

XINSIGHT: eXplainable Data Analysis Through The Lens of Causality

PINGCHUAN MA, Hong Kong University of Science and Technology, Hong Kong SAR

RUI DING*, Microsoft Research, China

SHUAI WANG, Hong Kong University of Science and Technology, Hong Kong SAR

SHI HAN, Microsoft Research, China

DONGMEI ZHANG, Microsoft Research, China

In light of the growing popularity of Exploratory Data Analysis (EDA), understanding the underlying causes of the knowledge acquired by EDA is crucial. However, it remains under-researched. This study promotes a transparent and explicable perspective on data analysis, called *eXplainable Data Analysis* (XDA). For this reason, we present XINSIGHT, a general framework for XDA. XINSIGHT provides data analysis with qualitative and quantitative explanations of causal and non-causal semantics. This way, it will significantly improve human understanding and confidence in the outcomes of data analysis, facilitating accurate data interpretation and decision making in the real world. XINSIGHT is a three-module, end-to-end pipeline designed to extract causal graphs, translate causal primitives into XDA semantics, and quantify the quantitative contribution of each explanation to a data fact. XINSIGHT uses a set of design concepts and optimizations to address the inherent difficulties associated with integrating causality into XDA. Experiments on synthetic and real-world datasets as well as a user study demonstrate the highly promising capabilities of XINSIGHT.

CCS Concepts: • **Information systems** → **Data analytics; Data management systems**; • **Mathematics of computing** → **Bayesian networks**.

ACM Reference Format:

Pingchuan Ma, Rui Ding, Shuai Wang, Shi Han, and Dongmei Zhang. 2023. XINSIGHT: eXplainable Data Analysis Through The Lens of Causality. *Proc. ACM Manag. Data* 1, 2, Article 156 (June 2023), 35 pages. <https://doi.org/10.1145/3589301>

1 INTRODUCTION

Exploratory data analysis (EDA) is key to acquiring insight from data and facilitating analysis towards decision making [38, 28]. With the advent of the digital age, the information explosion phenomenon [8] makes it difficult for users to justify and rely on knowledge and conclusions from EDA. To ease the cognitive process, data explanations are proposed to deliberate data facts (e.g., query outcomes) and enhance user comprehension [17]. In this paper, we term such a process as *eXplainable Data Analysis* (XDA), which advances data analysis by providing users with effective

*Corresponding author.

Authors' addresses: Pingchuan Ma, Hong Kong University of Science and Technology, Kowloon, Hong Kong SAR, pmaab@cse.ust.hk; Rui Ding, Microsoft Research, Beijing, China, juding@microsoft.com; Shuai Wang, Hong Kong University of Science and Technology, Kowloon, Hong Kong SAR, shuaiw@cse.ust.hk; Shi Han, Microsoft Research, Beijing, China; Dongmei Zhang, Microsoft Research, Beijing, China.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2836-6573/2023/6-ART156 \$15.00

<https://doi.org/10.1145/3589301>

explanations. By suggesting and justifying choices to alter outcomes, XDA helps users comprehend and trust phenomena emerging from data; as a result, it facilitates real-world decision making.

Explanations can be categorized as either causal or non-causal [22]. Causal explanations seek causal factors to explain an outcome. Fig. 1 depicts a hypothetical lung cancer dataset. Here, a patient's location (indicating regional tobacco control policy) and amount of stress have an impact on whether they would smoke. Then, smoking influences lung cancer's severity. The degree of severity further affects whether they would undergo surgery and the five-year survival rate. Here, smoking explains why a patient has high lung cancer severity (see Fig. 1(f)). In contrast, a non-causal explanation shows the results merely by statistical correlations. For example, surgery "explains" (more precisely, is relevant to) lung cancer severity (see Fig. 1(g)). Despite being helpful, this is not a causal explanation [45].

Existing data explanation tools (e.g., Tableau's Explain Data [13] in industry, Scorpion [54] and DIFF [1] in academia) often provide non-causal explanations [17]. Although valuable for data analysis, they may mislead users who want causal explanations. A well-known confusion, as noted in [23], is that Tableau's Explain Data reports that Massachusetts' low teenage pregnancy rate may explain this state's high ACT Math score. Such explanations are questionable. In comparison, causal explanations play a central role in human cognition [21, 39]. They enable users to make counterfactual thinking and actionable decisions. For instance, quitting smoking reduces lung cancer severity whereas cancelling surgery does not.

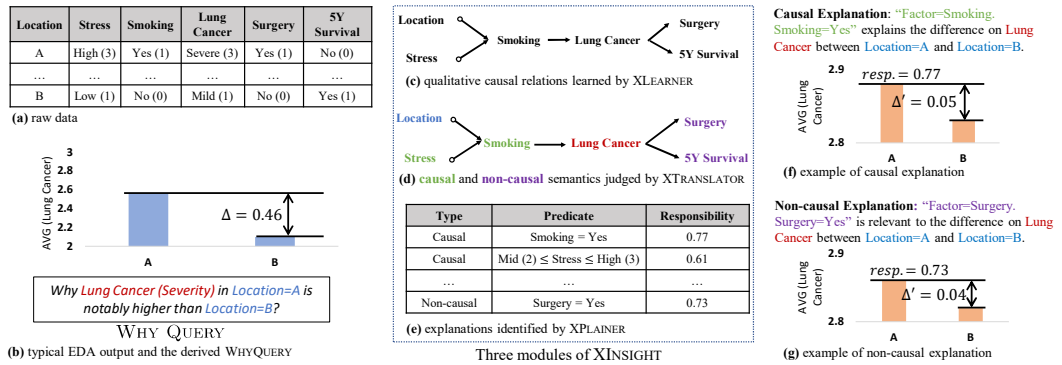


Fig. 1. Illustrative example of XINSIGHT.

According to Pearl's general causality [40], causal knowledge is typically represented by causal graphs. Each node in a causal graph represents a random variable in data, and each directed edge between parent and child nodes denotes a cause-effect relation. Causal knowledge primarily conveys *qualitative explanations* [50], such as smoking causes lung cancer. To further enable *quantitative explanations*, it is necessary to quantify the contribution of each input to the output. This way, we quantify smoking's contribution to lung cancer and compare it with other factors. Halpern's actual causality [19, 18], and essentially its adaptation, DB Causality [30], provide an elegant formulation of this concept.

This paper proposes XINSIGHT as a unified, causality-based XDA framework that qualitatively and quantitatively answers WHY QUERY raised by users. Considering the following WHY QUERY:

Example 1.1: The dataset in Fig. 1(a) depicts the patient information in a country. An analyst observes an interesting data fact: "the average severity of lung cancer of patients in location A is much higher than in location B" and then raises a WHY QUERY (Fig. 1(b)).

for which Fig. 1(f) and (g) illustrate two explanations provided by XINSIGHT. Each explanation is flagged as either a “causal” or “non-causal” explanation, and is composed of both qualitative and quantitative sub-explanations.

Example 1.2: An explanation (Fig. 1(f)) deems “Smoking” to be a qualitative causal factor of lung cancer severity and highlights “Smoking=Yes” and its responsibility as a quantitative sub-explanation.

XINSIGHT includes three modules, XLEARNER, XTRANSLATOR, and XPLAINER to gradually form explanations. XLEARNER first automatically discovers a causal graph \mathcal{G} from data (Fig. 1(c)). Then, given a WHY QUERY (Fig. 1(b)) with the target (i.e., the measure “Lung Cancer”) and the context (i.e., the breakdown dimension “Location”), XTRANSLATOR enumerates each remaining variable on \mathcal{G} and decides if it is causal or non-causal to the target under the context. Fig. 1(d) shows causal variables (i.e., those that can potentially provide causal explanations) in green and non-causal variables (i.e., those that can potentially provide non-causal explanations) in purple. Last, XPLAINER quantifies how well each variable answers WHY QUERY by searching possible predicates on the variables that are the most responsible, as shown in Fig. 1(e). Despite the promising capability of XINSIGHT, concretizing each module is challenging. We brief the challenges and our solutions in the following.

XLEARNER. Most real-world datasets are collected irrespective of causal sufficiency [43]. In other words, not all causally relevant variables are available in the dataset. Furthermore, real-world datasets often contain deterministic relations in the form of Functional Dependency (FD), especially when they have materialized from relational databases. These FDs may violate the faithfulness assumption [12], which is crucial for many causal discovery algorithms. To address these challenges, we establish a theory to propose an FD-induced graph \mathcal{G}_{FD} . XLEARNER uses \mathcal{G}_{FD} to select a subset of variables for standard causal discovery where the selected variables do not trigger faithfulness violations induced by FDs. It adopts FCI [57] to address causal insufficiency and synergistically combine the result of FCI with the causal relations entailed by \mathcal{G}_{FD} .

XTRANSLATOR. The translation from causal primitives (the structural relations in the causal graph) into XDA semantics (e.g., whether a variable provides causal explanations) is under-explored. Given a WHY QUERY (with a target and a context), it is unclear how to determine if a variable X can explain the target given the context, and, moreover, if X provides causal or non-causal explanations. XTRANSLATOR characterizes various causal primitives (e.g., m -separation, ancestor/descendant relations) from a causal graph and provides a taxonomy to translate them into XDA semantics.

XPLAINER. DB causality is primarily designed for data provenance, which usually provides tuples as explanations. Contrarily, we note that predicate-level explanations shall be more desirable for XDA scenarios. Moreover, computing the responsibility of explanations with DB causality is NP-complete in general [33]. XPLAINER adapts DB causality to XDA by using predicate-level explanations with the conciseness consideration and also significantly reduces the computing cost with theoretical guarantees. In summary, we make the following contributions:

- We propose XINSIGHT, a unified and causality-based framework for XDA. XInsight features adequate (by distinguishing causal from non-causal) and comprehensive (with qualitative and quantitative) explanations.
- XINSIGHT consists of three modules, XLEARNER, XTRANSLATOR and XPLAINER, each of which is meticulously designed to address technical challenges and deliver efficient analysis. XLEARNER learns the causal graph from causally insufficient data in the presence of FD-induced faithfulness violations, XTRANSLATOR translates causal primitives into XDA semantics, and XPLAINER efficiently provides quantitative explanations via an adaptation of DB causality to meet the needs of XDA scenarios.

- Empirically, we conduct thorough experiments on public data, production data, and synthetic data via quantitative experiments and human evaluations. The results are very encouraging.

Open Source and Real-world Adoption. We release our code at [36]. XPLAINER has been integrated into Microsoft Power BI to explain increase/decrease in data [37].

2 PRELIMINARY

2.1 Data Model and Query

Multi-Dimensional Data. Let $D := \{X_1, \dots, X_n\}$ represents multi-dimensional data comprising n attributes. In XINSIGHT, we assume that records of D are drawn independently from an identical distribution without selection biases (i.e, i.i.d. assumption) such that each attribute is a (random) variable. Here, selection bias is a preferential selection of units in data analysis [5]. A variable is either categorical or numerical. In accordance with previous works [28, 11], we denote a categorical variable as *dimension* and a numerical variable as *measure*. Multi-dimensional data is commonly represented as a spreadsheet in our context. For relational data, we anticipate taking a materialized provenance table [24] as input.

Aggregation and Discretization on Measure. Given a measure M , users may perform aggregation operations (such as SUM and AVG in SQL) over a set of realizations of M . In some cases, measures are processed in the form of a dimension (e.g., use measures for explanations), which necessitates discretization. It transforms numerical values into several discrete bins that form a categorical variable.

Filter. In this paper, filter is the basic unit of data operations. Given a multi-dimensional data D and a dimension X , a filter $p_i = \{X = x_i\}$ (e.g., “Smoking=Yes”) implies an equality assertion to X such that the value of X shall equal x_i .

Predicate. The disjunction of filters applied on the same dimension is a predicate. Given the dimension X , the predicate $P(x_1, \dots, x_k)$ is a set containment assertion $\{X = x_1 \vee \dots \vee X = x_k\} \equiv \{p_1, \dots, p_k\}$. On a discretized measure, a predicate is an assertion on ranges. A filter is a special case of a predicate. For clarity, we represent a general predicate with a capital P and a filter with a lower-case p .

Subspace. A subspace is a conjunction of filters on disjointed dimensions. Given multi-dimensional data D , a subspace corresponds to a subset of rows satisfying the conditions of all filters. If two subspaces only differ in one filter, they are regarded as *siblings*. The term *Context* refers to the variables of two sibling subspaces, where the *background variables* are the variables with the shared filters and the *foreground variables* are the variables with the different filters. In the following example, we provide a simple instantiation.

Example 2.1: Consider the preceding dataset in Fig. 1(a). $s = \{\text{Location} = A \wedge \text{Lung Cancer} = \text{Severe}\}$ represents the subspace denoting all patients in “Location=A” with severe lung cancer. All patients in “Location=A” with severe lung cancer and all patients in “Location=B” with severe lung cancer form a pair of sibling subspaces. Here, “Location” is the foreground variable and “Lung Cancer” is the background variable.

Selection. We use the following notation to represent the selection procedure over multi-dimensional data D . The subset of data after the selection operation is defined as D_{p_i} , D_P , or D_s , where p_i is a filter, P is a predicate, and s is a subspace. We define $D - D'$ as the rows remaining in D after removing those from D' .

WHY QUERY and Explanation. As illustrated in Fig. 1(b), the user would issue a WHY QUERY to XINSIGHT for explanation. We formally define WHY QUERY as follows.

DEFINITION 2.1 (WHY QUERY). *Given a multi-dimensional data D , a user launches aggregate query $\text{agg}()$ on a target measure M under two sibling subspaces s_1, s_2 . **WHY QUERY** is defined as $\Delta_{s_1, s_2, M, \text{agg}}(D) = \text{agg}_M(D_{s_1}) - \text{agg}_M(D_{s_2})$. For brevity, we use $\Delta(D)$ as the shorthand of $\Delta_{s_1, s_2, M, \text{agg}}(D)$. W.l.o.g., we assume Δ is always non-negative.*

Example 2.2: As shown in Fig. 1(b), we concretize the **WHY QUERY** Δ with the AVG aggregate on the target “Lung Cancer” over two sibling subspaces $s_1 = \{\text{Location} = A\}$, $s_2 = \{\text{Location} = B\}$, denoting the difference in average lung cancer severity in “A” and “B”.

Indeed, explaining the difference between two aggregate queries is one prevalent data analysis task. Identifying the cause in data difference constitutes the basis of many data explanation applications, such as outlier explanation and data debugging [17]. In accordance with prior works [54, 1, 17], we concretize the problem of XDA by concentrating on the explanation of *data difference*. The following form is used to provide explanations in response to **WHY QUERY**.

DEFINITION 2.2 (EXPLANATION). *Given a **WHY QUERY**, an explanation is represented by the following triplet*

$$\text{explanation} := \langle \text{type}, \text{predicate}, \text{responsibility} \rangle \quad (1)$$

where $\text{type} \in \{\text{causal}, \text{non-causal}\}$ denotes whether the explanation is causal or non-causal, the predicate is the content of the explanation, and responsibility, a score ranging from 0 to 1 quantifies the extent to which the explanation explains the given **WHY QUERY**.

Example 2.3: Fig. 1(e) lists several explanations to the **WHY QUERY**. Fig. 1(f)-(g) visualize two of them. Fig. 1(f), as a causal explanation, depicts that “Smoking=Yes” causes the lung cancer severity difference in Location A and B with a responsibility of 0.77.

Single- vs. Multi-Dimensional Explanation. For conciseness and clarity, we anticipate that each explanation reflects one aspect contributing to the outcome when explaining the **WHY QUERY**. We recommend adopting a single-dimensional explanation in XINSIGHT due to its unambiguous causal semantics, although it is feasible to extend an explanation as multi-dimensional using the Cartesian product. The joint causal semantics of several variables, however, could be obscure. Furthermore, multiple single-dimensional explanations (e.g., Fig. 1(e)) suffice to represent a multi-dimensional case.

Functional Dependency (FD). Functional dependency relations are common in multi-dimensional data. In a relational database, among the attributes, there may exist primary keys and foreign keys. Therefore, after materialization, the resulting multi-dimensional data may have functional dependencies. A functional dependency between X and Y is represented by $X \xrightarrow{\text{FD}} Y$. FD, as a deterministic relation among two variables, deems a form of reliable knowledge. This research focuses on one-to-one and one-to-many FDs. We present a simple exemplary dataset that contains FDs.

