

Decision-Making & Modelling Under Uncertainty (DMU)

Oliver Maclaren (oliver.maclaren@auckland.ac.nz)

[10 lectures / tutorials]

◦ Decision-making under uncertainty [5/10]

- ↳ Basic concepts
- ↳ Risk, probability, utility
- ↳ Statistical: extended setup
 - ↳ formulation & empirical risk approx.
 - ↳ minimax & Bayes
- ↳ Tutorial sheet

◦ Modelling under uncertainty { models of risk & intervention [5/10] }

- ↳ probability, graphical models, & independence
- ↳ causal interpretations of graphical models
- ↳ stochastic process models (esp. Markov)
- ↳ simulation & estimation tools
- ↳ Tutorial sheet

Lecture 5: Modelling under risk

→ representation & independence for probabilistic models

In the next few lectures we will look at models of, or incorporating, risk/probability

- these may represent 'frequentist' probabilities ('out there in the world') or 'belief' (Bayesian) probabilities
 - using Bayes' theorem ≠ being a 'Bayesian'
 - 'Bayesian' means using probability to represent belief
- See Wasserman →

From Larry Wasserman (author of 'All of Statistics'):

1. Some Obvious (and Not So Obvious) Statements

Before I go into detail, I'll begin by making a series of statements.

Frequentist Inference is Great For Doing Frequentist Inference.
Bayesian Inference is Great For Doing Bayesian Inference.

Frequentist inference and Bayesian Inference are defined by their goals, not their methods.

A Frequentist analysis need not have good Bayesian properties.
A Bayesian analysis need not have good frequentist properties.

Bayesian Inference \neq Using Bayes Theorem

Bayes Theorem \neq Bayes Rule

Bayes Nets \neq Bayesian Inference

Frequentist Inference is not superior to Bayesian Inference.
Bayesian Inference is not superior to Frequentist Inference.
Hammers are not superior to Screwdrivers.

Confidence Intervals Do Not Represent Degrees of Belief.
Posterior Intervals Do Not (In General) Have Frequency Coverage Properties.

Saying That Confidence Intervals Do Not Represent Degrees of Belief Is Not a Criticism of Frequentist Inference.
Saying That Posterior Intervals Do Not Have Frequency Coverage Properties Is Not a Criticism of Bayesian Inference.

Some Scientists Misinterpret Confidence Intervals as Degrees of Belief.
They Also Misinterpret Bayesian Intervals as Confidence Intervals.

Mindless Frequentist Statistical Analysis is Harmful to Science.
Mindless Bayesian Statistical Analysis is Harmful to Science.

See :

<https://normaldeviate.wordpress.com/2012/11/17/what-is-bayesianfrequentist-inference/>

→ The tools we discuss may be
useful for either (or something else!)

Probability models : representation & computation

[See summary sheet for background]

→ In principle a risk/probability model can be captured by a joint distribution over a collection of random variables

(see later for stochastic process models)

→ However, in practice it is difficult to represent & compute directly with joint distributions

Aside: densities & mass functions

- For discrete random variables, each outcome has finite probability & so eg $P(X=1)$ makes sense
 - The function giving the probabilities of each outcome, $P(X=x)$, is the probability mass function (pmf) & we write $\boxed{P(x) = P(X=x)}$
 - For continuous random variables each outcome typically only has 'infinitesimal' probability & for finite probabilities we have to consider eg intervals: $P(1 \leq X \leq 2)$ etc.
- This is given by $\int_1^2 p(x) dx$ etc.
for probability density function (pdf) $p(x)$ (or $p_x(x)$ or $f_x(x)$)

Aside: densities & mass functions

we can relate these considering (semi-formally) the probability mass of a continuous outcome as $\underset{x \in \text{dx}}{\sim} \int_p(x) dx$

$$\sim \text{pdf differential}$$

The discrete case is of course just:

$$\sim \underset{\text{pmf}}{P(X=x)} = P(x)$$

Importantly: In our context, most key relationships can be expressed in terms of ' $p(x)$ ' for both continuous & discrete vars, regardless of whether discrete (pmf) or continuous (pdf).

→ We will hence call both of these 'probability functions' & let context decide if $p(x)$ is a pdf or pmf.

Illustration : representation troubles

Consider N binary random variables

$$x_i, i=1, \dots, N$$

i.e each x_i can take values in $\{0, 1\}$

The joint distribution over these

is given by

$$P(x_1=x_1, \dots, x_N=x_N) = P(x_1, x_2, \dots, x_N)$$

which requires specifying 2^N

possible probabilities (one per input combination)

→ hard as N grows!

