

# Causal Relational Learning

## ABSTRACT

Causal inference is at the heart of empirical research in majority of sciences and social sciences, and is critical for making sound data-driven decisions. The gold standard in causal inference is performing *controlled experiments*, which is not always feasible due to ethical, legal, or cost constraints. As an alternative, inferring causality from *observational data* has been extensively used in statistical studies in public policy or social sciences. However, the existing methods critically rely on a restrictive assumption that the population of study consists of *homogeneous units* that can be represented as a single flat table. In contrast, in many real-world settings, the study domain consists of heterogeneous units with complex relational structure, where the data is naturally represented as multiple related tables, and therefore the data consists now of an entire relational database as opposed to a single table. In this paper, we present a formal framework for causal inference from such relational data, propose a declarative language called CaRL for capturing users' assumptions and specifying causal queries using simple Datalog-like rules, and develop an underlying inference engine that answers a suite of causal queries considering *relational* and *isolated* effects of the applied treatment. We give extensive experimental evaluations on synthetic and real data, and illustrate the applicability of our method on estimating the causal effect of institutional prestige on the acceptance of papers under single-blind and double-blind review processes.

## 1 INTRODUCTION

The importance of causal inference for taking informed policy decisions has been long recognised in health, medicine, social sciences, and other domains. However, today's decision-making systems typically do not go beyond *predictive analytics* and thus fail to answer questions such as "What will happen to revenue if the price of X is lowered?". While predictive analytics has achieved remarkable success in diverse applications, it is based on fitting a model to observational data based on associational patterns [7]. Causal inference, on the other hand, goes beyond associational patterns to the process that generates the data, thereby enabling analysts to reason about *interventions*, *policies*, and *counterfactuals*.

The gold standard in causal analysis is performing *randomized controlled experiments*, where the *subjects* or *units* of the study are assigned randomly to a treatment or a placebo (withheld from the treatment). The difference of the outcome variable between the treatment group and the control group is called *average treatment effect*, and represents the causal

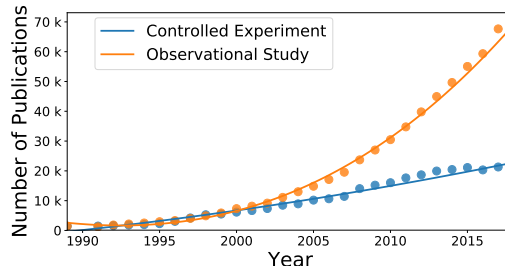


Figure 1: Number of publications on observational studies v.s. controlled experiments (obtained from SemanticScholar [10]).

effect. However, controlled experiments are not always feasible due to ethical, legal, or cost constraints. An attractive alternative that has been used in statistics, economics, and social studies is to simulate controlled experiments by *Observational Studies* using any data available, called *observational data*. While we can no longer assume that the treatment has been randomly assigned, under appropriate assumptions we can still learn causal relationships. Two well established frameworks exist for this purpose, Rubin's Potential Outcome Framework [9] and Pearl's Causal Models [5]. Both have been extensively studied in the literature and used in various applications, and there is an increased interest in observational studies, as can be seen in Fig 1, which shows the result of our quick analysis on SemanticScholar, revealing an increased interest in observational studies compared to controlled experiments. Babak: [The last sentence is too long] Dan: [agreed; it should capture concisely, and with the balance the statements "this is the real trend in causality today" and "it's a a result that we obtained by doing a quick analysis using SemanticScholar"]

However, all causal frameworks critically rely on the assumption that the units of study are sampled from a population consisting of homogeneous units that can be represented as a single flat table. In many real-world settings, the study domain consists of *heterogeneous units* that have a *complex relational structure*. In other words, the data usually consists of an entire relational database, and is not a single table. Standard notions used in causal analysis, such as units, casual DAG, covariates, no longer apply to relational data, prohibiting us from adopting existing causal inference frameworks to relational domains. We illustrate these challenges with an example.

*Example 1.1. (REVIEWDATA) OpenReview [4] is a collection of paper submissions and their reviews to several conferences, mostly in ML and AI. What makes it interesting is that it contains the review scores for both accepted and rejected papers. Scopus [11] is a large, well maintained database of peer-reviewed literature, including scientific journals, books and conference proceedings. We crawled and integrated these two sources to produce a relational database, which we show in a simplified form in Fig. 2. Data sources like this represent a treasure trove of information for the leadership of scientific conferences and journals. For example, they can help answer questions like “does double blind achieve its desired effect”? “Does increasing (or decreasing) the page limit affect the quality of the papers and how?”. To help in real-life decision, discovering associations is not sufficient, but instead decision makers would like to know if there exists a causal effect. For example, suppose a conference is currently requiring double blind submissions, and the leadership is questioning its effectiveness. If we revert to single blind submissions, would that represent an unfair advantage for authors **from top universities and companies (a.k.a. “prestigious institutions”)** authors? *Dan: [don’t overcomplicate it. We have already given two examples above, now we need to narrow down to the question the we will eventually answer in the experiments. Whether it’s about prestige of authors or prestige of institutions, that’s exactly how we should formulate the question here.]* Given a dataset like in Fig. 2, one can run a few SQL queries and check whether senior authors consistently get better reviews at single blind conferences than at double blind conferences, but this can only prove or disprove correlation, and not causation. Alternatively, one can try to apply the Neyman-Rubin causal model [9], but that requires us to present the data as a single table of independent units. If we do this naively on our dataset (e.g. computing the universal table [14]), then we create interference effects and contagion effects, *Babak: [joining data sets do not create contagion and interference, these are properties of the data generative model (how the domain works)]* *Dan: [you are right. But the point to make here is that the relational data captures correctly the generative model, but if we just join them and re-define units based in the joint data, then there is interference and contagion. Do you agree?]* both of which prohibit standard causal analysis. For instance, prestige of an author not only influences their acceptance rate, but also has a spill-over effect on the acceptance rate of their co-authors; this is called interfere. Authors’ qualifications can be contagious over a course of time, meaning that if a junior author collaborates frequently with a senior person then the overall quality of his/her paper may increase.*

In this paper, we propose a declarative framework for *Causal Relational Learning*, a foundation for inferring causal

Person	Author	Submission	Submitted	Conference
Bob	Bob	s1	s1	DB
Carlos	Eva	s2	s2	AI
Eva	Eva	s3	s3	AI
	Eva	s3		
	Carlos	s3		

Prestige	Qualifications	Score	Blind?
Bob	Expe- rience	s1	DB
Carlos	h-index	s2	AI
Eva		s3	Single
			Double

Expe- rience	h-index
Bob	10
Carlos	8
Eva	2

Score
s1
s2
s3

**Figure 2: Fragment of a relational database schema, obtained by integrating OpenReview and Scopus. (The data shown is not real.)**

inference from relational domains. Specifically, we propose Relational Causal Models, which extend Pearl’s Causal Models from representing causal relationship between attributes of homogeneous units, to heterogeneous units with complex relational structures. We extend Pearl’s do-operator to relational causal models, to capture complex *relational interventions* and *relational causal queries*. We develop algorithms for detecting a sufficient set of *covariates* that should be adjusted for to remove confounding effects. We develop novel techniques such as *embedding relational data* for transforming relational data into a flat table amenable for easy causal inference.

*Dan: [I think that this paragraph still needs work to explain our contributions]* At the core of the relational causal model is a declarative language that specifies a complex causal structure in relational domains, and allows uses to formulate causal queries. As oppose to predictive analysis, domain expert knowledge plays a critical role and is fundamentally required for causal inference [6]. In our declarative language called CaRL users can specify their background knowledge and their causal queries about relational domains with simple and intuitive Datalog-like rules.

*Dan: [We should have one contribution per section.]* We make the following contributions: (1) We propose Relational Causal Molds, which lift Pearl’s Causal Models to relational domains. (2) We propose CaRL a declarative language for capturing users’ assumptions and queries using simple Datalog-like rules. (2) We develop an underlying inference engine that automatically detects a sufficient set of *covariates* that should be adjusted for to remove confounding effects. (3) We develop novel techniques such as *embedding* to transform relational data into a flat table amenable to easy causal inference. (4) we answers a suite of causal queries considering relational and isolated effects. (5) We evaluate CaRL on synthetic and real data, and illustrate the applicability of our method on estimating the causal effect of institutional

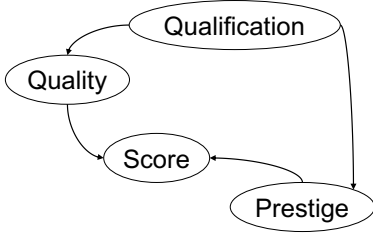


Figure 3: A simple causal DAG.

prestige on the acceptance of papers under single-blind and double-blind review processes.

## 2 BACKGROUND ON CAUSAL ANALYSIS

We briefly review here some basic concepts in causal analysis, covering both Pearl’s casual DAGs and the Neyman-Rubin model, following [5, 9].

**Causal Models.** A probabilistic causal model is a tuple  $M = \langle U, V, F, Pr_U \rangle$ , where  $U$  is a set of exogenous variables that cannot be observed,  $V$  is a set of observable or endogenous variables, and  $F = (F_X)_{X \in V}$  is a set of *non-parametric structural equations*  $F_X : Dom(Pa_V(X)) \times Dom(Pa_U(X)) \rightarrow Dom(X)$ . Here  $Pa_U(X) \subseteq U$  and  $Pa_V(X) \subseteq V - \{X\}$  are called the exogenous parents and endogenous parents of  $X$  respectively; and  $Pr_U$  is a joint probability distribution on the exogenous variables  $U$ . Intuitively, the exogenous variables  $U$  are not known, but we know their probability distribution, while the endogenous variables are completely determined by their parents (exogenous and/or endogenous).

**Causal DAG.** The probabilistic causal model is associated to a *causal graph*,  $G$ , whose nodes are the endogenous variables  $V$ , and whose edges are all pairs  $(Z, X)$  such that  $Z \in Pa_V(X)$ ; we write  $Z \rightarrow X$  for an edge; it is usually assumed that  $G$  is acyclic and then it is called a Causal DAG. In other words, the causal DAG hides the exogenous variables (since we can’t observe them anyway) and instead captures their effect by defining a probability distribution  $Pr_V$  on the endogenous variables. This is possible under the *causal sufficiency* assumption<sup>1</sup>. The formula for  $Pr_V$  is the same as that for a Bayesian network:

$$Pr(V) = \prod_{X \in V} Pr(X|Pa(X)) \quad (1)$$

We will only refer to endogenous variables in the rest of the paper and drop the subscript  $V$  from  $Pa_V$  and  $Pr_V$ . Fig. 3 shows a simple example of a causal graph: the Score (of a

<sup>1</sup>The assumption requires that, for any two variables  $X, Y \in V$ , their exogenous parents are disjoint and independent  $Pa_U(X) \perp\!\!\!\perp Pa_U(Y)$ . When this assumption fails, one adds more endogenous variables to the model to expose their dependencies.

paper) is affected by its Quality and by the Prestige of the author (assuming the reviews are single blind); the latter are both affected by the authors’ Qualification. In this paper we will assume that the causal DAG is known<sup>2</sup>, and one uses some observational data in order to learn the conditional probability distribution (1). Notice a fundamental limitation of traditional causal frameworks. The units need to be uniform and independent. In other words the data consists of a uniform set of tuples, called *units*, with four attributes (Qualification, Quality, Prestige, Score); each row represents one paper, its (unique?) author, and the outcome. Situations like multiple authors for the same paper, or multiple submissions by the same authors cannot be represented in this model and, in fact, would violate fundamental assumptions made by the model.

