

Causality

Talk 2: Directed acyclic graph (DAG) models

Yao Zhang

Intelligent Information Processing Research Group,
Faculty of Electrical Engineering and Computer Science,
Ningbo University

Note: The following slides are primarily adapted from the course materials¹.

Nov 26, 2025

¹C. Heinze-Deml. Causality. URL: <https://stat.ethz.ch/lectures/ss21/causality.php>. 



Directed acyclic graph (DAG) models

Causality

Christina Heinze-Deml

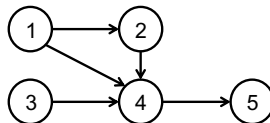
Spring 2021

Today

- Graph terminology
- Directed acyclic graph (DAG) models
- Markov properties
- d-separation

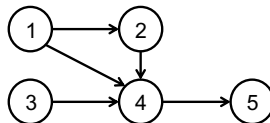
Graph terminology

- A graph $G = (V, E)$ consists of vertices (nodes) V and edges E
- There is **at most one edge** between every ordered pair of vertices
- Two vertices are **adjacent** if there is an edge between them
- If all edges are directed ($i \rightarrow j$), the graph is called **directed**
- A **path** between i and j is a sequence of **distinct** vertices (i, \dots, j) such that **successive vertices are adjacent**
- A **directed path** from i to j is a path between i and j where all edges are pointing towards j , i.e., $i \rightarrow \dots \rightarrow j$



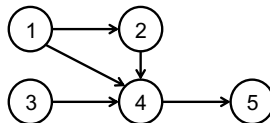
Graph terminology

- A **cycle** is a path (i, j, \dots, k) plus an edge between k and i
- A **directed cycle** is a directed path (i, j, \dots, k) from i to k , plus an edge $k \rightarrow i$
- A **directed acyclic graph (DAG)** is a directed graph without directed cycles



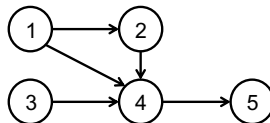
Graph terminology

- If $i \rightarrow j$, then i is a **parent** of j , and j is a **child** of i
- If there is a directed path from i to j , then i is an **ancestor** of j and j is a **descendant** of i
- Each vertex is also an ancestor and descendant of itself
- The sets of parents, children, descendants and ancestors of i in G are denoted by **$pa(i, G)$** , **$ch(i, G)$** , **$desc(i, G)$** , **$an(i, G)$**
- We omit G if the graph is clear from the context



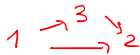
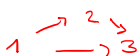
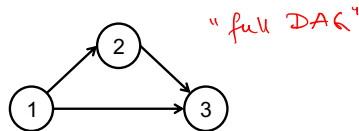
Graph terminology

- We write sets of vertices in bold face
- The previous definitions are applied disjunctively to sets
 - Example: $\text{pa}(\mathcal{S}) = \bigcup_{k \in \mathcal{S}} \text{pa}(k)$
- The non-descendants of \mathcal{S} are the complement of $\text{desc}(\mathcal{S})$:
 $\text{nondesc}(\mathcal{S}) := V \setminus \text{desc}(\mathcal{S})$



Graph terminology

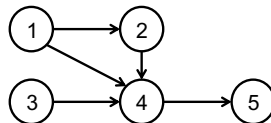
- We call G **fully connected** if all pairs of nodes are adjacent
- How many possibilities for a fully connected DAG?



$p!$ possibilities ; $p = |V|$

DAGs and random variables

- Each vertex represents a random variable:
vertex i represents random variable X_i
- If $A \subseteq V$, then $X_A := \{X_i : i \in A\}$



- Edges denote relationships between pairs of variables
(we will make this more precise)

Factorization of the joint density

- We can connect a distribution with density f to a DAG in the following way:
- We always have:

$$f(x_1, \dots, x_p) = \underbrace{f(x_1)} \underbrace{f(x_2|x_1)} \dots f(x_p|x_1, \dots, x_{p-1}) \quad \text{"chain rule"}$$

