

Chapter 3

Two key questions

We can roughly categorise the problems studied within causal inference into two groups, causal effect estimation and causal discovery. In causal effect estimation we assume (at least implicitly) that key aspects of the causal graph are known. The goal is then to estimate the effect of an intervention or range of interventions in the system. Causal effect estimation is implicit in countless studies in economics, social science and epidemiology of everything from the effect of education on earnings [26], diet on cancer [17] and breastfeeding on intelligence [66] to the effect of pet ownership on survival after a heart attack [43]. Almost every time someone runs a regression model the key quantity of interest is a causal effect. Given how it underlies so much of our scientific progress, there is a enormous potential in properly understanding when we can draw causal conclusions, exactly what assumptions are required to do so and how we can best leverage those assumptions to infer as much information as we can from our data.

Causal discovery aims to leverage much broader assumptions to learn the structure of causal graphs from data. This is critical in fields where we are generating a lot of data but have limited theoretical knowledge from which to draw on to determine how variables are related to one another. Causal discovery algorithms are being applied in bioinformatics [15, 107, 94, 4, 121, 47, 117, 123], medical imaging [95] and climate science [125]. An effective and generalisable approach for causal discovery would be a major step towards the automation of the scientific endeavour.

3.1 Causal effect estimation

Estimating causal effects from observational data comes down to determining if and how we can write expressions for the interventional distributions of interest in terms of observational quantities, which can be measured. We did this on an ad-hoc basis to resolve the examples discussed in chapter 2. In this chapter we describe a principled approach to mapping observational quantities to interventional ones and discuss some of the key issues involved in estimating such expressions from finite sample data. We assume the basic structure of the graph is known. That is, we assume that we can draw a network containing (at a minimum):

- the target/outcome variable we care about,
- the focus/treatment variables on which we are considering interventions,
- any variables which act to confound two or more of the other variables we have included, and
- any links between variables we have included.

Some of these variables may be latent in that the available data does not record their value, however their position in the network is assumed to be known. For example, consider estimating the impact of schooling on wages. Some measure of inherent ability could influence both the number of years of schooling people choose to pursue and the wages they receive. Even if we have no data to directly assess people's inherent ability we must include it in the graph because it influences two of the variables we are modelling.

How can the structure of the causal graph be leveraged to compute interventional distributions from observational ones? Given the graph corresponding to the observational distribution, the graph after any intervention can be obtained by removing any links into variables directly set by the intervention. The joint interventional distribution is the product of the factors associated with the interventional graph, as given by the truncated product formula 2.3. If there are no latent variables the interventional distribution of interest can be obtained by marginalising over the joint (interventional) distribution. However, if there are latent variables the joint interventional distribution will contain terms that cannot be estimated from the observed data.

The key to estimating causal effects in the presence of latent variables lies in combining the assumption of how an intervention changes the graph, encoded by the truncated product formula, with information the graph structure provides about conditional independencies between variables. By leveraging conditional independencies we can effectively localise the effect of an intervention to a specific part of a larger graph. This gives rise to the do-calculus [87]. The do calculus consists of three rules. They are derived from the causal information encoded in a causal network and the properties of d-separation and do not require any additional assumptions other than that of specifying the causal network.

3.1.1 Independence in Bayesian networks: D-separation

Many causal algorithms are based on leveraging the independence properties encoded in Bayesian networks. Therefore, in this section, we briefly review the key properties of Bayesian networks. A more thorough introduction (including proofs) can be found in [74]. Recall that a Bayesian network is a way of representing the joint distribution over its variables in a way which highlights conditional independencies between them.

Theorem 6. (*Local Markov condition*) *Given a Bayesian network G with nodes $X_1 \dots X_N$, each variable X_i is independent of its non-decedents given its parents in G for all distributions $P(X_1 \dots X_N)$ that are compatible with G .*

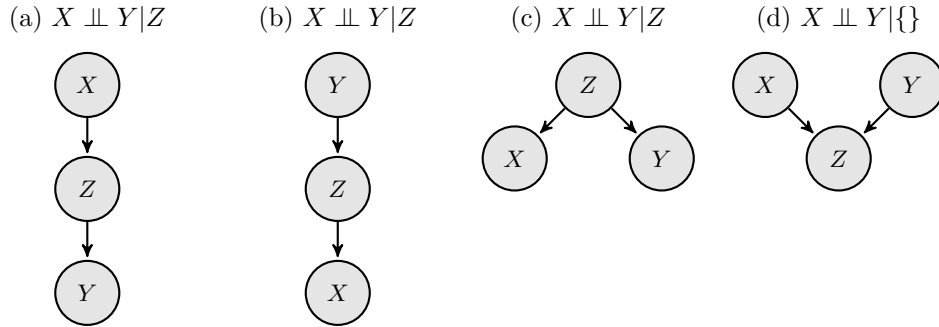
The set of conditional independence relations given by the local Markov condition can enforce additional independencies that also hold in all distributions that are compatible with G . D-separation is an algorithm that extends the local Markov property to find these additional independencies. It provides us with a simple way of reading from a network if a given conditional independence statement is true in all distributions compatible with that network.

The statement that X is conditionally independent of Y given Z implies that if we know Z , learning the value of Y gives us no additional information about X . From a graphical perspective you can think of this as Z blocks the flow of information from X to Y in the network. Figure 3.1 shows all possible network paths from a variable X to Y via Z . In figures (a) to (c) the path is blocked if we condition on Z and unblocked otherwise. In figure (d) the path is unblocked if we condition on Z and blocked otherwise.

The structure in figure 3.1d is referred to as a collider or v-structure. The somewhat counter-intuitive result that conditioning on Z introduces dependence between X and Y is called the *explaining away phenomena*. As an example, consider a scholarship available to female or disad-

vantaged students. Let X be gender, Y be family background and Z receipt of the scholarship. There are roughly equal numbers of boys and girls in both poor and wealthy families so X and Y are independent. However, if we know a student is receiving a scholarship then learning that they are male increases the probability that they are disadvantaged.

Figure 3.1: All possible two edge paths from X to Y via Z



Definition 7 (unblocked path). A path from X to Y is a sequence of edges linking adjacent nodes starting at X and finishing at Y , $(X, V_1, V_2 \dots V_k, Y)$. It is unblocked if every triple, $X - V_1 - V_2, V_1 - V_2 - V_3, \dots, V_{k-1} - V_k - Y$ in the path is unblocked (each triple will belong to one of the cases in figure 3.1)

Definition 8 (d-separation). The variables \mathbf{X} are d-separated from \mathbf{Y} given \mathbf{Z} in the network G if, there are no unblocked paths from any $X \in \mathbf{X}$ to any $Y \in \mathbf{Y}$ after conditioning on \mathbf{Z} .

Theorem 9 (d-separation and conditional independence). *If a set of variables \mathbf{Z} d-separates \mathbf{X} and \mathbf{Y} in a Bayesian network G then $(\mathbf{X} \perp\!\!\!\perp \mathbf{Y}|\mathbf{Z})$ in all distributions P compatible with G . Conversely, if \mathbf{X} and \mathbf{Y} are d-connected (not d-separated) given \mathbf{Z} then it is possible to construct a distribution P' that factorises over G in which they are dependent.*

Theorem 9 says that independencies implied by d-separation on a graph hold in every distribution that can be factored over that graph and that if $(\mathbf{X} \perp\!\!\!\perp \mathbf{Y}|\mathbf{Z})$ in *all* distributions that can be factored over G then they are d-separated in G . If we denote the independencies implied by d-separation in a graph by $\mathcal{I}(G)$ and the set of independencies in a distribution by $\mathcal{I}(P)$ then $\mathcal{I}(G) \subseteq \mathcal{I}(P)$.

