

An overview of the methodologies of causal discovery

This paper is part of our tech reports series. For more papers in the series see causalens.com/research. In our tech reports we share some of the most significant concepts in the causality literature and we point our readers to the most important publicly accessible academic references. The topics are chosen based on our frequent discussions with data science practitioners, and the reports are written primarily for a technical audience. While much of our original research forms part of our proprietary technology which we do not publicly disclose, these reports are part of our commitment to contribute to the wider research community and share the benefits of Causal AI.

Introduction

Until recently, discovering cause-and-effect relationships involved conducting a carefully controlled experiment or else relying on raw human intuition. Technological breakthroughs mean that AI can now help with causal discovery. Causal AI autonomously discovers causes in observational data, while also boosting human intuition and experimentation.

Causal knowledge

Causal knowledge of a system is formalized in terms of a *structural causal model* (SCM). An SCM $\{\mathcal{U}, \mathcal{V}, \mathcal{E}\}$ is composed of exogenous variables \mathcal{U} , endogenous variables \mathcal{V} , and a set of structural equations \mathcal{E} which describe the functional relationships between variables [Pea09]. Specifically, a structural equation relates the value of a target variable X with the values of those variables that have a directed edge terminating at X — its *parent* variables $\text{pa}(X)$. That is, $X = f(\text{pa}(X), \epsilon)$ for some function f , where ϵ is a noise variable reflecting the potentially stochastic nature of this relationship.

There are techniques for estimating these functional equations. However, in this brief report we focus on understanding the causal graph structure of the system. From a probabilistic perspective, this is reflected in the factorization of the joint probability across all variables $X_i \in \mathcal{U} \cup \mathcal{V}$:

$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i | \text{pa}(X_i)) \quad . \quad (1)$$

The key distinction between the two sets of variables is that the values of the exogenous variables \mathcal{U} are set by mechanisms extrinsic to the system of interest while endogenous variables \mathcal{V} are determined by variables (either in \mathcal{U} or \mathcal{V}) intrinsic to the system. For example, a model designed to predict share prices $V \in \mathcal{V}$ of companies in the United States may benefit from including an extrinsic variable reflecting the federal interest rate $U \in \mathcal{U}$.

The guiding aim of Causal AI is to understand the causes and effects amongst variables of interest in a system. The correlation between a pair of variables, X and Y , may

be due to one of three possible causal models:

- $X \rightarrow Y$ (X causes Y ; this can be an indirect causal relationship)
- $X \leftarrow Y$ (Y , directly or indirectly, causes X)
- $X \leftarrow Z \rightarrow Y$ (there is a common cause of both X and Y) .

Causal AI goes beyond statistical machine learning in distinguishing between these three scenarios leading to a more robust and flexible model of a system. There is no correlation without causation¹ and Causal AI acquires a more complete understanding of a system by identifying its causal structure.

How causal knowledge is acquired

Causal knowledge can be acquired in three, mutually beneficial, ways:

- through **interventional experimentation**.
- **causal discovery** from observational data.
- by integrating the knowledge of **domain experts**.

The technology established at causaLens facilitates all three approaches, which we briefly review in the following sections.

Interventional experimentation

Interventions form the basis of randomized controlled trials, the gold standard in clinical trials of new drugs, and of scientific experimentation more broadly. Consider two possible causal scenarios $X \rightarrow Y$ and $Y \rightarrow X$, relating variables X and Y . Only one of these is a true reflection of a system of interest and we would like to perform an experiment to determine which. Probabilistically, these scenarios correspond to two possible factorizations of the joint:

$$p(X, Y) = p(Y|X)p(X) \tag{2}$$

$$p(X, Y) = p(X|Y)p(Y) \tag{3}$$

Suppose we intervene and fix the X variable to a particular value $X = x$. Thus, we have performed a controlled experiment on the system. Consider the resulting factorized probabilities:

$$p(Y | \text{do}(X = x))p(\text{do}(X = x)) = p(Y | \text{do}(X = x)) \tag{4}$$

$$p(\text{do}(X = x) | Y)p(Y) = p(Y) \tag{5}$$

¹This is known as Reichenbach's *common cause principle*.

which we have identified using the *do*-operation from causal calculus [Pea09]. Note that the intervention $\text{do}(X = x)$ results in two different conditional probability distributions depending on the causal model. These distributions may be compared with experimental results in order to infer the true underlying causal structure. This simple case establishes the general principle that information regarding the causal structure of a system may be acquired by performing explicit *interventions* on the system. In the first scenario $X \rightarrow Y$, the distribution of Y varies with X . However, in the second $Y \rightarrow X$, it does not. Statistical machine learning does not distinguish between these two possibilities. Unfortunately, interventions may be either impossible or unethical and so we require a toolkit to perform causal discovery from data alone.

Causal discovery

Consider linear regression, which can be utilized to estimate the magnitude of the relationship between X and Y . Linear regression may be conceived as a structural equation with the independent variable $X \in \mathcal{U}$, the dependent variable $Y \in \mathcal{V}$, and the function $f \in E$ such that $y = f(x)$ is a linear function. Critically however, one must rely on domain knowledge or perform interventions in order to specify which variable is dependent. In contrast, causal discovery extracts the causal direction between variables automatically from data. Broadly-speaking, such methods are divided into two classes referred to as *constraint-based* and *score-based* [GZS19].

- **Constraint-based algorithms.** A causal structure implies a set of independence relations between variables. That is, if X and Y do not share a direct functional relationship or have a common parent variable, then they are independent. Constraint-based methods follow this logic by performing a sequence of statistical tests to determine the dependencies between variables and then constructing a causal graph by specifying the causal direction between dependent variables according to admissible orientation rules. Notable techniques such as the PC algorithm and Fast Causal Inference (FCI) fall into this category.
- **Score-based algorithms:** Rather than constructing a causal graph from local statistical dependencies, score-based methods search the space of graphs directly by evaluating the degree to which each graph can satisfactorily explain the observed data. The NOTEARS algorithm is one such approach which has inspired a proliferation of related techniques that leverage deep learning [VCB21]. Recent work at causaLens [KS21] identified a failure mode of NOTEARS which may be particularly problematic in applied settings.

Domain expertise

Humans are adept at identifying causal relationships. Translating such knowledge into structural equations can be a fruitful step in developing a formal theoretical model of

a system. Domain knowledge may be extracted under uncertainty or partial specification if the human is only partly knowledgeable of a certain phenomenon. Such partial understanding may still be used as a “prior” for further analysis. A key feature of the causaLens platform is to support the facility for domain experts to contribute their knowledge or belief regarding the causal structure of a target system, and to work in a collaborative feedback-cycle with Causal AI systems.

About causaLens

causaLens are the pioneers of Causal AI — a giant leap in machine intelligence. We build Causal AI powered products that empower humans to make superior decisions. We are creating a world in which humans can trust machines with the greatest challenges in the economy, society, and healthcare.

Bibliography

- [GZS19] Clark Glymour, Kun Zhang, and Peter Spirtes. “Review of causal discovery methods based on graphical models.” In: *Frontiers in genetics* 10 (2019), page 524.
- [KS21] Marcus Kaiser and Maksim Sipos. “Unsuitability of NOTEARS for Causal Graph Discovery.” In: *arXiv* (2021), page 2104.05441v1.
- [Pea09] J Pearl. *Causality*. 2009.
- [VCB21] Matthew J Vowels, Necati Cihan Camgoz, and Richard Bowden. “D’ya like DAGs? A Survey on Structure Learning and Causal Discovery.” In: *arXiv preprint arXiv:2103.02582* (2021).