility for the protected and non-protected groups.

Definition 4.4 (Utility of a prescription rule - overall, protected, non-protected). Given a database instance D with schema \mathbb{A} , an outcome variable O, a causal model \mathcal{G}_D on \mathbb{A} , a protected group p defined by a pattern \mathcal{P}_p , and a prescription rule $r = (\mathcal{P}_{grp}, \mathcal{P}_{int})$,

(1) the overall utility of r is defined as:

$$utility(r) := CATE_{G_D}(\mathcal{P}_{int}, O \mid \mathcal{P}_{grp})$$
 (2)

(2) the utility of r for the protected group p is defined as:

$$utility_p(r) := CATE_{G_D}(\mathcal{P}_{int}, O \mid \mathcal{P}_{grp} \land \mathcal{P}_p)$$
 (3)

(3) the utility of r for the non-protected group p is defined as:

$$utility_{\bar{p}}(r) := CATE_{G_D}(\mathcal{P}_{int}, O \mid \mathcal{P}_{grp} \land \mathcal{P}_{\neg p})$$
 (4)

The subscript \mathcal{G}_D denotes that the CATE is estimated using the causal model, and we drop the subscript when it is clear from context. If $Coverage(r) = Coverage(P_{grp}) = \emptyset$, i.e., if the rule does not apply to any tuple in D, then we assume that utility(r) = 0; similarly $utility_p(r) = 0$ if $Coverage(\mathcal{P}_{grp} \wedge \mathcal{P}_p) = \emptyset$, and $utility_{\bar{p}}(r) = 0$ if $Coverage(\mathcal{P}_{grp} \wedge \mathcal{P}_{\neg p}) = \emptyset$.

The goal of prescription rules is to improve the outcome O as desired. If the goal is to increase the outcome O (e.g., increase salary), we discard rules with negative utility, as they do not help achieve this objective. Similarly, if the aim is to decrease the outcome, we ignore rules with negative utility. Throughout the paper, without loss of generality, we assume that the goal is to increase the outcome, thereby focusing on maximizing utility.

Prescription to individuals when multiple rules apply: When dealing with a ruleset R, it is possible for multiple rules to apply to the same subpopulation. Specifically, if two rules $r_i = (\mathcal{P}_{grp}^i, \mathcal{P}_{int}^i)$ and $r_j = (\mathcal{P}_{\texttt{grp}}^j, \mathcal{P}_{\texttt{int}}^j) \in R$ share a non-empty intersection between their coverage, namely Coverage ($\mathcal{P}_{\mathsf{grp}}^i$) \cap Coverage ($\mathcal{P}_{\mathsf{grp}}^j$) \neq \emptyset , then the subpopulation defined by the pattern $\mathcal{P}_{grp}^i \wedge \mathcal{P}_{grp}^j$ will have more than one rule. In our definition below for utility of a ruleset, we refrain from applying more than one rule to a subpopulation for two reasons. First, two rules may conflict with each other. For instance, if one rule suggests individuals above 25 to earn a Ph.D., while another recommends women over 20 pursue an MBA, women above 25 would receive conflicting recommendations. Second, CATE is known to be non-monotonic [109], implying that appending a predicate to an intervention pattern can either increase or decrease the CATE value. Therefore, employing multiple rules simultaneously for a subpopulation might yield a utility gain smaller than the individual rules. Hence when multiple rules apply to a tuple in D, we assume that only one is chosen by the decision-maker.

Utility of a ruleset: For a prescription ruleset *R*, we use its *expected* utility on D as the utility of R.

DEFINITION 4.5 (EXPECTED UTILITY OF A RULESET). The expected utility of a prescription ruleset R is defined as the average maximum utility of an individual from Coverage(R) from the rules in R that applies to the individual, i.e.,

$$ExpUtility(R) = \frac{1}{n} \sum_{t \in CoverAGE(R)} \max_{r \in R_t} (utility(r))$$
 (5)

where $R_t \subseteq R$ denotes the set of rules covering the tuple t, and n = |D|. Note that if a rule does not apply to a tuple, its utility is zero, so the sum above is also over all $t \in D$.

Given a protected pattern \mathcal{P}_p , the expected utility for the protected and non-protected groups are defined as follows:

$$ExpUtility_p(R) = \frac{1}{n_p} \sum_{t \in COVERAGE_p(R)} \min_{r \in R_t} (utility(r))$$
 (6)

$$\begin{aligned} & ExpUtility_p(R) & = & \frac{1}{n_p} \sum_{t \in CoveRAGE_p(R)} \min_{r \in R_t} (utility(r)) & (6) \\ & ExpUtility_{\bar{p}}(R) & = & \frac{1}{n_{\bar{p}}} \sum_{t \in CoveRAGE_{\bar{p}}(R)} \max_{r \in R_t} (utility(r)) & (7) \end{aligned}$$

where $Coverage_p(R)$ denotes the set of protected individuals covered by R and $n_p = |COVERAGE_p(R)|$ (similarly $COVERAGE_{\bar{p}}()$ and $n_{\bar{p}}$).

Note the difference between formulas (6) and (5, 7). Since we do not assume any restriction on which rule is chosen for a tuple when multiple rules apply, we do a conservative worst-case analysis on fairness. We assume that protected individuals choose the worst possible rule, while the rest choose the best possible one. This ensures that the expected utility for the protected group in reality (irrespective of the rule chosen for each protected tuple) will be at least as high as the expected utility from the least beneficial relevant prescription rule for this group.

Size of a Prescription Ruleset

The size of a prescription ruleset is the number of rules in that set, denoted by size(R). Ideally, we want to find a small-size ruleset. The intuition is that, the fewer the rules in a set, the easier it is to understand the suggested interventions. Suppose we want to find a ruleset with high utility without specifying a constraint on the size. The following lemma shows that the best strategy is to return the optimal rule that applies to each individual. That is, to maximize utility, prescribing a personalized rule for each individual may lead to the best utility. Specifically, we can show that for every rule $r = (\mathcal{P}_{grp}, \mathcal{P}_{int})$, there exists a subgroup $g' \subseteq Coverage(\mathcal{P}_{grp})$ and a intervention $\mathcal{P}_{t'}$ s.t the utility of $r' = (g', \mathcal{P}_{t'})$ is greater than that of the original rule r.

Lemma 4.1. Given a rule $r=(\mathcal{P}_{grp},\mathcal{P}_{int})$, there exists a rule $r'=(\mathcal{P}_{g'},\mathcal{P}_{t'})$ $s.t \mathcal{P}_{q'} \subset Coverage(\mathcal{P}_{grp}) \ and \ utility(r') \geq utility(r).$

This property implies that the number of prescription rules in the optimal solution is O(|D|), making it impractical to implement in real-world scenarios. For instance, consider a policy enacted by a government official to allocate healthcare resources based on patient data. If the number of rules scales linearly with the size of the dataset, it would become infeasible to apply the policy effectively across a large population. Therefore, we limit the number of recommended rules. One approach is to impose a strict limit on the number of rules selected. However, pre-setting this constraint often requires tuning to balance utility and comprehensibility. Therefore, we incorporate the number of rules as an objective, considering rulesets with fewer rules to be more desirable, as was done in [46].

Coverage Constraints

We consider two types of coverage constraints: group coverage, where the goal is to find a solution that covers a predefined fraction of protected individuals and a certain fraction of the entire population, and rule coverage, where every selected rule must cover a certain fraction of the population and protected individuals.

Group Coverage Given two thresholds θ , $\theta_p \in [0, 1]$, we say that a ruleset R satisfies the group coverage constraint if R covers at least a θ fraction of the population, and a θ_p fraction of the protected subpopulation. Formally, both conditions are satisfied: (i) $Coverage(R) \ge \theta \cdot |D|$, (ii) $Coverage_p(R) \ge \theta_p \cdot |\mathcal{P}_p(D)|$, where $Coverage_p(R)$ denotes the number of covered protected individuals by *R*.

Rule Coverage Given two thresholds θ , $\theta_p \in [0, 1]$, we say that a ruleset R satisfies the rule coverage constraint if every rule $r \in R$ covers at least a θ fraction of the population, and at least a θ_p fraction of the protected subpopulation. Formally, both of the following conditions hold: (i) For every $r \in \mathbb{R}$: $coverage(r) \ge \theta \cdot |D|$, (ii) For every $r \in \mathbb{R}$: $coverage_p(r) \ge \theta_p \cdot |\mathcal{P}_p(D)|$, where $coverage_p(r)$ denotes the number of covered protected individuals by r.

4.6 Fairness Constraints

We study two definitions of fairness: statistical parity (SP) [56], and bounded group loss (BGL) [3]. Those definitions are based on equivalent notions of fairness for regression tasks [3]. We next provide an extension for these definitions to causal estimates.

Group and individual fairness are two key concepts in algorithmic fairness [12, 29, 91]. Group fairness aims to ensure that different groups receive similar outcomes. Individual fairness focuses on treating similar individuals similarly, meaning that if two individuals are alike in relevant aspects, they should receive similar outcomes. Both approaches aim to reduce bias, with the choice of which approach to adopt depending on the specific context. Next, we present four types of fairness constraints: SP and BGL, each of which can be applied to ensure group or individual fairness.

4.6.1 Statistical parity. In SP, the goal is to ensure that the gain in the utility of a protected individual is similar to that of any individual from the non-protected group.

Group Fairness: Intuitively, if we randomly sample a protected individual, the expected gain should be almost the same as that of an individual from the non-protected group. Formally:

 $|\text{ExpUtility}_{p}(R) - \text{ExpUtility}_{\bar{p}}(R)| \le \epsilon$, where $\epsilon > 0$ is a threshold.

Individual Fairness: Individual fairness says that the expected gain of every protected individual is similar to that of an individual from the non-protected group. That means that the expected utility of each rule $r \in R$ on a protected individual should be similar to that of an individual from the non-protected group. Formally, for every $r \in R$, $|utility_{\bar{R}}(r)-utility_{\bar{R}}| \le \epsilon$, where $\epsilon > 0$ is a threshold.