If $\mathcal{I}(G) = \mathcal{I}(P)$ then G is called a perfect map for P . However, it is possible to construct distributions that do not have a perfect map, that is they contain conditional independencies that cannot be represented by d-separation. A particular case in which this occurs is when there are deterministic relationships between variables. If we have a Bayesian network G in which we specify that some nodes are deterministic we cannot conclude that if X and Y are d-connected then there exists a distribution P' consistent with G in which they are dependent. This does not conflict with theorem 9 as *consistent* in this setting requires that P' both factorises over G and satisfies the specified the deterministic relations between variables. This subtlety led to confusion in assessing what independencies hold between counterfactuals via twin networks [87, 97] and demonstrates the caution required in using d-connecteness to assert lack of independence. D-separation can be extended to compute the additional independencies implied by a graph in which certain nodes are known to be deterministic [48].

3.2 The Do Calculus

The do-calculus is a set of three rules [86] that can be applied to simplify the expression for an interventional distribution. If by repeated application of the do-calculus, along with standard probability transformations, we can obtain an expression containing only observational quantities then we can use it to estimate the interventional distribution from observational data. Let $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$ and \mathbf{W} be disjoint sets of variables in a causal graph G . We denote the graph G after the an intervention $do(\mathbf{X})$, which has the effect of removing all edges into variables in \mathbf{X} , as $G_{\overline{\mathbf{X}}}$

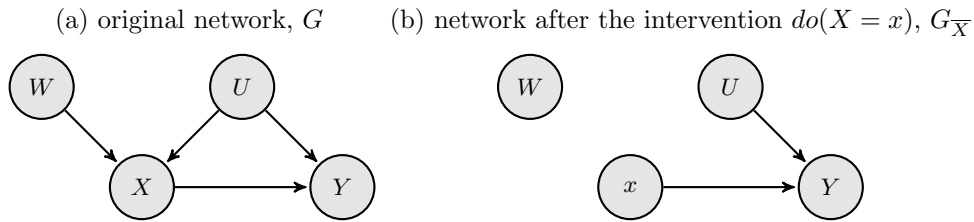
Rule 1: (adding or removing evidence)

Rule 1 allows us to remove (or insert) observational evidence from the right hand side of a conditional interventional distribution. It follows directly from the fact that the relationship between d-separation in a network and independence in the corresponding probability distribution still applies after an intervention.

If $(\mathbf{Y} \perp\!\!\!\perp \mathbf{W} | \mathbf{Z}, \mathbf{X})$ in $G_{\overline{\mathbf{X}}}$:

$$P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{Z} = \mathbf{z}, \mathbf{W} = \mathbf{w}) = P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{Z} = \mathbf{z}) \quad (3.1)$$

Figure 3.2: Rule 1 example. $(Y \perp\!\!\!\perp W | X)$ in $G_{\overline{X}} \implies P(Y | do(X), W) = P(Y | do(X))$



Rule 2: (exchanging actions with observations)

Rule 2 describes when conditioning on $\mathbf{X} = \mathbf{x}$ and intervening, $do(\mathbf{X} = \mathbf{x})$, have the same effect on the distribution over \mathbf{Y} . Let $G_{\underline{\mathbf{X}}}$ denote the causal graph G with all edges *leaving* \mathbf{X} removed.

If $(\mathbf{Y} \perp\!\!\!\perp \mathbf{X} | \mathbf{W})$ in $G_{\underline{\mathbf{X}}}$:

$$P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{W}) = P(\mathbf{Y} | \mathbf{X} = \mathbf{x}, \mathbf{W}) \quad (3.2)$$

The intuition behind this is that interventional distributions differ from observational ones due to the presence of indirect paths between X and Y . Observing a variable X provides information about Y both directly and indirectly, by changing our belief about the distribution of the parents of X . However setting X tells us nothing about its parents and therefor affects Y only via direct paths out of X . Removing edges *leaving* X removes all the direct paths out of X . If X is then independent of Y (conditional on W), that indicates there are no indirect paths. This implies conditioning on X is equivalent to setting X (given W).

Equation 3.2 does not cover cases where acting on one set of variables allows us to replace acting on another set with conditioning (see figure 3.4). The general form of rule 2 is given in equation 3.3.

If $(Y \perp\!\!\!\perp X | W, Z)$ in $G_{\underline{X}\bar{Z}}$:

$$P(Y | do(X = x), do(Z = z), W) = P(Y | X = x, do(Z = z), W) \quad (3.3)$$

Figure 3.3: An example of rule 2 with a single intervention $(Y \perp\!\!\!\perp X | W)$ in $G_{\underline{X}}$ $\implies P(Y | do(X), W) = P(Y | X, W)$. In this example, observing X provides information about Y both directly and indirectly, because knowing X tells us something about W which also influences Y . If we condition on W , we block this indirect path.

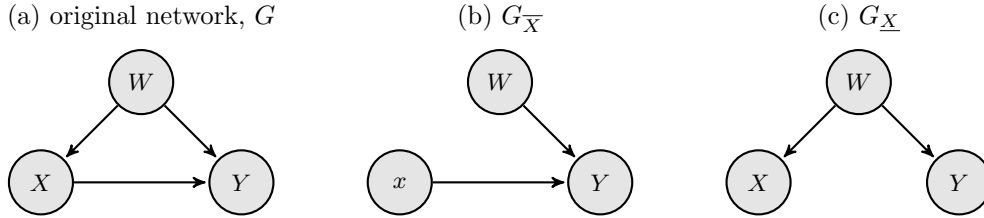
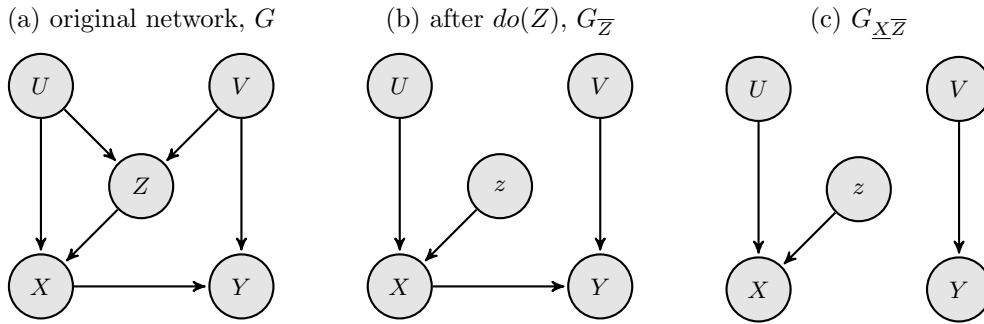


Figure 3.4: An example of applying equation 3.3. In this case $(Y \perp\!\!\!\perp X | Z)$ in $G_{\underline{X}\bar{Z}}$ $\implies P(Y | do(X = x), do(Z = z)) = P(Y | X = x, do(Z = z))$. Observing, rather than intervening, on Z would not have allowed us to exchange $do(X = x)$ for $X = x$. Conditioning on Z does block the indirect path $X - Z - V - Y$ but opens $X - U - Z - V - Y$.