Example 2.4: Let **CityInfo** be a dataset with three attributes (i.e., City, State, Country). It has three FDs, namely, $\text{City} \xrightarrow{\text{FD}} \text{State}$, $\text{State} \xrightarrow{\text{FD}} \text{Country}$, and $\text{City} \xrightarrow{\text{FD}} \text{Country}$.

FD-Induced Graph. Given a multi-dimensional data D and its functional dependencies, the FD-induced Graph $\mathcal{G}_{\text{FD}} := (V, E)$, where $V := \{X_i \mid \forall X_i \in D\}$ and $E := \{(X_i, X_j) \mid \text{if } X_i \xrightarrow{\text{FD}} X_j\}$. We assume \mathcal{G}_{FD} to be acyclic. Cycles in \mathcal{G}_{FD} imply redundant attributes; in such cases, we retain only one of them to ensure acyclicity.

2.2 Causal Discovery with Latent Variables

This section presents terminology essential to causal discovery with latent variables, such as the representation of causal graphs under causal insufficiency and typical assumptions in causal discovery.

Causal Sufficiency. Causal discovery aims to learn the causal relations from the observational data. Most causal discovery algorithms assume a sufficient observation of the underlying data generating process [51]. Formally, a set of variables X is said to be causally sufficient if there is no hidden variable $Z \notin X$ that is causing more than one variable in X . In other words, it assumes that latent confounders — the shared causes among two or multiple variables — do not exist. However, the process used to acquire real-world data does not provide such guarantees, thereby often yielding causally insufficient observations. Hence, the causal discovery procedure is compromised by the spurious association between two variables sharing a latent confounder. We present an example below.

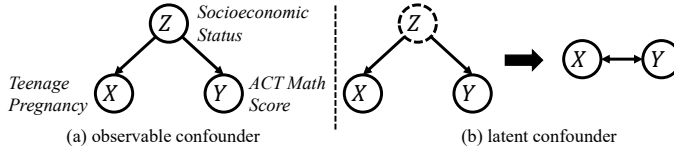


Fig. 2. Examples of observable and latent confounders.

Example 2.5: Consider the hypothetical causal graph in Fig. 2(a) where the socioeconomic status (Z) simultaneously causes teenage pregnancy (X) and their ACT math scores (Y). The socioeconomic status, however, does not appear in the dataset. This absence yields an insufficient observation (the left-side causal graph in Fig. 2(b)), which further results in a spurious association [41] of teenage pregnancy and ACT math scores (the bidirected edge in Fig. 2(b)).

Hence, popular directed acyclic graphs are not expressive enough to represent these subtle relations. This necessitates the Maximal Ancestral Graph [51], which is introduced shortly.

Notation and Terminology. Recall that we assume the dataset is *i.i.d.* with potential latent confounders and does not contain selection bias. Maximal Ancestral Graph (MAG) forms the standard representation of causal graphs in this setting. We now introduce important concepts of graphical models and properties of MAG.

A directed mixed graph \mathcal{G} is a graphical model that contains nodes X and two types of edges, including directed (\rightarrow) and bidirected (\leftrightarrow). There is at most one edge between any two nodes. For each directed edge $X \rightarrow Y$, X is a *parent* of Y and Y is a *child* of X . X and Y are *adjacent* if there is an edge (either directed or bidirected) between them. A *path* \mathcal{P} is a sequence of distinct nodes (X_1, \dots, X_k) where X_i and X_{i+1} are adjacent in \mathcal{G} for all $1 \leq i < k$. A path $\mathcal{P} = (X_1, \dots, X_k)$ is directed if X_i is a parent of X_{i+1} for all $1 \leq i < k$. X is an ancestor of Y if there exists a directed path from X to Y and Y is a descendant of X accordingly. Given a path (X_1, \dots, X_k) , a non-endpoint node X_i is a *collider* if there are arrowheads pointing to X_i from both X_{i-1} and X_{i+1} . Below, we list all possible cases of a collider.

Example 2.6: Given (X_{i-1}, X_i, X_{i+1}) , X_i is a collider if and only if a) $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$, or b) $X_{i-1} \leftrightarrow X_i \leftarrow X_{i+1}$, or c) $X_{i-1} \rightarrow X_i \leftrightarrow X_{i+1}$, or d) $X_{i-1} \leftrightarrow X_i \leftrightarrow X_{i+1}$. In Fig. 1(c), Smoking is a collider of Location and Stress since “Location $\circ \rightarrow$ Smoking $\leftarrow \circ$ Stress”, where \circ represents an undetermined edge endpoint.

A path (X, W_1, \dots, W_k, Y) is said to be *blocked* by $Z \subseteq X \setminus \{X, Y\}$ if there exists a node $W_i \in \{W_1, \dots, W_k\}$ such that a) W_i is not a collider but a member of Z , or b) W_i is a collider but not an ancestor of any nodes of Z . We now introduce *m-separation* and MAG.

DEFINITION 2.3 (M-SEPARATION [57]). X, Y are *m-separated* by Z (denoted by $X \perp\!\!\!\perp_{\mathcal{G}} Y \mid Z$) if all paths between X, Y are blocked by Z .

Example 2.7: Consider the causal graph in Fig. 1(c) where “Smoking” blocks the only path between “Location” and “Lung Cancer”. Hence, “Smoking” *m-separates* “Location” and “Lung Cancer” (denoted by $\text{Lung Cancer} \perp\!\!\!\perp_{\mathcal{G}} \text{Location} \mid \text{Smoking}$).

DEFINITION 2.4 (MAXIMAL ANCESTRAL GRAPH [57]). A directed mixed graph is called a MAG if a) it contains no directed cycles or almost directed cycles and b) for each pair of non-adjacent nodes, there exists a set of nodes that *m-separates* them. A directed cycle refers to the case where $X \rightarrow Y \rightarrow \dots \rightarrow X$ and an almost directed cycle refers to the case where $X \rightarrow Y \rightarrow \dots \rightarrow Z \leftrightarrow X$.

Then, the Global Markov Property (GMP) is developed to provide a probabilistic interpretation of *m-separation*.

DEFINITION 2.5 (GLOBAL MARKOV PROPERTY [51]).

$$X \perp\!\!\!\perp_{\mathcal{G}} Y \mid Z \Rightarrow X \perp\!\!\!\perp Y \mid Z \quad (2)$$

As aforementioned, *m-separation* indicates that all paths between X and Y are “blocked” by Z . Hence, it is intuitive that, if X and Y are *m-separated*, their statistical correlation is also “blocked” by Z . The term *conditional independence* (i.e., $X \perp\!\!\!\perp Y \mid Z$) depicts this absence of statistical correlation. Statistically, $X \perp\!\!\!\perp Y \mid Z$ implies that $P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$, which can be empirically examined using statistical hypothesis tests (e.g., χ^2 tests).

Example 2.8: Consider the dataset in Fig. 1(a). According to GMP, the *m-separation* in Ex. 2.7 implies that, for the dataset in Fig. 1(a), “Location” and “Lung Cancer” are conditionally independent given “Smoking” in a statistical sense.

With GMP, we can deduce statistical conditional independence in data from *m-separations*. Note that only data is available when performing causal discovery. Hence, we need to invert GMP and establish a connection from data distribution to the graphical structure. Faithfulness assumption establishes such connection.

DEFINITION 2.6 (FAITHFULNESS [51]).

$$X \perp\!\!\!\perp Y \mid Z \Rightarrow X \perp\!\!\!\perp_{\mathcal{G}} Y \mid Z \quad (3)$$

According to faithfulness, if we observe that two variables are conditionally independent by a set of variables in data, then they are *m-separated* by the same set of variables on the causal graph. Faithfulness and GMP together establish the equivalence between conditional independence and *m-separation* and they form the key to causal discovery. In addition, we define *skeleton* as follows.

DEFINITION 2.7 (SKELETON). The *skeleton* \mathcal{S} of a MAG \mathcal{G} is the undirected graph obtained by removing all arrowheads from \mathcal{G} .

Constraint-based Causal Discovery. Constraint-based approaches are the standard solution to causal discovery. With the faithfulness assumption, these methods exploit the conditional independence derived from data and gradually establish a MAG \mathcal{G} . \mathcal{G} is consistent with all *m-separations* entailed by conditional independence. However, there may exist multiple MAGs that are equally consistent with the *m-separations* and not distinguishable, which is called Markov equivalence class, denoted by $[\mathcal{G}]$. It is worth noting that these feasible MAGs share the same skeleton while differing in direction on certain edges. These MAGs are therefore summarized into a compact representation called Partial Ancestral Graph (PAG) with some undetermined edge endpoints.

Table 1. Four types of edges in PAG. Circle represents undetermined edge endpoint (can be either an arrowhead or tail).

Edge	Causal Semantics
$X \rightarrow Y$	X is a cause of Y .
$X \leftrightarrow Y$	neither X nor Y is a cause of each other but they share a latent common cause.
$X \circ \rightarrow Y$	1) X is a cause of Y ; or 2) neither X nor Y is a cause of each other but they share a latent common cause.
$X \circ \circ Y$	1) X may be a cause of Y ; or 2) Y may be a cause of X ; or 3) neither X nor Y is a cause of each other but they share a latent common cause.

DEFINITION 2.8 (PARTIAL ANCESTRAL GRAPH [57]). Let $[\mathcal{G}]$ be a Markov equivalence class of a MAG \mathcal{G} . A PAG for $[\mathcal{G}]$ is a graph \mathcal{P} with three possible edge endpoints (namely, tail, circle and arrowhead; and hence four kinds of edges: \rightarrow , \leftrightarrow , $\circ \rightarrow$, $\circ \circ$) such that 1) \mathcal{P} shares the same adjacencies with \mathcal{G} (and any member of $[\mathcal{G}]$), and 2) every non-circle edge endpoint indicates an invariant edge endpoint in $[\mathcal{G}]$.

The second condition in Def. 2.8 implies that an edge associates a tail “ \circ ” or arrowhead “ \rightarrow ” endpoint, if and only if it is invariant in all $\mathcal{G} \in [\mathcal{G}]$. Table 1 lists the semantics of edges.

Example 2.9: Location $\circ \rightarrow$ Smoking in Fig. 1 (c) implies that “Location” is a cause of “Smoking” or they share a latent confounder.

We clarify that the FCI algorithm [51], as a typical constraint-based approach, consists of two phases. The skeleton of $[\mathcal{G}]$ is first learned by assuming faithfulness (i.e., the FCI-SL phase of the FCI algorithm). Then, the undirected edges are subsequently oriented according to a set of orientation rules (i.e., the FCI-Orient phase of the FCI algorithm). Finally, the PAG is returned; see full details of the FCI algorithm in Supplementary Material. However, soon we will show that the faithfulness assumption can be violated by FD relations. In this paper, we focus on establishing a theory and proposing a solution to tackle this unique challenge that arises in data analysis scenarios. That is, our XLEARNER calibrates the FCI algorithm to correctly handling FDs (see details in Sec. 3.1).

3 XINSIGHT

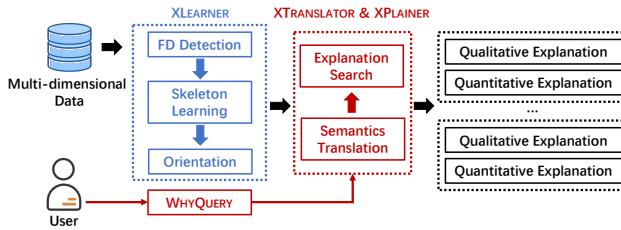


Fig. 3. Workflow of XINSIGHT. Offline phase is marked in blue and online phase is marked in red.

XINSIGHT delivers a unified framework for XDA with three modules. The workflow of XINSIGHT is shown in Fig. 3. First, given a multi-dimensional data D , XLEARNER pre-learns a causal graph \mathcal{G} from data in the offline phase (blue-annotated in Fig. 3). Then, in the online phase (red-annotated in Fig. 3), upon receiving a WHY QUERY, XTRANSLATOR identifies variables that have the potential to give either causal or non-causal explanations based on the causal primitives in \mathcal{G} . Finally, XPLAINER examines each identified variable with potential and decides the optimal explanation for the given WHY QUERY. After applying XPLAINER to all variables with potential, XINSIGHT yields a set of explanations (with qualitative sub-explanations and quantitative sub-explanations). By decoupling XINSIGHT into an offline phase and an online phase, heavy-weight computations are performed beforehand, and only light-weight computations are needed in the online phase, allowing for a

rapid response to a user's query. In the following, we elaborate on the design of each module. Due to page limits, we present proofs and theoretical discussion in the Supplementary Material.

3.1 XLEARNER

XLEARNER aims to learn a causal graph \mathcal{G} from multi-dimensional data D in the presence of latent confounders. The primary obstacle is that learning the skeleton of \mathcal{G} requires the faithfulness assumption (see Sec. 2.2), which may be violated by FDs in D . Below, we show how contradictory causal structures can be induced when being agnostic to FDs.

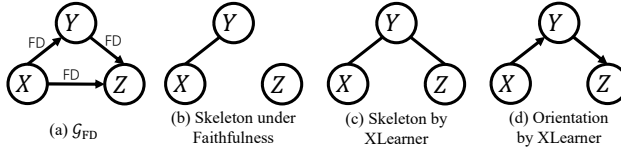


Fig. 4. Illustration of FD-induced faithfulness violation.

Example 3.1: We first consider the **CityInfo** dataset described in Ex. 2.4 and the corresponding FD-induced graph in Fig. 4(a), where X denotes city, Y denotes state, and Z denotes country. By definition of conditional independence (i.e., $X \perp\!\!\!\perp Y \mid Z \iff P(X, Y \mid Z) = P(X \mid Z)P(Y \mid Z)$), we have $Y \perp\!\!\!\perp Z \mid X$ and $X \perp\!\!\!\perp Z \mid Y$. The definition of faithfulness implies $Y \perp\!\!\!\perp_{\mathcal{G}} Z \mid X$ and $X \perp\!\!\!\perp_{\mathcal{G}} Z \mid Y$, given $Y \perp\!\!\!\perp Z \mid X$ and $X \perp\!\!\!\perp Z \mid Y$. Z is non-adjacent to both X and Y according to the m -separation definition (see the induced graph in Fig. 4(b)). Consequently, Z is an isolated node in Fig. 4(b). We have $Y \perp\!\!\!\perp_{\mathcal{G}} Z$ and GMP further implies that $Y \perp\!\!\!\perp Z$, which contradicts $Y \not\perp\!\!\!\perp Z$ entailed by $Y \xrightarrow{\text{FD}} Z$. Indeed, the skeleton is not consistent with any MAGs that are on the top of it.

Table 2. Comparing different causal discovery algorithms. ✓ denotes “support” whereas ✗ denotes “no support”.

Alg.	Orientation	FD-induced Faithfulness Violation	Causal Insufficiency
PC [51]	✓	✗	✗
FCI [57]	✓	✗	✓
REAL [12]	✗	✓	✗
XLEARNER	✓	✓	✓

As aforementioned in Sec. 2, the violations of causal sufficiency and faithfulness (induced by FDs) are common in the data analysis scenarios. However, they are addressed separately in the literature, as reviewed in Table 2. XLEARNER focuses on addressing both challenges simultaneously. Fig. 4(c)-(d) show the skeleton and orientation by XLEARNER, which are compliant with intuition.

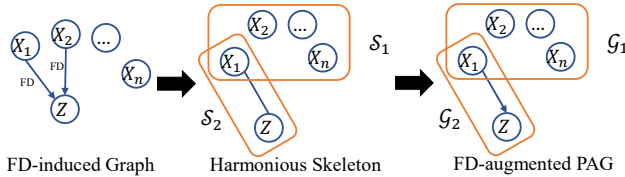


Fig. 5. Running example of XLEARNER.

Overall, XLEARNER tackles the problem in three stages. We outline the workflow of XLEARNER in Alg. 1 and present an example. Then, we elaborate on the design of XLEARNER.

Algorithm 1: XLEARNER procedure.

