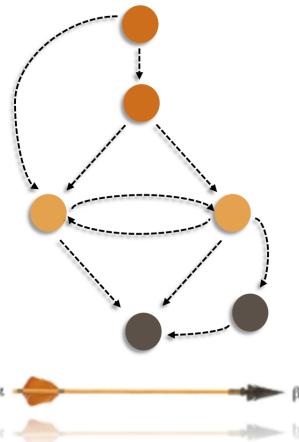


A Primer on Causal Diagram Learning

Interpreting Causation from the Causal Discovery Perspective

Xuanzhi CHEN

January 9, 2024, Chinese Lunar New Year



*"As X-rays are to the surgeon,
causal diagrams are for causation."*

— Judea Pearl
Father of Bayes Nets and Causal Models

Contents

1 Abstract	1
2 Introduction	2
2.1 What I Ready to Tell You in This Paper	2
2.2 Frequently Asked Questions	2
3 Revolutionary Causation	3
3.1 Is Human Influence to Blame for Global Warming?	3
3.1.1 Application: Counterfactuals	3
3.1.2 Foundation: Use Diagrams to Unravel Causation Scenarios	9
3.2 Can Vaccinal Treatments Save Lives of Loved Ones?	12
3.2.1 Application: Intervention	12
3.2.2 Foundation: Use Diagrams to Deduce Causal Inquiry	17
4 Inferred Causation	21
4.1 Causal Diagram Learning = Data + Restriction	21
4.1.1 Why Does the Equation Harbor Causal Significance?	21
4.1.2 Restrictions Give Rise to Identification	26
4.2 Challenges: Unmeasured Common Cause	32

1 Abstract

Causal diagram learning, namely the inference of "structural" causation from raw data (e.g. drawing an arrow (cause → effect) to specify a causal structure), usually serves as a blueprint that can be applied to mathematically model the real world causality.

Majoring in computer science, I tended to develop algorithms that are able to learn the causal diagram given the statistical patterns from data. I was pondering, however, I have spent much time in coding algorithms and fine tuning experiment results; whereas I actually know little about the profoundly broad ideas underneath the philosophy of *Causal Discovery*. In some ways, algorithms are productions and implementations of the idea of causality. From a bigger perspective, it is also pivotal to understand the reason why pioneers in causation have been dedicated to inferring the causal diagram from data.

In light of my research experiences, is it possible for me to write a paper that attempts to provide an initial interpretation of the pivotal roles of "causal diagram learning" by standing upon a broad view of causation?

Since I was motivated by fundamental ideas of causality from a popular science book named "*The Book of WHY*", I attempt to further create connections in the book to other celebrated books by leaders in causation. Centering around subjects of "causal diagram learning", three of celebrated books herein include: *Causation, Prediction, and Search* [2001] (Peter Spirtes, Clark Glymour); *Causality* [2008] (Judea Pearl); *Elements of Causal Inference* [2017] (Jonas Peters, Bernhard Schölkopf, etc). Hence, I want to show that ideas of *Causal Discovery* are deeply rooted in the past by standing on the shoulders of these giants.

Finally, ideas of causal diagram learning introduced in this paper is the tip of the iceberg. Furthermore, I hope that points of view from the work of the giants in causal science might be beneficial to future progression of AI, cognitive neuroscience, and other related fields.

Videos of the relative presentation can be found online¹.

Acknowledgement

The book "*Introduction to Causal Inference*" by *Brady Neal*, with its awesome course, deserves acknowledgement for forming the "spirit" of this paper: popularization of *causal science*.

I am grateful for the *DMIR lab*² (*Data Mining and Information Retrieval*) that offered me research opportunities in *Causal Discovery*, with special thanks to *Ruichu Cai* of the lab director and *Dongning Liu* of the dean in School of Computer, GDUT.

I owe a great debt to my advisor *Wei Chen*. I would not have completed this article during my busy graduation season but for remembering the encouragement from her when I first started studying causation two years ago — "Do it, just have your own interest of research and your own rhythm of lifetime".



Figure 1: Four of the celebrated "causal books" by leaders in causal science, providing the relevant ideas (as to causal diagram learning) that will be discussed in this paper.

NOTES¹:

YouTube: link TBD.
Bilibili: link TBD

REFERENCE²:

The DMIR Lab was established in Guangdong, China, dedicating to building the lab into an influential research center at home and abroad. The DMIR lab have been focusing on non-temporal and temporal causal inference over the last decade. The lab also has collaborated with front institutions such as CMU (Carnegie Mellon University) on causal discovery study.

— extracted from lab website

2 Introduction

2.1 What I Ready to Tell You in This Paper

The content of this paper are mainly divided into two sections: revolutionary causation and inferred causation¹. Partial ideas from celebrated causal books (*Causation, Prediction, and Search*; *Causality*; *Elements of Causal Inference*; *The Book of WHY*) that are beneficial to interpret causation in the paper are woven into the whole sections. In subsections, the content further consists of the "introductory part" and the "formalization part".

The first section will starts from the top of the "Ladder of Causation" (a central metaphor of *The Book of WHY*). Notions of *counterfactuals* and *intervention*, the paradigm of human causal thinking, are then introduced into two (simplified) applications of causal relation analysis in significant issues such as climate change and COVID-19. I attempt to have my readers' attention on fundamental roles that causal diagrams are playing therein.

After nearly a first half of the paper introducing the "downstream" capability of a given causal diagram, the second section will retrospectively focus on how to learn the causal diagram in the "upstream" task — namely the "*Causal Discovery*". I am trying to interpret the causal significance of the abstract principle behind classical causal inference. To approach concrete *causal discovery* methodology, we ultimately need to get a sense where specific restrictions are necessary to be understood.

2.2 Frequently Asked Questions

Couples of questions about the "instrumental function" of this paper are sorted into the FAQ form as the following.

Q: What relevant knowledge is required to understand causation?

A: In this paper, the introductory part mainly serves as a connection with *The Book of WHY*. If you have read the book, I hope there is no other prerequisites² required for this paper.

Q: How do you choose celebrated causal books based on the topic?

A: *The Book of Why* and *Causality* are representative work of Judea Pearl — the father of causality modeling. Authors and institutions of two rest books are major contributors in fields of *Causal Discovery*. For instance, you might probably heard of the TETRAD program developed by Carnegie Mellon University.

Q: What is the basic keynote when you write this paper?

A: I aim at popularization of causal science at first. It is more like an informing article if focusing on the "introductory part"³.

Q: Any other additional instructions?

A: Please forgive for my limitation of working proficiency in English writing (mandarin Chinese as the native language). I am sorry if some of the mathematical demonstration might be less rigorous and clear as well. I will try my best to illustrate the ideas that I deem are worthy in the rest of the paper.

REFERENCE¹:

Moral of "revolutionary causation" is inspired by the possible "second causal revolution" combining with the machine learning and AI system in the next decade. In the book of *Causality*, "inferred causation" refers to the fundamental intuition of causation analysis (causal discovery).

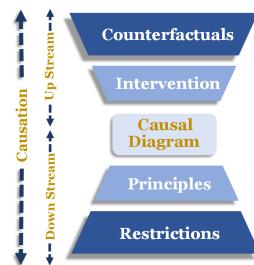


Figure 2: The outline of interpreting causation in this paper.

NOTES²:

However, basic probability and some topics of machine learning and statistics maybe used for the formalization part to elaborate details.

NOTES³:

Though the truth is that when I wrote the introductory part, I did virtually complete the "formalization part" beforehand.

3 Revolutionary Causation

3.1 Is Human Influence to Blame for Global Warming?

Debatable statements on the human responsibility for climate change have been existing for a long time. Why does focusing on describing a statement in front of major issues such as global warming become so important?

One way to effectively raise the public awareness may be enabling people to genuinely understand the cause-and-effect relationship behind the significant issues. In that sense, it may naturally require us to be able to describe the causation in **an obvious and concise way** in the first place, making it more acceptable and easy remembering for publics¹.

3.1.1 Application: Counterfactuals

Our brains are naturally the master at describing causation, and one of the remarkable capabilities to approach this is by **counterfactuals** thinking. Let us start our story with backgrounds lining up with the context² in *The Book of WHY*.

Conceptual Causation: Necessity and Sufficiency Causation

A heatwave occurring in France during August in 2003 has claimed the lives of hundreds of people. Dr. Allen, a meteorological scientist, used a metric called the "*fraction of attributable risk*" (FAR) to quantify how human influence or nature force can risk to global warming. Hoping to issue a scientific statement towards the public, Allen wrote the sentence at that point:

"It is 'very likely' that 'over half' the FAR of European summer temperature anomalies is attributable to human influence."

Scientist Allen essentially wanted to make an attribution — analysis of causes. the terms "very likely" and "over half", however, might seem a bit contradictory³. Since "*over half the FAR of human influence*" is a relatively **vague probability** (e.g. 60%, 70%, or 80%?), it might contrast with the extreme of that "*very likely the FAR*" (e.g. > 95%). What if instead, using the native accent of causation to describe this attribution probability?

Tricky senses exhibited from that statement are actually the results from tricky interaction between two of aspects of causation, namely **necessity causation** and **sufficiency causation**. When it comes to analyze the probability, this notion of causal relations yields respectively the **probability of necessity** (PN) and the **probability of sufficiency** (PS).

To this end, we can view as simple as an "imagination training" how PN and PS evaluate the human influence.

Take PN, given the fact that some people had lost their lives during the global heatwave, let us imagine a *counterfactual* world and ask ourselves: how much would we accept the likelihood that these people should have survived if the lethal heatwave did not happen? I suppose virtually we all tend to accept that likelihood — Our experience tells us that the lethal heatwave is

REFERENCE¹:

This motivation is in accord with what pearl has said in *The Book of WHY* that, quote, "Can ordinary people learn to understand the difference between necessary and sufficient causes? This is a nontrivial question. Even scientists sometimes struggle."

REFERENCE²:

I herein follow the example in "Counterfactuals", Chapter 8 of the book.

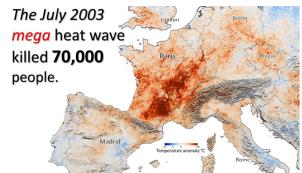


Figure 3: The 2003 lethal heatwave that occurred in France.

NOTES³:

On top of that, the term "European summer" indicates a broader metric of the average temperature in Europe over the entire summer, rather than specifying a singular scenario in France during August. We will further consider this "scenario based" angle in the second half.



Probability of Necessity (PN)

Figure 4: Those people should have survived if global warming did not happen. Normally it is unnecessary to concern about its threats, if not for the increasingly severe issue of greenhouse gases emission.

an extreme event that normally is not necessary to be placed concerns for its threats. But, the concern has become increasingly common only due to the greenhouse gases produced by humans. For this reason, we must assign a high value of *PN*, specifying a non-negligible responsibility of ourselves.

Manipulating a "*counterfactual*" imagination of *PN* based on our experience, is our unique talent; in terms of *PS* in the following, I would like to say it arguably involves our instinct.

Again, if we were asked to envision a more general case, for instance, about how strongly we do believe that global warming has the devastating power to sufficiently result in death. Based on our intuition, we would acknowledge the lethal effect of global warming without thinking twice, which also means to assign a high value of *PS*. Just think of how vulnerable the people in poverty are — people who are in lack of shelters and water will be impossible to resist when climate disaster happens.

Through the above "imagination training", we have seen how we are capable of drawing on our past experience and common sense in a rapid way. We have the gift: to articulate the analysis of causes in **an obvious and concise way**.

Conceptual Causation: Analysis of Causes Based on Scenarios

Recall Dr. Allens' formal statement that I mentioned at the beginning. Actually, the native language of causation can retain its simplicity but without compromising the accuracy implied by that statement.

Mathematically speaking, saying "*the very likely the FAR of human influence*" is essentially equal⁴ to saying "*the high PN of human influence*"; Interestingly, claiming "*the over half the FAR*" is still having something to do with claiming "*the high PS*". Consisting with our "imagination", both *PN* and *PS* are substantially in a high value, whereas the relatively **vague probability** of "*over half*" is mislead by the un conspicuous evidence. Which is to say, such deadly heatwave events in reality is rare (once-in-a-lifetime event), making it difficult to collect meteorological data over such a long period of time. The result? Scientists actually "calculate" a seemingly low value of *PS* without sufficiently evident data, diluting the effect of global warming and concluding our inadequate responsibility on heatwave victims.

In other words, it is possible that sufficiency causation can turn out to be "**insufficiency causation**⁵", as we dig deeper into a specific scenario such as "once-in-a-lifetime global heatwave". Let us add more details on the early example and imagine a more specific scenario. Suppose that both human influence and nature force, for instance, are causes of global warming. Assume that the industrial production associating with greenhouse gases emission (represented as "human influence") will impact upon government's political resistance on dealing with global warming. Also assume that the dense smoke caused by the widespread mega fires (represented as "nature force") will impede emission migration, which then aggravates the climate



Probability of Sufficiency (*PS*)

Figure 5: Global heatwave would have the devastating power that can sufficiently result in death, if the heatwave were occurring now.

TIPS⁴:

Under mild assumptions, I will show how *PN* becomes calculable, and how it is equal to the *FAR* metric in the formalization part.

REFERENCE⁵:

This is a motivation that readers can connect to the topic "insufficiency of necessary causation" in the book *Causality*, Chapter 10.

crisis. The controversy point is, though the government tends to impose restrictions on industrial production for its greenhouse gases emission issue, emerging techniques such as drones and satellites manufactured by the industries can also protect our climate via detecting signals of emission migration by wildfire — rendering the government to reconsider to relax the restrictions (and thus (indirectly) relax the restrictions on greenhouse gases emission, making the ambiguous human's responsibility for climate change).

Namely, neither of the causes can be ignored; yet neither of them takes sufficient "productivity" of its bad effects. If we must attribute to a single actual factor, which one should we pick⁶?

An analogical notion of "**sustenance-based causation**" will shed light on how to use environmental information as "the frame of reference" to define **actual causation**. In order to avoid "**insufficiency causation**" that focuses on the causality of "production" (able to bring effects), we can imagine another type of capability of "sustenance" (able to relatively maintain effects). This capability is described relative to the frame of reference in which the environmental information play an role, and the idea yields the *sustenance-based causation*, which amounts to be a **weak version of sufficiency causation**. The weakness refers to an alternative measurement of causal sufficiency, where we go from whether the cause can sufficiently "produce" its effects or not, to whether the cause can sufficiently "sustain" its effects.

For example, the nature force such as widespread mega fire would sometimes fail to pose threats towards global warming because human influence are powerful enough to immediately detect and extinguish the fire. In contrast, almost nothing could prevent the all-time human influence except for ourselves (to reduce greenhouse gases emission). **That asymmetry**, where mega fire is less likewise to continuously "sustain" (maintain) its causation towards climate devastation than human, suggests that human influence is exactly the *actual cause* to blame for global warming. Despite of the fact that occurrence of mega fire is time to coincide with industrial production activities (e.g. making drones and satellites), we can envision an accidental "freeze" of human industrial production to test whether mega fire could "sustain" its destructiveness or not. Picturing the "accidents" always help us understand causation better.

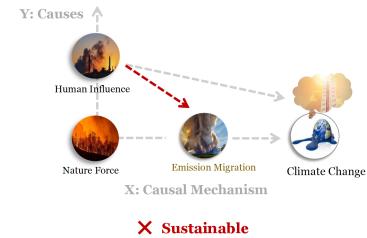
Aside for supplanting *sufficiency causation* with *sustenance causation*, we can add the connotation of *necessity causation* as well. We can describe the capability to "sustain necessary responsibility" as: global warming would not happen but for human influence, the same is truth when we have "frozen" mega fire. At the heart of it, *actual causation* encompass⁷ together the strength of *necessity causation* and the *sufficiency causation*, reaching a more comprehensive analysis of causes.



Figure 6: The causal diagram for a specific scenario: human influence, nature force, and climate change.

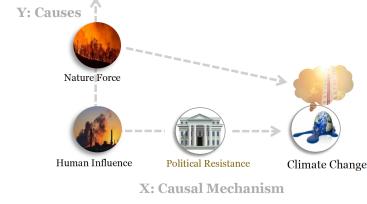
NOTES⁶:

If we pick greenhouse gases emission by human as an actual cause for climate change, it seems reasonable to argue that technologies by human also prevent climate from pollution by mega fire; That in turn means, if we deem air pollution by mega fire is the actual cause, then we find that its causal effects can be neutralized by human technologies as well.



