



HYPER: Hypothetical Reasoning With What-If and How-To Queries Using a Probabilistic Causal Approach

Sainyam Galhotra*
University of Chicago
sainyam@uchicago.edu

Amir Gilad*
Duke University
agilad@cs.duke.edu

Sudeepa Roy
Duke University
sudeepa@cs.duke.edu

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

ABSTRACT

What-if (provisioning for an update to a database) and how-to (how to modify the database to achieve a goal) analyses provide insights to users who wish to examine hypothetical scenarios without making actual changes to a database and thereby help plan strategies in their fields. Typically, such analyses are done by testing the effect of an update in the existing database on a specific view created by a query of interest. In real-world scenarios, however, an update to a particular part of the database may affect tuples and attributes in a completely different part due to implicit semantic dependencies. To allow for hypothetical reasoning while accommodating such dependencies, we develop HYPER, a framework that supports what-if and how-to queries accounting for probabilistic dependencies among attributes captured by a probabilistic causal model. We extend the SQL syntax to include the necessary operators for expressing these hypothetical queries, define their semantics, devise efficient algorithms and optimizations to compute their results using concepts from causality and probabilistic databases, and evaluate the effectiveness of our approach experimentally.

CCS CONCEPTS

• **Theory of computation** → **Incomplete, inconsistent, and uncertain databases**; • **Computing methodologies** → **Causal reasoning and diagnostics**.

KEYWORDS

hypothetical reasoning, causal inference, what-if, how-to

ACM Reference Format:

Sainyam Galhotra, Amir Gilad, Sudeepa Roy, and Babak Salimi. 2022. HYPER: Hypothetical Reasoning With What-If and How-To Queries Using a Probabilistic Causal Approach. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD '22)*, June 12–17, 2022, Philadelphia, PA, USA. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3514221.3526149>

*Both authors contributed equally to this research.

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 ACM 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.

SIGMOD '22, June 12–17, 2022, Philadelphia, PA, USA

© 2022 Association for Computing Machinery.

ACM ISBN 978-1-4503-9249-5/22/06...\$15.00

<https://doi.org/10.1145/3514221.3526149>

1 INTRODUCTION

Hypothetical reasoning is a crucial element in decision-making and risk assessment in business [24, 47, 53], healthcare [39, 40], real estate [19], etc. Such analysis is split by previous work into two categories: what-if analysis and how-to analysis. What-if analysis [9, 29, 31] is usually meant for testing assumptions and projections on a particular outcome by allowing users to pose queries about hypothetical updates in the database and examining their effect on a query result. Users detail a specific hypothetical scenario whose effect they wish to examine on their view of choice and the system computes the view as if the update has been performed in the database. On the other hand, how-to analysis [33, 35] has the reverse goal; users specify a target effect that they want to achieve and the system computes the appropriate hypothetical updates that have to be performed in the database to fulfill the goal.

EXAMPLE 1. Consider a simplified version of the Amazon product database [28] shown in Figure 1 describing product details and product reviews. Each tuple has a unique tuple identifier next to it for clarity. Now, consider an analyst who wants to examine the effect of laptop prices on their Amazon ratings. She may ask “what would be the effect of increasing the price of Asus laptops by 10% on their average ratings?”. This what-if query asks about the effect of the hypothetical update on the database (increasing the Price) on a specific view (average Rating). She may also be interested in “what fraction of Asus laptops would have rating more than 4.0 if their price drops by \$100?” or “What would be the average sentiment in the reviews for cameras if their color was changed to red?”. A different analyst may also be interested in maximizing the average rating of laptops reviews by changing their price. She may ask “how to maximize the average rating of laptops and cameras by updating the price of laptops so that it will not drop below 500 and increase above 800, and will be at most 100 away from its original value?” or “How to increase average sentiment in the reviews for cameras by changing their color?”. Both queries are forms of hypothetical reasoning that can assist analysts and decision-makers in gaining insights about their products and their marketing strategies.

Multiple works in the database community have studied hypothetical reasoning. A substantial part of these [7, 16–18, 33, 35] has focused on provenance updates and view manipulation as a main component for answering such queries. Therein, hypothetical updates are captured by changing values in the provenance and, thus, updating the view generated by the query of interest. However, in many real-world situations, due to complex probabilistic causal dependencies between attributes of tuples that are relationally connected, updating an attribute of a tuple has collateral effects on other attributes of the same tuple, as well as attributes of other

	PID	Category	Price	Brand	Color	Quality
p_1	1	Laptop	999	Vaio	Silver	0.7
p_2	2	Laptop	529	Asus	Black	0.65
p_3	3	Laptop	599	HP	Silver	0.5
p_4	4	DSLR Camera	549	Canon	Black	0.75
p_5	5	Sci Fi eBooks	15.99	Fantasy Press	Blue	0.4

(a) Product

	PID	ReviewID	Sentiment	Rating
r_1	1	1	-0.95	2
r_2	2	2	0.7	4
r_3	2	3	-0.2	1
r_4	3	3	0.23	3
r_5	3	5	0.95	5
r_6	4	5	0.7	4

(b) Review

Figure 1: Amazon product database

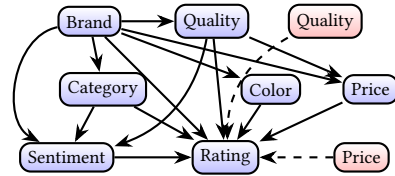


Figure 2: A graph showing the dependencies between the attributes in the database in Figure 1. Blue nodes are attributes of the same tuple and the red node is an attribute of a different tuple. A dashed edge denotes a dependency between attributes of different tuples

tuples. Such dependencies cannot be expressed and captured by provenance. We illustrate with an example.

EXAMPLE2. *Reconsider Example 1. The provenance of the average rating of Asus laptops will not change if the price of the laptops is augmented. Similarly, for the how-to query, the provenance of the average rating of laptops and cameras will not be affected by the change in price. Thus, previous work in databases fails to account for the collateral effect that increasing the price of a laptop may have on the user's ratings. Note that due to our lack of knowledge about the underlying process that leads to the user's ratings, we may only reason about the probabilistic effect of increasing the price on user's ratings. Figure 2 gives an intuitive description of potential dependencies between the attributes of the database in Figure 1. For example, changing the Price of a laptop may affect its Rating (denoted as the edge from the blue Price node to the blue Rating node in Figure 2). Furthermore, increasing the Price of Asus laptops may affect the Rating of Vaio laptops and vice versa (denoted as the edge from the red Price node to the blue Rating node in Figure 2). In general, a directed edge stands for an effect of the outbound node on the inbound node, e.g., Price affects Rating. Accounting for such dependencies is crucial for sound hypothetical reasoning.*

In this paper, we propose a novel probabilistic framework for hypothetical reasoning in relational databases that accounts for collateral effects of hypothetical updates on the entire data. Our system, **HYPER** (**H**ypothetical **R**easoning), allows users to ask complex relational what-if and how-to queries using a SQL-like declarative language. The underlying inference mechanism, then, internally accounts for the probabilistic causal effect of hypothetical updates and computes probabilistic answers to such hypothetical queries. Our framework brings together techniques from probabilistic databases [6, 15], and recent advancements in inference from relational data [46, 52, 54], to provide a principled approach for computing complex what-if

and how-to queries from relational databases. Specifically, **HYPER** relies on causal reasoning to capture background knowledge on probabilistic causal dependencies between attributes and interprets hypothetical updates as real-world actions that potentially affect the other attributes.

Our framework supports a rich class of what-if queries that involve joins and aggregations to support complex real-world what-if scenarios in relational domains. **HYPER** captures what-if queries through a novel model that can accommodate complex probabilistic dependencies, and computes their results efficiently by employing optimizations from probabilistic databases and causal inference. In addition, **HYPER** supports complex how-to queries and frames them as an optimization problem on the search space of consistent what-if queries, and searches for a hypothetical update that optimizes the desired query result. **HYPER** employs an efficient routine to solve this optimization problem, by expressing it as an Integer Program (IP) that can be efficiently handled using the existing IP solvers.

Our main contributions can be summarized as follows:

- We propose a formal probabilistic model for hypothetical what-if and how-to queries in relational domains that combines notions from probabilistic databases and causality. Our model assigns a probability to each possible world [15] that can be obtained after a hypothetical update according to the underlying probabilistic causal dependencies. We further define a probabilistic possible world semantics for complex what-if and how-to queries that support joins and aggregations.
- We develop a declarative language that extends the standard SQL syntax with new operators that capture hypothetical reasoning in relational domains and allow users to succinctly formulate complex probabilistic what-if and how-to queries.
- Evaluating hypothetical queries in a naive manner can be inefficient due to the need to iterate over all possible worlds, or explore the space of all possible hypothetical updates. To address these, we develop a suite of optimizations that allows **HYPER** to efficiently evaluate hypothetical queries:
 - We use the model of block-independent databases [41], i.e., the database can be partitioned into blocks of tuples, where the tuples in different blocks are independent, meaning there are no causal dependencies between the tuples across different blocks (without background knowledge, we assume tuple independence). We then show that what-if queries can be evaluated independently within each block and the results can be combined to get the result over the entire database.
 - We further show that under some assumptions complex what-if queries in relational domains can be evaluated using the existing techniques in causal inference and machine learning.
 - We frame how-to queries as an optimization problem and develop an efficient mechanism to solve this optimization problem, by expressing it as an Integer Program (IP) that can be efficiently handled using the existing IP solvers.
- We perform an extensive experimental evaluation of **HYPER** on both real and synthetic data. On real datasets, we show that the query output by **HYPER** matches the conclusions from prior studies in fair and explainable AI [23]. On synthetic datasets, we show that **HYPER**'s query output is accurate as compared to other baselines. Running time analysis shows that both what-if and how-to components of **HYPER** are highly efficient.

2 PROBABILISTIC UPDATES IN HYPER

In this section we describe our notations and then define the probabilistic hypothetical update model in HYPER (Section 2.1) that serve as the basis for probabilistic what-if and how-to queries in the following sections. Then in Section 2.2, we review necessary concepts from probabilistic causal models [38] that capture the propagation of the effect of an update through other attributes due to underlying dependencies between them and succinctly defines the probability distribution after updates.

Notations. Let D be a standard multi-relational database; we use D for both schema and instance (as a set of tuples) where it is clear from the context. For each relation R in D , $\text{Attr}(R)$ denotes the set of attributes of R and $\mathbf{A} = \cup_{R \in D} \text{Attr}(R)$ denotes the set of attributes in D . For attributes A appearing in multiple relations, we use $R.A$ for disambiguation. For an attribute $A \in \mathbf{A}$, $\text{Dom}(A)$ denotes the domain of A ; $A_i[t] \in \text{Dom}(A_i)$ denotes the value of the attribute A_i of the tuple t . We assume that each relation R has a (primary) key, that can be a single or a combination of multiple attributes. For easy reference, we annotate each tuple with a unique identifier as demonstrated by the identifiers p_i, r_j in Figure 1. We assume each relation can be modeled as a set of tuples (set semantics) and, for a relation R , we use the notation $t \in R$ to denote a tuple in R .

For the purpose of hypothetical updates, a subset of attributes that can change values directly or indirectly in tuples is referred to as **mutable attributes**, the other attributes are **immutable attributes**. The attribute that is updated in hypothetical updates is called the **update attribute**, and the effect is measured on an **output attribute** as specified by the user. The update and output attributes are always mutable, and the key attributes are always immutable.