- A set of variables $X_{\text{pa}(j)}$ is said to be **Markovian parents** of X_j if it is a minimal subset of $\{\underline{X_1}, \dots, \underline{X_{j-1}}\}$ such that $f(\underline{x_j|x_1, \dots, x_{j-1}}) = f(\underline{x_j|x_{\text{pa}(j)}})$
 - **Note:** Markovian parents depend on the chosen ordering of the variables

Factorization of the joint density

- We always have:

$$f(x_1, \dots, x_p) = f(x_1)f(x_2|x_1) \dots f(x_p|x_1, \dots, x_{p-1})$$

- A set of variables $X_{\text{pa}(j)}$ is said to be **Markovian parents** of X_j if it is a minimal subset of $\{X_1, \dots, X_{j-1}\}$ such that $f(x_j|x_1, \dots, x_{j-1}) = f(x_j|x_{\text{pa}(j)})$
- Then

$$f(x_1, \dots, x_p) = \prod_{j=1}^p f(x_j | \underline{x_{\text{pa}(j)}})$$

"factorization property"

- We can draw a DAG accordingly
- The distribution is said to **factorize according to this DAG**

Notes week 2 - I

Consider (X_1, X_2, X_3) and suppose that $X_1 \perp\!\!\!\perp X_3 \mid X_2$

is the only (conditional) independence:

$$f(x_1 \mid x_2, x_3) = f(x_1 \mid x_2)$$

$$f(x_3 \mid x_1, x_2) = f(x_3 \mid x_2)$$

Then $f(x_1, x_2, x_3) = f(x_1) f(x_2 \mid x_1) f(x_3 \mid x_1, x_2)$ ↘ simplify
 $= f(x_1) f(x_2 \mid x_1) f(x_3 \mid x_2)$

DAG: $1 \rightarrow 2 \rightarrow 3$

Or $f(x_3, x_2, x_1) = f(x_3) f(x_2 \mid x_3) f(x_1 \mid x_2, x_3)$ ↘ simplify
 $= f(x_3) f(x_2 \mid x_3) f(x_1 \mid x_2)$

DAG: $3 \rightarrow 2 \rightarrow 1$

Or $f(x_1, x_3, x_2) = f(x_1) f(x_3 \mid x_1) f(x_2 \mid x_1, x_3)$ cannot simplify further
DAG: $1 \rightarrow 3 \rightarrow 2$

Note: Markovian parents depend on the chosen ordering of the variables

Factorization of the joint density

- A distribution can factorize according to several DAGs
- Every distribution factorizes according to a full DAG
 - **Note:** there are $p!$ possibilities
- Sometimes a distribution factorizes according to a sparse DAG
 - I.e., a DAG with few edges
 - E.g. first-order Markov chain:
 - $f(x_1, \dots, x_p) = f(x_1)f(x_2|x_1) \dots f(x_p|x_1, \dots, x_{p-1}) = f(x_1)f(x_2|x_1) \dots f(x_p|x_{p-1})$
 - DAG: $1 \rightarrow 2 \rightarrow \dots \rightarrow p$

DAG models

- A **DAG model** or **Bayesian network** is a combination (G, P) , where G is a DAG and P is a distribution that factorizes according to G
- DAG models can be used for various purposes:
 - Estimating the joint density from low order conditional densities
 - Reading off conditional independencies from the DAG
 - Probabilistic reasoning (expert systems)
 - Causal inference

DAG models

- A **DAG model** or **Bayesian network** is a combination (G, P) , where G is a DAG and P is a distribution that factorizes according to G
- DAG models can be used for various purposes:
 - Estimating the joint density from low order conditional densities
 - **Reading off conditional independencies from the DAG**
 - Probabilistic reasoning (expert systems)
 - Causal inference

Reading off conditional independencies: Markov property

- First-order Markov models: the future is independent of the past given the present

$$1 \rightarrow 2 \rightarrow \dots \rightarrow (t-1) \rightarrow t \rightarrow (t+1)$$

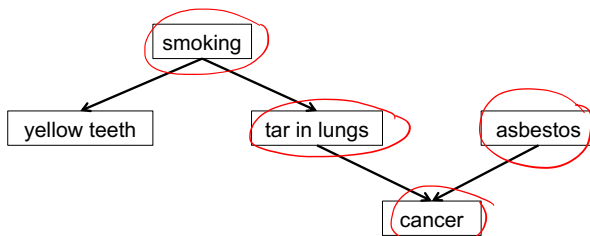
$$X_{t+1} \perp\!\!\!\perp \{X_{t-1}, X_{t-2}, \dots, X_1\} \mid X_t$$