**d-Separation and Markov compatibility.** A common inference question in a causal DAG is how to determine whether a CI  $(X \perp\!\!\!\perp Y|Z)$  holds. A sufficient criterion is given by the notion of d-separation, a syntactic condition  $(X \perp\!\!\!\perp Y|_d Z)$  that can be checked directly on the graph.  $Pr$  and  $G$  are called *Markov compatible* if  $(X \perp\!\!\!\perp Y|_d Z)$  implies  $(X \perp\!\!\!\perp Y|Z)$ ; if the converse implication holds, then we say that  $Pr$  is *faithful* to  $G$ . The following is known:

PROPOSITION 2.1. *If  $G$  is a causal DAG and  $Pr$  is given by Eq.(1), then they are Markov compatible.*

**Interventions and do operator.** An *intervention* represents, intuitively, actively setting an endogenous variable to some fixed value, and observing the effect. Pearl [5] introduced for this purpose the *do*-calculus. Formally, an intervention consists of setting a variables  $U$  to some values  $U = \mathbf{u}$ , and it defines the probability distribution  $Pr(V|do(U = \mathbf{u}))$  given by Eq. 1 where we remove all factors  $Pr(X|Pa(X))$  where  $X \in U$ . In other words we modify the causal DAG by removing all edges entering the variables on which we intervene; of course, this is different from conditioning,  $Pr(V|U = \mathbf{u})$ . Pearl has an extensive discussion on the rationale of the *do*-calculus and describes several equivalent formulas for estimating  $Pr(V|do(U = \mathbf{u}))$  from observed distribution.

**Average Treatment Effect (ATE).** The scope of Neyman-Rubin’s causal model [9] is to compare the effect of a binary treatment variable  $T$  on some response variable  $Y$ , which is often measured by the following quantify known as *Average Treatment Effect (ATE)* and expressed as follows in our notation:

$$ATE(Y, T) = \mathbb{E}[Y|do(T = 1)] - \mathbb{E}[Y|do(T = 0)] \quad (2)$$

Much of the literature on the Neyman-Rubin causal model in statistics is concerned with efficient estimation of *ATE* from observational data.

<sup>2</sup>There exists an extensive literature on learning the causal DAG from data.

### 3 RELATIONAL CAUSAL MODEL

As discussed in previous sections, the traditional frameworks for causal inference are unsuitable for capturing complex causal relationships in relational domains comprising heterogeneous units, which may be different entities or relationships between those entities. In this section we propose a declarative language called *CaRL* (*Causal Relational Language*) that extends Pearl’s Causal Model and causal graphs to relational data by allowing the user to (1) specify assumptions and backgrounded knowledge on the interactions among heterogeneous units (Section 3.2), and (2) ask various causal queries (Section 3.3). We start with our data model that forms the basis for our language in Section 3.1.

#### 3.1 Data Model

**Schema.** A *relational causal schema* is a tuple  $S = (P, A)$ , where  $P = \{P_1, \dots, P_n\}$  is a standard relational schema that represents a set of *entities*  $E$  and their *relationships*  $R$ , hence  $P = E \cup R$ ; and  $A = \{A_1, \dots, A_k\}$  is a set of *attribute functions* with fixed arity, domain and range  $Range(A)$  that encode the descriptive values of the entities and their relationships.

**Skeleton and observed instance** Our language represents the causal relationships between the attribute functions given a set of entities and their relationships. Therefore, we refer to a standard relational database instance from schema  $P$  as a *relational skeleton*. A relational skeleton  $\Delta$  together with the values of a subset of observed attribute functions  $A$  is called an *observed relational instance*. Moreover, we distinguish between the attribute functions and the entities and their relationships by denoting the former with  $P(\cdot)$  and the latter with  $A[\cdot]$ .

*Example 3.1.* The relational causal schema corresponding to the relational data in Figure 2 is as follows:

$P = \text{People}(A), \text{Submission}(S), \text{Conference}(C), \text{Submitted}(S, C), \text{Author}(A, S)$   
 $A = \text{Prestige}[A], \text{Qualification}[A], \text{Score}[S], \text{Quality}[S], \text{Blind}[C]$

Here,  $\text{People}(A)$ ,  $\text{Submission}(S)$  and  $\text{Conference}(C)$  are entities; whereas  $\text{Authors}(A, S)$  and  $\text{Submitted}(S, C)$  are relationships. Moreover,  $\text{Prestige}[A]$  describes the prestige of the institution that the author is affiliated with measured by its ranking.  $\text{Qualification}[A]$  is a compact representation of an author’s *h-index* and experience.  $\text{Score}[S]$  describes the average score reviewers gave to a submission: 1 being perfect and 0 being the worst possible.  $\text{Quality}[S]$  describes the quality of a submission (which is unobserved). Finally,  $\text{Blind}(C)$  shows whether a conference review policy is single-blind or double-blind. We refer to the tables in the first row of Figure 2 as the relational skeleton of the observed instance.

#### 3.2 Specification Language and Semantics

Now we discuss the syntax of CaRL to encode backgrounded causal knowledge in relational domains. At the core of CaRL,

we have *relational causal rules* to capture the causal assumptions.

*Definition 3.2.* A *relational causal rule* (or simply a *rule*) over a relational schema  $S = (P, A)$  has the following form:

$$A[X] \Leftarrow A_1[X_1], \dots, A_k[X_k] \text{ WHERE } Q(Y) \quad (3)$$

Here  $A, A_1, \dots, A_k \in A$ ,  $Q$  is a Boolean conjunctive query (BCQ) on the predicates in  $P$ ,  $\cup_{i=1}^k X_i \subseteq Y$ , and  $X \subseteq Y$ . We call  $A[X]$  the *head* of the rule,  $A_1[X_1], \dots, A_k[X_k]$  the *body* of the rule, and  $Q(Y)$  the *condition*. Denote  $\phi_A$  an rule with an attribute function  $A$  in the head.

*Definition 3.3.* A *relational causal model*  $\Phi$  for a relational schema  $S = (P, A)$  is a collection of rules associated to  $A$ , i.e.,  $\Phi = \cup_{A \in A} \phi_A$ .

A rule is template for generating multiple *grounded rules*.

*Definition 3.4.* A rule  $\phi_A$  and a relational skeleton  $\Delta$  define a set of *grounded rules*  $\phi_A^\Delta$  obtained by substituting the variables  $Z = X \cup X_1 \cup \dots \cup X_k$  in  $\phi_A$  with a set of constants  $z$  (denoted by  $Y/z$ ), such that  $\Delta \models Q([Y/z])$ , i.e., when we substitute variables  $Y$  with constants  $z$  the Boolean query  $Q$  evaluates to true.

A relational causal model  $\Phi$  together with a relational skeleton  $\Delta$  defines a collection of grounded rules  $\Phi^\Delta \stackrel{\text{def}}{=} \cup_{\phi_A \in \Phi} \phi_A^\Delta$ . To each  $\Phi^\Delta$ , we associate a *grounded causal graph*  $G(\Phi^\Delta)$  with vertices consisting of the set of all grounded atoms in  $\Phi^\Delta$ , denoted  $A^\Delta$ , and directed edges  $(A_i(x_i), A_j(x_j))$  if  $A_i(x_i)$  appears in the head and  $A_j(x_j)$  appears in the body of a grounded rule in  $\Phi^\Delta$ . We denote  $A^\Delta \subseteq A^\Delta$  the set of all groundings of an attribute function  $A \in A$  in  $A^\Delta$ . We interpret a grounded causal graph  $G(\Phi^\Delta)$  as a standard probabilistic causal models characterized by a joint probability distribution  $\Pr(A^\Delta)$  and the following conditional probability distributions, for each  $A[x] \in A^\Delta$ .

$$\Pr(A[x] \mid \text{Pa}(A[x])). \quad (4)$$

*Example 3.5.* Consider the following rules that postulate a relational causal model for REVIEWDATA in Figure 2.

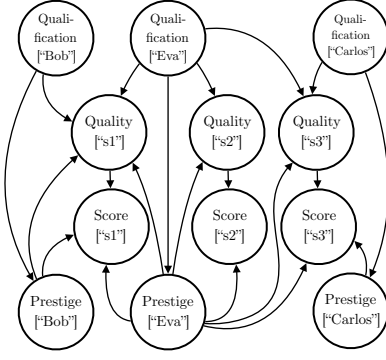
$$\text{Prestige}[A] \Leftarrow \text{Qualification}[A] \text{ WHERE } \text{Author}(A) \quad (5)$$

$$\text{Quality}[S] \Leftarrow \text{Qualification}[A], \text{Prestige}[A] \text{ WHERE } \text{Contributor}(A, S) \quad (6)$$

$$\text{Score}[S] \Leftarrow \text{Prestige}[A], \text{ WHERE } \text{Contributor}(A, S), \quad (7)$$

$$\text{Score}[S] \Leftarrow \text{Quality}[S] \text{ WHERE } \text{Submission}(S) \quad (8)$$

The rule (5) indicates that the qualification of an author causally affects their prestige due to the fact that prestigious institutions tend to hire most qualified candidates. The rule (6) stipulates that the quality of a submission is affected by the qualifications of its authors as well-as their prestige (presumably, authors from prestigious institutions have access



**Figure 4: grounded causal graph corresponded to the grounded rules in Example 3.5.**

to more resources). Finally, the rules (7) and (8) collectively stipulate that reviewers scores are based on the quality of a submission, but also potentially influenced by prestige of its authors. We implicitly assign a rule with empty body to the rest of functional attributes, meaning that they not causally influenced by the other attribute functions in the schema.

The relational causal model together with the toy relational skeleton instance in Figure 2 generates the following grounded rules:

$$\text{Prestige}["Bob"] \Leftarrow \text{Qualification}["Bob"] \quad (9)$$

$$\text{Prestige}["Carlos"] \Leftarrow \text{Qualification}["Carlos"] \quad (10)$$

$$\text{Prestige}["Eva"] \Leftarrow \text{Qualification}["Eva"] \quad (11)$$

$$\text{Quality}["s1"] \Leftarrow \text{Qualification}["Bob"], \text{Qualification}["Eva"] \quad (12)$$

$$\text{Quality}["s2"] \Leftarrow \text{Qualification}["Eva"] \quad (13)$$

$$\text{Quality}["s3"] \Leftarrow \text{Qualification}["Carlos"], \text{Qualification}["Eva"] \quad (14)$$

$$\text{Score}["s1"] \Leftarrow \text{Quality}["s1"], \text{Prestige}["Bob"], \text{Prestige}["Eva"] \quad (15)$$

$$\text{Score}["s2"] \Leftarrow \text{Quality}["s2"], \text{Prestige}["Eva"] \quad (16)$$

$$\text{Score}["s3"] \Leftarrow \text{Quality}["s3"], \text{Prestige}["Carlos"], \text{Prestige}["Eva"] \quad (17)$$

The rules in Example 3.5 result in the grounded causal graph shown in Figure 4. Note that our language allows for recursive rules:

$$\text{Qualification}[A] \Leftarrow \text{Qualification}[A'] \text{ WHERE } \text{Collaborator}(A, A') \quad (18)$$

Where,  $\text{Collaborator}(A, A')$  is a derived predicates that shows whether two authors were collaborators in the past. The rule (18) states that an author's qualifications has a feedback nature, i.e., the qualification of an author can influence the qualifications of his/her collaborators over time. Even though CaRL is equipped to handle cyclic causal dependency such as 18, using *equilibrium of discrete Markov chains* for brevity we discuss recursive rules only in the full technical report. [SR: do we want to mention cyclic rules? we have not discussed those details yet. How about saying out of scope of the current paper?]