4.6.2 Bounded group loss (BGL). : Fair regression with BGL minimizes the overall loss while controlling the worst loss in the protected group [3]. In our context, this translates to the following constraint: When selecting an individual from the protected group, the utility increase should exceed a specified threshold $\tau \ge 0$.

Group Fairness: We aim to ensure that the expected utility of a randomly sampled protected individual within Coverage(R) is above a given threshold τ . Formally, ExpUtility, $t(R) \ge \tau$.

Individual Fairness: We aim to ensure that the gain of every protected individual from Coverage(R) exceeds a threshold τ . Therefore, a ruleset R satisfies the individual loss constraint if the utility of every rule r on protected individuals is at least τ . Formally, for every rule $r \in R$, $utility_D(r) \ge \tau$.

4.7 The Prescription Ruleset Selection Problem

We are finally ready to present the problem we study in this paper. If we did not have *fairness or coverage constraints*, then our goal is to select a small-size perception ruleset with high expected utility. However, as demonstrated in the introduction, not considering fairness constraints may result in a ruleset that are only highly beneficial to a small, non-protected subset of the population. Therefore, we extend our problem definition to include coverage and fairness constraints. We can apply any of SP or BGL group or individual fairness constraints (Section 4.6), as well as rule or group coverage

constraints Section 4.5), along with no fairness or coverage constraints, resulting in 18 distinct problem variants. The choice of which constraints to apply is left to the user as it may be application-dependent, and is discussed below. To define the generic problem, we use $R \models \mathcal{F}$ and $R \models C$ to denote that a ruleset R satisfies a given fairness constraint \mathcal{F} and a given coverage constraint C respectively (if no constraints are given, these conditions are trivially satisfied). We assume that a set of candidate rules $\{r_i\}_{i=1}^{I}$ has been already mined and is available as an input to the problem.

DEFINITION 4.6 (PRESCRIPTION RULESET SELECTION UNDER FAIRNESS AND COVERAGE CONSTRAINTS). Given a database D, a causal model \mathcal{G}_D , an outcome attribute O, a fairness constraint \mathcal{F} , a coverage constraint C, and a collection of prescription rules $\{r_i\}_{i=1}^l$, a subset $R\subseteq \{r_i\}_{i=1}^l$ of prescription rules is called valid if (1) $R\models \mathcal{F}$ and (2) $R\models C$. The goal is to find a valid subset of rules $R^*\subseteq \{r_i\}_{i=1}^l$ s.t

$$R^* = argmax_{R \subseteq \{r_i\}_{i=1}^{l}} \left[\lambda_1 \cdot (l - size(R)) + \lambda_2 \cdot ExpUtility(R) \right]$$
 (8)

where λ_1, λ_2 are non-negative weights.

 λ_1 and λ_2 may be tuned by the user. The above optimization problem, as expected, is NP-hard even for simple variants, although some constraints are matroid constraints and therefore are amenable to greedy approaches (discussions and proofs in the Appendix). Therefore, we obtain efficient algorithms that work well in practice in Section 5 and experimentally demonstrate the effect of different constraints on the results in Section 6.

Since we have several options for fairness and coverage constraints, a natural question is which version to use. We observe that there is no one-size-fits-all solution and the best choice depends on the specific application. For instance, going for individual fairness gives a stronger fairness guarantee at the expense of possible lower utility to everyone. In addition, the complexity of different versions can vary. To assist in making this decision, we summarize the process through a decision tree that guides users in selecting the most suitable variant for their needs, presented in Figure 2. The decision to choose between SP or BGL fairness is left to the user. In Section 6, we present a case study that empirically compares the obtained rulesets under different problem variants.

5 THE FAIRCAP ALGORITHM

A brute-force approach, which considers all grouping and intervention patterns to form prescription rules, results in long runtimes (as we show in Section 7.3). Instead, we propose a more efficient algorithm, called FairCap (Fair CAusal Prescription), which avoids generating every possible prescription rule. FairCap can be adapted for any variant of the Prescription Ruleset Selection problem. For simplicity, we first describe FairCap for the case with SP group fairness and group coverage constraints. We then explain how it can be modified to accommodate other variants.

Our algorithmic framework is outlined in Algorithm 1. FairCap consists of three parts: (1) generating grouping patterns by using the Apriori algorithm [5] (line 2), (2) identifying promising intervention patterns for each grouping pattern by using a lattice traversal approach [8], and (3) finding a set of prescription rules using a greedy approach. We leverage existing solutions (e.g., [5, 8, 65, 109]) where

7

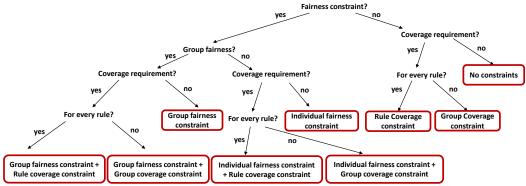


Figure 2: A decision tree for selecting the appropriate problem variant.

Algorithm 1: The FairCap algorithmic framework

```
\begin{array}{c} \textbf{input} : \text{A database relation } D, \text{ a protected group defined by the pattern } \mathcal{P}_p \\ & \text{and an outcome variable } O \\ \textbf{output} : \text{A set } \Phi \text{ of prescription rules.} \\ \textbf{1} \quad \Phi \leftarrow \emptyset; \\ \textbf{2} \quad \mathcal{G} \leftarrow \text{GetGroupingPatterns}(D, O); \\ \textbf{3} \quad \textbf{for } \mathcal{P}_g \in \mathcal{G} \text{ do} \\ \textbf{4} \quad \middle| \quad \mathcal{P}_t \leftarrow \text{GetIntervention}(\mathcal{P}_g, O, \mathcal{P}_p, D); \\ \textbf{5} \quad \middle| \quad \Phi \leftarrow \Phi \cup (\mathcal{P}_g, \mathcal{P}_p) \\ \textbf{6} \quad \Phi \leftarrow \text{ApplyGreedy}(\Phi, O, \mathcal{P}_p); \\ \textbf{7} \quad \textbf{return } \Phi \\ \end{array} \right. \\ \text{$ '' \  \  }
```

applicable, and develop novel techniques where necessary. Specifically, the first step follows the same approach as CauSumX [109], while the second and third steps introduce novel methods.

We show that the prescription rule selection problem is NP-hard even in simple settings (proofs deferred to the Appendix), and therefore developing effective heuristics considering several constraints is non-trivial. FairCap avoids generating all possible rules (as their number grows exponentially with the database size) and therefore does not perform an exhaustive search and may not return an optimal answer. If steps 1 and 2 were replaced by a brute-force approach that generates all rules, then a greedy approach for selecting a ruleset could approximate the optimal solution for certain problem variants, as the objective is a non-negative, monotone submodular function (even with a rule coverage or individual fairness constraints which are matroid constraints). However, other constraints are harder to satisfy. Future work will explore the complexity of various problem variants and establish theoretical bounds for finding approximate solutions.

Despite the fact FAIRCAP does not provide formal guarantees for the prescription ruleset selection problem, we emphasize that each selected rule represents an intervention that is statistically significant. Specifically, based on causal analysis, the expected utility reflects the anticipated average increase in the outcome for the specific subpopulation to which the rule applies.

5.1 Step 1: Mining Grouping Patterns

Considering every possible grouping pattern is infeasible as their number is exponential $(O(agrmax_{A_i \in \mathbb{A}} | dom(A_i)|^{|\mathbb{A}|})$. Instead, as done in previous work [65, 109], we utilize the Apriori algorithm [5] to generate candidate patterns. The Apriori algorithm gets a threshold τ , and ensures that the mined grouping patterns are present in at least τ tuples of D. The algorithm guarantees that each mined

pattern covers at least τ tuples from D, making them promising candidates for covering many tuples from D.

5.2 Step 2: Mining Intervention Patterns

Our next goal is to identify an intervention pattern \mathcal{P}_{int} for each mined grouping pattern \mathcal{P}_{grp} that maximizes utility (i.e., treatments with the highest CATE for \mathcal{P}_{grp}) while ensuring fairness to the protected group. Unlike step 1, this step requires novel techniques for finding treatments that are both fair and have high utility.

Since the number of potential intervention patterns for $\mathcal{P}_{\mathsf{grp}}$ can be large (exponential in $|\mathbb{A}|$), we employ a greedy lattice-traversal [8, 17] approach, inspired by [65, 109]. This allows us to materialize and assess the CATE only for promising patterns.

Concretely, the space of all intervention patterns can be represented as a lattice where nodes correspond to patterns and there is an edge between $\mathcal{P}_{\text{int}}^1$ and $\mathcal{P}_{\text{int}}^2$ if $\mathcal{P}_{\text{int}}^2$ can be obtained from $\mathcal{P}_{\text{int}}^1$ by adding a single predicate. This lattice can be traversed in a top-down fashion. Since not all nodes correspond to treatments with a positive CATE, we only materialize nodes if all their parents have a positive CATE. We note that this might lead the algorithm to overlook certain relevant intervention patterns. However, as shown in [109], combining patterns that exhibit a positive CATE is highly likely to result in an intervention with a positive CATE as well.

When a group fairness constraint is imposed, instead of searching for the treatment with the highest CATE, we search for the treatment that is "fair" by that it maximizes CATE for both protected and non-protected groups, while minimizing disparities.

To identify the most fair treatment, we define the *benefit* of an intervention pattern as follows. Intuitively, we penalize the treatment based on the difference between the utility for the non-protected group and the utility provided to the protected group. The larger the difference, the lower the benefit of the treatment. Formally, the benefit of a rule $r = (\mathcal{P}_{\text{grp}}, \mathcal{P}_{\text{int}})$ is defined as:

Formally, the benefit of a rule
$$r = (\mathcal{P}_{\mathsf{grp}}, \mathcal{P}_{\mathsf{int}})$$
 is defined as:
$$benefit(r) = \begin{cases} \frac{utility(r)}{1 + utility_{\bar{p}}(r) - utility_{\bar{p}}(r)}, & \text{if } utility_{\bar{p}}(r) \geq utility_{\bar{p}}(r) \\ utility(r), & \text{otherwise} \end{cases}$$

We implement two optimizations to improve efficiency: (i) we discard attributes that do not have a causal relationship with the outcome, since such attributes have no impact on CATE values. We can detect such attributes by utilizing the input causal DAG. (ii) The process of extracting intervention patterns for each grouping pattern can be performed in parallel since this procedure is dependent only on the grouping pattern.

