

Using Causal Graphs In Epidemiological Research

Dr David A. Stephens

Department of Mathematics & Statistics

david.stephens@mcgill.ca

www.math.mcgill.ca/dstephens/SISCER2025/



1. Probability distributions for observed data.
2. Independence and conditional independence.
3. Graphical modelling.
4. Cause and effect: structural models.
5. Paths and how to block them.
6. d-separation.

Key Challenge in Epidemiological Research

Key objective: causal conclusions from observational data

- Experimental studies:
 - Treatment assigned by the researcher, independent of confounding factors;
 - Causal statements possible.
- Observational studies:
 - Treatment assignment dependent on confounding factors;
 - Causal statements not possible ?

The objective of *causal inference* is to quantify the effect of an *intervention*: in a multi-variable system

- suppose we are able to *manipulate* (i.e. alter the value of) one of the variables separately from all other variables;
- we wish to report the impact of that manipulation on one or more of the other variables.

In many scientific enterprises, this is a primary objective.

We will collect data

$$\{(x_i, y_i, z_i), i = 1, \dots, n\}$$

which are observed values of the variables X , Y and Z .

- X - predictors, covariates, *confounders*
- Y - *outcome*, response
- Z - *treatment*, exposure

Causal Effects and Counterfactual Outcomes

The *causal effect* of Z on Y is the amount to which an *intervention* to change Z from z_0 to z_1 modifies Y .

- The *potential* or *counterfactual* outcome is denoted $Y(\mathbf{z})$, and is the outcome after an *intervention* to set $Z = \mathbf{z}$.

For example, if $Z \in \{0, 1\}$, then

$Y(0)$: outcome if intervention sets $\mathbf{z} = 0$ ('Untreated')

$Y(1)$: outcome if intervention sets $\mathbf{z} = 1$ ('Treated')

- The *observed* outcome Y is then

$$Y = (1 - Z)Y(0) + ZY(1) = \begin{cases} Y(0) & Z = 0 \\ Y(1) & Z = 1 \end{cases}.$$

Hypothetical data generating mechanism:

- individual brings their characteristics X ;
- for each z , the outcome $Y(z)$ is determined by X ;
- for *observed* treatment $Z = z$, we observe $Y = Y(z)$.

In an *experimental study* precisely the right kind of ‘intervention’ to study causal contrasts is made:

- we randomly assign Z , *independently* of X ;
- we *compare* the outcomes in the different groups indexed by different Z values.

Example: randomized controlled trials.

Observational studies

In an *observational study* we do *not* intervene to assign treatments, we observe it as part of the data collection process.

- we *cannot* treat Z as if it were independent of X ;
- groups with different Z may have *different* distributions of X , so these groups are *not directly comparable*.

To study cause and effect, we need to have an understanding of the *probabilistic* relationship between all the variables we observe.

To gain statistical insights, we need to build probability models.

Example: HIV Study

- C - CD4 count (continuous)
- D - assignment of new HIV drug (binary)
- H - underlying HIV severity (binary)
- S - symptoms (binary)
- U - follow up indicator (binary, $U = 0$ implies loss to follow-up)

We wish to examine the impact of the new drug on CD4 count:

Does intervening to change D affect outcome C ?

We need to describe how these variables vary jointly in the study.

A *joint* probability distribution

$$f_{X,Y,Z}(x, y, z)$$

describes how the data are generated. This model specifies

- the *marginal* distributions

$$f_X(x) \quad f_Y(y) \quad f_Z(z)$$

that describe how the variables behave individually,

- the *conditional* distributions such as

$$f_{X|Y}(x|y) \quad f_{X|Z}(x|z) \quad f_{Y|X}(y|x) \quad f_{Y|X,Z}(y|x, z) \quad f_{Y,Z|X}(y, z|x)$$

etc. that describe how the variables behave when one or more variable is *fixed*

We have the possible decompositions

$$f_{X,Y,Z}(x, y, z) = f_X(x)f_{Z|X}(z|x)f_{Y|X,Z}(y|x, z)$$

$$f_{X,Y,Z}(x, y, z) = f_Z(z)f_{Y|Z}(y|z)f_{X|Y,Z}(x|y, z)$$

and so on, for any ordering of the variables.

We can *always* consider this kind of sequential decomposition, which is termed a *chain rule factorization*.

Two random variables X, Z are *independent*

$$X \perp\!\!\!\perp Z$$

if and only if, for all values (x, z) ,

$$f_{X,Z}(x, z) = f_X(x)f_Z(z)$$

$$f_{Z|X}(z|x) = f_Z(z)$$

$$f_{X|Z}(x|z) = f_X(x)$$

i.e. knowledge of X *does not influence* our assessment of Z .