### 3.3 Query Language

This section describes different types of causal queries supported in CaRL. The semantics of the queries will be explained in Section 4.

*Average treatment effect (ATE).* In CaRL, queries about ATE (see (2)) of a binary attribute function  $T[X]$  on a response attribute function  $Y[X']$  expressed as follows:

$$Y[X'] \Leftarrow T[X]? \quad (19)$$

For example, the query

$$\text{Score}[A] \Leftarrow \text{Prestige}[A]? \quad (20)$$

compute the ATE of *Prestige* on *Score*, i.e., it compares papers' scores in two hypothetical worlds in which all authors are, and are not affiliated with prestigious institutions.

*Aggregated Response.* In CaRL one can define aggregated responses such as  $\text{AVG\_score}[A]$ , the average submission scores of an authors, and answer queries about ATE of a treatment on an aggregated outcome. For example the following query quantifies the ATE of *Prestige* on  $\text{AVG\_score}[A]$ :

$$\text{AVG\_score}[A] \Leftarrow \text{Prestige}[A]? \quad (21)$$

*Relational and Isolated Effects.* In relational domain units that are relationally connected can have a causal influence on each other. For example, Prestige of an authors not only influences their average submission scores but also their collaborators average submission scores. To measure such complex relational causal interactions CaRL support queries about *relational* and *isolated* effects, expressed as follows (will be defined rigorously in Section 4):

$$Y[X'] \Leftarrow T[X]? \text{ WHEN } \langle \text{cnd} \rangle \text{ PEERS TREATED} \quad (22)$$

where,  $\langle \text{cnd} \rangle$  is a condition with the following grammar:

$$\begin{aligned} \langle \text{cnd} \rangle \Leftarrow & \langle \text{LESS} \mid \text{MORE} \rangle \text{ THAN } k\% \mid \text{AT } \langle \text{MOST} \mid \text{LEAST} \rangle k \mid \\ & \text{EXACTLY } k \mid \text{ALL} \mid \text{NONE} \end{aligned} \quad (23)$$

For example, the query

$$\text{Score}[S] \Leftarrow \text{Prestige}[A]? \text{ WHEN ALL PEERS TREATED} \quad (24)$$

compute both isolated and relation effect prestige on reviewer scores if all coauthors of an author are from prestigious institutions. The isolated effect measures to what extent an author review score is influenced by their own prestige. Whereas, the relational effect quantifies to what extent the author review scores are influenced by the prestige of their collaborators.

### 3.4 Structural Homogeneity

Recall from Section 2 that in standard probabilistic causal models the conditional probability distribution  $\Pr(X|\mathbf{Pa}(X))$  for each  $X \in \mathbf{V}$ , captures the underlying structural equation  $F_X$  associated to  $X$  as a probabilistic mapping. Since, different



variables have different corresponding structural equations, it is natural to associate different conditional probability distribution to each variable. For instance in Figure 3 the structural equations corresponded to *Score* and *Prestige* models two different functions with different inputs and outputs. Also note that these conditional probability distributions are unknown and must be estimated from available data.

However, in CaRL assuming different structural equations, and thereby different conditional probability distribution  $\Pr(A[x] \mid \text{Pa}(A[x]))$  for each  $A[x] \in A^\Delta$ , is pathological. Here, the number of grounded atom are not fixed a priori and depends on the size of the skeleton (the relational database); we no longer have the attribute *Score*, but instead we have  $\text{Score}[s_1], \text{Score}[s_2], \dots$ . Under the unreasonable assumption that the structural equations associated to grounded atom are different, the estimation of the corresponding conditional probability distributions  $\Pr(A[x] \mid \text{Pa}(A[x]))$ , and hence causal inference becomes impossible. Therefore in this paper we make the following reasonable *Structural Homogeneity assumption*:

- **Structural Homogeneity:** All the grounded atoms  $A[x] \in A^\Delta$  of an attribute function  $A \in \mathbf{A}$  share the same structural equation, which models a function with variable input size that determines the value of  $A[x]$ , given its endogenous and exogenous parents.

For instance, in Example 3.5, we assume the underlying structural equation for submission’s scoring is the same for all submissions.

The structural homogeneity assumption, however, is not easily captured by associating the same conditional probability distribution corresponding to all grounded atoms  $A[x] \in A^\Delta$ . Different grounded atoms can have different number of parents in the grounded causal digram. For instance, consider the atoms  $\text{Score}[s_1]$  and  $\text{Score}[s_2]$  from (15) and (17). In this toy example,  $\text{Score}[s_1]$  depends on the *Prestige* of two authors “Eva” and “Bob”, whereas  $\text{Tutor}[s_2]$  depends on the *Prestige* of one author only: “Eva.” This requires a *unification* of the parents of the same type using another layer of functions  $\psi$  (using aggregates, or in general, an embedding, see Section ??).

Formally, instead of using (4), we assume the structural equations associated to the groundings of an attribute function  $A$  can be captured using the following conditional probability distribution shared by all  $A[x] \in A^\Delta$ :

$$\Pr(A[x] \mid \Psi^A(\text{Pa}(A[x]))) \quad (25)$$

where,  $\Psi^A$  is a collection of mappings that projects the parents of  $A[x]$  into a low-dimensional vector with fixed dimensionality for all  $A[x] \in A^\Delta$ . Intuitively, we assume that the mappings provide *sufficient statistics* for evaluation of the underlying structural questions corresponded to all  $A[x] \in A^\Delta$ .

Throughout this paper we assume  $\Psi_V^A$  is known and consists of set of mappings  $\{\psi_{A_1}^A, \psi_{A_2}^A, \dots\}$  where, each  $\psi_{A_i}^A$  is an *embedding function* that map a subset of  $\text{Pa}(A[x])$  consisting of the grounding of the attribute function  $A_i$  into a low-dimensional *embedding space* with fixed dimensionality. In section ?? we discuss a set of techniques to learn  $\Psi^A$  from data.

*Example 3.6.* Let us assume whether a submission is accepted or not is a function of the result of applying a mapping  $\psi_{\text{Prestige}}^{\text{Score}}$  to the vector of the *Prestige* of the submission authors, of the *Prestige* of the authors of the submission into a vector with a fixed dimension for all submissions, e.g., a vectors consisting of the average (weighted by the authors position in the submission) *Prestiges* and the number of authors. In our running example the vector  $\langle 1, 1 \rangle$  and  $\langle 0 \rangle$ , respectively corresponded to the prestige of authors of the submissions  $s_1$  and for  $s_2$ . Using average and size of the vector, the embedded vectors become  $\langle 1, 2 \rangle$  and  $\langle 0, 1 \rangle$ , respectively.

Note that in this paper we work on the common assumption of causal sufficiency and acyclicity of  $G(\Phi^\Delta)$ . The following factorization, which forms a basis for defining and identifying the effect of interventions (see Section 4 and 5), immediately implied from (25) and Markov compatibility of  $\Pr(A^\Delta)$  and the grounded causal graph  $G(\Phi^\Delta)$ .

$$\Pr(A^\Delta) = \prod_{A[x] \in A^\Delta} \Pr(A[x] \mid \Psi^A(\text{Pa}(A[x]))) \quad (26)$$

**Structural Heterogeneity.** In practice, relational data may consist of data collected from different domain, e.g., different organization, cooperation, countries, cities, etc that could lead to the violation of the Structural homogeneity assumption. For instance, its is reasonable to assume the the structural question associated to scoring is different across single-blind and double-blind conferences. Such situations can be expressed in CaRL by postulating different rule at *different granularity* in which the homogeneity assumption perceived to hold, same organization, all double-blind conferences, all conference in computer science, etc. For example, the heterogeneity of paper scoring can be captured using the following rules:

$$\begin{aligned} \text{SBblind\_Score}[S] &\Leftarrow \text{Prestige}[A] \text{ WHERE } \text{Contributor}(A, S) \\ \text{SBblind\_Score}[S] &\Leftarrow \text{Quality}[S] \text{ WHERE } \text{Submission}(S) \\ \text{DBblind\_Score}[S] &\Leftarrow \text{Quality}[S] \text{ WHERE } \text{Submission}(S) \end{aligned}$$

that show the generative models for scoring are different in double-blind and single-blind conferences. Note that CaRL supports a compact way for representing heterogeneous domains (to be included in full technical report).

In general for causal inference from relational data it is indispensable to assume that the homogeneity holds at some granularity level.

## 4 SEMANTIC FOR CAUSAL QUERIES

This section defines the semantics of the causal queries described in Section 3.3. Fix a relational causal schema  $S$ , a relational skeleton  $\Delta$ , and a relational causal model  $\Phi$  with a corresponding grounded causal graph  $G(\Phi^\Delta)$ . For an attribute function  $A \in \mathbf{A}$ , denote  $\mathbb{U}_A$  the set of all tuples of grounded entities  $\mathbf{x}$  such that  $A[\mathbf{x}] \in \mathbf{A}^\Delta$ . For example,  $\mathbb{U}_{Prestige}$  consists of all the authors, e.g., {"Bob", "Eva", "Carlos"}, whereas  $\mathbb{U}_{Score}$  consists of all the submissions, e.g., {"s1", "s2"}. We refer to each element  $\mathbf{x} \in \mathbb{U}_A$  as a *unit* of an attribute function  $A$ .

### 4.1 Treatment and response attribute function

In a causal analysis, we are given a pair of attribute functions  $T[\mathbf{X}], Y[\mathbf{X}'] \in \mathbf{A}$ , the goal is to estimate the effect of intervening on the attribute function  $T$  of the *treated units*  $\mathbb{U}_T$ , on the attribute function  $Y$  of the *response units*  $\mathbb{U}_Y$ , hence we call  $T$  the *treatment attribute function* and  $Y$  the *response attribute function*. For example, to study the effect of prestige on submission scores,  $Prestige[A]$  is the treatment and  $Score[S]$  is the response attribute function. Here, the treated units are  $\mathbb{U}_{Prestige}$  (all authors) and the response units are  $\mathbb{U}_{Score}$  (all submissions). For exposition, we assume  $Range(T) = \{0, 1\}$  and  $Range(Y) = \mathbb{R}$ .

Formally, given a set of treated units  $\mathbb{U}_T = \{\mathbf{x}_1, \mathbf{x}_2, \dots\}$  and a binary vector  $\vec{t} = (t_1, t_2, \dots)$ , we are interested in the effect of a set of interventions  $do(T(\mathbf{x}_i) = t_i)$  for all treated units  $\mathbf{x}_i$ , where each intervention replaces the NSE associated with  $T(\mathbf{x}_i)$  with a constant  $t_i$ ; in our example of the effect of prestige on score, the vector  $\vec{t}$  corresponds to a particular assignment of prestige to all authors, e.g., the vector  $\vec{1}$  identifies an intervention in which *hypothetically changes the authors' afflictions to a prestigious one*. By abuse of notation, we denote with  $do(T[\mathbb{S}] = \vec{t}_{\mathbb{S}})$  a set of interventions in which an arbitrary subset of treated units  $\mathbb{S} \subseteq \mathbb{U}_T$  receive  $\vec{t}_{\mathbb{S}}$  (with an implicit assumption on the order of the elements in the set  $\mathbb{S}$ ).