EXAMPLE3. In Figure 1a, the database has two relations *Product* and *Review* with keys $\{\text{PID}\}$ and $\{\text{PID}, \text{ReviewID}\}$ respectively. For example, suppose $\text{Dom}(\text{Price}) = [0, 500K]$. In tuple p_1 , $\text{Category}[p_1] = \text{Laptop}$ and $\text{Price}[p_1] = 999$ etc. The mutable attributes are *Price*, *Quality*, *Color*, *Rating*, and *Sentiment*, whereas *Brand* and *Category* are immutable. The update attribute is *Price* in relation *Product*, and the output attribute is *Rating* in relation *Review*.

We assume the update and output attributes do not appear in multiple relations, but as Example 3 illustrates, they can appear in two different tuples.

2.1 Probabilistic Hypothetical Updates

HYPER interprets hypothetical updates in terms of real-world interventions that potentially influence the value of other attributes in the data due to probabilistic dependencies between the attributes and tuples. To capture such probabilistic influence, we use the notion of *possible worlds* from the literature of probabilistic databases [15] as the set of all possible instances on the same schema with the same number of tuples in each relation that may contain different values in their mutable attributes from the appropriate domains.

DEFINITION 1 (POSSIBLE WORLDS). Let R in D be a relation, where in $\text{Attr}(R)$, A_1, \dots, A_m are immutable attributes (including keys) and B_1, \dots, B_ℓ are mutable attributes. For a tuple $t \in R$, a **possible world of tuple t** is the set (assuming values are associated with

corresponding attribute names for disambiguation)

$$PWD(t) = \{A_1[t], \dots, A_m[t], v_1, \dots, v_\ell : v_i \in \text{Dom}(B_i), i = 1 \text{ to } \ell\}.$$

The **possible worlds of relation R** is $PWD(R) = \times_{t \in R} PWD(t)$. The **possible worlds of a database D** is $PWD(D) = \times_{R \in D} PWD(R)$.

Next we define the notion of hypothetical updates.

DEFINITION 2 (HYPOTHETICAL UPDATES). A **hypothetical update** $U = u_{R,B,f,S}$ on a database D is a 4-tuple that includes a relation R in D containing the mutable update attribute $B \in \text{Attr}(R)$, a subset of tuples $S \subseteq R$ where the update will be applied, and a function $f : \text{Dom}(B) \rightarrow \text{Dom}(B)$ specifying the update for attribute $B[t]$ for tuples $t \in S$ to $f(B[t])$.

In other words, the hypothetical update $u_{R,B,f,S}$ forces all tuples in set S in relation R to take the value $f(B[t])$ instead of $B[t]$. In the what-if query in Example 1, intuitively, $R = \text{Product}$, S defines the set of Asus laptops, B is *Price*, and f increases the price by 10% (see Section 3.1 for details). This update, in turn, may change values of other mutable attributes in R or even mutable attributes in other relations R' in D through causal dependencies as discussed next in Section 2.2, eventually (possibly) changing the output attribute. These changes are likely not deterministic (e.g., changing price of a laptop does not change its reviews or their sentiments in a fixed way), therefore, we model the state of the database after a hypothetical update as a probability distribution called the *post-update distribution*.

DEFINITION 3 (POST-UPDATE DISTRIBUTION). Given a database D and an update $U = u_{R,B,f,S}$ (Definition 2), the **post-update distribution** is a probability distribution over possible worlds, i.e., $\text{Pr}_{D,U} : PWD(D) \rightarrow [0, 1]$ such that $\sum_{I \in PWD(D)} \text{Pr}_{D,U}(I) = 1$.

While the previous definition defines the post-update distribution in a generic form, there will be restrictions imposed by the hypothetical update as well as by its effect on the distribution of other attributes (e.g., for all possible worlds with non-zero probability, the value of attribute B for tuples $t \in S$ must be $f(B[t])$). We define this post-update distribution with the help of a probabilistic relational causal model in Section 2.2.

2.2 Causal Model for Probabilistic Updates

In this paper, we use causal modeling to capture probabilistic causal dependencies between attributes in relational domains, and to account for the collateral effect of hypothetical updates on other attributes. Specifically, HYPER rests on relational causal models, recently introduced in [46], which are briefly reviewed next.

Probabilistic Relational Causal Models (PRCM). A probabilistic relational causal model (PRCM) associated with a relational instance D is a tuple $(\epsilon, \mathcal{V}, \text{Pr}_\epsilon, \phi)$, where ϵ is a set of unobserved *exogenous (noise)* variables distributed according to Pr_ϵ , \mathcal{V} is a set of *endogenous ground*¹ variables associated with observed attribute values of each tuple $A[t]$, for all $A \in \text{Attr}(R)$, $t \in R$ and $R \in D$, and ϕ is a set of *structural equations*. The structural equations capture the *causal dependencies* among the attributes and are of the

¹The endogenous variables are called *ground* variables since in a PRCM the attribute $A[t]$ associated with each tuple t form the variables, generating multiple variables corresponding to the same attribute, in contrast to the standard probabilistic causal model [38] where each attribute or feature A forms a unique variable.

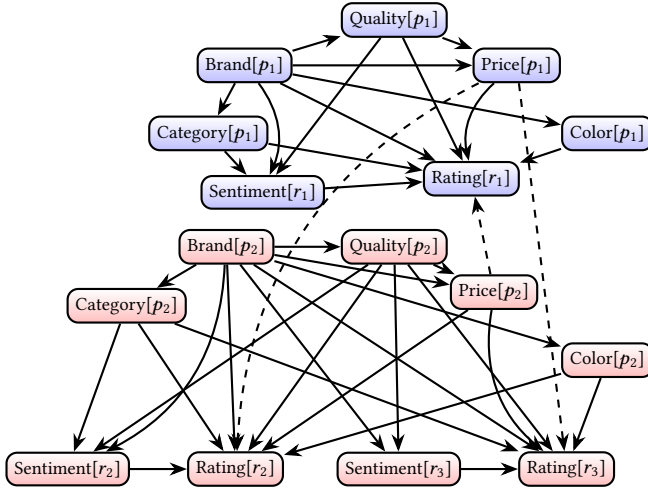


Figure 3: Part of the ground causal graph for the tuples in Figure 1. $A[p_i]$ ($A[r_j]$) represents the attribute A of tuple p_i (r_j). Blue nodes are related to p_1 , red nodes are related to p_2 , and dashed edges represent cross-tuple dependencies. Cross-tuple edges between Quality and Rating are dropped.

form $\phi_{A_i[t]} : \text{Dom}(Pa_{\mathcal{V}}(A_i[t])) \times \text{Dom}(Pa_{\mathcal{E}}(A_i[t])) \rightarrow \text{Dom}(A_i[t])$, where $Pa_{\mathcal{E}}(A_i[t]) \subseteq \epsilon$ and $Pa_{\mathcal{V}}(A_i[t]) \subseteq \mathcal{V} - \{A_i[t]\}$ respectively denote the exogenous and endogenous parents of $A_i[t]$. A PRCM is associated with a *ground causal graph* G , whose nodes are the endogenous variables \mathcal{V} and whose edges are all pairs (X, Y) (directed edges) such that $X \in \mathcal{V}$ and $Y \in Pa_{\mathcal{V}}(A_i[t])$. In this paper we assume the underlying causal model is acyclic. Due to uncertainty over the unobserved noise variables, the structural equations can be seen as a set of probabilistic dependencies² of the form $\Pr(A[t] \mid Pa_{\mathcal{V}}(A[t]))$ between the attributes. From now on, we will use $A[t]$ interchangeability to refer to both an attribute value and the ground variable associated with it.

EXAMPLE4. Reconsider the database in Figure 1 and the causal diagram in Figure 2. Part of its ground version w.r.t. the database is depicted in Figure 3, where the blue nodes are related to the tuple p_1 and the red nodes are related to the tuple p_2 . Cross-attribute dependencies within the same tuple are illustrated as solid edges and cross-tuple dependencies between the tuples are shown as dashed edges.

To be able to estimate the conditional probability distributions $\Pr(A[t] \mid Pa_{\mathcal{V}}(A[t]))$, for $t \in R$, from the relational instance D , we make the following assumptions that are common in causal inference from relational data [46, 52]. First, since $Pa_{\mathcal{V}}(A[t])$, the set of parents of $A[t]$ may have variable cardinality for each $t \in R$, we assume there exists a distribution preserving summary function ψ that projects $Pa_{\mathcal{V}}(A[t])$ into a fixed size vector such that $\Pr(A[t] \mid Pa_{\mathcal{V}}(A[t])) = \Pr(A[t] \mid \psi(Pa_{\mathcal{V}}(A[t])))$, for each $t \in R$. Second, we assume the conditional probability distributions $\Pr(A[t] \mid \psi(Pa_{\mathcal{V}}(A[t])))$ are the same for all $t \in D$, i.e., the conditional probability distributions $\Pr(A_i[t] \mid \psi(Pa_{\mathcal{V}}(A_i[t])))$ are independent of a particular $t \in R$ and can be readily estimated from

²Note that it is not necessary to have relational connections through database constraints like foreign key dependencies or functional dependencies for causal dependencies and vice versa.

D , hence we denote them by unified notation $\Pr_D(A_i \mid \psi(Pa_{\mathcal{V}}(A_i)))$. For more discussion on these assumptions, please see [46].

EXAMPLE5. Continuing Example 1, suppose we want to update attribute Price and examine its effect on Rating. Since each product has one price but several review ratings in Figure 1, we will summarize the Rating attribute into the Product table by, e.g., averaging the Rating for each product and price. Thus, for p_2 , we will have Price = 529 and Rating = Average(4, 2) = 3 (the average over tuples r_2 and r_3).

Post-update distribution by PRCM. We describe how the post-update distribution (Definition 3) is defined using a PRCM in HYPER. Given a relation R in D , an update attribute $B \in \text{Attr}(R)$, a hypothetical update $U = u_{R,B,f,S}$ (Definition 2) can be interpreted as an *intervention* that modifies the underlying PRCM and replaces the structural equation associated with the variables $B[t]$ for all $t \in S$ with the constant $f(B[t])$. Updating $B[t]$ propagates through all relations, tuples and attributes according to the underlying PRCM. The *post-update* state of a tuple $t' \in R'$ in a relation R' in D is the *solutions* to each ground variable $A[t']$, for $A \in \text{Attr}(R')$, in the modified set of structural equations. Now, the uncertainty over unobserved noise variables ϵ induces uncertainty over post-update states of all tuples t' captured by their post-update distribution on the possible worlds (Definition 1): $\Pr_{D,U}(\tau)$ for $\tau \in \text{PWD}(t')$, and in turn, the post-update distribution of the entire database $\Pr_{D,U}(I)$ for $I \in \text{PWD}(D)$. As we will show in Section 3.3, to answer what-if and how-to queries in HYPER, it suffices to estimate the post-update conditional distributions of the form $\Pr_{D,U}(Y = y \mid B = b, C = c)$, where $Y, B, C \in \text{Attr}(R)$, that measures the probabilistic influence of the update U on subset of tuples for which $B = b$ and $C = c$. It is known that if C satisfies a graphical criterion called *backdoor-criterion* (see Section 3.3) w.r.t. B and Y in the causal model G , then the following holds:

$$\Pr_{D,U}(Y = y \mid B = b, C = c) = \Pr_D(Y = y \mid B = f(b), C = c) \quad (1)$$

Where, the RHS of (1) can be estimated from D using standard techniques in causal inference and Machine Learning. Equation (1) also extends to multi-relation databases (see [22]).