Rule 3: (adding or removing actions)

This rule describes cases where the intervention $do(X = x)$ has no effect on the distribution of the outcome Y . A simple case of rule 3 is given in equation 3.4. If Y is independent of X in G after removing links entering X then can be no direct path from X to Y and any intervention on X will not affect Y .

if $(Y \perp\!\!\!\perp X)$ in $G_{\bar{X}}$:

$$P(Y | do(X = x)) = P(Y) \quad (3.4)$$

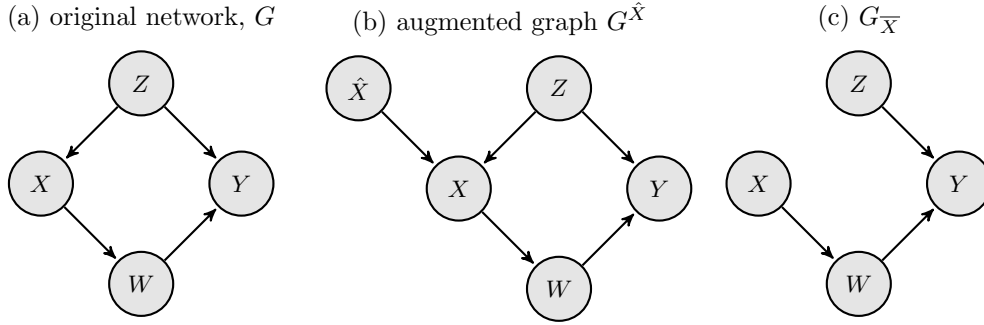
The general case of rule 3 is easier to state by explicitly representing the intervention in the graphical model. Let $G^{\hat{X}}$ denote the graph G after adding a variable \hat{X}_i as a parent of each variable $X_i \in \mathbf{X}$ (see figure 3.5b). The variable \hat{X}_i can be thought of as representing the mechanism by which X_i takes its value, either by being set via intervention or as a stochastic function of its other parents [74].

if $(Y \perp\!\!\!\perp \hat{X} | Z, W)$ in $G_{\bar{Z}}^{\hat{X}}$:

$$P(Y | do(Z = z), do(X = x), W = w) = P(Y | do(Z = z), W = w) \quad (3.5)$$

The statement that $\mathbf{Y} \perp\!\!\!\perp \hat{\mathbf{X}}$ (without conditioning on \mathbf{X}) implies that there is no unblocked path from \mathbf{X} to \mathbf{Y} in G which *includes* an arrow leaving \mathbf{X} . These are the only paths by which intervening in \mathbf{X} can effect \mathbf{Y} .

Figure 3.5: Example application of equation 3.5. $(Y \perp\!\!\!\perp \hat{X}|W, Z) \implies P(Y|do(X), W, Z) = P(Y|W, Z)$. We have to condition on Z because conditioning on W blocks the path $\hat{X} - X - W - Y$ but opens $\hat{X} - X - Z - Y$.



3.2.1 Identifiability

A natural question to ask is, given a set of assumptions about the causal graph, is it possible to estimate a given interventional distribution from observational data? This is the identifiability problem. It asks if we can obtain an unbiased point estimate for the causal query of interest in the infinite data limit. A query is non-parametrically identifiable if it is identifiable without assumptions about the functional form of the dependencies between variables in the graph.

Definition 10 (Non-parametric identifiability). Let G be a causal graph containing observed variables \mathbf{V} and latent variables \mathbf{U} and let $P(\cdot)$ be any positive distribution over \mathbf{V} . A causal query of the form $P(\mathbf{Y}|do(\mathbf{X}), \mathbf{W})$, where \mathbf{Y}, \mathbf{X} and \mathbf{W} are disjoint subsets of \mathbf{V} , is non-parametrically identifiable if it is uniquely determined by $P(\cdot)$ and G .

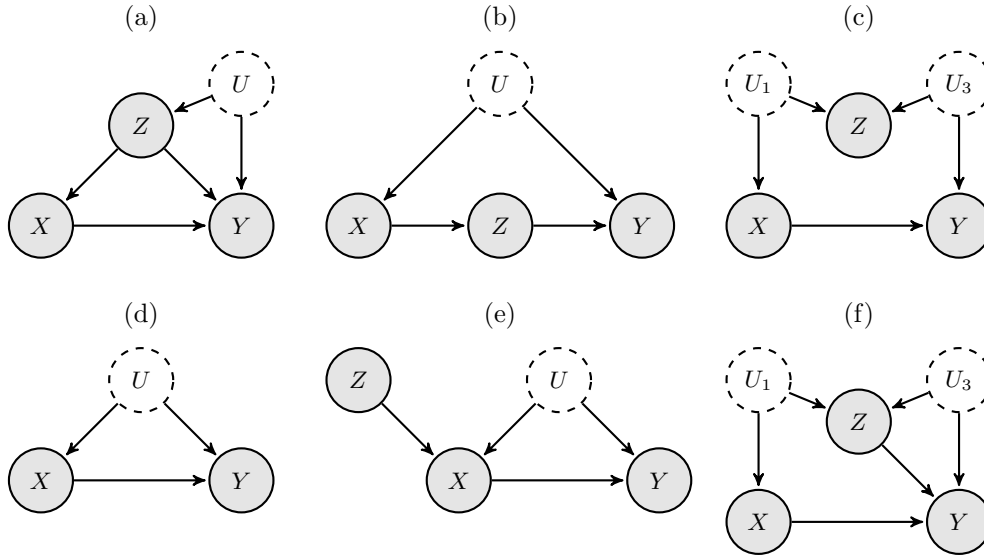
The question of non-parametric identifiability is solved. The do calculus is complete [113, 62]. A problem is identifiable if and only if the interventional distribution of interest can be transformed into term containing only observational quantities via repeated application of the do-calculus. There is a polynomial time algorithm [112] based on these properties that, for a given network and interventional (do-type) query, can:

1. determine if the query can be translated into an expression involving only distributions over observed variables. In other words, determine if the query is identifiable given the assumptions encoded by the network, and
2. if it is identifiable, return the required expression.

Figure 3.6 shows some examples of identifiable and non identifiable queries. I have created a javascript implementation of the identifiability algorithm [112] on which you can test your own queries <http://finnhacks42.github.io>.

Many interesting questions relating to identifiability remain open. What is the minimal (by some metric) additional information that would be required to make a non-identifiable query identifiable? What if we assume various restrictions on the functional form of the relationships between the variables? Some queries which are not identifiable non-parametrically can be identified by additional assumptions such as linearity. A complete algorithm for the problem of linear identifiability is yet to be found, despite a rich body of work [28, 124, 37].

Figure 3.6: Examples of identifiable and non-identifiable queries. In subfigures (a), (b) and (c) the causal query $P(Y|do(X))$ is identifiable. In subfigures (d), (e) and (f) it is not.



Although identifiability is a natural and important question to ask, it does not partition causal questions into solvable and unsolvable. Estimators for identifiable queries can be slow to converge and we may be able to obtain useful bounds on causal effects in cases where point estimates are not identified.