Here we emphasize the power of our framework for causal analysis on relational data with two key observations:

- **Non-uniform treated and response units.** *The treated and response units can correspond to different entities or relationships.* In our example, they correspond to entities for  $S$  and entities of  $(A, S)$ , but if we want to understand the intervention of prestige on a conference's average submission score, it could be two entities  $A$  and  $C$  that are not directly related.
- **Independent treated units.** *Each treated unit can be subjected to an intervention independently of all other treated units.* In our example, we can have a treatment vector  $\vec{1}$  where we force *all* to be afflicted with prestigious institutions, as well as another treatment vector

$(1, 0, 0, \dots)$  where we only force the "Bob" to be affiliated to a prestigious institution.

Our goal is to answer a set of queries that compare the *average response* of the units  $\mathbb{U}_Y$  to two alternative intervention strategies  $\vec{t}$  and  $\vec{t}'$  applied to the treated units  $\mathbb{U}_Y$  as we discuss below.

### 4.2 Average Treatment Effect (ATE) on attributes

The primary causal query in CaRL is that of average treatment effect (ATE) as in (19), which is defined as follows:

$$ATE(T, Y) \stackrel{\text{def}}{=} \sum_{\mathbf{x}' \in \mathbb{U}_Y} \frac{1}{m} (\mathbb{E}[Y[\mathbf{x}'] \mid do(T[\mathbb{U}_T] = \vec{0})] - \mathbb{E}[Y[\mathbf{x}'] \mid do(T[\mathbb{U}_T] = \vec{1})]) \quad (27)$$

Intuitively, ATE compares the expected response of the response units in two regimes of interventions in which all units receive treatment and all of them do not receive treatment respectively. For example,  $ATE(PRESTIGE, SCORE)$  compare papers' scores under two interventions in which all authors are and are not affiliated with prestigious institutions.

### 4.3 Extending attribute functions and schema with aggregates

In CaRL, one can extend the set of attribute functions  $\mathbf{A}$  by new aggregated attribute functions using one of the *aggregate rules* of the following forms. For  $A \in \mathbf{A}$ ,

$$\begin{aligned} AGG\_A[\mathbf{W}] &\Leftarrow A[\mathbf{X}] \text{ WHERE } Q(\mathbf{Z}) \\ AGG\_A[\mathbf{W}'] &\Leftarrow AGG\_A[\mathbf{W}] \text{ WHERE } Q(\mathbf{Z}) \end{aligned} \quad (28)$$

Here,  $\mathbf{Z} \supseteq \mathbf{X}' \cup \mathbf{W}$  and  $AGG$  is an aggregate function on  $\mathbf{A}$ , e.g., AVG (average) and VAR (variance). The new aggregated attribute functions  $AGG\_A$  are included in the extended attribute functions  $\mathbf{A}$  (for simplicity, we use  $\mathbf{A}$  for both given and extended attribute functions). The semantic of aggregated rules is defined similar to rules in Section 3, i.e., they define a set of grounded rules with corresponding vertices and edges in the grounded causal graph  $G(\Phi^\Delta)$  as described in Section 3. However, instead of a conditional probability distribution, a deterministic function  $AGG(Pa(AGG\_Y[\mathbf{w}]))$ , will be associated to each  $AGG\_Y[\mathbf{w}] \in AGG\_Y^\Delta$ . For example, the following aggregate rule defines the average acceptance rate for each author.

$$AVG\_Score[A] \Leftarrow Score[S] \text{ WHERE } Contributor(A, S) \quad (29)$$

Find the grounded causal graph that incorporates (29) in Figure 5.

#### 4.4 Average relational and isolated effects on attributes and aggregates

In order to formalize isolated and relational effects as described in Section 3.3, we need to establish a one-to-one correspondence between treated and response units by using aggregations carefully. To this end, first we define relational paths.

**Definition 4.1.** Given a relational causal schema  $S = (P, A)$  consisting of a set for entities and relations  $P = E \cup R$ , and a set of attribute functions  $A$  (possibly extended with aggregates), a *relational path* is a sequence of entities and relationships of the following form:

$$\mathcal{P} : E_1(X_1) \xleftrightarrow{R_1(X_1, X_2)} E_1(X_2) \cdots E_{\ell-1}(X_{\ell-1}) \xleftrightarrow{R_{\ell-1}(X_{\ell-1}, X_\ell)} E_\ell(X_\ell) \quad (30)$$

where,  $E_i(X_i) \in E$  and  $R_{i-1}(X_{i-1}, X_i) \in R$ , for  $i = 1, \dots, \ell$ .

For instance,  $Conference(C) \xleftrightarrow{Submitted(S, C)} Submission(S)$  is a relational path in our example. The treated and response units corresponding to a treatment and response attribute functions  $T$  and  $Y$  are said to be *relationally connected* if there exists a relational path  $\mathcal{P}$  that includes the entities or relationships that  $T$  and  $Y$  describe either as the endpoints in the path or as the labels of the edges at the ends of the path. For example, for  $T[X] = Prestige[A]$  and  $Y[X'] = Score[S]$ , the treatment is an attribute function of the entity  $Author(A)$ , the response is an attribute function of the relationship  $Contributor(A, S)$ , and the treated and response units are relationally connected by the following relational path:

$$Author(A) \xleftrightarrow{Contributor(A, S)} Submission(S) \quad (31)$$

In this paper, we make the natural assumption that the treated and response units are relationally connected by at least one relational path as in (30). Then the treated and response units can be unified using the aggregated response  $AGG\_Y[X]$  defined with the following aggregate rule that maps attribute  $Y$  in entity(s)  $X'$  to entity(s)  $X$ :

$$AGG\_Y[X] \Leftarrow Y[X'] \text{ WHERE } R_1(X_1, X_2), \dots, R_{\ell-1}(X_{\ell-1}, X_\ell) \quad (32)$$

For example, for unifying the treated and response units associated to  $T[X] = Prestige[A]$  and  $Y[X'] = Score[S]$ , the aggregate rule associated with the relational path in (31) coincides with (29).

Therefore, from now on, we assume that the response units  $\mathbb{U}_Y$  are the same as the treated unit  $\mathbb{U}_T$ . Henceforth, we simply refer to elements of  $\mathbb{U}_Y$  and  $\mathbb{U}_T$  as *units* and denote them with  $\mathbb{U}_{(T, Y)} = \mathbb{U}_T = \mathbb{U}_Y$ .

**Relational Peers.** Next, we define the notion of relational peers of a unit which is central to the notion of relational and isolated effects. Recall that the grounded causal graph

$G(\Phi^\Delta)$  is extended with vertices and edges corresponding to aggregated attributes as discussed in Section 4.3.

**Definition 4.2.** Given a treated attribute function  $T[X]$ , and a (possibly aggregated) response attribute function  $Y[X]$ , we define the *relational peers* of a unit  $\mathbf{x} \in \mathbb{U}_{(T, Y)}$  as a set of units  $\mathbb{P}(\mathbf{x}) \subseteq \mathbb{U}_{(T, Y)} - \{\mathbf{x}\}$  such that for each  $\mathbf{x}' \in \mathbb{P}(\mathbf{x})$ , there exists a directed path from  $T[\mathbf{x}']$  to  $Y[\mathbf{x}]$  in  $G(\Phi^\Delta)$ .

For example, for  $T[X] = Prestige[A]$  and  $Y[X'] = AVG\_Score[A]$ ,  $\mathbb{P}(\text{"Bob"}) = \{\text{"Eve"}\}$  and  $\mathbb{P}(\text{"Eve"}) = \{\text{"Bob"}, \text{"Carlos"}\}$ . In practice, the relational causal model is expected to form relational peers  $\mathbb{P}(\mathbf{x})$  such that they consist only of units that are in some *relational proximity* of  $\mathbf{x}$ , e.g., authors from the same institution, same research interests, etc.

The following quantity measures the expected response of a unit  $\mathbf{x} \in \mathbb{U}_{(T, Y)}$ , when it receives the treatment  $t$  and its relational peers receive the vector of treatments  $\vec{t} = (t_1, t_2, \dots)$ .

$$Y_{\mathbf{x}}(t, \vec{t}) \stackrel{\text{def}}{=} \mathbb{E}[Y[\mathbf{x}] \mid \text{do}(T[\mathbf{x}] = t), \text{do}(T[\mathbb{P}(\mathbf{x})] = \vec{t})] \quad (33)$$

In this paper, we assume  $\text{do}(T[\mathbb{P}(\mathbf{x})] = \vec{t})$  is a *well-defined intervention* for all units  $\mathbf{x}$ , i.e., it uniquely determines which relational peers of a unit would receive which treatment. For instance, this holds if  $\mathbb{P}(\mathbf{x})$  is of the same size for all  $\mathbf{x}$ , and either has a natural ordering or is ordering-invariant. However, as we discuss after (23), we allow several relaxations on the size and type on  $\vec{t}$  in our framework.

**Overall, Isolated and Relational Effects.** The query (22) in CaRL, compute the following quantities that compare the overall, isolated and relational effects of two alternative intervention strategies  $(t, \vec{t})$  and  $(t', \vec{t}')$ :

- The average *isolated causal effect*:

$$AIE(t; t' \mid \vec{t}) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{\mathbf{x} \in \mathbb{U}_{(T, Y)}} Y_{\mathbf{x}}(t, \vec{t}) - Y_{\mathbf{x}}(t', \vec{t}) \quad (34)$$

- The average *relational causal effect*:

$$ARE(\vec{t}; \vec{t}' \mid t) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{\mathbf{x} \in \mathbb{U}_{(T, Y)}} Y_{\mathbf{x}}(t, \vec{t}) - Y_{\mathbf{x}}(t, \vec{t}') \quad (35)$$

- The average *overall causal effect*:

$$AOE(t, \vec{t}; t', \vec{t}') \stackrel{\text{def}}{=} \frac{1}{n} \sum_{\mathbf{x} \in \mathbb{U}_{(T, Y)}} Y_{\mathbf{x}}(t, \vec{t}) - Y_{\mathbf{x}}(t', \vec{t}') \quad (36)$$

Intuitively, the isolated causal effect of a treatment fixes the treatment of the relational peers of a unit and compares its expected response under two treatment strategies assigned to the unit. On the other hand, the relational causal effect of a treatment fixes the treatment of a unit  $\mathbf{x}$  and compares its expected response under two treatment strategies assigned to its relational peers. For example, the relational effect of  $Prestige[A]$  on  $AVG\_Score[A]$  fixes the prestige of an author such as “Bob” and compares the counterfactual response  $AVG\_Score[\text{"Bob"}]$  under two regimes of interventions in



which the relational peers of “Bob”, e.g., “Eve”, receive two different treatment strategies, e.g., all of them have prestigious affiliation versus non-of them have prestigious affiliations. Notice that the overall causal effect is an extension of ATE (27) for two arbitrary treatment strategies. Indeed, ATE coincides with AOE(1,  $\vec{1}$  | 0,  $\vec{0}$ ), when the treated and response units are unified. The following proposition shows the connection between relational, isolated and overall effects.