We can consider *conditional independence*: say

$$Y \perp\!\!\!\perp Z \mid X$$

if and only if, *for all* (x, z, y)

$$f_{Y,Z|X}(y, z|x) = f_{Z|X}(z|x)f_{Y|X}(y|x)$$

i.e. if we *fix* $X = x$

- knowledge of Y *does not influence* our assessment of Z :

$$f_{Z|X,Y}(z|x, y) = f_{Z|X}(z|x).$$

- knowledge of Z *does not influence* our assessment of Y :

$$f_{Y|X,Z}(y|x, z) = f_{Y|X}(y|x)$$

Example

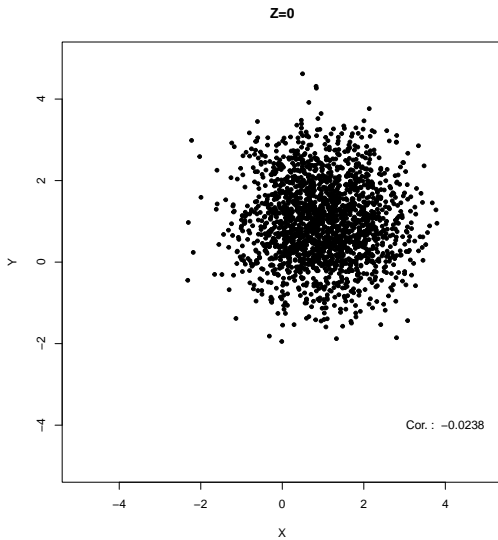
Three variables

- X and Y vary continuously,
- Z is binary.

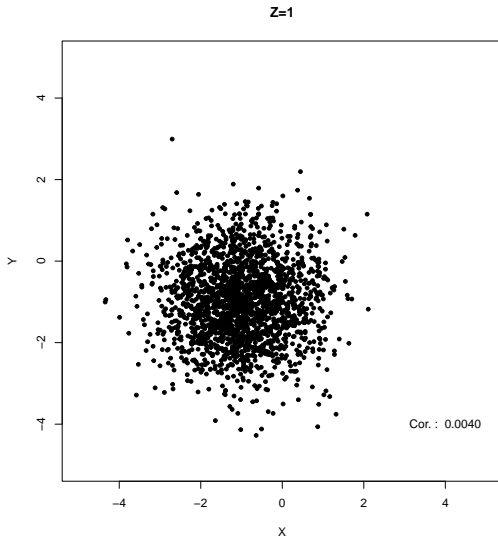
We can study the distribution of the data for X and Y

- for each level of Z separately,
- pooled over Z levels.

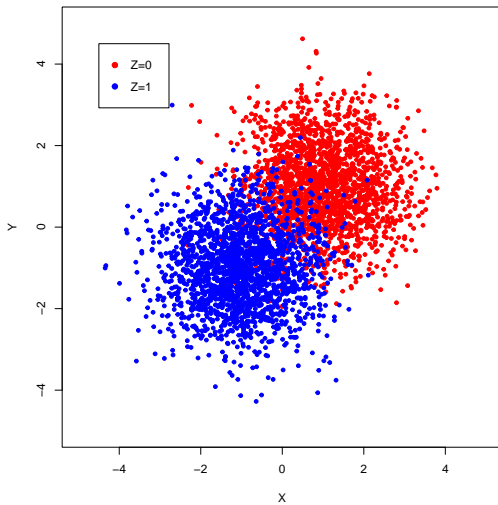
Example



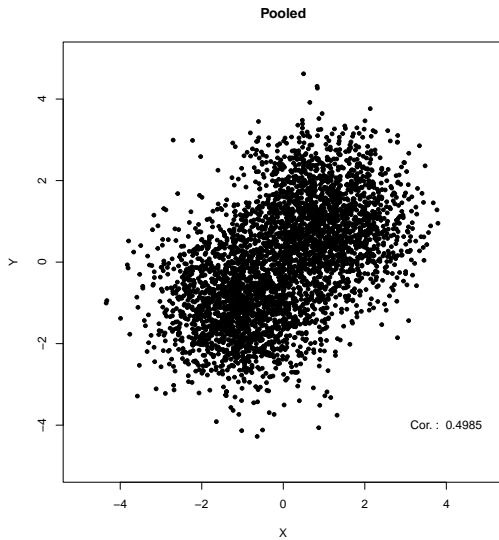
Example



Example



Example



Example

We see that

- for $Z = 0$ and $Z = 1$ separately, X and Y are *uncorrelated*;
- overall X and Y are *positively correlated*.

Thus X and Y are

conditionally unrelated given Z

but are

unconditionally related.

This illustrates that conditioning can remove (or *block*) dependence.

We have a chain rule factorization

$$f_{X,Y,Z}(x, y, z) = f_X(x)f_{Y|X}(y|x)f_{Z|X,Y}(z|x, y).$$

We might then *assume* the *conditional independence*

$$Z \perp\!\!\!\perp Y|X$$

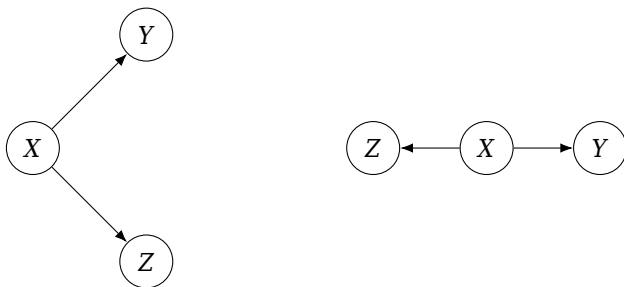
so that

$$f_{Z|X,Y}(z|x, y) = f_{Z|X}(z|x)$$

and so

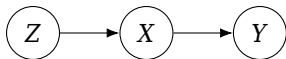
$$f_{X,Y,Z}(x, y, z) = f_X(x)f_{Y|X}(y|x)f_{Z|X}(z|x)$$

We can depict the conditional independence using a graph:



This type of graph is sometimes called a *fork*.