3.3 Estimation

3.3.1 Defining causal effects

So far we have described causal effect estimation in term of identifying the interventional distribution $P(Y|do(X))$ from observational data. This interventional distribution is in fact a family of distributions parameterised by the value, x , to which the treatment variable X is set. From a decision theoretic viewpoint, we can select an optimal action x by specifying a utility function $\mathcal{U} : y \in \mathcal{Y} \rightarrow \mathbb{R}$ that assigns a value to each outcome y and then selecting the action that maximises the expected utility.

$$x^* = \arg \max_x \mathbb{E}_{y \sim P(Y|do(X=x))} [\mathcal{U}(y)] \quad (3.6)$$

Frequently however, studies wish to define and estimate a causal effect without reference to a specific utility function. There are a variety of ways of defining causal effects that can be viewed as different ways of summarising the family of interventional distributions. For a binary treatment variable X , the average causal effect, ACE¹ is defined as:

$$ACE = \mathbb{E}[Y|do(X=1)] - \mathbb{E}[Y|do(X=0)] \quad (3.7)$$

Assuming the expectations in equation 3.7 are well defined, the ACE captures the shift in the mean outcome that arises from varying X . It does not capture changes in variance or higher

¹also referred to as the average treatment effect (ATE)

moments of the distribution. The ACE can be generalised to non-discrete interventions by considering the effect on the expectation of Y of an infinitesimal change in x . If X is linearly related to Y then the ACE is constant and equivalent to the corresponding coefficient in the linear structural equation model.

$$ACE(x) = \frac{d}{dx} \mathbb{E}[Y|do(X = x)] \quad (3.8)$$

The average causal effect is often introduced as the average over individual causal effects as discussed in section 2.2. Individual causal effects are deterministic and cannot be expressed as properties of the interventional distribution. However we can personalise the average causal effect by stating it with respect to some observed context. I will refer to this as the personalised causal effect (PCE).²

$$PCE(z) = \mathbb{E}[Y|do(X = 1), z] - \mathbb{E}[Y|do(X = 0), z] \quad (3.9)$$

In some cases we may be interested in the average causal effect for some sub-group of the population. A particularly common example of this is the average treatment effect of the treatment of the treated (ATT). This would be the key quantity of interest if we had to decide whether or not to continue providing a program or treatment for which we could not control the treatment assignment process.

$$ATT = \mathbb{E}_{z \sim P(Z|x=1)}[Y|do(X = 1)] - \mathbb{E}_{z \sim P(Z|x=1)}[Y|do(X = 0)] \quad (3.10)$$

Causal effects can also be written in terms of counterfactuals. The ACE is $\mathbb{E}[Y^1 - Y^0]$. We could estimate the ratio of expectations $\frac{\mathbb{E}[Y^1]}{\mathbb{E}[Y^0]}$ instead of the difference. However, the quantity $\mathbb{E}\left[\frac{Y^1}{Y^0}\right]$ depends on the joint distribution over the counterfactual variables (Y^1, Y^0) and thus cannot be computed from the interventional distribution.

Another way of conceptualising causal effects is as a property indicating the strength of the causal link between two variables. This notion is complex to formalise when the relationship between variables is non-linear. Suppose $Y = X \oplus Z$ with $P(Z = 1) = \frac{1}{2}$, the interventional distributions over X are identical after marginalising out Z . Janzing et al. [67] propose a number of postulates that a notion of causal strength could satisfy, demonstrate why previous measures fail these postulates and propose an alternative based on information flow.

3.3.2 Estimating causal effects by adjusting for confounding variables

Probably the two most frequently applied approaches to estimating causal effects from observational data are instrumental variables and adjusting for confounding factors. Instrumental variables correspond to the graph in figure 3.6e, which is not identifiable without parametric assumptions, however they can provide tight bounds. Adjusting for confounding equates to identifying a set of variables Z such that the ignorability assumption discussed in section 2.2 holds. This corresponds to a simple graphical test known as the backdoor criterion [87]. The setting is also referred to as unconfounded.

²This quantity is sometimes called the conditional average treatment effect (CATE), however that term is also used for the sample rather than population effect.

Figure 3.7

(a) There can be multiple valid adjustment sets. (b) Conditioning on Z opens the backdoor path $X - U_1 - Z - U_2 - Y$ from X to Y .



Theorem 11 (The backdoor criterion). [87] Let \mathbf{X} , \mathbf{Z} and \mathbf{Y} be disjoint sets of vertices in a causal graph G . If \mathbf{Z} blocks (see Definition 7) for every path from X_i to Y_j that contains a link into X_i , for every pair $(X_i \in \mathbf{X}, Y_j \in \mathbf{Y})$, and no node in \mathbf{Z} is a decedent of a node in \mathbf{X} then the backdoor criterion is satisfied and;

$$P(\mathbf{y}|do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y}|\mathbf{x}, \mathbf{z}) P(\mathbf{z}) \quad (3.11)$$

The backdoor criterion derives from rule 2 of the do-calculus. Selecting which covariates should be adjusted for to estimate a causal effect reduces to identifying a set which satisfies the backdoor criterion. There may be more than one valid adjustment set, see figure 3.7a. The seemingly simple problem of determining if a variable should be adjusted for when estimating causal effects has been the subject of substantial debate and controversy [88]. Adjusting for the wrong variables (even pre-treatment variables) can introduce or magnify bias, see figure 3.7b. Causal graphs and the back door criterion provide a clear mechanism for deciding which variables should be adjusted for. For a practical example, see the discussion in Schisterman et al. [109] on whether birth weight should be adjusted for to estimate the causal effect of smoking on neonatal mortality.

Given that a set of variables \mathbf{Z} satisfies the backdoor criterion (or equivalently the conditional ignorability assumption), the interventional distribution is asymptotically identifiable and can be estimated from equation 3.11. The expected value of Y after the intervention $do(X = x)$ is given by equation 3.12 and the average causal effect for a binary intervention $x \in \{0, 1\}$ is given by equation 3.13.

$$\mathbb{E}[Y|do(X = x)] = \mathbb{E}_{z \sim P(\mathbf{Z})} [\mathbb{E}[Y|x, \mathbf{z}]] \quad (3.12)$$

$$ACE = \mathbb{E}_{z \sim P(\mathbf{Z})} [\mathbb{E}[Y|1, \mathbf{z}] - \mathbb{E}[Y|0, \mathbf{z}]] \quad (3.13)$$

Assuming x and \mathbf{z} are discrete, equation 3.12, and thus the ACE, can be estimated by selecting the data for which $X = x$, stratifying by \mathbf{Z} , then computing the mean outcome within each stratum and finally weighting the results by the number of samples in each strata. However this approach is not workable for most real world problems with finite samples as the number of strata grows exponentially with the dimension of \mathbf{Z} and it cannot handle continuous covariates. There is a substantial body of work within in the statistics and econometrics literature on estimating

average causal effects assuming conditional ignorability, see Imbens [64] for a comprehensive review. The key approaches are based on matching on covariates, propensity score methods and regression. We now examine these approaches from a machine learning perspective.

In standard supervised learning, we have a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ assumed to be sampled i.i.d from an unknown distribution $P(\mathbf{x}, y) = P(\mathbf{x}) P(y|\mathbf{x})$. The goal is to select a hypothesis $h \in \mathcal{H} : \mathcal{X} \rightarrow \mathcal{Y}$ such that, on unseen data $\sim P(\mathbf{x}, y)$, $h(\mathbf{x})$ is close (by some metric) in expectation to y . In other words we wish to minimise the generalisation error $E_{out}(h)$,

$$E_{out}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim P(\mathbf{x}) P(y|\mathbf{x})} [L(h(\mathbf{x}), y)] \quad (3.14)$$

We cannot directly compute the generalisation error as $P(\mathbf{x}, y)$ is unknown, we only have access to a sample. We could search over \mathcal{H} and select a hypothesis $h^*(\mathbf{x})$ that minimises some loss function on the sample data.

$$E_{in}(h) = \frac{1}{n} \sum_{i=1}^n L(h(\mathbf{x}_i), y_i) \quad (3.15)$$

The VC-dimension of the hypothesis space provides (typically loose) bounds on the probability that $E_{out} \gg E_{in}$. However, in practice, the generalisation error is usually estimated empirically from a hold-out set of the sample that was not used to train the model, or via cross-validation.