X Sustainable

Figure 7: The nature force fails to "sustain" its causal effect.



✓ Sustainable

Figure 8: The human influence always "sustain" its causal effect, making it the actual cause in the specific scenario.

TIPS⁷:

In the formalization part, I will try to show a unity perspective, namely in form of the causal assertion, to combine the concepts as for necessity causation, sufficiency causation, and sustenance-based causation.

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of *Judea Pearl's* points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

Provided the example in the introductory part, let us start with introducing into the definition of *PN* relative to the causation of human influence ($X = x$)¹ on global warming ($Y = y$):

$$PN = P(Y_{x'} = y' \mid X = x, Y = y). \quad (1)$$

Where $Y_{x'} = y'$ is a standard "causation language", indicating " $Y = y'$ when intervening $X = x'$ (instead of observing $X = x$)". Following the context of "**structured-based counterfactuals**"², we specify $Y_{x'} = y'$ further as $Y_{M_{x'}} = y'$ with a *causal model* M containing other context-related variables $V \setminus \{X, Y\}$. We denote the origin model without intervention as M_x , and another "sub-model" $M_{x'}$ with intervention $X = x'$, as if we are often "imagine the another world" when thinking of *counterfactual causation*.

Here, the term "*counterfactuals*" implies the further evidences ($X = x, Y = y$) counterfactual to our imagination ($X = x', Y = y'$) in *PN* (Equ(*)). In cases without such evidences, we can connect it with the expression of *intervention causation* using the well-known *do calculus*³ over M_x :

$$P(y' \mid do(x')) = P(Y_{M_x} = y') = P(Y_{M_{do(X=x')}} = y') = P(Y_{x'} = y'). \quad (3)$$

Probability of intervention is fundamentally a general notion. This is because *intervention* expression is the average causality, actually marginalizing any background information (U) — the exogenous (unknown) factors over the cognitive causal model M . Based on this "background information" notion, *counterfactual* expression used to represent *PN* will further consider the evident background (namely the factual expression) over "specific scenarios" (\bar{U} that $\bar{U} \subset U$):

$$PN = \sum_u P(Y_{M_{x'}}(u) = y') := \sum_{\{u \mid Y_{M_{x'}}(u) = y', Y_{M_x}(u) = y, u \in \bar{U}\}} P_{M_{x'}}(u). \quad (4)$$

The formalism highlights the attribution or the analysis of causes based on that specific scenario $\bar{U} = u$. The emphasis on scenario evidence, compared to Equ(*), yields the "generic"⁴ "probability of assertion" on behalf of *PN*:

$$\text{Assertion : } P(Y_M(u) \mid e) = p. \quad (7)$$

Assigning $PA = p$ based on observational scenario evidence e , we characterize the intensity of our belief to the assertion.

The pivotal point is: **which causation relationship that an assertion is categorized** (e.g. *necessity causation*, *sufficiency causation*, *actual causation*), will essentially depends on how

OUTLINE:

- Structured-based Counterfactuals
- Causal Assertion

NOTES¹:

We denote $X = x$ and $X = x'$ as "with and without human influence", and $Y = y$ and $Y = y'$ as "whether global warming occurs or not".

REFERENCE²:

In *Causality*, Chapter 7.
"Structured-based" partially refers to a causal graph G_V with vertexes V associating with a structure M . We can see structures acting as the "bridge" to convey "background information" between a current world and a "counterfactual world" (will be explicitly expressed in Equ(4)).

REFERENCE³:

Notice that the notion of *do calculus* has always been exceptionally highlighted in The Book of WHY. Readers who have gotten a sense about $do(\cdot)$ *calculus* might know that the direct distinction between pure observation and causation lies in the action of "doing intervention":

$$P(y' \mid x') \neq P(y' \mid do(x')). \quad (2)$$

REFERENCE⁴:

Readers can compare the expression in Equ(7) with the standard form of *PN* and *PS* in text books:
Probability of Necessity:

$$P(Y_{x'} = y' \mid x, y). \quad (5)$$

Probability of Sufficiency:

$$P(Y_x = y \mid x', y'). \quad (6)$$

much evident information e we know about the scenario $U = u$. For instance, *necessary causation* entails the assertion $Y_M = Y_{M_{X=x'}} = y'$ in light of evidences $e = \{X = x, Y = y\}$ that are observed based on scenario u . This evident information in turn allow us to update our knowledge about the specific scenario $P(u) \rightarrow P(u|x, y)$. After explicitly specifying two of distinct causal models (e.g. M_x and the counterfactual one $M_{x'}$), applying the definition of conditional probability in Equ(1) will yield:

$$PN = \frac{\sum_u P(Y_{M_{x'}}(u) = y', X_{M_x}(u) = x, Y_{M_x}(u) = y)}{P(X = x, Y = y)}. \quad (8)$$

Where the joint statement of M_x and $M_{x'}$ becomes an ordinary event in a standard probability space—the space governed by associating background information $U = u^5$. By simplifying the constraint that $u \in \{u \mid Y_{M_{x'}}(u) = y', X_{M_x}(u) = x, Y_{M_x}(u) = y\}$, we ultimately obtain Equ(8) in form of $P(u)$:

$$PN = \frac{\sum_u P_{M_{x'}}(u) P_{M_x}(u)}{P_{M_x}(x, y)} = \sum_u P_{M_{x'}}(u) P_{M_x}(u|x, y). \quad (9)$$

The meaning by Equ(9) compared to Equ(1) is significant: the calculation of PN might be viewed as **abducting** the causal model from "central" information (X, Y) to "background" information (U) , then summing up weight of evidence (in form of $P_{M_x}(u|x, y)$) in the "current world" over the *necessity causation* assertions (in form of $P_{M_{x'}}(u)$) in the "counterfactual world"⁶.

Now notice that the moral is, the more the scenario evidence that we can hold on, the closer the *actual causation* that we are approaching. In order words, **we are more likely to make a powerful causal statement if we can include more evident information to specify the "scenario background"**⁷.

By breaking down the definition of *the probability of necessity causation*, let us briefly take a glimpse on *probability of actual causation*. We encompass⁸ necessity and sufficiency causal relations into a single assertion in Equ(10), along with a specific model-context named the **natural beam** M^u (a variant of causal model M , I will discuss it in the next Section). Roughly speaking, the *natural beam* M^u exactly characterizes "the frame of reference" of (environmental) background information u .

$$\text{Assertion : } P(Y_{M_{x'}}^u = y', Y_{M_x}^u = y \mid e) = p. \quad (10)$$

Notice that I simply omit the symbol u from $Y_{M_{x'}}^u(u)$ and $Y_{M_x}^u(u)$, because M^u is exactly tailored with respect to u .

Next in the second half, I would like to show you how to calculate PN (probability of necessity) following the "clue"⁹ at the beginning of the introductory part.

In form of *FAR* (fraction of attributable risk), it suggests that if human influence makes no sense to global warming, namely the two probabilities remain the same that $P(y \mid X = x) : P(y \mid X = x') = 1 : 1$, then it will lead to the *FAR* metric equalling to 0:

$$FAR = 1 - \frac{P(y \mid x')}{P(y \mid x)}. \quad (11)$$

NOTES⁵:

It might be slightly uncomfortable when initially encounter with the counterfactual symbol in joint probability formula (So am I). This largely due to the contrast that $Y_{x'} = y'$ and $Y_x = y$ are unlikely to be jointly true. However, we can relate to this by seeing the sharing background information $U = u$ as an "ordinary event". In that sense, the joint statement of $Y_{x'} = y'$ and $Y_x = y$ just characterizes how possibly the "information flow" stemming from $U = u$ can be conveyed (by causal diagrams) and finally result in "different branches" $Y_{x'} = y'$ and $Y_x = y$.

REFERENCE⁶:

The content on the left is following the standard three-steps of counterfactual inference (abduction-action-prediction) proposed in *Causality*.

NOTES⁷:

So far, I have highlighted the term "scenario background" over a causal model several times (the definition is explicitly expressed in Equ(4)). Personally, I think this is the fundamental philosophy of the causal model (as well the causal diagram learning).

NOTES⁸:

Recall to the sentences that we end the introductory part with:
"At the heart of it, *actual causation* encompass together the strength of *necessity causation* and the *sufficiency causation*, reaching a more comprehensive analysis of causes."

NOTES⁹:

Recall to the sentences in the introductory part:
"Mathematically speaking, saying 'the very likely the FAR of human influence' is essentially equal to saying 'the high PN of human influence'."

To see why *FAR* is essentially equal to *PN*, we need to turn *PN* from the *counterfactual* causation expression, into the ready-to-use statistic formula. Two of the important assumptions for the calculation are required: the ***exogeneity*** assumption and ***monotonicity*** assumption.

Correction is not causation. It is the unmeasured confounder that promotes us to establish the assumption to compensate for the estimation "gaps". X is exogenous relative to Y if the explicit generation function $Y_X := f_Y(X)$ possesses the independence:

$$M : X \perp\!\!\!\perp f_Y \Rightarrow X \perp\!\!\!\perp \{Y_x, Y_{x'}\}. \quad (12)$$

This makes sense since the functional mechanism f_Y where greenhouse gases will damage the climate is an established fact, no matter if human produces greenhouse gases or not. Associating with another mild assumption named the *consistency* (feeding X in $P(f_Y(X = x'))$) is consistent with conditioning X in $P(Y | X = x')$, we obtain the crucial conclusion by exogeneity:

$$P(Y_{x'} = y') = P(Y_{x'} = y' | X = x') = P(Y = y' | X = x'). \quad (13)$$

Meanwhile, the function f_Y mapping with four of the ***potential outcomes*** (Fig(9)) is monotonous if the "mapping trajectory" never changes down to $Y_x = y'$ given the X changing from x' to x — It is impossible that massive greenhouse gases produced by human ($X = x$) will never damage our climate ($Y = y'$):

$$(Y_{x'} = y) \wedge (Y_x = y') \Rightarrow \text{False}. \quad (14)$$

We now aim at unraveling the procedure of turning *PN* to the purely statistic expression. To begin with, we leverage the *monotonicity* assumption to rewrite *PN*, allowing us to equate $P(Y = y)$ with $P(Y_x = y)$ (*monotonicity* ensures there never exists the causation $P(Y_x = y')$ after observing $X = x$)¹⁰:

$$PN = \frac{P(Y_{x'} = y', Y_x = y, X = x)}{P(X = x, Y = y)}. \quad (16)$$

We further attain Equ(16) thanks to the independence in Equ(12) entailed by the *exogeneity* assumption:

$$PN = \frac{P(Y_{x'} = y', Y_x = y, X = x)}{P(X = x, Y = y)} \cdot P(X = x). \quad (17)$$

Then we obtain the following (simplified) equation where our primary attention boils down to the calculation of jointed distribution $P(y'_{x'}, y_x)$ in the numerator.

$$PN = \frac{P(y'_{x'}, y_x)}{P(y | x)}. \quad (20)$$

since the *monotonicity* assumption guarantees the only conversion from $y'_{x'}$ and $y_{x'}$ to y_x (given the X changing from x' to x). Therefore, we can obtain $P(y'_{x'}, y_x)$ directly by subtracting $P(y_{x'})$ from $P(y_x)$ ¹¹.

$$P(y'_{x'}, y_x) = P(y_x) - P(y_{x'}). \quad (21)$$

Via using the *exogeneity* again based on Equ(20) and Equ(13), we finally prove that *PN* is mathematically equal to *FAR*:

$$PN = \frac{P(y | x) - P(y | x')}{P(y | x)} = 1 - \frac{P(y | x')}{P(y | x)} = FAR. \quad (22)$$

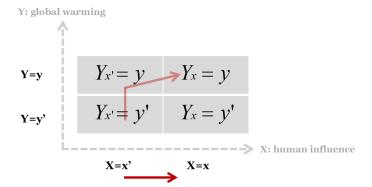


Figure 9: Illustration of the monotonicity assumption (over four of the potential outcomes).

NOTES¹⁰:

Notice that we can obtain Equ(15) by first applying the definition of conditional probability:

$$PN = \frac{P(Y_{x'} = y', X = x, Y = y)}{P(X = x, Y = y)}. \quad (15)$$

NOTES¹¹:

To be more formally, the capability of the monotonicity assumption herein is to tighten a lower bound. Seeing the jointed distribution as differently jointed solutions $Y_{x'} = y'$ or $Y_x = y$ for $Y_X(u)$ governed by a standard probability space over U , then in generally we have a sharp lower bound relating to "random events" $y'_{x'}$ and y_x that $P(y_x, y'_{x'})$:

$$\begin{aligned} &\geq \max[0, P(y_x) - (1 - P(y'_{x'}))] \\ &\geq \max[0, P(y_x) - P(y_{x'})]. \end{aligned} \quad (18)$$

Notice that

$$P(A, B) \geq \max[0, P(A) + P(B) - 1]. \quad (19)$$

3.1.2 Foundation: Use Diagrams to Unravel Causation Scenarios

Review the specific scenario we described in Fig(6), where the emission migration issue caused by mega fire could be detected and tackled by human technologies, making it slightly confusing to distinguish the actual responsibility for climate change. In the following, I will introduce how to leverage causal diagrams to define, transform, and solve the problem.

Graphical Causation: Analysis of Causes Based on Scenarios

In terminology, we depict the phenomenon in that singular circumstance as "the generic causal link (mega fire → climate change) has been "**preempted**" by other causes (namely human influence). Thus, mega fire is a legitimate but not actual cause, since the impact of a preemption is just about to invalidate its continuous causal effects to climate change, making the actual cause determination slightly elusive.

Thus, since the causal mechanism of nature force (e.g. causing issues of emission migration after the occurrence of mega fire) is possible to be preempted, it is insufficient to claim that the nature force is the actual cause.

Graphical Causation: Causal Beam as a Specific Causal Diagram

Interpreting by causal diagrams, the invalidation essentially amounts to expunge the arrows stemming from the cause that is preempted, which brings us to a revelation: In some of the singular scenarios, we can and should "slim down" a causal diagram by deleting some trivial arrows.

Notice that this makes sense because some causal relations genuinely cannot exist when preemption accidentally occurs in that scenarios. Let us take a bit addition to our vocabulary to describe this behavior: "slimming a causal graph down to a **causal beam**". The *causal beam*, less rigorously, suggests that a causal graph is projected to a subgraph relative to a singular scenario, serving as a "frame of reference" of the surroundings¹. Therefore, the function of constructing the causal beam lies in helping us to clearly figure out the singular cause or, may be exactly the actual cause.

For instance, consider another modification of the climate change example in which we have not idea about the weather condition background shown in Fig(11). Let us assume that the weather condition will influence the factors in the diagram (political resistance and emission migration) — Extremely bad weather (e.g. haze, blizzard, etc.) might hinder the political activities or natural process. It turns out that², "slimming down" the causal diagram in Fig (11) yields the causal beam in Fig (12). The remaining connections thereof (illustrated as red) partially represent the "strong" relation that would never disappear even in some singular circumstances. Now this "strong intensity" of causal relations, from the perspective of the causal diagram (causal beam), unravels the evidence in which human influence is the actual cause of climate change.



Figure 10: The causal diagram for a specific scenario: A causal preemption occurs.

REFERENCE¹:

Notice that we have introduced the sustainable-based causal relations in the previous section.

Accordingly, causal beam models the operation behind it: operations to, quote, "freeze trivial surroundings" (Causality, Chapter-10).

If we observe that the sustenance persisting, it hints the (strong) intensity of the actual cause, even in some singular cases.

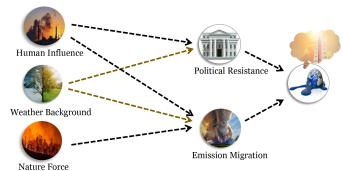


Figure 11: The causal diagram for a general scenario.

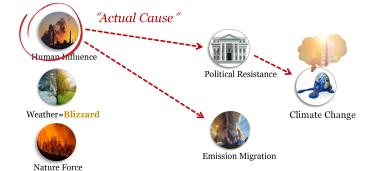


Figure 12: The causal beam for a specific scenario.

TIPS²:

In the formalization part, I will show how to slim down the causal diagram by tuning the responsive function entailed by the causal diagram.

