# Causal Inference and Probabilistic Graphical Models

MoeKid101

## 0   Properties of Conditional Independence

- Symmetry: $X \perp\!\!\!\perp Y | Z \Leftrightarrow Y \perp\!\!\!\perp X | Z$.

- Decomposition: $X \perp\!\!\!\perp Y, W | Z \Rightarrow X \perp\!\!\!\perp Y | Z$.

- Weak union: $X \perp\!\!\!\perp Y, W | Z \Rightarrow X \perp\!\!\!\perp Y | Z, W$.

- Contraction: $(X \perp\!\!\!\perp W | Z, Y) \wedge (X \perp\!\!\!\perp Y | Z) \Rightarrow X \perp\!\!\!\perp Y, W | Z$.

- Intersection: (for only positive distributions) $(X \perp\!\!\!\perp Y | Z, W) \wedge (X \perp\!\!\!\perp W | Z, Y) \Rightarrow X \perp\!\!\!\perp Y, W | Z$.

    *Proof.* From the two assumptions we have $\mathbf{Pr}\,[X|W, Z] = \mathbf{Pr}\,[X|Y, Z]$. So $\dfrac{\mathbf{Pr}\,[X, W|Z]}{\mathbf{Pr}\,[W|Z]} = \dfrac{\mathbf{Pr}\,[X, Y|Z]}{\mathbf{Pr}\,[Y|Z]}$.

    $$\sum_{w \in \mathrm{Val}(W)} \mathbf{Pr}\,[X, W|Z]\,\mathbf{Pr}\,[Y|Z] = \sum_{w \in \mathrm{Val}(W)} \mathbf{Pr}\,[X, Y|Z]\,\mathbf{Pr}\,[W|Z]$$

    Then we have $\mathbf{Pr}\,[X, Y|Z] = \mathbf{Pr}\,[X|Z]\,\mathbf{Pr}\,[Y|Z]$, i.e. $X \perp\!\!\!\perp Y | Z$. Symmetrically, $X \perp\!\!\!\perp W | Z$. $\square$

## 1   Representing Directed PGMs

The goal is to find a compact representation of probabilities so that it would not cost $2^n$ storage for merely $n$ binary random variables. The intuition is to utilizing conditional independence.

**Definition 1.** (**Naive Bayes Model**) *Assuming instances belong to a number of classes, i.e. $C \in \{c_1, ..., c_k\}$ for any sample and there are a number of features $X_1, ..., X_n$. The model is called a naive Bayes model iff. it satisfies the **naive Bayes assumption** that within each class, all features are conditionally independent $(X_i \perp\!\!\!\perp X - \{X_i\}|C)$.*

From the definition of Naive Bayes Model it is straightforward that this model is used for classification. $\mathbf{Pr}\,[C = c_i|x_1 \sim x_n] = \mathbf{Pr}\,[C = c_i] \prod_{j=1}^{n} \mathbf{Pr}\,[x_j|C = c_i]$ gives a very simple way to determine probability and confidence as well.

**Definition 2.** (**Bayesian Network Structures**) *The structure $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ of a Bayesian network is a DAG with $\boldsymbol{pa}_{X_i}$ denoting the parents of $X_i$ in $\mathcal{G}$ and $\boldsymbol{nd}_{X_i}$ denoting the vertices that are not descendants of $X_i$ in G. These networks satisfy **local independency** assumptions that $\forall X_i \in \mathcal{V}$, $X_i \perp\!\!\!\perp \boldsymbol{nd}_{X_i}|\boldsymbol{pa}_{X_i}$. These independency assumptions are denoted by $\mathcal{I}(P)$ where $P$ is the distribution of variables $X_1 \sim X_n$.*

**Definition 3.** (**I-Maps**) *$\mathcal{K}$ is any graph associated with independencies $\mathcal{I}(\mathcal{K})$. $\mathcal{K}$ is an I-map for a set of independencies $\mathcal{I}$ if $\mathcal{I}(\mathcal{K}) \subseteq \mathcal{I}$. (I-map ensures only correctness of $\mathcal{I}(\mathcal{K})$ but not sufficiency.)*

**Definition 4.** (**Factorization**) *If $\mathcal{G}$ is a BN graph with $\mathcal{V} = \{X_1 \sim X_n\}$ and $P$ is the distribution of these variables, then $P$ factorizes according to $\mathcal{G}$ when $P(X_1 \sim X_n) = \prod_{i=1}^{n} Pr\left[X_i|\boldsymbol{pa}_{X_i}\right]$.*