In the causal effect estimation under ignorability, we have training data $(\mathbf{x}_1, \mathbf{z}_1, y_1), \dots, (\mathbf{x}_n, \mathbf{z}_n, y_n)$ sampled i.i.d from $P(\mathbf{z}) P(\mathbf{x}|\mathbf{z}) P(y|\mathbf{x}, \mathbf{z})$. Estimating $\mathbb{E}[Y|do(\mathbf{X} = \mathbf{x})]$ corresponds to selecting a hypothesis $h \in \mathcal{H} : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Y}$ that minimises;

$$E_{out} = \mathbb{E}_{(\mathbf{x}, \mathbf{z}, y) \sim \delta(\mathbf{x} - \mathbf{x}') P(\mathbf{z}) P(y|\mathbf{x}, \mathbf{z})} [L_2(h(\mathbf{x}, \mathbf{z}), y)], \quad (3.16)$$

$$= \mathbb{E}_{(\mathbf{z}, y) \sim P(\mathbf{z}) P(y|\mathbf{x}, \mathbf{z})} [L_2(h(\mathbf{x}, \mathbf{z}), y)], \quad (3.17)$$

Johansson et al. [70] identified that this is equivalent to the covariate shift problem. If we let $\mathbf{v} = (\mathbf{x}, \mathbf{z})$ then we have training data sampled from $P_{train}\{\mathbf{v}\} P(y|\mathbf{v})$ where $P_{train}\{\mathbf{v}\} = P(\mathbf{z}) P(\mathbf{x}|\mathbf{z})$ but at test time the data will be sampled from $P_{test}\{\mathbf{v}\} P(y|\mathbf{v})$, where $P_{test}\{\mathbf{v}\} = \delta(\mathbf{x} - \mathbf{x}') P(\mathbf{z})$.³ With this connection to covariate shift in mind, let us return to regression, matching and propensity scores.

Regression

The regression approach is to learn a function that is a good approximation to the output surface $\mathbb{E}[Y|X, Z]$. Let $f_1(z) = \mathbb{E}[Y|X = 1, Z = z]$. The expectation of Y after the intervention $X = 1$ is then obtained by taking the expectation with respect to Z , $\mathbb{E}[Y|do(X = 1)] =$

³It is not obvious that the question of estimating causal effects under ignorability entirely reduces to covariate shift. Take the case where we have a binary intervention $x \in \{0, 1\}$. Suppose we learn $h(1, \mathbf{z}) = \mathbb{E}[Y|x = 1, \mathbf{z}] + g(\mathbf{z})$ and $h(0, \mathbf{z}) = \mathbb{E}[Y|x = 0, \mathbf{z}] + g(\mathbf{z})$, then the estimated average causal effect equals the true average causal effect for any function g , $\mathbb{E}[h(1, \mathbf{z}) - h(0, \mathbf{z})] = \mathbb{E}[Y|x = 1, \mathbf{z}] - \mathbb{E}[Y|x = 0, \mathbf{z}]$. More generally, if the goal is to select an optimal action x^* from a continuous space of possible interventions we need algorithms capable of leveraging any structure in the relationship between x and y as well as a means of focusing the loss on regions of the sample likely to affect x^* .

$\mathbb{E}_{z \sim P(Z)} [\mathbb{E}[Y|X=1, z]]$. We can learn a parametric regression model $\hat{f}_1(z)$ via empirical risk minimisation.

$$\hat{f}_1(z) = h_1(z; \hat{\theta}_{obs}), \text{ where } \hat{\theta}_{obs} = \arg \min_{\theta \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = 1\} L(h_1(z_i; \theta), y_i) \right] \quad (3.18)$$

This estimator is consistent with respect to the observational distribution. As the sample size tends to infinity, $\hat{\theta}_{obs}$ approaches the parameter within the hypothesis space that minimises the expected loss given data sampled from the observational distribution.

$$\lim_{n \rightarrow \infty} \hat{\theta}_{obs} = \arg \min_{\theta \in \Theta} \mathbb{E}_{(z,y) \sim P(z|x=1) P(y|x=1,z)} [L(h_1(z; \theta), y)] \quad (3.19)$$

If the model is correctly specified such that $f_1(z) = h_1(z; \theta^*)$ for some $\theta^* \in \Theta$ then the empirical risk minimisation estimate is consistent with respect to the loss over any distribution of Z [122], including the interventional one.

$$\lim_{n \rightarrow \infty} \hat{\theta}_{obs} = \theta^* = \arg \min_{\theta \in \Theta} \mathbb{E}_{(z,y) \sim P(z) P(y|x=1,z)} [L(h_1(z; \theta), y)] \quad (3.20)$$

The average causal effect can then be estimated by,

$$\hat{\tau}_{reg} = \sum_{i=1}^n \left(\hat{f}_1(z_i) - \hat{f}_0(z_i) \right) \quad (3.21)$$

Regression thus has a direct causal interpretation if the parametric model is correctly specified and the covariates included form a valid backdoor adjustment set for the treatment variable of interest in the corresponding structural equation model.

Propensity scores

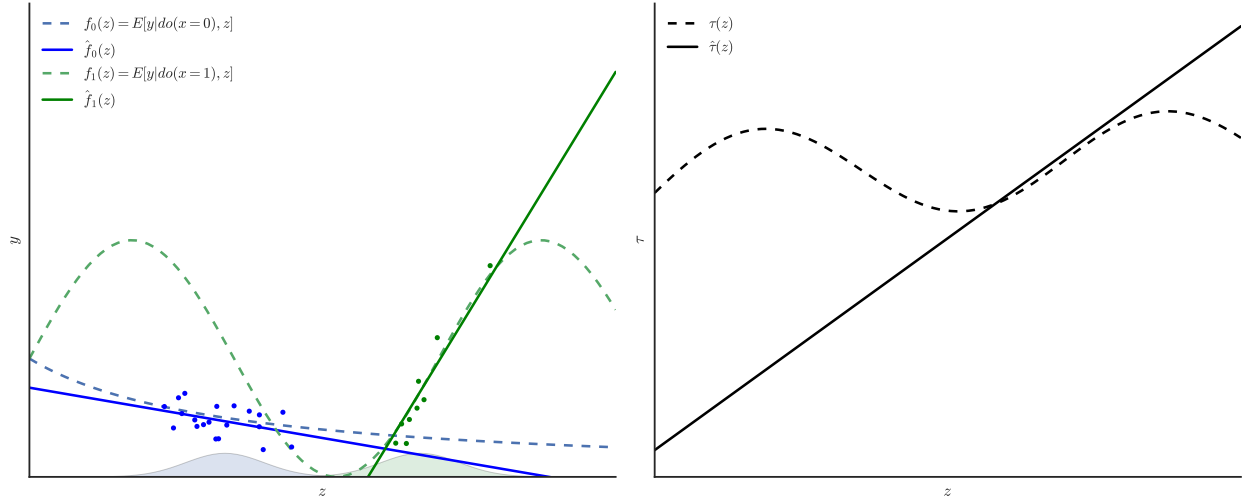
If the parametric model is miss-specified then the parameter that minimises the loss depends on the distribution from which the covariates z are sampled. The model learned by ERM could perform very well in a validation set (which estimates the generalisation error over the observational distribution of (x, z)) but yield very poor estimates of the causal effect, see figure 3.8.