5.3 Step 3: A Greedy Approach

The final step involves finding a solution from the rules mined in Steps 1 and 2. We propose a greedy algorithm that optimizes the problem's objectives. Intuitively, the algorithm operates as follows: at each iteration, it selects the next best rule that maximizes expected utility, benefit (as defined in Section 5.2), and coverage. Once the coverage constraints are met, the focus shifts to maximizing benefit and utility. The algorithm stops when the additional gain becomes negligible, as the number of rules is not predetermined.

Formally, the next best rule is determined as follows. Let $\{r_j\}_{j=1}^l$ denote the candidate rules and R_i is the ruleset selected in the first i iterations. The score of a rule r w.r.t R_i is defined as:

$$score(r) = Coverage(R_i \cup \{r\}) + benefit(R_i \cup \{r\}) + ExpUtility(R_i \cup \{r\})$$

The next best rule r_{i+1} to add in case the coverage constraints are not met yet is defined as: $r_{i+1}^* = argmax_{r_{i+1} \in \{r_j\}_{j=1}^l \setminus R_i} score(r_{i+1})$ In case the coverage constraints are met, ignore the coverage term. The algorithm stops at the first iteration i where the score of the selected rule r_i falls below a predefined threshold, indicating that the marginal gain from r_i is negligible.

5.4 Adjustments to Other Variants

We explain how FAIRCAP can be adjusted to solve other problem variants. We set the Apriori's threshold to ensure that each mined grouping pattern covers a sufficient number of individuals when a rule coverage constraint is imposed (step 1). Without coverage constraints, the Apriori threshold can be set to any value.

Without fairness constraints, in Step 2, the goal is to identify the intervention with the highest CATE value (as was done in [109]). Whit an individual fairness constraint, each rule must satisfy this constraint, so we only select interventions that are guaranteed to meet the constraint while maximizing CATE (step 2).

In case a group BGL fairness constraint is imposed, we define the benefit of a rule $r = (\mathcal{P}_{\sf grp}, \mathcal{P}_{\sf int})$ as follows. Intuitively, we penalize the treatment based on the difference between the minimum required utility for the protected group and the utility provided to the protected group by this treatment. The larger the difference, the lower the benefit of the treatment. Formally:

$$benefit(r) = \begin{cases} \frac{utility(r)}{1+\tau-utility_p(r)}, & \text{if } \tau \geq utility_p(r) \\ utility(r), & \text{otherwise} \end{cases}$$

where τ is the threshold for the BGL fairness constraint. This benefit definition is applied in Step 2 of the algorithm to identify fair and effective treatments for the mined grouping patterns.

Runtime complexity analysis: The maximum number of rules in a database D with attributes $\mathbb A$ is bounded by $|D|^{|\mathbb A|}$ (considering both grouping and intervention patterns and the active domain of attributes), which is polynomial in terms of *data complexity*, assuming a fixed schema [98]. The final greedy step is also polynomial in the number of rules considered. Additional operations, such as calculating CATE values, are polynomial in D, leading to worst-case polynomial data complexity. As we demonstrate in Section 7.3, our algorithm is capable of efficiently handling large datasets.

Table 3: Examined datasets.

Dataset	Tuples	Atts	Mut Atts	Protected Group
SO	38K	20	10	People from countries with a low GDP (21.5% of the data)
German	1000	20	15	Single Females (9.2% of the data)

6 CASE STUDY

The objective of this case study is to evaluate the impact of various constraints on the solution. We analyze two datasets, (1) German Credit (German in short) and (2) Stack Overflow (SO in short), each with a corresponding protected group, and assess the rules chosen by FairCap under different constraints. We aim to understand how these constraints influence coverage, utility, and disparities (for fairness) between protected and non-protected groups. We present example chosen rules under different configurations. We chose the rules by randomly picking one from each category (one that favors the protected group, one that favors the non-protected, and another that is more balanced). The full lists of rules are available in [1].

Datasets & protected groups. We examine two commonly used datasets: (1) Stack Oveflow (SO) [2], as described in Example 1.1. Here, the goal is to increase salary. (2) German Credit [9], which contains details of bank account holders, including demographic and financial information. Here, the goal is to increase the credit score (binary). The corresponding causal DAG was constructed using [108]. The datasets' statistics are presented in Table 3. The protected groups were selected to represent subgroups where the desired outcome was relatively low and sufficiently large to ensure the discovery of statistically significant rules. The protected group in Stack Overflow is defined as individuals from countries with a low GDP, which constitutes 21.5% of the data (the GDP attribute is categorical in this dataset). In the German data, the protected group is defined as single females, which constitutes 9.2% of the data.

Default parameters. Unless otherwise specified, the threshold of the Apriori algorithm is set to 0.1. For the SO dataset, the coverage thresholds are set to 0.5. The threshold for the SP and BFL fairness constraint is set at \$10k. For the German dataset, the coverage thresholds are set at 30% and the fairness thresholds are set at 0.1. This configuration allows for the generation of multiple rules.

The results are shown in Table 4, illustrating the trade-off between utility, coverage, and fairness. Without constraints, the expected utility is substantially higher, but this comes at the expense of greater disparities between protected and non-protected groups (as indicated by the unfairness score — the difference between the expected utility of protected and non-protected). In the examined scenarios, coverage for both groups was achieved without constraints, but other protected group definitions may require them. Stack Overflow. Observe that while the expected utility for both protected and non-protected groups reaches its highest value in the no-constraints variant, the unfairness score is very high. This indicates that achieving SP fairness requires compromising on the expected utility for both protected and non-protected groups. Interestingly, rule coverage and individual fairness are difficult to achieve, as most rules fail to meet these criteria. This leads to lower expected utility for all groups. On the other hand, group coverage and fairness constraints are easier to satisfy, as they offer more

Table 4: Comparison of Solutions in Terms of Size, Coverage, Expected Utility and Unfairness. IDS and FRL were used to either (i) replace step 1 of FAIRCAP to find grouping patterns; (ii) replace step 2 of FAIRCAP to find intervention patterns.

<u> </u>							
Stack Overflow (SP fairness)	# rules	coverage	coverage pro	exp utility	exp utility non-pro	exp utility pro	unfairness
No constraints	20	99.91%	99.98%	32634.2	32626.98	18432.66	14194.32
Group coverage	20	99.84%	99.88%	32597.02	32595.1	18340.29	14254.81
Rule coverage	10	99.99%	99.99%	22301.77	22292.02	16604.92	5687.1
Group fairness	8	97.52%	97.81%	27870.77	27998.47	17998.66	9999.81
Individual fairness	20	99.99%	99.99%	28014.58	28256.35	14241.07	14015.28
Group coverage, Group fairness	11	97.95%	98.85%	27934.76	28144.58	18145.23	9999.35
Rule coverage, Group fairness	12	99.96%	99.89%	22284.1	22279.93	16594.77	5685.16
Group coverage, Individual fairness	20	99.74%	99.88%	28057.78	28284.25	15128.91	13155.34
Rule coverage, Individual fairness	13	99.99%	99.99%	18591.41	18606.68	12797.15	5809.53
IDS (IF clause as grouping pattern)	16	100%	100%	29770.43	29988.1	16440.82	13547.28
IDS (IF clause as intervention pattern)	16	100%	100%	27763.89	27714.9	16888.1	10826.8
FRL (IF clause as grouping pattern)	9	99.5%	98.85%	27777.43	27782.3	18891.22	8891.08
FRL (IF clause as intervention pattern)	9	100%	100%	28999.22	28997.8	16453.8	12544
German Credit (BGL fairness)	# rules	coverage	coverage pro	exp utility	exp utility non-pro	exp utility pro	unfairness
No constraints	17	100.0%	100.0%	0.39	0.39	0.27	0.12
Group coverage	18	100.0%	100.0%	0.39	0.39	0.3	0.09
Rule coverage	6	96.0%	100.0%	0.31	0.31	0.3	0.01
Group fairness	18	100.0%	100.0%	0.39	0.39	0.3	0.09
Individual fairness	20	100.0%	100.0%	0.37	0.37	0.23	0.14
Group coverage, Group fairness	6	100.0%	100.0%	0.36	0.37	0.31	0.06
Rule coverage, Group fairness	3	90.0%	100.0%	0.29	0.29	0.31	-0.02
Group coverage, Individual fairness	20	100.0%	100.0%	0.37	0.37	0.23	0.14
Rule coverage, Individual fairness	8	96.8%	100.0%	0.29	0.29	0.23	0.06
IDS (IF clause as grouping pattern)	12	100%	100%	0.35	0.35	0.3	0.05
IDS (IF clause as intervention pattern)	12	100%	100%	0.34	0.34	0.24	0.1
FRL (IF clause as grouping pattern)	13	100%	100%	0.26	0.26	0.21	0.05
FRL (IF clause as intervention pattern)	13	100%	100%	0.3	0.3	0.23	0.07

flexibility by allowing the selection of some unfair rules alongside those specifically designed for the protected group.

3 Selected Rules out of 11 for SO (SP group fairness):

- \triangleright (S1_a) For individuals aged 24-34, pursue an undergraduate major in CS (exp utility protected: 10,292, exp utility non-protected: 22,586).
- \triangleright (S1_b) For individuals with 6-8 years of coding experience, work with a computer 9 12 hours a day. (exp utility protected: 17,161, expe utility non-protected: 19,254).
- \triangleright (S1_c) For males whose parents have a secondary school education, work as back-end developers (exp utility protected: 51,542, exp utility non-protected: 46,354).