**Definition 5.** (**Bayesian Network**) *A Bayesian network is $\mathcal{B} = (\mathcal{G}, P)$ where $P$ factorizes over $\mathcal{G}$.*

**Theorem 1.** *(**Connecting I-map and factorization**) $\mathcal{G}$ is an I-map for $P$ iff. $P$ factorizes according to $\mathcal{G}$.*

*Proof.* For a DAG $\mathcal{G}$, there is always a topological ordering. So wlog, assume $\mathcal{G}$ satisfies such ordering. Then we know that in $X_1, X_2, ..., X_n$, any edge $X_i \to X_j$ satisfies $i < j$. Consequently, $\mathbf{pa}_{X_i} \subset \{X_1 \sim X_{i-1}\}$, and let $\mathbf{np}_{X_i} = \{X_1 \sim X_{i-1}\} - \mathbf{pa}_{X_i} \subset \mathbf{nd}_{X_i}$. We have

$$P(X_1 \sim X_n) = \prod_{i=1}^{n} \mathbf{Pr}\left[X_i | X_1 \sim X_{i-1}\right] = \prod_{i=1}^{n} \mathbf{Pr}\left[X_i | \mathbf{pa}_{X_i}, \mathbf{np}_{X_i}\right] = \prod_{i=1}^{n} \frac{\mathbf{Pr}\left[X_i, \mathbf{np}_{X_i} | \mathbf{pa}_{X_i}\right] \mathbf{Pr}\left[\mathbf{pa}_{X_i}\right]}{\mathbf{Pr}\left[\mathbf{pa}_{X_i}\right] \mathbf{Pr}\left[\mathbf{np}_{X_i} | \mathbf{pa}_{X_i}\right]}$$

$$= \prod_{i=1}^{n} \mathbf{Pr}\left[X_i | \mathbf{pa}_{X_i}\right]$$

Conversely, we try proving $\mathbf{Pr}\left[X_i | \mathbf{pa}_{X_i}\right] = \mathbf{Pr}\left[X_i | \mathbf{pa}_{X_i}, \mathbf{nd}_{X_i}\right]$. Notice that we can always ignore the descendants of $X_i$, and in the induced subgraph $\mathcal{G}_i$ of $\mathcal{G}$ by $\{X_i\} \cup \mathbf{pa}_{X_i} \cup \mathbf{nd}_{X_i}$, $X_i$ can always be the last term in topological ordering.

$$\mathbf{Pr}\left[X_i | \mathbf{pa}_{X_i}, \mathbf{nd}_{X_i}\right] = \frac{\mathbf{Pr}\left[X_i, \mathbf{pa}_{X_i}, \mathbf{nd}_{X_i}\right]}{\mathbf{Pr}\left[\mathbf{pa}_{X_i}, \mathbf{nd}_{X_i}\right]} = \frac{\prod_{j=1}^{k} \mathbf{Pr}\left[X_j' | \mathbf{pa}_{X_j'}\right]}{\prod_{j=1}^{k-1} \mathbf{Pr}\left[X_j' | \mathbf{pa}_{X_j'}\right]} = \mathbf{Pr}\left[X_k' | \mathbf{pa}_{X_k'}\right]$$

Here $X_k' = X_i$. The second equality holds because they are identical terms in different DAGs. □

For a Bayesian network, we can infer other dependencies and independencies according to the given structure $\mathcal{G}$. Considering basic dependency relationships, we have the following four empirical conclusion:

- $X$ and $Y$ are dependent if $(X, Y) \in \mathcal{E}$ or $(Y, X) \in \mathcal{E}$.

- $X$ and $Y$ are *likely* dependent if there exists a path between $X$ and $Y$.

- $X$ and $Y$ are conditionally independent given $Z$ if **(1)** there is only one path between $X$ and $Y$; **(2)** $Z$ intercepts the path.

- $Y$ and $Z$ are independent conditional on $X$ if **(1)** $X$ is a common cause of $Y$ and $Z$; **(2)** there exists only one path between $Y$ and $Z$ in the undirected graph.

- $X$ and $Y$ are unconditionally independent but conditionally dependent given $Z$ or any descendants of $Z$ if **(1)** $Z$ is the collision node between $X$ and $Y$; **(2)** there exists only one path between $X$ and $Y$ in the undirected graph.

With these basic rules dealing with chains, forks and colliders, we can use **d-separation** as a formal expression of the above four empirical conclusions to measure in any graph whether variable nodes are connected (or dependent, termed as **d-connected** whose contrary is **d-separated**). Notice that being d-connected doesn't necessarily mean dependent, but being d-separated means independent.

