

Causal inference cheat sheet

Author: Juvid Aryaman

Last compiled: January 10, 2021

This document is a summary of useful/interesting results in causal inference, mainly from [Pearl \(2009\)](#).

1. Basic probability

- Law of total probability: $P(A) = \sum_i P(A, B_i)$ (a.k.a. marginalizing over B)
- Chain rule of probability: $P(A, B) = P(A|B)P(B)$
- Thus, $P(A) = \sum_i P(A|B_i)P(B_i)$
- Expectation: $E(g(X)) = \sum_x g(x)P(x)$
- Conditional mean: $E(X|Y) = \sum_x xP(x|y)$
- Variance: $\sigma_X^2 = E[(X - E(X))^2]$
- Covariance: $\sigma_{XY} = E[(X - E(X))(Y - E(Y))]$
- Correlation coefficient: $\rho_{XY} = \sigma_{XY}/(\sigma_X\sigma_Y)$
- Regression coefficient of Y on X : $r_{YX} = \rho_{XY}\sigma_Y/\sigma_X = \sigma_{XY}/(\sigma_X^2)$ (for the equation $Y = r_{YX}X + c + \mathcal{N}(0, \sigma^2)$)
- Conditional independence: $(X \perp\!\!\!\perp Y|Z) \iff P(x|y, z) = P(x|z)$
- Partial correlation $\rho_{XY \cdot Z}$: The correlation between residuals e_X and e_Y resulting from the linear regression of X with Z and Y with Z , respectively.

2. Bayesian networks

Let a **graph** G consist of a set of **vertices** (or **nodes**) V and a set of **edges** E that connect some pair of vertices. Each edge in a graph can be either directed, undirected, or bidirected. Bidirected edges will subsequently be used to denote unobserved common causes, or **confounders**. Let a **path** be a sequence of edges such that each edge starts with the vertex ending in the preceding edge. A path may go either along or against the arrows of a directed graph. Directed graphs may include cycles (e.g. $X \rightarrow Y$, $Y \rightarrow X$), which represent mutual causation or feedback processes, but not self-loops (e.g. $X \rightarrow X$).

The recursive decomposition of the joint distribution into parents which characterises Bayesian networks is

$$P(x_1, \dots, x_n) = \prod_i P(x_i | pa_i). \quad (2.1)$$

2.1. d -separation (blocking) in Bayesian networks

Let an arrow $pa_j \rightarrow x_j$ denote part of a directed acyclic graph (DAG) in which variables are represented by nodes, and arrows are drawn from each node of the parent set PA_j towards the child node X_j .

Definition 2.1. d -separation A path p is d -separated (or blocked) by a set of nodes Z if and only if

1. p contains a chain $i \rightarrow m \rightarrow j$ or a fork $i \leftarrow m \rightarrow j$ such that the middle node m is in Z , or
2. p contains a collider $i \rightarrow m \leftarrow j$ such that the middle node m is not in Z and such that no descendant of m is in Z

A set Z d -separates X from Y if and only if Z blocks **every** path from a node in X to a node in Y

Theorem 2.1. Probabilistic implications of d -separation Consequently, if X and Y are d -separated by Z in a DAG G , then $(X \perp\!\!\!\perp Y|Z)$ in every distribution compatible with G . Conversely, if X , Y , and Z are not d -separated by Z in a DAG G then X and Y are dependent conditional on Z in almost all distributions compatible with G (assuming no parameter fine-tuning).

3. Functional causal models

A functional causal model consists of a set of equations of the form

$$x_i = f_i(pa_i, u_i), \quad i = 1, \dots, n \quad (3.1)$$

where pa_i are the set of variables (parents) that directly determine the value of X_i (i.e. its direct causes) and U_i represents errors (or “disturbances”) due to omitted factors. Eq.(3.1) is called a causal model if each equation represents the process by which the *value* (not merely the probability) of variable X_i is selected.

When some disturbances U_i are judged to be dependent, it is customary to denote such dependencies in a causal graph with double-headed arrows. If the causal diagram is acyclic, then the corresponding model is called *semi-Markovian* and the values of the variables X are uniquely determined by those of the variables U . If the error terms U are jointly independent, the model is called *Markovian*.

Linear structural equation models obey

$$x_i = \sum_{k \neq i} \alpha_{ik} x_k + u_i, \quad i = 1, \dots, n \quad (3.2)$$

In linear models, pa_i corresponds to variables on the r.h.s. of the above equation where $\alpha_{ik} \neq 0$.

3.1. Counterfactuals in functional causal models: An example

Consider a randomized clinical trial, where patients are/are not treated $X \in \{0, 1\}$. We also observe whether the patients die after treatment $Y \in \{0, 1\}$. We wish to ask the question: did the patient die *because of* the treatment, *despite* the treatment, or *regardless* of the treatment.

Assume $P(y|x) = 0.5$, and therefore $P(y, x) = 0.25$ for all x and y . We can write two models with the same joint distribution

Model 1 (treatment no effect):

$$x = u_1 \quad (3.3)$$

$$y = u_2 \quad (3.4)$$

$$P(u_1 = 1) = P(u_2 = 1) = \frac{1}{2} \quad (3.5)$$

Model 2 (treatment has an effect):

$$x = u_1 \quad (3.6)$$

$$y = xu_2 + (1 - x)(1 - u_2) \quad (3.7)$$

$$P(u_1 = 1) = P(u_2 = 1) = \frac{1}{2} \quad (3.8)$$

Let Q =fraction of deceased subjects from the treatment group who would not have died had they not taken the treatment. In model 1, $Q = 0$ since X has no effect on Y . In model 2, subjects who died ($y = 1$) and were treated ($x = 1$) must correspond to $u_2 = 1$. If $u_2 = 1$ then the only way for $y = 0$ is for $x = 0$. I.e. if you are a patient for whom $u_2 = 1$ then the only way not to die is to not take the treatment, so the treatment caused your death. So $Q = 1$.

