

---

# Causal Inference using Causal Graphs in Epidemiologic Context

---

Ehsan KARIM

Department of Statistics, UBC

[ehsan@stat.ubc.ca](mailto:ehsan@stat.ubc.ca)

December 11, 2010

## Contents

<b>1</b>	<b>Prelude</b>	<b>1</b>
1.1	Motivation . . . . .	1
1.2	Objective . . . . .	1
1.3	Organization . . . . .	2
<b>2</b>	<b>Causal Diagrams</b>	<b>2</b>
2.1	Terminologies . . . . .	3
2.2	d-separation Rules . . . . .	5
2.2.1	Unconditional Rule . . . . .	5
2.2.2	Conditional Rule . . . . .	5
2.3	Statistical Relationships Implied by DAGs . . . . .	5
2.3.1	Compatibility Rule . . . . .	6
2.3.2	Weak Faithfulness Rule . . . . .	6
2.4	Illustration . . . . .	6
<b>3</b>	<b>Identifying Confounding using DAGs</b>	<b>7</b>
3.1	Back-door Criterion . . . . .	7
3.1.1	Criterion or Rule . . . . .	7
3.1.2	Steps for Identifying Confounding . . . . .	7
3.2	Illustration . . . . .	8
3.2.1	Identifying Confounding using Back-door Criterion . . . . .	8
3.2.2	Comparing Approaches . . . . .	8
<b>4</b>	<b>Control or Conditioning</b>	<b>10</b>
4.1	Control in the Design Stage . . . . .	10
4.1.1	Random Mechanism . . . . .	10
4.1.2	Exchangeability . . . . .	10
4.1.3	Matching . . . . .	10
4.2	Control in the Analysis Stage . . . . .	11
4.2.1	Stratification . . . . .	11
4.2.2	Regression Modelling . . . . .	11
4.2.3	Propensity Score . . . . .	11
<b>5</b>	<b>DAGs and Longitudinal Settings</b>	<b>11</b>
5.1	Scenario I . . . . .	12
5.2	Scenario II . . . . .	13

<b>6 Conclusion</b>	<b>14</b>
<b>Appendix</b>	<b>15</b>
<b>A Causation</b>	<b>15</b>
A.1 Philosophical Perspective . . . . .	15
A.2 Counterfactuals . . . . .	15
A.2.1 Causality and Counterfactual Models . . . . .	15
A.2.2 Relation with Philosophical Definition . . . . .	16
A.2.3 Probabilistic Causality . . . . .	16
A.2.4 General Methodology to estimate Causal Effect . . . . .	17
A.2.5 Illustration . . . . .	18
A.2.6 Confounding in terms of Counterfactuals . . . . .	20
<b>B Confounding and Relevant Concepts</b>	<b>22</b>
B.1 Conventional Definition of Confounding . . . . .	22
B.2 Residual Confounding . . . . .	23
B.3 Over-adjustment . . . . .	23
B.4 Sufficient Set of Confounders . . . . .	23
B.5 Ignorability . . . . .	23
B.6 Confounding and Collapsibility . . . . .	24
<b>C Bibliographic Notes</b>	<b>25</b>

# 1 Prelude

## 1.1 Motivation

In most of the scientific studies, exploration of causation has been the ultimate goal of research. Whenever a model is constructed for the purpose of prediction, researchers implicitly provide the model a causal interpretation unconsciously. Yet, most of the research results are based on correlation or some other measure of association among various variables under consideration. This is mostly because statisticians are mostly equipped with tools that can measure association. But association does not imply causation. Simple association results may sometimes be misleading or inadequate in assessing the relationship between variables. This is especially true in the field of epidemiology, while evaluating disease-exposure relationships. Finding out the cause of a health related outcome is usually the focus. Statistical association is merely an intermediate step in the process.

Confounding is an important issue in any Epidemiologic study since it distorts the relationship between the exposure variable and the outcome variable. Using the standard measures of association, researchers try to identify whether the resulting association is due to only the exposure variable of interest, or whether there exists some extraneous variables i.e., confounders that impose spurious association. Many attempts have been taken to define confounding properly. Traditional definition provided by many textbooks fail to capture the complete idea of confounding. The idea of counterfactuals are better ways frame the concept of confounding. However, although theoretically this approach is very useful to illustrate what confounding really is, implementation of these counterfactual models is not practically viable in a real world setting. Therefore, researchers often resort to graphical models to identify causal relationships.

## 1.2 Objective

The mini-lecture is focused on causal inference using causal diagrams. We also see the use of this approach in identifying the obstacle in the process of causal inference, such as, confounding. Instead of introducing all the conditional notations and mathematical jargons behind it, various epidemiologic examples will be used to explain the ideas. The ultimate goal is to provide a gentle introduction to the concepts of causal graphs.

### 1.3 Organization

In order to explain the ideas of causal diagrams, at first, the preliminary terminologies will be defined in §2.1. These graphs abide by various basic rules. They will be stated in §2.2. A few assumptions or rules are necessary to connect these graphical rules with statistical association, which will be stated in §2.3.

One major use of this approach is to find the situation where confounding exists. This will be explained, illustrated and compared with conventional methods in §3.

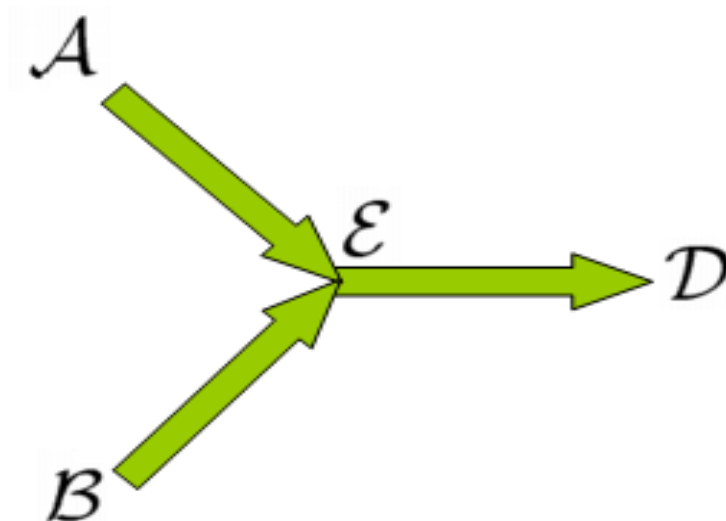
Conditioning or controlling for some variables in probabilistic sense is dealt by this graphical approach. However, to physically implement these conditioning, several approaches can be utilized. Some practical choices will be discussed in §4.

Examples of causal graphs related to longitudinal measurement from observational study will be shown in §5 without much technical details.

Also, the ideas of “causation” and “confounding” are pre-requisites in understanding the idea of causal inference. However, to make sure the focus of this talk does not shift from the main issue of this lecture, which is “causal graphs”, these external but relevant ideas are placed separately in the appendix. The framework of causation will be explained in §A in terms of philosophical perspective and counterfactual models. Then confounding will be defined in §B.1 using conventional and counterfactual approaches as well. Some miscellaneous relevant ideas will also be discussed in §B.

## 2 Causal Diagrams

Path diagrams have been used since the beginning of twentieth century to informally summarize various assumptions of models of statistical association. In 1980s, formal use of probabilities were integrated with these diagrams to express dependence or independence among variables from a structure of joint probability distribution. This opened the door of probabilistic inference using graphical approach incorporating model uncertainties. Probability theory is used causal connections sine it provides means to draw inference from the observed information. A decade later, these graphs was recognized as a tool for causal inference. These graphs enabled researcher to tell whether the causal effects can be estimated from available data, or what else variable or information is required to make a valid inference.

**Figure 2.1:** Explaining the Basics of a Causal Graph

Eventually, Pearl [2000] formalized these ideas in an organized fashion by developing directed acyclic graph (DAG) theory and relevant rules that helps us make causal inference. Causal DAGs are very intuitive and flexible tool to give insight about the process of cause-effect causal mechanism. These diagrams are well equipped to handle more than one cause and effect.

## 2.1 Terminologies

An overview of the components of causal graph theory will require readers to be familiar with the following terminologies. Figure 2.1 will be referred frequently to explain the definitions.

**Node** Each variable under consideration ( $A$ ,  $B$ ,  $E$  and  $D$  in figure 2.1) defines a node in the graph.

**Edge or Arc** Edge is the arrow ( $\rightarrow$  or  $\leftarrow$  or  $\leftrightarrow$ ) that connects causal variable to its effect. Statistical dependence is usually represented by connectivity of edges.