**Definition 6.** (**Blocking**) *For DAG $\mathcal{G}$, a trail (not necessarily a path) $p$ in $\mathcal{G}$ is blocked by a set of nodes $Z$ iff.*
*(1) $p$ contains a chain $A \to B \to C$ or a fork $A \leftarrow B \to C$ and $B \in Z$;*
*(2) $p$ contains a collider $A \to B \leftarrow C$ and none of $B$ and its descendants are in $Z$.*

**Definition 7.** (**d-Separated**) *$X$ and $Y$ as two sets of nodes are d-separated conditional on $Z$ if $Z$ blocks every trail between $X$ and $Y$ in $\mathcal{G}$, annotated d-sep$_\mathcal{G}(X; Y | Z)$. The set of independence assertions corresponding to d-separation is **global Markov independencies** defined by $\mathcal{I}(\mathcal{G}) = \{(X \perp\!\!\!\perp Y | Z) | \text{d-sep}_\mathcal{G}(X; Y | Z)\}$.*

**Theorem 2.** *(**Soundness of d-separation**) If $P$ factorizes over $\mathcal{G}$, then $\mathcal{I}(\mathcal{G}) \subseteq \mathcal{I}(P)$.*

*Proof.* The proof utilizes the soundness of undirected separation and correspondence between Bayesian networks and Markov networks. For any d-sep$_\mathcal{G}(X;Y|Z)$, we can construct $\mathcal{H} = \mathcal{M}[\mathcal{G}^+[U]]$ and obtain sep$_\mathcal{H}(X;Y|Z)$. By the soundness of undirected separation, $X \perp\!\!\!\perp Y|Z$ in the Gibbs distribution $P_{\mathcal{G}^+[U]}$ factorizing $\mathcal{H}$. Then we will only have to prove $P_{\mathcal{G}^+[U]}(U^*) = P_\mathcal{B}(U^*)$. Notice the Gibbs distribution has exactly the same factors as $\mathbf{Pr}\left[X_i|\mathbf{pa}_{X_i}\right]$, and $U \in U^* \Rightarrow \mathbf{pa}_U \subseteq U^*$,

$$P_\mathcal{B}(U^*) = \sum_{X_i \in \mathcal{V}-U^*} \prod_{i=1}^n \mathbf{Pr}\left[X_i|\mathbf{pa}_{X_i}\right] = \left(\prod_{X_i \in U^*} \mathbf{Pr}\left[X_i|\mathbf{pa}_{X_i}\right]\right)\left(\sum_{X_i \in \mathcal{V}-U^*} \prod_{X_i \in \mathcal{V}-U^*} \mathbf{Pr}\left[X_i|\mathbf{pa}_{X_i}\right]\right)$$

$$= P_{\mathcal{G}^+[U]}(U^*)\left(\sum_{X_i \in \mathcal{V}-U^*} \mathbf{Pr}\left[\mathcal{V}-U^*|U^*\right]\right) = P_{\mathcal{G}^+[U]}$$

Now we have proven that d-sep$_\mathcal{G}(X;Y|Z) \Rightarrow X \perp\!\!\!\perp Y|Z$. So $\mathcal{I}(\mathcal{G}) \subseteq \mathcal{I}(P)$. $\qquad\square$

**Theorem 3.** *(Completeness of d-separation)* *For any $(X, Y, Z)$, if not d-sep$_\mathcal{G}(X;Y|Z)$, then $\neg(X \perp\!\!\!\perp Y|Z)$ in $P$ factorizing over $\mathcal{G}$ (completeness).*

# 2 Representing Undirected PGMs

The Bayesian models aren't always useful. For example, when we want $A \perp\!\!\!\perp C|\{B, D\}$ and $B \perp\!\!\!\perp D|\{A, C\}$ without extraneous independencies, Bayesian model fails, due to its directions. Therefore, in undirectional PGMs, we use symmetric values to estimate the affinities between variables.