Consequence 0: joint probability distributions are insufficient for counterfactual computation

Consequence 1: stochastic causal models are insufficient for counterfactual computation

Consequence 2: functional causal models are sufficient to define and compute counterfactual statements.

3.2. General method to compute counterfactuals

Given evidence $e = \{X_{obs}, Y_{obs}\}$, to compute probability of $Y = y$ under hypothetical condition $X = x$ apply the following steps:

1. Abduction: Update the probability of disturbances $P(u)$ to obtain $P(u|e)$
2. Action: Replace the equations corresponding to variables in the set X by the equations $X = x$
3. Prediction: Use the modified model to compute the probability $Y = y$.

4. Causal Bayesian networks

Definition 4.1. Causal effect Given two disjoint sets of variables X and Y , the **causal effect** of X on Y , denoted as $P(y|\hat{x})$ or $P(y|do(x))$, is the probability of $Y = y$ by deleting all equations from Eq.(3.1) where variables X are on the l.h.s., and substituting $X = x$ in the remaining equations. This corresponds to mutilating the DAG such that all arrows pointing directly to X are removed. **Amputation is the difference between seeing and doing.**

For an atomic intervention, we get the *truncated factorization* formula

$$P(x_1, \dots, x_n | \hat{x}'_i) = \begin{cases} \prod_{j \neq i} P(x_j | pa_j) & \text{if } x_i = x'_i \\ 0 & \text{if } x_i \neq x'_i \end{cases} \quad (4.1)$$

The $j \neq i$ denotes the removal of the term $P(x_i | pa_i)$ from Eq.(2.1) (i.e. amputation). A $do(x_i)$ is a severely limited sub-space of the full joint distribution, since the distribution only has support where the intervention variable x_i is equal to its particular intervention value x'_i , rather than a continuum of values in Eq.(2.1).

Multiplying and dividing by $P(x'_i | pa_i)$ yields

$$P(x_1, \dots, x_n | \hat{x}'_i) = \begin{cases} P(x_1, \dots, x_n | x'_i, pa_i) P(pa_i) & \text{if } x_i = x'_i \\ 0 & \text{if } x_i \neq x'_i \end{cases} \quad (4.2)$$

Marginalization of the above leads to the following theorem.

Adjustment for direct causes Let PA_i denote the set of direct causes of variable X_i , and let Y be any set of variables disjoint of $\{X_i \cup PA_i\}$. The causal effect of $do(X_i = x'_i)$ on Y is

$$P(y | \hat{x}'_i) = \sum_{pa_i} P(y | x'_i, pa_i) P(pa_i) \quad (4.3)$$

where $P(y | x'_i, pa_i)$ and $P(pa_i)$ are preintervention probabilities. This is called “adjusting for PA_i ”.

Identifiability Causal quantities are defined relative to a causal model M , not the joint distribution $P_M(v)$ over the set of observed variables V . Non-experimental data provides information about $P_M(v)$ alone, and several graphs can give rise to the same $P_M(v)$. Thus, not all quantities are unambiguously **identifiable** from observational data, **even with infinite samples**. Added assumptions by specifying a particular M can provide enough details to compute quantities of interest without explicating M in full.

Theorem 3.2.5: Given a causal diagram G of any Markovian model in which a subset of variables V are measured, the causal effect $P(y|\hat{x})$ is identifiable whenever $\{X \cup Y \cup PA_X\} \subseteq V$. I.e. *all parents of the cause are necessary to estimate the causal effect.*

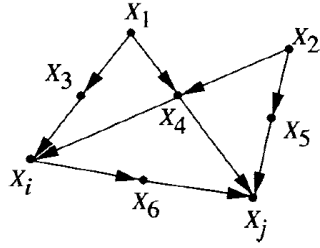


Figure 3.4 A diagram representing the back-door criterion; adjusting for variables $\{X_3, X_4\}$ (or $\{X_4, X_5\}$) yields a consistent estimate of $P(x_j | \hat{x}_i)$. Adjusting for $\{X_4\}$ or $\{X_6\}$ would yield a biased estimate.

Figure 1. Example of the back-door criterion

5. Inferring causal structure

- IC algorithm is for inferring causal structure given observational data when there are no latent variables
- IC* algorithm is for inferring causal structure given observational data when there are latent variables. The PC algorithm is apparently more contemporary (see [Spirtes et al. \(2010\)](#))
- There are local criteria for potential cause and genuine cause
- Spurious association: X and Y are spuriously associated if they are dependent in some context and there exists a latent common cause, as exemplified in the structure $Z_1 \rightarrow X \rightarrow Y \leftarrow Z_2$
- NOTEARS ([Zheng et al., 2018](#)) casts the structure learning problem as a continuous optimization problem over real matrices to avoid the superexponential combinatorial explosion with number of variables.

6. Adjusting for confounding bias

When seeking to evaluate the effect of one factor (X) on another (Y), we should ask **whether** we should *adjust* for possible variations in other factors (Z , known as “covariates”, “concomitants” or “confounders”). This becomes apparent in **Simpson’s paradox**: any statistical relationship between two variables may be reversed by including additional factors in the analysis.

6.1. The back-door criterion

This criterion demonstrates how confounders that *affect* the treatment variable can be used to facilitate causal inference.

Definition 6.1. Back-door criterion A set of variables Z satisfy the back-door criterion relative to an ordered pair of variables (X_i, X_j) in a DAG G if:

1. no node in Z is a descendant of X_i ; and
2. Z blocks every path between X_i and X_j that contains an arrow into X_i

Similarly, if X and Y are two disjoint subsets of nodes in G , then Z satisfies the back-door criterion relative to (X, Y) if it satisfies the criterion relative to any pair (X_i, X_j) such that $X_i \in X$ and $X_j \in Y$.