**Path** A path is a sequence of edge or arrows (say,  $\rightarrow \dots \rightarrow \dots \leftarrow$ ), connecting nodes or variables regardless of the directions of the arrows. There exist a path between  $D$  and  $A$ . The absence of a path between two variables would mean that there is no causal relationship between them.

**Direct Cause** If there is one arrow between  $\mathcal{E}$  and  $\mathcal{D}$ , then the relationship is direct in figure 2.1. These variables are then called adjacent or neighbors.

**Indirect Cause** If there is a sequence of arrows between  $\mathcal{A}$  and  $\mathcal{D}$ , then the relationship is indirect.

**Descendant or child** Descendant of  $\mathcal{E}$  will be the factors following it, say  $\mathcal{D}$ .

**Ancestor or parent** Ancestor of  $\mathcal{E}$  will be the factors prior to it, say,  $\mathcal{A}$ . The set of all parents of a given variable  $\mathcal{E}$  is denoted by  $pa[\mathcal{E}]$ . In this case,  $pa[\mathcal{E}] = [\mathcal{A}, \mathcal{B}]$ .

**Exogenous or external** The variable that has no parent is an exogenous variable. Here  $\mathcal{A}$  and  $\mathcal{B}$  are exogenous variables. Otherwise, it is endogenous or internal variable (such as,  $\mathcal{E}$  and  $\mathcal{D}$ ).

**Sink or terminal node** The variable that has no children, say  $\mathcal{D}$ .

**Essential Variables** If two or more factors are being caused by another (parent) factor, then we need to include that other (parent) factor as well. That variable does not have to be observed, i.e., we can also include a variable that is not observed. Otherwise, it is not essential to include all the causes of each variables.

**Directed Path** This is a path traced out entirely along arrows tail-to-head. If there is a directed path from  $\mathcal{A}$  to  $\mathcal{D}$ , then  $\mathcal{A}$  is parent of  $\mathcal{D}$ , and  $\mathcal{D}$  is a child of  $\mathcal{A}$ .

**Directed Acyclic Graph** There must not be any cycles or feedback loop between variables, directly or indirectly. If  $\mathcal{E}$  causes  $\mathcal{D}$ , then  $\mathcal{D}$  can not cause  $\mathcal{E}$  in these graphs. Same variables can not be cause and effect at the same time. That is, no variable can be its own parent or its own child. Therefore, figure 2.1 is a Directed acyclic graph (DAG). A DAG is a way of expressing certain constraints on the joint distribution of a set of variables or nodes. This is an elegant way to represent a complete causal structure that explains all the dependencies among variables.

**Intercept or Mediate** A variable that is in the path, but not at the end, say,  $\mathcal{E}$  in this case.

**Collider** If arrows from  $\mathcal{A}$  and  $\mathcal{B}$  both point to another factor  $\mathcal{E}$ , then  $\mathcal{E}$  is said to be a collider. The collider  $\mathcal{E}$  is a common effect of  $\mathcal{A}$  and  $\mathcal{B}$ . Whether a variable is a collider, depend on the path. For example, in the figure,  $\mathcal{E}$  is a collider in the path  $\mathcal{A} \rightarrow \mathcal{E} \leftarrow \mathcal{B}$ , but not a collider in the path  $\mathcal{A} \rightarrow \mathcal{E} \rightarrow \mathcal{D}$ .

**Open or Active Path** A path is open at non-collider. Therefore, the path  $\mathcal{A} \rightarrow \mathcal{E} \rightarrow \mathcal{D}$  is open.

**Close or Inactive Path** A path with a collider is closed. Here, the path  $\mathcal{A} \rightarrow \mathcal{E} \leftarrow \mathcal{B}$  is closed. This is also an undirected path and  $\mathcal{E}$  intercepts the path.

## 2.2 *d*-separation Rules

### 2.2.1 Unconditional Rule

Two variables  $\mathcal{A}$  and  $\mathcal{B}$  are *d*-separated if there is no open path between them. If there is an open path, we say that they are *d*-connected. If  $\mathcal{A}$  and  $\mathcal{B}$  are *d*-separated, then it implies:

- $\mathcal{A}$  and  $\mathcal{B}$  are unassociated
- $\mathcal{A}$  and  $\mathcal{B}$  may be marginally independent

Say, in a more complicated situation,  $\mathcal{F} \leftarrow \mathcal{G} \rightarrow \mathcal{H} \leftarrow \mathcal{I} \rightarrow \mathcal{J}$ , path is closed at  $\mathcal{H}$ . Although paths are open elsewhere ( $\mathcal{F} \leftarrow \mathcal{G}$  and  $\mathcal{I} \rightarrow \mathcal{J}$ ), due to a closed path at  $\mathcal{H}$ , we can say,  $\mathcal{F}$  and  $\mathcal{J}$  are *d*-separated.

### 2.2.2 Conditional Rule

A path is closed or blocked by conditioning on a set of variables  $\mathbf{S}$  if either of the following conditions hold:

1. We have a non-collider  $\mathcal{E}$  in the path (say,  $\mathcal{A} \rightarrow \mathcal{E} \rightarrow \mathcal{D}$ ), and non-collider  $\mathcal{E}$  is in  $\mathbf{S}$  (a path that is conditioned by a non-collider  $\mathcal{E}$ ). In this case, unconditionally, path is open at  $\mathcal{E}$ , but conditioning on  $\mathcal{E}$  closes this path. Then such conditioning removes  $\mathcal{E}$  as a source of association between  $\mathcal{A}$  and  $\mathcal{D}$ . The the path between  $\mathcal{A}$  and  $\mathcal{D}$  will be conditionally *d*-separated.
2. We have a collider  $\mathcal{E}$  in the path (say,  $\mathcal{A} \rightarrow \mathcal{E} \leftarrow \mathcal{B}$ ), and neither the collider, not any of its descendent is in  $\mathbf{S}$  (a path that is not conditioned by any collider or collider-descendent). Then the path is blocked by  $\mathcal{E}$ , and  $\mathcal{A}$  and  $\mathcal{B}$  will be *d*-separated. Conversely, if we put  $\mathcal{E}$  or any of its child in  $\mathbf{S}$ , a path will open via  $\mathcal{E}$ , and this will establish association between  $\mathcal{A}$  and  $\mathcal{B}$ , i.e., they will be conditionally *d*-connected.

This conditional *d*-separation criteria basically says that conditioning on a variable reverses the status of that variable on the path.

## 2.3 Statistical Relationships Implied by DAGs

It can be seen that  $\mathcal{E}$  is a collider in the path  $\mathcal{A} \rightarrow \mathcal{E} \leftarrow \mathcal{B}$ , since both  $\mathcal{A}$  and  $\mathcal{B}$  point into  $\mathcal{E}$ . Therefore, ( $\mathcal{A}$  and  $\mathcal{B}$ ) are unconditionally *d*-separated. However, since  $\mathcal{E}$  blocks the path



by being a collider, when they are conditioned on  $\mathcal{E}$ , path of ( $\mathcal{A}$  and  $\mathcal{B}$ ) are conditionally d-connected. But to illustrate the the statistical implication of this, we need the following two rules that links between causal graphs and statistical association:

### 2.3.1 Compatibility Rule

If ( $\mathcal{A}$  and  $\mathcal{B}$ ) are d-separated by  $\mathcal{E}$ , then ( $\mathcal{A}$  and  $\mathcal{B}$ ) will be statistically independent.

### 2.3.2 Weak Faithfulness Rule

If ( $\mathcal{A}$  and  $\mathcal{B}$ ) are d-connected conditional on  $\mathcal{E}$ , then independence of ( $\mathcal{A}, \mathcal{B} | \mathcal{E}$ ) cannot be assumed. In other words, if independence of ( $\mathcal{A}, \mathcal{B} | \mathcal{E}$ ) is assumed, then there exists no d-connecting path between ( $\mathcal{A}$  and  $\mathcal{B}$ ) conditional on  $\mathcal{E}$ .

There are some stronger, alternate versions and extensions of these rules that are often used; such as, ‘causal Markov assumption (CMA)’ and ‘faithfulness or stability’, but they are beyond the scope of this lecture.