Thus, even representation (writing down a model) is hard for probability models with many variables

→ motivates →

Independence for 'modularity'

- Random variables X & Y are called independent if

$$\boxed{P(x, y) = P(x)P(y)}$$

\Leftrightarrow

$$\boxed{P(y|x) = P(y)}$$

\Leftrightarrow

$$\boxed{P(x|y) = P(x)}$$

(recall
 $P(x=x) = P(x)dx$
notation)

We write this as $\boxed{Y \perp\!\!\!\perp X}$ which is equivalent to $\boxed{X \perp\!\!\!\perp Y}$

- This also applies to sets of random vectors like $\{x_1, x_2, \dots, x_n\}$

→ we can write $x_1 \perp\!\!\!\perp \{x_2, \dots, x_n\}$ etc to express that x_1 is independent of the other RVs

→ Independence assumptions make models more modular & hence tractable (see illustration...but first →)

Mutual independence of x_1, x_2, \dots, x_n

is expressed as:

$$x_i \perp\!\!\!\perp \{x_j\}_{j \neq i} \text{, for all } x_i$$

where $\{x_j\}_{j \neq i}$ is the set of
all RVS except x_i .

Random variables X & Y are called

'conditionally independent given Z '

$$\text{if } | P(y|x, z) = P(y|z) |$$

or, equivalently:

$$| P(x|y, z) = P(x|z) |$$

or, equivalently:

$$| P(x, y, z) = P(x|z)P(y|z)P(z) |$$

This is written: $| X \perp\!\!\!\perp Y | z |$

$$(\Leftrightarrow Y \perp\!\!\!\perp X | z)$$

We can also consider e.g. $X \perp\!\!\!\perp Y | \{z_1, \dots, z_n\}$

$\rightarrow X$ & Y independent given the set
of random vars z_1, \dots, z_n

Modularity: illustration

Suppose in our example that the N
binary variables are mutually independent

then:

$$P(x_1, x_2, x_3, \dots, x_N)$$

=

$$P(x_1) \cdot P(x_2) \cdot \dots \cdot P(x_N)$$

=

$$\prod_{i=1}^N P(x_i)$$

product version of sum $\sum_{i=1}^n$

\rightarrow only need $2N$ values to be given

instead of 2^N . Much less when N big!

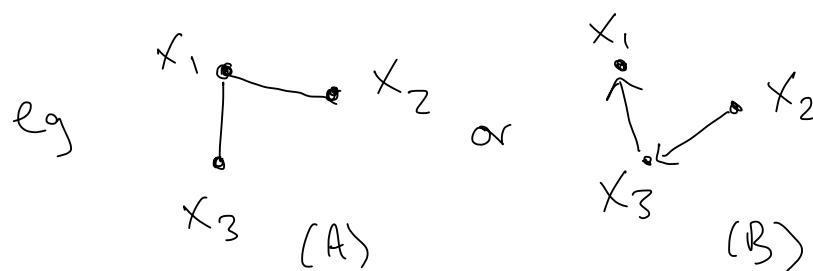
\rightarrow While mutual independence is rare, in practice
conditional independencies are not!

\rightarrow also leads to computational improvements

\rightarrow also conceptual/modelling advantages

Representing conditional independences

- Graphical models are a popular way to represent conditional independencies
- These use graphs consisting of vertices representing variables & edges representing (in)dependencies

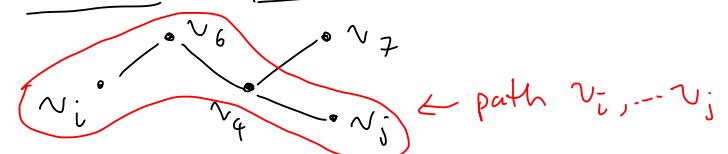


- Two common types of graph are
 - {'undirected': edges have no direction (A)}
 - {'directed': edges have a direction (B)}

Aside:

- A graph $G = (V, E)$ is defined formally as a structure containing
 - a set of vertices V &
 - a set of edges E between vertices.