Background knowledge on causal DAG. While in this paper we assume the underlying causal model is available, HYPER is designed to work with any level of background knowledge. If the causal DAG is not available, HYPER assumes a canonical causal model in which all attributes affect both the output and the updated attribute. In other words, HYPER assumes (1) holds for $C = \text{Attr}(R)$, i.e., all attributes are considered in the backdoor set in Equation 1, ensuring that the ground truth backdoor set is a subset of $\text{Attr}(R)$. We also examine this case experimentally in Section 5.

3 PROBABILISTIC WHAT-IF QUERIES

In this section we describe the syntax of probabilistic what-if queries supported by HYPER (Section 3.1), describe their semantics as expected value from the post-update distribution on possible worlds (Section 3.2), and present efficient algorithms and optimizations to compute the answers to what-if queries (Section 3.3).

3.1 Syntax of Probabilistic What-If Queries

A what-if query has two parts (see Figure 4):

- The required **USE** operator in the first part defines a single table as the **relevant view** with relevant attributes including the update and the output attribute to be used in the second part. The


```

USE RelevantView As
  (SELECT T1.PID,T1.Category,T1.Price,T1.Brand,
    AVG(Sentiment)As Sent,AVG(T2.Rating)As Rtn
  FROM Product As T1,Review As T2
  WHERE T1.PID = T2.PID
  GROUPBY T1.PID,T1.Category,T1.Price,T1.Brand)
WHEN Brand = 'Asus'
UPDATE(Price) = 1.1 × PRE(Price)
OUTPUT AVG(POST(Rtn))
FORPRE (Category) = 'Laptop' ANDPRE (Brand) = 'Asus'
AND POST(Sent) > 0.5

```

Figure 4: What-if query asking “If the prices of all Asus products is increased by 10%, what would the effect on average ratings of Asus laptops having average sentiments in the reviews > 0.5 after the update?”

USE operator can simply mention the table name if no transformation is needed, and both update and output attributes belong to this table (e.g., ‘USE Review’). Otherwise, a standard SQL query within the USE operator can define this relevant view as discussed below.

- The second part includes the new operators for hypothetical what-if queries supported by HYPER: the required UPDATE and OUTPUT operators for specifying the update and outcome attribute from the relevant view, and optional WHEN and FOR operators.

The second part takes as input the relevant view, denoted \mathcal{V}^{rel} (named as *RelevantView* in Figure 4), as defined by the required **USE operator** in the first part containing all relevant attributes, and therefore does not mention any table name for disambiguation in its operators. Recall that a hypothetical update in HYPER is of the form $U = u_{R,B,f,S}$, where the updated attribute $B \in \text{Attr}(R)$ in D , and is changed for all tuples $t \in S$ in R according to the function f (Definition 2). In the what-if query, the relevant view \mathcal{V}^{rel} defined by the first part combines the update and outcome attributes (Price and Rating in Figure 4) along with other attributes used in the second part. In particular, the SQL query defining \mathcal{V}^{rel} includes the update attribute B in the SELECT clause along with the key of R (here PID), and other attributes from R and (in aggregated form) from other relations in D that are used in the second part of the query. A group-by is performed on the attributes coming from relation R . Note that the first part always outputs a view having the same number of tuples as in R , which is ensured as the SELECT and Group By clauses include the key of R .

The required **UPDATE operator** mentions the update attribute B along with the function f . HYPER allows hypothetical update functions f of the form $\text{Update}(B) = \langle \text{const} \rangle$, $\text{Update}(B) = \langle \text{const} \rangle \times \text{PRE}(B)$, and $\text{Update}(B) = \langle \text{const} \rangle + \text{PRE}(B)$, where $\langle \text{const} \rangle$ is a constant specified by the user (here 1.1 models a 10% price increase). **PRE**(A) and **POST**(A) respectively denote the value of an attribute A before the hypothetical update (i.e., as given in the database instance D) and after the update according to the PRCM (see Sections 2.2 and 3.2); except in the operator as ‘UPDATE(B)’ which defines updating the value of B , **PRE** is assumed by default if **PRE** or **POST** is not explicitly mentioned in the query. **UPDATE** is always performed w.r.t. the **PRE** value of an attribute, rather than

the **POST** value which is the result of the update. The optional SQL query in the USE operator defining the relevant view can only have **PRE** values of attributes, so **PRE** is omitted in the query. Note that for immutable attributes A , $\text{PRE}(A) = \text{POST}(A)$.

The optional **WHEN operator** specifies the set S in Definition 2; any valid SQL predicate can be used here that is defined for each tuple in the relevant view \mathcal{V}^{rel} , and allows selection of a subset of tuples from \mathcal{V}^{rel} , e.g., $A = \langle \text{const} \rangle$, $A \in (\text{SELECT} \dots \text{AS } A \dots)$ etc. If the WHEN operator is not specified, we assume $S = R$ and the hypothetical update is applied to all tuples in R . Since the update is applied to the original attribute values, it can only use **PRE**(A) value for an attribute A , and therefore **PRE** is omitted.

The required **OUTPUT operator** mentions the output attribute Y (here, Rtn) on which we want to measure the effect of the hypothetical update. If Y belongs to another table $R' \neq R$, the SQL query in the USE operator describes how R and R' are combined in the join condition, and a SQL aggregate operator aggr_1 (SUM , AVG , COUNT) is used to aggregate Y (here $\text{AVG}(T2.Rating)$) to have a unique value for each tuple in R identified by its key in the relevant view. Note that the effect of an update is output as a single value, so another SQL aggregate operator aggr is used in the OUTPUT operator (here again AVG). If the user wants to measure effects on different subsets of tuples, it can be achieved by the use of the optional FOR operator described below. The OUTPUT operator can only use **POST**(A) values of attributes after the update.

The output specified in the OUTPUT operator is computed only considering the tuples in the relevant view \mathcal{V}^{rel} that satisfy the conditions in the optional **FOR operator** (details in Section 3.2). If no FOR operator is provided, all tuples in \mathcal{V}^{rel} are used to compute the output. FOR can contain both **PRE**(A) and **POST**(A) values of attributes, and **PRE** can be optionally provided for clarity. Further, like WHEN, any valid SQL predicate can be used that is defined on individual tuples in relevant view \mathcal{V}^{rel} .

EXAMPLE 6. Consider the what-if query statement shown in Figure 4. It checks the effect of hypothetically updating the price by 10% (UPDATE) on Brand = ‘Asus’ (WHEN). The effect is measured on their average of average ratings (OUTPUT) – the first average on ratings of the same type of Asus products, and the second average is on different types of Asus products, but only for Category = ‘Laptop’ (i.e., does not include phones for instance), and where the post-update average sentiment is still above 0.5. Since Rating and Sentiment come from the Review table whereas the update attribute Price belongs to the Product table, they are aggregated in the SQL query in the USE operator for each Product tuple.

HYPER supports multiple updates in a what-if query with attributes B_1, B_2, \dots , e.g., $\text{UPDATE}(Price) = 500 \text{ AND } \text{UPDATE}(Color) = \text{Red}$, provided there are no paths from any $B_i[t]$ to any $B_j[t']$ for any two tuples t, t' – a fact that we will use in Section 4 for how-to queries; we discuss other extensions in Section 7. Here, we discuss single-attribute updates for simplicity.

3.2 Semantics of Probabilistic What-If Queries

Here we define the semantics of what-if queries described in Section 3.1 as the *expected* value of the output attribute over possible worlds consistent with a what-if queries.

The operators in the what-if queries are evaluated in this order: $\text{USE} \rightarrow \text{WHEN} \rightarrow \text{UPDATE} \rightarrow \text{FOR} \rightarrow \text{OUTPUT}$.

(1) The USE operator outputs the relevant view \mathcal{V}^{rel} that contains all relevant attributes for the what-if query by a standard group-by SQL query.

(2) The WHEN operator takes \mathcal{V}^{rel} as input, and defines the set S in the update $U = u_{R,B,f,S}$. Suppose this operator uses an SQL predicate μ_{WHEN} defined on a subset of attributes of \mathcal{V}^{rel} . Then the output of the WHEN operator is the view $\mathcal{V}^{rel}_w = \{t \in \mathcal{V}^{rel} : \mu_{\text{WHEN}}(t) = \text{true}\}$. Note that in both USE and WHEN operators, the pre-update values (PRE values are assumed by default) from the given database D are used.

(3) Then the 'UPDATE $B = f(\text{PRE}(B))$ ' operation is applied to the tuples $t \in \mathcal{V}^{rel}_w$ on attribute B . As described in Section 2.2, this update is equivalent to modifying the structural equation $\phi_B[t]$ in the PRCM by replacing them with a constant value $f(\text{PRE}(B))$. Due to uncertainty induced by the noise variables, at this point, we get a set of possible worlds $PWD(D)$ (Definition 1) along with a post-update distribution $\text{Pr}_{D,U}$ on $PWD(D)$ induced by the update U . Clearly, some possible worlds I have $\text{Pr}_{D,U}(I) = 0$, e.g., if for a tuple t in relation R of I such that t corresponds to a tuple in \mathcal{V}^{rel}_w with the same key, $B[t] \neq f(\text{PRE}(B[t]))$.

(4 and 5) For the FOR and OUTPUT operators, for a possible world $I \in PWD(D)$ obtained from the previous step. Let \mathcal{V}_I^{rel} be the output of the SQL query in the USE operator on I . Suppose the predicate in the FOR operator is μ_{FOR} , which may include $\text{PRE}(A)$ and $\text{POST}(A')$ values for different attributes A, A' . For every tuple t (in any relation in D) and attribute A , consider two values of $A[t]$: $\text{PRE}(A[t])$ of t in D and $\text{POST}(A[t])$ of t in I (some values remain the same in PRE and POST, e.g., if A is immutable or if there is no effect of updating B for S tuples on A). Using these values, we evaluate the predicate μ_{FOR} , and using tuples from R that satisfy this predicate, we compute the aggregate aggr_Q ($\text{Avg}(\text{Rating})$) in Figure 4) mentioned in the OUTPUT operator using their values in I (i.e., POST values).

This aggregate aggr_Q is computed on attribute values $Y[t]$ for $t \in \mathcal{V}_I^{rel}$, where Y itself can be an aggregated attribute $Y = \text{aggr}_{\text{USE}}(Y')$ if it is coming from a different relation than the one containing the update attribute as defined by the SQL query in the USE operator (in Figure 4, $Y = \text{Rtnng}$, $\text{Rtnng} = \text{Avg}(\text{Review.Rating})$), and both aggr_Q and aggr_{USE} are Avg). Hence, when a possible world is fixed, the what-if query answer is computed as follows:

DEFINITION 4 (WHAT-IF QUERY RESULT ON A POSSIBLE WORLD). Given a what-if query Q and a database D , the answer to Q on a given possible world $I \in PWD(D)$ is the aggregate aggr_Q over $Y_I[t]$ values using the notations above:

$$\text{val}_{\text{whatif}}(Q, D, I) = \text{aggr}(\{Y_I[t] : \mu_{\text{FOR}}(t) = \text{true}, t \in \mathcal{V}^{rel}\}) \quad (2)$$

where $Y_I[t]$ denotes the value of attribute Y for tuple t in the possible world I . Here t is tuple in the relevant view \mathcal{V}^{rel} and therefore corresponds to a unique tuple in relation R .