The other common type of graph component is a *chain*



which implies the factorization

$$f_Z(z)f_{X|Z}(x|z)f_{Y|X}(y|x)$$

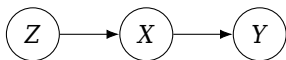
and the conditional independence

$$Y \perp\!\!\!\perp Z|X$$

That is, there are two ways the conditional independence

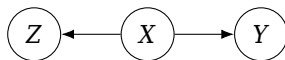
$$Y \perp\!\!\!\perp Z|X$$

could be represented



Chain

$$f_Z(z)f_{X|Z}(x|z)f_{Y|X}(y|x)$$

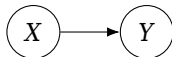


Fork

$$f_X(x)f_{Z|X}(z|x)f_{Y|X}(y|x)$$

- Nodes \textcircled{X} , \textcircled{Y} , \textcircled{Z} denote the variables;
- Edges with *arrows* indicate the nature of dependence in the chain rule factorization;
- *Directed* arrows specify the conditional independence assumptions;

- Nodes without *incoming* edges are *founders*;

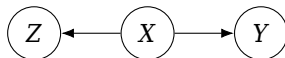


corresponds to

$$f_X(x)f_{Y|X}(y|x)$$

- Nodes with only *outgoing* edges act to *block* dependence.

For example, in



so that

$$f_{X,Y,Z}(x,y,z) = f_X(x)f_{Y|X}(y|x)f_{Z|X}(z|x)$$

it follows that

$$Z \perp\!\!\!\perp Y|X.$$

However, it also follows that, in general

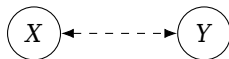
$$Y \not\perp\!\!\!\perp Z$$

(recall the earlier scatterplots)

Causal Graphs

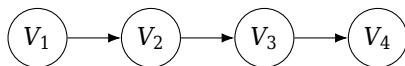
- *Nodes* or *vertices*, V_1, V_2, \dots , represent variables.
- *Edges*, E_1, E_2, \dots , represent dependencies.
- Two nodes are *adjacent* if there is an edge between them.
 - edges can be *directed*, denoted using arrows, or *undirected*;
 - if all edges are directed, the graph is directed.

Note: we can use 'bidirected' (edges with an arrow at each end) to indicate general dependence between two variables



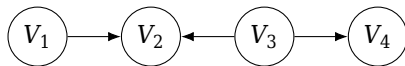
although these will be less important in causal settings.

- A *path* is a sequence of edges that connects two nodes;
 - a *directed* path is a path where the directions of arrows on edges are obeyed



Directed path from V_1 to V_4

whereas an *undirected* path is a path that is not directed.

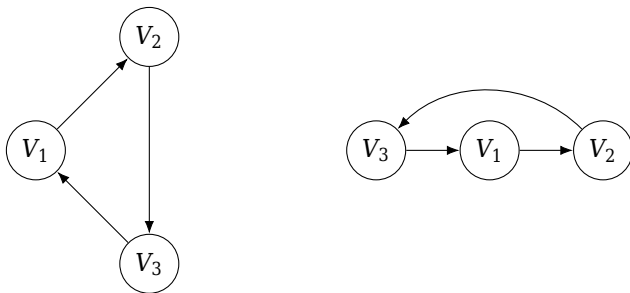


Undirected path from V_1 to V_4

- two nodes are *connected* if a path exists between them, and *disconnected* otherwise.

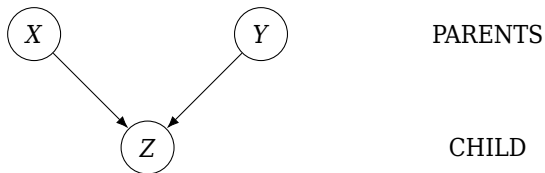
Causal Graphs

In general, a graph may contain *cycles*, that is, directed paths that *start* and *end* at the *same* node.



Directed acyclic graph (DAG): a directed graph with no cycles.

Colliders: Suppose we have the DAG



$$f_{X,Y,Z}(x, y, z) = f_X(x)f_Y(y)f_{Z|X,Y}(z|x, y)$$

In this DAG, we have $X \perp\!\!\!\perp Y$:

$$f_{X,Y}(x,y) = f_X(x)f_Y(y)$$

e.g. X and Y represent the scores on two dice rolled independently, Z is the total score

$$Z = X + Y.$$

We might observe

$$X = 2, Y = 3 \implies Z = 5.$$

However, conditioning on $Z = z$

$$f_{X,Y|Z}(x,y|z) \neq f_X(x|z)f_Y(y|z)$$

in general. Equivalently

$$f_{Y|Z,X}(y|z,x)$$

depends on the value of x .