Theorem 6.1. Back-door adjustment If a set of variables Z satisfies the back-door criterion relative to (X, Y) , then the causal effect of X on Y is identifiable and is given by

$$P(y|\hat{x}) = \sum_z P(y|x, z)P(z). \quad (6.1)$$

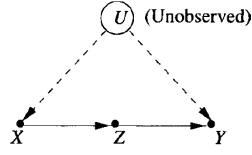


Figure 3.5 A diagram representing the front-door criterion. A two-step adjustment for Z yields a consistent estimate of $P(y | \hat{x})$.

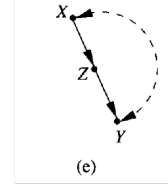


Figure 2. (Left) Example of the front-door criterion. The path $X \leftarrow U \rightarrow Y$ denotes an unobserved (latent) unobserved common cause. (Right) This is often represented as a **bi-directed path**.

This corresponds to partitioning the population into groups that are homogeneous relative to Z , assessing the effect of X on Y in each homogeneous group, and then averaging the results. Conditioning in this way means that the observation $X = x$ cannot be distinguished from an intervention $do(x)$.

6.2. The front-door criterion

This criterion demonstrates how confounders that are *affected by* the treatment variable can be used to facilitate causal inference.

Definition 6.2. Front-door A set of variables Z satisfy the front-door criterion relative to an ordered pair of variables (X, Y) if:

1. Z intercepts all directed paths from X to Y ;
2. there is no unblocked back-door path from X to Z ; and
3. all back-door paths from Z to Y are blocked by X .

Theorem 6.2. Front-door adjustment If Z satisfies the front-door criterion relative to (X, Y) and if $P(x, z) > 0$, then the causal effect of X on Y is identifiable and is given by

$$P(y|\hat{x}) = \sum_z P(z|x) \sum_{x'} P(y|x', z) P(x') \quad (6.2)$$

Conditions (2) and (3) of the front-door definition are overly restrictive: e.g. nested combinations of back-door and front-door conditions are permissible: see Section 7 for a more general set of conditions.

7. Do-calculus

The back-door and front-door criteria do not provide a complete set of rules for when/how causal effects can be computed. Do-calculus sidesteps the need for algebraic manipulation and provides a complete set of inference rules by which probabilistic sentences involving interventions and observations can be transformed into other such sentences, allowing a method of deriving/verifying claims about interventions. The aim is to compute causal effect expressions for $P(y|\hat{x})$ where Y and X are subsets of variables. When $P(y|\hat{x})$ can be reduced to an expression involving observable probabilistic quantities, we say that the causal effect of X on Y is **identifiable**.

7.1. Notation

- $G_{\overline{X}}$ = graph obtained by deleting from G all arrows pointing into nodes in X
- $G_{\underline{X}}$ = graph obtained by deleting from G all arrows pointing out of nodes in X
- $G_{\overline{X}\underline{Z}}$ = graph obtained by deleting from G all arrows pointing into nodes in X and out of nodes in Z
- $P(y|\hat{x}, z) := P(y, z|\hat{x})/P(z|\hat{x})$, meaning the probability of observing $Y = y$ given an *intervention* $X = x$ and an *observation* $Z = z$

7.2. Rules

Theorem 7.1. Rule 1 *Insertion/deletion of observations:*

$$P(y|\hat{x}, z, w) = P(y|\hat{x}, w) \quad \text{if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{X}}}. \quad (7.1)$$

Theorem 7.1 is a reaffirmation of d -separation (Section 2.1) as a valid test for conditional independence in the distribution resulting from $do(X = x)$. The rule follows from the fact that deleting equations from the system $(G_{\overline{X}})$ does not introduce any dependencies among the remaining disturbance terms.

Theorem 7.2. Rule 2 *Action/observation exchange:*

$$P(y|\hat{x}, \hat{z}, w) = P(y|\hat{x}, z, w) \quad \text{if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{X}\underline{Z}}}. \quad (7.2)$$

Theorem /refthm:do-calc-act-obs-ex provides a condition for an external intervention $do(Z = z)$ to have the same effect on Y as the passive observation $Z = z$. The condition amounts to $\{X \cup W\}$ blocking all back-door paths from Z to Y (in $G_{\overline{X}}$), since $G_{\overline{X}\underline{Z}}$ retains all (and only) such paths.

Theorem 7.3. Rule 3 *Insertion/deletion of actions:*

$$P(y|\hat{x}, \hat{z}, w) = P(y|\hat{x}, w) \quad \text{if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{X}, \underline{Z(W)}}} \quad (7.3)$$

where $Z(W)$ is the set of Z -nodes that are not ancestors of any W -node in $G_{\overline{X}}$.

Theorem 7.3 provides conditions for introducing (or deleting) an external intervention $do(Z = z)$ without affecting the probability of $Y = y$. The validity of this rule stems from simulating the intervention $do(Z = z)$ by the deletion of all equations corresponding to the variables in Z (hence $G_{\overline{X}\underline{Z}}$).

Completeness A quantity $Q = P(y|do(x), z)$ is identifiable if and only if it can be reduced to a do -free expression using the above 3 rules.

7.3. Identifiability

A causal effect $q = P(y_1, \dots, y_k|\hat{x}_1, \dots, \hat{x}_m)$ is identifiable in a model characterised by a graph G if there exists a finite sequence of transformations conforming to one of the three rules in Section 7.2 that reduces q into a standard (i.e. “hat”-free) probability expression involving observed quantities. Figure 3 provides a set of graphical conditions; if any one is satisfied then $P(y|\hat{x})$ is identifiable, and satisfying at least one of the conditions is necessary for $P(y|\hat{x})$ to be identifiable. I.e. $P(y|\hat{x})$ is unidentifiable then no finite sequence of inference rules reduces $P(y|\hat{x})$ to a hat-free expression. Figure 3 can also be used to define an algorithm for deriving a closed-form expression for control queries in terms of observable quantities, an implementation in R is in the package `causaleffect`, see [Tikka and Karvanen \(2017\)](#) and a Jupyter Notebook example [here](#).