We show above the three example rules selected under group fairness constraint. The first rule $S1_a$ is more advantageous for the non-protected group, the second $(S1_b)$ benefits both protected and non-protected groups similarly, while the third rule $(S1_c)$ is more beneficial for the protected group. Overall, all these rules together satisfy the group fairness requirement. We also present three example rules selected under individual fairness constraints. In this case, all rules $(S2_a, S2_b, S2_c)$ are nearly equally beneficial for both groups, but the overall expected utility is lower. Finally, consider the three example rules selected with no constraints. Here, all rules $(S3_a, S3_b, S3_c$ in the figure below) favor the non-protected group, highlighting the importance of including fairness constraints.

3 Selected Rules out of 20 for SO (SP individual fairness):

- ▶ (S2_a) For males aged 25-34 with no dependents, pursue a bachelor's degree (exp utility protected: 16,158, exp utility non-protected: 18,134).
- ▶(S2_b) For individuals aged 18 -24, work as back-end developers. (exputility protected: 12,664, exputility non-protected: 14,101).
- \triangleright (S2_c) For individuals with dependents, pursue an undergraduate major in CS (exp utility protected: 16,124, exp utility non-protected: 17,138).

3 Selected Rules out of 20 for SO (no fairness constraints):

- \triangleright (S3_a) For White aged 25-34 with dependents, work with computer 9-12 hours a day and work as back-end developers (exp utility protected: 11,147, exp utility non-protected: 32,248).
- \triangleright (S3_b) For males aged 35-44 with dependents, work as back-end developers. (exp utility protected: 11,189, exp utility non-protected: 40,207).
- ► (S3c) For students, pursue an undergraduate major in CS (exp utility for protected: 12,126, exp utility for non-protected: 22,174).

German. While the expected utility for both protected and non-protected peaks in the no-constraints variant, the unfairness score is relatively high. This suggests that achieving BGL fairness necessitates compromising utility for both groups. Notably, to reduce the unfairness, it is feasible to impose either a rule coverage constraint or a rule coverage constraint combined with group fairness. We show three rules selected under BGL group fairness constraints below. Since we are focusing on BGL fairness, which considers only the minimal gain for the protected group without regard for the gains of the non-protected group, we still observe a disparity between the two, even with a fairness constraint in place.

3 Selected Rules out of 20 for German (group BGL fairness):

- ho (G1_a) For people aged 24-30 with 0-2 dependents, maintain a minimum balance of 200 DM in the checking account and pursue skilled employment (exp utility protected: 0.26, exp utility non-protected: 0.35).
- ho (G1_b) For people seeking a loan to purchase furniture or equipment, maintain a minimum balance of 200 DM in the checking account (exputility protected: 0.38, exp utility non-protected: 0.29).
- ho(G1_c) For people seeking a loan for an unspecified purpose, maintain a minimum balance of 200 DM in the checking account and own a house. (exp utility protected: 0.54, exp utility non-protected: 0.41).

7 EXPERIMENTAL EVALUATION

We present an experimental evaluation that evaluates FAIRCAP effectiveness and efficiency. We aim to address the following questions: Q1: How does the quality of our generated rulesets compare to that of existing methods? Q2: What is the efficiency of FAIRCAP and how is it affected by various data and system parameters?

7.1 Experimental Setup

FAIRCAP was implemented in Python, and is publicly available in [1]. CATE values computation was performed using the DoWhy library [87]. The generated rules were translated into natural language using simple, manually constructed templates. We perform experiments on CloudLab [19] xl170 machines (10-core 2.4 GHz CPU, 64 GB RAM). The datasets, protected groups, and default parameters considered are the same as those described in Section 6.

Baselines. We compare FAIRCAP with the following baselines: 1. CauSumX: CauSumX [109] is designed to find a summarized causal explanation for group-by-avg SQL query results. When applied directly to the datasets, it can be viewed as a solution to our problem with only an overall coverage constraint. 2.IDS [46] is a framework for generating Interpretable Decision Sets for prediction tasks. IDS incorporates parameters restricting the percentage of uncovered tuples and the number of rules. These parameters were assigned the same values in our system. 3. FRL: The authors of [15] proposed a framework for creating Falling Rule Lists (FRLs) as a probabilistic classification model. FRLs comprise if-then rules with antecedents in the if-clauses and probabilities of the desired outcome in the then-clauses, ordered based on associated probabilities.

Since IDS and FRL assume a binary outcome, we binned the salary variable in SO using the average value. To address fairness considerations, we run the baseline algorithms twice (excluding Brute-Force): Once on the entire dataset to obtain a set of rules applicable to the entire population, and again solely on the tuples belonging to the protected population to generate rules specifically tailored for them. We report the number of rules generated by the baselines, their coverage, and runtime. To compare the expected utility, we proceed as follows: The rules generated by IDS and FRL are prediction rules (e.g., IF owning a house = YES, THEN credit score = 1). As such, these rules do not provide an intervention to improve outcomes. We, therefore, treat the IF clauses in two manners: (1) IF clauses as the selected grouping patterns and then apply step 2 (Section 5.2) of FAIRCAP to determine the intervention patterns; (2) IF clauses as the selected intervention patterns, where the grouping pattern is the entire data.

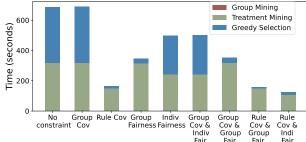


Figure 3: Runtime by-step of the FAIRCAP algorithm (SO)
7.2 Quality Evaluation (Q1)

We compare the set of rules chosen by each baseline and FAIRCAP. Stack Overflow. As discussed in Section 6, prescription rules selected without fairness constraints, similar to the behavior of CauSumX, were significantly more advantageous for non-protected. The rules generated by IDS do not suggest interventions to improve outcomes. For example, one rule states that if Country = Turkey and Age = 18-24 years, then the expected salary is low (with the outcome binned). Another key distinction is that these rules are not causal, as they are based on correlations in the data. For example, one rule indicates that if the years coding = 0-2 and Sexual Orientation = Gay or Lesbian, then the expected salary is low. Similarly, rules generated by FRL do not propose interventions to improve outcomes and are not causal. For example, one rule states that if Country = US and Sexual Orientation = Straight or Heterosexual, then the expected salary is high. In contrast, FAIRCAP generates interventions aimed at improving the outcome by leveraging causal relationships. It also allows users to impose fairness constraints, ensuring that the protected group benefits from these interventions. **German.** Here again, with no fairness constraint (akin to CauSumX), the selected rules were mostly beneficial for the non-protected. Here again, the rules generated by IDS are not causal and do not offer an intervention. For example, one of the rules suggested that single females at the age of 35-41 are unlikely to get a loan. As before, the rules generated by FRL are also not causal and do not propose ways to improve the credit risk score. For example, one rule suggests that if a person has lived in a house for 4-7 years, their credit risk score is likely to be high. Another rule states that if the purpose of the loan is to buy a used car, the credit risk score is also likely to be high. Clearly, these rules rely on correlations in the data rather than causal relationships. In contrast, FAIRCAP generated a ruleset that offers interventions to improve the credit risk score based on causal relationships. Example selected rules are shown in Section 6.

We report the solution size, coverage, expected utility for protected and non-protected, and the unfairness of the rulesets generated using IDS and FRL (as explained in Section 7.1). The results are presented in **Table 4**. Notably, the expected utility for both protected and non-protected groups across both datasets is generally lower than that achieved by FairCap. FairCap consistently delivers higher expected utility for both groups and a smaller difference between these values. This indicates that our approach to mining grouping and intervention patterns is more effective than relying on these algorithms for the same purpose. However, we note that the rules in IDS and FRL had different objectives (prediction accuracy) and had to be adapted for quantitative comparison using our measures.

exp utility pro Stack Overflow (SP fairness) # rules coverage coverage pro exp utility exp utility non-pro unfairness Group SP (2.5K) 97.82% 99.0% 20973.55 20772.77 18275.44 2497.33 4 Group SP (5K) 7 97.31% 98.24% 22805.52 23069.98 18071.12 4998.86 Group fairness (10K) 8 97.52% 97.81% 27870.77 27998.47 17998.66 9999.81 Group SP (20K) 20 99.88% 99.94% 32671.11 32664.45 18423.64 14240.81 Individual SP (2.5K) 20 99 95% 99 98% 24070.94 24433.55 12784.62 11648.93 Individual SP (5K) 20 99.99% 99.99% 25526.1 25911.22 15327.21 10584.01 Individual SP(10K) 20 99.99% 99.99% 28014.58 28256.35 14241.07 14015.28 Individual SP (20K) 20 99.51% 99.63% 29984.0 29966.29 14929.7 15036.59

Table 5: Comparison of Solutions in Terms of Fairness

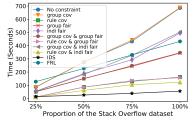


Figure 4: Runtime as a function of the dataset size (SO)

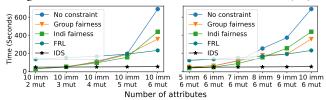


Figure 5: Runtime as a function of number of mutable and immutable attributes for SO with statistical parity

7.2.1 Robustness to the Causal DAG. The quality of the generated rules may depend on the accuracy of the underlying causal DAG. To evaluate this, we examine the impact of different causal DAGs on the rules. The causal DAGs considered are as follows: (1) 1-layer Indep DAG: A causal DAG where all attributes are independent of each other and only impact the outcome. This setting similar to the scenario where all the causal graph is ignored. (2) 2-layer Mutable DAG: A simplified DAG where immutable attributes affect the mutable attributes, which impact the outcome variable. In this graph, all immutable attributes act as confounders but do not directly impact the outcome. (3) 2-layer DAG: A simplified DAG where all variables affect the outcome but the mutable attributes are also confounded by all immutable attributes. (4) PC DAG: A causal DAG generated by the PC causal discovery algorithm [90].