Input: Multi-dimensional Data D , FD-induced graph \mathcal{G}_{FD}
Output: FD-augmented PAG \mathcal{G}

```

1 // stage 1: detect and preclude  $X_{\text{FD}}$  (Sec. 3.1.1)
2  $\mathcal{S}_2 \leftarrow (V, \emptyset)$ ;
3 Topologically sorting nodes in  $\mathcal{G}_{\text{FD}}$  and record depth as  $d(X_i)$ ;
4 while  $\mathcal{G}_{\text{FD}}$  has non-root nodes do
5    $X \leftarrow \operatorname{argmax}_{X \in \mathcal{G}_{\text{FD}}.V} d(X)$ ;
6    $Y \leftarrow \operatorname{argmin}_{Y \in \text{Pa}(\mathcal{G}_{\text{FD}}, X)} |Y|$ ;
7   add edge  $(X, Y)$  in  $\mathcal{S}_2$ ;
8   remove  $X$  and all connected edges from  $\mathcal{G}_{\text{FD}}$ ;
9 end
10 // stage 2: standard PAG learning
11  $\mathcal{S}_1 \leftarrow \text{FCI-SL}(D, \mathcal{G}_{\text{FD}}.V)$ ;
12  $\mathcal{G}_1 \leftarrow \text{FCI-Orient}(\mathcal{S}_1)$ ;
13 // stage 3: orient  $\mathcal{S}_1$  and generate  $\mathcal{G}$  (Sec. 3.1.2)
14 foreach  $(X \xrightarrow{\text{FD}} Y) \in \mathcal{G}_{\text{FD}}.E$  do
15   if  $X, Y$  is adjacent in  $\mathcal{S}$  then orient  $X \rightarrow Y$  on  $\mathcal{G}^2$ ;
16 end
17 generate  $\mathcal{G}$  concatenating  $\mathcal{G}^1$  and  $\mathcal{G}^2$ ;
18 return  $\mathcal{G}$ ;

```

Example 3.2: Consider the FD-induced graph \mathcal{G}_{FD} shown in Fig. 5. In the first stage, XLEARNER uses \mathcal{G}_{FD} to identify variables (e.g., X_1 and Z in Fig. 5) that may trigger faithfulness violations. Then, the skeleton \mathcal{S}_2 is built upon a harmonious assumption instead of faithfulness over X_1 and Z . In the second stage, the FCI algorithm (skeleton learning and orientation) is only conducted over variables that comply with the faithfulness assumption. Hence, the skeleton \mathcal{S}_1 and the PAG \mathcal{G}_1 are identified accordingly. In the third stage, we orient \mathcal{S}_2 to generate an FD-augmented PAG \mathcal{G}_2 . By concatenating \mathcal{G}_1 and \mathcal{G}_2 , the resultant (FD-augmented) PAG \mathcal{G} is obtained.

Comparison with FCI. Comparing with the FCI algorithm [51, 57], XLEARNER for the first time reconciles functional dependency (FD) and the faithfulness assumption for causally insufficient data within the harmonious skeleton framework. In that sense, it can learn causal graphs from real-world data adequately. As validated in Sec. 4.3, XLEARNER learns more accurate causal graphs than the FCI algorithm. Second, it uses FDs to provide a more complete orientation to the underlying causal graph. Hence, compared to the FCI algorithm, it leverages the knowledge from FDs to enforce a more precise causal graph with less undetermined edges. In sum, we deem that XLEARNER enhances the FCI algorithm from the theoretical perspective, and also addresses obstacles in the real-life adoption of the FCI algorithm.

3.1.1 Skeleton Learning with FD (lines 1–9 of Alg. 1). Ding et al. point out that in the presence of FD relations, faithfulness assumption can be violated thus we can at most obtain a *harmonious skeleton* [12]. However, the original theory of *harmonious skeletons* is established under causal sufficiency. Here, we further generalize the *harmonious skeleton* for causally insufficient systems:

DEFINITION 3.1 (HARMONIOUS SKELETON). A skeleton \mathcal{S} is said to be harmonious w.r.t. a joint probability distribution P if 1) there exists a MAG \mathcal{G} sharing the same adjacencies of \mathcal{S} , 2) P satisfies GMP to \mathcal{G} , and 3) any subgraph of \mathcal{S} does not satisfy the previous two conditions.

Def. 3.1 entails three properties of \mathcal{S} . First, since there exists a MAG \mathcal{G} on top of the skeleton \mathcal{S} , there exists a set of nodes that m-separates any non-adjacent nodes. Second, if two nodes (e.g., X, Y) are m-separated by Z , then $X \perp\!\!\!\perp Y \mid Z$. These two conditions imply that X and Y are non-adjacent in \mathcal{S} , if and only if there exists a set of nodes Z such that $X \perp\!\!\!\perp Y \mid Z$. The last condition implies the minimality of \mathcal{S} , which is commonly assumed [43]. When two graphs $\mathcal{G}, \mathcal{G}'$ are equally compatible with the data, we would prefer the simpler one. We now show the construction of \mathcal{S} , which begins with a basic case and generalizes to arbitrary structures.

THEOREM 3.1. *Let Z be a sink node (i.e., all edges of Z are oriented to Z) in \mathcal{G}_{FD} . $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ is a harmonious skeleton if 1) \mathcal{S}_1 is a harmonious skeleton over $X \setminus Z$ and \mathcal{S}_2 contains only one edge $X_i - Z$ where X_i can be any node connected to Z in \mathcal{G}_{FD} .*

Example 3.3: Fig. 5 presents an example of Thm. 3.1. Connecting the sink node Z with one of its parents X_1 in the \mathcal{G}_{FD} yields a skeleton \mathcal{S}_2 . If we can learn a harmonious skeleton \mathcal{S}_1 on the remaining nodes (X_1, \dots, X_n), Thm. 3.1 ensures that concatenating $\mathcal{S}_1, \mathcal{S}_2$ produces a harmonious skeleton over all variables.

According to Thm. 3.1, if Z has more than one parent, multiple harmonious skeletons exist (note that Z can connect to any one of its parents). In practice, we connect Z to the parent node with the lowest cardinality. Given a FD-induced graph, we recursively apply Thm. 3.1 to identify sink nodes and derive the corresponding harmonious skeleton \mathcal{S}_2 until all FDs are properly resolved. Then, we can apply the standard skeleton learning algorithm over the remaining nodes. The procedure is shown in lines 1–9 in Alg. 1.

THEOREM 3.2. *The skeleton of Alg. 1 is harmonious.*

Alg. 1 first constructs an empty skeleton \mathcal{S} that shares the same nodes as \mathcal{G}_{FD} (line 2). At line 3, we topologically sort the \mathcal{G}_{FD} nodes (note that \mathcal{G}_{FD} is a DAG). In each iteration (lines 5–8), we pick the deepest node and apply Thm. 3.1 to connect X to one of its parents (in \mathcal{G}_{FD}) Y in the skeleton. We use the parent node with the lowest cardinality as Y (line 6), as it usually aligns with human intuition.

Example 3.4: Consider the **CityInfo** dataset in Ex. 2.4. Alg. 1 identifies the correct skeleton as City – State – Country in Fig. 4(c).

For root nodes, since there are no FDs and thus the faithfulness assumption holds, we employ the standard FCI algorithm (lines 10–12) to infer the PAG \mathcal{G}_1 . After \mathcal{S}_2 being oriented to \mathcal{G}_2 (see Sec. 3.1.2), we concatenate them to form \mathcal{G} (line 17). Thm. 3.2 proves that the skeleton of \mathcal{G} is also harmonious after the concatenation.

3.1.2 Orientation (lines 13–16 of Alg. 1). Classical constraint-based causal discovery algorithms decide the direction of edges based on a set of orientation rules. These rules orient undirected edges on skeletons (i.e., $\circ-\circ$) based on a set of criteria, including conditional independence and some graphical structural relations (e.g., discriminating path) [57]. These rules are applied iteratively until no more orientations can be made. However, we argue that an FD itself reflects a causal relation to a good extent, of which the reason is twofold.

ANM Perspective on FD-related Edges. We anticipate incorporating the discrete additive noise model (ANM) [42] for orienting FD-related edges. The main theory of ANM implies that if an asymmetric ANM $Y = f(X) + N_Y$ exists from X to Y and N_Y is independent of X , then X causes Y . By FD, we note that, if $X \xrightarrow{\text{FD}} Y$ in \mathcal{G}_{FD} , an ANM construction from X to Y naturally exist with noise term $N_Y = 0$. On the other hand, an ANM construction from Y to X exists only in very rare cases, as determined by the identifiability of the discrete ANM (see Thm. 4.6 in [43]). In light of this, we hypothesize that $X \xrightarrow{\text{FD}} Y$ in \mathcal{G}_{FD} implies causation of $X \rightarrow Y$.

FCI Perspective on FD-related Edges. The rules in FCI may be unreliable due to the faithfulness violations by FDs. However, an FD itself is generally more reliable, which describes deterministic relations. More importantly, the directions from the FDs are compatible with the result of the FCI on the variables excluding FD-related variables. That is, incorporating ANM would not violate GMP.

We implement the above hypothesis in our orientation algorithm (lines 13–16 in Alg. 1). We examine, for each FD relation that is also adjacent in \mathcal{S}_2 , whether the edge is oriented as \rightarrow (lines 13–15). We note that, by incorporating ANM, the augmented graph is more informative and represents an overcomplete graph w.r.t. the ground-truth MAG’s Markov equivalence class, exhibiting greater precision than causal graphs learned only by rules.

3.2 XTRANSLATOR

A causal graph does not directly reveal if a variable adequately explains a WHY QUERY, nor does it directly reflect if the variable features a causal or non-causal explanation. Bridging this gap requires a translation from causal primitives to XDA semantics. To illustrate, we start with a WHY QUERY under AVG. We then show how to generalize the main result into SUM and other aggregates.

Principle of Explainability. Given a WHY QUERY Δ where $agg = \text{AVG}$, a variable X is said to have *No Explainability* if $X \perp\!\!\!\perp M \mid F \cup B$, where M is the target measure, F is the foreground variable, and B are background variable(s). In the subsequent discussion, we omit B for the ease of presentation without loss of generality.

A WHY QUERY in XDA requires us to observe the difference between aggregates on M within two subspaces. The conditional independence of $X \perp\!\!\!\perp M \mid F$ implies that $\mathbb{E}(M \mid F, X) = \mathbb{E}(M \mid F)$. Hence, $\Delta(D) = \Delta(D_{X=X})$ in the large sample limit for all feasible filters in X . If X is *conditionally independent* of M given F , X is simply impossible to offer explanations to the WHY QUERY. Thus, this principle imposes a restriction on possible variables that have the potential to provide explanations. In particular, we derive the following restriction.

PROPOSITION 3.1. *If X has explainability, M, X are not m -separated by F in the causal graph G .*

Proposition 3.1 illustrates the chance of pruning variables for which it is impossible to provide explanations. Table 3 further depicts the translation from causal primitives to XDA semantics. In XTRANSLATOR, a variable X is first confirmed to have explainability if X, M are not m -separated by F in G (1st row in Table 3). In addition, XTRANSLATOR also categorizes whether X is causal or non-causal according to Table 3. Overall, X provides a causal explanation if it is explainable and a cause (① and ② in Table 3) or a possible cause (③ and ④ rows in Table 3) of M . We show how the causal graph identified by XLEARNER is translated.

Example 3.5: Given the dataset in Fig. 1(a), XLEARNER identifies the corresponding causal graph in Fig. 1(c). With the WHY QUERY in Fig. 1(b), XTRANSLATOR translates the causal graph into the XDA semantics in Fig. 1(d). “Smoking” and “Stress Level” can be used to *causally* explain “Lung Cancer”. And, other variables (e.g., “Surgery”) are deemed non-causal explanations (last row in Table 3).

Table 3. Translating causal primitives to XDA semantics.

Rule	Path	Causal Primitive	XDA Semantics
①	$X \rightarrow F \rightarrow M, \dots$	m -separated	no explainability
②	$X \rightarrow M$	parent	causal explanation
③	$X \rightarrow \dots \rightarrow M$	ancestor	causal explanation
④	$X \circ \rightarrow M$	almost parent	causal explanation
⑤	$X \circ \rightarrow \dots \circ \rightarrow M$	almost ancestor	causal explanation
⑥	others	N/A	non-causal explanation

Extension to SUM. The above formulation over *explainability* is established on AVG aggregates. In the following, we discuss the implications of our formulation on SUM aggregates. If X has no

explainability, $X \perp\!\!\!\perp M \mid F$. When we enforce $X = x$, $\Delta(D_{X=x})$ can merely be affected by the number of rows where $X = x$ in two sibling subspaces (namely a COUNT-based explanation) instead of a causal relation between X and M (see detailed formulation in Supplementary Material). This may be valid for explanations; nevertheless, it is typically inconsistent with the common intuition of data analysis and may not align user expectations regarding explanations (i.e., a variable explains the target). Such COUNT-based explanation is unconventional and is thus less of a concern.

Semantics Consistency. Following the above discussion, we clarify that a variable may play different roles in various aggregates. However, in our current design, XTRANSLATOR focuses primarily on variables with strong connections to M , which are more likely to provide desirable explanations. Therefore, the semantics of a variable are consistent across different aggregates. As clarified in **Principle of Explainability** above, it is appropriate for pruning uninformative variables from general aggregates, and we do not observe notable issues in practice. We leave designing more comprehensive translation rules for future research.

3.3 XPLAINER

XLEARNER and XTRANSLATOR together provide a coarse-grained, variable-level qualitative explanation to a WHY QUERY. For instance, “Smoking” is a causal explanation for the differences in severity of “Lung Cancer” in Locations A and B. To go one step further, XPLAINER provides predicate-level quantitative explanations to answer WHY QUERY (e.g., “Smoking=Yes” explains the difference with the responsibility of 0.77 in Fig. 1(f)). XPLAINER is on the basis of a well-establish framework, DB causality [32] (an extension of actual causality). To ease reading, below we first provide a recap of the notations defined in Sec. 2.1. We then rewrite the formulation of DB causality in the context of XINSIGHT in Sec. 3.3.1.

Recap of Notations. We refer to a dataset as D , a filter as a lowercase p , and a set of filters as an uppercase P . The subset of D satisfying p (or P) is represented by D_p (or D_P , respectively). We use $D - D_p$ as the complement of D_p in D . By default, $\Delta(D)$ represents the WHY QUERY over the dataset D . Likewise, for arbitrary $D' \subseteq D$, $\Delta(D')$ represents the difference between the aggregated values of two sibling subspaces inside D' .

DEFINITION 3.2 (DB CAUSALITY [32]). *Given a multi-dimensional data D and WHY QUERY Δ , let t be a tuple in D . t is called a counterfactual cause to Δ , if $\Delta(D) > \epsilon$ and $\Delta(D - \{t\}) \leq \epsilon$, where ϵ is a user-defined threshold. t is called an actual cause to Δ , if there exists a contingency $\Gamma \subseteq D$ such that t is a counterfactual cause for $D - \Gamma$ (i.e., $\Delta(D - \Gamma - \{t\}) \leq \epsilon < \Delta(D - \Gamma)$).*

DEFINITION 3.3 (DB RESPONSIBILITY [32]). *Suppose P is an actual cause to WHY QUERY Δ and Γ ranges over all valid contingencies for P . The responsibility of P is defined as $\rho_P = \frac{1}{1 + \min_{\Gamma} |\Gamma|}$, where $|\Gamma|$ denotes the number of tuples in the contingency.*

DB causality is appealing as it offers both a normalized measure (responsibility $\in (0, 1]$) and a contingency. First, when the responsibility is close to 1, it implies that the tuple is more accountable for the outcome, and when it hits 1, it is totally responsible. Second, the minimal contingency reflects the additional influential factors that, together with the tuple, are fully responsible for the outcome. The two elements form a quantitative explanation and it is useful for users to understand why the difference exists.

3.3.1 Adaption. DB causality was originally designed for data provenance. As pointed out in [31], tuple-level explanations are usually too fine-grained for data analysis scenarios. An individual tuple usually has too little effect on the highly aggregated outcome of a large dataset. Recalling the example in Fig. 1, users would expect to know that “Smoking=Yes” causes high “Lung Cancer” severity rather than an individual patient being the cause of the high “Lung Cancer” severity. This

necessitates predicate-level explanations which are easier to understand and frequently used in data analysis scenarios, and by many data explanation tools [54, 1]. Motivated by this, we make three adaptations over DB causality (namely, W-CAUSALITY, W-RESPONSIBILITY and conciseness) to support XDA. We formulate W-CAUSALITY as follows.