Then the final value of the what-if query is the expected query result on all possible worlds of D :

DEFINITION 5 (WHAT-IF QUERY RESULT). Given a what-if query Q and a database D , the result of $Q(D)$ is the expected value of

$\text{val}_{\text{whatif}}(Q, D, I)$ over all possible worlds $I \in PWD(D)$, using the post-update probability distribution $\text{Pr}_{D,U}$:

$$\begin{aligned} \text{val}_{\text{whatif}}(Q, D) &= \mathbb{E}_{I \in PWD(D)} [\text{val}_{\text{whatif}}(Q, D, I)] \\ &= \sum_{I \in PWD(D)} \text{val}_{\text{whatif}}(Q, D, I) \cdot \text{Pr}_{D,U}(I) \end{aligned} \quad (3)$$

3.3 Computation of What-If Queries

The semantics presented in Section 3.2 does not directly lead to an efficient algorithm to compute the answer to what-if queries by Definition 5, since (1) the number of possible worlds can be exponential in the size of the database D , and (2) computation of post-update distribution $\text{Pr}_{D,U}$ is non-trivial. In this section, we present our algorithm for computing what-if query answers that use two key ideas to address these challenges: (a) Instead of computing the what-if query over the entire database, we decompose it into smaller problems and compute modified queries on subsets of tuples that are 'independent' of each other (as fewer tuples make the computation more efficient). Then we combine the results to get the result of the original query over the entire database. (b) To compute the distribution $\text{Pr}_{D,U}$ needed for estimating the query result, we use techniques from the *observational causal inference* and the *graphical causal model* literature [38] when the post-update distribution is determined by a PRCM.

Decomposing the computation. The decomposition, and subsequently the composition of answers, is achieved by the use of *block-independent databases* and *decomposable aggregate functions* supported by HYPER (SUM, COUNT, AVERAGE) described below.

Block-independent database decomposition. We adapt the notion of block-independent database model that has been used in *probabilistic databases* [14, 41] and *hypothetical reasoning* [30]. First, we need the notion of independence in our context. We say that two tuples $t, t' \in D$ are **independent** if there are no paths in the ground causal graph G (ref. Section 2.2) between $A[t]$ and $A'[t']$ for any two attributes A, A' .

Given a database D and a PRCM with a ground causal graph G , $\mathcal{B} = \{D_1, \dots, D_\ell\}$ is called a **block-independent decomposition of D** if (i) $\{D_1, \dots, D_\ell\}$ forms a partition of D , i.e., each $D_i \subseteq D$, $\bigcup_{i=1}^\ell D_i = D$, and $D_i \cap D_j = \emptyset$ for $i \neq j$, and (ii) for each $t \in D_i$ and $t' \in D_j$ where $i \neq j$, t and t' are independent. Note that these tuples t and t' can come from the same or different relations of D .

We compute block-independent decomposition of database D given a causal graph G as follows. The block decomposition process performs a topological ordering of the nodes in the causal graph and then performing a DFS or BFS on it, and is therefore linear in the size of the causal DAG. The causal DAG has at most $n \times k$ nodes where n is the number of tuples in D and $k = |\text{Attr}(D)|$. In particular, the decomposition does not depend on the structure or complexity of the query. Block-independent decomposition provides an optimization in our algorithms; in the worst case, all tuples may be included in a single block.

EXAMPLE 7. Consider the causal graph of the PRCM (Figure 3) defined on the database presented in Figure 1. The procedure first performs a topological sort of the nodes. For example, in Figure 2, the node $\text{Brand}[p_1]$ is first, and then the node $\text{Quality}[p_1]$ etc. Then, the algorithm performs a BFS to detect the connected components

of the graph which are all tuples belonging to the same category, along with their reviews. The block-independent decomposition of the database D in Figure 1 is then $\mathcal{B} = \{D_1, D_2, D_3\}$ where $D_1 = \{p_1, p_2, p_3, r_1, r_2, r_3, r_4, r_5\}$, $D_2 = \{p_4, r_6\}$, and $D_3 = \{p_5\}$ corresponding to laptops, camera, and books along with their reviews.

Decomposable functions. The aggregate functions supported by HYPeR are *decomposable* as defined below, which allows us to combine results from each block after a block-independent decomposition to compute the answer to a what-if query. Since the immutable attributes include keys that are unchanged in all possible worlds $I \in PWD(D)$ of D (Definition 1), given a block-independent decomposition \mathcal{B} of D , we will use the corresponding decomposition \mathcal{B}_I of I where the same tuples identified by their keys go to the same blocks in \mathcal{B} and \mathcal{B}_I . The aggregate functions $f_{Q,D}, f'_{Q,D}$ below map a set of tuples to a real number whereas g maps a set of real numbers to another real number.

DEFINITION 6 (DECOMPOSABLE AGGREGATE FUNCTION). Given a database D , a block-independent decomposition $\mathcal{B} = \{D_1, \dots, D_\ell\}$ of D , a what-if query Q , and any possible world $I \in PWD(D)$ of D , an aggregate function $f_{Q,D}$ is **decomposable** if there exist aggregate functions $f'_{Q,D}$ and g such that:

- $f_{Q,D}(I) = g(\{f'_{Q,D}(D_i) \mid \forall D_i \in \mathcal{B}_I\})$ where \mathcal{B}_I is the block partition of I corresponding to \mathcal{B} ,
- $\alpha g(\{x_1, \dots, x_l\}) = g(\{\alpha x_1, \dots, \alpha x_l\})$, $\forall \alpha \geq 0$, and
- $g(\{x_1, \dots, x_l\}) + g(\{y_1, \dots, y_l\}) = g(\{x_1 + y_1, \dots, x_l + y_l\})$

When the aggregate function $aggr$ given in Equation (2): $\text{val}_{\text{whatif}}(Q, D, I) = \text{aggr}(\{Y_I[t] : \mu_{\text{FOR}}(t) = \text{true}, t \in \mathcal{V}^{\text{rel}}\})$ is decomposable, we show that the computation can be performed on the blocks \mathcal{B}_I and then aggregated to compute $\text{val}_{\text{whatif}}(Q, D, I)$. We note that every supported aggregate function in this paper (SUM, AVG, COUNT) is decomposable. We demonstrate this for AVG below.

EXAMPLE 8. Reconsider the what-if query in Figure 4. Suppose the database can be partitioned into blocks by Category as demonstrated in Example 7. In this case, $aggr = \text{AVG}$ and $Y = \text{Rtngr} = \text{AVG}(\text{T2.Rating})$, and for any $I \in PWD(D)$, $\text{val}_{\text{whatif}}(Q, D, I) = \text{AVG}(\{\text{Rtngr}_I[t] \mid t \in \mathcal{V}^{\text{rel}}, \text{Category}[t] = \text{Laptop}, \text{Brand}[t] = \text{Asus}, \text{Post}(\text{Senti}[t]) > 0.5\})$. We use the standard formula for decomposing average: $\text{AVG}(D) = \frac{1}{|D|} \sum_{i=1}^{\ell} \text{SUM}(D_i)$. For each block $D_i \in \mathcal{B}_I$, $f'_{Q,D}(D_i) = \frac{1}{|D_i|} \text{SUM}(\{\text{Rtngr}_I[t] \mid t \in \mathcal{V}_I^{\text{rel}} \cap D_i, \text{Category}[t] = \text{Laptop}, \text{Brand}[t] = \text{Asus}, \text{Post}(\text{Senti}[t]) > 0.5\})$. Here, $g = \text{SUM}$, and SUM satisfies the properties in Definition 6.

In the proof of the following proposition, we leverage the ability to marginalize the distribution $\text{Pr}_{D,U}$ over the possible worlds of the database D (Definition 3) given a what-if query Q to get a distribution and a set of possible worlds for any block $D_i \in \mathcal{B}$, which we denote by $PWD(D_i) \subseteq PWD(D)$. $PWD(D_i)$ are all instances where all tuples $t' \notin D_i$ remain unchanged and all mutable attributes of $t \in D_i$ get all possible values from their respective domains. We further denote $PWD(D_i)$ as the set of possible worlds of D_i that only includes the tuples in D_i ; i.e., $PWD(D_i)$ is the projection of $PWD(D)$ on D_i . All proofs are deferred to the full version [22] due to space constraints.

PROPOSITION 1 (DECOMPOSED COMPUTATION). Given a database D , its block-independent decomposition $\mathcal{B} = \{D_1, \dots, D_\ell\}$, and a what-if query Q whose result on a possible world $I \in PWD(D)$ is $\text{val}_{\text{whatif}}(Q, D, I) = \text{aggr}(\{Y_I[t] : \mu_{\text{FOR}}(t) = \text{true}, t \in \mathcal{V}^{\text{rel}}\})$ (Definition 4), if $aggr$ is a decomposable function, i.e., if there exist functions g and $f'_{Q,D}$ according to Definition 6, then

$$\text{val}_{\text{whatif}}(Q, D) = g(\{\text{val}_{\text{whatif}}(Q', D_i) \mid \forall D_i \in \mathcal{B}\}) \quad (4)$$

where Q' is the same query as Q with $f'_{Q,D}$ replacing $aggr$ and

$$\text{val}_{\text{whatif}}(Q', D_i) = \mathbb{E}_{I_j \in \overline{PWD}(D_i)} [\text{val}_{\text{whatif}}(Q', D_i, I_j)] \quad (5)$$

Computing results with causal inference. We show the connection between the what-if query results and techniques in observational causal inference. This connection will allow us to compute the results for each block as given in Equation (5). Specifically, we show how the computation in each block is done by the post-update probabilities, which we further reduce to pre-update probabilities.

PROPOSITION 2 (CONNECTION TO CAUSAL INFERENCE FOR COUNT). Given a database D with its block independent decomposition \mathcal{B}_D , a block $D_i \in \mathcal{B}_D$, a ground causal graph G , a what-if query Q' where $\text{Agg} = \text{COUNT}$, and the FOR operator is denoted by $\mu_{\text{FOR}} = \vee_k (\mu_{\text{FOR,PRE}}^k \wedge \mu_{\text{FOR,POST}}^k)$, and $\mu_{\text{FOR,PRE}}$ and $\mu_{\text{FOR,POST}}$ are defined with the PRE and POST operators, respectively, the following holds.

$$\text{val}_{\text{whatif}}(Q', D_i) = \sum_{t \in D_i} \left(\sum_k \left(\text{Pr}_{D_i, U}(\mu_{\text{FOR,POST}}^k(t) = \text{true} \mid \mu_{\text{FOR,PRE}}^k(t) = \text{true}) \right) \right)$$

In this equation, $\text{Pr}_{D_i, U}(\mu_{\text{FOR,POST}}^k(t) = \text{true} \mid \mu_{\text{FOR,PRE}}^k(t) = \text{true})$ denotes the sum of probabilities of all possible worlds of D_i such that the tuple t that satisfied $\mu_{\text{FOR,PRE}}^k(t) = \text{true}$ before the update U also satisfies $\mu_{\text{FOR,POST}}^k(t)$ after the update.

The proof of the proposition relies on the fact that a FOR clause can be represented as a CNF of PRE and POST conditions (as shown in [22]). Proposition 2 assumes $\text{Agg} = \text{COUNT}$, however, a similar result for $\text{Agg} = \text{SUM}/\text{AVG}$ can be found in the full version [22].

Estimating the probability values. The expression in Proposition 2 relies on the post-update distribution to evaluate conditional probability of certain attribute values. For example, we need a way to estimate $\text{Pr}_{D, U}(A_i = a_i \mid A_j = a_j, \mu_{\text{WHEN}})$ when $aggr = \text{COUNT}$. Our goal is to find a way to estimate these probability values from the input database D , assuming we have a PRCM.