The results are depicted in Table 6. We report the expected utility as computed over the different causal DAGs. We observe that the expected utility remains similar for the Stack overflow dataset, demonstrating robustness towards the choice of causal dag. The results show some variability in German credit. However, the PC DAG and the original causal DAG are the most accurate (as they are based on the data distribution and domain knowledge) and achieve the highest coverage and expected utility.

7.3 Scalability Evaluation (Q2)

Breakdown analysis by step. Figure 3 shows the runtime comparison of FairCap for different problem settings. Observe that using rule coverage constraint has the lowest runtime because it helps to

prune rules which do not satisfy the coverage constraint. Employing rule coverage with individual fairness is the fastest among all settings, while no constraint setting takes the longest time. The time taken by the group mining phase is less than 2 seconds across all setups, and is therefore not visible in the plot. The intervention mining phase (Step 2) is the most inefficient phase, which takes around 6 mins for the unconstrained setting. The running time of these components aligns with our time complexity analysis (Section 5). Due to space restrictions, we do not present the corresponding plot for German dataset. All conclusions remain the same but the overall running time is $\approx 10\times$ faster due to its smaller size.

The running time of FairCap and the baselines is comparable. FRL is an order of magnitude slower than IDS because it uses a Bayesian modeling approach to simultaneously select a subset of rules and determine their optimal order, which involves solving a computationally intensive combinatorial problem. In contrast, IDS leverages submodular optimization on an unordered set of rules, significantly reducing the size of the search. We now analyze the impact of system parameters and data size on performance.

Data Size. Figure 4 compares the running time of FairCap and the baselines for varying dataset sizes. We observe that the time taken by FairCap and the baselines increases linearly for most of the settings, with FairCap demonstrating a runtime comparable to IDS under certain configurations. We also observed that the quality of rules returned by sampling 25% of the data points is comparable with the rules returned by using the whole dataset. Therefore, sampling-based optimizations can help to reduce the running time from 11 min to less than 2 min for the unconstrained setting and less than a minute with fairness constraints.

Number of Attributes. Figure 5 shows the runtime of FAIRCAP while increasing the number of mutable and immutable attributes. On increasing the number of mutable attributes, the number of intervention patterns increases exponentially while on increasing immutable attributes, the number of grouping patterns increases exponentially. Therefore, both have a similar impact on runtime. IDS and FRL do not distinguish between mutable and immutable attributes and there the runtime increases slightly due to an increase in the number of attributes, as more rules are considered.

In the following, we omit the results for the IDS and FRL baselines, as these parameters do not impact their runtime.

Fairness Threshold. Table 5 presents the results for varying ϵ for group and individual fairness. We observe that the unfairness of the returned solution increases with the increase in ϵ . Additionally, the overall expected utility increases but the expected utility of the

Table 6: Metrics Comparison with different Causal DAGs.

Stack Overflow (SP group fairness + group coverage)	# rules	coverage	coverage pro	exp utility	exp utility non-pro	exp utility pro	unfairness
Original causal DAG	11	97.95%	98.85%	27934.76	28144.58	18145.23	9999.35
1-Layer Indep DAG	11	98.38%	98.38%	28110.19	28117	18117.45	9972
2-Layer Mutable DAG	10	97.7%	98.4%	28198.59	28193.09	18193.23	9999.86
2-Layer DAG	10	98.47%	98.87%	28106.4	28211.17	18211.4	9999.77
PC DAG	10	97.7%	98.4%	28198.59	28193.09	18193.23	9999.86
German Credit (BGL group fair-	# rules	coverage	coverage pro	exp utility	exp utility non-pro	exp utility pro	unfairness
ness + group coverage)							
ness + group coverage) Original causal DAG	6	100.0%	100.0%	0.36	0.37	0.31	0.06
	6 12	100.0% 100%	100.0% 100%	0.36 0.31	0.37 0.31	0.31 0.29	0.06 0.02
Original causal DAG	-						
Original causal DAG 1-Layer Indep DAG	12	100%	100%	0.31	0.31	0.29	0.02

protected individuals decreases. This result matches our intuition as highly unfair rules are selected for higher values of ϵ . We also notice that the greedy algorithm satisfies the group fairness constraint in all scenarios (unfairness is always less than the desired threshold).

For individual fairness, the overall utility increases monotonically with ϵ . However, the rate of growth for individual fairness is slower than that of group fairness. One interesting observation about individual fairness is that when all rules have statistical parity difference less than 2500, the overall unfairness is still around 11K. This sudden increase in unfairness when considering multiple fair rules together is because we evaluate the upper bound of unfairness by taking the difference between max utility of unprotected and min utility of protected individuals. On manual inspection, we observed that all rules are indeed individually fair.

Coverage Threshold. With the change in coverage thresholds, we do not observe major difference in the overall results because the majority of the rules exhibit very high coverage (Table 4).

Apriori Threshold. We observe that increasing the Apriori threshold τ leads to a reduction in the number of grouping patterns considered, and thus to a decrease in runtime. However, our findings indicate that higher τ values lead to a decrease in both utility and fairness. Based on our findings, we recommend using a default value of 0.1, which provides satisfactory results in terms of coverage, utility, fairness and runtime.

8 LIMITATIONS AND FUTURE WORK

FAIRCAP generates actionable, causal-based recommendations to improve a target outcome while incorporating coverage and fairness constraints. To the best of our knowledge, this is one of the first works in this direction, and several directions of future work remain. In this section, we discuss some of the current limitations of FAIRCAP and future directions.

Generation and usage of rules by FAIRCAP. FAIRCAP can be used to recommend actions for different subpopulations toward optimizing a target. As an example scenario, a policymaker may select the target outcome and the parameters for coverage and fairness constraints (which may be iteratively varied based on the

application). FairCap then generates a prescription ruleset as recommended actions for different subpopulations. The current framework assumes that the policymaker is trustworthy, will not misuse the rules, and will publish the relevant recommendations for each subpopulation. However, it is important to note that if not all rules are provided to all subpopulations, disparities among subpopulations may increase. In addition, the generated rules may not impact all individuals receiving the recommendation in the same way. The gain in objective may vary across different subpopulations. For example, an increase in \$10k revenue may have varied impacts in different countries, depending on the cost of living and purchasing power. Addressing these will be interesting future work.

Considering constraints, costs, and resources in rule generation. The current framework does not account for the cost of interventions. Some interventions may be impractical (e.g., pursuing a bachelor's degree in CS for someone who already holds a PhD in CS) or vary significantly in cost (e.g., moving to the US versus learning Python). Further, the generated rules do not consider global constraints, e.g., if the targeted outcome is the salary in an institution, there may be a budget. Future research will incorporate intervention costs to generate budget-constrained rules and address finite resource allocation scenarios to account for cases where the population size that can achieve improved outcomes is limited.

Extension to multi-table data. FAIRCAP currently supports a single-relation database without dependencies among tuples to ensure compliance with the SUTVA assumption [79] (discussed in Section 3). However, this assumption breaks down even in single-table databases with tuple dependencies. In single-table settings, intervention and grouping patterns are straightforward to define. Extending these definitions to multi-table databases, where grouping attributes and interventions may originate from different tables, introduces a significant challenge. This complexity arises due to many-tomany relationships and cross-table patterns. Previous work, such as [24, 82], has extended causal models to handle multi-table data, but they have not explicitly targeted recommendations for subgroups. Expanding our framework to support multi-relational databases with complex dependencies remains an important direction for future research. Notably, prior work leveraging causal inference [53, 81, 107, 109] has also primarily focused on single-table settings.

Robustness of rules. The generated rules may be influenced by several factors, including the method used to evaluate causal effects, the thresholds set for the constraints, the overall quality of the data, and the quality of the causal DAG. In this work, we assume that the causal DAG is provided as part of the input, with the responsibility for validating its correctness resting on the policymaker. Nonetheless, the causal DAG only needs to specify causal dependencies between variables without detailing the nature of those dependencies. Developing methods that are robust to inaccuracies in the DAG is an important direction for future work.

Explainability and prescriptive causal nature of rules. While if-then rules for prediction or causation are considered explainable or interpretable in the literature [15, 35, 46, 71, 97], we note that no additional explanations or justifications come with the rules mined by FairCap. Generating meaningful explanations to describe how the rules impact the outcome and the variability of the outcome within various sub-populations is deferred to future work.

To conclude, observational causal analysis is the main foundation for any *prescription* or *recommendation* beyond predictions, when a randomized controlled trial is not possible due to cost, ethics, or feasibility issues. However, the analysis depends on assumptions (ignorability, causal DAG) that may not hold in a scenario and one should know the assumptions and limitations of these claims. How the rules should be used in practice considering practical and fairness aspects is a general direction of future work.

ACKNOWLEDGMENTS

This work was partially supported by the NSF awards IIS-2008107 and IIS-2147061, and a grant from Infosys. Additional funding was provided by the Henry and Marilyn Taub faculty for computer science at the Technion.