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of Judea Pearl's points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

The heart of this section is about the *causal beam*: slimming down a causal graph based on a singular scenario. From the SCMs¹ perspective, this modification equals to shutting down or downsizing the responsive functions f involved in the original causal model M , if given a singular background $U = u$:

$$f(PA, U) \rightarrow f(PA^*, U = u) \rightarrow f^u(PA^*). \quad (23)$$

Here $PA^* \subset PA$ indicates weaker "listening" relations in u , leading to the construction of f^u showing the newly functional response in a modified² causal model M^u — the *causal beam*.

Alternatively, suppose $PA = S \cup S'$ and $PA^* = PA \setminus S' = S$. While freezing some variables in the state of $S = S(u)$, f^u is constructed by eliminating S' from PA that remain f^u trivial:

$$\exists U = u, \forall S' = s', \sigma(f^u(S, S' = s')) = 0. \quad (24)$$

Therefore, **actual cause** (AC), namely the sustainable-based causal relations in which analysis of **necessary cause** (NC) and **sufficient cause** (SC) are converged³, can be clearly defined over the *causal beam* M^u :

$$AC(X \rightarrow Y) : [Y_{M_x^u} = y] \vee [Y_{M_{x'}^u} = y']. \quad (25)$$

Assertions in Equ(*) being true would indicate X is the *actual cause* of Y . Consider further with uncertainty $P(U = u)$, the **probability of actual causation** (PA) in light of evidence e is the average summing up against the evidence weight $P(u | e)$:

$$PA = \sum_u P(Y_{M_x^u} = y, Y_{M_{x'}^u} = y') P(u | e). \quad (26)$$

To see how to apply actual causation analysis, let us back to the example context in the introductory part. Fig (13) shows the graphical structure among variables whose responsive functions (represented as AND-OR logic) are exhibited in Equ(27):

$$M = \begin{cases} z_1 := x_1 \wedge \neg U \\ z_2 := x_2 \wedge (U \vee x'_1) \\ y := z_1 \vee z_2 \end{cases}. \quad (27)$$

For the unknown weather ($U = u$ or $U = u'$), we assume being in blizzard ($U = u'$) with evidences ($e = \{X_1 = x_1, X_2 = x_2\}$). we resort to the causal beam $M^{u'}$ yielded in the following:

$$M^{u'} = \begin{cases} z_1 := x_1, z_2 := x'_1 \\ y := z_1 \end{cases}. \quad (28)$$

Since less mega fires occurred in blizzard (u'), only link $x_1 \rightarrow z_1 \rightarrow y$ is distinctive ($M_X^{u'} : Y_{x_1} = y, Y_{x'_1} = y'$). We have to grant that humans are actually to be blamed for the climate change.

OUTLINE:

- Natural Beam and Causal Beam
- Probability of Actual Causation

NOTES¹:

Causal graphs are vague representation of the structure causal models (SCMs), in the sense that an arrow over a graph reflects a relationship of "listening" over a SCM. That word "listening" essentially refers to "the responsive functions f ". Variables that render f nontrivial are often called the parents (PA).

NOTES²:

For simplicity, here I omit some details in building up causal beam. In fact, how to restrict a modified model depends on the varying degree of "freezing". We call it the "**natural beam**" if we freeze all the variables.

NOTES³:

Remember that when discussing the sufficiency and necessity causal relations in the original causal model M , we implicitly denote expressions $Y_{M_x} = y$ and $Y_{M_{x'}} = y'$ as $Y_x = y$ (SC) and $Y_{x'} = y'$ (NC).

NOTES⁴:

Notice that the formula in Equ(26) is the expansion relative to the one shown in Equ(7) and (10), where we started discussing the general idea of causal effect assertion.

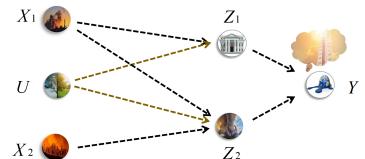


Figure 13: The causal diagram (variable formalization) in the example.

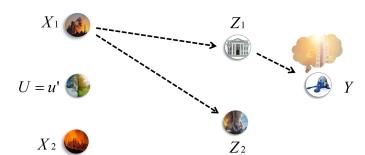


Figure 14: The causal beam (variable formalization) in the example.

Causal Books References

Reminder: "Partial ideas from celebrated causal books (*Causation, Prediction, and Search*; *Causality*; *Elements of Causal Inference*; *The Book of WHY*) that are beneficial to interpret causation are woven into the paper."

The Book of WHY

- *Chapter-8. Counterfactuals — Mining Worlds That Could Have Been*
 - Adaptation of the global heatwave example.
 - Drawing the notion of necessity and sufficiency causation into the notion of sustenance-based causation, associating with the actual cause analysis in the book *Causality*.

Causality: Models, Reasoning, and Inference

- *Chapter-7. The Logic of Structure-Based Counterfactuals*
 - 7.1 Structural Model Semantics
 - * Probabilistic causal model based on latent background variables.
 - 7.5 Structural Versus Probabilistic Causality
 - * Differences between intervention and counterfactuals.
 - * The (partial) reason why the monotonicity assumption is further required.
- *Chapter-9. Probability of Causation: Interpretation and Identification*
 - 9.2 Necessary and Sufficient Causes: Conditions of Identification
 - * Definition: probability of necessity and sufficiency, PNS.
 - * Definition: the exogeneity assumption and monotonicity assumption.
 - * Relations between PN and PNS.
 - * Relations between PN (PS) and excess risk ratio (ERR).
 - * The nonzero bound widths for PNS, implying that probabilities of causation cannot be defined uniquely in non-Laplacian models (The reason why the exogeneity and monotonicity assumptions is further required).
 - 9.3 Examples and Applications
 - * Counterfactual-based causal assertion for analysis of causes.
- *Chapter-10. The Actual Cause*
 - 10.1 Introduction: the Insufficiency of Necessary Causation
 - * Adaptation of the desert traveler example.
 - * Insufficiency in necessity causation.
 - 10.2 Production, Dependence, and Sustenance
 - * Relations between causal assertions and philosophical thinking.
 - 10.3 Causal Beams and Sustenance-Based Causation
 - * Definition: natural beam and causal beam.
 - * Definition: sustenance-based causation and actual cause determination.

3.2 Can Vaccinal Treatments Save Lives of Loved Ones?

In the previous section, I introduced the analysis of causes. When we find out a cause, we might also wonder the procedure of its inherent causal mechanisms¹ — ***causal effect*** (Fig(15)). For example, not being vaccinated probably causes infections by coronavirus. However, what is the mechanism by which the vaccine prevents infection?

A significant causal mechanism, or causal effect, is normally understood by ***intervention***. In this section, I will discuss the connotation of fundamental types of causal effects, which is woven into my grandfather's story struggling with COVID-19.

3.2.1 Application: Intervention

It appears inevitably to be a sense of immorality or unrealistic when it comes to talking causality seriously: It is undoubtedly immoral to impose virus infections upon patients (namely the *intervention*); yet it is also impossible to observe the treatment effect again in reality since man has only one life.

The point is, however, the infeasibility in implementing the causal *intervention* does not necessarily mean the infeasibility in modeling causal effects. Going through the following, I hope we can get a sense in which theoretically modeling causal effect is as significant as, if not more, than practically conducting causal treatments in real world.

Conceptual Causation: Inquire about Causal Effect

Many of people experienced losing love ones during the pandemic over the last couple of years. COVID-19 has as well taken away the life of my dear grandfather.

Accompany with memorial regret lingering within our mind, there might still be questions left with unclear truth of causal-effect behind our losing love ones. In fact, moral of the causality power largely depends on how we are capable of imagination, bringing us insights into seemly tricky questions. Envisioning *counterfactual* scenarios, we are exceptionally familiar with such a paradigm to attach causality: what effect it would be if we can turn back the clock, and start over another treatment? In terms of my case, could my grandfather have survived from coronavirus if he received getting vaccines in early time?

If I were placed in the other parallel worlds and encouraged my grandfather to receive vaccination in time, how much would the vaccine finally save my grandfather's life? Essentially, I herein wish to aware the (*total*) ***causal effect*** of vaccination — The (*total*) ***causal effect*** should immediately be divided into two layers: ***directed causal effect*** and ***indirected causal effect***.

For example, if I attempt to observe the (vaccine's) directed effect by controlling every treatments (into lose their efficacy) except for getting a vaccine, then I am seeking for the ***controlled directed effect (CDE)***. However, a plethora of control over other factor(s) could sometimes lead to pitfalls. In Fig(18), what if the



Figure 15: The second version of the question "WHY": Why? What is the mechanism by which vaccine prevent infection?

REFERENCE¹:

Namely, we wish to better understand the connection between a known cause and a known effect. In *The Book of WHY*, Chapter 9, understanding the causal mechanism is described as answering the "second" version of the "WHY" question (the first version is the analysis of causes).

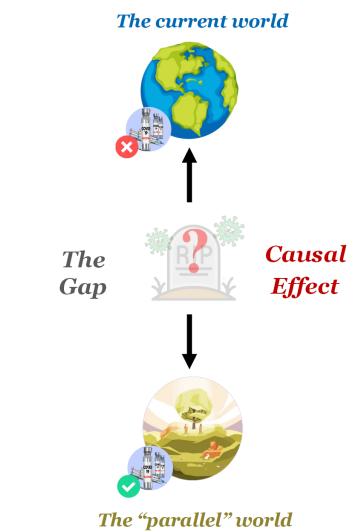


Figure 16: Moral of the causality power — the counterfactual thinking we discussed in the previous section. Generally speaking, the gap between the current world and the parallel (counterfactual) world is namely the causal effect, where the difference in the two different worlds therein lies in whether my grandfather receiving the vaccine treatment or not.

physiological mechanism is that a dose of vaccine is functioning by simultaneously proliferating the mediated factor(s) M (such as enzymes) to function as well? It certainly sounds absurd for my control to disobey the mechanism by "killing" enzymes just to malfunction them. Actually, I am more likely wanting a docilely natural control: curbing the physiological mechanism by continuously maintaining the enzyme at the level that is before vaccination. If after merely an injection of vaccine (without breaking the balance of the other mediated mechanisms), it can significantly save my grandfather's life — we will call it relatively the ***nature directed effect (NDE)*** (illustrated in Fig(19-(b))).

Unfortunately, the concrete medical mechanism appear to be slightly more complicated. Suppose that the abundant enzymes are primarily functioning by coordinately producing a spike level of antibodies to prevent the disease. Namely, merely an injection of vaccine with the normal level of the enzyme amount is actually not able to generate enough antibodies, and furious symptoms would still be rampant without enough antibodies. That means, the ***(natural) directed effect*** can be sometimes insignificant, whereas this brings us to focus more on the the crucial indirected effect.

Suppose that now I only consider the indirected effect (of the enzyme). Once the indirected mediated mechanism makes sense, I can imagine resorting to an elixir that can stimulate my grandfather's body to produce the enzymes in an amount as much as a vaccination level. In that sense, a strong outcome of the "***nature indirected effect (NIE)***" will imply my grandfather's miraculous recovery from COVID-19, even without injecting a vaccine beforehand (illustrated in Fig(20)).

Feeling frustrated again, it might turn out that the ***NIE*** stays insignificant similarly, if antibodies produced by the enzyme itself are still not enough to protect my grandfather from diseases.

Nonetheless, should we just deem that getting vaccines is never work at all? Not quiet. The tricky point is, again, the amount of antibodies that are promoted to release only by the vaccine or the enzyme will fail to reach the intensity threshold to help my grandfather dodge COVID-19 (though vaccine or enzyme does play the physiological role in their chemical reactions). It makes the question boiling down to a common setting in reality: ***non-linearity*** — causal effects of vaccines and causal effects of enzymes are fundamentally indispensable.

That is, "a unit of" intervening treatment (e.g. from not getting vaccines to getting vaccines) does not linearly or proportionally associated with "a unit of" recovery (e.g. from death to survival), unless a certain medically intensity threshold has been meet.

So what happens? Due to the ***non-linearity*** (e.g. common restrictions such as the medication threshold of the effective antibodies level), we recognize the seemly intuitive "principle of additivity" — "*the total causal effect of the vaccine equals to the directed effect of the vaccine itself, and 'plus' the indirected effect of the vaccine via triggering beneficial enzymes*" — just cannot

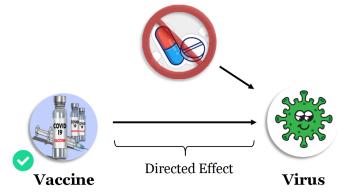


Figure 17: If the COVID-treatments involve taking the certain pill, then its efficacy should be controlled in order to measure the pure effect of the vaccine treatment.

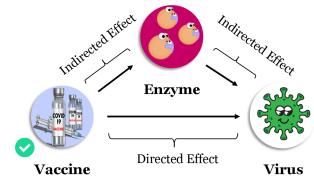


Figure 18: The assumptive physiological mechanism of vaccination. The vaccine treatment might result in an assistive "enzyme treatment".

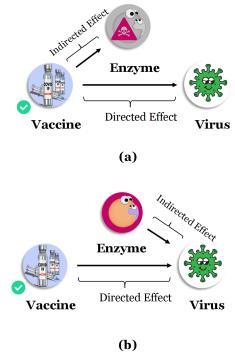


Figure 19: (a) An unnatural control by malfunctioning the enzyme; (b) naturally maintain the enzyme's amount that in the level before receive the vaccine (compare the "number" of enzyme with Fig(18)).

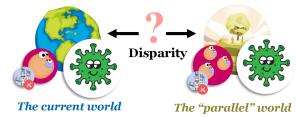


Figure 20: The intuition of the NIE. Notice that neither in the current world nor the parallel world my grandfather were supposed to receive the directed vaccination. In order to evaluate indirected enzyme efficacy, I presume there were an elixir to simulate the amount of enzymes to a vaccine-level in the parallel world. Finally, we can calculate the NIE by calculate the disparity in the two different worlds — the disparity of the probabilities where the COVID symptoms are still rampant.

work. How could we learn a tangible causal-effect given either the directed effect or the indirected effect is both insignificant? However, with the way of thinking flipping around, namely the *counterfactual* thinking, surprisingly we can redescribe it from another perspective — "**the principle of subtraction**":

"the total causal effect of the vaccine as well equals to the directed effect of the vaccine itself, and 'minus' the potentially indirected losses of the vaccine — the losses that without the potential enzymes' assistance triggered by vaccines."

Merits of this thinking perspective bring us to conjure up the other picture (Fig(21)), where the losses that are without the indirected effect of enzyme treatment can be characterized by an "inverse" *nature indirected effect (NIE)* (Fig(21-b)). The meaning entailed by this "inverse" version of the causal effect is that, instead of transferring from the previous treatment (time= t_0) to the hindsight treatment (time= T), namely from without enough amount of enzymes to the abundant enzymes, the perspective of *counterfactual* thinking instructs us to (hypothetically and straightforwardly) start with the hindsight treatment (time= T) in which the abundance of enzymes are exactly proliferated by vaccination. Then, we turn back the clock, and start over another "treatment" — namely without the vaccine treatment (and without the enough amount of enzymes).

Notice that this *counterfactual* thinking (suppose receiving the vaccination beforehand) brings the effect asymmetry shown in Fig(21-b): Only when the vaccine efficacy is involved, could we observe the significant *non-linearity* change of the levels of enzymes (e.g. under the threshold and above the threshold). This significant change amounts to the significant disparity between two of the different setting (in two of the different worlds), further contributing to an significant result of the "inverse" *NIE* — The inverse *NIE* equals to -1 whereas the original *NIE* equals to 0. Incorporating this result into the "principle of subtraction", we will discover the *total causal effect* actually equals to 1 — receiving the vaccine could have worked!

Finally, to close the end of the story where my grandfather had struggled with COVID-19, I want to discuss more a bit about our thinking way of causality in treatment effects. Not feeling slightly weird the next time when hearing of "*total causal effect = directed effect – indirected effect*", instead of "*total causal effect = directed effect + indirected effect*". We know that the "subtraction"² essentially indicates to evaluate treatments by the *counterfactual* thinking. The strength of it gives rise to a general "subtraction principle" — an "envisioning-parallel-world" way of thinking to understand causality in *intervention*.

May my dear grandfather rest in peace in his heaven world.

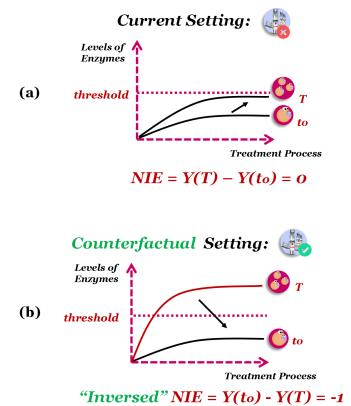


Figure 21: The inverse version of the nature indirected effect (NIE) that fits to the non-linearity context, which is a special type of the intervention in light of counterfactuals thinking.