Assorted facts on identifiability

- Whilst a causal effect is not identifiable for every joint distribution of variables if this condition is broken, it might be for *some* probability densities. For example, an instrumental variable can yield a causal effect identifiable in a linear model in the presence of a bow pattern (Fig. 3.7A of Causality), but will not be generally identifiable (see Section 3.5 of Causality).
- If $P(y|\hat{x})$ is identifiable, then if a set of nodes Z lies on a directed path from X to Y , then $P(z|\hat{x})$ is also identifiable (lemma 4.3.4).

1. *There is no back-door path from X to Y in G ; that is, $(X \perp\!\!\!\perp Y)_{G_{\bar{X}}}$.*
2. *There is no directed path from X to Y in G .*
3. *There exists a set of nodes B that blocks all back-door paths from X to Y so that $P(b|\hat{x})$ is identifiable. (A special case of this condition occurs when B consists entirely of nondescendants of X , in which case $P(b|\hat{x})$ reduces immediately to $P(b)$.)*
4. *There exist sets of nodes Z_1 and Z_2 such that:*
 - (i) *Z_1 blocks every directed path from X to Y (i.e., $(Y \perp\!\!\!\perp X | Z_1)_{G_{\bar{Z}_1\bar{X}}}$);*
 - (ii) *Z_2 blocks all back-door paths between Z_1 and Y (i.e., $(Y \perp\!\!\!\perp Z_1 | Z_2)_{G_{\bar{X}\bar{Z}_1}}$);*
 - (iii) *Z_2 blocks all back-door paths between X and Z_1 (i.e., $(X \perp\!\!\!\perp Z_1 | Z_2)_{G_{\bar{X}}}$);*
and
 - (iv) *Z_2 does not activate any back-door paths from X to Y (i.e., $(X \perp\!\!\!\perp Y | Z_1, Z_2)_{G_{\bar{Z}_1\bar{X}(Z_2)}}$). (This condition holds if (i)–(iii) are met and no member of Z_2 is a descendant of X .)*

(A special case of condition 4 occurs when $Z_2 = \emptyset$ and there is no back-door path from X to Z_1 or from Z_1 to Y .)

Figure 3. Graphical conditions for identification of causal effect (Theorem 4.3.1 Causality). Satisfying at least one renders the causal effect $P(y|\hat{x})$ identifiable, whereas satisfying none implies unidentifiability of the causal effect.

- **Complete identifiability condition** A sufficient condition for identifying the causal effect $P(y|do(x))$ is that there exists no bi-directed path (i.e. a path composed entirely of bi-directed arcs, see Fig. 2) between X and any of its children. Prior to applying this criterion, all nodes which are not ancestors of Y are deleted from the graph (i.e. only consider nodes which are on pathways from X to Y).

8. Actions, plans, and direct effects

Pearl defines two kinds of intervention:

- **Act:** An intervention which results from a reactive policy, deriving from an agent's beliefs, disposition, and environmental inputs (or the "outside")
- **Action:** An intervention which results from a deliberative policy, deriving from an agent's free will (or the "inside"; meditative traditions might not draw such a bright line between these two classifications as a description of physical reality, but it is no doubt a useful distinction for reasoning about the future when conscious agents are involved)

8.1. Conditional actions and stochastic policies

In general, interventions may involve complex policies in which X is made to respond according to e.g. a deterministic functional relationship $x = g(z)$, or more generally through a stochastic relationship whereby X is set to x with probability $P^*(x|z)$.

Let $P(y|do(X = g(z)))$ denote the distribution of Y prevailing under the deterministic policy $do(x =$

$g(z)$). Then,

$$\begin{aligned} P(y|do(X = g(z))) &= \sum_z P(y|do(X = g(z)), z)P(z|do(X = g(z))) \\ &= \sum_z P(y|\hat{x}, z)|_{x=g(z)}P(z) \\ &= E_z[P(y|\hat{x}, z)|_{x=g(z)}]. \end{aligned} \quad (8.1)$$

Hence, the evaluation of the outcome of an intervention under a complicated conditional policy $x = g(z)$ amounts to being able to evaluate $P(y|\hat{x}, z)$. The equality $P(z|do(X = g(z))) = P(z)$ stems from the fact that Z **cannot** be a descendant of X : in other words, **one cannot define a coherent policy of action for X based on an (indirect) effect of X because actions change the distributions of their effects!** (Aside: I suppose one might argue about whether an agent has any choice over the form of $g(z)$)

Similarly, let $P(y)|_{P^*(x|z)}$ denote the distribution of Y prevailing under the stochastic policy $P^*(x|z)$ – i.e. given $Z = z$, $do(X = x)$ occurs with probability $P^*(x|z)$. Then,

$$P(y)|_{P^*(x|z)} = \sum_x \sum_z P(y|\hat{x}, z)P^*(x|z)P(z). \quad (8.2)$$

Since $P^*(x|z)$ is specified externally, it is again the case that $P(y|\hat{x}, z)$ is sufficient for the identifiability of any stochastic policy which shapes the distribution of X by the outcome of Z .

8.2. Identification of dynamic plans

A **control problem** consists of a DAG with vertex set V partitioned into four disjoint sets $V = \{X, Z, U, Y\}$ where

- X = the set of control variables (exposures, interventions, treatments, etc.)
- Z = the set of observed variables, often called **covariates**
- U = the set of unobserved (latent) variables, and
- Y = an outcome variable

We are interested in settings where we have gathered data $\mathcal{D} = \{X, Z, Y\}$ for previous agents making actions X . The problem is, given a new instance of the system (e.g. a new patient whom we seek to treat), can we estimate the outcome of $\{do(x_1), \dots, do(x_n)\}$ using only the observational data \mathcal{D} . See Section 4.4.1 of Causality for a specific motivating example.