**Definition 8.** (**Factor**) *For the set $D \subset \mathcal{X}$ of rvs, a factor is $\phi : Val(D) \to \mathbb{R}$. $D$ is the scope of $\phi$, denoted Scope$[\phi]$.*

**Definition 9.** (**Factor Product**) *For disjoint sets $X, Y, Z$, suppose there are two factors $\phi_1(X, Y)$ and $\phi_2(Y, Z)$. The factor product $\psi = \phi_1 \times \phi_2 : Val(X, Y, Z) \to \mathbb{R}$ is defined by $\psi(X, Y, Z) = \phi_1(X, Y) \cdot \phi_2(Y, Z)$.*

Here factors describe local relationships between variables, which resembles the conditional probability used in Bayesian model. And similarly to the case of Bayesian model, we define $P(X_1 \sim X_n)$ as the product of multiple factors. Notably, joint distributions and conditional distributions in Bayesian models are both instances of factors.

**Definition 10.** (**Gibbs Distribution**) *A distribution $P_\Phi$ is a Gibbs distribution parameterized by factors $\Phi = \{\phi_1(D_1) \sim \phi_K(D_K)\}$ if $P_\Phi(X_1 \sim X_n) = \frac{1}{Z}\prod_{j=1}^m \phi_j(D_j)$ with $Z$ as a normalizing constant.*

**Definition 11.** (**Factorization**) *A Gibbs distribution $P_\Phi$ factorizes over a Markov network $\mathcal{H}$ if $\forall k \in [K]$, $D_k$ is a complete subgraph of $\mathcal{H}$. Here the factors parameterizing the Markov network are **clique potentials**.*

Since every complete subgraph $D_k$ is the subset of some maximal clique, we can let restrict for every factor that the scope can only be maximal subgraphs. However, in this case, the required number of parameters might exceed the original representation. For instance, a complete graph encoding pairwise interactions requires $O(n^2)$ entries, but maximal clique potentials require $O(2^n)$ entries.

**Definition 12.** (**Factor Reduction**) *To model conditional probabilities $\mathbf{Pr}[...|U = u]$ where $U \subseteq \mathcal{V}$, the reduction $\phi[U = u]$ or $\phi[u]$ with $U \subseteq Y$ is $\phi[u](y') = \phi(y', u)$; while for $U \nsubseteq Y$, $\phi[u] = \phi[U \cap Y = u']$.*

**Definition 13.** (**Reduced Gibbs Distribution**) *For $P_\Phi$ parameterized by $\Phi = \{\phi_1 \sim \phi_K\}$ and context $u$, the reduced Gibbs distribution is $P_\Phi[u]$ defined by $\Phi[u] = \{\phi_1[u] \sim \phi_K[u]\}$.*

**Definition 14.** (**Reduced Markov Network**) *$\mathcal{H}$ is a Markov network over $\mathcal{V}$ and $U = u$ is a context. The reduced Markov network is the induced subgraph $\mathcal{H}[u]$ from $\mathcal{H}$ by $\mathcal{V} - U$.* Here conditioning eliminates edges, while in Bayesian networks conditioning also creates edges.

**Definition 15.** (**Blocking**) *For a Makov network structure $\mathcal{H}$, a path $p$ is blocked by $Z \subseteq \mathcal{V}$ if there exists $X_i \in p \cap Z$.*

**Definition 16.** (**Separation**) *$Z$ separates $X$ and $Y$ in $\mathcal{H}$ (denoted sep$_\mathcal{H}(X;Y|Z)$) if all paths between $X$ and $Y$ are blocked by $Z$. The corresponding global independency assertions are $\mathcal{I}(\mathcal{H}) = \{(X \perp\!\!\!\perp Y|Z)|sep_\mathcal{H}(X;Y|Z)\}$.*

**Theorem 4.** *(Soundness of separation)* *For distribution $P$ and Markov network structure $\mathcal{H}$, if $P$ is a Gibbs distribution factorizing over $\mathcal{H}$, then $\mathcal{I}(\mathcal{H}) \subseteq \mathcal{I}(P)$.*

*Proof.* Soundness is equivalent to saying $\text{sep}_{\mathcal{H}}(X;Y|Z) \Rightarrow X \perp\!\!\!\perp Y|Z$ under the condition of factorization.

**Lemma 5.** *Suppose $\mathcal{X} = X \cup Y \cup Z$, then $X \perp\!\!\!\perp Y|Z$ iff. $P(\mathcal{X}) = \phi_1(X, Z)\phi_2(Y, Z)$.* *(We should observe that when $Z$ is not conditioned on, this decomposition won't lead to independence.)*

   Using the lemma, we first prove soundness with assumption $\mathcal{X} = X \cup Y \cup Z$. Suppose $(X, Y) \in X \times Y$ satisfies $\{X, Y\} \in D_j$, then $X$ and $Y$ are not separated because $D_j$ is a complete subgraph. Consequently, any $D_j$ contains either $X \cup Z$ or $Y \cup Z$. Then we can group $D_j$ into those don't contain $Y$ (and thus $D_j \subseteq X \cup Z$) and those don't contain $X$, which allows for a decomposition into $P(\mathcal{H}) = f(X, Z)g(Y, Z)$.