**PROPOSITION 4.3.** The average overall effect can be decomposed into the average isolated and average relational effects as follows:

$$\begin{aligned} \text{AOE}(t, \vec{t}; t', \vec{t}') &= \text{AIE}(t, t' | \vec{t}) + \text{ARE}(\vec{t}, \vec{t}' | t') \\ &= \text{AIE}(t, t' | \vec{t}') + \text{ARE}(\vec{t}, \vec{t}' | t) \end{aligned} \quad (37)$$

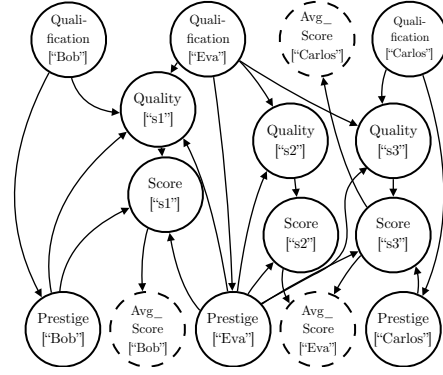
**Relaxations on the treatment vector to the peers.** In the definitions of the average isolated, relational, and overall causal effects in (34), (35), and (36), we do not need the treatment vectors  $\vec{t}, \vec{t}'$  applied to the peers to have the same size although they are applied to all units  $\mathbf{x}$ , as well as we do not need all units  $\mathbf{x}$  to have the same number of peers in  $\mathbb{P}(\mathbf{x})$ . As the grammar defined in (23) describes, we can assign treatments to “at least/most  $k$  or  $k\%$ ” neighbors, and that is well-defined for all units  $\mathbf{x}$  even if they do not have the exact same number of peers in  $\mathbb{P}(\mathbf{x})$ . On the other hand, for such conditions, we do need to assume that the effects of interventions to the peers are *ordering-invariant*, e.g., the intervention can be applied to any of the  $k$  peers (with possible truncations for smaller peer sets) in  $\mathbb{P}(\mathbf{x})$ .

## 5 ANSWERING CAUSAL QUERIES

The query answering component of CaRL consists of the *covariate detection* (Section 5.1) and *covariate adjustment* (Section 5.2). The goal of covariate detection is to find a sufficient set of covariates that should be adjusted for to remove confounding effects. Then in the process of covariate adjustment, the data is transformed into a flat single-table format such that adjustment can be performed using standard methods.

### 5.1 Covariate Detection

Given a treatment and a response attribute functions  $T[\mathbf{X}]$  and  $Y[\mathbf{X}']$ , the quantity  $\Pr(Y[\mathbf{x}'] = y | \text{do}(T[\mathbb{S}] = \vec{t}_{\mathbb{S}}))$  for a response unit  $\mathbf{x}' \in \mathbb{U}_Y$  and a subset of treated units  $\mathbb{S} \subseteq \mathbb{U}_T$ , is primitive to all causal queries defined in Section 4. For example, ATE (27) considers  $\mathbb{S} = \mathbb{U}_T$ , whereas AIE (34), ARE (35), and AOE (36) consider  $\mathbb{S} = \{\mathbf{x}\} \cup \mathbb{P}(\mathbf{x})$ . Our goal is to use the assumptions encoded in a GCD  $G(\Phi^\Delta)$  to rewrite  $\Pr(Y[\mathbf{x}'] | \text{do}(T[\mathbb{S}] = \vec{t}_{\mathbb{S}}))$  into an equivalent probabilistic formula in terms of a set of carefully chosen *observed covariates* to adjust for. If such rewriting exists, then



**Figure 5: Extending the GCD in Fig 4 with the aggregated attribute  $\text{AVG\_Score}[A]$ .**

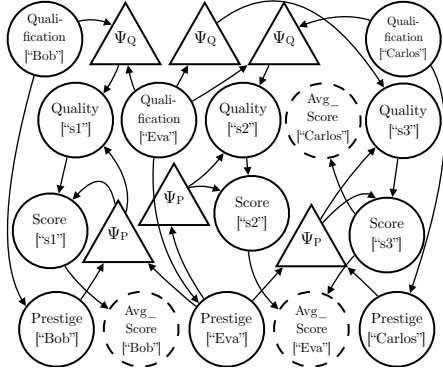
$\Pr(Y[\mathbf{x}'] = y | \text{do}(T[\mathbb{S}] = \vec{t}_{\mathbb{S}}))$  can be estimated from data, thereby the causal queries can be answered.

First, we extend a GCD  $G(\Phi^\Delta)$  to account for the embedding functions  $\Psi^A = \{\psi_{A_1}^A, \psi_{A_2}^A, \dots\}$  for each  $A \in \mathbf{A}$ , that map the parents  $\text{Pa}(A(\mathbf{x}))$  of  $A(\mathbf{x})$  in GCD  $G(\Phi^\Delta)$  into a low-dimensional embedding space as given in (??). For instance, in Example 5.1,  $\psi_{\text{Prestige}}^{\text{Score}}$  embeds the *Prestige* of the (possibly varying number of) authors of submissions into vectors with a fixed dimension for all submissions. The motivation is to enable the search for a low-dimensional embedding representation of covariates, which in turn enables efficient query answering from finite data (see Section 5.2). To this end, with abuse of notation, we extend the set of attribute functions  $\mathbf{A}$  with a collection of *relational embedding* attribute functions  $\Psi^A = \{\psi_{A_1}^A[X], \psi_{A_2}^A[X], \dots\}$  for each  $A \in \mathbf{A}$ , such that for each  $\mathbf{x} \in \mathbb{U}_A$ ,  $\psi_{A_i}^A[\mathbf{x}]$  corresponds to the result of applying the mapping  $\psi_{A_i}^A$  to a subset of  $\text{Pa}(A(\mathbf{x}))$  consisting of the grounding of  $A_i$ . Then, we augment a GCD  $G(\Phi^\Delta)$  by inserting the grounding of the relational embeddings as intermediate vertices between a ground atom and its parents.

**Example 5.1.** In Example 3.6, the relational relational embedding attribute  $\psi_{\text{Prestige}}^{\text{Score}}(S)$  now corresponds to the result of applying the function  $\psi_{\text{Prestige}}^{\text{Score}}$  to map the *Prestige* attributes of the authors of a submission (of the *Author* entities) to a new attribute of the corresponding *Submission* entities (see Figures ??, 4, 5 for the input GCDs, and Figure 6 for the output GCD with augmented attributes).

We show the following theorem in the technical report.

**THEOREM 5.2 (RELATIONAL ADJUSTMENT FORMULA).** Given an augmented GCD  $G(\Phi^\Delta)$ , the quantity  $\Pr(Y[\mathbf{x}'] = y | \text{do}(T[\mathbb{S}] = \vec{t}_{\mathbb{S}}))$



**Figure 6: Augmenting the GCD in Fig 5, for clarity**  $\psi_{Qualifications}^{Quality}[S]$  represented as  $\psi_Q$  and  $\psi_{Prestige}^{Quality}[S]$  and  $\psi_{Prestige}^{Score}[S]$  compactly represented as  $\psi_P$

$\vec{t}_S$ ) is given by the following *relational adjustment formula*:

$$\Pr(Y[x'] = y \mid \text{do}(T[S] = \vec{t}_S)) = \sum_{z \in \text{Dom}(Z)} \Pr(Y[x'] = y \mid Z = z, T[S'] = \vec{t}_{S'}) \Pr(Z = z) \quad (38)$$

Where  $S' \subseteq S$  is such that for each  $x \in S'$  there exists a directed path from  $T[x]$  to  $Y[x']$  in  $G(\Phi^A)$ , and  $Z$  is set of vertices in  $G(\Phi^A)$  corresponded to the groundings of a subset of  $A_{Obs}$  (the observed attribute functions) such that:

$$(Y[x'] \perp\!\!\!\perp \bigcup_{x \in S} \text{Pa}(T[x]) \mid_{G(\Phi^A)} \bigcup_{x \in S} T[x], Z) \quad (39)$$

Notice that the relational adjustment formula ignores the treated units in  $S \setminus S'$ , simply because,  $\Pr(Y[x'] = y \mid \text{do}(T[S] = \vec{t}_S)) = \Pr(Y[x'] = y \mid \text{do}(T[S'] = \vec{t}_{S'}))$ . The following observation immediately follows:

**OBSERVATION 5.1.** *The  $d$ -separation condition (39) always satisfied for  $Z = \bigcup_{x \in S} \text{Pa}(T[x])$  (any subset of vertices  $d$ -separates itself from the rest of the vertices). Hence, adjusting for the joint distribution of the parents of the treated atoms is always sufficient for identification. However, if some of the parents are unobserved or redundant for adjustment, one can use the condition (39) to select a minimal set of observed covariates that are sufficient for identification.*

We illustrate with an example.

**Example 5.3.** To compute  $ATE(Prestige, Score)$  in our toy example, we need to compute quantities of the form

$$\Pr(\text{Score}[s] = 1 \mid \text{do}(\text{Prestige}[\{\text{"Bob"}, \text{"Eve"}, \text{"Carlos"}\}] = \vec{1})) \quad (40)$$

where we intervene on all three authors in the example. By applying Theorem 5.2 and Observation 5.1 for submission

$s = \text{"s}_1\text{"}$  and  $Z = \{\text{Qualifications}[\text{"Bob"}], \text{Qualifications}[\text{"Eve"}]\}$  we obtain,

$$(40) = \sum_{z \in \text{Dom}(Z)} \Pr(\text{Score}[\text{"s}_1"] = 1 \mid Z = z, \text{Prestige}[\{\text{"Bob"}, \text{"Eve"}\}] = \vec{1}) \Pr(Z = z) \quad (41)$$

Similarly, for  $s = \text{"s}_2\text{"}$  and  $Z = \{\text{Qualifications}[\text{"Eve"}]\}$  we obtain,

$$(40) = \sum_{z \in \text{Dom}(Z)} \Pr(\text{Score}[\text{"s}_2"] = 1 \mid Z = z, \text{Prestige}[\{\text{"Eve"}\}] = \vec{1}) \Pr(Z = z) \quad (42)$$

For estimating  $ATE(Quality, Score)$  (assuming quality is observed), by applying (39) we obtain that for each submission  $s$ ,  $\Pr(\text{Score}[s] = 1 \mid \text{do}(\text{Prestige}[U] = \vec{1}))$  can be estimated by adjusting for the embedded attribute functions  $Z = \{\psi_{Prestige}^{Score}[s], \psi_{Qualifications}^{Quality}[s]\}$ .

For estimating  $ATE(Quality, AVG\_Score)$  (the effect on average acceptance rate of an author) we need to estimate  $\Pr(AVG\_Score[A] = y \mid \text{do}(\text{Prestige}[U] = \vec{1}))$ , for each author. According to Theorem (39), this can be done by adjusting for the joint distribution of the qualifications of *all* of their past coauthors, which is potentially very high-dimensional.

---

#### Algorithm 1: Constructing unit-table.

---

**Input:** An augmented GCD  $G(\Phi^A)$ , a treated and outcome attribute functions attribute function  $T[X]$  and  $Y[X']$ .