REFERENCE²:

A piece of anecdote from Judea Pearl, quote,
"I was extremely thrilled to see this subtraction principle emerging from the analysis of Total Causal Effect (TCE), despite the nonlinearity of the equations." .

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of *Judea Pearl's* points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

In fields of biological and medical studies, researches can refer to **average treatment effects (ATE)**, the (positive) efficacy ($Y = y$) of the associated treatment(X), by measuring the average disparity between participants in different assigned groups (e.g. $X = x$ or $X = x'$):

$$ATE = P(y \mid X = x) - P(y \mid X = x'). \quad (29)$$

There is a primary hurdle, however, preventing our estimation from the "true causal efficacy" in practical implementations. Reasons behind often involve the confounding that normally comes from "unknown" factors during treatment experiments:

$$|P(Y = y \mid do(x)) - P(Y = y \mid x)| > 0. \quad (30)$$

The inequality $P(Y = y \mid do(x)) \neq P(Y = y \mid x)$ described by *intervention (do calculus)*¹ suggests that, in practice, additional "known" factors z might be helpful to shrink or even eliminate the gap via mathematical transformation Φ :

$$\Phi_z : \underbrace{P(Y = y \mid do(x))}_{\text{theoretical}} \mapsto \underbrace{(P(Y = y \mid x, Z = z))}_{\text{practical}}. \quad (31)$$

Equ(31) shows the whole picture. But since we mainly focus on modeling theoretical causal effects, I will use $P(Y = y \mid do(x))$ as the building-block symbol for the following demonstration².

Assuming binary variables, we start with the well-known **ACE (average causal effect)**, the general causal-effect estimation that measures the average change corresponding with the different intervening treatment effects $P(y \mid do(X))$:

$$ACE := P(y \mid do(X = x)) - P(y \mid do(X = x')). \quad (32)$$

In the pandemic example, rampant symptoms by disastrous COVID amounts to putting a low probability of escaping from infection ($Y = y$) without the previous (time= t_0)³ vaccine protection ($X_{t_0} = x'$). In a hindsight situation (time= T), we presume participants were vaccinated ($X_T = x$). We then ask could they stay survival (assigning a high probability $Y = y$) by assessing the causal effects $X \rightarrow Y$ after the intervening $X_T = x$ versus the previously $X_{t_0} = x'$:

$$ACE_{X \rightarrow Y}^{t_0 \rightarrow T} := P(y \mid do(X_T)) - P(y \mid do(X_{t_0})). \quad (33)$$

In the basic formula of ACE, we then immediately introduce the mediating mechanism W (e.g. enzymes spured by vaccines)

OUTLINE:

- Average Treatment Effects
- Controlled Directed Effect
- Natural Directed Effect
- Natural Indirected Effect

NOTES¹:

In causation, recall that the intervention expression is often formally denoted as $do(x)$ (e.g. doing a treatment instead of seeing a treatment).

NOTES²:

But please keep in mind the endeavor where we always need to maneuver around Φ_z from the theoretical "causal estimation" into the practical "statistical estimation".

NOTES³:

To put it into our context, let us introduce the concept about "time" to enhance our understanding as to the causal effect.

and fully control it, obtaining a relatively "sensitive" estimation named *CDE* (*control directed effect*):

$$CDE_{X \rightarrow Y}^{t_0 : T} := P(y | do(X_T), do(\mathbf{w})) - P(y | do(X_{t_0}), do(\mathbf{w})). \quad (34)$$

Slightly different from *ACE*, measuring the sensitivity implies that the *CDE* of $X \rightarrow Y$ is obtained when all other causal path (e.g. $X \rightarrow W \rightarrow Y$) are cut off by $do(\mathbf{w})$ ⁴. However, it might sometimes be a absurd compulsion because *CDE* hints us to select which mechanism (e.g. $W = ?$) to be "controlled". We may not be willing to make such a choice since a wrong control over mechanism can probably result in a more problematic issue.

So to speak, there is reasonable to obtain a more natural estimation: measuring the anticipated change of Y under the "pure shifting"⁵ from $do(X_{t_0})$ into $do(X_T)$, which will retain the "previous preference"⁶ of the mediate mechanism W related to X_{t_0} . With merely a directed change of the vaccine treatment X , we further define the *nature directed effect* (*NDE*) as:

$$NDE_{X \rightarrow Y}^{t_0 : T} := P(y_{w_{t_0}} | do(X_T)) - P(y_{w_{t_0}} | do(X_{t_0})). \quad (35)$$

The probability $P(y_{w_{t_0}})$ that we interested in reflects the likelihood to dodge COVID ($Y = y$) in a hypothetical situation where all other variables W remains the previous value $W = W(X_{t_0})$ (except for X , because we were exactly curious about what things would be after a switch up to X_T against X_{t_0}).

Notice that we choose to use the notion of *structure-based counterfactuals* (e.g. $P(Y_{M_{do(w)}} = y' | X_{t_0})$, recall M denoted as the *causal model*) since *intervention* (e.g. $P(Y = y' | do(\mathbf{w}))$) is a general notion for controlling over (every) scenarios $do(\mathbf{w})$; whereas we prefer to the specific scenario ($do(\mathbf{w}) = do(w(x'))$)⁷. We emphasized this in the formalization part in Section 3.1.1.

Once again, leveraging the expressive *counterfactual* notion, it seems that we are able to analogically define the *NIE* (*nature indirected effect*) exhibited in Equ(37) compared to Equ(35):

$$NIE_{X \rightarrow Y}^{t_0 : T} := P(y_{x_{t_0}} | do(W_T)) - P(y_{x_{t_0}} | do(W_{t_0})). \quad (37)$$

Naturally we want to measuring the anticipated change of Y under the "pure shifting" from $do(W_{t_0})$ into $do(W_T)$ by retaining the "previous state" of treatment $X = X_{t_0}$. However, notice the fact that W is actually depended on X that $do(W_T) = W(do(X = X_T))$. This would contrast with the *intervention* involved in the *counterfactual* expression that $y_{x_{t_0}} = y_{do(X=X_{t_0})}$. To fix this this mild fault, we can rewrite Equ(37) as⁸:

$$NIE_{X \rightarrow Y}^{t_0 : T} := P(y_{w_T} | do(X_{t_0})) - P(y_{w_{t_0}} | do(X_{t_0})). \quad (38)$$

As mentioning at the end of the introductory part, we can obtain the *total causal effect* (*TCE*) via a more general "principle of subtraction" to combining together *NDE* and *NIE*:

$$TCE_{X \rightarrow Y}^{t_0 : T} := \left(NDE_{X \rightarrow Y}^{t_0 : T} \right) - \left(NIE_{X \rightarrow Y}^{T : t_0} \right). \quad (39)$$

I need to highlight the term $NIE_{X \rightarrow Y}^{t_0 : T}$ (instead of $NIE_{X \rightarrow Y}^{T : t_0}$) that exactly specifies "the losses that are without account for the indirected effect". I would like to say that silently swapping the "time point" essentially implies our *counterfactual* thinking.

NOTES⁴:

Notice that again, the symbol $do(\mathbf{w})$ indicates to confine W into constants by physical intervention $do(W = w)$, not by conditioning $see(W = w)$.

NOTES⁵:

Recall the content in the introductory part that "It certainly sounds absurd for my control to disobey the mechanism by 'killing' enzymes just to malfunction them."

"Actually, I am more likely wanting a docilely natural control."

NOTES⁶:

"Retaining the previous preference" is explicitly represented as $w_{t_0} = W(X_{t_0}) = w(x')$, with the "previous" value $X_{t_0} = x'$ before shifting.

NOTES⁷:

To supplement to this, *NDE* can be intuitively reduced to the weighted average of *CDE* under certain "no confounding" assumptions (discussed in literature):

$$\sum_w (CDE_{X \rightarrow Y}^{t_0 : T}(w)) \cdot P(w | do(X_{t_0})) \quad (36)$$

REFERENCE⁸:

Equ(38) is the standard mathematical form of the *NIE* highlighted in The Book of WHY; and we have seen its intuition in the introductory part (presume an elixir that can stimulate the enzymes from the normal level w_{t_0} to the vaccination level w_T). Personally, Equ(38) might be the most representative one to mathematically interpret how to deal with causal intervention about intermediate mechanism.

3.2.2 Foundation: Use Diagrams to Deduce Causal Inquiry

Review the versatile causal-effect inquiry we introduced in the previous section, in the following, let us briefly redescribe them but from the perspective of causal diagrams — how to execute different "operations" on the causal diagram in light of different requirements for inquiry? Notice that this is significant, as the permission for answering the causal inquiry through graph operations will imply the feasibility for answering the causal inquiry through statistical estimation in the real world.

Graphical Causation: Inquire about Causal Effect

Formalism of inquiries boils down to asking *directed causal effects* and *indirected causal effectss*. For each category, there might be a "controlled" version and a "nature" version. We can execute different types of "surgery" for the causal diagram to gain some intuition of these inquires. Literally, two "scalpels" relative to the graph's nodes are called (in Judea Pearl's words) "*to hold it on constant*" and "*to tweak it on compulsion*".

Take Figure (22) and (23), I herein take another unknown common factor into consideration for the illustration of "graphical surgery". "Tweaking the node" in a causal diagram refers to enforce an intervention, thus in the figures the arrows pointing to the node that is waiting to be tweaked are deleted from causal diagrams — the only arrow is our operation of "tweaking it". Meanwhile, "hold constant the node" is similar to "control", but implicitly involves the "*counterfactual*" information. This is because the most natural way to keep the (other) factors unchanged (while tweaking the main node) is just to render them to be *counterfactual* to the situation that after intervention.

Importantly, notice that the graphical surgery results in a ***causal subgraph***, which is the "generic format" as to how the causal diagram can answer causal inquiries. Difference "shapes" of the *causal subgraph* imply the versatile scenarios — the different intervention based on the different causal inquiries.

Graphical Causation: Inquire about Causal Effect

Graph operations, as a friendly language of causation, are extremely good at "simulating" experiments. The following pieces involve testing content. Strongly reliable causal formula would still be deduced (not matter how poor the data we have) to answer our inquiries without massive and costly experiments, only if we pass specific graph tests with flying colors.

We end this part with an example from *The Book of Why*, where Pearl left us with a vague meditation about that muse whisper "Try the do-calculus". You can also type a keyword to search online for the video that best fits your document.

Pivotal roles played by causal diagrams lie in driving as an "inference engine" to deduce the reliable causal formula. Conversely, the prerequisite of such deductive reliability exactly relies on satisfying a "graph test". Certainly, one could deduce

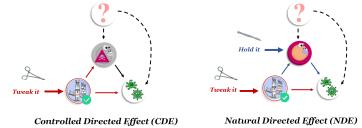


Figure 22: The "graphical surgery" relative to the directed causal effect. Compared to the previous situation where vaccination was not received, "tweak" the node (vaccine) on compulsion is meant to intervene a vaccination treatment and to observe the change brought by the flow (red) of causal information (namely arrows in the causal diagram). Therefore, it seems unnatural that the flow of causal information is stuck at the controlled node (enzyme). Instead, we tend to just hold constant the information (blue) of the node (enzyme) to be the level that is before tweaking the node (vaccine).

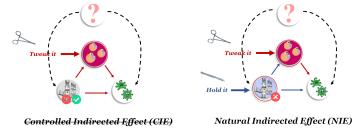


Figure 23: The "graphical surgery" relative to the indirected causal effect. Notice that we implicitly skip the notion like "controlled indirected effect". This is on purpose since such a notion does not exist. Roughly speaking, vaccine (directed effect) is exactly the main cause cause of enzymes (indirected effect). Thus the indirected effect will not exist if the directed effect were shut down (controlled) beforehand.

Margin Note: Recall that we have mentioned; emphasis on theoretically modeling; ...

Margin Note: Readers who might be familiar with the *front-door criterion*. Viewing it as a "sequential decisions", ...

the same formula by hand using the rudimentary *do operation* principles. However, image a machine, embedding with concise representation of a huge causal graph, can compute deductions automatically and answer our causal inquiries swiftly. It would be more appealing and exhilarating, would it?

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of Judea Pearl's points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

In fact, versatile estimations¹ of causal effect by enforcing an intervention $do(x)$ can end up converging upon an atomic deduction tool: **back-door adjustment**².

$$P(y \mid do(x)) = \sum_z P(y \mid x, z) \cdot P(z). \quad (40)$$

Generally, when many intervening objects $do(x_1), \dots, do(x_d)$ are sequentially in estimation (the **dynamic plan**), an "advanced version" named **sequential back-door adjustment**³ is available:

$$P(y \mid do(\mathbf{x}_{1:d})) = \sum_{z_1, \dots, z_d} P(y \mid \mathbf{x}_{1:d}, \mathbf{z}_{1:d}) \cdot \prod_{k=1}^d P(z_k \mid \mathbf{x}_{1:k-1}, \mathbf{z}_{1:k-1}). \quad (41)$$

Readers might come across a bit confusion.

I mean, these extra variables z cannot occur in a vacuum. But should we just roughly select z from $V \setminus \{x, y\}$? How could we select a sequence $\mathbf{z}_{1:d} = \langle z_1, \dots, z_d \rangle$ that matches⁴ a plan $do(\mathbf{x}_{1:d}) = \langle do(x_1), \dots, do(x_d) \rangle$?

The whole picture is, since we know it should connect an action $do(x)$ with a "surgery" of the graph \mathcal{G} , then analogously we can marry a plan ($do(\mathbf{x}_{1:d})$) with a series of "surgery" ($\pi_{\mathcal{G}}$):

$$\pi_{\mathcal{G}} = \langle \mathcal{G}_{\underline{X}_1 \overline{X}_{2:d}}, \dots, \mathcal{G}_{\underline{X}_{1:k} \overline{X}_{k+1:d}}, \dots, \mathcal{G}_{\underline{X}_{1:d}} \rangle. \quad (42)$$

Tune on each $\mathcal{G}^{(k)}$, the answer as to the above question is: seeking for z_k over that $\mathcal{G}^{(k)}$ just by applying the well-known **back-door criterion** — a set of graph testing rules. Roles of the graph is distinctive: Test whether $do(x_k)$ can be transition to purely statistical expression by controlling z_k (Equ(41)).

It is also similar to use the *do operation* rules (discussed in Section(*)). Back to our previous "three lines of do-calculus" example aiming at $P(y \mid do(\mathbf{x}_{1:2}))$. We start turning $do(x_1)$ to x_1 as whether receiving the first-dose is fully controlled by ourselves:

$$P(y \mid do(\mathbf{x}_{1:2})) = P(y \mid x_1, do(x_2)). \quad (44)$$

Then introducing m (enzyme) over the marginal probability yields:

$$P(y \mid x_1, do(x_2)) = \sum_m P(y \mid x_1, do(x_2), m) P(m \mid x_1, do(x_2)). \quad (45)$$

Ultimately we turn $do(x_2)$ to x_2 in the front part as the second-dose effect is now fully determined by enzyme m , and delete $do(x_2)$ in the back as the $do(x_2)$ does not function on m .

$$P(y \mid do(\mathbf{x}_{1:2})) = \sum_m P(y \mid x_1, x_2, m) P(m \mid x_1). \quad (47)$$

OUTLINE:

- Dynamic Plans
- Sequential Back-Door Criterion
- *Do* Operation

NOTES¹:

average causal effect (ACE), control directed effect (CDE), nature directed effect (NDE), nature indirected effect (NIE), total (causal) effect (TCE), to name a few.

NOTES²:

Recall the mathematical transformation Φ shown in Equ(31) that helps eliminate the gap between causation and correlation. Now the back-door adjustment is exactly a basic transformation in practical causation.

NOTES³:

Notice that here I just refer to it directly in Equ(41), with a glimpse to see how it mirrors the formula in Equ(40).

NOTES⁴:

In other words, how could we match the pair variables in each of the "phrase (k)" given a plan? (denote variables in each of the "phrase" as $do(x_k)$ and z_k)

Margin Note:

Readers can review that "surgery symbols" \mathcal{G} in ...