   When $\mathcal{X} = X \cup Y \cup Z \cup U$, then there doesn't exist a path passing through $U$ while not passing through $Z$. In the induced subgraph of $\mathcal{H}$ by $X \cup Y \cup U$, $X$ and $Y$ are not connected. Consequently, there exists partition $U = U_1 \cup U_2$ such that in $\mathcal{H}$, there doesn't exist a path from $X \cup U_1$ to $Y \cup U_2$ not passing through $Z$. Then the lemma is applied for $X \cup U_1$ and $Y \cup U_2$ given $Z$. $\qquad\square$

**Theorem 6.** *(Completeness of separation)* *If $X$ and $Y$ are not separated given $Z$, then they are conditionally dependent in some distribution $P$ factorizing over $\mathcal{H}$.*

*Proof.* Since $X$ and $Y$ are not separated, we can find a minimal path $p : X = U_0 \to U_1 \to ... \to U_k = Y$ where minimal means $\forall j \neq i \pm 1$, there doesn't exist an edge $U_i - U_j$, then any clique $C_i$ containing $U_i - U_{i+1}$ must be different. Without loss of generality, we assume all variables are binary-valued, and $\phi_i(C_i) = W$ (very large) when $u_i = u_{i+1}$, otherwise $\phi_i(C_i) = 1$. Then the constructed probability yields dependent $X$ and $Y$. $\qquad\square$

   Similar to the Bayesian models, we can also strengthen the conclusion to be: for almost all distributions except for some measure-zero distributions $P$ that factorizes over $\mathcal{H}$, we have $\mathcal{I}(P) = \mathcal{I}(\mathcal{H})$.

## 2.1   Converting Bayesian Networks into Markov Networks

   There are two approaches, one given probability $P_{\mathcal{B}}$, we represent $P_{\mathcal{B}}$ as a Markov network (find minimal I-map for distribution); the other given graph $\mathcal{G}$, we represent independencies in $\mathcal{G}$ using undirected graph $\mathcal{H}$ (find minimal I-map for independencies).

   When given $P_{\mathcal{B}}$ parameterized over $\mathcal{G}$, this $P_{\mathcal{B}}$ is exactly a Gibbs distribution. Moreover, conditioning on evidence $E = e$, the distribution is exactly the reduced Gibbs distribution with context $E = e$.

   When building Markov network from independencies $\mathcal{G}$, we want for each $X_i$ to be connected to $\mathbf{pa}_{X_i}$ and $\mathbf{pa}_{X_i}$ forms a complete subgraph. The resulting graph is called moralized graph.

**Definition 17.** (**Moralized Graph**) *The moral graph $\mathcal{M}[\mathcal{G}]$ is the undirected graph containing an edge $(X, Y) \in \mathcal{E}_{\mathcal{M}[\mathcal{G}]}$ if $(X, Y) \in \mathcal{E}_{\mathcal{G}}$ or $(Y, X) \in \mathcal{E}_{\mathcal{G}}$ or $X, Y$ are the parents of the same node.*

**Theorem 7.** *The moralized graph $\mathcal{M}[\mathcal{G}]$ is a minimal I-map (thus expressing maximal independencies) for $\mathcal{G}$.*

*Proof.* Consider the **Markov blanket** (the parents, children and other parents of children) of $X$, denoted $\text{MB}_{\mathcal{G}}(X)$. Then $\text{MB}_{\mathcal{G}}(X)$ is exactly the set d-separating $X$ from other variables in $\mathcal{G}$ (no subset has this property, requiring parents of children because conditioning on common children opens the path). $\qquad\square$

**Definition 18.** (**Induced Subgraph of Upward Closure**) *For $U \subseteq V$ in $\mathcal{G}$, $U^* = U \cup \text{Ancestors}_U$ is the upward closure. $\mathcal{G}^+[U]$ is the induced subgraph from $\mathcal{G}$ by $U^*$.*

**Theorem 8.** *For disjoint nodes $X, Y, Z$ in a Bayesian network $\mathcal{G}$, if $U = X \cup Y \cup Z$, let $\mathcal{H} = \mathcal{M}[\mathcal{G}^+[U]]$, then $\text{d-sep}_{\mathcal{G}}(X;Y|Z) \Leftrightarrow \text{sep}_{\mathcal{H}}(X;Y|Z)$.*