**Output:** The unit table  $D(Y, H_T, H_Z)$

```

1 for  $x' \in \mathbb{U}_Y$  do
2    $\mathbb{U}'_T \leftarrow$  A minimal subset of  $\mathbb{U}_T$  such that there exists a
   directed path in  $G(\Phi^A)$  from  $T[x]$  to  $Y[x']$ , for all  $x \in \mathbb{U}'_T$ 
3    $Z \leftarrow$  A minimal set of vertices in  $G(\Phi^A)$ 
   that satisfies  $d$ -separation statement in Eq (39)
4    $\psi_T \leftarrow \psi_Y^T(\langle T[x_1], \dots, T[x_{|\mathbb{U}'_T|}] \rangle)$ 
5    $\Psi_Z \leftarrow \Psi_Y^Z(Z)$ 
6   Insert the tuple  $(Y[x], \psi_T[x], \Psi_Z[x])$  to unit table  $D$ 
```

---

## 5.2 Covariate Adjustment

This section addresses the following challenges that arise in estimating the causal queries in Section 4 using the relational adjustment formula (38) (see Section 5.1).

(1) *High dimensional distributions.* The first challenge is that the relational adjustment formula (38) consists of two potentially high-dimensional distributions  $\Pr(Y[x] = y \mid Z = z, T[S] = \vec{t}_S)$  and  $\Pr(Z = z)$  that are unknown and must be estimated from limited data (see Example 5.3).

(2) *Variability of treatment and covariate vectors.* The causal queries need to compute *averages* across all response units. Hence, we need to estimate  $\Pr(Y[x'] = y \mid \text{do}(T[S] = \vec{t}_S))$  for all  $x' \in \mathbb{U}_Y$ . This is challenging since Theorem 5.2 can lead to

Unit	Outcome (Y)	Embedded Authors' Treatments ( $\psi_T^Y$ )		Embedded Authors' Covariates ( $\psi_Z^Y$ )	
Submission ID	Score	Prestige (AVG)	Authors (COUNT)	Experience (AVG)	H-index (AVG)
$S_1$	0.8	0.5	2	6	30
$S_2$	0.15	0.2	1	2	10
$S_3$	0.2	0.5	2	5	15

**Table 1: The unit table for  $T[X] = \text{Prestige}[A]$  and  $Y[X'] = \text{AVG\_Score}[A]$ . The unit table for  $T[X] = \text{Prestige}[A]$  and  $Y[X'] = \text{Score}[A]$ .**

Unit	Outcome (Y)	Embedded Collaborators' Treatments ( $\psi_T^Y$ )		Embedded Collaborators' Covariates ( $\psi_Z^Y$ )	
Author ID	AVG_Score	Prestige (AVG)	Centrality (COUNT)	Experience (AVG)	H-index (AVG)
Bob	1	1	1	2	10
Carlos	0	1	1	2	10
Eve	0.5	2	0.5	9	35

different probabilistic formulas for different response units. Then even under the homogeneity assumption, the corresponding probabilistic formulas must be estimated for each  $\mathbf{x}' \in \mathbb{U}_Y$ , which is not feasible. For instance, in Example 5.3 (41) and (42) must be estimated separately.

To address the aforementioned issues, similar to Section 3, we use a set of embedding functions  $\psi_T^Y$  and  $\psi_Z^Y$ , respectively to project the treatment and covariates vectors into a low-dimensional embedding space with *fixed dimensionality* and for all response units. Note that it implies from the homogeneity assumption in Section 3 that the embedding representation of the treated and covariates are *identically distributed*. This enables us to transform a relational instance to a low-dimensional flat-table that suffices for query answering.

**Unit-Table.** Given a GCD  $G(\Phi^A)$ , a treatment and response attribute functions  $T[X]$  and  $Y[X']$ , we use Algorithm 1 to construct a *unit table*, which is a standard relation with schema  $D(Y, \psi_T^Y, \psi_Z^Y)$ , consists of tuples  $(Y[\mathbf{x}'], \psi_T^Y[\mathbf{x}'], \psi_Z^Y[\mathbf{x}'])$  for each response unit  $\mathbf{x}' \in \mathbb{U}_Y$ , where  $\psi_T^Y[X'], \psi_Z^Y[X']$  (with abuse of notation) are relational embedded attribute functions corresponded to the result of applying  $\psi_T^Y, \psi_Z^Y$ , respectively to the treatment and covariate vectors.

*Example 5.4.* Table ?? shows the unit-table corresponding to  $T[X] = \text{Prestige}[A]$  (the *submissions* constitute the response units) and  $Y[X'] = \text{Score}[S]$ , whereas Table 1 shows the unit-table associated with  $T[X] = \text{Prestige}[S]$  and  $Y[X'] = \text{AVG\_Score}[A]$  (the *Authors* constitute the response units). In these two tables, simple mappings as average and the number of treated units are used for embedding.

By rewriting the RHS of the relational adjustment formula (38) in term of the attributes of the unit table and  $\vec{t}_{\mathbf{S}'}^e$ , the embedded representation of the treatment assignment  $\vec{t}_{\mathbf{S}'}^e$ , i.e.,  $\vec{t}_{\mathbf{S}'}^e = \psi_T^Y(\vec{t}_{\mathbf{S}'}^e)$ , we obtain:

$$\sum_{z \in \text{Dom}(\psi_Z^Y)} \Pr(Y = y \mid \psi_Z^Y = z, \psi_T^Y = \vec{t}_{\mathbf{S}'}^e) \Pr(\psi_Z^Y = z) \quad (43)$$

Now, all causal queries in Section 4 can be estimated using standard techniques in statistics for estimating (43) using the unit table  $D$ .

Dataset	Att. [#]	Rows [#]	Embedding Choice			Query Ans.
			Predefined	RNN	Forest	
REVIEWDATA	7	6K	0.2	480	2.5	180
SYNTHETICDATA	7	300K	11	24,000	x	400

**Table 2: Runtime in seconds for experiments.**

Note that the idea of embedding addresses both issues of high dimensionality and variable size of the treatments and covariates corresponding to the response units, thereby makes the estimation of all causal queries feasible. However, (43) only approximates (38), hence the quality of the answers depend on whether the embeddings preserve sufficient statistics. In Section ??, we develop several techniques to learn this embeddings from data.

## 6 EXPERIMENTS

This section presents experiments that evaluate the feasibility and efficacy of CaRL. We aim to answer the following questions: **Q1:** What is the end-to-end performance of CaRL? In particular, can the system (1) avoid confusing correlation with causation; and (2) consider the relational structure in multi-relational settings rather than naively applying traditional causal inference methods on the *universal table* obtained by joining all base relations? **Q2:** What is the sensitivity of query answering component of CaRL to different embedding methods? Additionally, we report the performance of CaRL in Table 2.

### 6.1 Setup

This section describes the datasets and causal queries we used in the experiments.

*Real Dataset.* The REVIEWDATA consists of 2,075 papers submitted for review across 10 computer science conferences and workshops. Each submission is associated with a number of referee reviews, as well as an acceptance or rejection decision. About half of all submissions are double-blind, while in the other half the names of the authors are revealed. All submissions have been unblinded after the conferences concluded. The dataset also contains an authors table, with the citation count, h-index, publishing experience (in years)

and university ranking for each of the 4490 authors that contributed to a paper in the dataset.

*Synthetic Dataset.* **Babak:** [**@Harsh**] To answer both questions we needed the ground truth, so we generated SYNTHETICDATA data according to a slightly modified version of the relational causal model in Example 1.1. . Specifically, ... explain how you generated the relational skeleton .... and then ... given them the ground truth.

*Causal Queries.* We used both REVIEWDATA and SYNTHETICDATA to answer the following causal queries.

$$\text{Score}[S] \Leftarrow \text{Prestige}[A] ? \quad (44)$$

$$\text{Score}[S] \Leftarrow \text{Prestige}[A] ? \text{ WHEN LESS THAN 35\% } \langle \text{PEERS TREATED} \rangle \quad (45)$$

where, (44) quantifies *ATE* of *Prestige* on *Score*, and (45) quantifies the *ATE*, *ARE*, *AIE* effect of having less than and more than one third of the prestigious collaborator on submissions scores (cf. 4). The threshold of 1/3 was as chosen as it resulted maximum correlation between prestige and review score. However, in addition to *ATE*, *ARE*, *AIE* and *AOE* as defined in Section 4, we also report these quantities as per each level of covariates  $Z = z$ . We denote the corresponding conditional quantities, respectively as *CATE*, *CARE*, *CAIE* and *CAOE*. For instance *CATE* quantifies the *ATE* of prestige on submission scores as per each level of authors qualifications (the covariates). The motivation is to explore the variability of causal effects across different sub populations. Also note we used Pearson’s correlation coefficients to measure the association between the *Prestige* and *Score* and *AVG\_Score*.

Note that for estimating the *ATE*, since the REVIEWDATA is sparse (*i.e.*, there are only a few submissions with and without prestigious authors), we further aggregated the treatment vector associated to each authors (as in Figure1). Specifically, we considered submissions to be treated if 1/3 of the authors are from prestigious universities and untreated otherwise.).

**Babak:** [**@harsh: is there any sub-population for which the effect is significantly large? Then we should report that.**]

*Answering Queries and Embeddings.* For query answering we used Algorithm 1.1 for creating unit tables. The tables were generated using the same relational causal model as in Example 1.1. Then the unit tables used for estimating the relational adjustment formula as in (43) for both treatment vectors the queries concern about.. Note that we used Random Forest Regression for estimating the underling conditional probability distribution. For embeddings we used basic for aggregation functions such as mean and median together with the cardinality of the vectors to account for the underling topology of the relational skeleton (e.g, number of authors,

number of co authors). We also the vector of  $K$  moments for embedding, *i.e.*, the vector consist of mean, variance, skewness, etc., when we learn  $k$  automatically from data. In general, when the vectors are ordered or have temporal pattern the basic aggregations and moment summarization lead to information loss, which would negatively impact the quality of query answer. However, for experiment all these methods turned out to be effective. In future we plan to use techniques such as convolutional graph neural networks and Recurrent neural networks.

CaRL

covariate detection and creating the corresponding the unit table. For estimating the conditional probability distribution involved in the probabilistic formulas (43) we used Random Forest regression.

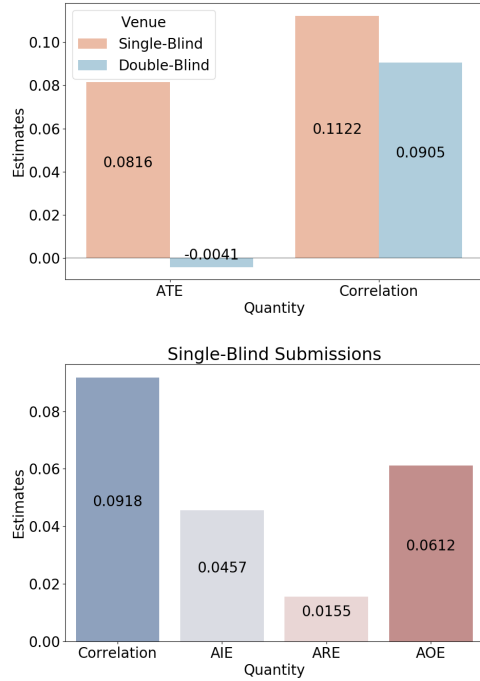
## 6.2 Q1: End-To-End Results

For this experiment we report the answers to our causal queries on both real and synthetic data. We used the the embedding method for computing these results.