For example, if we observe $Z = 5$, and we *know* $X = 2$, then *we know with certainty* that $Y = 3$.

That is,

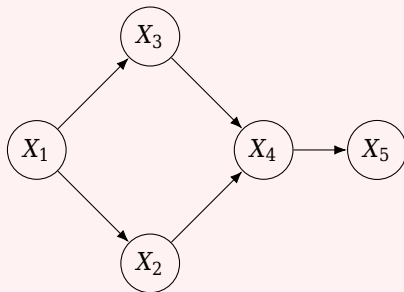
$$X \perp\!\!\!\perp Y$$

but

$$X \not\perp\!\!\!\perp Y \mid Z$$

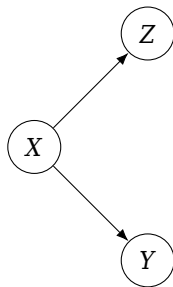
Conditioning on Z *induces dependence*; the node is termed a *collider*.

Example: Factorization



$$f_{X_1}(x_1)f_{X_2|X_1}(x_2|x_1)f_{X_3|X_1}(x_3|x_1)f_{X_4|X_2,X_3}(x_4|x_2,x_3)f_{X_5|X_4}(x_5|x_4)$$

When we write



what precisely does the symbol \longrightarrow mean ?

A *structural* interpretation states that we

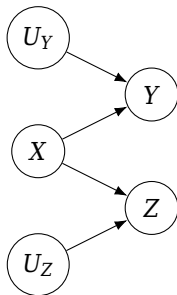
- generate X independently,
- generate Y and Z independently as functions of the realized X , for example

$$Y = 3X$$

$$Z = 4X + 9$$

Structural models and equations

X, U_Z, U_Y
independent



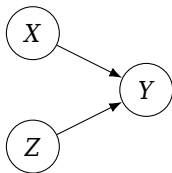
$$Y = g_1(X, U_Y)$$

$$Z = g_2(X, U_Z)$$

For example

$$Y = X + U_Y$$

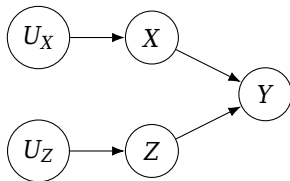
$$Z = X + U_Z$$



$$Y = g(X, Z)$$

Fixing $X = x$ and $Z = z$ fixes $Y = g(x, z)$.

Structural models and equations



$$X = g_1(U_X)$$

$$Z = g_2(U_Z)$$

$$Y = g(X, Z)$$

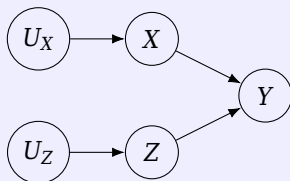
If we know $X = x$ and $Z = z$, then we do not need to know the values of U_X and U_Z to determine Y . That is

$$Y \perp\!\!\!\perp (U_X, U_Z) \mid (X, Z).$$

We can interpret *causation* in terms of these functions.

- X *causes* Y if it appears in the function, g , that assigns Y 's value;
- X *causes* Y if, in the graph representing the joint distribution, there is a *directed path* from X to Y ;
- X is a *direct cause* of Y if there is an arrow from X to Y .

Note



$$X = g_1(U_X)$$

$$Z = g_2(U_Z)$$

$$Y = g(X, Z)$$

so that

$$Y = g(X, Z) = g(g_1(U_X), g_2(U_Z))$$

so both (X, Z) and (U_X, U_Z) can be interpreted as causes of Y .

- X and Z are direct causes,
- U_X and U_Z are indirect causes.

Note

We will proceed by assuming that in a practical setting, the structural relationship and the corresponding causal graph is *known* before any analysis can be carried out.

- Usually in practice this requires expert knowledge;
- Learning the causal graph from data is a hard problem.

Note

If we obtain all variables simultaneously, it is *not possible* to learn which of the possible factorizations is the data generating one.

For example, if we simply observe (X, Y) jointly, we cannot distinguish

$$f_X(x)f_{Y|X}(y|x) \quad \text{from} \quad f_Y(y)f_{X|Y}(x|y)$$

i.e. does X cause Y or does Y cause X ?

Note

In order for X to cause Y , we must have that X *precedes* Y temporally.

The structural equations form the variables on the left hand side from the variables on the right hand side

$$Y = g(X, Z)$$

that is, we *first* generate X and Z , and *then* generate Y .

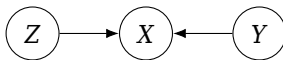
That is, there must be a *temporal* ordering.

To assess whether

$$Y \perp\!\!\!\perp Z \quad \text{or} \quad Y \perp\!\!\!\perp Z \mid X$$

for any distribution compatible with the DAG, we must assess whether there is any way for ‘information’ to ‘flow’ between Z and Y , maybe once X has been accounted for.

First, recall the *collider* graph

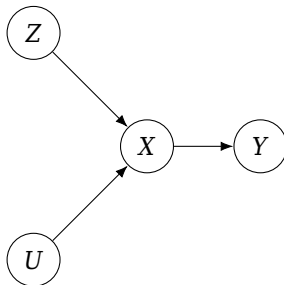


X is a collider on the path between Z and Y . Therefore

$$Y \perp\!\!\!\perp Z \quad \text{but} \quad Y \not\perp\!\!\!\perp Z \mid X$$

Note that a *directed path* from one node to another *cannot* contain a collider.