*Proof.* First assume $\text{d-sep}_{\mathcal{G}}(X;Y|Z)$ and $\neg\text{sep}_{\mathcal{H}}(X;Y|Z)$, then all trails connecting $X$ and $Y$ in $\mathcal{G}$ are blocked by $Z$ but there exists a path $p_0 \in \mathcal{H}$ satisfying $X$ is connected to $Y$ without passing $Z$. There are two possibilities, either being a new edge that doesn't exist in $\mathcal{G}$ or $X \leftrightsquigarrow A \to B \leftarrow C \leftrightsquigarrow Y$ with $B$'s descendants not in $Z$. The latter case implies that $B$ is an ancestor of $X$ (otherwise $B \notin \mathcal{H}$), which means the trail $X \leftarrow B \leftarrow C \leftrightsquigarrow Y$ is blocked by $Z$. This procedure is done infinitely, which is clearly impossible. The former case is $X \leftrightsquigarrow A - B \leftrightsquigarrow Y$ with $A$ and $B$ being parents of $C$, which is reduced to the latter case. Contradiction!

Then assuming $\text{sep}_{\mathcal{H}}(X;Y|Z)$, every path in $\mathcal{H}$ passes through $Z$. Moreover, all $A \to Z \leftarrow B$ are illegal without other presence of $Z$ in any path because it creates shortcut $A - B$ in $\mathcal{H}$. Consequently, d-sep$_{\mathcal{G}^+[U]}(X;Y|Z)$ holds. However, if there exists a trail passing through $G - G^+[U]$, then it is $X \leftrightsquigarrow A \to p' \leftarrow B \leftrightsquigarrow Y$ (otherwise vertices in $p'$ is in $\mathcal{H}$). Therefore there always exist $A' \to B' \leftarrow C'$ in $p'$, which is blocked. Therefore, d-sep$_{\mathcal{G}}(X;Y|Z)$ holds. □

## 2.2 Converting Markov Networks into Bayesian Networks

**Definition 19.** (**Chordal Graph**) *For any loop $X_1 - X_2 - ... - X_k - X_1$, a chord is an edge connecting nonconsecutive $X_i$ and $X_j$. An undirected graph $\mathcal{H}$ is chordal if any loop with $k \geq 4$ has a chord.*

**Theorem 9.** *Suppose $\mathcal{G}$ is minimal I-map for $\mathcal{H}$, then $\mathcal{G}$ can have no immoralities, and consequently $\mathcal{G}$ is chordal.* (In inversed conversion, we have also eliminated immoralities)

The conclusion is that only chordal graphs can represent independencies without loss.

# 3 Causal Model

Causal chains typically have multiple latent variables that aren't observed in reality. In purely probabilistic graphical models, the correlation contributed by latent variables can always be represented by simply calculating the joint distribution or conditional distribution, despite leading to more complicated models. However, in causal models, causal relationships and correlations due to confounding factors are significantly different.

An important type of latent variables is **selection bias**. For example, when trying to model the causality between athletic activity participation and high GPA, one might easily find a negative correlation. However, the reason might be the student did well in neither sports nor grades tend not to respond the survey.

**Definition 20.** (**Exogenous Variables**) *The set of variables that are determined outside the model and act as the driving force of the model.*

**Definition 21.** (**Endogenous Variables**) *The set of variables that are influenced by other variables inside the model.*

**Definition 22.** (**Cause / Direct Cause**) *$X$ is a cause of $Y$ if $X$ is a direct cause of $Y$ or the cause of any cause of $Y$.*

A **structural causal model** (SCM) has two sets of variables, exogeous variables $U$, endogenous variables $V$ and a set of stochastic functions $f$ used to determine the values of endogenous variables based on other variables (when variables in $U$ are determined, $f$ gives all distributions in $V$). Each SCM is associated with a **graphical causal model** represented by a Bayesian network structure $G = (\mathcal{V}, \mathcal{E})$ where $(X, Y) \in \mathcal{E}$ represents direct causality. Consequently, in a graphical causal model $G$, $\mathbf{Pr}[x_1, ..., x_n] = \prod_i \mathbf{Pr}[x_i|pa_i]$ where $X_1, ..., X_n$ are variable nodes.