## 2.4 Illustration

Using these simple rules, and the basic rules of DAG, from figure 2.1, we can make the following statements:

- $\mathcal{E}$  is a direct cause of  $\mathcal{D}$ , and  $\mathcal{A}$  and  $\mathcal{B}$  are indirect causes.
- $\mathcal{A}$  and  $\mathcal{B}$  both directly causes  $\mathcal{E}$
- $\mathcal{A}$  and  $\mathcal{B}$  have no causal relationship between themselves.
- $(\mathcal{E}, \mathcal{D})$ ,  $(\mathcal{A}, \mathcal{D})$ ,  $(\mathcal{B}, \mathcal{D})$ ,  $(\mathcal{A}, \mathcal{E})$ ,  $(\mathcal{B}, \mathcal{E})$  are statistically dependent since directed paths exists between each of them.
- $(\mathcal{A}, \mathcal{B})$  are statistically independent, since a collider  $\mathcal{E}$  exists in their path (from the unconditional d-separation rule).
- $(\mathcal{A}, \mathcal{D})$  are statistically independent, conditional on  $\mathcal{E}$ , the non-collider in this path (from the conditional d-separation rule).
- $(\mathcal{A}, \mathcal{B})$  are statistically dependent, conditional on  $\mathcal{E}$ , the collider in this path (from the conditional d-separation rule).

### 3 Identifying Confounding using DAGs

While assessing  $\mathcal{E}$  and  $\mathcal{D}$  relationship, any undirected open path is an indication of bias, and hence named as biasing path for the particular effect. Although these paths do not represent the effects of  $\mathcal{E}$  on  $\mathcal{D}$ , they surely can contribute to the measure of association of  $\mathcal{E}$  and  $\mathcal{D}$ . The only case which graphical approach considers as unbiased is when the directed paths are open from  $\mathcal{E}$  to  $\mathcal{D}$ . Bias can creep into the analysis in the form of confounding, selection bias or over-adjustment (discussed in §B.3).

Identify the conditions where a relationship between a cause and an effect is confounded, is one of the most lucrative feature of graphical approach. It can do so using a criterion called ‘back-door criterion’:

#### 3.1 Back-door Criterion

The rule that detects confounding from a causal diagram is called the ‘back-door criterion’. It should be noted that, this rule does not define a confounder. It merely describes the situation when the statistical association between two variables would be confounded. Thus, we should regard this rule as a identifier of a precise condition for lack of confounding situation.

##### 3.1.1 Criterion or Rule

Usually, undirected paths from  $\mathcal{E}$  to  $\mathcal{D}$  are known as back-door, when the arrow points into  $\mathcal{E}$ . In Figure 3.1,  $\mathcal{E} \leftarrow \mathcal{C} \rightarrow \mathcal{D}$  is a backdoor path from  $\mathcal{E}$  to  $\mathcal{D}$ . These undirected paths are then biasing paths and non-causal association might arise between  $\mathcal{E}$  and  $\mathcal{D}$ . A set of variables  $\mathcal{F}$  is said to satisfy the back-door criterion relative to  $(\mathcal{E}, \mathcal{D})$  if

1. no variable in  $\mathcal{F}$  is a child of  $\mathcal{E}$ , and
2. every path between  $\mathcal{E}$  and  $\mathcal{D}$  that contains an arrow into  $\mathcal{E}$  is blocked by  $\mathcal{F}$ .

##### 3.1.2 Steps for Identifying Confounding

To determine the existence of confounding in a given situation, the following steps should be executed:

1. in the first step delete all the arrows starting from  $\mathcal{E}$
2. in the second step, find any unblocked back-door path from  $\mathcal{E}$  to  $\mathcal{D}$ .

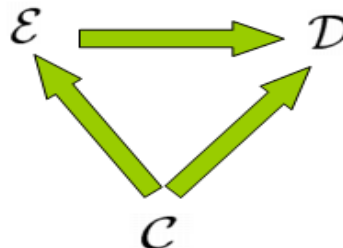
If no such path exists, then, by this criterion, we can say that there exists no confounding in that situation.

**Figure 3.1:** Simple Causal relationship and Confounding

Simple Causal relationship



Confounding Situation



## 3.2 Illustration

### 3.2.1 Identifying Confounding using Back-door Criterion

Let us consider a very complicated scenario of figure 3.2. Here we can see that  $\mathcal{C}$  fails the back-door criterion relative to  $(\mathcal{E}, \mathcal{D})$ , since  $\mathcal{C}$  is a collider in the path  $\mathcal{E} \leftarrow \mathcal{F}_3 \leftarrow \mathcal{F}_1 \rightarrow \mathcal{C} \leftarrow \mathcal{F}_2 \rightarrow \mathcal{F}_5 \rightarrow \mathcal{D}$ . After conditioned on it, the back-door path ( $\mathcal{E} \leftarrow \mathcal{F}_3 \leftarrow \mathcal{F}_1 \rightarrow \mathcal{C} \leftarrow \mathcal{F}_2 \rightarrow \mathcal{F}_5 \rightarrow \mathcal{D}$ ) is blocked, but the biasing path ( $\mathcal{E} \leftarrow \mathcal{C} \rightarrow \mathcal{D}$ ) is still open. Now, let us include a non-collider  $\mathcal{F}_3$  with  $\mathcal{C}$  in  $\mathbf{S}$ . Conditioning on a non-collider  $\mathcal{F}_3$  will alter the status of this node. Now,  $\mathbf{S} = (\mathcal{F}_3, \mathcal{C})$  blocks all back-door paths, and this is a lack of confounding situation. Therefore, we should be able to estimate the causal effects of  $\mathcal{E}$  on  $\mathcal{D}$  without any interruption.

### 3.2.2 Comparing Approaches

Conventional methods of identifying confounder existed well before this DAG approach. One might be wondering why a new approach was necessary to develop in the presence of an existing one. Researchers showed that there are some situations where conventional definition fails to identify confounding, but DAG approach successfully pin-points confounding. The following example was mentioned by Oakes and Kaufman [2006] and many others. The illustration is in figure 3.3. Let us denote

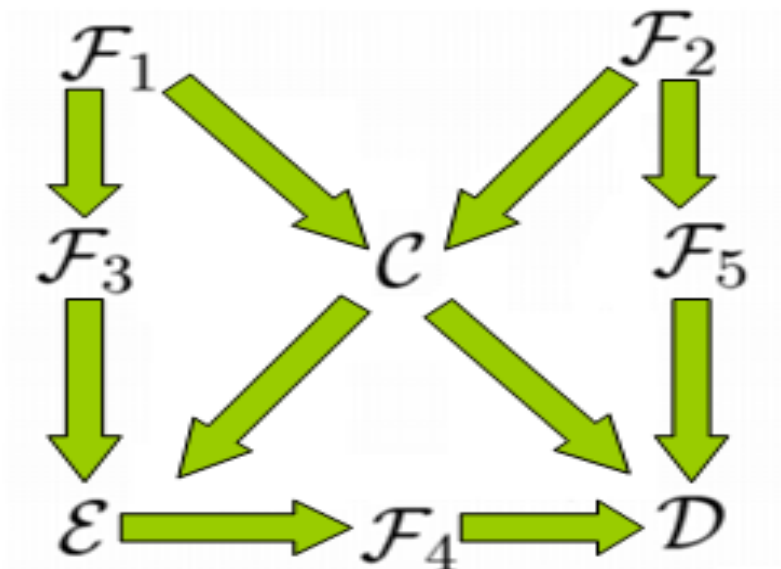
$\mathcal{E}$  = Respondent's Low education

$\mathcal{F}_1$  = Family income during childhood (unobserved)

$\mathcal{F}_2$  = Mother's genetic diabetes risk (unobserved)

$\mathcal{C}$  = Mother had diabetes

Figure 3.2: Application of Back-door Criterion



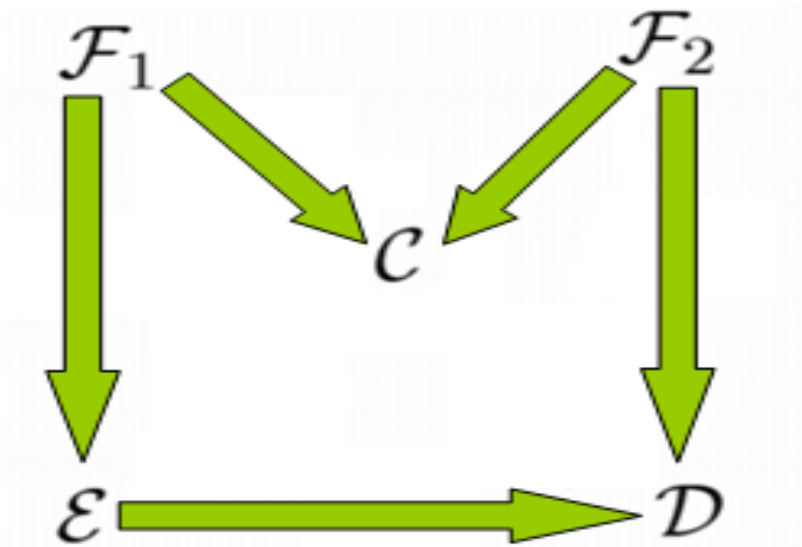
$\mathcal{D}$  = Respondent had diabetes

We are interested in the  $\mathcal{E}$ ,  $\mathcal{D}$  relationship.