The notion of being a collider is *path-specific*: for example



- X is a *collider* on path $Z \rightarrow X \rightarrow U$
- X is *not a collider* on path $Z \rightarrow X \rightarrow Y$.

Unconditional d-separation

Consider a general path (directed or undirected) between Z and Y .

The path is *open* (or *unblocked*) if there is *no collider* on the path;

- if there is a collider, the path is *closed* (*blocked*).

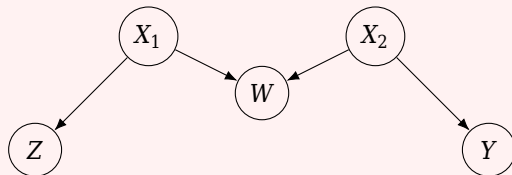
Z and Y are *d-separated* if there is *no open path between them*;

If there is an open path, Z and Y are *d-connected*.

- this path must comprise *chains* or *forks* only

Example: Diabetes example (Rothman et al. p 188)

- X_1 family income
- X_2 genetic risk
- W parental diabetes
- Z low educational attainment
- Y diabetes of subject



Example: Diabetes example (Rothman et al. p 188)

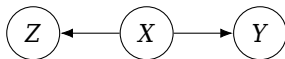
Z and Y are d-separated; there is one path between Z and Y , but it is blocked by the collider W .

$$f_{X_1}(x_1)f_{X_2}(x_2)f_{W|X_1,X_2}(w|x_1,x_2)f_{Z|X_1}(z|x_1)f_{Y|X_2}(y|x_2)$$

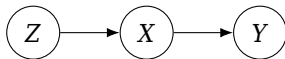
and Z and Y are *independent*.

Conditional d-separation

For a *non-collider* X : *conditioning* on X :



$$Z \perp\!\!\!\perp Y \mid X$$

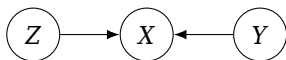


$$Z \perp\!\!\!\perp Y \mid X$$

Conditioning *blocks* the path.

Conditional d-separation

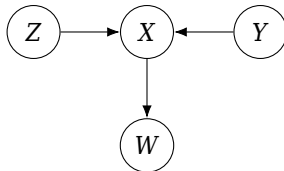
For a *collider* X : conditioning on X *opens* the path



$$Z \not\perp\!\!\!\perp Y \mid X$$

Conditional d-separation

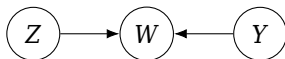
Consider the following DAG:



and conditioning on *a descendant*, W , of X :

$$\begin{aligned} f_{Z,Y,W}(x, y, w) &= f_Z(z)f_Y(y) \int f_{X|Z,Y}(x|z, y)f_{W|X}(w|x) dx \\ &= f_Z(z)f_Y(y)f_{W|Z,Y}(w|z, y) \end{aligned}$$

Therefore we have that



$$Z \not\perp\!\!\!\perp Y \mid W$$

and so W *is a collider* in the reduced graph.

Therefore

- (i) conditioning on a *non-collider* X *blocks* the path at X ;
- (ii) conditioning on
 - a *collider* X or
 - a *descendant* W of X*opens* the path at X ;

Consider two nodes X and Y with possibly several open paths connecting them. Suppose S is a set of variables.

- S *blocks* a path if, after conditioning on S , the path is *closed*;
- S *unblocks* a path if after conditioning the path is *open*;
- If S blocks *every path*, then X and Y are *d-separated* by S .

- If S d-separates X and Y , then

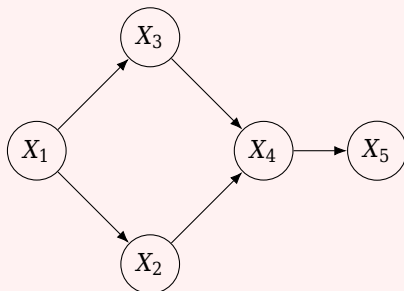
$$X \perp\!\!\!\perp Y \mid S,$$

so that

$$f_{X|Y,S}(x|y,s) \equiv f_{X|S}(x|s) \quad \forall (x,y,s).$$

X and Y are *conditionally independent* given S .

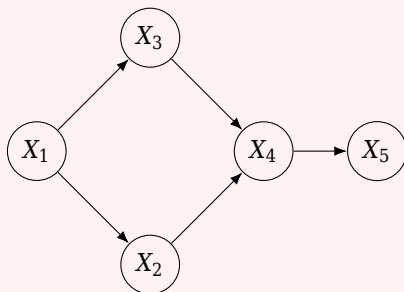
Example:



$\{X_2\}$ and $\{X_3\}$ are d-separated by $\{X_1\}$, and $X_2 \perp\!\!\!\perp X_3 \mid X_1$.

- there are two paths between X_2 and X_3 ;
 - $X_2 \rightarrow X_1 \rightarrow X_3$: blocked by conditioning on X_1 .
 - $X_2 \rightarrow X_4 \rightarrow X_3$: blocked by the collider at X_4 , and $X_4 \notin \{X_1\}$.