Usually the goal of statistical studies are to study the effect of **interventions** (through random controlled experiments), but such experiments aren't always practical. Intervening on a variable $\mathbf{Pr}[Y = y|do(X = x)]$ and conditioning on a variable $\mathbf{Pr}[Y = y|X = x]$ are different in that intervening changes the system, through which other variables are changed, however in conditioning, we are simply narrowing the focus to a subset of the system. Specifically, intervention changes the graphical causal model while conditioning doesn't. Mathematically, intervening a set of variables $S = s$ are expressed as

$$\mathbf{Pr}_m[X_i|pa_i] = \begin{cases} \mathbf{Pr}[X_i|pa_i], & X_i \notin S \\ 1, & X_i \in S \wedge x_i = s_i \\ 0, & X_i \in S \wedge x_i \neq s_i \end{cases}$$

Considering the existence of latent variables, in most cases we can obtain only marginal distribution over observable variables. Consequently, the problem lies in finding out what intervention queries are **identifiable** using only conditional probabilities from observed variables.

## 3.1 Identifying Causality

**Definition 23.** (**Augmented Causal Model**) *In $\mathcal{G}$ we define for every variable $Z$ the **decision variable** $\hat{Z} \in \{\epsilon\} \cup Val(Z)$ determining whether we intervene $Z$ and the intervention value. Then resulting augmented graph is $\mathcal{G}^\dagger$. $\mathcal{G}^\dagger_{\bar{Z}}$ is the graph obtained from $\mathcal{G}^\dagger$ with every $Z \in \mathbf{Z}$ having only parent $\hat{Z}$.*

**Theorem 10.** (***Intervention Query Simplification Rule 1***) *For a causal model over $\mathcal{G}$, $Pr\,[Y|do(Z), X, W] = Pr\,[Y|do(Z), X]$ if d-sep$_{\mathcal{G}^\dagger_{\bar{Z}}}(W; Y|Z, X)$.*

*Proof.* Trivial. $\mathbf{Pr}_{\mathcal{G}}\,[Y|do(Z), X] = \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{Z}}}\,[Y|Z, X]$. $\qquad\square$

**Theorem 11.** (***Intervention Query Simplification Rule 2***) *For a causal model over $\mathcal{G}$, $Pr\,[Y|do(Z), do(X), W] = Pr\,[Y|do(Z), X, W]$ if d-sep$_{\mathcal{G}^\dagger_{\bar{Z}}}(\hat{X}; Y|Z, X, W)$.*

*Proof.* Notice that the d-separation condition doesn't hold when there exists causal chain $Y \rightsquigarrow X$. D-separation $\Rightarrow$ all trails $p : \hat{X} \to X \leftrightsquigarrow Y$ are blocked by $X$. Therefore, $p : \hat{X} \to X \to ... \to U \leftarrow V \leftrightsquigarrow Y$. Since $\mathbf{Pr}\,[V|X] = \frac{1}{\mathbf{Pr}[X]}\sum_u \mathbf{Pr}\,[V]\,\mathbf{Pr}\,[X]\,\mathbf{Pr}\,[U|V, X] = \mathbf{Pr}\,[V]$, all marginal distributions of $Y$'s ancestors (when $X$ are observed) are unchanged compared to intervening $X$ (where those marginal distributions are naturally unchanged because none of those nodes are descendants of $X$). $\qquad\square$

**Theorem 12.** (***Intervention Query Simplification Rule 3***) *For a causal model over $\mathcal{G}$, $Pr\,[Y|do(Z), do(X), W] = Pr\,[Y|do(Z), W]$ if d-sep$_{\mathcal{G}^\dagger_{\bar{Z}}}(\hat{X}; Y|Z, W)$.*

*Proof.* Similarly, all trails connecting $\hat{X}$ and $Y$ are blocked by $\varnothing$. That is, $p : \hat{X}... \to U \leftarrow ...Y$ (here $U$ can be $X$). The same reasoning shows that changing causal mechanism of $X$ changes nothing about $Y$. $\qquad\square$

Notably, despite being complete (and thus stated as three axioms in **do-calculus**), these three rules have limited applicability, while the more general rules apply under **back-door criterion** or **front-door criterion**.