DEFINITION 3.4 (W-CAUSALITY). *Given a multi-dimensional data D , an attribute of interest X and WHY QUERY Δ , let $P \subseteq \bigcup p_i$ be a predicate in D , where $\bigcup p_i$ denotes the set of all possible filters on X . P is called a counterfactual cause of Δ , if $\Delta(D) > \epsilon$ and $\Delta(D - D_P) \leq \epsilon$, where ϵ is a user-defined threshold. P is deemed to be an actual cause of Δ , if there is a contingency $\Gamma \subseteq \bigcup p_i$ such that P is a counterfactual cause for $D - D_\Gamma$ (i.e., $\Delta(D - D_\Gamma - D_P) \leq \epsilon < \Delta(D - D_\Gamma)$), where $P \cap \Gamma = \emptyset$.*

From Tuples to Predicates. Def. 3.4 transforms the tuple and contingency into two predicates. This way, explanations as well as contingencies constitute a form of intervention over the multi-dimensional data. When a contingency Γ is applied, it indicates that, if the events related to Γ do not happen, then the events related to P are fully responsible for Δ . This adaption in turn entails another adaption to the responsibility for predicate-level explanations.

DEFINITION 3.5 (W-RESPONSIBILITY). *Suppose P is an actual cause to WHY QUERY Δ and Γ range over all valid contingencies for P . The responsibility of P is defined as $\rho_P = \frac{1}{1 + \min_\Gamma |\Gamma|_W}$, where $|\Gamma|_W$ is defined as $\max(\frac{\Delta(D - D_P) - \Delta(D - D_P - D_\Gamma)}{\Delta(D)}, 0)$. We let $\rho_P = 0$ if P is not an actual cause.*

W-RESPONSIBILITY. Instead of using the number of rows in D_Γ as $|\Gamma|_W$, Def. 3.5 employs the truncated difference in Γ over Δ to measure the importance of P . In particular, $\Delta(D - D_P) - \Delta(D - D_P - D_\Gamma)$ can be deemed as *first-order finite backward difference* to the function $\Delta(\cdot)$ at the point of $D - D_P$ and Γ is the step size. This supplies a simple and intuitive way to understand to what extent Γ plays an important role in reducing the difference. The large difference imposed by Γ implies a low importance of explanation P to Δ ; because the reduction in Δ is primarily caused by Γ instead of P itself. Therefore, the responsibility of P is measured by a valid contingency Γ^* (such that $\Delta(D - D_P - D_\Gamma) \leq \epsilon$ and $\Delta(D - D_\Gamma) > \epsilon$) with minimal difference on Δ .

Conciseness. Using responsibility as the sole criterion is not sufficient in practical data analysis scenarios [17]. Typically, a *concise* explanation is preferable. Therefore, given an attribute of interest X , we formulate the optimal explanation of X as follows.

$$\operatorname{argmax}_{P \subseteq \bigcup p_i} \rho_P - \sigma|P| \quad (4)$$

where $\bigcup p_i$ is the set of all possible filters in X , $|P|$ is the number of filters in P , and $\sigma|P|$ ($\sigma > 0$) forms a conciseness regularization. In practice, we would prefer $\sigma = 1/m$ such that when all filters are picked, the score is zero.

Table 4. Different search solutions in XPLAINER. FP is false positive and FN is false negative.

Solution	Complexity	Optimality
Brute-force Search	$O(2^m)$	Optimal
Approx. Search (SUM)	$O(m \log m)$	Moderated FP; Negligible FN
Approx. Search (AVG)	$O(m^2)$	Moderated FP&FN

3.3.2 Optimization. As pointed out in [6, 7], computing responsibility (i.e., ρ_P) is intractable. Furthermore, solving the optimization problem in Eqn. 4 is itself difficult given 2^m search space (m is the number of filters in X). We characterize the performance of different solutions in Table 4. First, the brute-force search is the most accurate and general method for arbitrary aggregates, despite being very slow. The explanation discovered by brute-force search is exactly the optimal explanation. In this paper, we design two approximate solutions for SUM and AVG, respectively. In particular, we first show the existence of a linearithmic approximated solution when the aggregation

is SUM. This solution also has negligible false negatives, as theoretically guaranteed by a lemma on the completeness. Furthermore, we present a heuristics-based solution for AVG with quadratic complexity. Moreover, this solution should be applicable for other aggregate functions with a mild downgrade in optimality. Our evaluation (Sec. 4.4) shows that both approximations are tight and efficient in comparison to brute-force search.

Optimization for SUM. Given the additive property of SUM (i.e., $\Delta(D_{P_1} + D_{P_2}) = \Delta(D_{P_1}) + \Delta(D_{P_2})$), we obtain the following proposition to prune the search space.

PROPOSITION 3.2. *If P^* is the optimal explanation of Eqn. 4, $\forall p \in P^*, \Delta(D_p) > 0$.*

According to Proposition 3.2, the search algorithm can omit filters with a non-positive Δ_i (i.e., $\Delta(D_{p_i})$). Recall that Eqn. 4 seeks the optimal explanation. When the aggregate function is SUM, we only need to focus on filters with a reasonably high Δ_i without losing optimality and we define such filters as *canonical filters*.

DEFINITION 3.6 (CANONICAL FILTER AND PREDICATE). *Without loss of generality (w.l.o.g.), given a WHY QUERY Δ and an attribute of interest X , let filters $\{p_1, \dots, p_m\}$ of X be ordered by Δ_i (i.e., $\Delta(D_{p_i})$) such that $\Delta_1 \geq \dots \geq \Delta_m$. We let p_1, \dots, p_j be canonical filters if*

$$\Delta(D) - \sum_{i=1}^j \Delta_i \leq \epsilon < \Delta(D) - \sum_{i=1}^{j-1} \Delta_i \quad (5)$$

$P^C = \{p_1, \dots, p_j\}$ is called a canonical predicate and $\tau = \sum_{i=1}^j \Delta_i$.

With canonical filters and a corresponding canonical predicate P^C , we observe that P^C manifests good properties. First, P^C is the minimal counterfactual cause entailed by Eqn. 5. Our construction of canonical predicates guarantees completeness.

PROPOSITION 3.3 (COMPLETENESS). *For SUM, given a WHY QUERY Δ , an attribute of interest X and corresponding canonical predicate P^C , there exists an optimal explanation $P^* \subseteq P^C$.*

The completeness proposition (Proposition 3.3) allows us to only focus on canonical filters when searching for the optimal explanation without loss of optimality. More importantly, the canonical predicate also allows us to efficiently identify actual causes and the corresponding valid contingencies.

THEOREM 3.3. *For SUM, given a WHY QUERY Δ , an attribute of interest X and corresponding canonical predicate P^C , $\forall P \subset P^C$, P is an actual cause and $\bar{P} = P^C - P$ is a valid contingency.*

The advantages of Thm. 3.3 are twofold. First, we can directly confirm valid explanations without exhaustive enumerations. Second, by the property of \bar{P} , we bound P 's responsibility (ρ_P).

THEOREM 3.4. *For SUM, given a WHY QUERY Δ , an attribute of interest X and corresponding canonical predicate P^C , the W-RESPONSIBILITY ρ_P of $P \subset P^C$ satisfies*

$$\frac{1}{1 + \frac{\tau - \Delta(D_P)}{\Delta(D)}} \leq \rho_P \leq \frac{1}{2 - \frac{\Delta(D_P) + \epsilon}{\Delta(D)}} \quad (6)$$

When $\Delta(D_P) \ll \tau$ and $0 < \tau \leq \Delta(D)$, $\frac{1}{1 + \tau - \Delta(D_P)} \approx \frac{1 + \tau + \Delta(D_P)}{(1 + \tau)^2}$ and the corresponding approximation error rate $E < 0.25$.

Thm. 3.4 provides a way to efficiently approximate responsibility with theoretical guarantees. In that sense, we can compute responsibility immediately and alleviate searching the minimal

Algorithm 2: XPLAINER For AVG**Input:** WHY QUERY Δ , threshold ϵ , consiseness parameter σ **Output:** (near) optimal explanation P^*

```

1  $P^C \leftarrow \emptyset$ ;
2 foreach  $r = 1, \dots, \min(m, \frac{1}{\sigma})$  do
3   if  $\Delta(D - D_{P^C}) \leq \epsilon$  then break ;
4   else
5      $\bar{P} \leftarrow \{p_1, \dots, p_m\} - P^C$ ;
6     if homogeneous then
7        $S \leftarrow \{p_i \mid p_i \in \bar{P}, \Delta_i > \Delta(D - D_{P^C})\}$ ;
8        $p^* \leftarrow \operatorname{argmin}_{p \in S} \Delta(D - D_{P^C} - D_p)$ ;
9     else
10       $p^* \leftarrow \operatorname{argmin}_{p \in \bar{P}} \Delta(D - D_{P^C} - D_p)$ ;
11    end
12     $P^C \leftarrow P^C \cup \{p^*\}$ ;
13  end
14 end
15 if  $\Delta(D - D_{P^C}) > \epsilon$  then return  $\perp$ ;
16 foreach  $k \in 1, \dots, |P^C|$  do
17    $P_k \leftarrow$  top-k filters of  $P^C$ ;
18    $\Gamma_k \leftarrow P^C - P_k$ ;
19   compute  $\rho_{\hat{P}_k}$  with  $\Gamma_k$ .
20 end
21 return  $\operatorname{argmax}_k \rho_{\hat{P}_k} - \sigma|P_k|$ ;

```

contingency. Let $\hat{\rho}_P = \frac{1+\tau+\Delta(D_P)}{(1+\tau)^2}$, we can rewrite the objective function in the following form.

$$\hat{\rho}_P - \sigma|P| = C_1 + C_2 \times \sum_{p_i \in P} (\Delta_i - C_3) \quad (7)$$

Here, C_1, C_2, C_3 are constants. Then, the optimal explanation to Eqn. 7 is straightforward:

$$P^* = \{p_i \mid \Delta_i > C_3\} \quad (8)$$

where $C_3 = \frac{\sigma\Delta(D)}{(1+\frac{\sigma}{\Delta(D)})^2}$. The complexity is $O(m \log(m))$ (primarily in sorting filters for generating canonical predicates).

Optimization for AVG. In terms of AVG, it is generally much more challenging due to the absence of the additive characteristics on $\Delta(D_P)$. Therefore, the majority of the preceding propositions are not applicable. Having said that, we find the causal graph gives considerable opportunities to prune unnecessary computations.

DEFINITION 3.7 (HOMOGENEOUS SIBLING SUBSPACE). Given sibling subspaces s_1, s_2 (with foreground variable F and background variables B), an attribute X and the causal graph G , s_1, s_2 are homogeneous on X if X, F are m -separated given B on G .

PROPOSITION 3.4. For a homogeneous AVG, given a WHY QUERY Δ , an attribute of interest X , a predicate $P \subseteq \bigcup p_i$ and a filter $p \in P$, if $\Delta(D_p) > \Delta(D_P)$, then $\Delta(D_P - D_p) < \Delta(D_P)$.

To practically address the search problem of AVG (Eqn. 4), we rely on greedy-based heuristics with a pruning strategy enabled by Proposition 3.4. The algorithm is outlined in Alg. 2.

The high-level idea behind Alg. 2 is similar to the one for SUM, which attempts to construct a canonical predicate P^C such that P^C forms a counterfactual cause, each subset $P \subset P^C$ of the canonical predicate constitutes an actual cause, and the complement set $P^C - P$ is a valid contingency. Unlike SUM, however, Alg. 2 does not ensure the optimality of the resulting explanation, primarily due to the incompleteness of the canonical predicate (Proposition 3.3) under AVG. Recall that ρ_P ranges from 0 to 1 in Eqn. 4. The optimal explanation contains at most $1/\sigma$ filters (otherwise, $\rho_P - \sigma|P| < 0$). Hence, the canonical predicate P^C shall contain at most $1/\sigma$ or m (i.e., the number of filters in the attribute).

Alg. 2 employs a greedy strategy to construct P^C progressively. It starts with an empty canonical predicate (line 1) and inserts one filter in each iteration (lines 2–13). Before insertion, it checks whether P^C is a canonical predicate (line 3) and terminates the loop if P^C is already valid. Otherwise, it picks the remaining filters that were not chosen in earlier iterations as candidates (line 5) and inserts the filter that minimizes the difference $\Delta(D - D_{P^C} - D_{p_i})$ at the highest magnitude into P^C (lines 6–12). When homogeneity holds and $\Delta_i \leq \Delta(D - D_{P^C})$, Alg. 2 prunes p_i according to Proposition 3.4 (lines 7–8). Note that Δ_i is invariant throughout the loop; thus it only needs to be queried once. In general cases where homogeneity does not hold, Alg. 2 has to enumerate all possible filters in \bar{P} (line 10). If we cannot obtain a valid canonical predicate (i.e., a counterfactual cause to Δ) after the loop, Alg. 2 terminates and outputs \perp , indicating that it fails to find the optimal explanation within the attribute (line 15). Empirically, we do not observe such rare cases. When the canonical predicate P^C is obtained, $\forall k = 1, \dots, |P^C| - 1$, the top- k filters of $P_k \subseteq P^C$ is a valid actual cause and the complement set $\Gamma_k = P^C - P_k$ is a valid contingency. According to the termination condition in the above loop (line 3), $\Delta(D - D_{P_k}) > \epsilon$. In addition, according to the definition of canonical predicate $\Delta(D - D_{P^C}) = \Delta(D - D_{P_k} - D_{\Gamma_k}) \leq \epsilon$, Γ_k is a valid contingency to P_k . Therefore, we compute the approximated responsibility ρ_{P_k} by using its lower bound deduced by Γ_k (line 19). After enumerating each k , Alg. 2 returns P_k such that $\rho_{P_k} - \sigma|P_k|$ is maximized (line 21). In summary, the first loop (lines 2–14) is of quadratic complexity regardless of homogeneity and the second loop is linear (lines 16–20). The total complexity is $O(m^2)$.

4 EVALUATION

In this section, we evaluate XINSIGHT to answer the following three research questions (RQs):

- (1) **RQ1: End-To-End Performance.** How can XINSIGHT facilitate end users in explainable data analysis?
- (2) **RQ2: XLEARNER Evaluation.** Does XLEARNER effectively recover causal relations from observational data?
- (3) **RQ3: XPLAINER Evaluation.** Does XPLAINER accurately and efficiently yield explanations?¹

4.1 Datasets & User Study Setup

To the best of our knowledge, there is no real-world benchmark with manually labeled query/explanation pairs. To deliver a scientific evaluation, we conduct experiments on ① public datasets collected from previous works, ② real-world data collected from a production environment for user study and human evaluation, and ③ synthetic data with ground-truth explanations. The detailed steps for generating synthetic datasets are given in the Supplementary Material. We make necessary preprocessing before feeding to XINSIGHT (e.g., remove missing values).

① **Flight Delay (FLIGHT).** We use the flight delay dataset from [49] to explore the causes of flight delays in US airports. After preprocessing, the resulting dataset contains 17 variables, including the weather conditions of departure airports (temperature, humidity, visibility, rain, etc.), flight carrier,

¹The correctness of XTRANSLATOR has been rigorously discussed in Sec. 3.2.

flight time (month, quarter, year, day of week and hour) and two variables indicating flight delays, DelayMinute (continuous) and Delay>15min (binary).

① **Hotel Booking (HOTEL).** The hotel booking dataset [3] is frequently used for demonstrating data analysis methods. It contains 119,390 observations from two hotels. Each observation depicts the booking status (e.g., “room type”, “reservation status”, and “is canceled”) of a guest.

② **Web Service Behavior Dataset (WEB).** The dataset is collected from a web service’s production environment. It contains 29 columns and 764 rows, where each row is a list of binary values. The first 28 columns describe user behaviors on the web service (e.g., whether he clicks a specific button), which are collected by a logging module. The last column indicates whether the user was blocked for publishing malicious content (i.e., “IsBlocked”), which is annotated by cybersecurity experts. These behaviors are known to exhibit strong and clear causal relations, making it appropriate for testing XINSIGHT in real-world scenarios.

③ **Synthetic Data A (SYN-A).** Ground-truth causal graphs are unattainable in practice and it is common to generate random graphs and then sample observational data from this graphical model. We generate MAGs with 10 to 150 variables (141 distinct scales in total). For each scale, we synthesize five random graphs and the associated datasets, resulting in 705 (141×5) datasets. Each dataset is injected with different amount of FDs.