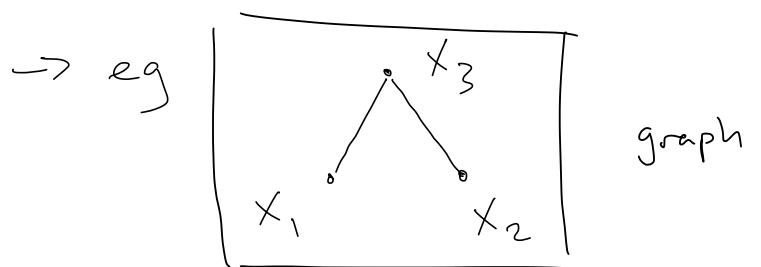
- We also call vertices 'nodes'
- We can consider an edge $e \in E$ as a tuple (v_1, v_2) for some $v_1, v_2 \in V$
- A path between two vertices v_i, v_j is a sequence of vertices v_i, \dots, v_j , starting at v_i & ending at v_j , such that there is an edge between each consecutive pair of vertices



Undirected graphs : independence interpretation

These represent purely probabilistic (cf causal-see later) independencies using the rule:

Given a set of variables X_1, \dots, X_N , we do not draw an edge between a pair X_i & X_j if $\boxed{X_i \perp\!\!\!\perp X_j \mid \text{(the rest)}}$



represents $\boxed{X_1 \perp\!\!\!\perp X_2 \mid X_3}$ independence

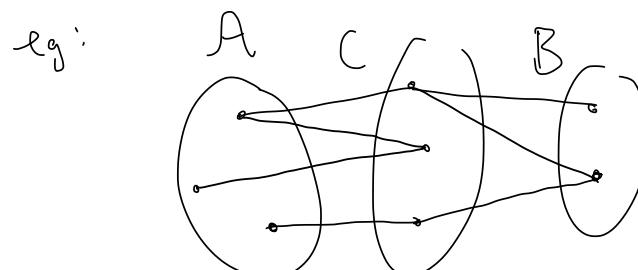
i.e. $\boxed{P(X_1 \mid X_2, X_3) = P(X_1 \mid X_3)}$ simplification of prob. function.
etc

Undirected separation (u-separation)

The graph property of u-separation relates sets of vertices (nodes) in an undirected graph as follows:

Given three sets of vertices $A, B, C \subseteq \mathcal{V}$, we say C separates A & B if

all paths from elements of A to B pass through C



Sets of variables & independence:

The graph property of separation translates to the probabilistic property of independence:

If:

$\{z_1, \dots, z_p\}$ separates $\{x_1, \dots, x_n\}$ & $\{y_1, \dots, y_m\}$

then:

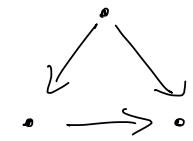
$$\{x_1, \dots, x_n\} \perp\!\!\!\perp \{y_1, \dots, y_m\} \mid \{z_1, \dots, z_p\}$$

Called the 'global Markov property'

What about directed graphs?



Directed graphical models



- Also called 'Bayesian networks'
(misnomer: not really 'Bayesian')
- Designed to incorporate information about both 'causality', &/or directional ideas, & probability

Recall: 'correlation \neq causation'

\uparrow
probability
concept \uparrow
not probability
concept ?

- First we will look at how they encode (probabilistic) conditional independence

\hookrightarrow different to undirected graphs

\hookrightarrow can use just as alternative tool for representing conditional independencies

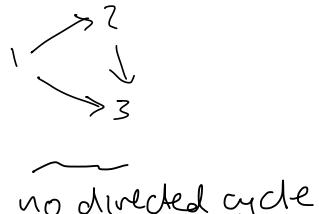
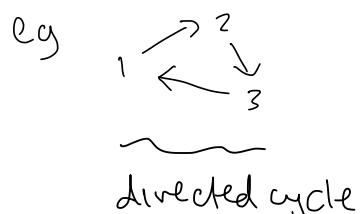
\hookrightarrow or, also as causal modelling tool

Directed graphical models

- Edges now have directions (are arrows) : \rightarrow
- we will restrict attention to acyclic directed graphs or
 $\boxed{\text{DAGs}}$ (directed acyclic graphs)

DAG? No directed cycles, where:

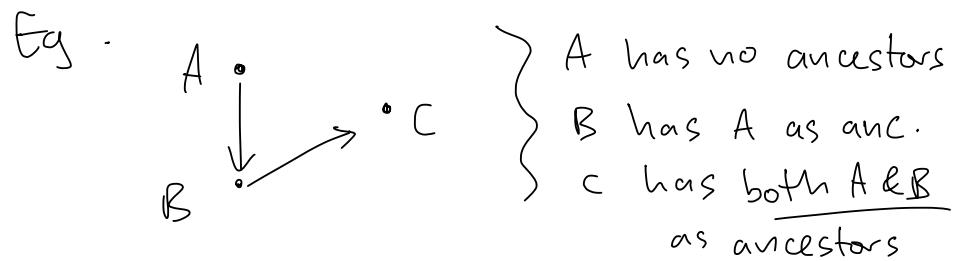
- A cycle is a path from a node back to itself
- A directed cycle is a cycle of directed edges where they all have the 'same' direction along the path



Directed independence?