Here, family income during childhood,  $\mathcal{F}_1$  have impact on both mother had diabetes,  $\mathcal{C}$  and respondent's low education,  $\mathcal{E}$ , since as the individual was poor during his childhood, he had less chance in education,  $\mathcal{E}$ . Also, since mother was poor, her poverty increased her risk of having diabetes,  $\mathcal{C}$  (known from study results). Therefore,  $\mathcal{C}$  and  $\mathcal{E}$  are related.

Now, mother's genetic diabetes risk  $\mathcal{F}_2$  have impact on both mother had diabetes,  $\mathcal{C}$  and respondent had diabetes,  $\mathcal{D}$ , therefore,  $\mathcal{C}$  and  $\mathcal{D}$  are also related. However, mother's diabetes,  $\mathcal{C}$  should not be effected by respondent's low income,  $\mathcal{E}$  or respondent's diabetes,  $\mathcal{D}$ . Thus, by conventional definition of confounding as described in §B.1,  $\mathcal{C}$  is a confounder.

Notice that, we see that there is only one path in between low education,  $\mathcal{E}$  and own diabetes,  $\mathcal{D}$  (other than the direct path). When we check  $\mathcal{C}$  using the back-door criterion, we see that  $\mathcal{C}$  already blocks the path, since mother's diabetes  $\mathcal{C}$  is a collider. Rather, if we want to adjust for mother's diabetes,  $\mathcal{C}$ , then the path is unblocked, inducing spurious statistical association. Therefore, we see that in graphical criteria, we can achieve better results, because it is easy to identify collider there. Conditioning on a collider is often very misleading, and we should be cautious about that.

**Figure 3.3:** Superiority of Graphical Method

## 4 Control or Conditioning

Once the set of variables  $\mathcal{C}$  are identified as confounder, we need to control them by conditioning on those variables  $\mathcal{C}$ . Usually this means, controlling either in the design stage, or in the analysis stage.

### 4.1 Control in the Design Stage

#### 4.1.1 Random Mechanism

In presence of randomization mechanism, the comparison groups are expected to be balanced and hence comparable with respect to all covariates. The word ‘expectation’ indicated availability of large sample sizes. In the studies, where randomization is not implemented, say in observational studies; then confounding is more likely to occur.

#### 4.1.2 Exchangeability

The idea of exchangeability is somewhat similar to counterfactuals (explained in §A.2.1). This assumption makes sure that the comparison groups are both from the same target population, and hence free from any systematic bias.

#### 4.1.3 Matching

Matching and restriction are also the tools to eradicate the potential imbalance between comparison groups. Too much matching variable, however, causes further problem.

## 4.2 Control in the Analysis Stage

Often situation arises when randomization is not feasible, or the data is already collected. In such situations, estimates of effects should be adjusted if existence of confounding is confirmed.

### 4.2.1 Stratification

The idea is to control the variation due to the confounder,  $\mathcal{C}$ . One need to stratify the target population according to confounder  $\mathcal{C}$  and make separate sub-groups, each of which shared common values of  $\mathcal{C}$ . For example, for confounders that take discrete values, we can create different strata for different values of  $\mathcal{C}$ . Then, when we measure the causal effects of  $\mathcal{E}$  on  $\mathcal{D}$  within each balanced strata, under the assumption of conditional randomization, the effect due to variation of the confounder  $\mathcal{C}$  will vanish.

### 4.2.2 Regression Modelling

A variety of multivariate regression analysis that defines the relationship between exposure and the response are used to control confounders. Say, logistic regression for binary dependent variable, conditional logistic regression for matched data or analysis of covariance approach in some other situations.

### 4.2.3 Propensity Score

Propensity score method is superior to usual stratification method or matching, due to its ability to account for a large number of confounders or control variables.

## 5 DAGs and Longitudinal Settings

A longitudinal data allows researchers to understand the temporal order of the variables under consideration. This enables them to determine whether changes of suspected exposure variable are happening before the changes in the outcome variable. Observing this sequence establishes the causal interpretation of the process. In particular, having access to longitudinal data is advantageous in observational studies in this respect.

However, this is not without a price. Several analytic issues creep in while analyzing longitudinal data, especially due to existence of time-varying exposure and confounding variables. To make it worse, these explanatory variables could interact with response process. This requires very careful assessment of identification of confounding.

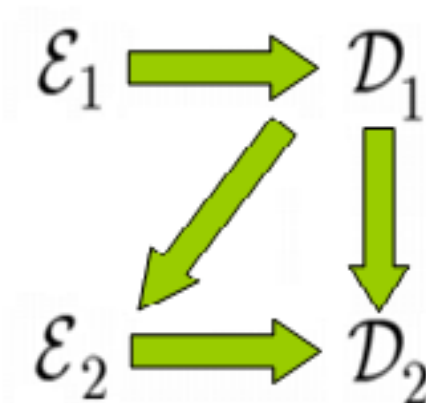
Since this lecture is merely an introduction to DAGs, simple examples will be presented here without much technical details. The motivation is to show how the theories of DAGs are applied in complicated scenarios of longitudinal measurement from observational study setting. Readers can refer to recent research such as Gran et al. [2010], Moodie and Stephens [2010], or textbooks such as Fitzmaurice [2009] (appendix §23.7), Diggle [2002] (chapter 12) for more elaborate examples and explanations. Followings are the simplified versions of those discussed in these references.

### 5.1 Scenario I

Let, for simplicity, there are two time points,  $t = 1$  and 2.  $\mathcal{E}_t$  is the treatment indicator (received treatment, or not) and  $\mathcal{D}_t$  is the disease outcome indicator (diseased, or not). In the two time points, measurements are recorded for both  $\mathcal{E}_t$  and  $\mathcal{D}_t$ .

Let us further assume that, after treatment status  $\mathcal{E}_1$  at time 1, the outcome  $\mathcal{D}_1$  at time 1 is measured. Subsequently, depending on  $\mathcal{D}_1$ , treatment status  $\mathcal{E}_2$  on time 2 was decided. Now, both disease status  $\mathcal{D}_1$  in the previous stage and treatment status  $\mathcal{E}_2$  in the second stage is associated with final disease status  $\mathcal{D}_2$  at time 2. Figure 5.1 shows a DAG to represent the time dependence situation.

**Figure 5.1:** Time dependent treatment and outcome for  $t = 1, 2$



Say, the aim is to determine the treatment effect of  $\mathcal{E}_1$  and  $\mathcal{E}_2$  on the final disease status  $\mathcal{D}_2$ . From the DAG, we can see that

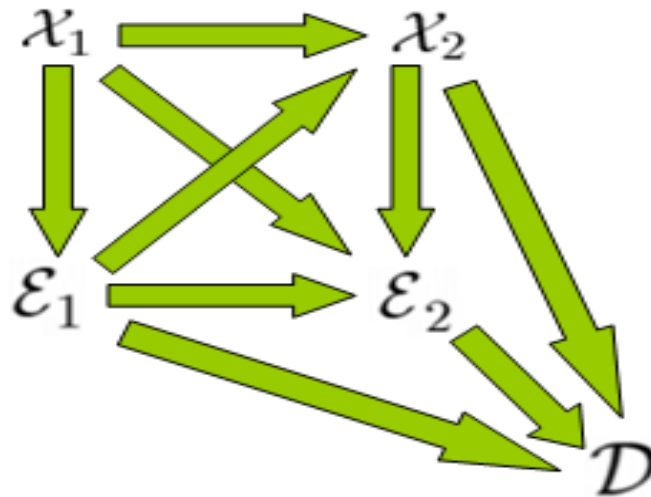
- with respect to treatment status  $\mathcal{E}_1$  at time 1, treatment status  $\mathcal{E}_2$  at time 2 is simply an intermediate variable, and causal effect of both can be simultaneously obtained.
- However, with respect to treatment status  $\mathcal{E}_2$  at time 2, a back-door is open at  $\mathcal{E}_2 \leftarrow \mathcal{D}_1 \rightarrow \mathcal{D}_2$ . Therefore, to find the treatment effect of  $\mathcal{E}_2$  on  $\mathcal{D}_2$ , confounding is

present here. Controlling  $\mathcal{D}_1$  will remove the effect of  $\mathcal{D}_1$  and  $\mathcal{E}_1$ , and causal pathway of  $\mathcal{E}_2$  to  $\mathcal{D}_2$  will be established.