③ **Synthetic Data B (SYN-B).** We follow the approach in Scorpion [54] to synthesize datasets for assessing XPLAINER. Each dataset includes a valid WHY QUERY and a ground-truth explanation to this difference. We generate 18 datasets with different difficulties.

User Study Setup. In addition to the experiments that will be launched shortly, we intend to determine the extent to which the results on WEB is correct and reasonable. Nonetheless, rendering professional judgments on explanations and causal claims require sufficient expertise in this domain, which makes gathering a large number of participants difficult. In this study, we recruit six domain experts for the WEB dataset; we confirm each expert can evaluate the explanations and causal claims with professionalism and high confidence. We organize the user study as follows:

- (1) **Participant Education.** We organize an education session for participants and demonstrate how to discern between causation and correlation. Then, a pilot study is conducted to confirm that participants have an adequate sense of causality.
- (2) **Explanation Assessment.** We raise four WHY QUERY and ask XINSIGHT to generate two explanations for each WHY QUERY. We then ask participants to give each explanation a score (between 0 to 5) based on their domain knowledge.
- (3) **Causal Claim Assessment.** Following [23], we collect eight edges connected to “IsBlocked”, transform these causal relations into human-comprehensible causal claims and ask participants to independently evaluate them (by labeling them as “reasonable,” “not reasonable,” or “unsure”).
- (4) **Follow-up Discussion.** Participants explain their decision and discuss the aggregated results.

4.2 RQ1: End-To-End Performance

We show that XINSIGHT generates plausible and intuitive explanations for diverse datasets (FLIGHT, HOTEL and WEB) and invite experts to assess the quality of explanations generated for WEB. In this experiment, we manually discover noticeable differences to form WHY QUERY, and ask XINSIGHT to supply the explanations. We also compare XINSIGHT’s outputs with naive correlation-based explanations. To ease presentation, we describe a WHY QUERY in human-readable natural language in the following paragraphs.

FLIGHT & HOTEL. We ask the following WHY QUERY on ①:

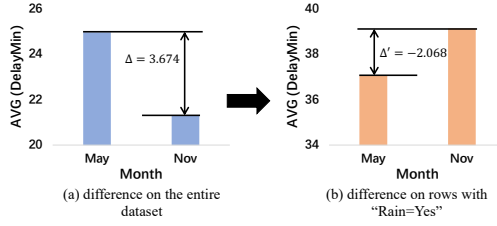


Fig. 6. Explanation of WHY QUERY on the FLIGHT dataset.

- (1) FLIGHT: *why* AVG(DelayMinute) in May (24.95 min) is notably higher than the one in November (21.28 min)?
- (2) HOTEL: *why* AVG(IsCanceled) (cancellation rate) in July (0.37) is notably higher than the one in January (0.30)?

For the first WHY QUERY, we observe that the duration of flight delay differs by month, particularly for May and November, which motivates us to ask XINSIGHT for explanations — what the cause of the flight delay difference is. XINSIGHT first learns a causal graph from data and identifies “rain” as a direct cause of DelayMinute. Then, XPLAINER finds that the difference is reversed ($\Delta = 3.674$ vs. $\Delta' = -2.068$) when the condition “rain=Yes” is enforced (Fig. 6). Thus, it returns “rain=Yes” as an explanation. We interpret the explanation as correct because 1) rain is a typical reason for flight delay, and 2) for most states, monthly precipitation in May is usually higher than in November. Hence, when only counting the rainy cases (by enforcing “rain=Yes”), the difference is eliminated.

For the second WHY QUERY, we observe that the cancellation rate varies by month of arrival. For instance, the cancellation rate in July is 0.37, which is higher than in January. Thus, we ask XINSIGHT for explanations. XINSIGHT identifies “LeadTime” (number of days between booking data and arrival date) as an indirect cause of “IsCanceled”. It also discovers that when enforcing “LeadTime ≤ 133 ”, the difference is reduced. This is intuitive. A longer “LeadTime” results in greater uncertainty about guests’ future schedules, leading to a higher cancellation rate. In January, LeadTime of most reservations is less than 133 days (91%). In contrast, there are far more early reservations (> 133 days) in July (48%), resulting in a higher cancellation rate. When these early reservations are excluded, the difference becomes much smaller.

Table 5. Results of explanation assessment. E_i and P_i stand for the i th explanation and the i th participant, respectively.

	E1	E2	E3	E4	E5	E6	E7	E8
P1	4	4	5	4	4	4	5	3
P2	4	4	4	4	3	4	3	4
P3	5	3	4	5	3	5	5	5
P4	3	4	5	4	4	3	3	4
P5	4	2	5	3	5	4	3	3
P6	5	4	5	5	5	4	5	5
mean	4.16	3.50	4.67	4.17	4.00	4.00	4.00	4.00
std	0.69	0.76	0.47	0.69	0.82	0.58	1.00	0.82

WEB. We report the results of the second phase in the user study (i.e., Explanation Assessment in Sec. 4.1) in Table 5. We view the results as encouraging, since nearly all responses are positive (≥ 3). Moreover, the average scores for seven out of eight explanations are ≥ 4 . We also investigated the explanation with the lowest score (E2 in Table 5). We find this is counter-intuitive but reasonable in retrospect. In the follow-up session, the discussion among participants also confirmed our finding. During the assessments, experts find many explanations inspiring and insightful, despite their familiarity with the dataset. It continuously increases their knowledge and help them design a better criteria for detecting malicious behavior.

Answer to RQ1: *XINSIGHT shows a promising end-to-end performance in explaining data differences. The user study validates that XINSIGHT achieves a respectable level of agreement with experts.*

4.3 RQ2: XLEARNER Evaluation

As a cornerstone of XINSIGHT, XLEARNER is crucial to the effectiveness of the entire pipeline. To answer **RQ2**, we run XLEARNER on **SYN-A** which has ground truth causal graphs and on a real-world dataset **WEB**. Since **WEB** does not associate a ground-truth causal graph, we assess the quality of causal relations by the user study.

Table 6. Overall comparison between XLEARNER and FCI.

Algo.	F1-Score	Precision	Recall
XLEARNER	0.88 ± 0.04	0.95 ± 0.03	0.82 ± 0.06
FCI	0.72 ± 0.05	0.92 ± 0.04	0.59 ± 0.06

Table 6 provides an overall comparison between XLEARNER and FCI on **SYN-A** datasets. We find that XLEARNER is more accurate than FCI in the presence of FDs. In particular, while FCI has comparable precision, XLEARNER has a much higher recall. This confirms our discussion on the implications of FDs in Sec. 3.1. The faithfulness violations mislead FCI to incorrectly refute true edges (thus yield a lower recall) while XLEARNER is aware of such faithfulness violations and handles them with the procedure in Sec. 3.1.

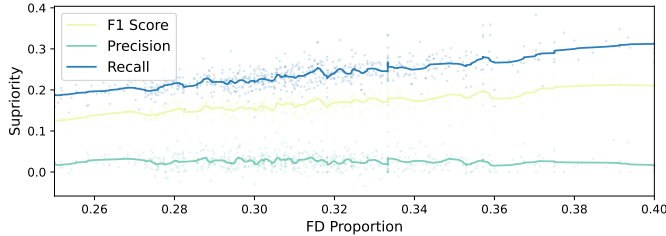


Fig. 7. Comparison by FD Proportion. The x-axis is the proportion of FDs in the causal graph. The y-axis is the superiority (determined by subtracting the FCI’s score from the XLEARNER’s score) of XLEARNER over FCI.

Since XLEARNER focuses primarily on FDs as the opposite of FCI, we further study how varying proportions of FDs in the causal graph affect XLEARNER’s performance. We report the superiority in terms of varied amounts of FDs in Fig. 7. Overall, we observe an increasing trend in XLEARNER performance (particularly for F1 and recall) as the FD proportion increases. More importantly, we observe that “superiority” increases as FDs increase. Recall, as noted in the caption of Fig. 7, that the superiority is computed by subtracting the FCI’s score from the XLEARNER’s score. Thus, we interpret that XLEARNER gradually outperforms the FCI algorithm with greater degree as the proportion of FDs grows.

In addition to the experiments on synthetic datasets, we also evaluate XLEARNER with the **WEB** dataset (Sec. 4.1). As aforementioned, this real-world dataset lacks a ground-truth causal graph. Evaluating the accuracy of an estimated causal graph is thus challenging, if not impossible. At this step, we involve human experts to assess the correctness of the identified causal relations in the third phase of our user study (i.e., Causal Claim Assessment in Sec. 4.1).

We report the results of our user study in Table 7: first, out of 48 responses (6 participants \times 8 questions), only three (6.3%) suggest that the causal claims are “Not Reasonable,” while 40 responses (83.3%) mark the causal claims as “Reasonable.” It indicates that the causal relations identified by XLEARNER correspond with expert knowledge in the majority of instances. Second, we investigate the claims

Table 7. User study. C_i stands for the i th causal claim.

	C1	C2	C3	C4	C5	C6	C7	C8
# Reasonable	6	4	4	6	6	4	5	5
# Not Sure	0	2	1	0	0	0	1	1
# Not Reasonable	0	0	1	0	0	2	0	0

that have been deemed “Not Reasonable” or “Unsure.” Encouragingly, we find that a notable proportion of causal claims are counter-intuitive yet *correct*. For instance, one causal claim states that “malicious intent would lead to more frequent configuration changes than benign intent.” In the causal claim assessment phase, one expert deemed it “Not Reasonable” and presumed that malicious users would keep a default configuration. In the follow-up session where they shared the independent assessments, this participant was persuaded and confirmed this causal relation as “Reasonable.”

Answer to RQ2: *As reflected by the carefully-designed quantitative experiments and the user study, XLEARNER generates plausible causal graphs that are consistent with expert knowledge.*

4.4 RQ3: XPLAINER Evaluation

Recall the descriptions of **SYN-B** in which different parameters result in datasets with varying degrees of difficulty. In this experiment, we explore the accuracy of XPLAINER on **SYN-B**.

Baseline. We compare XPLAINER with three baselines, namely, Scorpion [54], RSExplain [47] and BOExplain [27], which use predicates as explanations. Scorpion is an explanation engine for explaining outliers. It uses a metric called influence score to rank explanations, which considers the effect of the explanation between the outlier region and the hold-out region. RSExplain uses the concept of intervention to measure the effectiveness of an explanation. BOExplain is originally designed for explaining black-box machine learning models. When explaining WHY QUERY, it employs the inference score and the Bayesian optimization to find the optimal explanation. To launch an apple-to-apple comparison, all baselines are enforced to search over a set of pre-defined causal filters $\{p_1, \dots, p_m\}$ derived from the generation procedure of **SYN-B** (see details in Supplementary Material); all these filters have been confirmed to constitute legitimate causal explanations.

Metric. We use the top-ranked explanation of each baseline as its optimal explanation. We mark a method as “N/A” denoting timeouts (more than one hour to process). We report the F1 Score of filters in the explanation over the ground-truth explanation.

Different Dataset Sizes. To study the scalability of XPLAINER, we generate datasets of varying sizes and report the results in Table 8. Overall, we observe that XPLAINER is more accurate and efficient than all baselines across all the studied settings. This is encouraging and also reasonable, as XPLAINER uses many distinct characteristics of aggregation functions to optimize the search process, while other methods primarily treat them as a “black-box.” Scorpion and BOExplain often produce incomplete explanations, while RSExplain may frequently find extra spurious filters. We presume that this is because the objective function of Scorpion (and also BOExplain) is for explaining anomalies instead of WHY QUERY, whereas RSExplain is primarily designed for data provenance. In contrast, explanations provided by XPLAINER are seen as consistent with the ground truth.

XPLAINER is highly efficient particularly for high cardinality regimes, while both Scorpion and RSExplain run out of time when the cardinality exceeds 30 (see the bottom half of Table 8). BOExplain uses Bayesian optimization to search for explanations, and its accuracy downgrades as cardinality increases. Similarly, when iterating different #Rows (the top half of Table 8), XPLAINER also exhibits highly encouraging efficiency: XPLAINER takes on average 0.06 seconds to explain WHY QUERY whereas BOExplain (the second best) takes 13.17 seconds. In sum, we interpret from

Table 8. XPLAINER and baselines under various settings. ✓ denotes that F1=1.0 and the best metric is **high-lighted**.

#Rows (Cardinality=10)		10K	20K	50K	100K	500K	1M
XPLAINER (SUM)	F1 Score	✓	✓	✓	✓	✓	✓
	Time (sec.)	0.004	0.005	0.007	0.010	0.017	0.019
Scorpion (SUM)	F1 Score	0.5	0.5	0.5	0.5	0.5	0.8
	Time (sec.)	0.68	0.82	1.25	1.93	2.45	2.93
RSExplain (SUM)	F1 Score	0.75	0.75	0.75	0.75	0.75	0.75
	Time (sec.)	0.68	0.83	1.25	1.94	2.44	2.90
BOExplain (SUM)	F1 Score	0.8	0.8	0.5	0.5	0.5	0.8
	Time (sec.)	5.24	5.32	5.62	6.38	9.80	13.53
XPLAINER (AVG)	F1 Score	✓	✓	✓	✓	✓	✓
	Time (sec.)	0.016	0.019	0.026	0.038	0.052	0.063
Scorpion (AVG)	F1 Score	✓	✓	✓	✓	✓	✓
	Time (sec.)	0.59	0.67	0.90	1.29	1.69	2.01
RSExplain (AVG)	F1 Score	0.75	0.75	0.75	0.75	0.75	0.75
	Time (sec.)	0.58	0.66	0.90	1.28	1.68	1.95
BOExplain (AVG)	F1 Score	0.86	✓	0.86	✓	✓	0.8
	Time (sec.)	5.33	5.37	5.56	6.56	8.67	12.62
Cardinality (#Rows=100k)		10	15	20	30	50	100
XPLAINER (SUM)	F1 Score	✓	✓	✓	✓	✓	✓
	Time (sec.)	0.010	0.014	0.018	0.025	0.040	0.077
Scorpion (SUM)	F1 Score	0.5	0.5	0.5	N/A	N/A	N/A
	Time (sec.)	1.96	16.50	75.72	N/A	N/A	N/A
RSExplain (SUM)	F1 Score	0.75	0.75	0.75	N/A	N/A	N/A
	Time (sec.)	1.95	16.61	75.82	N/A	N/A	N/A
BOExplain (SUM)	F1 Score	✓	0.86	0.86	0.46	0.27	0.15
	Time (sec.)	6.28	8.71	11.17	15.44	25.44	48.73
XPLAINER (AVG)	F1 Score	✓	✓	✓	✓	✓	✓
	Time (sec.)	0.038	0.060	0.082	0.124	0.211	0.426
Scorpion (AVG)	F1 Score	✓	✓	✓	N/A	N/A	N/A
	Time (sec.)	1.27	10.58	47.90	N/A	N/A	N/A
RSExplain (AVG)	F1 Score	0.75	0.75	0.75	N/A	N/A	N/A
	Time (sec.)	1.28	10.59	47.91	N/A	N/A	N/A
BOExplain (AVG)	F1 Score	✓	0.86	0.5	0.5	0.27	0.14
	Time (sec.)	5.87	8.23	10.44	15.00	24.35	46.35

Table 8 that XPLAINER delivers highly encouraging accuracy and efficiency across different settings in comparison with the baseline methods.

Different $\mu^* - \mu$. The difference between μ^* , μ indicates the magnitude of Δ . To study the sensitivity of XPLAINER, we study how well XPLAINER performs under varying differences and compare it with baselines in Table 9. To clarify, $\mu^* - \mu = 5$ and $\mu^* - \mu = 10$ denote two relatively more challenging settings in Table 9, given the very subtle differences. On SUM aggregates, we find that all methods have difficulties in identifying explanations in those two challenging settings; still, XPLAINER yields the best results for both settings. Even on the most challenging setting ($\mu^* - \mu = 5$), XPLAINER finds highly accurate explanations whereas RSExplain is less accurate.

On AVG aggregates, XPLAINER and Scorpion both perform well on identifying the ground-truth explanations; XPLAINER is slightly better particularly for the most challenging setting when $\mu^* - \mu = 5$. Nevertheless, RSExplain and BOExplain are less accurate on AVG. Overall, we conclude that XPLAINER is more robust to subtle data differences while all other methods have difficulties in such challenging settings. We omit reporting the processing time here, since it has already been evidently explored in Table 8.