Let control variables be ordered $X = X_1, \dots, X_n$ such that every X_k is a non-descendant of X_{k+j} ($j > 0$) and let the outcome Y be a descendant of X_n . A **plan** is an ordered sequence $(\hat{x}_1, \dots, \hat{x}_n)$ of value assignments to the control variables. A **conditional plan** is an ordered sequence $(\hat{g}_1(z_1), \dots, \hat{g}_n(z_n))$ where $\hat{g}_k(z_k)$ means “set X_k to $\hat{g}_k(z_k)$ whenever $Z_k = z_k$ ”, where the support Z_k of each $g_k(z_k)$ must not contain any variables that descendants of X_k .

Theorem 8.1. Plan identification: the sequential back-door criterion. *The probability of the unconditional plan $P(y|\hat{x}_1, \dots, \hat{x}_n)$ is identifiable if, for every $1 \leq k \leq n$ there exists a set Z_k of covariates satisfying the following conditions:*

$$Z_k \subseteq N_k \quad (8.3)$$

where N_k is the set of observed nodes that are non-descendants of any element of $\{X_k, X_{k+1}, \dots, X_n\}$, and

$$(Y \perp\!\!\!\perp X_k | X_1, \dots, X_{k-1}, Z_1, \dots, Z_k)_{G_{\underline{X}_k, \overline{X}_{k+1}, \dots, \overline{X}_n}} \quad (8.4)$$

When these conditions are satisfied, the effect of the plan is given by

$$P(y|\hat{x}_1, \dots, \hat{x}_n) = \sum_{z_1, \dots, z_n} P(y|z_1, \dots, z_n, x_1, \dots, x_n) \times \prod_{k=1}^n P(z_k|z_1, \dots, z_{k-1}, x_1, \dots, x_{k-1}) \quad (8.5)$$

8.3. Direct and indirect effects

We are often concerned with the extent to which a variable affects another directly, rather than the total causal effect mediated through all other intervening variables. For example, in cases of sex discrimination, we may be interested in asking the direct effect of an applicant's sex on the outcome of an applicant's job application. In effect, we are concerned with the causal effect of variable X on Y while all other factors in the analysis are held fixed (*Ceteris paribus*).

Definition 8.1. Direct effect. *The direct effect of X on Y is given by $P(y|\hat{x}, \hat{s}_{XY})$ where \hat{s}_{XY} is the set of all variables in the model except X and Y*

Corollary 8.1. *The direct effect of X on Y is given by $P(y|\hat{x}, \hat{p}_{a_{Y \setminus X}})$ where $p_{a_{Y \setminus X}}$ is any realization of the parents of Y excluding X .*

It is sometimes meaningful to average the direct effect over all levels of $p_{a_{Y \setminus X}}$. To do this, we define the natural direct effect:

Definition 8.2. Natural direct effect. *The natural direct effect ($DE_{x,x'}(Y)$) is defined as*

$$DE_{x,x'}(Y) := E[Y(x', Z(x)) - E(Y(x))] \quad (8.6)$$

where $Z = p_{a_{Y \setminus X}}$, and $Y(x', Z(x))$ is the value that Y would attain under the counterfactual scenario of $X = x'$, but Z retaining the values under the setting $X = x$.

The natural direct effect involves probabilities of nested counterfactuals, and cannot generally be written in terms of the $do(x)$ operator. However, if certain assumptions of "no confounding" are deemed valid, the natural direct effect can be reduced to

$$DE_{x,x'}(Y) = \sum_z [E(Y|do(x', z)) - E(Y|do(x, z))]P(z|do(x)) \quad (8.7)$$

which is simply a weighted average of controlled direct effects.

We can also define the indirect effect which quantifies the influence of X on Y through all paths except for the direct path from $X \rightarrow Y$.

Definition 8.3. Indirect effect. *The natural indirect effect ($IE_{x,x'}(Y)$) is defined as*

$$IE_{x,x'}(Y) = E[Y(x, Z(x')) - E(Y(x))] \quad (8.8)$$

We can define

Definition 8.4. Total effect. *The total effect of X on Y is given by $P(y|do(x))$, namely, the distribution of Y while X is held constant at x and all other variables are permitted to run their natural course. Confusingly, we also sometimes denote the total effect $TE_{x,x'}(Y)$ as*

$$TE_{x,x'}(Y) := E[Y(x') - E(Y(x))] \quad (8.9)$$

TODO: Write $TE_{x,x'}(Y)$ in terms of $P(y|do(x))$?

Theorem 8.2. Relationship between total effect, direct effect, and indirect effect *The total effect of a transition is the difference between the direct effect of that transition and the indirect effect of the reverse transition*