**Definition 24.** (**The Back-Door Criterion**) *Given $X, Y$ in $\mathcal{G}$, $\mathbf{Z}$ satisfies back-door criterion relative to $X, Y$ if no node in $\mathbf{Z}$ is a descendant of $X$ in $\mathcal{G}$, and $\mathbf{Z}$ blocks every trail between $X$ and $Y$ that points to $X$ (these trails are called **back-door paths**). (Intuitively, $\mathbf{Z}$ blocks all spurious paths (those with arrows into $X$) and creates no new spurious paths (by not including descendants of $X$).)*

**Theorem 13.** *If $\mathbf{Z}$ satisfies back-door criterion relative to $X, Y$, then*

$$Pr\,[Y = y|do(X = x)] = \sum_z Pr\,[Y = y|X = x, Z = z]\,Pr\,[Z = z]$$

*Proof.*

$$\mathbf{Pr}_{\mathcal{G}}\,[Y|do(X)] = \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|\hat{X} = do(x)] = \sum_z \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|\hat{X} = do(x), Z = z]\,\mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Z = z|\hat{X} = do(x)]$$

Since $\mathbf{Z}$ are not descendants of $X$, there must exist a collider in the trail $p : \hat{X} \to X \leftrightsquigarrow Z$, or equivalently, d-sep$_{\mathcal{G}^\dagger_{\bar{X}}}(\hat{X}; Z|\varnothing)$. Therefore, applying the 3rd rule, $\mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Z = z|\hat{X} = do(x)] = \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Z = z] = \mathbf{Pr}_{\mathcal{G}}\,[Z = z]$.

Meanwhile, $\mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|\hat{X} = do(x), Z] = \sum_{x'} \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|X = x', \hat{X} = do(x), Z] = \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|X = x, \hat{X} = do(x), Z]$.

Note that $\hat{X} \perp\!\!\!\perp Y|X, Z$ because all trails connecting $\hat{X}$ and $Y$ have either $X$ as a collider (blocked by $Z$) or blocked by $X$ itself. Consequently,

$$\mathbf{Pr}_{\mathcal{G}}\,[Y|do(X = x)] = \sum_z \mathbf{Pr}_{\mathcal{G}^\dagger_{\bar{X}}}\,[Y|X = x, \hat{X} = \text{idle}, Z]\,\mathbf{Pr}_{\mathcal{G}}\,[Z] = \sum_z \mathbf{Pr}_{\mathcal{G}}\,[Y|X = x, Z]\,\mathbf{Pr}_{\mathcal{G}}\,[Z]$$

$\qquad\square$

Back-door adjustment are not applicable for $X \leftarrow U \to Y$ and $X \to Y$ because the backdoor path $X \to U \to Y$ can't be blocked by anything. However, if we introduce another observed variable $Z$ intercepting the edge $X \to Y$, we can measure the causality $\mathbf{Pr}\,[Y|do(X = x)]$.

**Definition 25.** (**The Front-Door Criterion**) $Z$ *satisfies the front-door criterion relative to* $X, Y$ *if* **(1)** $Z$ *intercepts all paths* $p : X \rightsquigarrow Y$; **(2)** *there is no unblocked back-door path from* $X$ *to* $Z$; **(3)** *all back-door paths from* $Z$ *to* $Y$ *are blocked by* $X$.

**Theorem 14.** *If* $Z$ *satisfies the front-door criterion relative to* $X, Y$ *and* $Pr[x, z] > 0$, *then*

$$Pr[Y|do(X = x)] = \sum_z Pr[Z = z|X = x] \sum_{x'} Pr[Y|X = x', Z = z] Pr[X = x']$$

*Proof.*

$$\mathbf{Pr}_{\mathcal{G}}[Y|do(X = x)] = \sum_z \mathbf{Pr}_{\mathcal{G}_{\bar{X}}^{\dagger}}[Y|Z = z] \mathbf{Pr}_{\mathcal{G}_{\bar{X}}^{\dagger}}[Z = z|X = x] = \sum_z \mathbf{Pr}_{\mathcal{G}_{\bar{Z}}^{\dagger}}[Y|Z = z] \mathbf{Pr}_{\mathcal{G}}[Z = z|X = x]$$

$$= \sum_z \mathbf{Pr}[Y|do(Z = z)] \mathbf{Pr}[Z = z|X = x]$$

$$= \sum_z \mathbf{Pr}[Z = z|X = x] \left( \sum_{x'} \mathbf{Pr}[Y|Z = z, X = x'] \mathbf{Pr}[X = x'] \right)$$

Here, the first equality requires **(1)** to omit $X = x$ as a condition in the first term. The second equality requires $\mathbf{Pr}_{\mathcal{G}_{\bar{X}}^{\dagger}}[Y|Z] = \mathbf{Pr}_{\mathcal{G}_{\bar{Z}}^{\dagger}}[Y|Z]$ and $\hat{X} \perp\!\!\!\perp Z|X$, which are satisfied by **(2)** and **(3)** respectively. □