Tightness of Approximation. In Sec. 3.3, we show the approximation of minimal contingency for computing responsibilities under SUM and AVG. In the following, we compare the tightness of the responsibility $\hat{\rho}$ computed by $\bar{P} = P^C - P$ to the true responsibility ρ computed by the minimal contingency P_{\min} via brute-force search. The approximation error is computed as $E = \frac{|\hat{\rho} - \rho|}{\rho}$. Recall

Table 9. XPLAINER and baselines with different $\mu^* - \mu$. ✓ denotes that the result is identical to the ground truth (F1=1.0).

$\mu - \mu^*$	5	10	15	30	50	100
XPLAINER (SUM)	0.86	✓	✓	✓	✓	✓
Scorpion (SUM)	0.50	0.50	0.50	0.50	0.50	0.50
RSExplain (SUM)	0.75	0.75	0.75	0.75	0.75	0.75
BOExplain (SUM)	0.50	0.86	0.80	0.80	0.80	✓
XPLAINER (AVG)	✓	✓	✓	✓	✓	✓
Scorpion (AVG)	0.80	✓	✓	✓	✓	✓
RSExplain (AVG)	0.75	0.75	0.75	0.75	0.75	0.75
BOExplain (AVG)	0.80	✓	0.86	0.86	0.80	✓

that we craft three filters that form the counterfactual cause in each dataset. In SUM, we can craft six ($\binom{3}{2}$) actual causes from the three filters. In AVG, since the canonical predicate of AVG only supports the first k filters as actual causes and the rest as contingencies, we pick the top-1 and top-2 filters as two actual causes and repeat the experiments on three datasets (2×3 in total). On the six actual causes of SUM aggregates, we find that the brute-force algorithm is $253.3\times$ slower than our approximated solution. More importantly, the approximation error is highly negligible with an average of 0.007. We also observe that the approximation error on AVG is slightly higher (0.066) and that our heuristics-based solution is $27.3\times$ faster. This result is reasonable, as the heuristics-based solution does not provide guarantees of accuracy and requires more queries than SUM.

Answer to RQ3: *XPLAINER shows high scalability to large datasets and also accurately generates explanations in very difficult settings. On a mild cost of precision, two approximation solutions of XPLAINER substantially improve efficiency.*

5 DISCUSSION

FD in Noisy Data. XINSIGHT only considers *deterministic* FDs. As illustrated in Ex. 3.1, taking deterministic FDs into account eliminates faithfulness violations. However, when the data is noisy, the FDs may be stochastic (e.g., probabilistic interpretation of FDs [58]), which is currently not considered in XINSIGHT. We clarify that considering only deterministic FDs deems a common setup shared by relevant works in this field [12, 30]. It remains unclear how noisy FDs may impact faithfulness. We leave this for future exploration.

Acquiring Causal Knowledge. Inferring causal relations is difficult. Typically, it needs a combination of domain knowledge [2], randomized experiments [52] and causal discovery [51, 57, 12, 29, 10]. XINSIGHT performs causal discovery from observational data due to its simplicity. Nevertheless, we envision users of XINSIGHT can combine additional sources for acquiring more accurate causal knowledge. In this paper, we explain several key obstacles of applying causal discovery to real-world data, including causal insufficiency [57] and FD-induced faithfulness violations [12, 30]. XLEARNER, for the first time, simultaneously addresses all of them.

Other Forms of Explanations. Currently, XINSIGHT employs predicates as the content of explanations, which is general enough for common data analysis scenarios. However, in some cases, explanations may be formed by the number of records in a database [13] or counterbalances [35]. Furthermore, when explaining changes in a whole time series [9], XINSIGHT may be not applicable. We leave integrating XINSIGHT into these scenarios for future work.

6 RELATED WORK

Data Explanation. Explaining an unexpected query outcome in database is a crucial phase in the lifecycle of data analysis. In general, an explanation aims to provide certain forms of patterns

that lead to the unexpected query outcome. Such patterns may be a set of predicates [54, 1, 4, 47], tuples [31], or counterbalances [35]. Scorpion is the most relevant work for XINSIGHT, which also provides explanations to aggregated queries [54]. In particular, it employs an influence score to quantify explanations and features a set of optimizations to reduce the cost of explanation search. Recently, many tools have attempted to enhance explanations with additional knowledge (e.g., join tables) about the underlying data. However, such additional knowledge does not imply causation — top ranked explanations could be rated low by human participants due to a lack of causal semantics [24]. These observations evidently show the necessity of integrating causality into XDA.

Causality in Database. Most works related to causality analysis in the database is on the basis of Halpern’s seminal framework on actual causality [19, 18]. It provides an elegant and natural way to reason about input-output relations. Its results not only highlight the output’s cause, but also provide a contingency describing how it is triggered. The adaption of actual causality in the database (i.e., DB causality) is widely used for data provenance [33], data explanation [47] and debugging [34, 14, 56, 20]. However, it has limitations when applied alone. On the one hand, as noted in [17], DB causality does not necessarily imply *true causation*. Indeed, it assumes that causal knowledge is already known, and focuses solely on quantitative explanations. On the other hand, considerable adaptations are required to make it applicable to XDA scenarios, as discussed in Sec. 3.3. We also notice other methods for quantifying explanations, such as sufficient score, necessity score, and average causal effect [48, 53]. Despite their usefulness, we design XPLAINER on top of actual causality because it is more understandable and general. Furthermore, without prior causal knowledge, none of these methods can imply true causation.

XDA vs. XAI. We note that XAI (explainable artificial intelligence) is parallel and complementary to XDA. Through the lens of data analysis, XAI aims to explain a prediction or model [46, 16], while XDA enhances EDA for understanding data facts. In addition, we also observe a line of research [55, 15, 26, 25] that identifies a subset of (training) data that is responsible for a prediction. While this line of research shares similar output format with XINSIGHT, it is essentially for explaining how model predictions are influenced by training/test data, a scenario that is orthogonal to our research.

7 CONCLUSION

This paper advocates XDA, a concept that ships comprehensive and in-depth explainability toward EDA. XDA offers either causal or non-causal explanations for EDA outcomes, from both quantitative and qualitative perspectives. We have also presented the design of XINSIGHT, a production framework for XDA over databases. Experiments and human evaluations reveal that XINSIGHT manifests highly encouraging explanation capabilities. XPLAINER has been integrated into Microsoft Power BI to explain increase/decrease in data.

ACKNOWLEDGMENTS

The authors would like to thank the anonymous reviewers for their valuable comments and suggestions. The author would also like to thank Siwen Zhu, Haidong Zhang, Zhitao Hou, Ruming Wang, Ziyu Wang, Long Ding, Kai Zhang, Jon Kay, and Dingkun Xie for helpful discussions and all our participants in the user study for their valuable feedback. The authors at HKUST were supported in part by RGC RMGS under the contract RMGS22EG02.

REFERENCES

- [1] Firas Abuzaid et al. “Diff: a relational interface for large-scale data explanation”. In: *The VLDB Journal* 30.1 (2021), pp. 45–70.

- [2] Bryan Andrews, Peter Spirtes, and Gregory F Cooper. “On the completeness of causal discovery in the presence of latent confounding with tiered background knowledge”. In: *International Conference on Artificial Intelligence and Statistics*. PMLR. 2020, pp. 4002–4011.
- [3] Nuno Antonio, Ana de Almeida, and Luis Nunes. “Hotel booking demand datasets”. In: *Data in brief* 22 (2019), pp. 41–49.
- [4] Peter Bailis et al. “Macrobase: Prioritizing attention in fast data”. In: *Proceedings of the 2017 ACM International Conference on Management of Data*. 2017, pp. 541–556.
- [5] Elias Bareinboim and Judea Pearl. “Controlling selection bias in causal inference”. In: *Artificial Intelligence and Statistics*. PMLR. 2012, pp. 100–108.
- [6] Leopoldo Bertossi. “Score-Based Explanations in Data Management and Machine Learning”. In: *International Conference on Scalable Uncertainty Management*. Springer. 2020, pp. 17–31.
- [7] Leopoldo Bertossi et al. “Causality-based explanation of classification outcomes”. In: *Proceedings of the Fourth International Workshop on Data Management for End-to-End Machine Learning*. 2020, pp. 1–10.
- [8] Michael Buckland. *Information and society*. MIT Press, 2017.
- [9] Yiru Chen and Silu Huang. “TSExplain: Surfacing Evolving Explanations for Time Series”. In: *Proceedings of the 2021 International Conference on Management of Data*. 2021, pp. 2686–2690.
- [10] Haoyue Dai et al. “ML4C: Seeing Causality Through Latent Vicinity”. In: *arXiv preprint arXiv:2110.00637* (2021).
- [11] Rui Ding et al. “Quickinsights: Quick and automatic discovery of insights from multi-dimensional data”. In: *Proceedings of the 2019 International Conference on Management of Data*. 2019, pp. 317–332.
- [12] Rui Ding et al. “Reliable and Efficient Anytime Skeleton Learning”. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 34. 06. 2020, pp. 10101–10109.
- [13] *Discover Insights Faster with Explain Data*. https://help.tableau.com/current/pro/desktop/en-us/explain_data.htm. 2022.
- [14] Anna Fariha, Suman Nath, and Alexandra Meliou. “Causality-guided adaptive interventional debugging”. In: *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020, pp. 431–446.
- [15] Lampros Flokas et al. “Complaint-Driven Training Data Debugging at Interactive Speeds”. In: *Proceedings of the 2022 International Conference on Management of Data*. 2022, pp. 369–383.
- [16] Sainyam Galhotra, Romila Pradhan, and Babak Salimi. “Explaining black-box algorithms using probabilistic contrastive counterfactuals”. In: *Proceedings of the 2021 International Conference on Management of Data*. 2021, pp. 577–590.
- [17] Boris Glavic, Alexandra Meliou, Sudeepa Roy, et al. “Trends in Explanations: Understanding and Debugging Data-driven Systems”. In: *Foundations and Trends® in Databases* 11.3 (2021), pp. 226–318.
- [18] Joseph Y Halpern. *Actual causality*. MIT Press, 2016.
- [19] Joseph Y Halpern and Judea Pearl. “Causes and explanations: A structural-model approach. Part I: Causes”. In: *The British journal for the philosophy of science* 56.4 (2005), pp. 843–887.
- [20] Zhenlan Ji, Pingchuan Ma, and Shuai Wang. “PerfCE: Performance Debugging on Databases with Chaos Engineering-Enhanced Causality Analysis”. In: *arXiv preprint arXiv:2207.08369* (2022).
- [21] Frank C Keil. “Explanation and understanding”. In: *Annu. Rev. Psychol.* 57 (2006), pp. 227–254.
- [22] Marc Lange. *Because Without Cause: Non-Casual Explanations In Science and Mathematics*. Oxford University Press, 2016.
- [23] Po-Ming Law et al. “Causal Perception in Question-Answering Systems”. In: *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 2021, pp. 1–15.

- [24] Chenjie Li et al. “Putting Things into Context: Rich Explanations for Query Answers using Join Graphs (extended version)”. In: *arXiv preprint arXiv:2103.15797* (2021).
- [25] Yanhui Li et al. “Training data debugging for the fairness of machine learning software”. In: *Proceedings of the 44th International Conference on Software Engineering*. 2022, pp. 2215–2227.
- [26] Jinkun Lin et al. “Measuring the Effect of Training Data on Deep Learning Predictions via Randomized Experiments”. In: *International Conference on Machine Learning*. PMLR. 2022, pp. 13468–13504.
- [27] Brandon Lockhart et al. “Explaining inference queries with bayesian optimization”. In: *Proceedings of the VLDB Endowment* 14.11 (2021), pp. 2576–2585.
- [28] Pingchuan Ma et al. “MetaInsight: Automatic Discovery of Structured Knowledge for Exploratory Data Analysis”. In: *Proceedings of the 2021 International Conference on Management of Data*. 2021, pp. 1262–1274.
- [29] Pingchuan Ma et al. “ML4S: Learning Causal Skeleton from Vicinal Graphs”. In: *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2022.
- [30] Ahmed Mabrouk et al. “An efficient Bayesian network structure learning algorithm in the presence of deterministic relations”. In: *ECAI 2014*. IOS Press, 2014, pp. 567–572.
- [31] Alexandra Meliou, Sudeepa Roy, and Dan Suciu. “Causality and explanations in databases”. In: *Proceedings of the VLDB Endowment* 7.13 (2014), pp. 1715–1716.
- [32] Alexandra Meliou et al. “Causality in databases”. In: *IEEE Data Engineering Bulletin* 33.ARTICLE (2010), pp. 59–67.
- [33] Alexandra Meliou et al. “The complexity of causality and responsibility for query answers and non-answers”. In: *Proceedings of the VLDB Endowment* 4.1 (2010), pp. 34–45.
- [34] Alexandra Meliou et al. “Tracing data errors with view-conditioned causality”. In: *Proceedings of the 2011 ACM SIGMOD International Conference on Management of data*. 2011, pp. 505–516.
- [35] Zhengjie Miao et al. “Going beyond provenance: Explaining query answers with pattern-based counterbalances”. In: *Proceedings of the 2019 International Conference on Management of Data*. 2019, pp. 485–502.
- [36] Microsoft. *Microsoft/reliableAI*. <https://github.com/microsoft/reliableAI>. 2022.
- [37] Microsoft. *Use the Analyze feature to explain fluctuations in report visuals*. <https://learn.microsoft.com/en-us/power-bi/consumer/end-user-analyze-visuals>. 2022.
- [38] Tova Milo and Amit Somech. “Automating exploratory data analysis via machine learning: An overview”. In: *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020, pp. 2617–2622.
- [39] Gregory L Murphy and Douglas L Medin. “The role of theories in conceptual coherence.” In: *Psychological review* 92.3 (1985), p. 289.
- [40] Judea Pearl. “Causal inference in statistics: An overview”. In: *Statistics Surveys* 3.none (2009), pp. 96–146.
- [41] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [42] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. “Causal inference on discrete data using additive noise models”. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 33.12 (2011), pp. 2436–2450.
- [43] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. The MIT Press, 2017.
- [44] David L Poole and Alan K Mackworth. *Artificial Intelligence: foundations of computational agents*. Cambridge University Press, 2010.
- [45] Mark Povich and Carl F Craver. *Because without Cause: Non-Causal Explanations in Science and Mathematics*. 2018.

- [46] Romila Pradhan et al. “Explainable AI: Foundations, Applications, Opportunities for Data Management Research”. In: *Proceedings of the 2022 International Conference on Management of Data*. 2022, pp. 2452–2457.
- [47] Sudeepa Roy and Dan Suciu. “A formal approach to finding explanations for database queries”. In: *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*. 2014, pp. 1579–1590.
- [48] Babak Salimi et al. “Causal relational learning”. In: *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020, pp. 241–256.
- [49] Babak Salimi et al. “Zaliql: causal inference from observational data at scale”. In: *Proceedings of the VLDB Endowment* 10.12 (2017), pp. 1957–1960.
- [50] Richard Scheines, Peter Spirtes, and Clark Glymour. “A qualitative approach to causal modeling”. In: *Qualitative simulation modeling and analysis*. Springer, 1991, pp. 72–97.
- [51] Peter Spirtes et al. *Causation, prediction, and search*. MIT press, 2000.
- [52] Sofia Triantafillou and Ioannis Tsamardinos. “Constraint-based causal discovery from multiple interventions over overlapping variable sets”. In: *The Journal of Machine Learning Research* 16.1 (2015), pp. 2147–2205.
- [53] David S Watson et al. “Local explanations via necessity and sufficiency: Unifying theory and practice”. In: *Uncertainty in Artificial Intelligence*. PMLR. 2021, pp. 1382–1392.
- [54] Eugene Wu and Samuel Madden. “Scorpion: Explaining Away Outliers in Aggregate Queries”. In: *Proceedings of the VLDB Endowment* 6.8 (2013).
- [55] Weiyuan Wu et al. “Complaint-driven training data debugging for query 2.0”. In: *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020, pp. 1317–1334.
- [56] Dong Young Yoon, Ning Niu, and Barzan Mozafari. “Dbsherlock: A performance diagnostic tool for transactional databases”. In: *Proceedings of the 2016 International Conference on Management of Data*. 2016, pp. 1599–1614.
- [57] Jiji Zhang. “On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias”. In: *Artificial Intelligence* 172.16–17 (2008), pp. 1873–1896.
- [58] Yunjia Zhang, Zhihan Guo, and Theodoros Rekatsinas. “A statistical perspective on discovering functional dependencies in noisy data”. In: *Proceedings of the 2020 ACM SIGMOD International Conference on Management of Data*. 2020, pp. 861–876.