To do so, we leverage the notion of **backdoor criterion** from causal inference [38]. A set of attributes C satisfies the backdoor criterion w.r.t. A_i and B if no attribute $C \in C$ is a descendant of A_i or B and all paths from B to A_i which contain an incoming edge into A_i are blocked by C . For example, in Figure 3, $\text{Brand}[p_1]$, $\text{Quality}[p_1]$, and $\text{Category}[p_1]$ satisfy the backdoor criterion with respect to $\text{Sentiment}[p_1]$ and $\text{Rating}[p_1]$. Using this criterion, we show (in the full version) that the element $\text{Pr}_{D, U}(A_i = a_i \mid B = b, C = c, A_j = a_j, \mu_{\text{WHEN}})$ in the query result expression in Proposition 2 can be estimated from Pr_D using the following calculations.

$$\text{Pr}_{D, U}(A_i = a_i \mid A_j = a_j, \mu_{\text{WHEN}}) =$$

$$\sum_{c \in \text{Dom}(C)} \text{Pr}_{D, U}(A_i = a_i \mid C = c, A_j = a_j, \mu_{\text{WHEN}}) \text{Pr}_D(C = c \mid A_j = a_j, \mu_{\text{WHEN}})$$

The first probability term can be simplified as:

$$\text{Pr}_{D, U}(A_i = a_i \mid C = c, A_j = a_j, \mu_{\text{WHEN}}) = \sum_{b \in \text{Dom}(B)} \text{Pr}_{D, U}(A_i = a_i \mid B = b, C = c, A_j = a_j, \mu_{\text{WHEN}}) \cdot \text{Pr}_D(B = b \mid C = c, A_j = a_j, \mu_{\text{WHEN}}).$$

```

Use (...) /* same as Figure 4 */
When Brand = 'Asus' AND Category = 'Laptop'
HowToUpdate Price, Color
Limit 500 ≤ Post(Price) ≤ 800 AND
    L1(PRE(Price), POST(Price)) ≤ 400
ToMaximize AVG(POST(Rtng))
For (PRE(Category) = 'Laptop' Or
    PRE(Category) = 'DSLR Camera') AND Brand = 'Asus'

```

Figure 5: How-to query asking “how to maximize the average rating of Asus laptops and cameras over the determined view by changing the price and/or color of Asus laptops so that it will not drop below 500 and increase above 800, and will be at most 400 away from its original value?”

This shows that the query output relies on $\Pr_{D,U}(A_i = a_i \mid B = b, C = c, A_j = a_j, \mu_{\text{WHEN}})$, which can be estimated from \Pr_D using equation (1). Using these probability calculations, we estimate the query output from the input data distribution \Pr_D . The equations require that we iterate over the values in the domain of \mathbf{B} and \mathbf{C} , which can be inefficient as the domain set size increases exponentially with the number of attributes in the set. However, the majority of the values in $\text{Dom}(\mathbf{C})$ would have zero-support in the database D , implying $\Pr_D(C = c \mid A_j = a_j, \mu_{\text{WHEN}}) = 0$ for $C = c$. Therefore, we build an index of values in $\text{Dom}(\mathbf{C})$ to efficiently identify the set of values that would generate a positive probability-value. This optimization ensures that the runtime is linear in the database size.

4 PROBABILISTIC HOW-TO QUERIES

How-to queries support *reverse data management* (e.g., [34]), and suggest how a given mutable attribute can be updated to optimize the output attributes subject to various constraints. In this section we describe the syntax of probabilistic how-to queries supported by HYPER (Section 4.1), describe their semantics (Section 4.2), and present algorithms to compute their answers (Section 4.3). How-to queries are computed by solving an optimization problem over several relevant what-if queries.

4.1 Syntax of Probabilistic How-To Queries

The syntax of how-to queries in HYPER is similar to that of what-if queries (see Figures 4 and 5, and Section 3.1). How-to queries have two parts. The first part uses the required **Use** operator and is identical to the **USE** operator in the what-if queries in its functionality – it defines the relevant view \mathcal{V}^{rel} that contains the key of the relation containing the update attribute, and includes all attributes used in the second part of the query; attributes coming from other relations are aggregated.

The optional **WHEN** and **FOR** operators have the same functions as in what-if queries. The **WHEN** operator specifies the set S on which an update $U = u_{R,B,f,S}$ can be applied, whereas the **FOR** operator defines the subset on which the effect is estimated. Like what-if queries, **WHEN** only includes pre-update values $\text{PRE}(A)$, whereas **FOR** can include both pre- and post-update values $\text{PRE}(A), \text{POST}(A)$.

The required **HowToUpdate** operator corresponds to the **UPDATE** operator of what-if queries, and uses $\text{PRE}(A)$, but instead

of specifying an attribute (or a set of attributes) to update, it specifies the set of mutable attributes that can be updated. In Figure 5, ‘HowToUpdate Price, Color’ states that any combination of these three attributes can be updated, and some attributes can be left unchanged as well. To ensure that the updates on these attributes are valid, our algorithms assume that, for any pair of the attributes mentioned in this clause A_1, A_2 , there are no paths in the ground causal graph of the PRCM between $A_1[t]$ and $A_2[t']$ for any $t, t' \in D$.

Possible **outputs of the how-to queries** are of these forms for each attribute A specified in the **HowToUpdate** operator: (i) $\text{UPDATE}(B) = \langle \text{const} \rangle$, (ii) $\text{UPDATE}(B) = \langle \text{const} \rangle \times \text{PRE}(B)$, (iii) $\text{UPDATE}(B) = \langle \text{const} \rangle + \text{PRE}(B)$, and $\text{UPDATE}(B) = \text{no change}$, where $\langle \text{const} \rangle$ is a constant found by our algorithms from the search space. One example output of this **HowToUpdate** query is

{Price: 1.1x, Color: no change}

stating the price should be increased by 10%, the color should be changed to red, and the category should not be changed.

The optional **Limit** operator states the constraints for optimization, i.e., it defines the conditions that restrict the post-update values of update attributes specified in the **HowToUpdate** operator for tuples in \mathcal{V}^{rel} that satisfy the **WHEN** operator. In particular, if an attribute A is numeric, its updates can be bounded by numeric limits, e.g., $l \leq \text{POST}(A) \leq h$, $\text{POST}(A) \leq \text{PRE}(A) + \langle \text{const} \rangle$, $\text{POST}(A) \leq \text{PRE}(A) \times \langle \text{const} \rangle$, etc., and if A is categorical or numeric, the user can specify the permissible values as a set, e.g., $\text{POST}(A) \text{ IN } (v_1, v_2, v_3)$. Furthermore, this operator allows users to specify the maximal or minimal $L1$ distance between the original attribute values ($\text{PRE}(A)$) and the updated ones ($\text{POST}(A)$) for attributes A in the **HowToUpdate** operator for the tuples satisfying the condition in the **WHEN** operator: $L1(\text{POST}(A), \text{PRE}(A))$ takes a vector of values V_u and $V_u[i]$ is an update value of the i ’th attribute mentioned in the **LIMIT** operator, and returns the normalized $L1$ distance between the original value vector and the vector of update values $|V_u - V_{orig}|$. The $L1$ operator helps model the **cost of an update** (with suitable weights) as some updates can be more expensive than the others.

Finally, the how-to query needs to include a required **ToMaximize** or **ToMinimize** operator, which specifies an aggregated value of an attribute from the relevant view \mathcal{V}^{rel} that is to be maximized or minimized using the updates on the attributes specified in the **HowToUpdate** operator. Only post-update values $\text{POST}(A)$ of attributes are allowed in **ToMaximize** and **ToMinimize**.

EXAMPLE9. Consider the query in Figure 5. It asks for the maximum value of the average value of *Rtng* (**HowToUpdate**) by updating the tuples with *Brand* = ‘Asus’, *Category* = ‘Laptop’ (**WHEN**). The attributes allowed to be updated are *Price*, *Color* (**HowToUpdate**). The update to the *Price* attribute is restricted to [500, 800], where distance between the original values and the updated values in this attribute has to be ≤ 400 . The average of *Rtng* is computed over the view defined by the **FOR** operator.

4.2 Semantics of Probabilistic How-To Queries

We next define the results of how-to queries in terms of what-if queries. Intuitively, every how-to query optimizes over a set of what-if queries, where each what-if query contains a possible update allowed in the how-to query. Assuming, without losing generality,

that the how-to query contains a ToMAXIMIZE operator, the result of the how-to query is then the what-if query that yields the maximum result of the output attribute in the ToMAXIMIZE operator of the how-to query, subject to the constraints on post-update values of attributes specified in the Limit operator.

DEFINITION 7 (CANDIDATE WHAT-IF QUERY). Given a how-to query Q_{HT} that includes (i) a ToMAXIMIZE operator of $\text{Agg}(\text{Post}(Y))$, (ii) a HowToUpdate operator with update attributes B_1, \dots, B_c , and (iii) a LIMIT operator that specifies permissible ranges \mathcal{R}_i and $L1(\text{PRE}(B_i), \text{POST}(B_i)) < \theta_i, \forall i \in [1, c]$ (if there are no constraints on the range for B_i , $\mathcal{R}_i = \text{Dom}(B_i)$ and if no L1 constraint is specified, $\theta_i = \infty$), a **candidate what-if query** is a what-if query Q_{WI} such that:

- The USE, WHEN, and FOR operators in Q_{WI} are identical to the ones in Q_{HT} ,
- Q_{WI} contains $\text{UPDATE } B_{j_1} = b_1, \dots, B_{j_k} = b_k$, where $\{j_1, \dots, j_k\} \subseteq \{1, \dots, c\}$, $b_i \in \mathcal{R}_{j_i}$, and $L1(\text{PRE}(B_{j_i}), \text{POST}(B_{j_i})) < \theta_{j_i}$.
- The OUTPUT operator in Q_{WI} specifies the attribute $\text{Agg}(\text{Post}(Y))$ from the ToMAXIMIZE operator in Q_{HT} .

This query is denoted as $Q_{WI}((B_{i_1}, b_1), \dots, (B_{i_c}, b_c))$. The **set of all candidate what-if queries** for a how-to query Q_{HT} is denoted by $Q_{\text{whatif}}(Q_{HT})$.

EXAMPLE 10. A candidate what-if query $Q_{WT}(\text{Price}, 500)$ for the how-to query depicted in Figure 5 is given below (USE operator is the same as that in Figure 4):

```
USE (...)
WHEN Brand = 'Asus' AND Category = 'Laptop'
UPDATE Price = 500
OUTPUT AVG(Post(Rating))
FOR (PRE(Category) = 'Laptop' OR
     PRE(Category) = 'DSLR Camera') AND Brand = 'Asus'
```

In particular, the update on the Price attribute is in $[500, 800]$ and satisfies the L1 distance since the original price of the Asus laptop is 529, and the rest of the query is identical to the query in Figure 5.

We now define the result of a how-to query that optimizes over the result of all candidate what-if queries.

DEFINITION 8 (HOW-TO QUERY RESULT). Given a database D and a how-to query Q_{HT} with a ToMAXIMIZE operator, the result of Q_{HT} is defined as $\text{argmax}_{Q_{WI} \in Q_{\text{whatif}}(Q_{HT})} \text{val}_{\text{whatif}}(Q_{WI}, D)$, where $\text{val}_{\text{whatif}}(Q_{WI}, D)$ denotes the result of the what-if query Q_{WI} on D as defined in Definition 5; ToMinimize is defined similarly.

We take the argmax of $Q_{\text{whatif}}(Q_{HT})$ since a how-to query asks about the manner in which the database needs to be updated and not about the result. This corresponds to the output we defined and demonstrated in Section 4.1. Definition 8 requires taking the maximum over a large set of candidate what-if queries, which can even be infinite if the domain is continuous. In the next section, we provide optimizations to make their computation feasible.