Example:

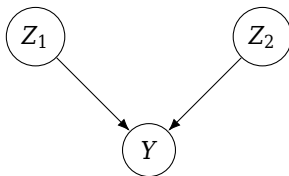


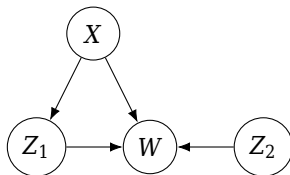
$\{X_2\}$ and $\{X_3\}$ are *not* d-separated by $\{X_1, X_5\}$:

- $X_2 \not\perp\!\!\!\perp X_3 \mid (X_1, X_5)$.
- X_5 is a descendant of collider X_4 ;

Conditioning on the common effect of two causes renders the two causes dependent;

- this is known as *selection bias* or *Berkson bias*
- it is the effect we observe in the collider graph

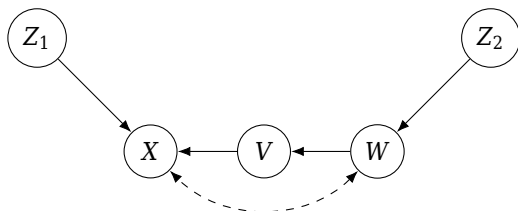




Here $Z_1 \perp\!\!\!\perp Z_2$: there are two paths to consider

- $Z_1 \rightarrow X \rightarrow W \rightarrow Z_2$
- $Z_1 \rightarrow W \rightarrow Z_2$

both blocked by collider W . Therefore $Z_1 \not\perp\!\!\!\perp Z_2 \mid \{W\}$.



Here $Z_1 \perp\!\!\!\perp Z_2$: there are two paths to consider

- $Z_1 \rightarrow X \rightarrow V \rightarrow W \rightarrow Z_2$ is blocked by the collider X .
- $Z_1 \rightarrow X \rightarrow W \rightarrow Z_2$ is blocked by the colliders X and W .

Therefore $Z_1 \not\perp\!\!\!\perp Z_2 \mid \{X, W\}$.

If X and Y are d-separated by S then

$$X \perp\!\!\!\perp Y \mid S$$

for *all* distributions compatible with the graph; conversely, if they are not d-separated, then X and Y are *dependent* given S for at least one distribution compatible with the graph.

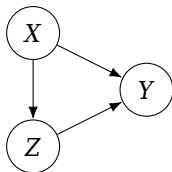
Intervening set the level of Z to z has the effect of

- removing all *incoming* arrows to Z
- switching the marginal for Z to the *degenerate distribution* $f_Z^*(.)$

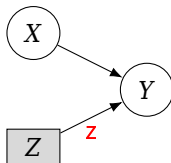
$$f_Z^*(z) = \mathbb{1}_{\{z\}}(z) \quad z \in \mathbb{R}.$$

That is, Z takes the value z with probability 1.

In the earlier example

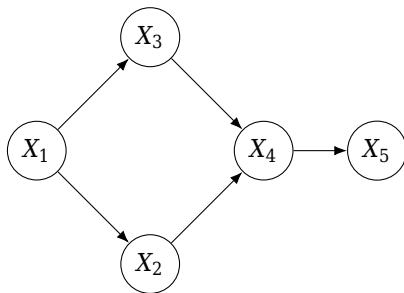


$$f_X(x)f_{Z|X}(z|x)f_{Y|X,Z}(y|x,z)$$



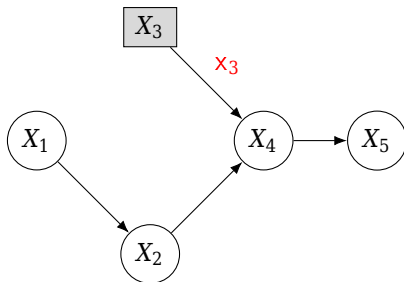
$$f_X(x)f_Z^*(z)f_{Y|X,Z}(y|x,z)$$

Consider the DAG



$$f_{X_1}(x_1)f_{X_2|X_1}(x_2|x_1)f_{X_3|X_1}(x_3|x_1)f_{X_4|X_2,X_3}(x_4|x_2,x_3)f_{X_5|X_4}(x_5|x_4)$$

Suppose we *intervene* to set $X_3 = \mathbf{x}_3$. The relevant DAG is



$$f_{X_1}(\mathbf{x}_1)f_{X_2|X_1}(\mathbf{x}_2|\mathbf{x}_1)f_{X_3}^*(\mathbf{x}_3)f_{X_4|X_2,X_3}(\mathbf{x}_4|\mathbf{x}_2,\mathbf{x}_3)f_{X_5|X_4}(\mathbf{x}_5|\mathbf{x}_4)$$

and X_1 is *no longer a cause* of X_3 .

We aim to understand the effect of Z on Y .

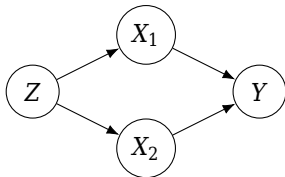
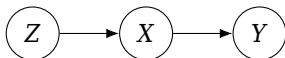
- An *open undirected path* between Z and Y allows for the association between Z and Y to be modified by the presence of other variables.