REFERENCES

- [1] [n. d.]. Git Repository.
- [2] 2021. 2021 Stackoverflow Developer Survey. https://insights.stackoverflow.com/survey/2021.
- [3] Alekh Agarwal, Miroslav Dudík, and Zhiwei Steven Wu. 2019. Fair regression: Quantitative definitions and reduction-based algorithms. In *International Conference on Machine Learning*. PMLR, 120–129.
- [4] Sameet Agarwal, Rakesh Agrawal, Prasad Deshpande, Ashish Gupta, Jeffrey F. Naughton, Raghu Ramakrishnan, and Sunita Sarawagi. 1996. On the Computation of Multidimensional Aggregates. In VLDB. Morgan Kaufmann, 506–521.
- [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] Abdullah Alomar, Pouya Hamadanian, Arash Nasr-Esfahany, Anish Agarwal, Mohammad Alizadeh, and Devavrat Shah. 2023. CausalSim: A Causal Framework for Unbiased Trace-Driven Simulation. In 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23). 1115–1147.
- [7] Sihem Amer-Yahia, Tova Milo, and Brit Youngmann. 2021. Exploring Ratings in Subjective Databases. In Proceedings of the 2021 International Conference on Management of Data. Association for Computing Machinery, New York, NY, USA, 62–75.
- [8] Abolfazl Asudeh, Zhongjun Jin, and HV Jagadish. 2019. Assessing and remedying coverage for a given dataset. In 2019 IEEE 35th International Conference on Data Engineering (ICDE). IEEE, 554–565.
- [9] Arthur Asuncion and David Newman. 2007. UCI machine learning repository.
- [10] Eftychia Baikousi, Georgios Rogkakos, and Panos Vassiliadis. 2011. Similarity measures for multidimensional data. In 2011 IEEE 27th International Conference on Data Engineering. IEEE, IEEE, 171–182.
- [11] Nicole Bidoit, Melanie Herschel, and Katerina Tzompanaki. 2014. Query-based why-not provenance with nedexplain. In Extending database technology (EDBT).
- [12] Reuben Binns. 2020. On the apparent conflict between individual and group fairness. In Proceedings of the 2020 conference on fairness, accountability, and transparency. 514–524.

- [13] Simon Caton and Christian Haas. 2024. Fairness in machine learning: A survey. Comput. Surveys 56, 7 (2024), 1–38.
- [14] Adriane Chapman and HV Jagadish. 2009. Why not?. In Proceedings of the 2009 ACM SIGMOD International Conference on Management of data. 523–534.
- [15] Chaofan Chen and Cynthia Rudin. 2018. An optimization approach to learning falling rule lists. In *International conference on artificial intelligence and statistics*. PMLR, 604–612.
- [16] Daniel Deutch, Nave Frost, and Amir Gilad. 2020. Explaining Natural Language query results. VLDB J. 29, 1 (2020), 485–508.
- [17] Daniel Deutch and Amir Gilad. 2019. Reverse-Engineering Conjunctive Queries from Provenance Examples. In Advances in Database Technology - 22nd International Conference on Extending Database Technology, EDBT 2019, Lisbon, Portugal, March 26-29, 2019, Melanie Herschel, Helena Galhardas, Berthold Reinwald, Irini Fundulaki, Carsten Binnig, and Zoi Kaoudi (Eds.). OpenProceedings.org, 277-288. https://doi.org/10.5441/002/edbt.2019.25
- [18] Daniel Deutch, Amir Gilad, Tova Milo, Amit Mualem, and Amit Somech. 2022. FEDEX: An Explainability Framework for Data Exploration Steps. Proc. VLDB Endow. 15, 13 (2022), 3854–3868. https://www.vldb.org/pvldb/vol15/p3854-gilad.pdf
- [19] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. 2019. The Design and Operation of CloudLab. In Proceedings of the USENIX Annual Technical Conference (ATC). 1–14. https://www.flux.utah.edu/paper/duplyakin-atc19
- [20] Ahmad-Reza Ehyaei, Amir-Hossein Karimi, Bernhard Schölkopf, and Setareh Maghsudi. 2023. Robustness implies fairness in causal algorithmic recourse. In Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency. 984–1001.
- [21] Michael D Ekstrand, Robin Burke, and Fernando Diaz. 2019. Fairness and discrimination in recommendation and retrieval. In Proceedings of the 13th ACM Conference on Recommender Systems. 576–577.
- [22] 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.
- [23] Uriel Feige. 1998. A threshold of ln n for approximating set cover. Journal of the ACM (JACM) 45, 4 (1998), 634–652.
- [24] Sainyam Galhotra, Amir Gilad, Sudeepa Roy, and Babak Salimi. 2022. Hyper: Hypothetical reasoning with what-if and how-to queries using a probabilistic causal approach. In Proceedings of the 2022 International Conference on Management of Data. 1598–1611.
- [25] 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.
- [26] Sainyam Galhotra, Karthikeyan Shanmugam, Prasanna Sattigeri, and Kush R Varshney. 2022. Causal feature selection for algorithmic fairness. In Proceedings of the 2022 International Conference on Management of Data. 276–285.
- [27] Yu Gan, Mingyu Liang, Sundar Dev, David Lo, and Christina Delimitrou. 2021. Sage: Practical and Scalable ML-Driven Performance Debugging in Microservices. In Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems. 135–151.
- [28] Ruoyuan Gao and Chirag Shah. 2020. Counteracting bias and increasing fairness in search and recommender systems. In Proceedings of the 14th ACM Conference on Recommender Systems. 745–747.
- [29] David García-Soriano and Francesco Bonchi. 2021. Maxmin-fair ranking: individual fairness under group-fairness constraints. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining. 436–446.
- [30] Manisha Girotra, Kanika Nagpal, Saloni Minocha, and Neha Sharma. 2013. Comparative survey on association rule mining algorithms. *International Journal of Computer Applications* 84, 10 (2013).
- [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] Yue Gong, Sainyam Galhotra, and Raul Castro Fernandez. 2024. Nexus: Correlation Discovery over Collections of Spatio-Temporal Tabular Data. Proceedings of the ACM on Management of Data 2, 3 (2024), 1–28.
- [33] Helga Gudmundsdottir, Babak Salimi, Magdalena Balazinska, Dan RK Ports, and Dan Suciu. 2017. A demonstration of interactive analysis of performance measurements with viska. In Proceedings of the 2017 ACM International Conference on Management of Data. 1707–1710.
- [34] Shubha Guha, Falaah Arif Khan, Julia Stoyanovich, and Sebastian Schelter. 2024. Automated data cleaning can hurt fairness in machine learning-based decision making. IEEE Transactions on Knowledge and Data Engineering (2024).
- [35] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Dino Pedreschi, Franco Turini, and Fosca Giannotti. 2018. Local rule-based explanations of black box decision systems. arXiv preprint arXiv:1805.10820 (2018).
- [36] Paul W Holland. 1986. Statistics and causal inference. Journal of the American statistical Association 81, 396 (1986), 945–960.

- [37] Zezhou Huang, Jiaxiang Liu, Haonan Wang, and Eugene Wu. 2023. The Fast and the Private: Task-based Dataset Search. arXiv preprint arXiv:2308.05637 (2023).
- [38] Guido W Imbens. 2024. Causal inference in the social sciences. Annual Review of Statistics and Its Application 11 (2024).
- [39] Shomik Jain, Vinith Suriyakumar, Kathleen Creel, and Ashia Wilson. 2024. Algorithmic Pluralism: A Structural Approach To Equal Opportunity. In The 2024 ACM Conference on Fairness, Accountability, and Transparency. 197–206.
- [40] Manas Joglekar, Hector Garcia-Molina, and Aditya Parameswaran. 2017. Interactive data exploration with smart drill-down. IEEE Transactions on Knowledge and Data Engineering 31, 1 (2017), 46–60.
- [41] Gagandeep Kaur and Shruti Aggarwal. 2013. Performance analysis of association rule mining algorithms. International Journal of Advanced Research in Computer Science and Software Engineering 3, 8 (2013), 856–58.
- [42] Samir Khuller, Anna Moss, and Joseph Seffi Naor. 1999. The budgeted maximum coverage problem. *Information processing letters* 70, 1 (1999), 39–45.
- [43] Alexandra Kim, Laks VS Lakshmanan, and Divesh Srivastava. 2020. Summarizing hierarchical multidimensional data. In 2020 IEEE 36th International Conference on Data Engineering (ICDE). IEEE, 877–888.
- [44] Been Kim, Cynthia Rudin, and Julie A Shah. 2014. The bayesian case model: A generative approach for case-based reasoning and prototype classification. Advances in neural information processing systems 27 (2014).
- [45] Trupti A Kumbhare and Santosh V Chobe. 2014. An overview of association rule mining algorithms. International Journal of Computer Science and Information Technologies 5, 1 (2014), 927–930.
- [46] Himabindu Lakkaraju, Stephen H Bach, and Jure Leskovec. 2016. Interpretable decision sets: A joint framework for description and prediction. In Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. 1675–1684.
- [47] Connor Lawless, Sanjeeb Dash, Oktay Gunluk, and Dennis Wei. 2023. Interpretable and fair boolean rule sets via column generation. *Journal of Machine Learning Research* 24, 229 (2023), 1–50.
- [48] Seokki Lee, Bertram Ludäscher, and Boris Glavic. 2020. Approximate summaries for why and why-not provenance (extended version). arXiv preprint arXiv:2002.00084 (2020).
- [49] Chenjie Li, Zhengjie Miao, Qitian Zeng, Boris Glavic, and Sudeepa Roy. 2021. Putting Things into Context: Rich Explanations for Query Answers using Join Graphs. In Proceedings of the 2021 International Conference on Management of Data. 1051–1063.
- [50] Jinyang Li, Yuval Moskovitch, Julia Stoyanovich, and HV Jagadish. 2023. Query Refinement for Diversity Constraint Satisfaction. Proceedings of the VLDB Endowment 17. 2 (2023), 106–118.
- [51] 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.
- [52] Yin Lou, Rich Caruana, Johannes Gehrke, and Giles Hooker. 2013. Accurate intelligible models with pairwise interactions. In Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining. 623– 631
- [53] Pingchuan Ma, Rui Ding, Shuai Wang, Shi Han, and Dongmei Zhang. 2023. Xinsight: explainable data analysis through the lens of causality. Proceedings of the ACM on Management of Data 1, 2 (2023), 1–27.
- [54] Ayan Majumdar and Isabel Valera. 2024. CARMA: A practical framework to generate recommendations for causal algorithmic recourse at scale. In The 2024 ACM Conference on Fairness, Accountability, and Transparency. 1745–1762.
- [55] Markos Markakis, An Bo Chen, Brit Youngmann, Trinity Gao, Ziyu Zhang, Rana Shahout, Peter Baile Chen, Chunwei Liu, Ibrahim Sabek, and Michael Cafarella. 2024. Sawmill: From Logs to Causal Diagnosis of Large Systems. In SIGMOD. 444–447.
- [56] Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. ACM computing surveys (CSUR) 54, 6 (2021), 1–35.
- [57] Alexandra Meliou, Wolfgang Gatterbauer, Joseph Y Halpern, Christoph Koch, Katherine F Moore, and Dan Suciu. 2010. Causality in databases. IEEE Data Engineering Bulletin 33, 3 (2010), 59–67.
- [58] Alexandra Meliou, Wolfgang Gatterbauer, Katherine F Moore, and Dan Suciu. 2009. Why so? or why no? functional causality for explaining query answers. arXiv preprint arXiv:0912.5340 (2009).
- [59] Alexandra Meliou, Wolfgang Gatterbauer, Katherine F Moore, and Dan Suciu. 2010. The complexity of causality and responsibility for query answers and non-answers. arXiv preprint arXiv:1009.2021 (2010).
- [60] Alexandra Meliou, Sudeepa Roy, and Dan Suciu. 2014. Causality and explanations in databases. Proceedings of the VLDB Endowment 7, 13 (2014), 1715–1716.
- [61] Zhengjie Miao, Qitian Zeng, Boris Glavic, and Sudeepa Roy. 2019. Going beyond provenance: Explaining query answers with pattern-based counterbalances. In Proceedings of the 2019 International Conference on Management of Data. 485–502.
- [62] Fatemeh Nargesian, Abolfazl Asudeh, and HV Jagadish. 2021. Tailoring data source distributions for fairness-aware data integration. Proceedings of the