4.3 Computation of How-to queries

The naive approach to compute the result of a how-to query by Definition 8 is inefficient as it evaluates a large number of candidate what-if queries. Instead, we model the problem of computing the result of how-to queries as an Integer Program (IP).

Denote by $U = \{B_1, \dots, B_c\}$ the set of update attributes in the HowToUpdate operator. For each attribute $B_i \in U$, we enumerate all permissible updates (denoted by S_{B_i}) and define an indicator variable δ_{b_i} for every b_i which denotes the potential updated value of attribute B_i . For example, S_{Price} can consist of the following updates: $\{1.1 \times \text{Pre}(\text{Price}), 1.2 \times \text{Pre}(\text{Price}), \dots, 2.5 \times \text{Pre}(\text{Price}), 100 + \text{Pre}(\text{Price}), 200 + \text{Pre}(\text{Price}), \dots, 500 + \text{Pre}(\text{Price}), 250, 300, \dots, 600\}$.

The elements of set S_A are defined such that all these updates satisfy the constraints mentioned in Limit operator. If the set of potential updates is continuous, we bucketize them so that we can treat their values as discrete. Given a set S_{B_i} and variables δ_{b_i} for all $b_i \in S_{B_i}$, we add a constraint for each attribute that $\sum_{b_i \in S_{B_i}} \delta_{b_i} \leq 1$ to ensure that at most one of the updates is performed. If δ_{b_i} is zero for all values in S_{B_i} , then B_i is not updated. Given this formulation, the corresponding what-if query is estimated as a linear expression by using Proposition 2 and training a regression function over the dataset D . Let this linear function be $\phi : \text{Dom}(U) \rightarrow O$, where O is the range of the output of candidate what-if queries. The following IP models the solution to the how-to query using the variables δ_{b_i} .

$$\text{argmax} \quad \phi(D, \sum_{b_1 \in S_{B_1}} \delta_{b_1} b_1, \dots, \sum_{b_c \in S_{B_c}} \delta_{b_c} b_c) \quad (6)$$

$$\text{subject to} \quad \sum_{b_i \in S_{B_i}} \delta_{b_i} \leq 1, \quad \forall i = 1 \text{ to } c \quad (7)$$

$$\delta_{b_i} \in \{0, 1\}, \quad \forall b_i \in S_{B_i}, \forall i = 1 \text{ to } c \quad (8)$$

In addition to these constraints, additional constraints are added to the IP based on the constraints in the Limit operator. Since all constraints and the objective function are linear equations, we leverage standard IP solvers to calculate the output of the HowToUpdate query³. Note that the number of constraints in the IP grows linearly with the number of attributes in U and the number of variables grows linearly in the number of possible updates for each attribute.

Extension to preferential multi-objective optimization. HYPER can be adapted to the settings where an user aims to optimize multiple objectives that are lexicographically ordered based on preference. Consider an ordered set of preferences p_1, \dots, p_t where each preference p_i is less important than p_j for $j < i$. In this case, we propose to solve IP iteratively as follows. First, we can solve the optimization problem for the first preference p_1 as described above, ignoring other preferences. In the second iteration, the identified objective value of the first objective is added as a constraint to maximize the second preference p_2 . In this way, all previously solved objectives are added as constraints while optimizing for a preference p_i . The solution to the last IP that optimizes for p_t where all other preferences are added as constraints is returned as the final solution to the preferential multi-objective optimization.

EXAMPLE 11. Consider the database in Figure 1 and a how-to query that aims to maximize the average ratings as a first priority and the average sentiment as a second priority. In the first IP, we will solve for the clause ToMAXIMIZE $\text{AVG}(\text{Post}(\text{Rtnng}))$, where Rtnng are the ratings. Suppose the maximum average rating we get is c . We then solve the IP for the clause ToMAXIMIZE $\text{AVG}(\text{Post}(\text{Sentiment}))$ and add the constraint that $(\text{AVG}(\text{Post}(\text{Rtnng})))$ will equal c .

³As an alternate formulation, our framework allows optimizing the cost (L1 distance between the original attribute and the updated value) while adding a constraint on the aggregated attribute. We discuss more details in the full version [22].

5 EXPERIMENTS

We evaluate the effectiveness of HYPER and its variants on various datasets and answer the following questions:

- (1) Do the results provided by HYPER make sense in real-world scenarios?
- (2) How does HYPER compare to other baselines for hypothetical reasoning when the ground truth is available?
- (3) How does the runtime of HYPER depend on query complexity and dataset properties like number of tuples, the causal graph structure, discretization of continuous attributes, and the number of attributes in different operators of the query?
- (4) How does combining a sampling approach with HYPER influences runtime performance and the quality of the results?

Our experimental study includes 5 datasets and 3 baselines that are either inspired by previous approaches or simulate the absence of a causal model. We provide a qualitative and quantitative evaluation of HYPER, showing that it gives logical results in real-world scenarios and achieves interactive performance in most cases.

Implementation and setup. We implemented the algorithms in Python. HYPER was run on a MacOS laptop with 16GB RAM and 2.3 GHz Dual-Core Intel Core i5 processor. We used random forest regressor [51] to estimate conditional probabilities.

5.1 Datasets and Baselines

We describe the datasets and baselines used in this section.

Datasets. The following datasets and causal models were used.

- The **Adult** income dataset [32] comprises demographic information of individuals along with their education, occupation, hours of work, annual income, etc. It is composed of a single table. We used the causal graph from prior studies [11].
- **German** dataset [20] contains details of bank account holders including demographic and financial information along with their credit risk. It composed of a single table and the causal graph was used from [11].
- **Amazon** dataset [28] is a relational database consisting of two types of tables, as described in Figure 1, and the causal graph is presented in Figure 2. We identified product brand from their description, used Spacy [2] for sentiment analysis of reviews and estimated quality score from expert blogs [1].
- **German-Syn** is a synthetically generated dataset using the same causal graph as German dataset [20]. It consists of a single table. We consider two different versions for our analysis, one with 20K records and the other with 1 million records.
- **Student-Syn** dataset contains two different tables (a) Student information consisting of their age, gender, country of origin and their attendance. (b) Student participation attributes like discussion points, assignment scores, announcements read, and overall grade. Each student was considered to enroll in 5 different courses and their overall grade is an average over respective courses. This data was generated keeping in mind the effect of attendance on class discussions, announcements and grade. The causal model has student age, gender and country of origin as the root nodes, which affect their attendance and other performance related attributes.

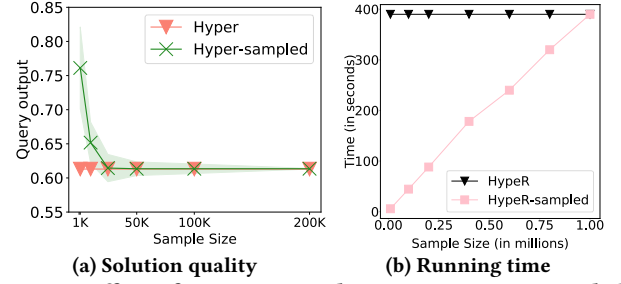


Figure 6: Effect of varying sample size on HYPER-sampled output and running time for German-Syn (1M) dataset

Variations. In the experiments, HYPER is run assuming that background knowledge about the causal graph is known a priori. Additionally, we consider the following variations.

- **HYPER-NB:** when no causal model is available, all attributes are assumed to affect the updated attribute and the output.
- **HYPER-sampled:** is an optimized version of HYPER that considers a randomly chosen subset of 100k records for the calculation of conditional probabilities of Proposition 2. The choice of sample size is discussed in Section 5.2

Baselines. We consider two different baselines of HYPER to evaluate hypothetical queries:

- **Indep** baseline inspired by previous work using provenance [16]: it ignores the causal graph and assumes that there is no dependency between different attributes and tuples.
- **Opt-HowTo** baseline for how-to analysis where we compute the optimal solution by enumerating all possible updates, evaluating what-if query output for each update and choosing the one that returns the optimal result.

5.2 HYPER and HYPER-sampled

We evaluate the effectiveness of HYPER with its variant HYPER-sampled to understand the tradeoff between quality and running time. Figure 6 compares the effect of changing the sample size on the quality of the output (Figure 6a) and running time (Figure 6b) by HYPER-sampled. Figure 6a shows that the standard deviation in query output of HYPER-sampled reduces with an increase in sample size and is within 1% of the mean whenever more than 100k samples are considered. In terms of running time, we observe a linear increase in time taken to calculate the output. Due to low variance of HYPER-sampled for 100k samples and reasonable running time, we consider 100k as the sample-size for subsequent analysis.

5.3 What-If Real World Use Cases

In this experiment, we evaluate the output of HYPER on different hypothetical queries on various real-world datasets. Due to the absence of ground-truth, we discuss the coherence of our observations with intuitions from existing literature.

German. We considered a hypothetical update offing attributes 'Status', 'Credit history', and 'housing' to their respective minimum and maximum values to evaluate the effect of these attributes on individual credit. Figure 7a demonstrates the query template, where X, x, X_2, x_2 are varied to evaluate the effect of different updates. Whenever status or credit history are updated to the maximum

USE D UPDATE(B) = b OUTPUTCOUNT ($Credit = Good$) FORPRE (A) = a

(a) What-if query (German dataset): What fraction of individuals will have good credit if B is updated to b ?

USE D UPDATE(B) = b OUTPUTCOUNT (*)
FOR POST($Income$) > 50k ANDPRE (A) = a

(b) What-if query (Adult dataset): How many individuals with attribute $A = a$ will have income $\geq 50K$ if B is updated to b ?

Figure 7: What-if queries for real world use cases

Table 1: Average Runtime in seconds for a COUNT what-if query. The time in (..) is by HYPER(-NB)-sampled, which has identical runtime as HYPER(-NB) on all other datasets with < 100k tuples.

Dataset	Att. [#]	Rows[#]	HYPER	HYPER-NB	Indep
Adult [32]	15	32k	45s	105s	3s
German [20]	21	1k	1.2s	12.5s	0.4s
Amazon [28]	5,3	3k, 55k	1.7s	10.5s	0.8s
Student-syn	3,6	10k, 50k	4.5s	12.3s	1.2s
German-Syn (20k)	6	20k	7.2s	22.45s	1.4s
German-Syn (1M)	6	1M	390s (44.5s)	1173s (132s)	73s

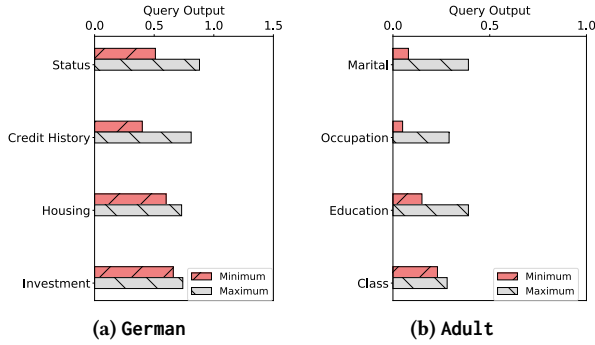


Figure 8: What-if query output for German and Adult datasets on updating each attribute to their min and max values; a larger gap denotes higher attribute importance.

value, more than 81% of the individuals have good credit. Similarly, updating these attributes to the minimum value reduces the credit rating of more than 30% individuals. On the other hand, updating other attributes like ‘housing’ and ‘investment’ affects the credit score of less than 20% individuals. Figure 8a presents the effect of updating these attributes to their minimum and maximum value. Larger gap in the query output for Status and credit history shows that these attributes have a higher impact on credit score. We also tested the effect of updating pairs of attributes and observed that **updating ‘credit history’ and ‘status’ at the same time can affect the credit score of more than 70% individuals**. These observations are consistent with our intuitions that credit history and account status have the maximum impact of individual credit.