A general approach to estimating the expectation of some function $f(\cdot)$ with respect to data from some distribution $P(\cdot)$, when we have data sampled from a different distribution $Q(\cdot)$ is importance sampling [59, 74].

$$\mathbb{E}_{v \sim P(v)} [f(v)] = \mathbb{E}_{v \sim Q(v)} \left[f(v) \frac{P(v)}{Q(v)} \right] \quad (3.22)$$

This importance weighting approach can be applied to the covariate shift/average causal effect problem by weighting the terms in the empirical risk minimisation estimator [122].

Figure 3.8: Parametric regression may yield poor estimates of causal effects if the model is misspecified, even if the model fits well over the domain of the training data. In this example, $P(Z|X=0) \sim N(\mu_0, \sigma_0^2)$ and $P(Z|X=1) \sim N(\mu_1, \sigma_1^2)$ with little overlap in the densities. If $X=0$ then $Y \sim N(f_1(x) = \sin(x), \sigma_y^2)$ and if $X=1$ then $Y \sim N(f_0(x) = \frac{1}{x+1}, \sigma_y^2)$. We estimate $f_1(z)$ from the sample in which $X=1$ (green points) and $f_0(z)$ from the sample for which $X=0$ (blue points). In both cases the linear model is a good fit to the data. However, the resulting estimate of the causal effect is very poor for the lower values of z .



$$\hat{\theta}_{iw} = \arg \min_{\theta \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = 1\} L(h_1(z_i; \theta), y_i) \frac{P(z_i) \delta(x_i - 1)}{P(z_i) P(x_i = 1|z_i)} \right] \quad (3.23)$$

$$= \arg \min_{\theta \in \Theta} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = 1\} L(h_1(z_i; \theta), y_i) \frac{1}{e(z_i)} \right], \quad (3.24)$$

where $e(z)$ is the propensity score, defined by [101];

$$e(z) \equiv P(x = 1|z) \quad (3.25)$$

The estimator in equation 3.23 is an example of a doubly robust estimator [108?]. Doubly robust methods are asymptotically unbiased as long as either the regression model h or propensity score e are correctly specified [100].

The propensity score can be used to estimate the average causal effect without specifying a regression model for $\mathbb{E}[Y|X, Z]$. Rosenbaum and Rubin [101] demonstrated that if the ignorability assumption is satisfied by conditioning on Z , then it is also satisfied by conditioning on $e(z)$. This allows for estimators based on stratifying, matching or regression on the propensity score rather than the covariates Z . Inverse propensity weighting can also be combined with empirical estimation of $\mathbb{E}[Y|X, Z]$ yielding the simple, albeit inefficient, estimator in equation 3.27 [64]. In some settings, such as stratified randomised trials [65] or learning from logged bandit feedback [18] the propensity score may be known. However in general, it must be estimated from data. Frequently this is done with a simple parametric model such as logistic regression, but a wide range of standard machine learning algorithms including bagging and boosting, random forests and neural networks can also be applied [12]. Lunceford et al. [79] review the theoretical properties of key propensity score based estimators, including stratification and inverse propensity weighting.

$$\mathbb{E}[Y|do(X = x)] = \mathbb{E}_{z \sim P(\mathbf{Z})} [\mathbb{E}[Y|x, \mathbf{z}]] = \mathbb{E}_{z \sim P(\mathbf{Z}|\mathbf{x})} \left[\mathbb{E}[Y|x, \mathbf{z}] \frac{1}{e(\mathbf{z})} \right] \quad (3.26)$$

$$\hat{\tau}_{ip} = \frac{1}{n} \sum_{i=1}^n \left(\frac{\mathbb{1}\{x_i = 1\} y_i}{e(\mathbf{z}_i)} - \frac{\mathbb{1}\{x_i = 0\} y_i}{1 - e(\mathbf{z}_i)} \right) \quad (3.27)$$

Matching

There is a straightforward connection between matching and regression for causal effect estimation. If $h \in \mathcal{H} \implies h + a \in \mathcal{H}$ for any constant a and \hat{f} is selected by minimising empirical risk with an L_2 loss then $\frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = 1\} \hat{f}_1(\mathbf{z}_i) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i = 1\} y_i$ ⁴, and equation 3.21 can be re-written as:

$$\hat{\tau}_{reg} = \sum_{i=1}^n \left[\mathbb{1}\{x_i = 1\} \left(y_i - \hat{f}_0(\mathbf{z}_i) \right) + \mathbb{1}\{x_i = 0\} \left(\hat{f}_1(\mathbf{z}_i) - y_i \right) \right] \quad (3.28)$$

This formulation of the regression estimator highlights the missing data aspect of casual effect estimation. For each instance, the regression models are used to estimate the counterfactual outcome had the instance received the alternate treatment. Matching estimates the counterfactual outcome for an instance from the outcome of *similar* instances that received a different treatment. Abadie and Imbens [1] analyse an estimator where both target and control instances are matched and the matching is done with replacement, let $j \in J_k(i)$ be the indices of the k instances closest to i by some metric $d(\mathbf{z}_i, \mathbf{z}_j)$ such that $x_i \neq x_j$.

$$\hat{\tau}_{match} = \sum_{i=1}^n \left[\mathbb{1}\{x_i = 1\} \left(y_i - \frac{1}{k} \sum_{j \in J_k(i)} y_j \right) + \mathbb{1}\{x_i = 0\} \left(\frac{1}{k} \sum_{j \in J_k(i)} y_j - y_i \right) \right] \quad (3.29)$$

This estimator is equivalent to equation 3.28 with k nearest neighbour regression. There are many variants of matching estimators utilising different distance metrics, matching with or without replacement (and in the latter case, greedy or optimal matching) and with or without discarding matches beyond some threshold [30, 102]. Although intuitive, matching estimators in general have poor large sample properties [2]. An exception is where the goal is to estimate the average treatment effect of treatment on the treated in settings where there is a large set of control instances (compared to treatment instances) [64].

The practical performance of the estimation approaches discussed in this section will depend on the sample size, dimensionality of the covariates, the complexity of the treatment assignment mechanism and output function, and the degree of prior knowledge available about these functions. A key difference between standard machine learning problems and causal effect estimation is that when estimating causal effects we cannot directly apply cross-validation or a hold-out set for model selection because we lack samples from the counterfactual.

The significance of this should not be underestimated. Cross-validation has allowed applied machine learning to succeed with a very atheoretical approach on the basis that we can identify

⁴[64] state this holds for most implementations

when a model is successful. With causal effect estimation there is no guarantee that a model that performs well at prediction (even out of sample) will accurately estimate the outcome of an intervention. Sugiyama et al. [122] propose inverse propensity weighted cross validation for the covariate shift problem. There is relatively little theory on model selection for estimating the propensity score. To achieve asymptotically unbiased estimates, the covariates should satisfy the backdoor criterion. It is also known that conditioning on instrumental variables, which directly influence X but not Y , increases variance without any reduction in bias and can increase bias if there are unmeasured confounding variables [130, 16, 89, 83]. With doubly robust estimators, one could apply an iterative approach, fitting a propensity score model, using the results for inverse propensity weighted cross-validation of the regression model and then selecting covariates for the propensity model on the assumption the estimated regression function was correct. I have found no examples of such an approach.

The performance of methods for causal effect estimation can be tested on simulated data [44, 137, 58, 36] or by comparing estimates from observational studies with the results from corresponding experiments [75, 42, 57, 56, 35, 116, 7]. Unfortunately there are a relatively small number of examples where comparable observational and experimental data are available. The results are mixed with later studies finding generally better alignment of results but it is hard to ascertain if this is due to improved methodological approaches or over-fitting to the available data-sets.