$$TE_{x,x'}(Y) = DE_{x,x'}(Y) - IE_{x',x}(Y) \quad (8.10)$$

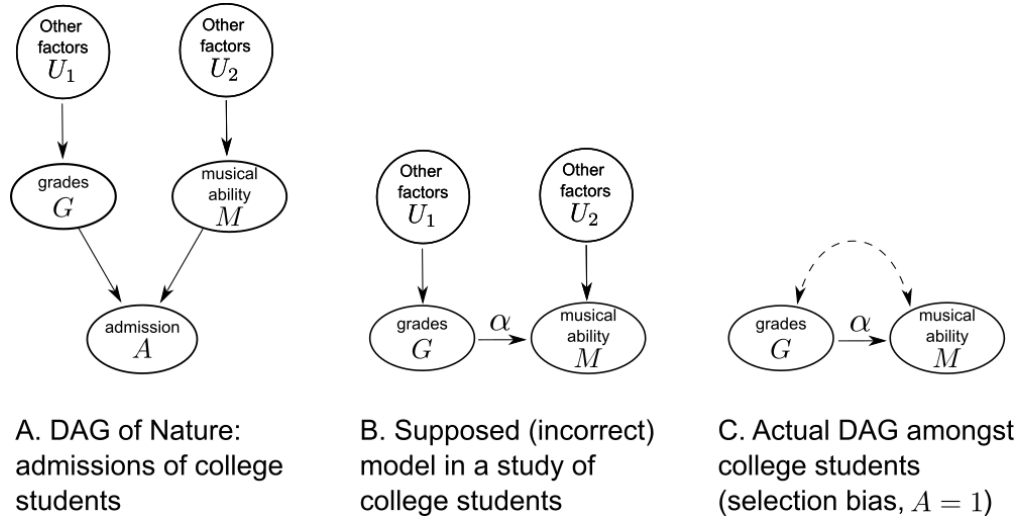


Figure 4. An example of selection bias in a study. A) Suppose, for a particular college, both grades and musical ability affect admission rate. B) Suppose investigators seek to understand the relationship between grades and musical ability, and measure G and M amongst students of the college in an attempt to estimate α ($\alpha = 0$ in reality). The investigators assume the DAG in (B) and discover a strong negative relationship between G and M . C) However, the investigators would discover that the equality in Eq.(9.3) does not hold, and therefore U_1 and U_2 are correlated (owing to a selection bias from $A = 1$, see Berkson's paradox). Hence the investigators must correct their DAG to include a latent common cause between G and M . In this case, α becomes unidentifiable since all of the causal effect can (and, here, is!) attributable to the latent common cause.

9. Causality and structural models

Let's rewrite Eq.(3.1) as

$$x_i = f_i(pa_i, \epsilon_i), \quad i = 1, \dots, n. \quad (9.1)$$

In general, for the **partial correlation** $\rho_{XY \cdot Z}$,

$$(X \perp\!\!\!\perp Y | Z) \implies \rho_{XY \cdot Z} = 0 \quad (9.2)$$

and therefore, in **any** Markovian model with DAG G , the partial correlation $\rho_{XY \cdot Z}$ vanishes whenever the nodes corresponding to the variables in Z d -separate node X from node Y in G , regardless of model parameters. Moreover, no other partial correlation vanishes, for all model parameters. **[Q: Not sure if this is general or only for linear models]**

Theorem 9.1. Test for correlation between error terms For any two non-adjacent variables X and Y , where Y is **not** a parent of X , a sufficient test of whether ϵ_X and ϵ_Y are uncorrelated is if the following equality holds:

$$E[Y|x, do(S_{XY})] = E[Y|do(x), do(S_{XY})] \quad (9.3)$$

where S_{XY} stands for (any setting of) all variables in the model excluding X and Y . If ϵ_X and ϵ_Y are uncorrelated then we are justified in having the absence of a bidirected arc between X and Y . I.e. the omitted factors which directly affect X , ϵ_X , are independent of the omitted factors which directly affect Y , ϵ_Y .

Note that **selection bias** can arise when two uncorrelated factors have a common effect that is omitted from the analysis but influences the selection of samples for the study, see Fig. 4. Hence bidirected arcs should be assumed to exist, by default, between any two nodes in a diagram – since

they at worst compromise the identifiability of model parameters. They should be deleted only by well-motivated justifications, such as the unlikely existence of a common cause, and the unlikely existence of selection bias.

9.1. Exogeneity and instrumental variables

Definition 9.1. General Exogeneity Let X and Y be two sets of variables, and let λ be any quantity which may be computed from a structural model M (structural, statistical, etc.) in a theory T . We say that X is exogenous relative to (Y, λ, T) if λ is identifiable from the conditional distribution of $P(y|x)$, that is, if

$$P_{M_1}(y|x) = P_{M_2}(y|x) \implies \lambda(M_1) = \lambda(M_2) \quad (9.4)$$

for any two models M_1 and M_2 satisfying theory T .

Definition 9.2. Instrumental variable For the parameter $\lambda = P(w|do(z))$ where W and Z are two sets of variables in Y , then if X is exogenous relative to (Y, λ, T) then

$$P_{M_1}(z, w|x) = P_{M_2}(z, w|x) \implies P_{M_1}(w|do(z)) = P_{M_2}(w|do(z)). \quad (9.5)$$

for any two models M_1 and M_2 satisfying theory T . Under these conditions, we call X an instrumental variable for the causal effect $P(w|do(z))$

Corollary 9.1. If X is an instrumental variable for $P(w|do(z))$, then $P(w|do(z))$ is identifiable from $P(z, w|x)$. I.e. X renders a causal effect $P(w|do(z))$ identifiable which X itself does not directly participate in.

An alternative definition of exogeneity may also be given, in terms of the correlation between errors and variables:

Definition 9.3. Error-based exogeneity A variable X is exogenous relative to $\lambda = P(y|do(x))$ if X is independent of all errors ϵ that influence Y , except those mediated by X .

9.2. Linear structural equation models

Linear structural equation models (SEMs) obey

$$x_i = \sum_{k \neq i} \alpha_{ik} x_k + \epsilon_i, \quad i = 1, \dots, n \quad (9.6)$$

If, in Eq.(9.6)

$$\epsilon \sim \mathcal{N}(\mu, \Sigma) \quad (9.7)$$

then X_i will also be multivariate normal, and the SEM will be entirely determined by the set of correlation coefficients ρ_{ij} . For a linear SEM

$$\rho_{XY \cdot Z} = 0 \iff (X \perp\!\!\!\perp Y | Z). \quad (9.8)$$

9.2.1. Interpretation

Definition 9.4. Structural equations An equation

$$y = \beta x + \epsilon \quad (9.9)$$

is said to be structural if it is to be interpreted as follows: In an ideal experiment where we control X to x and any other set Z of variables (not containing X or Y) to z , the value of Y is given by $\beta x + \epsilon$, where ϵ is not a function of the settings x and z . **[I think: ϵ can be arbitrarily distributed, potentially correlated with X , but assumed to have $E[\epsilon] = 0$.]**

The equality sign in structural equations has a different behaviour to algebraic equality signs. In the context of observations, the equality sign in Eq.(9.9) behaves symmetrically between X and Y : e.g. observing $Y = 0$ implies $\beta x = -\epsilon$. In contrast, in the context of interventions, the equality sign in Eq.(9.9) behaves asymmetrically between X and Y : e.g. setting $Y = 0$ tells us nothing about the relationship between x and ϵ .