This is known as a *biasing* path.

- By 'association', we mean some form of *correlation*.
- Usually association is estimated using regression,

Graphical representation of bias

- The association between Z and Y is *unbiased* for the effect of Z on Y if the only open paths between them are *directed paths*.

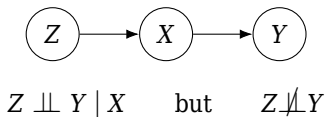


Consider a set of nodes S :

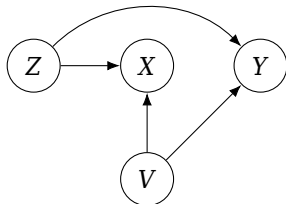
- S is *sufficient* to control bias in the association between Z and Y if after conditioning on S the remaining *open* paths between Z and Y are precisely the *directed* paths between Z and Y ;
- S is *minimally sufficient* if it is the *smallest* sufficient set.

Note: Conditioning on descendants of Z

(i) *blocks* directed paths



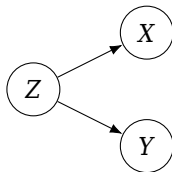
(ii) may *create* paths that lead to *biasing* of the effect of Z on Y.



In this graph,

- the only open path between Z and Y is the *direct* path;
- conditioning on X opens a *biasing* path.

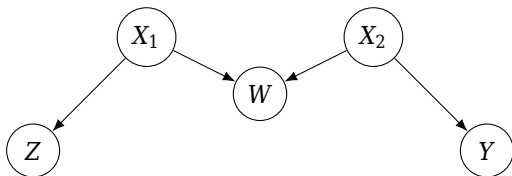
(iii) may be unnecessary in statistical terms: for example



In this graph, conditioning on X will not affect bias.

Graphical representation of bias

Undirected paths from Z to Y are termed *backdoor* paths (relative to Z) if they *start* with an arrow pointing *into* Z .



The only path from Z to Y is a backdoor path; however, it is not open because of the collider W .

Before conditioning

- *all biasing* paths in a DAG are backdoor paths, and
- all *open* backdoor paths are biasing paths.

To obtain an unbiased estimate of the effect of Z on Y , all backdoor paths between Z and Y must be *blocked*.

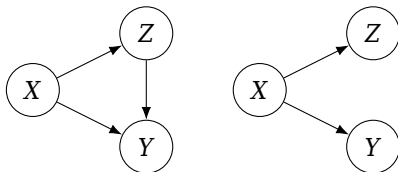
Set S satisfies the *backdoor criterion* with respect to Z and Y if

- (i) S *contains no descendant* of Z , and
- (ii) there is *no open backdoor path* from Z to Y after *conditioning* on S .

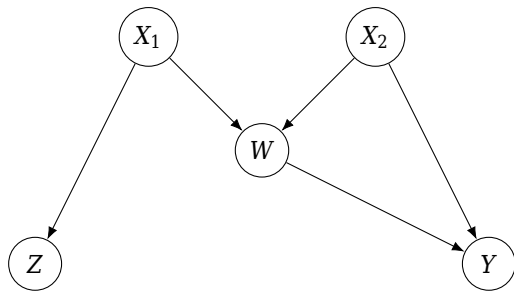
A *confounding path* between Z and Y is

- (i) a *biasing* path (that is, an *undirected open path*) that
- (ii) *ends* with an arrow into Y.

Variables on a confounding path are termed *confounders*.



X is a confounder in both cases.



W is a collider on the undirected path from Z to Y

Path 1: $Z \rightarrow X_1 \rightarrow W \rightarrow X_2 \rightarrow Y$

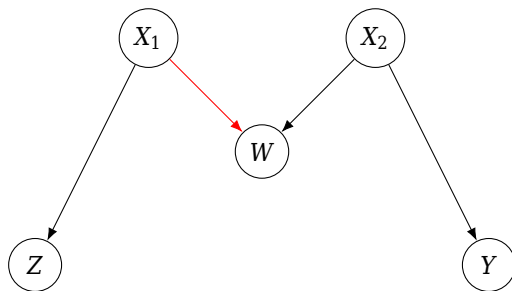
and hence this path is blocked.

However unconditional on W , the effect of Z on Y is confounded by the backdoor path

$$\text{Path 2: } Z \rightarrow X_1 \rightarrow W \rightarrow Y.$$

Conditioning on W alone opens Path 1, therefore to block both paths, we need to condition on

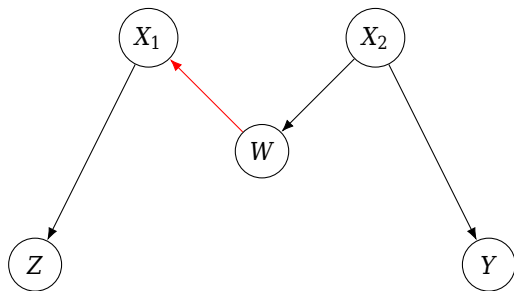
$$S \equiv \{W, X_2\}.$$



Conditioning on W *opens* the confounding path. Therefore $Z \not\perp\!\!\!\perp Y$ (as there is no open path between them), but

$$Z \not\perp\!\!\!\perp Y \mid W$$

Further conditioning on either $\{X_1\}$ or $\{X_2\}$ blocks the path.