Margin Note:

Readers who are familiar with the *d-separation* criteria (we would briefly induce in Section(*))...

$$(y \perp\!\!\!\perp x_k \mid z_k)_{\mathcal{G}^{(k)}} \quad (43)$$

Margin Note:

If we further combine all the "phrase" $do(x_k)$ over a "series level", then this is essentially how the *sequential back-door criterion* works (Equ(*)).

Margin Note:

$$\sum_{\emptyset, m} P(y \mid x_1, x_2, \emptyset, m) P(m \mid x_1, \emptyset). \quad (46)$$

Causal Books References

Reminder: "Partial ideas from celebrated causal books (*Causation, Prediction, and Search*; *Causality*; *Elements of Causal Inference*; *The Book of WHY*) that are beneficial to interpret causation are woven into the paper."

The Book of WHY

- *Chapter-1. The Ladder of Causation*
 - Intuition as to the causation in the smallpox-vaccine example.
- *Chapter-5. The Smoke-Filled Debate: Clearing the Air*
 - Basic ideas on randomized controlled trial (RCT).
 - Hazards of unobserved confounding.
- *Chapter-7. Beyond Adjustment: The Conquest of Mount Intervention*
 - The problem of sequential decisions and the "three lines" do-calculus, associating with the dynamic plan and the sequential back-door criterion in the book *Causality*.
- *Chapter-9. Mediation: The Search for a Mechanism*
 - Why is not "total effect = directed effect + indirected effect"?
 - From the "controlled causal effect" into the "natural causal effect".
 - Judea Pearl : "When I managed to strip the natural-causal-effect formula from all of its counterfactual representation, it was the greatest thrills in my life."

Causality: Models, Reasoning, and Inference

- *Chapter-3. Causal Diagrams and the Identification of Causal Effects*
 - 3.2 Intervention in Markovian Models
 - * Understanding intervention from the Bayesian perspective.
 - 3.3 Controlling Confounding Bias
 - * Definition: the back-door adjustment.
 - 3.4 A Calculus of Intervention
 - * The symbol of intervention that used in the sub-graph test for the sequential back-door criterion.
- *Chapter-4. Actions, Plans, and Direct Effects*
 - 4.4 The Identification of Plans
 - * Definition: the sequential back-door criterion and the back-door criterion.
 - * Relations between the calculation of the natural causal effect and the sequential back-door adjustment.
 - * Relations between the sequential decision and dynamic plan.
 - * Graphical language of the do-calculus.
 - 4.5 Direct and Indirect Effects
 - * Formalization: "controlled causal effect" and "natural causal effect".
 - * Definition: the mediation formula.

4 Inferred Causation

4.1 Causal Diagram Learning = Data + Restriction

Conceptual ideas of *causal models*, such as **causal attribution** and **causal effect**, are partially discussed in the first half of this paper, whereas causal diagrams can be viewed as a "fine-grained" reflection of *causal models*. Merits of the diagrams lie in the "arrow" that implies the simplicity of **causal significance** and the visibility of **inducing bias**, deserving our discussions as to what is a "learnable causal diagram".

4.1.1 Why Does the Equation Harbor Causal Significance?

Learnable causal graphs refer to "recovering" causal structures from statistical data, which is why we literally get the name of this standard task: **Causal Discovery**. Moral of "recovering" hints a "losing data generation" we are seeking for — the mechanism embedded in the *structure-based causal model*¹.

Causal Mechanism and Structure Causal Models

"Data-driven" learning methods might sometimes be unstable since they are statistically focusing on the data patterns that stay on the surface of probability distributions. What if behind the probability distribution, we assume there exists an invariant mechanism that governs how the data is generated?

Take physics, we know force is the cause of changing the state of motion. The principle about *Newton's laws of motion* is invariant, remaining independent of how large or how small the force imposed on an object (hold constant the mass). Similarly, if there exists an inherently embedded "law of causality" that is dominating how our world is functioning, we tend to believe such a causal mechanism is independent of its cause. The idea behind this intuition is the celebrated principle "**Independent Causal Mechanism**".² It imposes restrictions on the data-driven approaches to recognize patterns with *causal significance*.

Let us take a step back, say, applying this principle over multiple factors, we will have a topological structure permitting multiple invariant mechanisms. Notice that this is the heart of the *structure-based causal models* since the invariance implies that each mechanism is an "**autonomous module**".² Keeping this in mind, I want to further borrow the words from the book *Causation, Prediction, and Search*: "*Intervention and counterfactuals are obtained with a suitable metaphysical gyration.*" I believe that the "suitable gyration" is partially come from "gyrating the input for a causal diagram", if the causal model behind the causal diagram entails the *autonomous* mechanism.

From my perspective, I think this is the key of forming the world view of causal diagrams. Think of events weaving our lifetime, for instance, as a topological structure based on their causal relations. Obviously factors such as social wealth and political power are crucial to resulting in a majority of events (e.g. business success); whereas tracing back to fundamental factors, we might sense that some factors, such as inher-

REFERENCE¹:

"Structured-based" originates from one of the top-level notion named "structured-based counterfactuals" shown in *Causality*, Chapter 7. We can see the structure of causal diagrams acting as a "bridge" to convey generic information (what we typically call as the "background information") between a current world and a "counterfactual world".

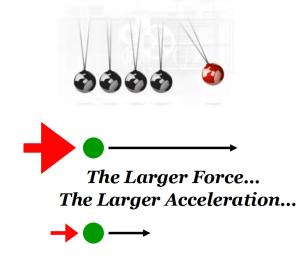


Figure 24: Newton's laws of motion, which is an established and invariant mechanism in physics.

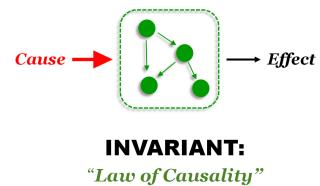


Figure 25: The expectant "laws of causality" implied by the well-known assumption of Independent Causal Mechanism (ICM).

REFERENCE²:

Assuming invariant data generation mechanisms gives an arrow (within a diagram) an unambiguous meaning of "governance". Readers can find out these ideas in the book *Elements of Causal Inference*.

ent characteristics, potential circumstances, or even ultimately one's fortune, are substantially some of the causes that largely influence our destinies. Plus, it seems that, the more we close to the "fundamental top factors" of causal relations, the weaker the connection among factors might be. Wealth and power might be closely related, whereas one's good and bad luck over his or her lifetime might be independent against other factors.

Interestingly, if we accept the *causal significance* from this moral, I would like to say that the rest of the content as to causal diagram learning (including mathematical theory) is, more or less, having something to do with this fundamental ideas.

Causal Markov Framework

Let us directly start to discuss the celebrated ***Causal Markov*** framework with an assertion from Judea Pearl's book *Causality*:

*"If we accept that causation **cannot** be inferred from statistics alone, then the Markovian equivalent models is inevitable".*

Allow me to throw out the terminologies there: "*Markovian*" describes a certain property that is able to identify a certain (unique) diagram; while the term "*equivalent*" partially refers to the several diagrams sharing that same property, meaning they are unidentifiable. Understanding the *Causal Markov* condition is a great way to grasp the ins and outs of mainstream *causal discovery* frameworks. My favourite metaphor is — "When the world of causation casts a shadow, it left with the property of **conditional independence**³ in the world of statistic".

The *Causal Markov* condition will become insightful in an analogic thought-experiment in which (causal) arrow connections can be viewed as a pattern of the "flow of water". Suppose I am a maintenance plumber, and my work is to help check the direction of the flow of water within a pipe. I will manage this by "locking" or "blocking" the junctions along the pipe, and observing whether there is any change of the state of water flow (Figure (18)). Unfortunately, if chances are, directions of the flow run in the three other ways (Figure (19)), I have no idea to tell which one it is. The moral is to imagine the "*conditional independence*" as if "locking" a pipe junction that curbs the water dependence along the pipe. Therefore, we master something like "*causal-information-flow*" just as if we are controlling the dependence of a pipe's water-flow. However, since we are merely work on the outer structure of a pipe, knowing nothing inside the pipe, "*the inaccurate check about the water flow direction is inevitable.*"

For another thought-experiment, envision a computer with a single push of the power button to start, allowing both the corresponding mouse and the keyboard to work normally (e.g. mouse \leftarrow button \rightarrow keyboard). At the first glimpse, it seems slightly contrast with the *Causal Markov* framework: Despite of the total control of the button, the odds are that the keyboard is able to work more normally while also observing the mouse working normally. Say, consider the following possible cases: (i) normally, the keyboard and mouse work when the button on; (ii) normally, the keyboard and mouse are stilled when the button

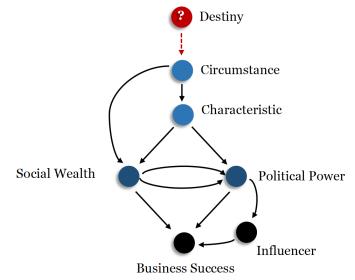


Figure 26: Toy causal diagram describing one's lifetime events as a topological structure.

REFERENCE³:

As for conditional independence, there is a popular example of alarm-smoke-fire in *The Book of WHY*. Given a causal diagram "fire \rightarrow smoke \rightarrow alarm", one knows that there is essentially no causation between fire and an alarm (namely they are independent), given the condition that smog is observed beforehand.

As Judea Pearl said at the beginning of his book "*Causality*" that, quote, "Conditional independence is the heart in causal modelling."

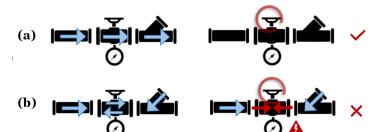


Figure 27: Check the direction of the flow of water by turning off the valve at the pipe junction.

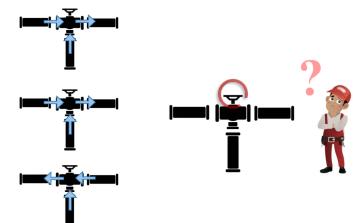


Figure 28: The three ways of water flow that cannot be determined, because of the same state of the pipe after turning off the valve.



Figure 29: Moral of the "pipe story": causal-information-flow in diagrams.

off; (iii) Unnormally, the keyboard and mouse are stilled when the button on, implying something (e.g. inside circuit) broken. It is the (iii) case that confuses us — given the control (the button is on), the working states of the keyboard and mouse are dependent (relevant): working or not working simultaneously.

In fact, it is the (iii) scenario that lulls us into a false sense of indeterminism. In other words, "two" of the manipulations (on and off) of the button do not fully determine the matchable "three" of the possible states of the computer. This phenomenon is conform to the characteristics of a **indeterministic system**⁴. In a "microscopic" system (Figure(21-b)), however, these devices are fundamentally "deterministic" when we totally control the inside circuit (rather than the outer button). We refer to this micro-level system as the **pseudo-indeterministic system**. The term "pseudo" hints that the computer can only be stuck by "**exogenous**" (unnatural) errors (e.g. the circuit cannot work in my computer, but it works well in yours). What we should learn from this thought-experiment is: The *Causal Markov* framework only holds in *deterministic systems* or *pseudo-indeterministic systems*; it cannot hold in *indeterministic systems*.

Coming to an end, do you still remember the metaphor? Though "casting a shadow" offers a few characteristics to take a closer look to causation, yet, it leave a distance for us to the elusively real appearance of causation.

The Key: Noise Perturbation

To point out a clue, I will introduce a causal assumption named **independent noise** implied by the *Independent Causal Mechanism (ICM)*. Virtually, it is consistent with the *independent exogenous errors*⁵ that I mentioned above. Notice that when the errors over a system is *exogenous*, it will consequently produce the *pseudo-indeterministic system*, which is exactly giving rise to the **Causal Markov** systematic natures⁶. Thus, we see how *causal discovery* frameworks have been emphasised from different views, while ultimately boiling down to a quintessence of *causal significance*: perturbation of noise.

So, why is an independently perturbative noise (sounds like a bad thing) so imperative for causality? Personally I believe it entails a type of metaphorical evidence of "how God create the world". Whenever there is an assignment of the certain cause from God, there is randomly an assignment of a mild perturbation — the noise where God encapsulates the potential uncertainty of the remainder. In other words, an "unexpected perturbation" is necessary since, compared to the omniscient God, it reflects our ignorance about the world and our destiny.

Following that logic, if one has entirely controlled all certain factors via **manipulation**, **condition**, and **intervention**, then we believe that not any correlation should have left because the remainders are insignificant perturbative noises. **Therefore, this idea has become the heart of modern causal diagram learning frameworks**. As romantic words I would like to say:

"Causation is given birth from perturbation."

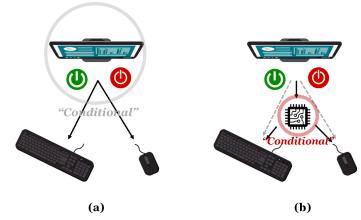


Figure 30: Examples as to the "indeterministic system" (a) and the "pseudo-indeterministic system" (b).

REFERENCE⁴:

Terminologies (including the following) are corresponding to ideas in the book *Causation, Prediction, and Search*. One conclusion I borrow herein, on the other side, involves the fact that the Causal Markov condition will be naturally satisfied within a deterministic system or pseudo-indeterministic system (discussing in the following).

NOTES⁵:

The term "exogenous errors" are technically referred as to "random (namely independent) noises" that are normally the hidden factors beyond one's control. Thus, (unknown) variables which characterize this phenomenon are called as the exogenous variables over a system.

REFERENCE⁶:

The pseudo-indeterministic system naturally yields the Causal Markov assumption. This has been proved in *Causation, Prediction, and Search*.



Figure 31: How God create the world? Does God play dice (to make some noise) but essentially conquer the causation over our destiny?

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of *Judea Pearl's* points, quote, "*it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics.*" Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

Remembering at first, I was seemly told that causal models are particular forms inherited from dominating models such as *Bayesian Network (BN)* and *Structure Equation Models (SEMs)*. A bit surprised, however, I later recognize that *BN* and *SEMs* have fundamentally derived from their own causality definition.

Firstly, we interpret the *Causal Markov* condition from the construction of *BN*, whose skeleton commonly represents for the *direct acyclic graph (DAG)*. Readers who are familiar with *BN* may know that an "engine" to infer the arrow between cause (*C*) and effect (*E*) over *BN* lies in the celebrated *Bayesian formula*:

$$P(C|E) = \text{Bayesian}(P(E|C), P(E)). \quad (48)$$

The inversion property though, seemly blurs the *BN*'s causal semantics, permitting it to compute "unrealistic probabilities" $P(\text{Cause} | \text{Effect})$ where the effect happens before the cause. Naturally, the arrow within the network should have entailed causality on its own, meaning that not only *BN* is empowering probability across the *DAG*, but in consistent with our cognition where causal relations in real world prone to be sparse¹.

Speaking of sparsity, in graphical fashion to model the sparsity, we avoid the fully connected *DAG* or *BN* by using:

- the symbol " $\perp\!\!\!\perp$ " describing *C* and *E* are "separated";

and projecting the sparsity relatively to probabilities with:

- the symbol " $\perp\!\!\!\perp$ " implying independence between *C* and *E*.

Modeling the sparsity property amounts to seek for a connection, from a graphical perspective to a probabilistic one, of modeling the independence property:

$$(C \perp\!\!\!\perp E)_G \Leftrightarrow (C \perp\!\!\!\perp E)_P. \quad (49)$$

Thus it has surreptitiously brought a celebrated criteria into our discussion: ***d-separation***, a criteria advances constructing connections between graph and probability². To approach this, the *d-separation* criteria needs to clarify it is who (*W*) that separate(s) *C* and *E* from each other, also depicting *W* as the one should be set in condition on in probabilities:

$$(C \perp\!\!\!\perp E | W)_G \Leftrightarrow (C \perp\!\!\!\perp E | W)_P. \quad (50)$$

More formally, on top of *BN*'s economical representation of joint probability functions leading to factorization commonly shown in Equation(51) (denote parents as *pa*):

$$P(V) = \prod_{V \in V} P(V | pa(V)). \quad (51)$$

OUTLINE:

- Bayesian Net and D-separation
- Structure Equation Models

NOTES¹:

Thus we literally should not draw an arrow between *C* and *E* in *BN*, if there is not undoubtedly causal relations between them.

REFERENCE²:

In Chapter 6 of the book *Elements of Causal Inference*, the notion of *d-separation* has all been a preferred introduction before talking about the (constraint-based) causal discovery.