Furthermore, the strongest empirical claim made by Eq.(9.9) is that

$$P(y|do(x), do(z)) = P(y|do(x)) \quad (9.10)$$

i.e. the statistics of Y remain invariant to the manipulation of Z under the condition of $do(x)$. In contrast, regression equations make no empirical claims whatsoever.

The operational definition of the structural parameter β in Eq.(9.9) is

$$\beta = \frac{\partial}{\partial x} E[Y|do(x)] \quad (9.11)$$

(since $E[Y|do(x)] = \beta x$). In words, β is the rate of change in the expectation of Y in an experiment where X is held at x by external control. This is true regardless of the correlation between X and ϵ in non-experimental studies (e.g. via another equation $x = \alpha y + \delta$).

As a consequence of the above, the operational definition of the error term ϵ is

$$\epsilon = y - E[Y|do(x)] \quad (9.12)$$

(again since $E[Y|do(x)] = \beta x$).

9.2.2. Estimation

Define the conditional variance $\sigma_{X|z}^2$, conditional covariance $\sigma_{XY|z}^2$, and the conditional covariance $\rho_{XY|z}$. For multivariate normal, $\sigma_{X|z}^2$, $\sigma_{XY|z}^2$, and $\rho_{XY|z}$ are all independent of the value of z . For the MVN, the **partial** variance $\sigma_{X \cdot Z}^2$, covariance $\sigma_{XY \cdot Z}$, and correlation $\rho_{XY \cdot Z}$ all coincide with the conditional variance, covariance, and correlation respectively (although this is not generally the case).

A **partial regression coefficient**, $r_{YX \cdot Z}$ is given by

$$r_{YX \cdot Z} = \rho_{YX \cdot Z} \frac{\sigma_{Y \cdot Z}}{\sigma_{X \cdot Z}} \quad (9.13)$$

and is equal to the coefficient of X in the linear regression of Y on X and Z . So, the coefficient of x in the regression equation

$$y = \alpha x + b_1 z_1 + \dots + b_k z_k \quad (9.14)$$

is

$$\alpha = r_{YX \cdot Z_1 Z_2 \dots Z_k} \quad (9.15)$$

Theorem 9.2. d -Separation in General Linear Models For any linear model Eq.(9.6), which may include cycles and bidirected arcs (i.e. dependent ϵ between different variables), $\rho_{XY \cdot Z} = 0$ if Z d -separates X from Y , where bidirected arcs between i and j are interpreted as a latent common parent $i \leftarrow L \rightarrow j$.

Theorem 9.2 provides a method for finding models in the context of linear SEMs: by searching over all $\rho_{XY \cdot Z}$, we can construct a DAG. Not all partial correlations need to be searched.

Definition 9.5. Basis Let S be a set of partial correlations. A basis B for S is a set of zero partial correlations where (i) B implies the zero of every element of S and (ii) no proper subset of B sustains such an implication.

An obvious choice of basis for a DAG D is

$$B = \{\rho_{ij \cdot pa_i} | i > j\} \quad (9.16)$$

where i ranges over all nodes in D and j ranges over all predecessors of i in any order that agrees with the arrows of D . More economical choice of basis exist, such as a Graphical Basis.

Theorem 9.3. Markov linear-normal equivalence *Two Markovian linear-normal models are observationally indistinguishable if every covariance matrix generated by one model can be parametrically generated by the other (covariance equivalent). Two such models are covariance equivalent if and only if their corresponding graphs have the same sets of zero partial correlations. Moreover, two such models are covariance equivalent if and only if they have the same edges and the same sets of v -structures. (I.e. arrows can be reversed as long as they do not alter v -structures)*

9.2.3. Parameter identifiability

Consider an edge $X \rightarrow Y$ in graph G , and let α be the path coefficient associated with that edge (i.e. the strength of the direct causal effect of X on Y). The regression coefficient in a linear model can, in general, be decomposed into

$$r_{YX} = \alpha + I_{YX} \quad (9.17)$$

where I_{YX} is independent of α , since I_{YX} is composed of other indirect paths connecting X and Y . If we remove the edge $X \rightarrow Y$ and observe that the resulting subgraph entails zero correlation between X and Y then $I_{XY} = 0$ and $r_{YX} = \alpha$, and hence α is identified. This idea is extended in the following theorem

Theorem 9.4. Single-door Criterion for Direct Effects *Let G be any path diagram in which α is the path coefficient associated with link $X \rightarrow Y$, and let G_α denote the diagram that results when $X \rightarrow Y$ is deleted from G . The coefficient α is identifiable if there exists a set of variables Z such that (i) Z contains no descendant of Y and (ii) Z d -separates X from Y in G_α . If Z satisfies these two conditions then, in a linear SEM, $\alpha = r_{YX \cdot Z}$. Conversely, if Z does not satisfy these conditions, then $r_{YX \cdot Z}$ is not a consistent estimand of α .*

Theorem 9.5. Back-door Criterion for Total Effects *For any two variables in a causal diagram G , the total effect of X on Y is identifiable if there exists a set of measurements Z such that*

1. *no member of Z is a descendant of X ; and*
2. *Z d -separates X from Y in the subgraph $G_{\underline{X}}$ formed by deleting all G arrows emanating from X*

If the two conditions are satisfied, then the total effect of X on Y in a linear SEM is given by $r_{YX \cdot Z}$.