**REVIEWDATA.** Figure 7(top), compares the *ATE* as in (44) and the naive correlation between the prestige and scores across the single and double-blind conferences. The results shows a significant correlation in both single and double blind conferences. However, after conducting a sound causal analysis using CaRL, it turned out that there is no significant causal effect double-blind conferences. However, the result shows a significant causal effect turned for single-blind conferences. This results suggest double-blind reviewing can be effective in reducing bias toward prestigious institutions

To further investigate the effect of prestige on review score in single-blind conferences, we used CaRL compute isolated and relational and overall effect of prestige on average review score as in (45). Figure 7(button) reveals that the isolated effect (*AIE*) is more significant than the relational effect (*ARE*), meaning that author’s own prestige have stronger effect on their average submission score than their collaborators prestige, which is expected. Furthermore, one can verify that approximately  $AOE = AIE + ARE$ , which confirms Theorem 4.3.

Notice the response variable underling the analysis in Figure 7(top) is the review score of each paper, whereas in Figure 7(bottom) is the average review scores for each author (corresponded to the tables in Figure 1). Therefore, the correlations in obtained in Figure 7(top) and Figure 7(bottom) are different. by naively interpreting these correlation as causation and might arrive at the wrong conclusion that double-blinding is not effective in reducing bias. While the validity

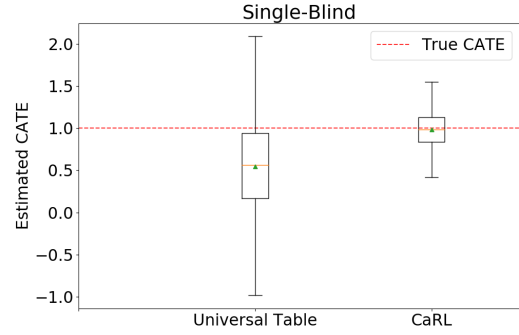


**Figure 7: Average treatment effects estimates and difference of sample means of treated and control units for single-blind and double-blind submissions.**

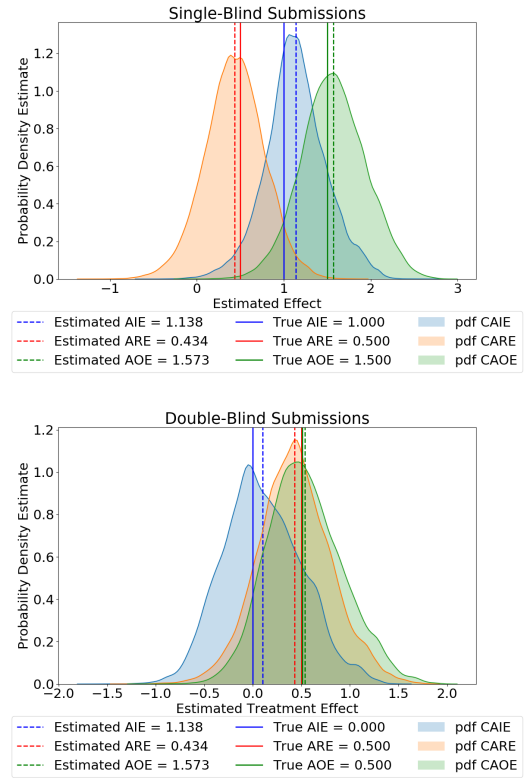
of results depend on the validity of the underlying assumption, are indeed insightful and goes way beyond associational analysis. We note that our results are in accordance with a series of controlled experiments that suggest double-blind reviewing reduces institutional prestige bias [3, 8, 12, 13].

**SYNTHETICDATA.** To further investigate the efficacy of CaRL, we used it to answer the same queries from SYNTHETICDATA for which we knew the ground truth. Figure 8, compares *CATE* (conditional ATE) of prestige on submission scores on SYNTHETICDATA computed using CaRL with the naive computation of CATE on the universal table obtained by joining all based relations. As shown, CaRL approximately recovers the ground truth. The variance is expected and is the product of limited data. However, naive causal analysis on the universal table leads to *incorrect ATE with huge variance*. This experiment reveals that ignoring the relational structure in relational domains can lead to incorrect conclusions.

Furthermore, Figure 9 compares the answers to both queries in (44) and (45) computed with CaRL with the ground truth. As shown in all cases CaRL approximately recovered the the ground truth. Note that the probability density functions shows **Babak: [ @Harsh ]**



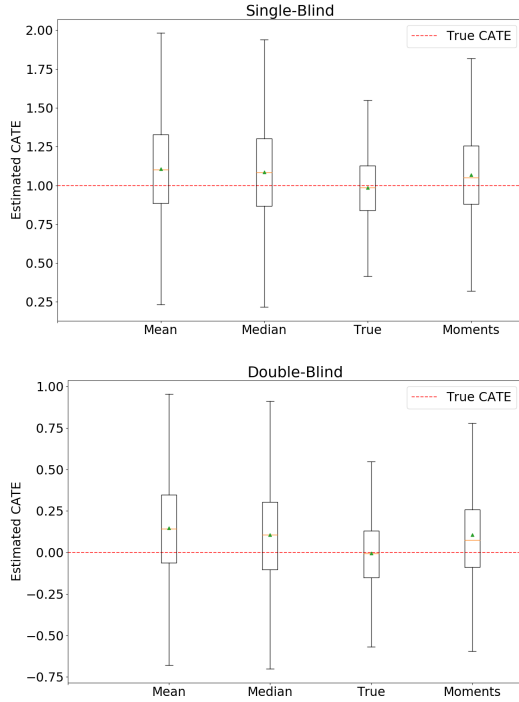
**Figure 8: Universal Table vs CaRL**



**Figure 9: Probability density function and corresponding averages for isolated treatment effect, relational treatment effect and total effects for (a) Single-Blind submission and (b) Double-blind submission.**

### 6.3 Q2: Sensitivity and Learnability of Embeddings

Figure 10 box-plots shows the distribution of estimated conditional average treatment effects (CATEs) under different choices of embedding functions for a) single-blind venue



**Figure 10: Comparing the sensitivity of the quality of query answer (CATE) to different choice of embeddings.**

submissions, and b) double-blind venue submissions. Here, “Mean” refers to using mean of the set as the hardcoded embedding of the set, “Median” refers to using median of the set as the hardcoded embedding of the set. “True” refers to using harmonic mean as the embedding for authors’ experience and weighted mean of the log of authors’ citation of each paper as the embedding, which are close to true generative model. The learned embedding uses optimal “Moments” of the set as the embedding. Here, we primarily observe that even the simple embeddings like mean or median are able to recover the true average treatment effect approximately. While the true embedding is able to recover almost the exact truth, the learned moment summarization embedding (which learns the optimal number of moments for accurate outcome prediction) also shows promising results. It is important to note that moment summarization is one of the most basic learning approach for set embedding and more complex approaches using recurrent neural networks or kernel density estimators can also be used. Across, both single and double-blind reviews estimates, we observe from Figure 10 that the estimated mean is almost approximately same for all methods while the variances for complex embedding is the smallest followed by the variance for the learned moment summarizing embedding.

Next, we generated the data using the second DGP where the review score has both a constant isolated effect of median prestige of the authors equal to 1 and a constant relational effect of mean prestige of the coauthors of the authors equal to 0.5.

## 7 RELATED WORK AND DISCUSSIONS

Statistical Relational Learning (SRL) aims at unifying logic and probability to model a joint probability distribution over relational data amenable for probabilistic reasoning [2]. A plethora of models such as PRMs [1], PSL, DAPER, PBNs, BLPs and MLNs has been proposed for SRL in the literature. While our proposed framework for Causal Relational Learning, CaRL, shares some similarities with some of these models, it is fundamentally motivated for causal reasoning. The closest SRL model to CaRL are PRMs and DAPER that extend the standard entity-relationship (ER) model to incorporate probabilistic dependencies. They can also be seen as a first-order extension of Bayesian Networks to relational data. However, there are two primary differences between CaRL and these models. First, PRMs are dependency models that essentially define a joint probability distribution over all possible completions of a partially specified instance. In contrast CaRL models the generative process of relational data. Here, the uncertainty is induced by hidden noise variables, *and prevails even when the entire relational data is available*. There has been prior attempt to use these models for causal inference from relational data, e.g.,

## 8 CONCLUSIONS

### TODOS:

- Do we need FKs and other constraints?
- Do we need to define ER model?



## REFERENCES

- [1] Lise Getoor, Nir Friedman, Daphne Koller, Avi Pfeffer, and Benjamin Taskar. Probabilistic relational models. *Introduction to statistical relational learning*, 8, 2007.
- [2] Lise Getoor and Ben Taskar. *Introduction to Statistical Relational Learning (Adaptive Computation and Machine Learning)*. The MIT Press, 2007.
- [3] Kanu Okike, Kevin T Hug, Mininder S Kocher, and Seth S Leopold. Single-blind vs double-blind peer review in the setting of author prestige. *Jama*, 316(12):1315–1316, 2016.
- [4] OpenReview. <https://openreview.net>.
- [5] Judea Pearl. *Causality*. Cambridge university press, 2009.
- [6] Judea Pearl et al. Causal inference in statistics: An overview. *Statistics Surveys*, 3:96–146, 2009.
- [7] Judea Pearl and Dana Mackenzie. *The book of why: the new science of cause and effect*. Basic Books, 2018.
- [8] Joseph S Ross, Cary P Gross, Mayur M Desai, Yuling Hong, Augustus O Grant, Stephen R Daniels, Vladimir C Hachinski, Raymond J Gibbons, Timothy J Gardner, and Harlan M Krumholz. Effect of blinded peer review on abstract acceptance. *Jama*, 295(14):1675–1680, 2006.
- [9] Donald B Rubin. *The Use of Matched Sampling and Regression Adjustment in Observational Studies*. Ph.D. Thesis, Department of Statistics, Harvard University, Cambridge, MA, 1970.
- [10] Semantic Scholar. [www.semanticscholar.org/](http://www.semanticscholar.org/).
- [11] Scopus. <https://www.scopus.com/>.
- [12] Richard Snodgrass. Single-versus double-blind reviewing: an analysis of the literature. *ACM Sigmod Record*, 35(3):8–21, 2006.
- [13] Andrew Tomkins, Min Zhang, and William D Heavlin. Reviewer bias in single-versus double-blind peer review. *Proceedings of the National Academy of Sciences*, 114(48):12708–12713, 2017.
- [14] Jeffrey D. Ullman. *Principle of Database Systems*. Pitman, 2nd edition, 1982.

## 9 SUPPLEMENTARY MATERIAL

This section includes the proof of theorems and propositions in the paper.

**PROOF OF PROPOSITION ??.** (Sketch) Markov compatibility  $\Pr(A^\Delta)$  and  $G(\Phi^\Delta)$  implied from the fact that  $\Pr(A^\Delta)$  satisfies *global Markov property*, i.e., for any set of ground atoms  $N$  that are non-descendants of a ground atom  $A[x]$  in  $G(\Phi^\Delta)$  it holds that  $(*) (A[x] \perp\!\!\!\perp N \mid \text{pa}(A[x]))$  (implied from the functional dependence between  $\text{pa}(A[x])$  and  $A[x]$  and the assumption of causal sufficiency. Furthermore, it follows from the assumption in Eq 25 that  $(**) (A[x] \perp\!\!\!\perp \text{Pa}(T[x]) \perp\!\!\!\perp \Psi^\Delta(\text{Pa}(A[x])))$ .