The idea of a DAG is to encode a type of conditional independence representing an 'ordered Markov factorisation' of the probability model

First define a 'predecessor' or 'ancestor' of a vertex v as any other vertex occurring 'before' v on a directed path connecting both vertices



A parent is an immediate ancestor (single edge path)

→ Exercise: Define descendants & children in obvious way!

Given these definitions, a DAG directly encodes conditional independencies of the form:

$$\boxed{x \perp \text{Pred}(x) \setminus \text{Pa}(x) \mid \text{Pa}(x)}$$

where

$\text{Pred}(x)$	= 'Predecessors of x '
$\text{Pa}(x)$	= 'Parents of x '
$\text{Pred}(x) \setminus \text{Pa}(x)$	= 'Non parent predecessors'

i.e.

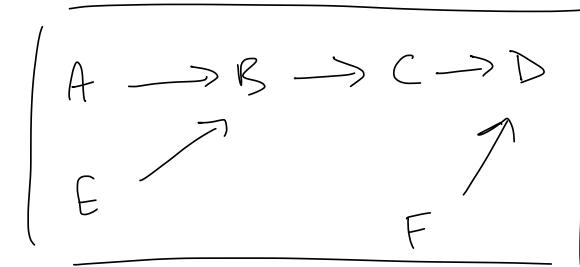
$\boxed{'x \text{ is independent of its non-parent predecessors, given its parents}'}$

(only the immediate predecessors, i.e. parents, matter)

$$\Rightarrow \boxed{p(x \mid \text{pred}(x)) = p(x \mid \text{pa}(x))}$$

Where $\text{pa}(x)$ stands for 'values of parents of x ' etc.

Eg :



Encodes independencies such as:

$$\boxed{c \perp \{a, e\} \mid b}$$

i.e.

$$\boxed{p(c \mid a, b, e) = p(c \mid b)}$$

etc.

Chain rule of probability

Given a collection of N random variables and any ordering of them: x_1, x_2, \dots, x_N

We can always write the joint distribution as

$$P(x_1, x_2, \dots, x_n) =$$

$$\boxed{P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1} | x_{n-2}, \dots, x_1) \dots P(x_2 | x_1) P(x_1)}$$

$$\begin{aligned} \text{eg } P(x, y, z) &= P(x|y, z) P(y|z) P(z) \\ &= P(y|x, z) P(x|z) P(z) \\ &= P(z|x, y) P(x|y) P(y) \\ &\vdots \\ &\text{etc.} \end{aligned}$$

Product decomposition

Combining the chain rule & the conditional independences implied by a DAG gives, for a set of N random variables x_1, \dots, x_N gives:

$$\boxed{P(x_1, x_2, \dots, x_N) = \prod_{i=1}^N P(x_i | \text{pa}(x_i))}$$

Notes:

- o this does not rely on x_1, x_2, \dots, x_N being 'in order' according to the DAG
→ this is why there is the ' $\text{pa}(x_i)$ ' in the expression
- o In general this allows us to specify the joint distribution as a product of simpler distributions similar to the case of mutual independence

Other independencies in DAGs

A given DAG tells us 'directly' about the dependencies of a node on its ancestors (past nodes)

→ It doesn't 'directly' tell us about the (probabilistic) dependencies of a node on its descendants!

$$\text{e.g. } X \rightarrow Y$$

tells us Y depends on X

$$\& \quad P(x, y) = P(y|x)P(x)$$

$\uparrow \quad \uparrow$
pa(y) no parents

But $P(x, y) = P(x|y)P(y)$ always] Prob. theory

$$\& \quad P(x|y) = \frac{P(y|x)P(x)}{P(y)} \neq P(x)$$

⇒ X is not indep. of Y in $X \rightarrow Y$

Causal vs probabilistic dependence?

- The DAG interpretation so far is purely probabilistic, not 'causal'
- We can add interpretational components so that $X \rightarrow Y$ also captures 'X causes Y'...
 BUT for now just focus on the probabilistic conditional (in)dependencies implied by a DAG

→ combine ordered Markov factorisation & probability theory
 --- OR use graph theory!