BN also acts as a "carrier" of the *conditional independence* relationship, if running through the graphical angle. Which is to say, shown in Equation(53), applying the *Markov* condition² (simply denoted as $Ma(\cdot)$) with respect to *BN* (denoted as graph \mathcal{G}) should yield a list of relationships of *conditional independence* among variables V (here we denote $V^* = V \setminus \{V_i, V_j\}$).

$$Ma(\mathcal{G}_V) = \{(V_i \perp\!\!\!\perp V_j | V_k) \mid V_i \neq V_j, V_k \in V^*\}. \quad (53)$$

The point is, back to the probability angle, we claim the ***Markov compatibility***³ property if the (statistical) *conditional independence* implied in the distribution $P(\mathcal{G}_V)$ perfectly coincide with the one implied by the *Markov* condition $Ma(\mathcal{G}_V)$:

$$(V_i \perp\!\!\!\perp V_j | V_k)_{Ma(\mathcal{G}_V)} \Leftrightarrow (V_i \perp\!\!\!\perp V_j | V_k)_{P(\mathcal{G}_V)}. \quad (54)$$

The *Markov compatibility* reflects a crucial idea behind modeling *BN* since $Ma(\mathcal{G}_V)$ characterizes the modeler's (causal) assumption (e.g. in form of the "sparsity" of graphs) towards the data generation procedure, which ultimately leading to (statistical) *conditional independence* over empirical distribution $P(\mathcal{G}_V)$.

Nevertheless, computers obviously do not bear knowledge about $Ma(\mathcal{G}_V)$ to further ascertain the compatibility. In fact, we teach machines and design algorithm to obtain them based on, again, the significant *d-separation* criterion. As a consequent result, the *d-separation* criterion naturally becomes the product of the *Markov* condition, and perhaps more significant, it is more reasonable to understand the *d-separation* criterion from the causality point of view. We determine the directions as to the "flow of causal information" (e.g. $* - *$ represents possible directions \rightarrow or \leftarrow) via (indirectly) determining the separation ability of some "junction"⁴ variables V_k :

$$V_i * - * V_k * - * V_j. \quad (55)$$

Secondly, the *Independence Causal Mechanism (ICM)* would start from the definition of *Structure Equation Models (SEMs)* shown in Equ(56). The term "structure" implies the noise N should be irrelevant with C . In other words, no matter the statistics relations between C and B , the (causal) coefficient β^5 is invariant, which exactly conforms to the notion of "invariant causal mechanism" and "independence noise assumption".

$$SEMs : E := \beta C + N. \quad (56)$$

In fact, the *Structure Causal Models (SCMs)* are non-parametric versions of the *SEMs* ($E := f(C) + N$). weaving the *SCMs* into the *ICM* framework, we naturally conclude:

$$(P(C) \perp\!\!\!\perp f) \text{ or } (P(C) \perp\!\!\!\perp P(N)). \quad (57)$$

Therefore, from the causality point of the view, I wish I can briefly introduce my readers these two models *BN* and *SEMs* retrospectively (or maybe that is the exact way we should be), so that readers might get a sense about how the modern causal discovery frameworks such as the *Causal Markov Condition* and the *Independence Causal Mechanism* have developed based on these conventional models.

NOTES²:

The *Markov* condition has also been used as the definition as *BN* (denote descendant as de):

$$X \perp\!\!\!\perp V' \setminus \{pa_X \cup de_X\} \mid pa_X. \quad (52)$$

NOTES³:

The *Markov compatibility* shown in Equation(54) further clarifies the connection between graph and probability based on Equation(49).

NOTES⁴:

Recall the "pipe-story" example in the introductory part. Given a pipe connection I-K-J, if we turn off the valve on the junction K, then we "separate" the water "dependence" between I and J.

REFERENCE⁵:

In the literature, the original definition of the coefficient β in *SEMs* is the "causal effect" (within a linear system). However, this fact has become slightly ambiguous since the manipulation $C = \frac{1}{\beta}(E - N)$, making the coefficient β becomes trivial in an absolute equation.

Borrowing the statement from Judea Pearl in the book *Causality*: "Now, speakers of the *SEMs* language are in search of its meaning."

4.1.2 Restrictions Give Rise to Identification

The previous section gives an overview of causal diagram learning frameworks, specifying why "*causal discovery*" is meaningful at the first place. To further apply these principles over data, we now need to focus on the "***identification***¹"— identify causal diagrams from observational data. The whole picture of this section is that, bonding particular restrictions to the fundamental principle will draw a (strict) line between what is causal and what is non-causal.

Restrict the Trait About Causal Dependence

We have already seen how *Causal Markov* framework builds a "matching pattern" between the world where causal graphs sit and world where our every-day events happens (e.g. We can associate the causal connection with the dependence of water flow in our daily lives). Yet this matching pattern is not perfect enough. Though the *Causal Markov* framework defines the "normal causal relations", it does not limit the emergence of "abnormal causal relations²". The causal diagram entailing the abnormal relations, however, certainly would not be the truly justified one for which we are seeking. Hence, when we require to ensure a "normal relation" and end up reaching a "perfect match pattern", some additional restrictions are necessary.

Here the restriction is: constraints of ***Causal Minimality*** and ***Causal Faithfulness***. In fact, these are two of the names that you would often hear when digging deeper into the *Causal Markov* condition.

To illustrate them straightforwardly, back to the thought-experiment context³ where we started introducing the *Causal Markov* condition. Suppose the pipe connecting three of the supply places ($A \rightarrow B \rightarrow C$). I will ensure the water supply in places A and C functioning normally via pipe $A-C$, after I worked on locking the valve located in B and consequently observing the flow of water from A to C . Such a simple logic, namely the water supply between two places are dependent whenever given a (normal) pipe, is where the "*Causal Minimality* restriction" have stood. From personal perspective, I prefer to think of the "*Minimality*" like this: I would not be delightful if there is nothing happen on pump $A-C$ given a locked valve at B — since initially I anticipate a "minimum amount" of checking workload, whereas now, there must be something abnormal occurs (e.g. the pipe that should have carried water might has something broken inside). In other words, the real situation does not fit well the "normal causation", unfortunately leading to some extra checking work (compared to the minimum checking work at the begin, in which "the causation functions well given a pipe").

Applying the moral about this plumbing-story into the other thought-experiment⁴, we expect the same picture where thing can function normally (Fig(22)). The *Causal Markov* condition ensures the working mechanism itself of either the mouse or the keyboard is "irrelevant" with each other — they are working or shutting down simultaneously just because someone con-

REFERENCE¹:

A "learnable" causal diagram from data should provide us an unambiguous way to interpret the data. We anticipate the data pattern is precisely represented by an unique causal diagram (because something unambiguous and unique means it can be used for precise prediction). In terms of graphs, we are expecting the particular nature of the causal diagram: Structure Identifiability. Relevant topics can be found in the book *Elements of Causal Inference*, Chapter 7.

NOTES²:

In order words, the *Causal Markov* framework will still even hold on, even giving mixtures of both reasonable and absurd causal relations over a causal diagram.

NOTES³:

Recall the thought-experiment in the previous section. Assuming my work as a plumber is to help check the direction of the pipe's water flow, I will check the state of the water flow by "locking" or "blocking" the junctions along the pipe.

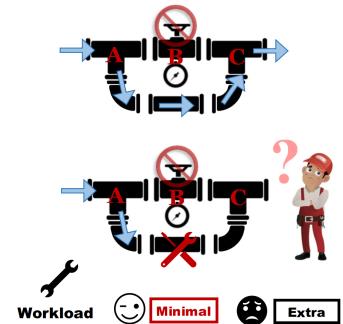


Figure 32: The workload would not be the minimum one if the relation of water dependence given a pipe does not fit the normal situation.

NOTES⁴:

Review the example I introduced in the previous section, where the *Causal Markov* condition ensure the working states of mouse and keyboard are "independent" if one has controlled the battery button.

trols the computer's power switch. The "*Causal Faithfulness*" restriction", on the other hand, ensure the "relevance" of the connection among the devices, making both mouse and keyboard normally depend on (or "listen to") someone's control.

While both of them mount constraints on dependence that a "normal regularity" should entails, the "*Faithfulness* restriction" is stronger than the "*Minimality* restriction". If possible, we tend to rule out everything absurd, abnormal, and "unfaithful" to our cognition. For instance, we do not expect that, sometimes, the flow of water would have maintain still even the valve has been released; we also do not expect that our mouse or keyboard get faults every so often just without causes (Fig(22)). These (stronger) expectations are all encapsulated into a general and thus stronger condition: the "*Causal Faithfulness* restriction".

The pivotal idea is: These restrictions, such as "*Faithfulness*" and "*Minimality*", are not concepts as obscure as it literally looks like. Conversely, they are relevantly on behalf of our normal cognition: anticipating the normal relations among things can perfectly match with the regular traits behind it. Hence, that relations, entailed by certain regularity's traits, are more likely to be the causal relations⁵, making the structure representing a regularity more likely to be an uniquely appropriate causal diagram—the causal diagram for us to identify from data.

Restrict the Complexity About Causal Mechanisms

We probably hear of the "***Bayesian Method***", the "inverse probability" analysis left by *Thomas Bayes* who was once a mathematics geek. In the following example⁶, let us briefly see the relationship between that "inverse probability" and our causal cognition. Suppose Bobby throws a ball toward a window, it is easy to predict the mechanism where the action of throwing may break the window most of the time—due to the "law of physics" (e.g. ball's mass, force, velocity, ect). Given a broken window, however, it seems hasty to making a deduction that it is Bobby who threw a ball toward the window. Though analysis of "*Bayesian inverse probability*" is designed to help break this **cognition asymmetry**, we should recognize the prompt by our causal cognition: the "inverse mechanism" is complicate. If the tricky backward mechanism does exist, it will require fitting multiple background conditions (Fig(25)). In contrast, when a general mechanism can virtually fit well, we tend to believe such a mechanism is more likely to entail causation, partially because the simple one captures the "law of causation".

Restricting the complexity indicates that, practically, learning a one-size-fits-all mechanism is more preferable. This does not necessarily means that the causal mechanism must be truly simple. In fact, we just tend to bear in mind a priori: the simpler a hypothetical mechanism is enough to fit the data, the more likely the mechanism can genuinely describe the causal relation⁷. After all, most of fundamental effective principles or guiding rules are invariably straightforward, wouldn't it?



Figure 33: The causal faithfulness restriction expects an absolute dependence (the arrow are highlighted as red), without any "accident" fault such as the "bug" occurring in the electric circuit.

NOTES⁵:

e.g. the normal relations of water dependence in the plumbing or electric component connection among the devices.

REFERENCE⁶:

Again, I briefly borrow an example context shown from *The Book of WHY*, Chapter 3.

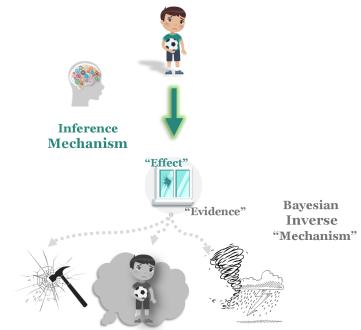


Figure 34: The cognition asymmetry: which of the mechanism is the causal one (simple or complicate)? Given a broken window, several things might be associated: cracked by a hammer, broke by throwing a ball, attacked by extreme weather.

REFERENCE⁷:

This is one of the important opinions of the book *Elements of Causal Inference* (Chapter 4).

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of *Judea Pearl's* points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

In order to yield a identifiable result, different restrictions substantially represents the intensity of different "standards" that we expect to reach to narrow the connection gap between probabilities and graphs. Fig (35) explicitly illustrates the range of relatively different standards amidst causal discovery, which also shown as different constraints from an "omniscient view".

Methods following this philosophy, such as the SGS algorithm and the *PC*¹ algorithm, discover the causal diagram by discovering the statistic traits of independence from the general hypothetical dependence (imagine a solid circle representing all the possible dependence in statistic). To connect it with graph, take *PC*, the methodology starts with a complete (undirected) graph with full connections and then conduct independence test to remove unnecessary connections (imagine a hollow circle representing the discovering independence is expending now). It stops until no connection can be removed anymore and turn to finalize causal directions — expecting remaining connections might perfectly match their causal significance from the true causal (directed) graph. In contrast, if it keeps on removing the connections that entails causation, then it will enlarge the "unnecessary independence" that does not be included by the *Causal Markov condition* (from an "omniscient view" (shown in Fig(35))), we can see when the expending hollow circle should finally stop at an exact "border" without further "breaking in" the vicinity range, which should be promised by the *Causal Minimality* and *Faithfulness* restrictions).

Though saying "remove unnecessary connections" is similar as saying "expend the hollow circle or enlarge the unnecessary independence", notice that the hollow circle can still continues to expend even if we keep unchanged the connections², which is why I draw the *Faithfulness* restriction at the outermost. As mentioning previously, the *Faithfulness* restriction ruling out every "unfaithful" dependence is stronger than the *Minimality*.

Thus, the range within that border can be viewed as reaching the standard of the *Markov* restriction from the direction of sufficiency; the range outside of that border is able to be seen as reaching the standard of the *Faithfulness* restriction and the *Minimality* restriction from the direction of necessity.

Combining together, two types of the restrictions expect an algorithm to end up with a narrow border with few "margins". Eventually, the narrower the "margin gap" is, the less likely the final graph equivalents can satisfy the same restriction, meaning reaching a perfect (unique) match.

OUTLINE:

- Minimality Assumption
- Faithfulness Assumption
- Statistical-Indistinguishability
- Faithful-Indistinguishability

REFERENCE¹:

These are state-of-art constraint-based causal algorithms that readers can learn more in the book *Causation, Search, and Prediction*.

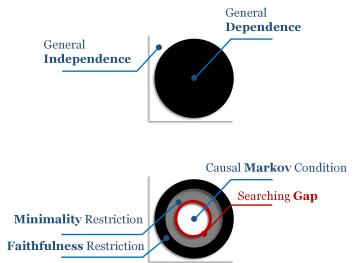


Figure 35: The range of relatively different standards amidst causal discovery. Given initially general dependence and independence from data, "Markov constraints" specify the (conditional) independence that implies the causal significance; and the "Minimality and Faithfulness constraints" ensure the remaining dependence. Searching by the causal discovery algorithms, final estimated causal graphs satisfying these constraints are involved in the margin (gap), which are also referred to as "Markovian Equivalents".

NOTES²:

This is possible since "accidents" or "unexpected noise" might impact the statistic independence test (e.g independence test over mixing data can sometimes trigger "abnormal" independence, which is not implied by the separation in the graph).

Additionally, the *Minimality* condition puts a constraint on BN, which confines the network connection to reach "the most economical" level. When it satisfies the *Minimality* condition, edges can not be taken away from BN since that operation would "maximize" the degree of *conditional independence* beyond the range yielded by the *Markov* condition.

$$\Gamma[Ma(\mathcal{G}_V)] \leq \Gamma[Ma(\bar{\mathcal{G}}_V)], \quad (58)$$

where $Ma(\cdot)$ simply denotes applying the *Causal Markov* condition, Γ denotes the volume of independence relations entailed by the *Causal Markov* condition.

The *Causal Markov* framework constructs the elementary connections between the probability and the graph (the same as constructing BN based on the *Markov* condition). We typically need to impose the additional *Minimality* limitation to narrow down the selecting range, and hopefully ends up with nearly a one-versus-one connection.

Interestingly, flipping the case around yields the standards of "**unidentifiability**": the *Minimality* and *Faithfulness* unidentifiability³. However, the *Minimality* unidentifiability is stronger than the *Faithfulness* one. We can clearly aware this by drawing a number axis in Fig(36).

Why do I repeatedly discuss different intensity relative to the "standards" in terms of identifiable and unidentifiable cases? If we agree to this matching methodology applied by constraint-based algorithms, it is more reasonable to accept that most of the unidentifiability raised by unmeasured common causes (namely the latent confounder, will be discussed in the next section) are at the end of the other spectrum. So aside for the identifiable causal diagrams, we also need to contrive the other expressive causal diagrams to represent unidentifiability.

Strictly speaking, discovering independence largely lies in discovering the so-called "**V-structures**" — "sub-components" contributing to finalizing causal directions. Attain equivalents is essentially common since we cannot absolutely determine all the directions⁴. However, by learning how intense each restriction is, we bear knowledge about characterizing causal graphs into equivalent classes (based on varying degree of restrictions).

This philosophy, from my perspective, is the quintessence of independence-test-based causal discovery. It is a beauty that shrinking the gap between graph and probability step-by-step; Yet one might still be wandering whether it is possible to make a "huge step" leading to more directed identification, which brings us to the next restriction: causal mechanism complexities.