## 5.2 Scenario II

In the observational survival analysis setting, to estimate the treatment effect, a Cox's proportional hazards model is very popular, especially if there is some delay in receiving treatment. On the top of that, if there exists time dependent confounder, the situation gets even more complicated and marginal structural model is very popular. To ensure comparability, some weighting scheme is used, such as inverse probability of treatment. However, in presence of time dependent confounders, the estimates of these weights sometime can very unstable. This is why, Gran et al. [2010] used a sequential Cox approach to estimate the causal effect of the treatment for HIV in the Swiss HIV cohort study.

**Figure 5.2:** Time dependent treatment and covariate for  $t = 1, 2$  and outcome



The idea was to split the data in a monthly basis, and assuming effect of confounder (say, CD4 cell count, which is associated with both treatment and outcome) being constant within such interval. This way, the analyst only has to deal with time dependent treatment status.

To simplify the idea, let us just consider two time points,  $t = 1, 2$ , and like the previous scenario, in the two time points, measurements are recorded for both treatment status  $\mathcal{E}_t$  (treatment received, or not at time  $t$ ) and  $\mathcal{X}_t$  (observed covariate at time  $t$ ).  $\mathcal{D}$  is the ultimate outcome (AIDS or death).

Here, we can see covariate in time 1,  $\mathcal{X}_1$  is associated with treatment status in time



1,  $\mathcal{E}_1$ . Covariate at time 2,  $\mathcal{X}_2$  depend on treatment status in time 1,  $\mathcal{E}_1$ . Then again,  $\mathcal{X}_2$  also depends on the covariate in the previous time point  $\mathcal{X}_1$ . This covariate at time 2,  $\mathcal{X}_2$  then affects treatment status in time 2,  $\mathcal{E}_2$ , which is also associated with  $\mathcal{X}_1$  and  $\mathcal{E}_1$ .  $\mathcal{E}_1$ ,  $\mathcal{E}_2$  and  $\mathcal{X}_2$  eventually influences the final outcome  $\mathcal{D}$ . Figure 5.2 shows a DAG to represent this scenario.

Now, we want to estimate the effect of the treatment status of time 1,  $\mathcal{E}_1$  on final disease status  $\mathcal{D}$  from the first interval or first month. From the causal diagram, we can see that there exists a back door from  $\mathcal{E}_1$  to  $\mathcal{D}$ , which is  $\mathcal{E}_1 \leftarrow \mathcal{X}_1 \rightarrow \mathcal{X}_2 \rightarrow \mathcal{D}$ . Then we can easily identify that conditioning on  $\mathcal{X}_1$  closes the back-door. We now have multiple paths to go from  $\mathcal{E}_1$  to  $\mathcal{D}$ , such as  $\mathcal{E}_1 \rightarrow \mathcal{E}_2 \rightarrow \mathcal{D}$ ,  $\mathcal{E}_1 \rightarrow \mathcal{X}_2 \rightarrow \mathcal{D}$  and  $\mathcal{E}_1 \rightarrow \mathcal{D}$ , and none of them are noncausal or non-directed. Since we are only interested in treatment effect of  $\mathcal{E}_1$  on  $\mathcal{D}$  (while dealing with first interval only), artificial censoring is imposed on those who starts treatment after some delay period.

In a more complicated scenario, we have multiple time points  $t = 1, 2, 3, 4, \dots$ , we subsequently deal with the corresponding time intervals, and average the treatment effect obtained in these intervals to report the aggregated treatment effect.

## 6 Conclusion

Causal inference is interrupted by many factors, such as confounding, missing data, selection bias, model specification, measurement error, etc. In this mini-lecture, only the issue of confounding was focused. Formulation of statistical theory to cover this source of bias has long been a challenge. Development of precise formulation, such as causal diagrams has opened a new door in identifying confounding situations. This mini-lecture was merely a gentle introduction to this approach. Those who are interested in pursuing further, might consider the references mentioned in the bibliographic notes in §.

## Appendix

### A Causation

The idea of causation is a very old one. Some definitions are provided below following historical hierarchy:

#### A.1 Philosophical Perspective

In simplest term, “cause” is that which produces an effect (outcome or result). Causes are conditions that have role in producing occurrence of effect. Philosopher Hume presented two set of definitions [Hume and Beauchamp, 2000]:

1. Cause is an event followed by another (effect), i.e., events labeled ‘cause’ are followed by the events labeled ‘effect’. That means the first event is sufficient for the occurrence of second event. The conditions of this assumption can be illustrated as follows: an event  $\mathcal{E}$  is considered cause of another event  $\mathcal{D}$ , if both  $\mathcal{E}$  and  $\mathcal{D}$  occurs, and the sequence is such that whenever  $\mathcal{E}$  occurs, then  $\mathcal{D}$  also occurs. Not only the cause  $\mathcal{E}$  occurs before effect  $\mathcal{D}$ , but the two events should be closely connected in terms of time or space.
2. Without the first event (cause), the second (effect) would never happen, i.e, if the first event does not happen, then the second event will cease to exist. That means the first event is necessary for the occurrence of second event. This definition states that an event  $\mathcal{E}$  is considered cause of another event  $\mathcal{D}$ , if  $\mathcal{E}$  had not occurred, then  $\mathcal{D}$  would not have occurred as well.

These two definitions imply that cause  $\mathcal{E}$  is a sufficient and necessary condition for  $\mathcal{D}$ .

In modern era, MacMahon and Pugh [1970] equates association with cause given the fact that, if the cause is altered, then the effect will also changed.

#### A.2 Counterfactuals

A definition of causation using counterfactual notations was provided by Neyman [1923]:

##### A.2.1 Causality and Counterfactual Models

The difference in counterfactual mechanism and the actual mechanism of events is that in the actual mechanism, cause  $\mathcal{E}$  and its effect  $\mathcal{D}$  happens, whereas in counterfactual mechanism, cause  $\mathcal{E}$  and its effect  $\mathcal{D}$  are absent. This gives us basis to find ‘the magnitude

of difference of the effect  $\mathcal{D}$  in presence and absence of  $\mathcal{E}$  (by comparing  $\mathcal{D}$  from actual mechanism to  $\mathcal{D}$  from counterfactual mechanism). This enables us to tell the impact of cause  $\mathcal{E}$ . Hence, concept of counterfactual is very relevant while making a causal statement.

### A.2.2 Relation with Philosophical Definition

This definition closely resembles the second definition of Hume, as discussed in §A.1. In the second definition, it is mentioned that if  $\mathcal{E}$  did not occur,  $\mathcal{D}$  would not occur. This is a condition that is counter to fact, since it is assumed what would happen in absence of  $\mathcal{E}$ , but not really observed. On the contrary, the effect is truly observed, once the cause is present. That is why this concept is called counterfactual.

### A.2.3 Probabilistic Causality

So far we have considered causality as deterministic. This would imply that whenever  $\mathcal{E}$  occurs,  $\mathcal{D}$  must also occur. Eells [1991] makes a compelling point in first chapter of his book using the example of “Smoking causes Lung Cancer” for extending causality from determinism, and advocates for population-level probabilistic causation in these situations. In this extension, occurrence of  $\mathcal{E}$  has an impact over the probability of occurrence of  $\mathcal{D}$ ,  $P(\mathcal{D})$ . Population-level probabilistic causation has little to say about individual-level causation, rather deals with frequencies in repeated sampling from a population. Therefore, causal theory involves techniques that are similar to those involved in explains relationship in probabilistic theory.

Hence, in an ideal experimental setting, we measure “causal effect of an exposure on the outcome” as follows:

1. we run an experiment at first on an experimental unit (a single patient) with the exposure  $\mathcal{E}$ , and record  $\mathcal{D}$ . Then under identical condition, we again run the experiment without the exposure  $\mathcal{E}$ , and record  $\bar{\mathcal{D}}$ . This means, the exposure condition was determined, before we observed the outcome.
2. We measure the observed difference in outcome under the two experimental conditions.
3. We repeat the same experiment on different experimental units under consideration, and record  $\mathcal{D}$ ,  $\bar{\mathcal{D}}$  under  $\mathcal{E}$ ,  $\bar{\mathcal{E}}$  respectively, holding all other factors fixed.
4. The average difference in response over all the experimental units under condition gives us the desired measure.