Adult. This dataset has been widely studied in the fairness literature to understand the impact of individual’s gender on their income. It has a peculiar inconsistency where married individuals report total household income demonstrating a strong causal impact of marital status on their income [45, 50, 56]. We ran a hypothetical what-if query to analyze the fraction of high-income individuals when everyone is married (Figure 7b). We observed that 38% of the individuals have more than 50K salary. Similarly, **if all individuals were unmarried or divorced, less than 9% individuals**

have salary more than 50K. This wide gap in the fraction of high-income individuals for two different updates of marital status demonstrate its importance to predict household income. Figure 8b shows the effect of updating the attributes with the minimum or the maximum value in their domain. These observations match the observations of existing literature [23], where marital status, occupation and education have the highest influence on income.

Amazon. We evaluated the effect of changing price of products of different brands on their rating. When all products have price more than the 80th percentile, around 32% of the products have average rating of more than 4. **On further reducing the laptop prices to 60th and 40th percentiles, more than 60% of the products get an average rating of more than 4.** This shows that reducing laptop price increases average product ratings. Among different brands, we observed that Apple laptops have the maximum increase in rating on reducing laptop prices, followed by Dell, Toshiba, Acer and Asus. These observations are consistent with previous studies on laptop brands [3], which mention Apple as the top-quality brand in terms of quality, customer support, design, and innovation.

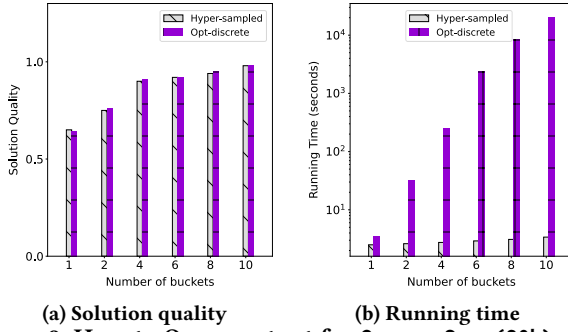
5.4 Solution Quality Comparison

In this experiment, we analyzed the quality of the solution generated by HYPER with respect to the ground truth and baselines over synthetic datasets. The ground truth values are calculated using the structural equations of the causal DAG for the synthetic data.

What-if. For the German-Syn (1M) dataset, Figure 10a presents the output of a query that updates different attributes related to individual income and evaluates the probability of achieving good credit. For all attributes, HYPER-sampled, and HYPER-NB estimate the query output accurately with an error margin of less than 5%. In contrast, Indep baseline ignores the causal structure and relies on correlation between attributes to evaluate the output. Since, the individuals with high status are highly correlated with good credit, Indep incorrectly outputs that updating Status would automatically improve credit for most of the individuals.

For the Student-Syn dataset, Figure 10b presents the average grade of individuals on updating different attributes that are an indicator of their academic performance. In all cases, HYPER-sampled and HYPER-NB output is accurate while Indep is confused by correlation between attributes and outputs noisy results. In addition to these hypothetical updates, we considered complex what-if queries that analyzed the effect of assignment and discussion attributes on individuals that read announcements and have high attendance. In these individuals, we observed that improving assignment score has the maximum effect on the overall grade of individuals.

How-to. For the German-Syn (20k) dataset, we considered a how-to query that aims to maximize the fraction of individuals receiving good credit. We provided Status, Savings, Housing and Credit amount as the set of attributes in the HowToUpdate operator. HYPER-sampled returned that updating two attributes i) account status, and ii) housing attributes is sufficient to achieve good credit. This showed that updating a single attribute would not maximize the fraction of individuals with good credit. We evaluated the ground truth (Opt-HowTo) by enumerating all possible update queries and used the structural equations of the causal graph to evaluate the



(a) Solution quality (b) Running time
Figure 9: How-to Query output for German-Syn (20k) with varying number of buckets.

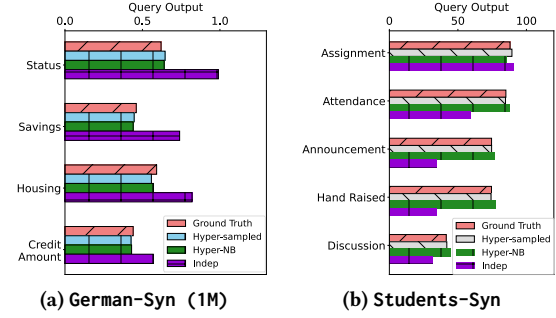
post-update value of the objective function for each update. We identified that HYPER’s output matches the ground truth update.

For the Student-Syn dataset, we evaluated a how-to query to maximize average grades of individuals with a budget of updating atmost one attribute. HYPER-sampled returned that improving individual attendance provide the maximum benefit in average grades. This output is consistent with ground truth calculated by evaluating the effect of all possible updates (Opt-HowTo).

Effect of discretization. HYPER bucketizes all continuous attributes before solving the integer program. In this experiment, we evaluate the effect of number of buckets on the solution quality and running time on a modified version of German-Syn (20k) dataset that contains continuous attributes. We partitioned the dataset into equiwidth buckets and compared the solution returned by HYPER-sampled and the optimal solution calculated after discretization (Opt-discrete) with the ground truth (OptHowTo). Figure 9a compares the quality of HYPER-sampled and Opt-discrete as a ratio of the optimal value. We observe that the solution quality improves with the increase in the number of buckets and the returned solution is within 10% of the optimal value whenever we consider more than 4 buckets. The time taken by Opt-discrete increases exponentially with the number of buckets, while the time taken by HYPER-sampled does not increase considerably as the number of variables in the integer program depends linearly on the number of buckets. This shows that running HYPER-sampled over a bucketized dataset leads to competitive quality in reasonable amount of time.

5.5 Runtime Analysis and Comparison

In this section, we evaluate the effect of different facets of the input on the runtime of HYPER. HYPER comprises two steps: (a) creating the aggregate view on which the query should be computed (done using a join-aggregate query), and (b) training regression functions to calculate conditional probability in the calculation of query output. This training is performed over a subset of the attributes of the view computed in the previous step. Training a regression function is more time-consuming than computing the aggregate view in step (1). Therefore, HYPER is as scalable as prior techniques for regression (we use a random forest regressor from the sklearn package). Hence the parameters we consider include (1) database size, (2) backdoor set size (see Section 3.3), and (3) query complexity. Since the effect of (1), (2) on the runtime of what-if query evaluation is directly translated to an effect on the runtime of how-to query evaluation, for how-to queries, we focus on the



(a) German-Syn (1M) (b) Students-Syn
Figure 10: What-If Query output.

effect of the number of attributes in the HowToUpdate operator which will change the optimization function ϕ (see Section 4.3).

What-if: database size. Table 1 presents the average running time to evaluate the response to a what-if query. To evaluate the effect of database size on running time, we varied the number of tuples in German-Syn dataset from 10K to 1M. Figure 11 compares the average time taken by HYPER, HYPER-sampled with Indep for five different What-If queries and Opt-HowTo for How-to queries. We observed a linear increase in running time with respect to the dataset size for all techniques except HYPER-sampled. The increase in running time is due to the time taken to train a regressor which is used to estimate conditional probabilities for query output calculation. To answer a hypothetical query, aggregate view calculation requires less than 1% of the total time. The majority of the time is spent on calculating the query output using the trained regressor. Therefore, the time taken by HYPER-sampled does not increase considerably when the dataset size is increased beyond 100K.

What-if: backdoor set size. This experiment changed the background knowledge to increase the backdoor set from 2 attributes to 6 attributes. The running time to calculate expected fraction of high credit individuals on updating account status increased from 7.2 seconds when backdoor set contains age and sex to 22.45 seconds when the backdoor set contains all attributes.

What-if: query complexity. In this experiment, we synthetically add multiple attributes in the Student-syn dataset and the different operators of the query to estimate their on running time.

On adding multiple attributes in the Use operator, the time taken to compute the relevant view increases minutely. For Student-Syn, Use operator was evaluated in less than 0.5 seconds when 5 different attributes are added from other datasets. The increase in these attributes do not affect the running time of subsequent steps unless the attributes in FOR operator increase.

We now compare the effect of adding attributes in the FOR operator of a Count query. Adding conditions involving Pre value of attributes increases the number of attributes used to train the regressor, which increases the running time (Figure 12a). Running time increased from 4.2 seconds when FOR operator is empty to 12.1 seconds and 17.7 seconds when it contains 5 and 10 attributes, respectively. In contrast, Indep is more efficient as it does not use additional attributes to compute query output. However, if the added attribute is in the backdoor set, then the output is evaluated faster. To understand the effect of adding such attributes, we considered a query where the backdoor set contained 10 attributes. To evaluate

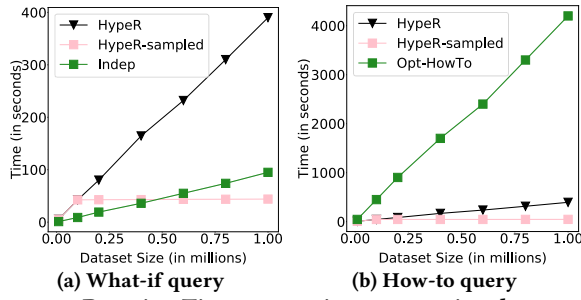


Figure 11: Running Time comparison on varying dataset size for German-Syn dataset averaged over five different queries.

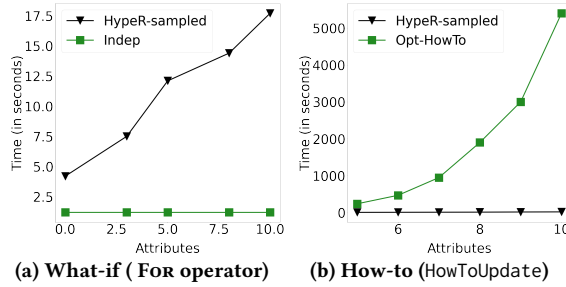


Figure 12: Running Time comparison on varying number of attributes in different operators for Student-Syn dataset.

the output, HYPER required 49.7 seconds and the time reduced to 7.4 seconds when 5 conditions on these attributes are added.

How-to: query complexity. Figure 12b presents the effect of the number of attributes in HowToUpdate operator on the time taken to process the query. Increasing attributes leads to a linear increase in the number of variables in the integer program (IP). Therefore, the time taken by HYPER-sampled increases from 7 seconds for 5 attributes to 20 seconds for 10 attributes. In contrast, Opt-HowTo considers all possible combinations of attribute values in the domain of attributes in the HowToUpdate operator. It takes around 4 minutes for 5 attributes and more than 90 minutes for 10 attributes, justifying the usefulness of the IP based formulation.