3.4 Causal Discovery

We now move to the much more general problem of learning a causal graph from observational data. In this setting we make much broader assumptions about the structure of the graph. For example, that it is acyclic or that we have no unmeasured confounding variables. We do not assume the existence or directions of any links between the variables. Amazingly, it is possible to infer some aspects of causal structure with such general assumptions. The set of conditional independence in a non-experimental data set indicates some causal structures are more likely than others. In addition, there can be subtle asymmetries in the relationship between the joint distribution of cause and effect and the distributions of cause given effect and effect given cause. These clues are the key to causal discovery algorithms.

Causal discovery is a much grander goal than causal effect estimation given a known causal network. Arguably, if achieved, it would equate to the automation of scientific discovery. We need simply supply our algorithm with a vast collection of variables (regardless of their relevance to the problem) and it would learn the causal structure and from that allow us to estimate the effects of any intervention we cared to make. Unfortunately, causal discovery is very hard. Even with the assumption that the causal graph is acyclic and there are no latent variables, the number of possible graphs grows exponentially with the square of the number of variables.

In the next sections we briefly survey the key approaches to causal discovery. We roughly divide the methods into those based on those that exploit the connection between the conditional independencies in a joint distribution and the structure of a causal model and those that leverage assumptions about the functional form of the relationships between cause and effect.

3.4.1 Conditional independence based methods

One general approach is to look for clues about the structure of the network in the conditional independence relations in the distribution. For any Bayesian network, G , (causal or not) we can read off conditional independencies in the joint distribution from the structure of the network. If

a set of variables Z d-separates X and Y in G then $(X \perp\!\!\!\perp Y|Z)$ in the distribution P . However, we want to work in the other direction, from conditional independence in the distribution to the structure of the network. This requires that we assume the reverse condition: $(X \perp\!\!\!\perp Y|Z)$ in P must imply Z d-separates X and Y in G . This assumption, commonly referred to as **faithfulness** ??, says there are no additional independence relations that are satisfied in P but not in all distributions P' that are compatible with G . Stating that P is faithful to G is equivalent to G is a **perfect map** [85] for P .

Faithfulness is an assumption. It does not always hold and we cannot verify it from the observational data we wish to use for causal inference. However, most distributions generated by a causal Bayesian network will be faithful to that network. For faithfulness to be violated, different causal effects must exactly balance one-another out. For example, consider a simple binary variable model of chocolate consumption, income and obesity (figure). If the coefficients in the conditional probability tables are just right then the direct effect of chocolate on obesity will exactly balance the indirect effect through income and obesity will appear independent of chocolate consumption. However, this independence is not stable. It would disappear under a small perturbation to any of the parameters. In discrete systems, violations correspond to the solutions to polynomial equations over values in the CPD tables and thus are a space of measure zero with respect to all possible distributions associated with the graph [74].

Given the faithfulness assumption, our causal discovery problem reduces to finding the set of Bayesian networks that have exactly the dependency structure as we observe in P . This set can also be referred to as the Markov equivalence class compatible with P .

Without hidden common causes

The strong assumption that there are no hidden variables that cause two or more variables in \mathbf{V} significantly reduces the 'search space' of Bayesian networks we must consider. This assumption is referred to as causal sufficiency [118].

We will begin with a brute force algorithm (described as the SGS algorithm in Spirtes et al. [118] and IC algorithm in Pearl [87]). While it is impractical for all but the smallest of networks, it demonstrates key concepts that also underlie the more useful and complex algorithms we will discuss later.

The SGS (or IC) Algorithm

Input: A distribution P , over variables \mathbf{V} , that was generated by and is faithful to an (unknown) Bayesian network G

Output: A partially directed network that represents the Markov equivalence class of G

1. Join all pairs of vertices $(a, b) \in \mathbf{V}$ with an un-directed link to form a complete graph.
 2. For each link $a - b$ search for a set $\mathbf{S}_{a,b} \subseteq \mathbf{V} \setminus \{a, b\}$ that renders a and b conditionally independent. If such a set (including the empty set) exists then a and b cannot be directly connected in G so delete the link.
 3. For all pairs of non-linked variables (α, β) with a common neighbour, c , if $c \notin \mathbf{S}_{\alpha,\beta}$, then c must be a collider in the path α, c, β so add arrows to direct the links $\alpha - c$ and $\beta - c$ towards c .
 4. Recursively try to orient any edges that remain un-directed to avoid creating cycles (because they are not there by assumption) and additional colliders (because any colliders were found in step 3).
-

The SGS algorithm utilises the fact that a collider structure (figure 3.1d) induces a distinct conditional independence relation. Assuming you have a consistent conditional independence test, it converges to return a partially directed network that represents the Markov equivalence class for the generating causal model. Unfortunately the number of conditional independence tests required for step 2 grows exponentially (in the worst case) with the number of variables. Not only that, but for each edge that is in the true network, the algorithm will always tests all other possible subsets of variables. If the assumption that there are no hidden common causes or that the distribution is faithful are violated, step 3 of the SGS algorithm can produce double headed arrows.

The PC algorithm Spirtes et al. [118] modifies step 2 of the SGS algorithm to utilise the fact that if two variables (a, b) are conditionally independent given some set, they will also be conditionally independent given a set that contains only variables adjacent to a or b . It also checks for low order conditional independence relations before higher order ones. This allows it to exploit any sparsity in the true network, leading to much better average case performance [118] (although the worst case, where the true network is complete, is still exponential). With finite data, the order in which the links are considered can change the output (unlike for SGS). The effect of wrongly removing a link early on flows through to later conditional independence tests by changing which nodes are considered adjacent.

The PC Algorithm

Input: A distribution P , over variables $\mathbf{V} = \{V_1 \dots V_k\}$, that was generated by and is faithful to an (unknown) Bayesian network G

Output: A partially directed network that represents the Markov equivalence class of G

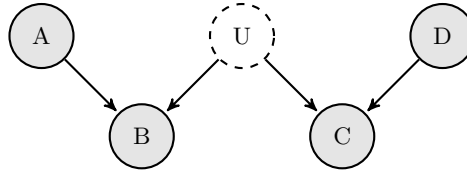
1. As for SGS
 2. **for** each link $a - b$:
 - $n = 0$
 - $\mathbf{A}_{a,b} = \{A_1 \dots A_j\}$ be the set of nodes adjacent to a and/or b
 - while** a and b are connected and $n < j$:
 - if** any subset of size n of \mathbf{A} makes a and b conditionally independent:
 - delete the link
 - $n = n + 1$
 3. as for SGS
 4. as for SGS
-

The PC algorithm also returns a set of Markov equivalent networks consistent with the distribution. Since we have assumed there are no hidden variables, for any single graph in this set we can calculate causal effects from the truncated product formula 2.3. We can then bound the true causal effect by combining the results for the all the networks. This procedure is the IDA algorithm [81] and has been found to outperform standard regularisation techniques at finding causal effects in a high-dimensional yeast gene expression data set [80]. An implementation is available in the R package [71]

With hidden variables

There are an number of difficulties in extending the approach of the last section to deal with the case where there are latent variables. With an unknown number of hidden variables there

Figure 3.9: A distribution faithful to this DAG is not faithful to any DAG over the variables $\{A, B, C, D\}$ after marginalising over U .



are infinity many possible structures to search over. In addition, the space of causal networks is not closed under marginalisation. If we have a distribution that $P'(\mathbf{O}, \mathbf{U})$ generated by and is faithful to a network G the distribution $P(\mathbf{O})$, that results from marginalising over \mathbf{U} , may not be faithful to any Bayesian network (see figure 3.9). The key to constraining the space of possible models is that many latent structures are equivalent (under transforms of the hidden variables).