### A.2.4 General Methodology to estimate Causal Effect

The basic methodology in epidemiologic setting is described as follows with appropriate notations. Here the causes are called treatments and effects are, say disease status or outcome. For simplicity, the underlying assumption is that any treatment is possible to be assigned on any unit, and outcome of any unit only depends on the particular treatment that was applied in that unit.

**Treatment** Suppose that a vector of  $J + 1$  treatments is  $\mathbf{x} = (x_0, x_1, \dots, x_J)'$  applied at one point in time. The unit  $i$  receives exactly one of the  $J + 1$  treatments. Also, the treatment  $x_0$  is the reference treatment. This can be either placebo, or any standard treatment. The rationale of having such treatment is to evaluate other treatment's effect compared to this one.

**Potential outcome** For unit  $i$ , the corresponding outcome vector is  $\mathbf{y}_i = (y_{0i}, y_{1i}, \dots, y_{Ji})'$ . This vector  $y_i$  is actually potential outcome, whereas, in actuality, outcome  $Y_i = y_{ji}$  is  $x_j$  treatment is administered to unit  $i$ .

**Counterfactual** If we assign  $x_i$  and record outcome, and then again assign  $x_k$  in that particular unit, then outcome will be due to both of  $x_i$  and  $x_k$ , and separating the effects will not be straightforward. That is why, we use this counterfactual model, where only one treatment is assigned on one unit. Only  $y_{ji}$  will be observed from the vector after applying  $x_i$ , whereas the remainder will be unobserved (contrary-to-fact).

**Homogeneity Assumption** To get the outcome corresponding to all the treatments, we assume that each unit is homogeneous and causal effect of a particular treatment will be the same on any unit. However, in an individual level, such assumption is not realistic. However, on a group level, the average causal effect of two groups from the same population should be the same. Therefore, we have different homogenous sub-groups and we apply different treatments on those homogenous subgroups.

**Causal Effect** We define causal effects in terms of counterfactuals. We say that the receiving treatment  $x_j$  instead of  $x_0$  would have an effect on  $Y_i$  if  $y_{ji} \neq y_{0i}$ . If we want to estimate the 'causal effect',  $C.E. = Y_{ji} - Y_{0i}$ ,

- We say that  $y_{ji} - y_{0i}$  is the effect of  $Y_i$  of receiving  $x_i$  instead of  $x_0$ .
- If the counterfactuals are strictly positive,  $y_{ji}/y_{0i}$  or  $\log(y_{ji}) - \log(y_{0i})$  also measures the causal effect.

Since this is defined on an individual level, in the next step we generalize it to group level.

**Inference** Inference about population effects can be obtained as differences in average population response under a particular treatment from that of reference treatment. That is, the ‘Average Causal Effect’,  $A.C.E. = E(Y_{ji}) - E(Y_{0i})$  for treatment  $j$  could be estimated from  $n$  population units as:  $\sum_{i=1}^n (y_{ji} - y_{0i})/n = \sum_{i=1}^n y_{ji}/n - \sum_{i=1}^n y_{0i}/n = \bar{y}_j - \bar{y}_0$ . The assumption would be that treatment assignment is independent of outcome, and randomization prevails.

**Randomization Assumption** If randomization fails (if the treatment assignment follows a systematic pattern, e.g, if almost all subjects with higher age gets treatment  $j$ , and almost all subjects with lower age gets treatment  $m$ , so that there is a covariate imbalance in the comparison groups), this estimate may be biased. Allocating treatment using a random mechanism may avoid this pitfall. If random mechanism is followed in a large sample, it is less probable to systematically favoring higher or lower values of covariates, hence the covariate balance could be established in the two comparison groups.

**Observational Studies** The formulation is appropriate for randomized trials. For observational studies, we need to have more assumptions, such as ignorability, which will be discussed in §B.5.

**Criticism** The primary objection is that it deals with unobserved objects, hence making assumptions empirically untestable.

### A.2.5 Illustration

For the purpose of illustration, one example is provided here, where we only have two comparison groups: people drinking coffee and people not drinking coffee. Let that we are interested about pancreatic cancer, and we want to see whether drinking coffee has any causal effect on it or not. Therefore, we organize our experiment as follows:

**Exposed Group** Let us select a sample, and expose each subject in the sample (of size  $n$ ) to coffee drinking, and follow them for appropriate amount of time, and determine whether they developed pancreatic cancer or not. Note that, the coffee exposure happened before the pancreatic cancer was observed. This temporal ordering is absolutely necessary for making causal inference.

**Unexposed Group** Then, hypothetically, go back in time and make each subject coffee abstainer, follow them for same amount of time and determine whether they developed pancreatic cancer or not. These would be counterfactuals. Logically, we can not really do so. Instead, ideally we choose subjects with identical conditions and

run the experiment without the exposure coffee to get the results from the second group.

**Compare Results** The possible results are provided in Table 1. From this table, group *I* develops pancreatic cancer irrespective of their coffee drinking habit. Number of subjects in this sub-group is  $n \pi_{11}$ . For group *IV*, its the opposite. Number of subjects in this sub-group is  $n \pi_{00}$ . Differences due to coffee drinking is really visible in group *III* and *IV*, with frequency  $n \pi_{10}$  and  $n \pi_{01}$  respectively.

Therefore, assuming these proportions are based on a large sample, we have,

**Table 1:** Possible hypothetical result

Sub-sample	$\mathcal{E}$	$\bar{\mathcal{E}}$	Proportion
<i>I</i>	$\mathcal{D}$	$\bar{\mathcal{D}}$	$\pi_{11}$
<i>II</i>	$\mathcal{D}$	$\bar{\mathcal{D}}$	$\pi_{10}$
<i>III</i>	$\bar{\mathcal{D}}$	$\mathcal{D}$	$\pi_{01}$
<i>IV</i>	$\bar{\mathcal{D}}$	$\bar{\mathcal{D}}$	$\pi_{00}$

$$P(\mathcal{D}|\mathcal{E}) = \pi_{11} + \pi_{10}$$

$$P(\mathcal{D}|\bar{\mathcal{E}}) = \pi_{11} + \pi_{01}$$

Using these probabilities, we can calculate causal effects as follows:

$$\text{Risk Difference} = P(\mathcal{D}|\mathcal{E}) - P(\mathcal{D}|\bar{\mathcal{E}}) = \pi_{11} + \pi_{10} - (\pi_{11} + \pi_{01}) \quad (\text{A.1})$$

$$\text{Risk Ratio} = \frac{P(\mathcal{D}|\mathcal{E})}{P(\mathcal{D}|\bar{\mathcal{E}})} = \frac{\pi_{11} + \pi_{10}}{\pi_{11} + \pi_{01}} \quad (\text{A.2})$$

$$\text{Odds Ratio} = \frac{\frac{P(\mathcal{D}|\mathcal{E})}{1-P(\mathcal{D}|\mathcal{E})}}{\frac{P(\mathcal{D}|\bar{\mathcal{E}})}{1-P(\mathcal{D}|\bar{\mathcal{E}})}} = \frac{\frac{\pi_{11}+\pi_{10}}{1-(\pi_{11}+\pi_{10})}}{\frac{\pi_{11}+\pi_{01}}{1-(\pi_{11}+\pi_{01})}} \quad (\text{A.3})$$

Here, the null hypothesis is  $H_0 : \pi_{01} = \pi_{10}$ . When we can not reject the  $H_0$ , that does not mean that coffee has no causal effect on the sample; this simply means that drinking coffee does not change population risk of pancreatic cancer.

**Assumption** We need to make sure about assumptions such as randomization, exchangeability, etc. Hypothetically, without these assumptions, its is possible that while sample size is small, some cells in Table 1 might be empty, which will in turn cause difficulty in formulating equations (A.1), (A.2), (A.3).

### A.2.6 Confounding in terms of Counterfactuals

An excellent way to explain confounding is to illustrate it in terms of a counterfactual model.

**Notations** Let us say that the target population is  $\Omega$ . Let us denote

- the treatment of interest as  $x_1$ , corresponding outcome as  $y_1$ . The marginal distribution of outcome is  $F_{\Omega}(y_1)$ . When  $x_1$  treatment is assigned, the summary parameter  $\mu$  is equal to  $\mu_{\Omega 1}$ .
- the reference treatment as  $x_0$ , corresponding outcome as  $y_0$ . The marginal distribution of outcome is  $F_{\Omega}(y_0)$ . When  $x_0$  treatment is assigned, the summary parameter  $\mu$  is equal to  $\mu_{\Omega 0}$ .