Theorems 9.4 and 9.5 are special cases of a more general scheme. In order to identify any **partial effect**, as defined by a select bundle of causal paths from X to Y , we must find a set Z of measured variables that block all non-selected paths between X and Y . For linear models, the partial effect is equal to the regression coefficient $r_{YX \cdot Z}$.

Some direct effects require evaluation of a broader causal effect first, in order to extract the direct effect of interest (see Fig. 5). The parameter α cannot be directly estimated with Theorem 9.4 because of the confounder, or its constituents (since it has none). Instead, we may apply Theorem 9.5 twice like so:

$$P(Y|\hat{z}) = r_{YZ} = \alpha\beta \quad (9.18)$$

$$P(X|\hat{x}) = r_{YX} = \beta \quad (9.19)$$

$$\alpha = E(Y|\hat{x}) = \frac{r_{YZ}}{r_{YX}} = \frac{P(Y|\hat{z})}{P(X|\hat{z})} \quad (9.20)$$

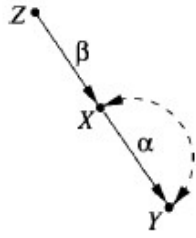


Figure 5.9 Graphical identification of α using instrumental variable Z .

Figure 5. Example of inference of a direct effect by evaluating a broader causal effect, and extracting the effect of interest using an instrumental variable.

A

	Combined	E	$\neg E$	Recovery Rate
(a)	Drug (C)	20	20	40
	No drug ($\neg C$)	16	24	40
		36	44	80
	Males	E	$\neg E$	Recovery Rate
(b)	Drug (C)	18	12	30
	No drug ($\neg C$)	7	3	10
		25	15	40
	Females	E	$\neg E$	Recovery Rate
(c)	Drug (C)	2	8	10
	No drug ($\neg C$)	9	21	30
		11	29	40

Figure 6.1 Recovery rates under treatment (C) and control ($\neg C$) for males, females, and combined.

B

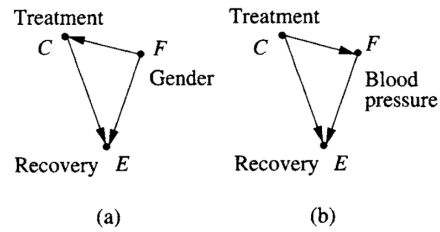


Figure 6. Example of Simpson's paradox. A. The recovery rate in the pooled population is greater with the drug, but for both the male and female populations the recovery rate is greater without the drug. B. Two causal models which are observationally equivalent capable of generating the data in (A). In (A), $P(E|do(C)) = \sum_F P(E|C, F)P(F)$ (Theorem 6.1) which averages over sub-populations, whereas in (B), $P(E|do(C)) = P(E|C)$ (Theorem 7.2) which uses the pooled data.

[TODO PROVE:] In a linear SEM, to evaluate the partial effect of X on Y along a single path, take the product of all path coefficients for each edge for all observed variables in the model. To evaluate the total effect for multiple branches, sum the partial effects across each branch (consisting of endogenous edges).

10. Simpson's paradox

Simpson's paradox is a reversal effect observed in sub-populations.

Definition 10.1. Simpson's paradox *The phenomenon whereby an event C increases the probability of E in a given population p and, at the same time, decreases the probability of E in every sub-population of p . In other words, if F and $\neg F$ are two complementary properties describing two sub-populations, then it is possible to encounter the equalities:*

$$\begin{aligned}
 P(E|C) &> P(E|\neg C) \\
 P(E|C, F) &< P(E|\neg C, F) \\
 P(E|C, \neg F) &< P(E|\neg C, \neg F).
 \end{aligned}
 \tag{10.1}$$

Note: Simpson's paradox is not strictly a paradox because it does not involve any contradiction.

An example dataset displaying Simpson's paradox is in Fig. 6 for a drug treatment (C), recovery (E), for sub-populations of gender (F). The question, therefore, is: "what is the total effect of the drug on recovery?". The solution to this question depends on the causal assumptions we bring to bear on the

problem: two different causal models which are capable of generating the data yield different answers for the quantity $P(E|do(C))$, one of which averages over the subpopulations, and the other using pooled data.

Theorem 10.1. Sure-thing Principle *An action C that increases the probability of an event E in each sub-population F must also increase the probability of E in the population as a whole, provided that the action does not change the distribution of the sub-populations. In other words, for dichotomous sub-populations F , if*

$$\begin{aligned} P(E|do(C), F) &< P(E|do(\neg C), F) \\ P(E|do(C), \neg F) &< P(E|do(\neg C), \neg F) \end{aligned} \tag{10.2}$$

then

$$P(E|do(C)) < P(E|do(\neg C)) \tag{10.3}$$

Theorem 10.1 follows from do-calculus, and gives the intuitive result that if the drug in Fig. 6A harms both men and women then one should not administer the drug if one does not know the patient's gender, despite the fact that the observational conditional densities obey $P(E|C) > P(E|\neg C)$. The "paradox" in the example arises because the males, who recover (regardless of the drug) more often than the females, are also more likely than the females to use the drug.

References

- Pearl, J., 2009 *Causality*. Cambridge university press.
- Spirtes, P., C. Glymour, R. Scheines, and R. Tillman, 2010 Automated search for causal relations: Theory and practice .
- Tikka, S. and J. Karvanen, 2017 Identifying causal effects with the R package causaleffect. J. Stat. Softw .
- Zheng, X., B. Aragam, P. K. Ravikumar, and E. P. Xing, 2018 Dags with no tears: Continuous optimization for structure learning. Advances in Neural Information Processing Systems **31**: 9472–9483.