8 SUPPLEMENTARY MATERIAL

8.1 Fast Causal Inference (FCI) Algorithm [51, 57]

We first introduce concepts and notations in addition to the one described in Sec. 2.2. Then, we outline the FCI algorithm in Alg. 3 and Alg. 4.

DEFINITION 8.1 (UNSHIELDED TRIPLE). *In a graph, a triple (X, Y, Z) is an unshielded triple if X and Z are non-adjacent, X and Z are adjacent, and Y and Z are adjacent.*

DEFINITION 8.2 (POSSIBLE-D-SEP). *In a graph, $\text{Possible-D-SEP}(X, Y)$ is the set of nodes Z such that there is an undirected path \mathcal{P} between X and Z and for each subpath $S \ast \ast W \ast \ast T$ of \mathcal{P} one of the following conditions holds.²*

- (1) W is a collider; or,
- (2) W is not marked as a non-collider and S, W, T are a triangle. (A triangle is a set of three nodes all adjacent to one another).

DEFINITION 8.3 (EXT-D-SEP). *In a graph, $\text{Ext-D-SEP}(X, Y)$ is the union of $\text{Possible-D-SEP}(X, Y)$ and $\text{Possible-D-SEP}(Y, X)$.*

DEFINITION 8.4 (DISCRIMINATING PATH). *In a MAG, a path between X and Y , $\mathcal{P} = (X, \dots, W, V, Y)$, is a discriminating path for V if*

- (1) \mathcal{P} includes at least three edges;
- (2) V is a non-endpoint node on \mathcal{P} , and is adjacent to Y on \mathcal{P} ; and
- (3) X is not adjacent to Y , and every node between X and V is a collider on \mathcal{P} and is a parent of Y .

DEFINITION 8.5 (UNCOVERED PATH). *In a PMG (partial mixed graph), a path $\mathcal{P} = (V_0, \dots, V_n)$ is said to be uncovered if for every $1 \leq i \leq n-1$, V_{i-1} and V_{i+1} are not adjacent, i.e., if every consecutive triple on the path is unshielded.*

DEFINITION 8.6 (POTENTIALLY DIRECTED PATH). *In a PMG (partial mixed graph), a path $\mathcal{P} = (V_0, \dots, V_n)$ is said to be potentially directed (abbreviated as p.d.) from V_0 to V_n if for every $0 \leq i \leq n-1$, the edge between V_i and V_{i+1} is not into V_i or out of V_{i+1} .*

DEFINITION 8.7 (CIRCLE PATH). *In a PMG (partial mixed graph), a circle path $\mathcal{P} = (V_0, \dots, V_n)$ is a special case of p.d. path, where each edge on \mathcal{P} is $\circ - \circ$.*

8.2 Principle of Explainability under AVG and SUM

In AVG, we notice that $\text{AVG}_M(s_1) \approx \mathbb{E}(M \mid F = f_1)$ and $\text{AVG}_M(s_2) \approx \mathbb{E}(M \mid F = f_2)$ asymptotically. $\text{AVG}_M(s_1 \wedge X = x) \approx \mathbb{E}(M \mid F = f_1, X = x) = \mathbb{E}(M \mid F = f_1) = \text{AVG}_M(s_1)$. Similarly, $\text{AVG}_M(s_2 \wedge X = x) \approx \text{AVG}_M(s_1)$. Hence, $\Delta(D) \approx \Delta(D_{X=x})$.

In SUM, we notice that $\text{SUM} = \text{COUNT} \times \text{AVG}$ and $\text{COUNT}_M(s_1) = N \times P(F = f_1)$, $\text{COUNT}_M(s_2) = N \times P(F = f_2)$, where N is a constant indicating the number of rows in the data. Then, $\Delta = N(P(F = f_1)\mathbb{E}(M \mid F = f_1) - P(F = f_2)\mathbb{E}(M \mid F = f_2))$. $\Delta(D_{X=x}) = N(P(F = f_1, X = x)\mathbb{E}(M \mid F = f_1, X = x) - P(F = f_2, X = x)\mathbb{E}(M \mid F = f_2, X = x)) = N(P(F = f_1, X = x)\mathbb{E}(M \mid F = f_1) - P(F = f_2, X = x)\mathbb{E}(M \mid F = f_2))$.

8.3 Proof of Thm. 3.1

LEMMA 8.3.1. *If $X \xrightarrow{\text{FD}} Y$, then $Y \not\perp\!\!\!\perp X$, and for any other variable set Z , $Z \perp\!\!\!\perp Y \mid X$.*

²* denotes a wildcard endpoint; either $-$ (tail), \rightarrow (arrowhead) or $\circ -$ (circle) can be matched.

Algorithm 3: FCI-SL**Input:** Data D **Output:** Skeleton G

```

1 initialize a complete undirected graph  $Q$  with  $n$  nodes;
2  $n \leftarrow 0$ ;
3 repeat
4   repeat
5     select an ordered pair of adjacent nodes  $X, Y$  such that  $Neighbor(X) \setminus \{Y\}$  has
       cardinality  $\geq n$ , and a subset  $S \subseteq Neighbor(X) \setminus \{Y\}$  of cardinality  $n$ , and, if
        $X \perp\!\!\!\perp Y \mid S$ , delete the edge between  $X$  and  $Y$  from  $Q$ , and record  $S$  in  $Sepset(X, Y)$ 
       and  $Sepset(Y, X)$ ;
6   until all ordered pairs of adjacent nodes  $X$  and  $Y$  such that  $Neighbor(X) \setminus \{Y\}$  has
       cardinality  $\geq n$  and all subsets  $S$  in  $Neighbor(X) \setminus \{Y\}$  have been tested for  $d$ -separation;
7    $n \leftarrow n + 1$ ;
8 until for each ordered pair of adjacent nodes  $X, Y$ ,  $Neighbor(X) \setminus \{Y\}$  has less than  $n$ 
   neighbors;
9 let  $F'$  be the undirected graph from the above step and orient each edge as  $\circ-\circ$ ;
10 foreach unshielded triple  $(X, Y, Z) \in F'$  do
11   if  $Y$  is not in  $Sepset(X, Z)$  then
12     orient  $X \ast \ast Y \ast \ast Z$  as  $X \ast \rightarrow Y \leftarrow \ast Z$ ;
13   end
14 end
15 foreach adjacent pair  $(X, Y) \in F'$  do
16   if there exists a subset  $S \subseteq Ext\text{-}D\text{-}SEP(X, Y)$  such that  $X \perp\!\!\!\perp Y \mid S$  then
17     delete the edge between  $X$  and  $Y$  from  $F'$ ;
18   end
19 end
20 let  $G$  be the undirected graph of  $F'$ ;
21 return  $G, Sepset$ 

```

PROOF. Because $X \xrightarrow{FD} Y$, $P(Y \mid X) = I_{Y=f(X)}$. Since $|X|, |Y| > 1$, the following inequality holds.

$$P(XY) = P(X)P(Y \mid X) = P(X)I_{Y=f(X)} \neq P(X)P(Y) \quad (9)$$

Therefore, $Y \not\perp\!\!\!\perp X$. Given a variable set Z ,

$$\begin{aligned}
 P(YZ \mid X) &= \frac{P(XYZ)}{P(X)} \\
 &= \frac{P(XZ)I_{Y=f(X)}}{P(X)} \\
 &= P(Z \mid X)I_{Y=f(X)} \\
 &= P(Z \mid X)P(Y \mid X)
 \end{aligned} \quad (10)$$

Therefore, $Z \perp\!\!\!\perp Y \mid X$. □

LEMMA 8.3.2. If $X \xrightarrow{FD} Y$ and $Z \perp\!\!\!\perp X \mid W$, then $Z \perp\!\!\!\perp Y \mid W$, where Z and W are two disjoint variable set other than X and Y .

Algorithm 4: FCI-Orient**Input:** Skeleton S , Sepset $Sepset$ **Output:** PAG G

```

1 let  $G$  be the graph sharing all adjacencies of  $S$  and orient each edge as  $\circ\circ$ ;
2 foreach unshielded triple  $(\alpha, \beta, \gamma) \in G$  do
3   if  $\beta$  is not in  $Sepset(\alpha, \gamma)$  then
4     | orient  $\alpha \ast \beta \ast \gamma$  as  $\alpha \ast \rightarrow \beta \leftarrow \ast \gamma$ ;
5   end
6 end
7 repeat
8    $\mathcal{R}1$ : if  $\alpha \ast \rightarrow \beta \circ \ast \gamma$ , and  $\alpha$  and  $\gamma$  are not adjacent, then orient the triple as  $\alpha \ast \rightarrow \beta \rightarrow \gamma$ ;
9    $\mathcal{R}2$ : if  $\alpha \rightarrow \beta \ast \gamma$  or  $\alpha \ast \rightarrow \beta \rightarrow \gamma$ , and  $\alpha \ast \circ \gamma$ , then orient  $\alpha \ast \circ \gamma$  as  $\alpha \ast \rightarrow \gamma$ ;
10   $\mathcal{R}3$ : if  $\alpha \ast \rightarrow \beta \leftarrow \ast \gamma$ ,  $\alpha \ast \circ \theta \circ \ast \gamma$ , and  $\alpha$  and  $\gamma$  are not adjacent, and  $\theta \ast \circ \beta$ , then orient  $\theta \ast \circ \beta$ 
    as  $\theta \ast \rightarrow \beta$ ;
11   $\mathcal{R}4$ : if  $u = (\theta, \dots, \alpha, \beta, \gamma)$  is a discriminating path between  $\theta$  and  $\gamma$  for  $\beta$ , and  $\beta \circ \ast \gamma$ ;
    then if  $\beta \in Sepset(\theta, \gamma)$ , orient  $\beta \circ \ast \gamma$  as  $\beta \rightarrow \gamma$ ; otherwise orient the triple  $(\alpha, \beta, \gamma)$  as
     $\alpha \leftrightarrow \beta \leftrightarrow \gamma$ ;
12 until none of the orientation rules applies;
13 repeat
14   $\mathcal{R}5$ : for every (remaining)  $\alpha \circ \circ \beta$ , if there is an uncovered circle path  $p = (\alpha, \gamma, \dots, \theta, \beta)$ 
    between  $\alpha$  and  $\beta$  s.t.  $\alpha$  and  $\beta$  are not adjacent, and  $\gamma$  and  $\theta$  are not adjacent, then orient
     $\alpha \circ \circ \beta$  and every edge on  $p$  as undirected edges ( $-$ );
15   $\mathcal{R}6$ : if  $\alpha - \beta \circ \ast \gamma$  ( $\alpha$  and  $\gamma$  may or may not be adjacent), then orient  $\beta \circ \ast \gamma$  as  $\beta \ast \gamma$ ;
16   $\mathcal{R}7$ : if  $\alpha - \beta \circ \ast \gamma$ , and  $\alpha$  and  $\gamma$  are not adjacent, then orient  $\beta \circ \ast \gamma$  as  $\beta \ast \gamma$ ;
17   $\mathcal{R}8$ : if  $\alpha \rightarrow \beta \rightarrow \gamma$  or  $\alpha \circ \beta \rightarrow \gamma$ , and  $\alpha \circ \rightarrow \gamma$ , orient  $\alpha \circ \rightarrow \gamma$  as  $\alpha \rightarrow \gamma$ ;
18   $\mathcal{R}9$ : if  $\alpha \circ \rightarrow \gamma$  and  $p = (\alpha, \beta, \theta, \dots, \gamma)$  is an uncovered p.d. path from  $\alpha$  to  $\gamma$  such that  $\beta$ 
    and  $\gamma$  are not adjacent, then orient  $\alpha \circ \rightarrow \gamma$  as  $\alpha \rightarrow \gamma$ ;
19   $\mathcal{R}10$ : suppose  $\alpha \circ \rightarrow \gamma$ ,  $\beta \rightarrow \gamma \leftarrow \theta$ ,  $p_1$  is an uncovered p.d. path from  $\alpha$  to  $\beta$  and  $p_2$  is an
    uncovered p.d. path from  $\alpha$  to  $\theta$ . Let  $\mu$  be the node adjacent to  $\alpha$  on  $p_1$  ( $\mu$  could be  $\beta$ ),
    and  $\omega$  be the node adjacent to  $\alpha$  on  $p_2$  ( $\omega$  could be  $\theta$ ). If  $\mu$  and  $\omega$  are distinct, and are
    not adjacent, then orient  $\alpha \circ \rightarrow \gamma$  as  $\alpha \rightarrow \gamma$ ;
20 until none of the orientation rules applies;
21 return  $G$ 

```

PROOF. Let $\pi_y = \{x \mid f(x) = y\}$. Therefore, $\{X \in \pi_y\} = \{Y = y\}$. Since $Z \perp\!\!\!\perp X \mid W$,

$$\begin{aligned}
 P(Y = y, Z \mid W) &= P(X \in \pi_y, Z \mid W) \\
 &= P(X \in \pi_y \mid W)P(Z \mid W) \\
 &= P(Y = y \mid W)P(Z \mid W)
 \end{aligned} \tag{11}$$

Therefore, $Z \perp\!\!\!\perp Y \mid W$. □

LEMMA 8.3.3. In a MAG \mathcal{G} over variables $\{X_1, \dots, X_n\}$, if $X_1 \rightarrow X_2$ and no other edges connect X_2 with other vertices, when $X_2 \perp\!\!\!\perp_{\mathcal{G}} Y \mid U$ for any $Y \in \{X_3, \dots, X_n\}$, $U \subset \{X_3, \dots, X_n\}$, $Y \notin U$, $X_1 \perp\!\!\!\perp_{\mathcal{G}} Y \mid U$.

PROOF. If $X_2 \perp\!\!\!\perp_{\mathcal{G}} Y \mid U$, according to the definition of m-separation, each path between X_2 and Y satisfies one of the following two conditions: 1) there exists $X_p \rightarrow X_s \rightarrow X_q$, $X_p \leftarrow X_s \rightarrow X_q$, or

$X_p \leftarrow X_s \leftarrow X_q$ where $X_s \in U$. 2) there exists $X_p \rightarrow X_s \leftarrow X_q$ where X_s or any of its descendants does not belong to U . Because $X_1 \notin U$, $X_1 \neq X_s$. For the first case, X_s blocks X_1 and Y given that there is only one edge connecting to X_2 . For the second case, X_1 cannot be a collider, because $X_1 \rightarrow X_2$. Therefore, U also blocks X_1 and Y on this path. In summary, $X_1 \perp\!\!\!\perp_{\mathcal{G}} Y \mid U$. \square

With above lemmas, now we prove Thm. 3.1.

PROOF. We first construct a MAG \mathcal{G} on the top of the skeleton, by adding an arrowhead from X_1 to Z (\mathcal{G}_2). Because \mathcal{S}_1 is learnt from $\{X_1, \dots, X_n\}$ where faithfulness assumption holds, this part is harmonious. For Z , there are two types of m-separation in \mathcal{G} . Here, we prove that each type of m-separation satisfies GMP to data distribution P_V .

Type 1: $Z \perp\!\!\!\perp_{\mathcal{G}} X_{j \neq 1} \mid X_1$

According to Lemma. 8.3.1, $X_1 \xrightarrow{\text{FD}} Z$ implies that $Z \perp\!\!\!\perp_{\mathcal{G}} X_{j \neq 1} \mid X_1$.

Type 2: $Z \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U (j \neq 1, U \cap \{X_1, X_j\} = \emptyset)$

By Lemma. 8.3.3, $Z \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$ implies that $X_1 \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$. Because \mathcal{S}_1 is harmonious, $X_1 \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$ implies that $X_1 \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$. According to Lemma. 8.3.2, $X_1 \xrightarrow{\text{FD}} Z$ and $X_1 \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$ imply $Z \perp\!\!\!\perp_{\mathcal{G}} X_j \mid U$.

Minimality is satisfied, since removing edge from X_1 to Z will make $X_1 \perp\!\!\!\perp Z$ which contradicts P_V . \square

8.4 Proof of Thm. 3.2

PROOF. We prove Thm. 3.2 by mathematical induction. In the sense, Alg. 1 returns a harmonious skeleton when \mathcal{G}_{FD} has arbitrary number of non-root vertices. Denote the number of non-root vertices as $s \geq 0$.