- VLDB Endowment 14, 11 (2021), 2519-2532.
- [63] Fatemeh Nargesian, Abolfazl Asudeh, and H. V. Jagadish. 2022. Responsible Data Integration: Next-generation Challenges. In Proceedings of the 2022 ACM SIGMOD International Conference on Management of data.
- [64] Eliana Pastor, Elena Baralis, and Luca de Alfaro. 2023. A hierarchical approach to anomalous subgroup discovery. In 2023 IEEE 39th international conference on data engineering (ICDE). IEEE, 2647–2659.
- [65] Eliana Pastor, Luca De Alfaro, and Elena Baralis. 2021. Looking for trouble: Analyzing classifier behavior via pattern divergence. In Proceedings of the 2021 International Conference on Management of Data. 1400–1412.
- [66] Judea Pearl. 2009. Causal inference in statistics: An overview. (2009).
- [67] Dana Pessach and Erez Shmueli. 2022. A review on fairness in machine learning. ACM Computing Surveys (CSUR) 55, 3 (2022), 1–44.
- [68] Alireza Pirhadi, Mohammad Hossein Moslemi, Alexander Cloninger, Mostafa Milani, and Babak Salimi. 2024. Otclean: Data cleaning for conditional independence violations using optimal transport. Proceedings of the ACM on Management of Data 2, 3 (2024), 1–26.
- [69] Drago Plecko and Elias Bareinboim. 2023. Causal fairness for outcome control. Advances in Neural Information Processing Systems 36 (2023), 47575–47597.
- [70] Drago Plecko and Elias Bareinboim. 2024. Causal fairness analysis. ECAI (2024).
- [71] Romila Pradhan, Jiongli Zhu, Boris Glavic, and Babak Salimi. 2022. Interpretable data-based explanations for fairness debugging. In Proceedings of the 2022 International Conference on Management of Data. 247–261.
- [72] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. 2020. Fairbatch: Batch selection for model fairness. arXiv preprint arXiv:2012.01696 (2020).
- [73] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. 2023. Improving Fair Training under Correlation Shifts. In International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA (Proceedings of Machine Learning Research, Vol. 202), Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (Eds.). PMLR, 29179–29209.
- [74] 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.
- [75] Sudeepa Roy. 2022. Toward interpretable and actionable data analysis with explanations and causality. Proc. VLDB Endow. 15, 12 (2022), 3812–3820.
- [76] 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.
- [77] 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.
- [78] Donald Bruce Rubin. 1971. The use of matched sampling and regression adjustment in observational studies. Ph. D. Dissertation. Harvard University.
- [79] Donald B Rubin. 2005. Causal inference using potential outcomes: Design, modeling, decisions. J. Amer. Statist. Assoc. 100, 469 (2005), 322–331.
- [80] Omer Sagi and Lior Rokach. 2021. Approximating XGBoost with an interpretable decision tree. *Information Sciences* 572 (2021), 522–542.
- [81] 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.
- [82] Babak Salimi, Harsh Parikh, Moe Kayali, Lise Getoor, Sudeepa Roy, and Dan Suciu. 2020. Causal relational learning. In Proceedings of the 2020 ACM SIGMOD international conference on management of data. 241–256.
- [83] 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.
- [84] Aécio Santos, Aline Bessa, Fernando Chirigati, Christopher Musco, and Juliana Freire. 2021. Correlation sketches for approximate join-correlation queries. In Proceedings of the 2021 International Conference on Management of Data. 1531–1544.
- [85] Gayatri Sathe and Sunita Sarawagi. 2001. Intelligent rollups in multidimensional OLAP data. In VLDB. 307–316.
- [86] Holger Schielzeth. 2010. Simple means to improve the interpretability of regression coefficients. Methods in Ecology and Evolution 1, 2 (2010), 103–113.
- [87] Amit Sharma and Emre Kiciman. 2020. DoWhy: An End-to-End Library for Causal Inference. arXiv preprint arXiv:2011.04216 (2020).
- [88] Jessie J Smith and Lex Beattie. 2022. RecSys Fairness Metrics: Many to Use But Which One To Choose? arXiv preprint arXiv:2209.04011 (2022).
- [89] Seamus Somerstep, Ya'acov Ritov, and Yuekai Sun. 2024. Algorithmic Fairness in Performative Policy Learning: Escaping the Impossibility of Group Fairness. In The 2024 ACM Conference on Fairness, Accountability, and Transparency. 616– 630
- [90] Peter Spirtes, Clark Glymour, and Richard Scheines. 2001. Causation, prediction, and search. MIT press.
- [91] Julia Stoyanovich, Bill Howe, and Hosagrahar Visvesvaraya Jagadish. 2020. Responsible data management. Proceedings of the VLDB Endowment 13, 12

(2020).

- [92] Hao Sun, Evan Munro, Georgy Kalashnov, Shuyang Du, and Stefan Wager. 2021. Treatment allocation under uncertain costs. arXiv preprint arXiv:2103.11066 (2021).
- [93] Tanmay Surve and Romila Pradhan. 2024. Example-based Explanations for Random Forests using Machine Unlearning. arXiv preprint arXiv:2402.05007 (2024).
- [94] Ki Hyun Tae, Hantian Zhang, Jaeyoung Park, Kexin Rong, and Steven Euijong Whang. 2024. Falcon: Fair Active Learning using Multi-armed Bandits. Proc. VLDB Endow. 17, 5 (2024), 952–965.
- [95] Yuchao Tao, Amir Gilad, Ashwin Machanavajjhala, and Sudeepa Roy. 2022. DPXPlain: Privately Explaining Aggregate Query Answers. Proc. VLDB Endow. 16, 1 (2022), 113–126. https://www.vldb.org/pvldb/vol16/p113-tao.pdf
- [96] Balder ten Cate, Cristina Civili, Evgeny Sherkhonov, and Wang-Chiew Tan. 2015. High-level why-not explanations using ontologies. In Proceedings of the 34th ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems. 31-43.
- [97] Jasper van der Waa, Elisabeth Nieuwburg, Anita Cremers, and Mark Neerincx. 2021. Evaluating XAI: A comparison of rule-based and example-based explanations. Artificial intelligence 291 (2021), 103404.
- [98] Moshe Y. Vardi. 1982. The Complexity of Relational Query Languages (Extended Abstract). In Proceedings of the Fourteenth Annual ACM Symposium on Theory of Computing (San Francisco, California, USA) (STOC '82). ACM, New York, NY, USA, 137–146. https://doi.org/10.1145/800070.802186
- [99] Manasi Vartak, Sajjadur Rahman, Samuel Madden, Aditya Parameswaran, and Neoklis Polyzotis. 2015. Seedb: Efficient data-driven visualization recommendations to support visual analytics. In VLDB, Vol. 8. NIH Public Access, 2182.
- [100] Stefan Wager and Susan Athey. 2018. Estimation and inference of heterogeneous treatment effects using random forests. J. Amer. Statist. Assoc. 113, 523 (2018), 1228–1242.
- [101] Tong Wang and Cynthia Rudin. 2022. Causal rule sets for identifying subgroups with enhanced treatment effects. INFORMS Journal on Computing 34, 3 (2022), 1626–1643.
- [102] Eugene Wu and Samuel Madden. 2013. Scorpion: Explaining away outliers in aggregate queries. (2013).
- [103] Yu Xie, Jennie E Brand, and Ben Jann. 2012. Estimating heterogeneous treatment effects with observational data. Sociological methodology 42, 1 (2012), 314–347.
- [104] Hongyu Yang, Cynthia Rudin, and Margo Seltzer. 2017. Scalable Bayesian rule lists. In International conference on machine learning. PMLR, 3921–3930.
- [105] Brit Youngmann, Sihem Amer-Yahia, and Aurélien Personnaz. 2022. Guided Exploration of Data Summaries. Proc. VLDB Endow. 15, 9 (2022).
- [106] Brit Youngmann, Michael Cafarella, Amir Gilad, and Sudeepa Roy. 2024. Summarized Causal Explanations For Aggregate Views (Full version). arXiv:2410.11435 [cs.DB] https://arxiv.org/abs/2410.11435
- [107] Brit Youngmann, Michael Cafarella, Yuval Moskovitch, and Babak Salimi. 2023. On Explaining Confounding Bias. 2023 IEEE 39th International Conference on Data Engineering (ICDE) (2023).
- [108] Brit Youngmann, Michael Cafarella, Babak Salimi, and Anna Zeng. 2023. Causal Data Integration. arXiv preprint arXiv:2305.08741 (2023).
- [109] Brit Youngmann, Michael J. Cafarella, Amir Gilad, and Sudeepa Roy. 2024. Summarized Causal Explanations For Aggregate Views. Proc. ACM Manag. Data 2. 1 (2024), 71:1–71:27.
- [110] Meike Zehlike, Francesco Bonchi, Carlos Castillo, Sara Hajian, Mohamed Megahed, and Ricardo Baeza-Yates. 2017. Fa* ir: A fair top-k ranking algorithm. In Proceedings of the 2017 ACM on Conference on Information and Knowledge Management. 1569–1578.
- [111] Meike Zehlike, Ke Yang, and Julia Stoyanovich. 2022. Fairness in ranking, part i: Score-based ranking. Comput. Surveys 55, 6 (2022), 1–36.
- [112] Meike Zehlike, Ke Yang, and Julia Stoyanovich. 2022. Fairness in ranking, part ii: Learning-to-rank and recommender systems. Comput. Surveys 55, 6 (2022), 1–41
- [113] Hantian Zhang, Ki Hyun Tae, Jaeyoung Park, Xu Chu, and Steven Euijong Whang. 2023. iFlipper: Label Flipping for Individual Fairness. Proc. ACM Manag. Data 1, 1 (2023), 8:1–8:26.
- [114] Mengdi Zhang and Jun Sun. 2022. Adaptive fairness improvement based on causality analysis. In Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering. 6–17
- [115] Xiaozhong Zhang, Xiaoyu Ge, Panos K Chrysanthis, and Mohamed A Sharaf. 2021. Viewseeker: An interactive view recommendation framework. Big Data Research 25 (2021), 100238.
- [116] Qiankun Zhao and Sourav S Bhowmick. 2003. Association rule mining: A survey. Nanyang Technological University, Singapore 135 (2003), 18.