The factorization in Eq 26 immediately follows from the following independence obtained by applying the contraction axiom in Graphoid to  $(*)$  and  $(**)$ .

$$A(x) \perp\!\!\!\perp N \mid \Psi^\Delta(\text{Pa}(A[x])) \quad (46)$$

□

**PROOF OF THEOREM 5.2.** (Sketch) Since  $\Psi^\Delta$  are deterministic functions the ground random variables  $A^\Delta$ , they are also random variables, hence we can define the joint probability distribution  $\Pr(A^\Delta, \Psi^\Delta)$ . Also note that that the parents of

an atom in the augmented GCD  $\hat{G}(\Phi^\Delta)$  corresponded to the embedded parents of the same node in GCD  $G(\Phi^\Delta)$ .

Since each atomic intervention  $do(T[x_i] = t_i)$  modifies the augmented GCD  $\hat{G}(\Phi^\Delta)$  by removing the parents of  $T[x_i]$  from  $\hat{G}(\Phi^\Delta)$  (implied from the factorization Eq 26), the post intervention distribution  $\Pr(A^\Delta, \Psi^\Delta \mid do(T[\mathbb{S}] = \vec{t}_\mathbb{S}))$  can be obtained from the pre-intervention (observed) distribution  $\Pr(A^\Delta)$  by removing all factors  $\Pr(A(x), \mid \text{Pa}(A(x)))$ , from  $\Pr(A^\Delta, \Psi^\Delta)$  (cf. Eq 26), hence we obtain the following:

$$\Pr(A^\Delta, \Psi^\Delta \mid do(T[\mathbb{S}] = \vec{t}_\mathbb{S})) = \frac{\Pr(A^\Delta)}{\prod_{x \in \mathbb{S}} \Pr(A(x) \mid \Psi^\Delta(\text{Pa}'(A(x))))} \quad (47)$$

The following factorization implied by the chain rule of probability:

$$\Pr(A^\Delta, \Psi^\Delta) = \Pr\left(\bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])\right) \Pr\left(T[x_1] \mid \bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])\right) \quad (48)$$

$$\Pr\left(T[x_2] \mid T[x_1], \bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])\right) \quad (49)$$

$\vdots$

$$\Pr\left(T[x_i] \mid \bigcup_{j=0}^{i-1} T[x_j], \bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])\right) \quad (50)$$

$\vdots$

$$\Pr(\overline{A^\Delta}, \overline{\Psi^\Delta} \mid \bigcup_{x \in \mathbb{S}} \text{Pa}(T[x]), \bigcup_{x \in \mathbb{S}} T[x]) \quad (51)$$

where  $\overline{A^\Delta} \cup \overline{\Psi^\Delta}$  consist of all ground atoms in  $A^\Delta \cup \Psi^\Delta$  except for  $\bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])$  and  $\bigcup_{x \in \mathbb{S}} T[x]$ . The above factorization is valid if  $\bigcup_{x \in \mathbb{S}} \text{Pa}(T[x])$  and  $\bigcup_{x \in \mathbb{S}} T[x]$  are disjoint, which is true in acyclic  $\hat{G}(\Phi^\Delta)$ .

The following implied from the factorization in Eq. 51, the generic conditional independence in Eq. 46, and Eq 47, and marginalization we obtain the following:

$$\begin{aligned} \Pr(Y[x'] = y \mid do(T[\mathbb{S}] = \vec{t}_\mathbb{S})) &= \\ \sum_{\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])} \Pr(Y[x'] = y \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]), T[\mathbb{S}] = \vec{t}_\mathbb{S}) &\Pr\left(\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])\right) \end{aligned} \quad (52)$$

Now, given a set  $Z \subseteq \Psi^\Delta$ , we can rewrite the RHS of Eq. 52 into the following equivalent formula:

$$\begin{aligned}
\text{RHS} = & \sum_{\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])} \sum_{z \in \text{Dom}(Z)} \Pr(Y[x'] = y \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]), T[\mathbb{S}] = \vec{t}_{\mathbb{S}}, Z = z) \\
& \Pr(Z = z \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]), T[\mathbb{S}] = \vec{t}_{\mathbb{S}}) \Pr(\bigcup_{x \in \mathbb{S}} \text{pa}(T[x]))
\end{aligned} \tag{53}$$

Now from the conditional independence in Eq.39 and the conditional independence statements  $T[x] \perp\!\!\!\perp Z, \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]) \mid \text{pa}(T[x])$  (implied from the generic conditional independence in Eq.46) for each  $x \in \mathbb{S}$ , Eq ?? can be simplified as follows:

$$\begin{aligned}
\text{RHS} = & \sum_{\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])} \sum_{z \in \text{Dom}(Z)} \Pr(Y[x'] = y \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]), T[\mathbb{S}] = \vec{t}_{\mathbb{S}}, Z = z) \\
& \Pr(Z = z \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x]), T[\mathbb{S}] = \vec{t}_{\mathbb{S}}) \Pr(\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])) \\
= & \sum_{\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])} \sum_{z \in \text{Dom}(Z)} \Pr(Y[x'] = y \mid T[\mathbb{S}] = \vec{t}_{\mathbb{S}}, Z = z) \\
& \Pr(Z = z \mid \bigcup_{x \in \mathbb{S}} \text{pa}(T[x])) \Pr(\bigcup_{x \in \mathbb{S}} \text{pa}(T[x])) \\
= & \sum_{z \in \text{Dom}(Z)} \Pr(Y[x] = y \mid Z = z, T[\mathbb{S}] = \vec{t}_{\mathbb{S}}) \Pr(Z = z)
\end{aligned} \tag{54}$$

This completes the proof.  $\square$

**9.0.1 Data-Generative Process.** We now briefly explain the data generative process (DGP) for the synthetic data experiments and highlight the different levels of complexity. We generate each institution  $i$  in the set of Institutes using the following generator:

$$\text{Prestige}_i \sim \text{exponential}(\lambda = 0.1).$$

Similarly, an author  $a$  in set of Authors is generated using the following generative process:

$$\begin{aligned}
\text{Experience}_a & \sim \text{Bernoulli}(1/3) \times \text{Exponential}(3/4) \\
& + \text{Bernoulli}(1/2) \times \text{Normal}(20, 4) \\
& + \text{Bernoulli}(1/4) \times \text{Normal}(35, 1) \\
\text{Gender}_a & \sim \text{Bernoulli}(1/2) \\
\text{Citation}_a & \sim \text{Bernoulli}(1/2) \times \text{Poisson}(100) \\
& + \text{Bernoulli}(1/3) \times \text{Poisson}(200) \\
& + \text{Bernoulli}(1/5) \times \text{Poisson}(500) \\
& + \text{Bernoulli}(1/16) \times \text{Poisson}(1000) \\
\text{Expertise}_a & \sim \text{Multinomial}(1, \{1/10 \dots 1/10\}_{10}).
\end{aligned}$$

For each author  $a$  we generate their affiliation to an institute using following process. We use this affiliated institute's

prestige as author's prestige in our experiments:

$$\begin{aligned}
\text{AfP}(a, i) & = \text{expit}\left(\frac{\text{Prestige}_i \text{Citation}_a}{10}\right) \\
& + \frac{305}{\text{Prestige}_i(\text{Citation}_a + 1)} \\
& - \frac{305}{10} \frac{\text{Prestige}_i}{\text{Citation}_a + 1} \\
& - \frac{305}{\text{Prestige}_i} \frac{\text{Citation}_a}{\text{Experience}_a} \times \frac{30}{\text{Experience}_a} \\
\text{Affiliation}_a & \sim \text{Dice-Roll}(N = |\text{Institutes}|, \\
& P = \{\text{AfP}(a, i)\}_{i \in \text{Institutes}}) \\
\text{Prestige}_a & = \text{Prestige}_{\text{Affiliation}_a}
\end{aligned}$$

Similarly, we generate each conference  $c$  in set of conferences as follows:

$$\begin{aligned}
\text{Impact}_c & \sim \text{Exponential}(0.1) \\
\text{Area}_c & \sim \text{Multinomial}(1, \{1/10 \dots 1/10\}_{10}) \\
\text{SingleBlind}_c & \sim \text{Bernoulli}(2/3).
\end{aligned}$$

Finally, for each paper  $p$  in the set of paper submissions, the number of authors are randomly drawn from  $1 + \text{Exponential}(2/5)$  and then for the purpose of simulation we let the authors for each paper be randomly chosen without replacement from the set of authors. We also let each paper be submitted to any conference uniformly at random. This is just to keep the generation process less complicated while driving the major point home. Other parameters for the paper are generated as follows:

$$\begin{aligned}
\text{Quality}_p & = \text{expit}\left(\sum_{a \in \text{Authors}_p} e^{-2a} \log\left(\frac{30 \times \text{Citation}_a}{\text{Experience}_a + 5} + 1\right)\right) \\
\text{Review}_p^{(1)} & = 3 + (\text{SingleBlind}_{\text{Conference}_p} \\
& \times (\text{median}(\{\text{Prestige}_a\}_{a \in \text{Authors}_p}) > 10)) \\
& + \frac{20}{3} \left(\text{Normal}\left(\log(\text{Quality}_p + 1)\right.\right. \\
& \left.\left. - \log(\text{Impact}_{\text{Conference}_p})/10, 0.1 * \text{Quality}_p\right)\right).
\end{aligned}$$

We can observe from the generative model of the review score of a paper that for a submission in single blind conference the treatment effect of the median prestige of the authors more than 10 is 1 while for a submission in a double-blind conference, the analogous treatment effect is 0. In this data generation setup, there is no relational effect. In the following section, we will show CaRL's ability to recover true treatment effect under different embeddings of paper's authors' features.

We define a second possible data generator for review score of the paper where we have a positive relational effect:

$$\begin{aligned}
\text{Review}_p^{(2)} & = 3 + (\text{SingleBlind}_{\text{Conference}_p} \\
& \times (\text{median}(\{\text{Prestige}_a\}_{a \in \text{Authors}_p}) > 10)) \\
& + \frac{20}{3} \left(\text{Normal}\left(\log(\text{Quality}_p + 1)\right.\right. \\
& \left.\left. - \log(\text{Impact}_{\text{Conference}_p})/10, 0.1 * \text{Quality}_p\right)\right) \\
& + \frac{1}{2} (\text{mean}(\{\text{Prestige}_{a'}\}_{a' \in (\bigcup_{a \in \text{Authors}_p} \text{Coauthor}_a \setminus \text{Authors}_p)}) > 10).
\end{aligned}$$

In the above mentioned DGP, the isolated treatment effect for single blind submissions is 1 while it is 0 for double blind submission. However, the relational treatment effect of the mean prestige of coauthors being more than 10 is 0.5 for both single-blind and double-blind submissions.

PROOF.

$$\begin{aligned}
\text{AOE}(t, \vec{t}; t', \vec{t}') &= \frac{1}{n} \sum_{\mathbf{x} \in \mathbb{U}_{(T, Y)}} Y_{\mathbf{x}}(t, \vec{t}) - Y_{\mathbf{x}}(t', \vec{t}') \\
&= \frac{1}{n} \sum_{\mathbf{x} \in \mathbb{U}_{(T, Y)}} Y_{\mathbf{x}}(t, \vec{t}) - Y_{\mathbf{x}}(t', \vec{t}) + Y_{\mathbf{x}}(t', \vec{t}) - Y_{\mathbf{x}}(t', \vec{t}') \\
&= \text{AIE}(t, t' \mid \vec{t}) + \text{ARE}(\vec{t}, \vec{t}' \mid t')
\end{aligned}$$

Similarly the second equality in (37) holds.  $\square$