Theorem 12. [127] *For every latent structure there is a dependency equivalent structure such that every latent (unobserved) variable is a root node with exactly two children .*

Since we only care about the causal relationships between observed variables, it is sufficient to search over networks where any hidden variables have no parents and directly cause two of the observed variables. Instead of representing hidden variables explicitly we can capture the necessary independence relations with a more general graphical model that supports bi-directed edges that play the role of a hidden confounding variable. These models, referred to as maximal ancestral graphs (MAGs) are closed under marginalisation and conditioning.

For any DAG with latent (and selection) variables there is a unique MAG [96]. This makes it possible to extend the PC algorithm to latent structures, resulting in the FCI algorithm [?]. The logic behind the algorithm is very similar. Certain structures are ruled out as a consequence of being inconstant with the observed conditional independence relations. The output is an equivalence class of MAGs, which can be represented graphically as a partial ancestral graph PAG [119]. Assuming there are no selection variables, the PAG can contain four types of link:

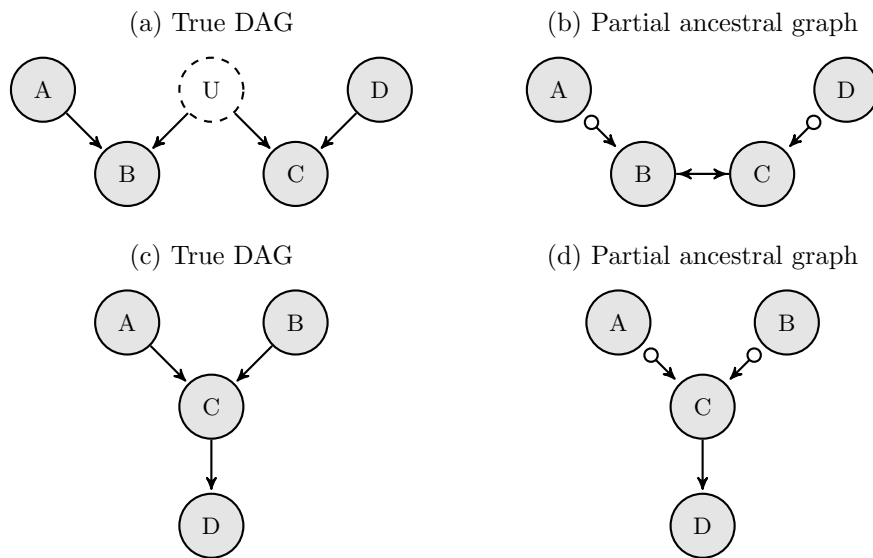
1. $X \rightarrow Y$, meaning X causes Y
2. $X \leftrightarrow Y$, meaning there is a latent variable that causes X and Y .
3. $X \circ \rightarrow Y$, either X causes Y or a latent variable causes both.
4. $X \circ - \circ Y$, either X causes Y or Y causes X or a latent variable causes both.

The circles indicate where it is ambiguous if there should be an arrowhead (IE where there is one in some MAGs and not in others in the equivalence class). Counter-intuitively it is sometimes possible to rule out or confirm the existence of a confounding variable and fully determine the causal type of a link (see examples in figure 3.10).

The FCI algorithm can be made complete such that it discovers all aspects of the true causal structure that are identifiable from the conditional independence relations of a distribution over observed variables and the faithfulness assumption [134]. More recently [32] have proposed the RFCI algorithm, which in some cases returns more ambiguous links than FCI but is substantially faster. [29] point out that the problem of learning sparse causal networks from data is not NP-hard and propose the FCI+ algorithm, that requires $O(N^{2(k+2)})$ conditional independence tests, where k is the maximum node degree over the observed variables.

Latent variables can create constraints on the marginal distribution over the observed variables that cannot be expressed in term of conditional independencies. These generalised constraints can be expressed and leveraged within nested Markov models [114, 111]

Figure 3.10: FCI examples: true graph and FCI output



All the algorithms discussed in this section rely on being able to perform conditional independence tests. This is non-trivial with high dimensional data. If the functional relationship between the variables is linear with Gaussian noise then the network represents a multivariate normal distribution and a pair of variables are conditionally independent if and only if the corresponding entry in the inverse correlation matrix is non-zero [74]. Where the functions are non-linear one can apply kernelised independence tests [51, 136]

3.4.2 Discovery with functional models

The algorithms we have considered so far return a Markov equivalence class. They cannot distinguish between two models that result in the same set of conditional independence relations. Consider the very simple case where there are only two variables and the possible causal structures are $X \rightarrow Y$ or $Y \rightarrow X$. These models have the same dependency structure but in one case $P(Y|do(X)) = P(Y|X)$ and in the other $P(Y|do(X)) = P(Y)$. No algorithm relying purely on conditional independence relations can separate these two cases.

Let us focus only on the two variable case $X \rightarrow Y$ or $Y \rightarrow X$. What possible clues could there be in the distribution $P(X, Y)$ that could indicate which causal model it was generated from? Recall the functional definition of causality (section 2.3). There are a number of assumptions about the form of the functions that can allow us to identify the causal direction: non-invertible functions, additive noise [61], post-non-linear additive noise [135] or linear models with non-Gaussian noise [60].

The causal direction can also be identified via a connection between causal discovery and semi-supervised learning [69]. Suppose we are trying to learn $P(Y|X)$. The goal of semi-supervised learning is to improve our estimate of $P(Y|X)$ by leveraging additional data sampled from $P(X)$. However, if the true causal model is $X \rightarrow Y$ then there is some function mapping values of X to Y which should be invariant to any changes in the input distribution $P(X)$. Therefore $P(X)$ should be independent of $P(Y|X)$ and semi-supervised learning should not perform any better than standard supervised learning. However if the true causal model is $Y \rightarrow X$ then variations in the $P(X)$ can result from both the input distribution over Y and the mapping from X to Y and semi-supervised learning could help. This assumption of independence of mechanism and input can also allow the identification of the causal direction in SEMs even where

the functions are deterministic and invertible [33]. Janzing et al. [68] leverage an information geometric viewpoint of the independence of mechanism and input to infer the causal direction between a pair of associated variables.

Instead of positing a functional restriction on the relationship between variables and then developing theory to exploit that assumption, [78] propose learning what the causal relationship looks like from data. They assume there will be a difference between the relationship of $P(X)$, $P(Y)$ and $P(X|Y)$ between $X \rightarrow Y$ versus $Y \rightarrow X$. Their algorithm requires a data set in which each row is itself a data set consisting of pairs of variables (x_i, y_i) with a label indicating the direction of causality between X and Y . They use a kernel mean embedding to represent the distributions $P(X)$, $P(Y)$ and $P(X|Y)$ as features for each individual sub-data set and train an algorithm to learn the direction of causality. Unfortunately we do not have a large collection of data sets where the causal direction is known to train such a model. Lopez-Paz et al. [78] instead use a simulated data set so their model will necessarily be based on the assumptions they make when generating the data. Nonetheless this approach makes it possible to rapidly construct a model from a wide range of possible assumptions, without doing a lot of theory to design a specific algorithm optimised to that setting.

[91] have extended results from the bi-variate case to the multivariate setting. They show that if we can come up with a condition that guarantees identifiability for the bi-variate case, we can extend that result to get the conditions under which the multivariate case is identifiable. They build on this to develop an algorithm that allows the construction of causal graphs based on the additive noise assumption.