9 MISSING PROOFS

Proof of Lemma 4.1. The utility of a rule r denotes the expected increase in outcome O when all individuals within the subgroup \mathcal{P}_q are

treated with \mathcal{P}_t .

$$utility(r) = \frac{1}{|\mathcal{P}_g|} \sum_{i \in coverage(\mathcal{P}_g)} utility_i(\mathcal{P}_t)$$

where $utilit y_i(\mathcal{P}_t)$ denotes the utility for tuple i with respect to treatment \mathcal{P}_t . Since utility(r) is an average over multiple different utilities, the utility will be higher than the expected value for certain tuples in $coverage(\mathcal{P}_g)$. Let $i^* = \arg\max utilit y_i(\mathcal{P}_t)$.

Consider a new prescription rule $r'(i^*, \mathcal{P}_t)$ which considers the same treatment \mathcal{P}_t for the tuple i^* . Therefore, $utility(r') = utility_i(\mathcal{P}_t) > utility(r)$.

9.1 Hardness Results

We next study the complexity of the Prescription Ruleset Selection problem under different constraint combinations. We show that Prescription Ruleset Selection is equivalent to optimizing a non-negative and monotone sub-modular function. Furthermore, the individual fairness constraint and rule coverage constraints are matroid constraints. Therefore, a greedy algorithm is appropriate approach to solve the problem.

Proposition 9.1. The optimization objective of Prescription Ruleset Selection problem is a non-negative submodular function.

Proof of Proposition 9.1. According to [46], the size objective is a non-negative and submodular function. Similarly, the expected utility—assuming each individual receives a single rule and selects the best option—is also a non-negative submodular function. As a result, the linear combination of these functions remains a non-negative submodular function, and maximizing it is known to be NP-hard [42].

Proposition 9.2. Individual fairness and rule coverage constraints are matroid constraint

PROOF OF PROPOSITION 9.2. We will show these constraints satisfy the following properties:

- (1) **Hereditary Property**: If *S* is an independent set, then every subset of *S* is also an independent set.
- (2) **Exchange Property**: If *S* and *T* are independent sets and |S| < |T|, then there exists an element $e \in T \setminus S$ such that $S \cup \{e\}$ is also an independent set.

These two properties ensure that the set system behaves like a matroid.

In our setting Start by specifying the ground set is all possible rules and what qualifies as an "independent set" is a subset of rules satisfying a constraint.

Individual Fairness. If a set of rules R satisfies the individual fairness constraint, this means each rule within R individually satisfies the constraint. Consequently, any subset $R' \subseteq R$ also upholds individual fairness. This further implies the exchange property, as any rule that satisfies individual fairness can be added to an individually fair set of rules while preserving individual fairness.

Rule coverage. If a set of rules R satisfies the rule coverage constraint, this means each rule within R individually satisfies the rule coverage constraint. Consequently, any subset $R' \subseteq R$ also satisfies the rule coverage constraint. This also implies the exchange property, as any rule that satisfies rule coverage can be added to a set of rules satisfying rule coverage while preserving the rule coverage constraints.

For the group coverage, we can show that merely finding a solution that satisfies the constraints, even without maximizing expected utility, is NP-hard via a reduction from the Set Cover problem [23].

Proposition 9.3. Prescription Ruleset Selection with a group-coverage constraint is NP-hard

PROOF of 9.3. In the decision version of the Set Cover problem, we are given a universe of elements $U = \{x_1, \ldots, x_{n'}\}$, a collection of m subsets $S_1, \ldots, S_{m'} \subseteq U$ and a number k. The question is whether there exists a cover of U of at most k' subsets.

In the decision version of Prescription Ruleset Selection we are searching for a set of rules R such that: $f(R) \ge \tau$, where:

$$f(R) = \lambda_1 \cdot (l - size(R)) + \lambda_2 \cdot ExpUtility(R)$$

such that R satisfies the group-coverage constraints, defined by the parameter θ . In this proof, we assume no protected group is given, namely the constraint requires that the selected ruleset R would cover at least a θ fraction of the population (i.e., the protected group to be the empty group).

Given an instance of the set cover problem, we build an instance of the Prescription Ruleset Selection problem as follows. We build a relation R with m'+1 attributes, $A = (A_1, \ldots, A_{m'}, O)$, and containing n'+m' tuples. For each element $x_i \in U$, we create a tuple t_i , such that $t_i[A_j] = 1$ iff $x_i \in S_j$. We further add m' tuples t_{S_j} such that $t_{S_j}[A_j] = 1$, $t_{S_j}[O] = 0$, and $t_{S_j}[A_p] = l \neq 0$ for all $p \neq j$ where l is a unique number not used anywhere else in an attribute of R. We set the outcome variable to be O.

Here, \mathcal{P}_g can be any predicate. Note that each set of tuples defined by a pattern can only have an outcome of 0, as the outcome of all tuples is 0. Therefore, the utility of all intervention patterns is 0. For Prescription Ruleset Selection, we further define the threshold for the group coverage constraint $\theta = \frac{n'+k'}{n'+m'}$. The underlying causal DAG, G, only contains the edges of the form $A_j \to O$ for all $1 \le j \le m'$. We claim that there exists a cover of U with at most k sets iff there exists a solution R to Prescription Ruleset Selection such that $f(R) \ge (l-k)$.

(⇒) Assume we have a collection S_{j_1}, \ldots, S_{j_k} such that $\cup_{j=j_1}^{J_k} S_j = U$. We show that there is a solution for Prescription Ruleset Selection as follows. For each S_{j_1} , we choose for the solution the pattern $\mathcal{P}_g^{J_i}: A_{J_i} = 1$. We show that $R = \{(\mathcal{P}_g^{J_1}, \emptyset), \ldots, (\mathcal{P}_g^{J_k}, \emptyset)\}$ is a solution to Prescription Ruleset Selection. The intervention pattern can be any pattern, as the utility of

every intervention pattern is 0. First, we note that all tuples of the form t_i are covered by at least one grouping pattern by their definition. For the m remaining tuples, we have coverage of at most k tuples. These are the tuples $t_{S_{j_i}}$ that have $A_{j_i} = 1$. Thus, the number of covered tuples is exactly n' + k' out of n' + m' tuples in R. If there are fewer than k tuples we can augment the original cover with arbitrary sets to obtain a cover of size k.

(⇐) Assume we have a solution R to Prescription Ruleset Selection with the aforementioned parameters. We show that we can find a solution to the set cover problem. First, note that since $f(R) \ge (l-k)$ and the expected utility is always 0, that means we have selected no more than k rules.

Suppose $R = \{(\mathcal{P}_g^{j_1}, \emptyset), \dots, (\mathcal{P}_g^{j_k}, \emptyset)\}$. We first claim that no grouping pattern that includes $A_i = 0$ in a conjunction can be included in R as such a pattern will not cover any tuple t_{S_i} since these tuples do not have an attribute with value 0 by definition (and any other number other than 1 will only cover a single tuple). Thus, the number of covered tuples will be $<\frac{n'+k'}{n'+m'}=\theta$, which would contradict the assumption that this is a valid solution to Prescription Ruleset Selection that satisfies the coverage constraint. Hence, all patterns are of conjunctions of $A_i = 1$. For each intervention pattern of the form $\mathcal{P}_g = \bigwedge_{j=i_a}^{i_b} (A_j = 1) \wedge (A_p = l)$, we choose an arbitrary attribute in the conjunction A_i if $(A_i = 1) \in \mathcal{P}_q$ and choose S_i for the cover. Finally, if there is an uncovered element x in U and R includes a pattern in of the form $\mathcal{P}_q = (A_i = l)$ where $l \neq 1$, we choose for the cover a set S that covers x arbitrarily. We claim that the chosen collection of sets is a cover of U. To see this, recall that we claimed that the coverage of R is at least n + k. If the coverage includes tuples of the form t_{S_i} , then each pattern covers a single tuple. Suppose these patterns are $\mathcal{P}_a, \ldots, \mathcal{P}_b$. When building the coverage, instead of these patterns, we add a set that covers elements not yet covered by existing patterns. Thus, there are at least b-a covered elements from U in addition to the n+k-(b-a)tuples covered by the patterns. Thus, the set cover we have assembled contains n - (b - a) + (b - a) = n elements and covers all elements in