Base Case. When $s = 0$, the skeleton is harmonious because all variables are under faithfulness assumptions.

Induction. Suppose the returned skeleton is harmonious for \mathcal{G}_{FD} when $s = n$. Now we prove the skeleton is still harmonious when we add a vertex X' to a vertex set $X \subseteq \mathcal{G}_{\text{FD}}.V$. Denote the new FD-induced graph as \mathcal{G}'_{FD} , the skeleton of \mathcal{G}_{FD} as \mathcal{S} , and the skeleton for \mathcal{G}'_{FD} as \mathcal{S}' . According to Alg. 1, $\mathcal{S}, \mathcal{S}'$ share the same vertices and edges (except X'). Given that \mathcal{S} is harmonious, \mathcal{S}' can be decomposed into two subgraphs, i.e., $\mathcal{S}_1, \mathcal{S}_2$, where $\mathcal{S}_1 = \mathcal{S}$ and \mathcal{S}_2 corresponds to X' and one edge from X' to one of X . Since $X'' \xrightarrow{\text{FD}} X'$, by Thm. 3.1, \mathcal{S}' is harmonious.

By the principle of mathematical induction, Thm. 3.2 holds. \square

8.5 FCI Rules on FD-related Edges

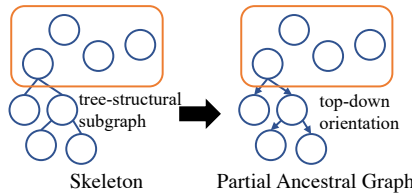


Fig. 8. Orientation on tree-structural subgraph.

Since we only consider one-to-one and one-to-many FDs, FD-induced vertices only have at most one parent node in \mathcal{G}_{FD} . Therefore, it forms tree-structural subgraphs in the skeleton starting from the root vertices of \mathcal{G}_{FD} (e.g., the left-hand side skeleton shown in Fig. 8). Intuitively, it would be

ideal if the orientation rules allow us to orient the tree-structural subgraphs in a top-down manner (e.g., the right-hand side of Fig. 8). In particular, if there is an edge heading to the root node of the tree-structural subgraph, we can assign \rightarrow (\rightarrow is defined in Table 1) as the edge from the root node to its child vertices and propagate the directions to the leaf node by applying Rule 1 of FCI recursively.

DEFINITION 8.8 (RULE 1 OF FCI [57]). *If $X * \rightarrow Y \circ \rightarrow Z$ and X, Z are non-adjacent, then orient $Y \rightarrow Z$.*³

However, it is not necessarily attainable if we cannot identify the direction on the edge heading to the root node (i.e., Y in Def. 8.8). As a result, we cannot derive any directions on the subgraphs. That is, all edges on the subgraph are $\circ \circ$, which denote edges with unknown directions. Such ambiguous cases are surely undesirable in XDA scenarios.

8.6 ANM on FD-related Edges

We note that, from viewpoint of conditional independence, functional dependency does not imply causal directions in all cases. However, we present that the counterexamples of such cases (i.e., contrasting direction between causal relation and functional dependency) are very rare in practice, if at all possible. Recall the theorem of “Identifiability of discrete ANMs” in [43].

THEOREM 8.1. *Assume that a distribution $P_{X,Y}$ allows for an ANM $Y = f(X) + N_Y$ from X to Y and that either X or Y has finite support. $P_{X,Y}$ allows for an ANM from Y to X if and only if there exists a disjoint decomposition $\bigcup_{i=0}^I C_i = \text{supp}X$, such that the following conditions a), b), and c) are satisfied:*

a) *The C_i ’s are shifted versions of each other*

$$\forall i \exists d_i > 0 : C_i = C_0 + d_i \quad (12)$$

and f is piecewise constant: $f \mid C_i \equiv c_i \forall i$.

b) *The probability distributions on the C_i s are shifted and scaled versions of each other with the same shift constant as above: For $x \in C_i$, $P(X = x)$ satisfies*

$$P(X = x) = P(X = x - d_i) \cdot \frac{P(x \in C_i)}{P(x \in C_0)} \quad (13)$$

c) *The sets $c_i + \text{supp}N_Y := \{c_i + h : P(N_Y = h) > 0\}$ are disjoint sets.*

By condition c), since $N_Y = 0$ by functional dependency, $\nexists x_i \in C_i, x_j \in C_j$ s.t. $f(x_i) = f(x_j)$. In other word, X are decomposed by the value of $f(x)$ and each $c_i \in f(X)$ forms a disjoint subset C_i of X . To admit condition a), each C_i is a shifted version of others. Therefore, each C_i should at least have an equal cardinality. The above two conditions imply that each $y \in Y$ corresponds to equal size of $x \in X$. Consider the aforementioned example where Country $\xrightarrow{\text{FD}}$ Continent. To satisfy the condition, each continent must have equal number of countries. Furthermore, by condition b), the probability of a country x in continent c_i in the database can be exactly scaled by $\frac{P(x \in C_i)}{P(x \in C_0)}$ the probability of another shifted country $x - d_i$. The conditions together exhibits a very rare scenario. Therefore, we ignore such cases in practice.

As consistent with [43], we consider it reasonable to infer that direction of ANM as causal. Therefore, we hypothesis that functional dependencies in \mathcal{G}_{FD} intrinsically imply causal directions.

³“*” is used as a wildcard symbol accepting either $\rightarrow, \leftrightarrow, \circ \rightarrow$.

8.7 Proof of Proposition 3.2

PROOF. Suppose there exists an optimal explanation P^* such that $\exists p_i \in P^*, \Delta_i \leq 0$. Let $P' = \{p \mid p \in P^*, \Delta_p > 0\}$ and Γ be the minimal contingency of P^* .

For SUM, since $\Delta(D_{P'}) \geq \Delta(D_{P^*})$, $\Delta(D) - \Delta(D_\Gamma) - \Delta(D_{P'}) < \Delta(D) - \Delta(D_\Gamma) - \Delta(D_{P^*}) \leq \epsilon$.

For homogeneous AVG,

$$\begin{aligned} \Delta(D - \Gamma - P') &= \frac{\sum_{p_i \in \{p_1, \dots, p_m\} - \Gamma - P'} a_i \Delta_i}{A_{D - \Gamma - P'}} \\ &< \frac{\sum_{p_i \in \{p_1, \dots, p_m\} - \Gamma - P^*} a_i \Delta_i}{A_{D - \Gamma - P^*}} \\ &= \Delta(D - \Gamma - P^*) \leq \epsilon \end{aligned} \quad (14)$$

Hence, Γ is also a valid contingency for P' . $\rho_{P'} \geq \rho_{P^*}$ and $|P'| < |P^*|$ contradict the fact that P^* is optimal. \square

8.8 Proof of Proposition 3.3

PROOF. When there are more than one optimal explanations, let P' be the one with the smallest predicate size (i.e., $|P|$) and Γ' be the corresponding minimal contingency. Otherwise, let P' and Γ' be the optimal explanation and corresponding minimal contingency, respectively. If $|P'| \geq |P^C|$, given that $\rho_{P'} \leq 1$ and $\rho_{P^C} = 1$, $\rho_{P'} - \sigma \rho_{P'} \leq \rho_{P^C} - \sigma \rho_{P^C}$. P^C is at least as optimal as P' . If $|P'| < |P^C|$, then let P'' be a predicate by replacing all non-canonical filters in P' with canonical ones. Then, $1 - d_{P''} \leq 1 - d_{P'} \leq \epsilon'$. Γ' is then a valid contingency for P'' . Therefore, $\rho_{P''} \geq \rho_{P'}$ and $|P''| = |P'|$. $\rho_{P''} - \sigma \rho_{P''} \geq \rho_{P'} - \sigma \rho_{P'}$. $P'' \subset P^C$ is at least as optimal as P' . In summary, there must exist an optimal explanation $P'' \subseteq P^C$. \square

8.9 Proof of Thm. 3.3

PROOF. Let $m_k = \sum_{i=1}^k d_{p_i}$. According to the definition of canonical predicate, $1 - m_j \leq \epsilon' < 1 - m_{j-1}$. This inequality implies that $1 - (m_j - d_{p_j}) + d_{p_j} \leq \epsilon' < 1 - (m_j - d_{p_j})$. Since $\bar{P} \subset P^C$ and $d_{p_1} \leq \dots \leq d_{p_j}$, $d_{\bar{P}} \geq d_{p_j}$. The following inequalities hold.

$$\begin{aligned} \epsilon' &\geq 1 - (m_j - d_{\bar{P}}) - d_{\bar{P}} = 1 - d_P - d_{\bar{P}} \\ \epsilon' &< 1 - (m_j - d_{p_j}) \leq 1 - (m_j - d_{\bar{P}}) = 1 - d_P \end{aligned} \quad (15)$$

Since $1 - d_P - d_{\bar{P}} \leq \epsilon' < 1 - d_P$, Thm. 3.3 holds. \square

8.10 Proof of Thm. 3.4

PROOF. $\rho_P = \frac{1}{1 + \min(|\Gamma|_W)} \geq \frac{1}{1 + |\bar{P}|_W} = \frac{1}{1 + m_j - d_P}$. Furthermore, since $1 - d_P - d_\Gamma \leq \epsilon'$, $d_\Gamma \geq 1 - d_P - \epsilon'$ and $\rho_P \leq \frac{1}{2 - d_P - \epsilon'}$. Thus, we derive the lower and upper bounds of ρ_P . When assuming $d_P \ll m_j$ and

$$\begin{aligned} 0 < m_j \leq 1, \quad \frac{1}{1 + m_j - d_P} &= \frac{1 + m_j + d_P}{(1 + m_j)^2 - d_P^2} \approx \frac{1 + m_j + d_P}{(1 + m_j)^2}. \text{ The approximation error rate } E = \frac{\frac{1}{1 + m_j - d_P} - \frac{1 + m_j + d_P}{(1 + m_j)^2}}{\frac{1}{1 + m_j - d_P}} \\ &= \frac{d_P^2}{(1 + m_j)^2} < \frac{m_j^2}{(1 + m_j)^2} \leq 0.25. \end{aligned} \quad \square$$

8.11 Proof of Proposition 3.4

To prove Proposition 3.4, we first introduce the following notations and propositions.

Notations. Given WHY QUERY Δ and corresponding sibling subspaces s_1, s_2 , let x_i, y_i be $agg_M(D_{s_1 \cap p_i})$ and $agg_M(D_{s_2 \cap p_i})$, respectively. We also let a_i, b_i be $|D_{s_1 \cap p_i}|$ and $|D_{s_2 \cap p_i}|$ denoting the rows of $D_{s_1 \cap p_i}$ and $D_{s_2 \cap p_i}$, respectively, and A_D, B_D be $|D_{s_1}|$ and $|D_{s_2}|$. Then, it is obvious that $\Delta(D_{p_i}) = x_i - y_i$

and we use Δ_i as a shorthand for $\Delta(D_{P_i})$. For AVG and $\Delta(D_P) = \sum_{P_i \in P} \Delta_i$, $\Delta(D)$ can be represented in the form of $\sum_1^m (\frac{a_i x_i}{A_D} - \frac{b_i y_i}{B_D})$.

PROPOSITION 8.1. *If a pair of sibling subspaces are homogeneous on X , then $\frac{a_1}{b_1} = \dots = \frac{a_m}{b_m}$.*

PROOF. For homogeneous AVG, X, F are m -separated given B . Hence, $X \perp\!\!\!\perp F \mid B$ and $P(X, F \mid B) = P(X \mid B)P(F \mid B)$. Recall that $a_i = P(X = x_i, F = f_1 \mid B = b)$ and $b_i = P(X = x_i, F = f_2 \mid B = b)$. We have $\frac{a_i}{b_i} = \frac{P(F=f_1|B=b)}{P(F=f_2|B=b)}$ which is a constant and invariant with respect to i . \square

PROPOSITION 8.2. *If P_1 and P_2 are two disjoint predicates on the same attribute, for homogeneous AVG, $\Delta(D_{P_1} + D_{P_2}) < \Delta(D_{P_1}) + \Delta(D_{P_2})$.*

PROOF. For homogeneous AVG, $\Delta(D_{P_1}) = \frac{\sum_{P_i \in P_1} a_i \Delta_i}{A_{P_1}}$, $\Delta(D_{P_2}) = \frac{\sum_{P_i \in P_2} a_i \Delta_i}{A_{P_2}}$ and $\Delta(D_{P_1} + D_{P_2}) = \frac{\sum_{P_i \in P_1 \cup P_2} a_i \Delta_i}{A_{P_1 \cup P_2}}$. Also, $A_{P_1 \cup P_2} = A_{P_1} + A_{P_2}$. Hence, $A_{P_1} < A_{P_1 \cup P_2}$ and $A_{P_2} < A_{P_1 \cup P_2}$. Therefore,

$$\begin{aligned} \Delta(D_{P_1}) + \Delta(D_{P_2}) &= \frac{\sum_{P_i \in P_1} a_i \Delta_i}{A_{P_1}} + \frac{\sum_{P_i \in P_2} a_i \Delta_i}{A_{P_2}} \\ &> \frac{\sum_{P_i \in P_1} a_i \Delta_i}{A_{P_1 \cup P_2}} + \frac{\sum_{P_i \in P_2} a_i \Delta_i}{A_{P_1 \cup P_2}} \\ &= \frac{\sum_{P_i \in P_1 \cup P_2} a_i \Delta_i}{A_{P_1 \cup P_2}} \\ &= \Delta(D_{P_1} + D_{P_2}) \end{aligned} \tag{16}$$

\square

Now, we prove Proposition 3.4.

PROOF.

$$\begin{aligned} \Delta_{D_P - D_{P_j}} &= \frac{\sum_{P_i \in P} a_i \Delta_i - a_j \Delta_j}{A_{D - P_j}} \\ &= \frac{\sum_{P_i \in P} a_i \Delta_i - a_j \Delta_j}{A_D - a_j} \\ &< \frac{\sum_{P_i \in P} a_i \Delta_i}{A_D} \\ &= \Delta_{D_P} \end{aligned} \tag{17}$$

The inequality in the above equation is true because $\frac{A}{B} < \frac{C}{D}$ implies $\frac{A-C}{B-D} < \frac{A}{B}$ where A, B, C, D are positive and $A > C, B > D$. \square

8.12 Generating Synthetic Data

③ **Synthetic Data A (SYN-A).** We use Erdős-Rényi model, a well-established random graph model, to synthesize random causal graphs of different scales and to produce datasets via forward sampling [44]. To simulate causally insufficient systems, we mask 5% variables at random and return the corresponding PAG (Partial Ancestral Graph) as the ground truth. We construct conditional probability tables based on a Dirichlet prior and generate two additional FD nodes on each leaf node. Afterwards, these functional dependencies are employed to build FD-induced graphs.

③ **Synthetic Data B (SYN-B).** In particular, we design a data generating process with three variables X, Y, Z , where X is a binary variable, Y is a categorical variable, and Z is a numerical variable. Different values in X first impact Y and then Y 's values further impact Z . When Y 's

realizations are equal to some specific values (i.e., $Y \in \{y_1, \dots, y_k\}$), Z would be sampled from a Gaussian distribution with a higher mean μ^* ; otherwise, Z is sampled from another Gaussian distribution with a lower mean $\mu < \mu^*$. Note that the higher k , the harder XPLAINER (and other tools) to comprehensively identify all filters to form the correct explanation. As a result, given different X , the aggregated result of Z differs. Recall Def. 2.1 where X and Z form the context and the target of a WHY QUERY. Then, we seek to extract the explanation from Y and the way Y impacts Z specified in the data generating process (i.e., $Y = y_1 \vee \dots \vee Y = y_k$) constitutes the ground truth of this WHY QUERY. We concretize the above data generating process with different parameters (e.g., cardinality and conditional probability table of Y and conditional distributions of Z) and yield 18 datasets of different difficulties. By default, we generate **SYN-B** datasets with 10,000 rows; the variable Y contains ten values, three of which would trigger abnormal Z (i.e., $Z \sim \mathcal{N}(\mu^*, 10)$, $\mu^* = 60$) while the others would produce normal $Z \sim \mathcal{N}(\mu, 10)$, $\mu = 10$. The parameters are on a par with the configuration in Scorpion.

Received October 2022; revised January 2023; accepted February 2023