Further conditional independencies

implied by a DAG: directed separation

To determine the other independencies implied by the directed Markov factorisation (encoded by the DAG) & the rules of probability, we can consider the graphical concept of

| 'd-separation' |

(for directed separation)

For this we will call any sequence of edges between two nodes, regardless of direction, a dependence path or just 'path':

$A \rightarrow B \leftarrow C \quad \{ A \rightarrow B, B \leftarrow C \text{ is a path}$

d-separation

{more complex than u-separation unfortunately! }

Definition: | path blocking |

A (dependence) path between any two nodes X, Y in a DAG is 'blocked' by a set of nodes B iff the path contains at least one sub-path ('junction') of the form:

1. $A \rightarrow B \rightarrow C$ (chain) where $B \in B$

or 2. $A \leftarrow B \rightarrow C$ (fork) where $B \in B$

or 3. $A \rightarrow D \leftarrow C$ (collider) where $D \notin B$

& $\text{des}(D) \not\subseteq B$

↑
descendants of D

d-separation:

If a set of nodes B blocks every path between two nodes X & Y then X & Y are said to be 'd-separated given B '

→ According to the standard 'directed Markov' interpretation of a DAG & standard probability theory, we have:

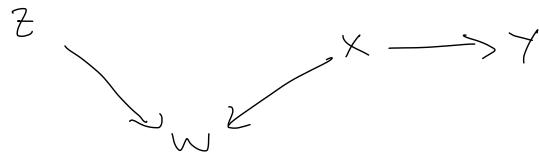
$$\begin{array}{|c|} \hline X \perp\!\!\!\perp Y \mid B \\ \hline (\Rightarrow) \\ \hline X \& Y \text{ d-separated given } B \\ \hline \end{array}$$

→ Furthermore:

d-separation can derive all the conditional independencies implied by the DAG!

Example

Consider:



- What can we say about the (in)dependence of Z & Y ?

→ we can 'condition' on the empty set to get 'unconditional' independence properties

→ Note that the dependence path between Z & Y contains the sub path $Z \rightarrow W \leftarrow X$, i.e. a 'collider'

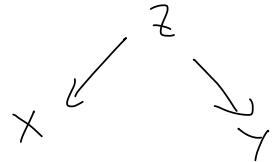
→ Hence $Z \perp\!\!\!\perp Y \mid \{\} \Leftrightarrow Z \perp\!\!\!\perp Y$

But $Z \not\perp\!\!\!\perp Y \mid \{W\}$ 'unblocks' path

Note: If two nodes are not d-separated, they are called d-connected & are (likely) 'not independent', i.e. ' $X \not\perp\!\!\!\perp Y$ '

Exercises

1. Consider



which of these is true:

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

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

2. Suppose in (1) we are

sampling from a population
& measuring variables:

Z = 'age'

X = 'retirement savings'

Y = 'number of wrinkles'

What can you say about the
relationship between retirement
savings & wrinkles?

How should you analyse this relationship
to understand the 'causal effects'?

3. Consider $A \rightarrow B \rightarrow C$, $A \leftarrow B \rightarrow C$

$$\& A \leftarrow B \leftarrow C$$

- Derive the implied conditional
independencies using d-separation

- What do you notice?

4. Consider: $A \rightarrow B \leftarrow C$

$$\text{Is } A \perp\!\!\!\perp C ?$$

$$\text{Is } A \perp\!\!\!\perp C \mid B ?$$

Suppose A represents 'height'
B represents 'basketball success'
C represents 'basketball shooting
skill'

In the above.

Interpret the associated conditional
independencies.

5. Further reading:

d-separation & associated
independence properties can be
extended to sets of variables

→ see Wasserman (2004) p.270
(attached)

Wasserman (2004) 'All of statistics'

17

Directed Graphs and Conditional
Independence

6. Further reading:

Look up 'Simpson's paradox'
& its relation to DAGs.

17.1 Introduction

A directed graph consists of a set of nodes with arrows between some nodes.
An example is shown in Figure 17.1.

Graphs are useful for representing independence relations between variables.
They can also be used as an alternative to counterfactuals to represent causal
relationships. Some people use the phrase **Bayesian network** to refer to a
directed graph endowed with a probability distribution. This is a poor choice
of terminology. Statistical inference for directed graphs can be performed using

7. Further reading:

'd-separation without tears' at:

<http://bayes.cs.ucla.edu/BOOK-2K/d-sep.html>

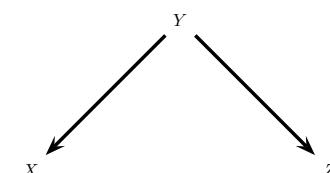


FIGURE 17.1. A directed graph with vertices $V = \{X, Y, Z\}$ and edges $E = \{(Y, X), (Y, Z)\}$.

frequentist or Bayesian methods, so it is misleading to call them Bayesian networks.