Conditioning on W *blocks* the confounding path. Therefore conditioning on any one of

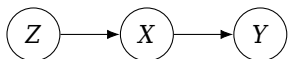
$$\{X_1\}, \{W\}, \{X_2\}$$

will block the path.

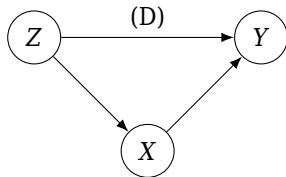
For the effect of Z on Y relative to X :

- *Direct effect*: A direct effect of Z on Y is the effect captured by a *directed* path from Z to Y that does not pass through X .
- *Indirect effect*: An indirect effect of X on Y that is captured by directed paths that pass through X .
 - X is termed an *intermediate* or *mediator* variable.

Direct and indirect effects

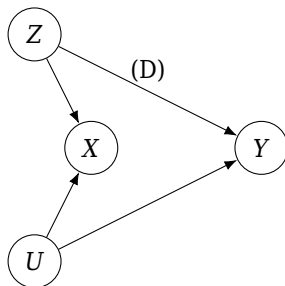


Indirect effect



Direct (D) & Indirect effect

X is a mediator of the indirect effect



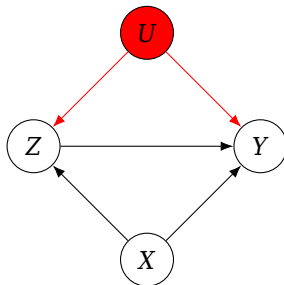
No indirect effect

Direct effect is not confounded

X is a collider, so there is no other open path from Z to Y .

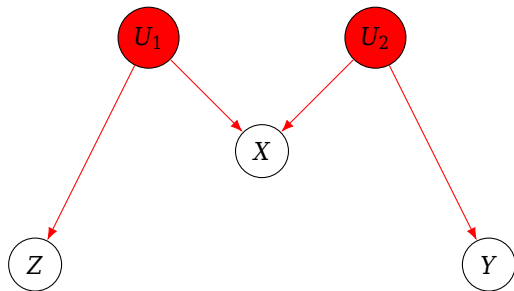
Unmeasured confounding

Suppose that in reality there is a further variable U that is a confounder, but is unmeasured in the observed data.



There is a hidden confounding path $Z \rightarrow U \rightarrow Y$. Conditioning on U is not possible, as we are unaware of its existence.

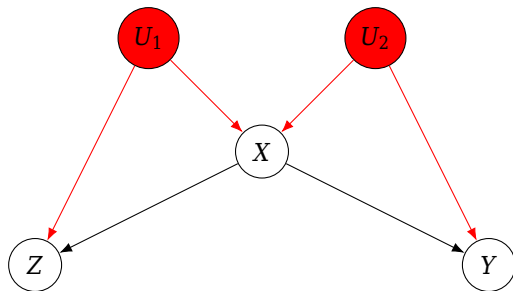
With two unmeasured confounders:



We have that X , Y and Z are *independent*; the (true but hidden) path between Z and Y is blocked at collider X .

Unmeasured confounding

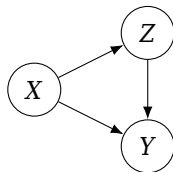
Suppose we condition on X :



In the modelled DAG, $Y \perp\!\!\!\perp Z \mid X$; however, conditioning on X *opens* the *hidden* path through U_1 and U_2 , so there is now an open biasing path.

This is sometimes referred to as the *M-bias* phenomenon.

We have seen that conditioning on variables can close biasing paths, allowing an unbiased assessment of the causal effect of Z on Y .



The open, undirected path

$$Z \rightarrow X \rightarrow Y$$

can be blocked by conditioning on X .

If all the variables are jointly Normally distributed, then this conditioning can be achieved by including X as a predictor in a linear regression model of Y on Z .

That is, we can fit the linear model where

$$\mathbb{E}[Y|X = x, Z = z] = \beta_0 + \beta_1 x + \psi z$$

and estimate the direct effect of Z on Y by estimating ψ .

Note

Blocking confounding paths (e.g. by conditioning) is not quite the end of the story.

Typically we need to utilize *parametric* inference, and there is usually a requirement that certain parametric models are *correctly specified*.

In a statistical formulation of a causal inference problem

1. We identify *treatment* Z and *outcome* Y
2. We form the *DAG* representing the relationships between Z and Y which contains other measured variables X .
3. The causal effect of Z on Y flows down *open* and *directed* paths from Z to Y ;
 - there may be a *direct* effect if there is an arrow from Z into Y ;
 - there may also be *indirect* effects if the directed path passes through *mediating* variables.

4. If there are *undirected* paths from Z to Y that are open, then these paths may induce *bias* in estimation of the effect of Z on Y .
5. In order to obtain unbiased estimation, the open undirected (biasing) paths must be *blocked*; typically this is done by *conditioning* on variables on those paths.
6. A *collider* node blocks a path; however, conditioning on the collider *opens* the path at that node.