Notice the fact, that we do admit different probability factorization to numerically describe a joint distribution $P(C, E)$:

$$P(C, E) = P(E | C)P(C) = P(C | E)P(E). \quad (59)$$

Whereas this equivalent essentially does not account for $P(C, E)$'s generation procedure. Behind the fixed and visible $P(C, E)$, we

REFERENCE³:

The two types of unidentifiability are referred to the "Strong-Statistical Indistinguishability (s.s.i.)" and the "Faithfulness Indistinguishability (f.i.)", which are introduced in the book *Causation, Prediction, and Search*, Chapter-4.

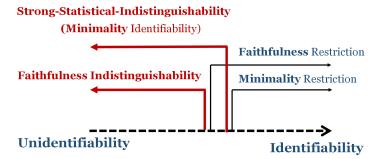


Figure 36: Illustration of different causal intensity (restriction) of identifiability and unidentifiability.

NOTES⁴:

If we accept that causation cannot be inferred from statistics alone, then the Markovian equivalent models is inevitable;

In another words, that "margin" we mentioned previously is always existing no matter how "narrow" it is.

tend to believe there exists an intangible class of mechanisms or functions $\{\mathcal{F}\}$ that describes how $P(C, E)$ is given birth to:

$$P(C, E) := \mathcal{F}(P(E | C), P(C)). \quad (60)$$

The moral is, mechanism \mathcal{F} seems "natural" from cause C to effect E , not vice versa⁵ (compared to the factorization in Equ(59)). Such "common sense" is the essence of restricting the class $\{\mathcal{F}\}$ into a "simple" spectrum because a reversing mechanism from E to C is prone to be unnatural, weird, and — complicated⁶.

We might then wish to learn physical formalism of \mathcal{F} . Recall our example in the introductory part that "most of time forcibly throwing a ball will result in a broken window". Here the words "most of the time" hints that, other factors (e.g. temperature, places, the window material) are trivial compared to the main cause (throwing a ball). This rough intuition gives a sense about how to select a concrete restricted mechanism \mathcal{F} by Equ(61):

$$\mathcal{F}(P(E | C), P(C)) := f(P(C)) + P(E | C), \quad (61)$$

where we break down $P(C)$ and $P(E | C)$ within an **additive causal model**, imposing certain functional mechanism f merely on $P(C)$. In the **additive model**, "other factors" enveloped by $P(E | C)$ are trivial with only the "shifting"⁷ impact on the "main cause" C . Since we expect $P(E | C)$ should be insignificant, the common **Gaussian** distribution is often assigned at $P(E | C)$.

Finally, notice that we do not require the form of f most the time — generally f is non-linear. However, if f is confined to be linear (typically a linear mechanism is slightly trivial than the non-linear one, e.g. exponential functions grows fast than linear functions), we might naturally want to impose restriction on $P(E | C)$, making it slightly non-trivial to reach a "balance". In fact, when f is a linear, we exactly assign the **non-Gaussian** distribution to $P(E | C)$ to maintain its causal identification.

These are namely two of the representative causal models in "functional-based" methodology: **ANMs** and **LiNGAM**. Readers can learn more about these "identifiable functional model" in the book *Elements of Causal Inference*, Chapter 4 and 7.

NOTES⁵:

Alternatively, the same class $\{f\}$ will be less persuasive when applying to the case in which C and E are swapped.

REFERENCE⁶:

In (functional-based) causation, one of the popular way to model the complexity involves the Kolmogorov complexity as an algorithmic information theory.

Meanwhile, this will also become natural from the machine learning point of view, where we limits our "fitting function" within a "reasonable" model class (e.g. polynomial function) in standard tasks such as regression and classification.

Relevant details can be found in the book *Elements of Causal Inference*, Chapter 7.

NOTES⁷:

e.g. non-additive models such as $\mathcal{F} := f(P(C)) \times P(E | C)$) might seem less intuitive, as the effects of $P(E | C)$ within \mathcal{F} might make the main cause C more susceptible to the other factors.

Causal Books References

Causation, Prediction, and Search: The following pieces involve testing content. Video provides a powerful way to help you prove your point.

- **5.4 New Algorithms:** The following pieces involve testing content. Video provides a powerful way to help you prove your point.
- **5.4 New Algorithms:** The following pieces involve testing content. Video provides a powerful way to help you prove your point.

4.2 Challenges: Unmeasured Common Cause

So far, the context of causal diagram learning has not involved the circumstance in which the variable system represented as the dataset can be essential "incomplete" — some variables in a complete system are just missing.

In other words, compared to an omniscient causal diagram reflected by a complete system, most of the time we need to face the challenge as to learning a relatively "smaller" causal diagram given a "smaller" dataset observed from the system.

Following one of the primary ideas¹ of the book *Causation, Prediction, and Search*, presence of the **unmeasured common cause (latent confounder)**, however, should arguably be the intriguing part of causal diagram learning since *unmeasured common causes* are virtually ubiquitous in practice, and it has long been a conundrum waiting for crack. As for the last section of this paper, I will attempt to provide some open perspectives relative to the problem.

Topic: Education, Experience, and Salary

Let us continue the discussion on the example "education, experience, and salary" from *The Book of Why*. Suppose we were conducting control experiments looking for a certain firm to see whether the higher education levels potentially contribute to an employee's salary. Meanwhile, it is known that experienced employees are also likely to get a high payment (experience → salary). Thus such a relevant factor should be controlled beforehand, ensuring unbiased measurements of the causal effect we are truly interested in (education → salary).

Given the candidates of employees, and before enforcing the controlled experiments, suppose we might not sure about the exact relation between the employee's education degree and the employee's working experience. The odds are, some candidates would refuse to engage further study in order to spend extra time in advancing working experience (education → experience). In contrast, it is also possible that the other candidates' excellent attachment in working experience just in turn hints that they are more excel at their study (experience → education).

Namely, the factor "experience" can stand for an intermediate factor (education → experience → salary); or it can also represent the **common cause (confounder)** of education and salary (education ← experience → salary). Thus, considering the difference, should our operation include measuring the "experience" factor before the experiments? Or should we just keep it as an *unmeasured common cause (latent confounder)*?²

Furthermore, if it were implied to directly ask people's salary, or if it were abstract to truly assess one's working experience (e.g. one's real competence reflected by the working experience is hard to be directly evaluated), then one of the approach is to design pieces of questionnaires to assembly characterize these factors we interested in. In other words, all the factor we truly care about have invisibly become the "superior" factors upon

REFERENCE¹:

The issue of the existence of latent confounder is the most important topic discussed in Chapter 6 and 7 in the book.

Actually, the book is exactly seeking for a better way for "prediction in causation" and "causal diagram search" when it comes to the context with latent confounders.

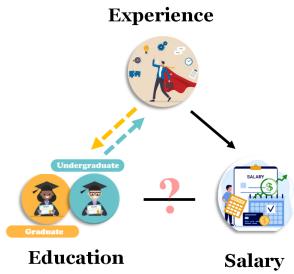


Figure 37: The popular example of "education, experience, and salary", with undetermined relations between "education" and "experience".

NOTES²:

To slightly remind, this choice before the controlled experiment becomes vital since "controlling the variable or not" should be cautiously treated when it comes to the causation with unmeasured common cause. We will see this opinion in the following.

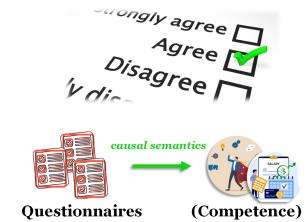


Figure 38: Use questionnaires to infer "superior" factors in the first place, because the "superior" factors (instead of the questionnaires) are in the causal semantics of cognition.

our cognition, and the interesting point is that these "superior" factors are now the *unmeasured common cause* with respect to the several questionnaires. The point lies in how could we infer the "superior" causal relations based on the "subordinate" information gathered from the questionnaires?

The Challenge of Latent Confounders: Horizontal and Vertical

Prior to learn about what adaptions could be made to fix the challenge, the essence of above examples can boil down to the following two perspectives:

- *Horizontal: Perception*³. Perception is often imperative than operation in causality. Does there exist a particular type of causal diagram which is capable of providing an elastic space to perceive the uncertainty from the possible unmeasured common causes?
- *Vertical: Elaboration*⁴. Low-level data measurement could be viewed as "multiple projections" from a single high-level (hidden) source that elaborates *causal semantic*. Does the "structure of projection" in turn leave us some insightful clues to detect the structure itself?

In terms of *perception*, we develop the concept from the causal path into the *inducing path*. As for *elaboration*, we would see a noble statistic constraint named the **Tetrad Constraint** beyond *Conditional Independence Constraints*. Above all, the following content we will introduce aims at providing powerful theoretical tools which serve as a guiding principle for the current causal discovery algorithms to tackle the challenge of causal diagram learning with *unmeasured common causes*.

From the Causal Path to the Inducing Path

Mentioned in the previous section, *conditional independence* is the heart of causal diagram learning framework. This framework, however, entails an implicit condition where all variables should be sufficiently observed. Lining up with the challenge of "unobserved variables" perception, we thus need to broaden the "horizontal boundary" that causal diagrams use to describe the connection among variables.

Yet, if not for *unmeasured common causes*, it does make sense to purely concentrate on "removing" the redundant dependence by the signals (causal significance) as to "conditional independence" (illustrated in Fig(41)), meaning we are confident that the remaining dependence will undoubtedly imply causation. Conversely, causation could not longer be simply obtained by "removing" the independence since the dependence is now a mixture of both causation and, the trickily spurious correlation from *unmeasured common causes*(illustrated in Fig(42)).

Roughly speaking, the "*conditional independence*" herein can be similarly attached by the controlled experiments that we mentioned in the previous few paragraphs. Thus, the general dependence between two variables "education" and "salary", for

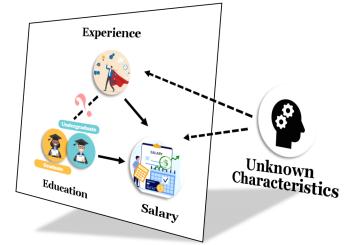


Figure 39: The challenge from the horizontal perspective: perception.

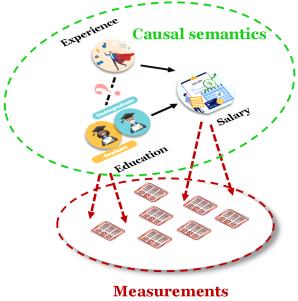


Figure 40: The challenge from the vertical perspective: elaboration.

REFERENCE³:

There is a similar opinion by Judea Pearl from *Causality, Models, Reasoning and Inference*:

"Generally, signal sensing is more fundamental to the notion of causation than manipulation".

REFERENCE⁴:

This "elaboration problem" with unobserved variables is introduced in the book *Causation, Prediction, and Search*, Chapter 10 and 11.

$$\text{Correlation} - \text{(Conditional Independence)} = \text{(Causal Dependence)} \quad (\text{Causal Path})$$

Figure 41: "Correlation is not causation." However, removing the conditional independence, the rest of the correlation is causation.

$$\text{Correlation} - \text{(Conditional Independence)} = \text{(Causal Dependence)} \quad (\text{Causal Path})$$

Figure 42: "Correlation is not causation", even the conditional independence has been removed, due to the existence of latent confounders.

example, might not be able to entail the causal significance (*conditional independence*) after the controlled experiments if considering the existence of *unmeasured common causes*. This possible uncertainty, can be described as there presents an *inducing path* between "education" and "salary".

In order to interpret the concept of the *inducing path*, let us back to our example. We describe a path as "inducing" because of the following dilemma:

- 1) Suppose we do not control the variable "experience", then we fail to infer causation between "education" and "salary" since the factor "experience" is possible to be an *unmeasured common cause* ($\text{education} \leftarrow \text{experience} \rightarrow \text{salary}$) leading to confounding bias.
- 2) Assume we do choose to control the variable. However, considering the other possible mechanism ($\text{education} \rightarrow \text{experience} \rightarrow \text{salary}$), conditioning on the factor "experience" will mistakenly "induce" spurious correlation.

Let us discuss more about the meaning of dilemma (2).

Suppose that the people who refuses to pursue the higher educational degree (in order to spend more time in accelerating the working experience) happens to be the people who owns the "first chance" to work in certain companies. However, we do not have any priori knowledge with respect to what might offer someone such a "first chance" to work. Perhaps someone just happens to have strong connections with the leading figures in the company they work, and that connections bring the person more potential opportunities to obtain high salary in a certain company ($\text{experience} \leftarrow \text{social connection} \rightarrow \text{salary}$). Notice that such a factor that has already beyond our control is the (real) *unmeasured common cause*.

The most tricky problem is, when we continue to conduct controlled experiments based on that "particular" group of people, we "induce" the spurious correlation stemmed from an unknown common cause ($\text{education} \leftarrow \text{social connection} \rightarrow \text{salary}$): people who obtains relatively the lower educational degree tend to achieve a high payment. But maybe the truth is — these are "particular" group of people who has the great opportunity and social connection to work early and ask for handsome income.

Thus, an *inducing path* can be activated by our control and simultaneously induces another *unmeasured common cause*. From the perspective of causal diagram (shown in figure(43)), this phenomenon amounts to activate a "V-shape connection"⁵, given the possible relation $\text{education} \rightarrow \text{experience}$.

Moving on, recall that in the previous section we have known that the ideation of causal diagram learning lies in uniquely matching the patterns between graph and probability. Meanwhile, the inevitable weakness involves the issue of equivalent classes, and we did not go any further about the representation of equivalent classes. But, when *unobserved common causes* come along, it should not be overlooked. This is because the attentions of representing the equivalent classes are gathered

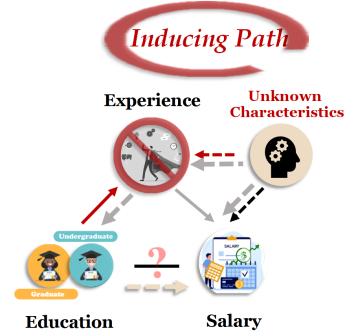


Figure 43: Illustration of an inducing path: $\text{education} \rightarrow \text{experience} \leftarrow \text{unknown characteristic} \rightarrow \text{salary}$. Notice that the "sub-path" $\text{education} \rightarrow \text{experience} \leftarrow \text{unknown characteristic}$ forms a "V-shape connection", which is prone to induce bias after controlling the collider "experience" in that "V-shape connection".

REFERENCE⁵:

The hazard of the "V-shape connection" (or V-structure) is highlighted in *The Book of WHY*, Chapter 5.

from the "independence" to the "dependence" where the mixture of both causation and spurious dependence only increases the matching diversity and complexity between graph and probability. Therefore, this is the reason why we need to define the other **inducing path graph equivalents**⁶ (described in Fig(44)), where the symbol $o \rightarrow$ that shows a circle at the end of edge is partially the element of the *inducing path* (representing the possible relations \rightarrow or \leftarrow).

Having the guiding model in mind in which the uncertainty triggered by *unmeasured common causes* is elastically captured, one could image that causal algorithms are driven to uncover an **information-maximization graph equivalents**⁶ based on the observed data (e.g. trying to determine the $o \rightarrow$ symbol as many as possible). The *graph equivalents* is an emblem of perceiving the causal knowledge before starting the next-step operation. Specifically, state-of-art algorithms such as the PC and the FCI algorithm (Fig(45)) are developed in light of the different "object type" of the causal diagrams (e.g. *Markov equivalent* or *inducing path graph equivalent*), further reducing computational cost and gaining popularity over the last two decades.

From CI Constraints to Tetrad Constraints

Take causal semantic, the "structure of projection" entailed by latent confounders formally stands for "**elaborated models**", "**indicator models**", or "**measurement models**"⁶. Good news is, insightful clues to causal inference are still available and the idea behind it is about the "trade-off" — the model constraints to simplify usually brings the statistics strength to identify.

We assume *linear elaborated models*, for instance, where we believe the salary levels are uniformly and proportionally varying with the employee's working experience. We also presume **causal purity**⁶. Like, simultaneously observing an employee performing well on the task in both the questionnaires A and B, should be essentially interpreted as the result from the working competence of that employee, instead of "questionnaire A being the cause of B" or vice versa. Such model restrictions including *linearity* and *purity* actually remain ubiquitous for the purpose of simplicity, hence bring us additional insights to identify the causal structure even with the presence of *latent confounders*.