Before getting into details about directed acyclic graphs (DAGs), we need to discuss conditional independence.

17.2 Conditional Independence

17.1 Definition. Let X , Y and Z be random variables. X and Y are conditionally independent given Z , written $X \perp\!\!\!\perp Y | Z$, if

$$f_{X,Y|Z}(x,y|z) = f_{X|Z}(x|z)f_{Y|Z}(y|z). \quad (17.1)$$

for all x , y and z .

Intuitively, this means that, once you know Z , Y provides no extra information about X . An equivalent definition is that

$$f(x|y,z) = f(x|z). \quad (17.2)$$

The conditional independence relation satisfies some basic properties.

17.2 Theorem. The following implications hold:¹

$$\begin{aligned} X \perp\!\!\!\perp Y | Z &\implies Y \perp\!\!\!\perp X | Z \\ X \perp\!\!\!\perp Y | Z \text{ and } U = h(X) &\implies U \perp\!\!\!\perp Y | Z \\ X \perp\!\!\!\perp Y | Z \text{ and } U = h(X) &\implies X \perp\!\!\!\perp Y | (Z, U) \\ X \perp\!\!\!\perp Y | Z \text{ and } X \perp\!\!\!\perp W | (Y, Z) &\implies X \perp\!\!\!\perp (W, Y) | Z \\ X \perp\!\!\!\perp Y | Z \text{ and } X \perp\!\!\!\perp Z | Y &\implies X \perp\!\!\!\perp (Y, Z). \end{aligned}$$

17.3 DAGs

A **directed graph** \mathcal{G} consists of a set of vertices V and an edge set E of ordered pairs of vertices. For our purposes, each vertex will correspond to a random variable. If $(X, Y) \in E$ then there is an arrow pointing from X to Y . See Figure 17.1.

¹The last property requires the assumption that all events have positive probability; the first four do not.

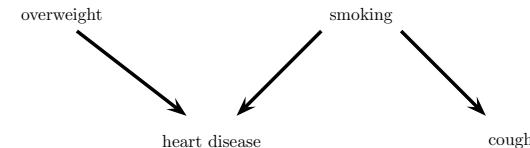


FIGURE 17.2. DAG for Example 17.4.

If an arrow connects two variables X and Y (in either direction) we say that X and Y are **adjacent**. If there is an arrow from X to Y then X is a **parent** of Y and Y is a **child** of X . The set of all parents of X is denoted by π_X or $\pi(X)$. A **directed path** between two variables is a set of arrows all pointing in the same direction linking one variable to the other such as:

$$X \longrightarrow \dots \longrightarrow Y$$

A sequence of adjacent vertices starting with X and ending with Y but ignoring the direction of the arrows is called an **undirected path**. The sequence $\{X, Y, Z\}$ in Figure 17.1 is an undirected path. X is an **ancestor** of Y if there is a directed path from X to Y (or $X = Y$). We also say that Y is a **descendant** of X .

A configuration of the form:

$$X \longrightarrow Y \longleftarrow Z$$

is called a **collider** at Y . A configuration not of that form is called a **non-collider**, for example,

$$X \longrightarrow Y \longrightarrow Z$$

or

$$X \leftarrow Y \leftarrow Z$$

The collider property is path dependent. In Figure 17.7, Y is a collider on the path $\{X, Y, Z\}$ but it is a non-collider on the path $\{X, Y, W\}$. When the variables pointing into the collider are not adjacent, we say that the collider is **unshielded**. A directed path that starts and ends at the same variable is called a **cycle**. A directed graph is **acyclic** if it has no cycles. In this case we say that the graph is a **directed acyclic graph** or **DAG**. From now on, we only deal with acyclic graphs.

17.4 Probability and DAGs

Let \mathcal{G} be a DAG with vertices $V = (X_1, \dots, X_k)$.

17.3 Definition. If \mathbb{P} is a distribution for V with probability function f , we say that \mathbb{P} is **Markov to \mathcal{G}** , or that \mathcal{G} represents \mathbb{P} , if

$$f(v) = \prod_{i=1}^k f(x_i \mid \pi_i) \quad (17.3)$$

where π_i are the parents of X_i . The set of distributions represented by \mathcal{G} is denoted by $M(\mathcal{G})$.