6 RELATED WORK

We review relevant literature in hypothetical reasoning in databases, probabilistic databases, and causality. *The main distinction of this paper from previous work is a framework that allows for hypothetical reasoning over relational databases using a post-update distribution over possible worlds that is able to capture both direct and indirect probabilistic dependencies between attributes and tuples using a probabilistic relational causal model.*

Previous work has focused on What-if and How-to analysis mainly in terms of provenance and view updates. Due to its practicality, and real applications like evaluating business strategies, there have been several works that developed support for *hypothetical what-if reasoning* in SQL, OLAP, and map-reduce environments [9, 29, 31, 37, 55]. What-if reasoning through provenance updates have been studied in [7, 16–18] to efficiently measure the direct effect of updating values in the database on a view created by the query. Nguyen et. al. [36] study the problem of efficiently performing what-if analysis with conflicting goals using data grids. Other

works have considered models for hypothetical reasoning in temporal databases [8, 27], where Arenas et. al. [8] focused on a logical model in which each transaction updates the database and the goal is to answer a query about the generated sequence of states, without performing the update on the whole database, and GreyCat [27] focused on time-evolving graphs. Christiansen et. al. [12] propose an approach that considers a single possible world and then modifies the query evaluation procedure within a logic-based framework. Another part of hypothetical reasoning is *how-to queries* which have been explored mostly in terms of provenance updates [33–35] that compute their results with hypothetical updates modeled as a Mixed Integer Program. MCDB [30] allows users to create an uncertain database that has randomly generated values in the attributes or tuples (that may be correlated with other attributes or tuples). These are generated using variable generation functions that can be arbitrarily complex. It then evaluates queries over this database using Monte Carlo simulations. Eisenreich et. al. [21] propose a data analysis system allowing users to input attribute-level uncertainty and correlations using histograms and then perform operations on the data such as aggregating or filtering uncertain values. We note that uncertainty in databases has been studied in previous work on *probabilistic databases* [4, 6, 14, 15, 48], where each tuple or value has a probability or confidence level attached to it, and in stochastic package queries [10] that allow for optimization queries on stochastic attributes. We adapt and use the concept of block-independent database model from probabilistic databases [14, 41] in this paper. The framework suggested in this paper uses a *probabilistic relational causal model* [46] to model updates as interventions and generate the post-update distribution that describes the dependencies between the attributes and tuples. There is a vast literature on *observational causal inference* on stored data in AI and Statistics (e.g., [5, 13, 25, 26, 38, 42–44, 49]), and we use standard techniques from this literature to compute query output.

7 CONCLUSIONS

We have defined a probabilistic model for hypothetical reasoning in relational databases. While the post-update distribution can stem from any probabilistic model, we focus here on causal models. We develop HYPER: a novel framework that supports what-if and how-to queries and performs hypothetical updates on the database, measures their effect, and computes the query results. Our framework includes new SQL-like operators to support these queries for testing different hypothetical scenarios. We prove that the query output can be computed using causal inference and we further devise optimizations. We show that our approach provides query results that are rational and account for implicit dependencies in the database. In future work, we plan to add support for multi-attribute updates consisting of dependent attributes and also account for database constraints and other semantic constraints. Extensions to cyclic dependencies of attributes in causal graphs is an intriguing future work. We also plan to develop an interactive UI, where users can pose and explore hypothetical queries.

ACKNOWLEDGMENTS

This work is partially supported by the NSF awards IIS-1552538, IIS-1703431, IIS-1814493, IIS-2008107, IIS-2112606, and IIS-2127309.

REFERENCES

- [1] Pcmag (<https://www.pcmag.com/>).
- [2] Spacy <https://spacy.io/>.
- [3] Top laptop brands in the world <https://www.globalbrandsmagazine.com/top-laptop-brands-in-the-world/>, 2021.
- [4] P. Agrawal, O. Benjelloun, A. D. Sarma, C. Hayworth, S. U. Nabar, T. Sugihara, and J. Widom. Trio: A system for data, uncertainty, and lineage. In *PVLDB*, pages 1151–1154, 2006.
- [5] J. D. Angrist, G. W. Imbens, and D. B. Rubin. Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455, 1996.
- [6] L. Antova, C. Koch, and D. Olteanu. Maybms: Managing incomplete information with probabilistic world-set decompositions. In *ICDE*, pages 1479–1480, 2007.
- [7] B. S. Arab and B. Glavic. Answering historical what-if queries with provenance, reenactment, and symbolic execution. In *USENIX*, 2017.
- [8] M. Arenas and L. E. Bertossi. Hypothetical temporal reasoning in databases. *J. Intell. Inf. Syst.*, 19(2):231–259, 2002.
- [9] A. Balmin, T. Papadimitriou, and Y. Papakonstantinou. Hypothetical queries in an OLAP environment. In *PVLDB*, pages 220–231, 2000.
- [10] M. Brucato, N. Yadav, A. Abouzied, P. J. Haas, and A. Meliou. Stochastic package queries in probabilistic databases. In *SIGMOD*, pages 269–283, 2020.
- [11] S. Chiappa. Path-specific counterfactual fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 7801–7808, 2019.
- [12] H. Christiansen and T. Andreassen. A practical approach to hypothetical database queries. In *DYNAMICS*, volume 1472, pages 340–355, 1998.
- [13] L. A. Cox Jr. Probability of causation and the attributable proportion risk. *Risk Analysis*, 4(3):221–230, 1984.
- [14] N. N. Dalvi, C. Ré, and D. Suciu. Probabilistic databases: diamonds in the dirt. *Commun. ACM*, 52(7):86–94, 2009.
- [15] N. N. Dalvi and D. Suciu. Efficient query evaluation on probabilistic databases. *Vldb J.*, 16(4):523–544, 2007.
- [16] D. Deutch, Z. G. Ives, T. Milo, and V. Tannen. Caravan: Provisioning for what-if analysis. In *CIDR*, 2013.
- [17] D. Deutch, Y. Moskovitch, and N. Rinetzy. Hypothetical reasoning via provenance abstraction. In *SIGMOD*, pages 537–554, 2019.
- [18] D. Deutch, Y. Moskovitch, and V. Tannen. Provenance-based analysis of data-centric processes. *Vldb J.*, 24(4):583–607, 2015.
- [19] H. Donner, K. Eriksson, and M. Steep. Digital cities: Real estate development driven by big data. Technical report, Working Paper. 2018. Available online: <https://gpc.stanford.edu...>, 2018.
- [20] D. Dua and C. Graff. UCI machine learning repository, 2017.
- [21] K. Eisenreich and P. Rösch. Handling uncertainty and correlation in decision support. In *Proceedings of the Fourth International VLDB workshop on Management of Uncertain Data (MUD 2010)*, volume WP10-04, pages 145–159, 2010.
- [22] S. Galhotra, A. Gilad, S. Roy, and B. Salimi. Hyper: Hypothetical reasoning with what-if and how-to queries using a probabilistic causal approach. *CoRR*, abs/2203.14692, 2022.
- [23] S. Galhotra, R. Pradhan, and B. Salimi. Explaining black-box algorithms using probabilistic contrastive counterfactuals. In *SIGMOD*, pages 577–590, 2021.
- [24] M. Golfarelli and S. Rizzi. What-if simulation modeling in business intelligence. *Int. J. Data Warehous. Min.*, 5(4):24–43, 2009.
- [25] S. Greenland. Relation of probability of causation to relative risk and doubling dose: a methodologic error that has become a social problem. *American journal of public health*, 89(8):1166–1169, 1999.
- [26] S. Greenland and J. M. Robins. Epidemiology, justice, and the probability of causation. *Jurimetrics*, 40:321, 1999.
- [27] T. Hartmann, F. Fouquet, A. Moawad, R. Rouvoy, and Y. L. Traon. Greycat: Efficient what-if analytics for data in motion at scale. *Inf. Syst.*, 83:101–117, 2019.
- [28] R. He and J. J. McAuley. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *WWW*, pages 507–517, 2016.
- [29] H. Herodotou and S. Babu. Profiling, what-if analysis, and cost-based optimization of mapreduce programs. *PVLDB*, 4(11):1111–1122, 2011.
- [30] R. Jampani, F. Xu, M. Wu, L. L. Perez, C. M. Jermaine, and P. J. Haas. MCDB: a monte carlo approach to managing uncertain data. In *SIGMOD*, pages 687–700, 2008.
- [31] L. V. S. Lakshmanan, A. Russakovsky, and V. Sashikanth. What-if OLAP queries with changing dimensions. In *ICDE*, pages 1334–1336, 2008.
- [32] M. Lichman. UCI machine learning repository, 2013.
- [33] A. Meliou, W. Gatterbauer, and D. Suciu. Bringing provenance to its full potential using causal reasoning. In *TaPP*, 2011.
- [34] A. Meliou, W. Gatterbauer, and D. Suciu. Reverse data management. *PVLDB*, 4(12):1490–1493, 2011.
- [35] A. Meliou and D. Suciu. Tiresias: the database oracle for how-to queries. In *SIGMOD*, pages 337–348, 2012.
- [36] Q. V. H. Nguyen, K. Zheng, M. Weidlich, B. Zheng, H. Yin, T. T. Nguyen, and B. Stantic. What-if analysis with conflicting goals: Recommending data ranges for exploration. In *ICDE*, pages 89–100, 2018.
- [37] S. Nieva, F. Sáenz-Pérez, and J. Sánchez-Hernández. HR-SQL: extending SQL with hypothetical reasoning and improved recursion for current database systems. *Inf. Comput.*, 271:104485, 2020.
- [38] J. Pearl et al. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.
- [39] B. Qureshi. Towards a digital ecosystem for predictive healthcare analytics. In *MEDES*, pages 34–41, 2014.
- [40] S. Ramakrishnan, K. Nagarkar, M. DeGennaro, K. Srihari, A. K. Courtney, and F. Emick. A study of the CT scan area of a healthcare provider. In *Proceedings of the conference on Winter simulation*, pages 2025–2031, 2004.
- [41] C. Ré and D. Suciu. Materialized views in probabilistic databases for information exchange and query optimization. In *Vldb*, pages 51–62, 2007.
- [42] D. W. Robertson. Common sense of cause in fact. *Tex. L. Rev.*, 75:1765, 1996.
- [43] J. Robins and S. Greenland. The probability of causation under a stochastic model for individual risk. *Biometrics*, pages 1125–1138, 1989.
- [44] D. B. Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- [45] B. Salimi, J. Gehrke, and D. Suciu. Bias in OLAP queries: Detection, explanation, and removal. In *Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018*, pages 1021–1035, 2018.
- [46] B. Salimi, H. Parikh, M. Kayali, L. Getoor, S. Roy, and D. Suciu. Causal relational learning. In *SIGMOD*, pages 241–256, 2020.
- [47] S. K. Singh and J. B. Lee. How to use what-if analysis in sales and operations planning. *The Journal of Business Forecasting*, 32(3):4, 2013.
- [48] D. Suciu. Probabilistic databases for all. In *PODS*, pages 19–31, 2020.
- [49] J. Tian and J. Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics and Artificial Intelligence*, 28(1-4):287–313, 2000.
- [50] F. Tramèr, V. Atlidakis, R. Geambasu, D. Hsu, J.-P. Hubaux, M. Humbert, A. Juels, and H. Lin. Fairtest: Discovering unwarranted associations in data-driven applications. In *IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 2017.
- [51] <https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>. Random forest regression – sklearn python library.
- [52] T. J. VanderWeele and W. An. Social networks and causal inference. *Handbook of causal analysis for social research*, pages 353–374, 2013.
- [53] Y. Zhang, H. Chen, H. Sheng, and Z. Wu. Applying hypothetical queries to e-commerce systems to support reservation and personal preferences. In *IDEAS*, pages 46–53, 2007.
- [54] E. Zheleva and D. Arbour. Causal inference from network data. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 4096–4097, 2021.
- [55] G. Zhou and H. Chen. What-if analysis in MOLAP environments. In *FSKD*, pages 405–409, 2009.
- [56] I. Zliobaite, F. Kamiran, and T. Calders. Handling conditional discrimination. In *Proceedings of the 2011 IEEE 11th International Conference on Data Mining*, page 992–1001, 2011.