One type of such additional benefits has a popular name "**vanishing tetrad difference**" or "**tetrad constraints**"⁶". Notice that the *conditional independence* constraints in form of " A is independent to D given the condition B and C ", for example, are insignificant given the elaborated model. That is because when A, B, C, D act like the measurement of questionnaires, they are all confounded by superior latent factors in the higher causal semantic level. However, the *tetrad constraints* move beyond it, implying different intriguing mathematical patterns such as " $AC - BD = 0$ " (vanished into zero) that are consistent with the different models under *linearity* and *purity* constraints. Fig(46) and (47) illustrate how to utilize these patterns to cluster the "subordinate" information to detect the "superior" factors.

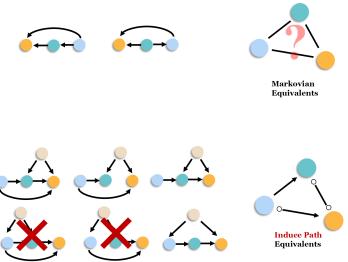


Figure 44: Different matching diversity between the Markov equivalent and the inducing path graph equivalent. The later one further consider the existence of latent confounder, and use flexible graph edge symbols to rule out some impossible cases in order to reduce matching diversity.

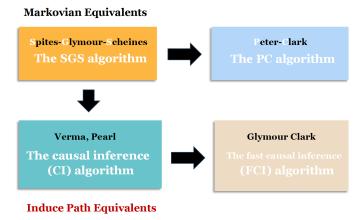


Figure 45: Development of the state-of-art causal diagram learning algorithms.

REFERENCE⁶:

These are technical but important terminology in causal diagram learning with latent confounders, referring to the book *Causation, Prediction, Search*, Chapter 6 , 10 and 11.

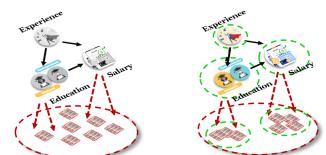


Figure 46: Intuition of the tetrad constraint. First step: using tetrad constraints to find out the cluster.

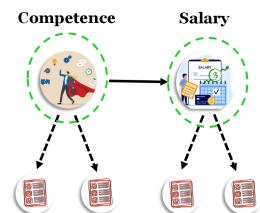


Figure 47: Intuition of the tetrad constraint. Second step: Infer causal relations among the superior factors based on tetrad constraints.

Mathematical Formalization

Reminder: I additionally attach this part alone to discuss the formalization as I fully convinced one of *Judea Pearl's* points, quote, "it is the formalization that eliminates the metaphysical controversy from causation by using elementary mathematics." Nonetheless, several formulas may be necessary to be included. I made it a relatively independent part so that readers interested more in general causal ideas can feel free to skip the following.

Suppose we have a (complete) vertex set V that satisfies **causal sufficiency**¹, then X and Y are not **d-separated**² by any subset Z of $V^* = V \setminus \{X, Y\}$ if and only if there is a **causal path** (CP) over V between X and Y .

$$CP(X, Y) \Leftrightarrow (X \not\perp\!\!\!\perp Y \mid \{Z \mid Z \subseteq V^*\})_G. \quad (62)$$

In contrast, assume that the observational set O violates **causal sufficiency**, in the sense that $O \subset V$. Again, X and Y are not **d-separated** by any subset Z of $O^* = O \setminus \{X, Y\}$ if and only if there is an **inducing path** (IP) over O between X and Y .

$$IP(X, Y) \Leftrightarrow (X \not\perp\!\!\!\perp Y \mid \{Z \mid Z \subseteq O^*\})_G. \quad (63)$$

The notion of *inducing path* therein brings a corresponding concept compared to *d-separation*: namely the **d-connection**. Notice that the notion of *d-connection* characterizes the general dependence (e.g. the causal dependence and the spurious dependence raised by *latent confounders*) associating with both the graphical and the probabilistic perspectives.

If one refers to the formal definition of *inducing path*, he or she would be told that two characteristics³ should be placed on the members of O^* along the IP that

- for every member $O_i \in O^*$ has $O_i \in (An(X) \cup An(Y))$,
- for every member $O_i \in O^*$ has $O_i = col(O_i)$.

Where we denote the collider as $col(\cdot)$ and the ancestor as $An(\cdot)$. From the above two characteristics we get a sense about why the *inducing path* has its name. Literally, condition 1) indicates that the ancestor relations (including the parent relations) are typically needed to be controlled (namely set into condition) to help us obtain possible conditional independence. However, condition 2) implies that the member we controlled will be exactly the collider, and the result of controlling a collider, is to "induce" possible spurious dependence that is yielded by the *unmeasured common causes*.

The graphical language for the *inducing path*, usually employed by *constraint-based* approaches, lies in developing the multiple connecting-symbols (Equ(64)) in order to characterize the uncertain information. Then, the uncertainty of the connections are expected to be narrowed down consequently:

$$\text{Edge Symbols : } \{o - o\} \geq \{o \rightarrow, \leftarrow o\} \geq \{\leftrightarrow\} \geq \{\leftarrow, \rightarrow\}. \quad (64)$$

Consideration of equivalent classes⁴ becomes inevitable since the phenomenon, in which a wide spectrum of **inducing graphs**

OUTLINE:

- Inducing Path and D-Connection
- Partially Inducing Path Graph
- Constraint-based Causal Discovery
- Possible-D-separation Set
- Definite Discriminating Path

REFERENCE¹:

In terminology, the variable system that is "complete" or "incomplete" (with or without missing variables) can be described as satisfying "causal sufficiency" and "causal insufficiency" respectively (See in the book *Causation, Prediction, and Search*).

NOTES²:

Recall that the d-separation criterion is naturally the product of the Markov condition. We discussed that in Section 4.1.

NOTES³:

It is imperative to understand the *inducing path* from an omniscient perspective. Which is, the *inducing path* is reasonably a continuous path over V^* , whereas only O^* are "visible" (observational) as well as that only O^* need to meet the characteristics.

NOTES⁴:

Recall the similar problem of the "Markovian equivalents" (without the latent confounder) we discussed in Section 4.1.

G' (the causal diagrams with the edge symbols in Equ(64)) share the same *d-connection*, tends to be common. The thing is: how could we generally summarize graphs from equivalent classes and depict a noble representative that covers whole uncertainty?

First of all, we need a **partially oriented inducing path graph** π over O that has the same vertices $V(\pi) = V(Eq(G'))$ ⁵ and the adjacency connection $ADJ(\pi) = ADJ(Eq(G'))$ with the inducing graphs in $Eq(G')$. Furthermore, not so strictly speaking, the π over O can be attained by "adding" together the $Eq(G')$ ⁶:

$$\pi := \sum_{G'_i \in Eq(G')} \Theta(G'_i), \quad (65)$$

Where we define the graph-aggregating function Θ as a set of constraint rules listed by *Causation, Prediction, and Search*. In this paper, I basically illustrate them as the following:

$$\Theta : \begin{cases} \leftrightarrow := (\rightarrow) + (\leftarrow), \\ o \rightarrow := (\rightarrow) + (\leftrightarrow), \\ o - o := (\rightarrow) + (\leftarrow) + (\leftrightarrow). \end{cases} \quad (66)$$

Keeping this in mind, in terms of causal discovery algorithms⁷, I attempt to synchronize the following introduction with and without causal sufficiency. What I hope is to offer a dynamic picture for algorithms' developments where my audience could pinpoint to the primary idea amid we tease out these comparable and improved versions.

When causal sufficiency comes along, the SGS algorithm starts from a complete graph and further prunes it (e.g. namely to "cut off" the (unnecessary) edge between any of the two variables X and Y in the graph) based on conditional independence implied by data.

$$Cut(X, Y) \Leftrightarrow X \perp\!\!\!\perp Y \mid \mathbf{V}^*. \quad (67)$$

The crucial point where the PC algorithm refines the original SGS algorithm lies in narrowing down the conditional set \mathbf{V}^* (required by conditional independence test) via replacing it as the subset $Sub(\cdot)$ of the adjacent variables $Ad(\cdot)$ of X and Y .

$$Cut(X, Y) \Leftrightarrow X \perp\!\!\!\perp Y \mid Sub(Ad(X, Y)). \quad (68)$$

The certain (minimal) subset resulting into the conditional independence serves as the "approximation" of the d-separation set $D-Sep(X, Y)$ between X and Y . Aside from cutting off edges according to the independence, the algorithm need to orient the edges base on dependence. Technically, given any of the triple $Tr(\cdot)$ of X, Z, Y over \mathbf{V} , a "V-structure" $X \rightarrow Z \leftarrow Y$ can be oriented (denoted as $Ori(X, Z, Y)$) if:

$$\forall (X, Z, Y) \subseteq Tr(\mathbf{V}) \Rightarrow (Ori(X, Z, Y) \Leftrightarrow Z \notin D-Sep(X, Y)). \quad (69)$$

Leveraging the $D-Sep(X, Y)$ searched by adjacent variables, from a Bayesian perspective, equally indicates the swift search for the Markov blanket $Mb(X)$ (e.g. parents $Pa(X)$ of X):

$$X \perp\!\!\!\perp \mathbf{V} \setminus Ad(X) \mid Ad(X) \Rightarrow X \perp\!\!\!\perp \mathbf{V} \setminus Mb(X) \mid Mb(X). \quad (70)$$

NOTES⁵:

Here $Eq(\cdot)$ denotes an unity summarizing all the equivalent classes of the inducing graphs.

NOTES⁶:

In fact, given these equivalent classes, it is also important to denote the "non-collider" (Z) in π as $X - o - o Z - o - o Y$ after the additive procedure.

REFERENCE⁷:

Popular constraint-based causal discovery algorithms introduced by the book *Causation, Prediction, and Search*, Chapter 5 and 6.

- SGS: the Spites-Glymour-Scheines algorithm;
- PC: the Peter-Clark algorithm (Spites Peter, Glymour Clark);
- CI: the causal inference algorithm (Verma, Pearl);
- FCI: the fast causal inference algorithm (Glymour Clark)

In the presence of latent confounders, however, the equality no longer holds and we need additional conditional variables acting as the ***possible-D-separation set***. The point is, instead of searching for all possible candidates, we could think of it in the opposite way: excluding the impossible candidates — based on the *inducing path*. Notice that this is exactly the core idea behind the CI and FCI algorithm.

$$Cut(X, Y) \Leftrightarrow X \perp\!\!\!\perp Y \mid Possible-D-Sep(X, Y). \quad (71)$$

As for the orientation over the *partial inducing path graph*, we have slightly learned the sophisticated "connection symbols" shown in Equ(64). Quite similar to the theoretical concept⁸ of the *inducing path*, we also have an intermediate concept emerging amid the practical execution of algorithms named *definite discriminating path*.

Notice that a *definite discriminating path* for M ($Disc(X, Y \mid M)$) refers to an inducing path ($IP(X, Y)$) between X and Y containing the "additionally discriminative" variable M . In other words, we are not sure whether the variable M will satisfy the two characteristics (discussed in a few previous paragraphs) of the inducing path ($IP(X, Y)$). Given the variables U consisting of such an "extensive" inducing path, namely

$$\exists M \in V^*, U \subseteq Disc(X, Y \mid M). \quad (72)$$

it implies that, roughly speaking, orienting the *V-structure* should be confined to a relatively small range of the variables U when it comes to the existence of latent confounders. Analogically to Equ(69), orienting an *partial inducing path graph* implies that

$$\forall (P, M, R) \subseteq Tr(U) \Rightarrow (Ori(P, M, R) \Leftrightarrow M \notin D-Sep(X, Y)). \quad (73)$$

When all of the observed variables O are susceptible to latent confounders $L = V \setminus O$, than none of conditional independence relation will exist:

$$\forall X, Y (X \not\perp\!\!\!\perp Y \mid \{Z \mid Z \subseteq O^*\})_G. \quad (74)$$

In this case, as we said the "vertical perspective", difference is mathematically about the tiny shifting of causal variables based on investigators' interest and access. It could have formally named the *full latent variables model*. Particularly, it consists of two parts: (linear) *structure equations model* and *measurement model*. For the latter one, it conjures up the *measurement purity assumption* in the tutorial part. And the most universality one, in terminology, named *latent-measured purity*

$$\forall L_i, L_j \in L, O^{(i)} := O(L_i) \Rightarrow O^{(i)} \cap O^{(j)} = \emptyset. \quad (75)$$

With the assumptions of linearity and *purity* and without being in form of conditional independence constraints, *vanishing tetrad difference* or *tetrad constraints* is the constraints about a handful of equations characterizing correlation (as second-order statistic) among measured variables $O = \{X_1, X_2, X_3, X_4\}$.

NOTES⁸:

Notice that the inducing path is a theoretically established concept. The algorithm does not know what is called the inducing path; it utilizes the property of the inducing path by defining and searching for the "discriminating path".

Similarly, the algorithm does not know what is called the D-separation set yielded by the Causal Markov condition. Instead, the algorithm approaches the relation of conditional independence by searching the "variable adjacent set" (in order to approximate the d-separation set).

Margin Note:

In other words, causal relations of the structure model present only among the latent variables

Latent Variables Model

Measurement Model + Structure Model

Figure 48: This is a marginal figure.

Margin Note: The following pieces involve $\rho_{13}\rho_{24} = \rho_{14}\rho_{23} = \rho_{12}\rho_{34}$

$$\begin{cases} \rho_{13}\rho_{24} - \rho_{14}\rho_{23} = 0, \\ \rho_{14}\rho_{23} - \rho_{12}\rho_{23} = 0, \\ \rho_{12}\rho_{34} - \rho_{13}\rho_{24} = 0. \end{cases} \quad (76)$$

Structures on the right side are typical instances to illustrate how their entailing tetrad constraints lead to identification of the latent variables model.

$$\rho_{13}\rho_{24} = \frac{\sigma_{13}\sigma_{24}}{\sigma_1\sigma_2\sigma_3\sigma_4}, \quad \rho_{14}\rho_{23} = \frac{\sigma_{14}\sigma_{23}}{\sigma_1\sigma_2\sigma_3\sigma_4}. \quad (77)$$

The following pieces involve testing content. Video provides a powerful way to help you prove your point.

$$\sigma_{13} = E[(x_1 - E(x_1))(x_3 - E(x_3))]. \quad (78)$$

and on the basis of linearity and purity we expend the generation of each measurement variable in the same pattern

$$\sigma_{13} = E\{[(\lambda_x\xi + \varepsilon_x) - E(\lambda_x\xi + \varepsilon_x)][(\lambda_y\xi + \varepsilon_y) - E(\lambda_y\xi + \varepsilon_y)]\}. \quad (79)$$

We uncover that constraint among measurement variables fundamentally equals to a fully mapping from the latent variable.

$$\sigma_{13} = \lambda_x\lambda_y\sigma_\xi^2 + (\sigma_{\varepsilon_x\varepsilon_y} + \sigma_{\xi\varepsilon_x} + \sigma_{\xi\varepsilon_y}) = \lambda_x\lambda_y\sigma_\xi^2. \quad (80)$$

The following pieces involve testing content. Video provides a powerful way to help you prove your point.

$$\lambda_x\lambda_y\sigma_\xi^2 \neq \lambda_x\lambda_y\sigma_{\xi_x\xi_y}. \quad (81)$$

Mathematical deduction implies that statistic about measurement variables ultimately boils down to statistic about latent variables, hence enabling **detection** of latent variables. Properties for linearity and purity of the generation model is the essence of deduction. The following pieces involve testing content. Video provides a powerful way to help you prove your point.

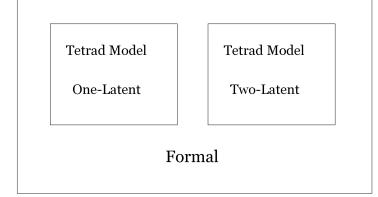


Figure 49: This is a marginal figure.

Margin Note: For algorithms to effectively search *vanishing tetrad difference*, elaborate concept such as *trek* and *choke point* are defined to establish *tetrad representation theorem*. Readers who feel interested in more details could refer to *Causation, Prediction, and Search*.

Causal Books References

Causation, Prediction, and Search: The following pieces involve testing content. Video provides a powerful way to help you prove your point.

- **5.4 New Algorithms:** The following pieces involve testing content. Video provides a powerful way to help you prove your point.
- **5.4 New Algorithms:** The following pieces involve testing content. Video provides a powerful way to help you prove your point.