17.4 Example. Figure 17.2 shows a DAG with four variables. The probability function for this example factors as

$$\begin{aligned} f(\text{overweight, smoking, heart disease, cough}) \\ &= f(\text{overweight}) \times f(\text{smoking}) \\ &\times f(\text{heart disease} \mid \text{overweight, smoking}) \\ &\times f(\text{cough} \mid \text{smoking}). \blacksquare \end{aligned}$$

17.5 Example. For the DAG in Figure 17.3, $\mathbb{P} \in M(\mathcal{G})$ if and only if its probability function f has the form

$$f(x, y, z, w) = f(x)f(y)f(z \mid x, y)f(w \mid z). \blacksquare$$

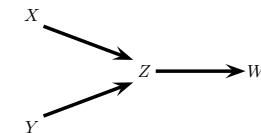


FIGURE 17.3. Another DAG.

The following theorem says that $\mathbb{P} \in M(\mathcal{G})$ if and only if the **Markov Condition** holds. Roughly speaking, the Markov Condition means that every variable W is independent of the “past” given its parents.

17.6 Theorem. A distribution $\mathbb{P} \in M(\mathcal{G})$ if and only if the following **Markov Condition** holds: for every variable W ,

$$W \perp\!\!\!\perp \widetilde{W} \mid \pi_W \quad (17.4)$$

where \widetilde{W} denotes all the other variables except the parents and descendants of W .

17.7 Example. In Figure 17.3, the Markov Condition implies that

$$X \perp\!\!\!\perp Y \quad \text{and} \quad W \perp\!\!\!\perp \{X, Y\} \mid Z. \blacksquare$$

17.8 Example. Consider the DAG in Figure 17.4. In this case probability function must factor like

$$f(a, b, c, d, e) = f(a)f(b|a)f(c|a)f(d|b, c)f(e|d).$$

The Markov Condition implies the following independence relations:

$$D \perp\!\!\!\perp A \mid \{B, C\}, \quad E \perp\!\!\!\perp \{A, B, C\} \mid D \quad \text{and} \quad B \perp\!\!\!\perp C \mid A \blacksquare$$

17.5 More Independence Relations

The Markov Condition allows us to list some independence relations implied by a DAG. These relations might imply other independence relations. Con-

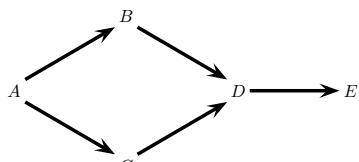


FIGURE 17.4. Yet another DAG.

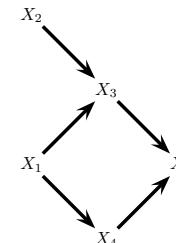


FIGURE 17.5. And yet another DAG.

sider the DAG in Figure 17.5. The Markov Condition implies:

$$\begin{aligned} X_1 \perp\!\!\!\perp X_2, \quad X_2 \perp\!\!\!\perp \{X_1, X_4\}, \quad X_3 \perp\!\!\!\perp X_4 \mid \{X_1, X_2\}, \\ X_4 \perp\!\!\!\perp \{X_2, X_3\} \mid X_1, \quad X_5 \perp\!\!\!\perp \{X_1, X_2\} \mid \{X_3, X_4\} \end{aligned}$$

It turns out (but it is not obvious) that these conditions imply that

$$\{X_4, X_5\} \perp\!\!\!\perp X_2 \mid \{X_1, X_3\}.$$

How do we find these extra independence relations? The answer is “d-separation” which means “directed separation.” d-separation can be summarized by three rules. Consider the four DAG’s in Figure 17.6 and the DAG in Figure 17.7. The first 3 DAG’s in Figure 17.6 have no colliders. The DAG in the lower right of Figure 17.6 has a collider. The DAG in Figure 17.7 has a collider with a descendant.

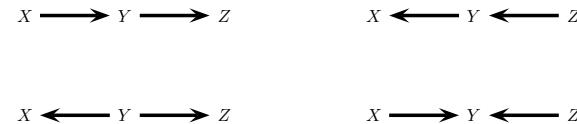


FIGURE 17.6. The first three DAG’s have no colliders. The fourth DAG in the lower right corner has a collider at Y.

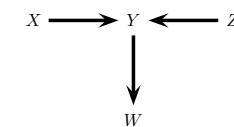


FIGURE 17.7. A collider with a descendant.