The same is true under another population  $\bar{\Omega}$ . We say that  $\bar{\Omega}$  is the reference population and we assume that  $\bar{\Omega}$  will be similar to  $\Omega$  (i.e., all the conditions are identical). All the items are shown in Table 2.

**Table 2:** Counterfactual Items

Population	$\Omega$		$\bar{\Omega}$	
Treatment	$x_1$	$x_0$	$x_1$	$x_0$
Outcome	$y_1$	$y_0$	$y_1$	$y_0$
Marginal Distribution	$F_{\Omega}(y_1)$	$F_{\Omega}(y_0)$	$F_{\bar{\Omega}}(y_1)$	$F_{\bar{\Omega}}(y_0)$
Summary parameter	$\mu_{\Omega 1}$	$\mu_{\Omega 0}$	$\mu_{\bar{\Omega} 1}$	$\mu_{\bar{\Omega} 0}$

**Objective** The objective is to determine the ‘effect’ of applying the exposure  $x_1$  on a summary parameter  $\mu$  of the distribution of outcome  $y$  in the population  $\Omega$ , compared to applying another exposure  $x_0$ . Therefore, one way to measure the causal effect of  $x_1$  relative to  $x_0$  should  $\mu_{\Omega 1} - \mu_{\Omega 0}$ .

**Difficulty in Practice** In reality, if we assign treatment  $x_1$  in the population  $\Omega$ , only  $\mu_{\Omega 1}$  will be observed, whereas  $\mu_{\Omega 0}$  will remain unobserved. Similarly, if we assign treatment  $x_0$  in the population  $\bar{\Omega}$ , only  $\mu_{\bar{\Omega} 0}$  will be observed, whereas  $\mu_{\bar{\Omega} 1}$  will remain unobserved. The observed objects are shown in Table 3. Therefore, the measure of treatment effect of  $x_1$  instead of  $x_0$  will be  $\mu_{\Omega 1} - \mu_{\bar{\Omega} 0}$ , where we replace  $\mu_{\Omega 0}$  with  $\mu_{\bar{\Omega} 0}$ .

### Defining Confounding

**Table 3:** Observed Items

Population	$\Omega$		$\bar{\Omega}$	
Treatment	$x_1$	$\times$	$\times$	$x_0$
Outcome	$y_1$	$\times$	$\times$	$y_0$
Marginal Distribution	$F_{\Omega}(y_1)$	$\times$	$\times$	$F_{\bar{\Omega}}(y_0)$
Summary parameter	$\mu_{\Omega 1}$	$\times$	$\times$	$\mu_{\bar{\Omega} 0}$

**Association**  $\mu_{\Omega 1} - \mu_{\bar{\Omega} 0}$  measures the association of treatments with outcomes across the populations.

**Causal**  $\mu_{\Omega 1} - \mu_{\Omega 0}$  measures the effect of treatment  $x_1$  in the population  $\Omega$ .

The essence of the whole analysis rests on the fundamental assumption that the separate unexposed group is comparable to the exposure group. From this formulation, confounding is said to exist if such assumption of comparability is void, i.e., when the risk in the unexposed group is different than that of the exposed group. In summary, if confounding exists, then association measure is a biased estimate of causal effect. That is, if  $\mu_{\Omega 0} \neq \mu_{\bar{\Omega} 0}$ , then the confounding exists.

### Measures of Confounding

1. One way to measure confounding is  $(\mu_{\Omega 1} - \mu_{\bar{\Omega} 0}) - (\mu_{\Omega 1} - \mu_{\Omega 0}) = \mu_{\Omega 0} - \mu_{\bar{\Omega} 0}$
2. When the outcome parameters are probabilities, an analogous measure is

$$\text{Confounding Risk Ratio} = \frac{\frac{\mu_{\Omega 1}}{\mu_{\Omega 0}}}{\frac{\mu_{\bar{\Omega} 1}}{\mu_{\bar{\Omega} 0}}} = \frac{\mu_{\Omega 0}}{\mu_{\bar{\Omega} 0}}$$

### Implications of this Formulation

**Outcome Parameter** Confounding depends on the outcome parameter. If  $\mu$  denotes 100 year survival probability due to a treatment, instead of 5 year survival probability, then we expect no confounding for that parameter.

**Target Population** Confounding depends on the target population. If  $\Omega$  is the target population, then  $(\mu_{\Omega 1} - \mu_{\bar{\Omega} 0})$  can be confounded for the effect  $(\mu_{\Omega 1} - \mu_{\Omega 0})$  iff  $\mu_{\Omega 0} \neq \mu_{\bar{\Omega} 0}$ . However, if  $\bar{\Omega}$  is the target population, then  $(\mu_{\Omega 1} - \mu_{\bar{\Omega} 0})$  can be confounded for the effect  $(\mu_{\bar{\Omega} 1} - \mu_{\bar{\Omega} 0})$  iff  $\mu_{\Omega 1} \neq \mu_{\bar{\Omega} 1}$ .

**Exception** Randomization does not exclude the possibility of confounding in a given experiment. If for any particular reason, say for small sample size, the covariate bal-



ance between the comparison groups are not established, then this kind of situation may occur.

**Illustration** Continuing the example of §A.2.5, say, male subjects in the sample are more likely to drink coffee than the female subjects. This would mean that we are more likely to find a male coffee drinker rather than a female coffee drinker, and hence the covariate balance is distorted. This also violates randomization assumption, since under randomization of treatment assignment, there should not be any bias or preference for male subject over the female subjects. Therefore, we can see that exposure to coffee,  $\mathcal{E}$  depends on gender,  $\mathcal{C}$ , and consequently, gender is also associated with pancreatic cancer,  $\mathcal{D}$ . But gender is neither determined by coffee drinking habit, nor the pancreatic cancer. This is exactly the textbook definition of confounding. Hence  $\mathcal{C}$  is a confounder.

## B Confounding and Relevant Concepts

### B.1 Conventional Definition of Confounding

The word ‘confounding’ originates from a Latin word ‘confundere’ – which means ‘mix together’. In a study aiming to explore association between an exposure (risk factor or cause) and the outcome (disease occurrence), confounding can occur when another exposure (extraneous to study variables) exists, which is associated with both the exposure and disease under consideration. Confounding occurs when the effects of two exposures have not been identified and separated, and we carry out our analysis based on one exposure only, excluding the other. When this extraneous factor is unequally distributed between the exposure subgroups, the estimates of effect in general is biased.

Three basic characteristics of a confounder is as follows [Szklo and Nieto, 2007]:

**Outcome Variable** The confounding variable must be causally related with the outcome variable.

**Exposure Variable** The confounding variable is associated with the exposure.

**External Variable** The confounding variable is not the intermediate variable. This characteristic requires it to be an external variable from the causal pathway between outcome and exposure, hence confounding variable must not be affected by exposure or outcome.

Rothman et al. [2008] defined confounding as a distortion in the estimated exposure effect that is due to difference in risk between exposed and unexposed that can not be

attributed to exposure.

## B.2 Residual Confounding

Residual confounding occurs in the following situations

- the categories of the confounder controlled for are too broad, resulting in imperfect adjustment (smaller categories are preferred for this reason),
- when some confounding variables remain unaccounted for, after some confounding variables have been controlled,
- when the confounder is mismeasured or misclassified. Like any other variable, confounder can be subject to misclassification. Misclassification of confounder makes it hard to control for that confounder.

## B.3 Over-adjustment

Over-adjustment should be avoided as well. When too many potential confounders are adjusted simultaneously, the main results become obscure. Even simple statistical tools, like logistic regression, are equipped to control for a lot of covariates, but the result out of this analysis will be less transparent to the audience, and more prone to errors that are not easily detectable. This, it is advised to control for confounder parsimoniously.

## B.4 Sufficient Set of Confounders

A set of variables are defined as a ‘sufficient’ for control of confounding if simultaneous stratification on all the variables are sufficient for estimating stratum-specific causal effects. Backdoor criterion can be used to identify this set. However, there is always possibility of such set being insufficient due to lack of measurement of an important variable. Sensitivity analysis should be performed to detect such unmeasured covariates.

## B.5 Ignorability

To obtain the unbiased estimate of causal effect in absence of randomization, the assumption of ignorability is necessary. The idea is that if we can find the set of covariates  $X$ , such that given those  $X$ , treatment assignment  $Z$  is independent of outcome  $Y$ , i.e.,  $Y \perp Z|X$ , then it is easy to show that  $E(Y_j|X = x, Z = j) = E(Y_j|X = x)$ , so that

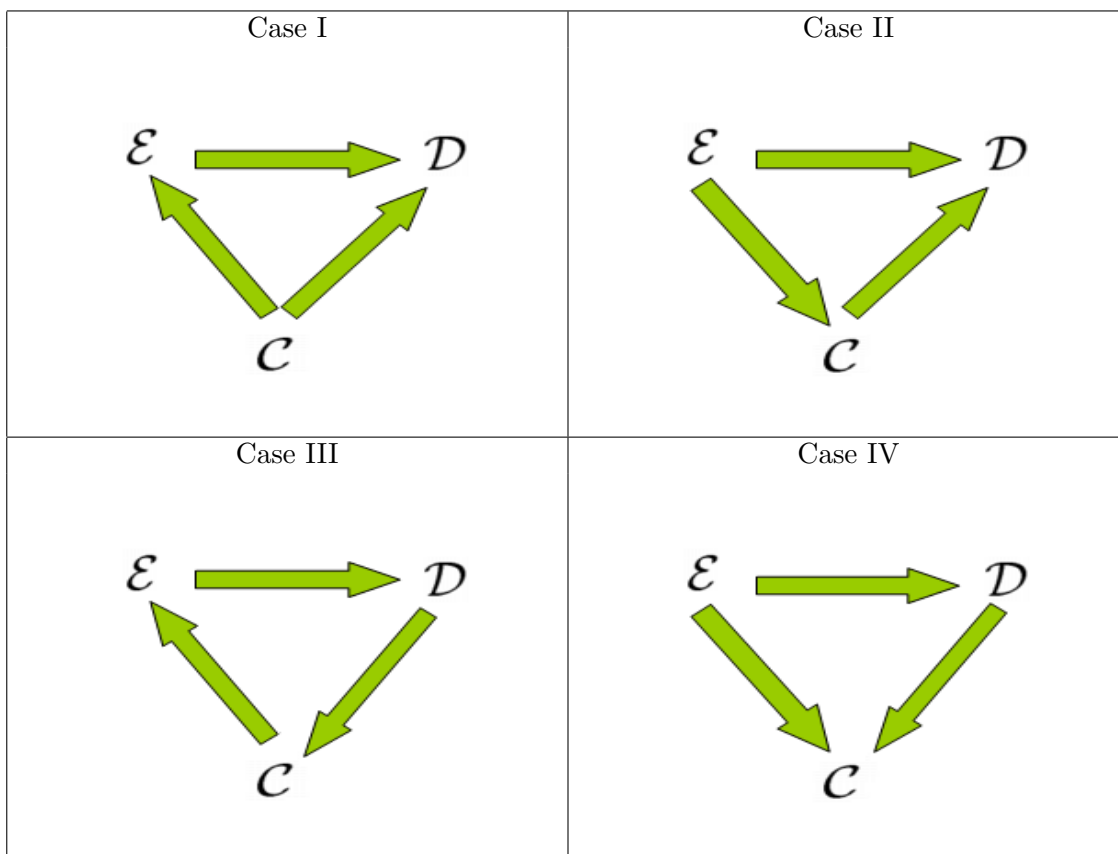
$$\begin{aligned} E(A.C.E.) &= \sum_x \left( E(Y_j|X = x) - E(Y_k|X = x) \right) P(X = x) \\ &= E(Y_j) - E(Y_k) \end{aligned}$$

Notice that, since  $y_j$  and  $y_k$  both can not be observed, this is another example of counterfactual concept.

## B.6 Confounding and Collapsibility

So far we have been introduced with the idea of stratification, which is used as a tool to control for confounding. However, stratification also introduces the possibility of collapsibility. By definition, noncollapsibility of a measure of association implies that the measure changes upon stratification by covariates. That is, different measure of associations results are obtained if the covariate is collapsed or ignored. The idea of collapsibility can easily be illustrated using Simpson's Paradox.

**Figure B.1:** Possible Cases of Relationship among exposure  $\mathcal{E}$ , outcome variable  $\mathcal{D}$  and an external variable  $\mathcal{C}$



If we use sufficient set of variables to control for confounding, then the problem of noncollapsibility will not occur as well. That means, the conditions of confounding and noncollapsibility are identical when the controlled covariates form a *sufficient set of control* (discussed in §B.4). This is the situation where confounding and noncollapsibility are equivalent, and this explains why these two different concepts are often not distinguished.

In terms of causal graphs, the idea of collapsibility does not distinguish between any cases shown in Figure B.1 since collapsibility is non-directional. Therefore, for case II, confounding and noncollapsibility will be different. The only case when they will be equivalent, is case I, where  $\mathcal{C}$  causes both  $\mathcal{E}$  and  $\mathcal{D}$ . Just like confounding situation, collapsibility of tables might be misleading, while conditioned by a collider i.e., case IV. Notice that case III is not a DAG, and hence causal graph cannot handle this case.

## C Bibliographic Notes

Gentle introductions to causal diagrams in context of epidemiology can be found in chapter 12 of Rothman et al. [2008], chapter 5 of Szklo and Nieto [2007] and chapter 8 of Jewell [2003]. From a longitudinal data analysis point of view, some might find chapter 9 and 10 of Gelman et al. [2007] helpful. For advanced causal graph theory, one could consult Pearl [2000], Oakes and Kaufman [2006] and Edwards [2000].

The ideas of causation is discussed in Parascandola and Weed [2001] and Fang [2003]. Neyman [1923] and Rubin [1974] discussed causal model in terms of counterfactuals. The basic ideas of confounding is discussed in Hume and Beauchamp [2000], MacMahon and Pugh [1970], Newman [2001] and Bonita et al. [2006].

A number of examples are provided in Fang [2003], Jewell [2003], and Greenland et al. [1999], some of which were re-used and extended in this lecture for illustrative purposes.

## References

- Bonita, R., Beaglehole, R., and Kjellström, T. (2006). *Basic epidemiology*. WHO.
- Diggle, P. (2002). *Analysis of longitudinal data*. Oxford University Press, USA.
- Edwards, D. (2000). *Introduction to graphical modelling*. Springer Verlag.
- Eells, E. (1991). *Probabilistic Causality*. Cambridge: Cambridge University Press.
- Fang, J. (2003). *Advanced Medical Statistics*. World Scientific Publishing Company.
- Fitzmaurice, G. (2009). *Longitudinal data analysis*. Chapman & Hall/CRC.
- Gelman, A., Hill, J., and Corporation, E. (2007). *Data analysis using regression and multilevel/hierarchical models*, volume 625. Cambridge University Press Cambridge.

- Gran, J., Røysland, K., Wolbers, M., Didelez, V., Sterne, J., Ledergerber, B., Furrer, H., von Wyl, V., and Aalen, O. (2010). A sequential Cox approach for estimating the causal effect of treatment in the presence of time-dependent confounding applied to data from the Swiss HIV Cohort Study. *Statistics in Medicine*.
- Greenland, S., Robins, J., and Pearl, J. (1999). Confounding and collapsibility in causal inference. *Statistical Science*, 14(1):29–46.
- Hume, D. and Beauchamp, T. (2000). *An enquiry concerning human understanding: a critical edition*. Oxford University Press, USA.
- Jewell, N. (2003). *Statistics for epidemiology*. CRC Press.
- MacMahon, B. and Pugh, T. (1970). *Epidemiology: principles and methods*. Little, Brown Boston.
- Moodie, E. and Stephens, D. (2010). Using Directed Acyclic Graphs to detect limitations of traditional regression in longitudinal studies. *International Journal of Public Health*, pages 1–3.
- Newman, S. (2001). *Biostatistical methods in epidemiology*. Wiley-Interscience.
- Neyman, J. (1923). On the application of probability theory to agricultural experiments: Essay on principles, section 9.(translated in 1990). *Statistical Science*, 5:465–480.
- Oakes, J. and Kaufman, J. (2006). *Methods in social epidemiology*. Jossey-Bass Inc Pub.
- Parascandola, M. and Weed, D. (2001). Causation in epidemiology. *British Medical Journal*, 55(12):905.
- Pearl, J. (2000). *Causality: Models, reasoning, and inference*. New York: Cambridge University Press.
- Rothman, K., Greenland, S., and Lash, T. (2008). *Modern epidemiology*. Lippincott Williams & Wilkins.
- Rubin, D. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688–701.
- Szklo, M. and Nieto, F. (2007). *Epidemiology: beyond the basics*. Jones & Bartlett Publishers.