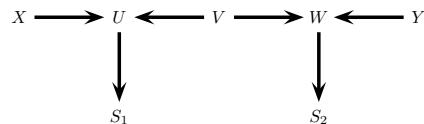


FIGURE 17.8. d-separation explained.

The Rules of d-Separation
Consider the DAGs in Figures 17.6 and 17.7.

1. When Y is not a collider, X and Z are **d-connected**, but they are **d-separated** given Y .
2. If X and Z collide at Y , then X and Z are **d-separated**, but they are **d-connected** given Y .
3. Conditioning on the descendant of a collider has the same effect as conditioning on the collider. Thus in Figure 17.7, X and Z are **d-separated** but they are **d-connected** given W .

Here is a more formal definition of d-separation. Let X and Y be distinct vertices and let W be a set of vertices not containing X or Y . Then X and Y are **d-separated given** W if there exists no undirected path U between X and Y such that (i) every collider on U has a descendant in W , and (ii) no other vertex on U is in W . If A , B , and W are distinct sets of vertices and A and B are not empty, then A and B are d-separated given W if for every $X \in A$ and $Y \in B$, X and Y are d-separated given W . Sets of vertices that are not d-separated are said to be **d-connected**.

17.9 Example. Consider the DAG in Figure 17.8. From the d-separation rules we conclude that:

- X and Y are d-separated (given the empty set);
- X and Y are d-connected given $\{S_1, S_2\}$;
- X and Y are d-separated given $\{S_1, S_2, V\}$.

17.10 Theorem. ² Let A , B , and C be disjoint sets of vertices. Then $A \perp\!\!\!\perp B \mid C$ if and only if A and B are d-separated by C .

²We implicitly assume that \mathbb{P} is **faithful** to \mathcal{G} which means that \mathbb{P} has no extra independence relations other than those logically implied by the Markov Condition.

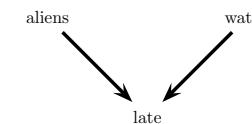


FIGURE 17.9. Jordan's alien example (Example 17.11). Was your friend kidnapped by aliens or did you forget to set your watch?

17.11 Example. The fact that conditioning on a collider creates dependence might not seem intuitive. Here is a whimsical example from Jordan (2004) that makes this idea more palatable. Your friend appears to be late for a meeting with you. There are two explanations: she was abducted by aliens or you forgot to set your watch ahead one hour for daylight savings time. (See Figure 17.9.) Aliens and Watch are blocked by a collider which implies they are marginally independent. This seems reasonable since — before we know anything about your friend being late — we would expect these variables to be independent. We would also expect that $\mathbb{P}(\text{Aliens} = \text{yes} \mid \text{Late} = \text{yes}) > \mathbb{P}(\text{Aliens} = \text{yes})$; learning that your friend is late certainly increases the probability that she was abducted. But when we learn that you forgot to set your watch properly, we would lower the chance that your friend was abducted. Hence, $\mathbb{P}(\text{Aliens} = \text{yes} \mid \text{Late} = \text{yes}) \neq \mathbb{P}(\text{Aliens} = \text{yes} \mid \text{Late} = \text{yes}, \text{Watch} = \text{no})$. Thus, Aliens and Watch are dependent given Late. ■

17.12 Example. Consider the DAG in Figure 17.2. In this example, overweight and smoking are marginally independent but they are dependent given heart disease. ■

Graphs that look different may actually imply the same independence relations. If \mathcal{G} is a DAG, we let $\mathcal{I}(\mathcal{G})$ denote all the independence statements implied by \mathcal{G} . Two DAGs \mathcal{G}_1 and \mathcal{G}_2 for the same variables V are **Markov equivalent** if $\mathcal{I}(\mathcal{G}_1) = \mathcal{I}(\mathcal{G}_2)$. Given a DAG \mathcal{G} , let $\text{skeleton}(\mathcal{G})$ denote the undirected graph obtained by replacing the arrows with undirected edges.

17.13 Theorem. Two DAGs \mathcal{G}_1 and \mathcal{G}_2 are Markov equivalent if and only if (i) $\text{skeleton}(\mathcal{G}_1) = \text{skeleton}(\mathcal{G}_2)$ and (ii) \mathcal{G}_1 and \mathcal{G}_2 have the same unshielded colliders.

17.14 Example. The first three DAGs in Figure 17.6 are Markov equivalent. The DAG in the lower right of the Figure is not Markov equivalent to the others. ■