- In DAG models, we have a similar (local) Markov property
- Let S be any collection of nodes. Then:

$$X_S \perp\!\!\!\perp X_{\text{nondesc}(S) \setminus \text{pa}(S)} \mid X_{\text{pa}(S)}$$

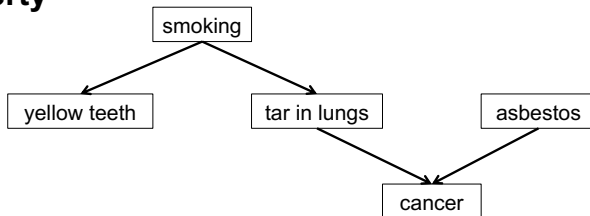
$$X_S \perp\!\!\!\perp X_{\text{nondesc}(S) \setminus \text{pa}(S)} \mid X_{\text{pa}(S)}$$

Example



- Take $S = \{\text{yellow teeth}\}$ and apply the local Markov property
- Then:
 - $\text{pa}(\text{yellow teeth}) = \{\text{smoking}\}$
 - $\text{nondesc}(\text{yellow teeth}) = \{\text{smoking, tar, cancer, asbestos}\}$
- Hence, **yellow teeth** $\perp\!\!\!\perp \{\text{tar, cancer, asbestos}\} \mid \text{smoking}$ in any distribution that factorizes according to this DAG

Markov property



- Is $\text{tar} \perp\!\!\!\perp \text{asbestos} \mid \text{cancer}$?
- The local Markov property cannot be used to read off arbitrary conditional (in)dependencies
 - For this we have **d-separation**

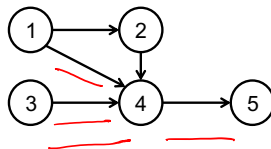
Graph terminology

- Need new terminology:
 - A non-endpoint node i is a **collider on a path** if the path contains $\rightarrow i \leftarrow$ (arrows collide at i)
 - Otherwise, it is a **non-collider on the path**

- Is 4 a collider in the given graph? \Leftarrow bad question

- 4 is a collider on the path (3, 4, 1)
- 4 is a non-collider (3, 4, 5)

\Rightarrow collider status is always relative to a path!



d-separation

- A **path** between i to j is **blocked** by a set S (not containing i or j) if at least one of the following holds:
 - There is a non-collider on the path that is in S ; or
 - There is a collider on the path such that neither this collider nor any descendants are in S
- A path that is not blocked is **active**
- If all paths between $i \in A$ and $j \in B$ are blocked by S , then **A and B are d-separated by S** . Otherwise they are **d-connected** given S .
- Denote d-separation by \perp

Global Markov property

- **Definition:**

A distribution P with density p satisfies the **global Markov property** with respect to a DAG G if:

$$A \text{ and } B \text{ are d-separated by } S \text{ in } G \Rightarrow X_A \perp\!\!\!\perp X_B \mid X_S \text{ in } P$$

- **Theorem** (Pearl, 1988):

A distribution P with density p satisfies the global Markov property with respect to G if and only if p factorizes according to G .

Notes week 2 - II

G :



A distribution P is said to satisfy

* the **global Markov property** wrt G if

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$

and

$$X_1 \perp\!\!\!\perp X_4 \mid X_2, X_3$$

* the **local Markov property** wrt G if

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$

$$X_4 \perp\!\!\!\perp X_1 \mid X_2, X_3$$

* the **Markov factorization property** wrt G if

$$p(x_1, x_2, x_3, x_4) = p(x_3) p(x_1 \mid x_3) p(x_2 \mid x_1) p(x_4 \mid x_2, x_3)$$

Faithfulness

- Given a DAG $G = (V, E)$, a distribution P on X_V is said to be **faithful** with respect to G if for all pairwise disjoint subsets A, B and S of V :

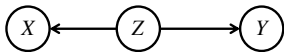
$$X_A \perp\!\!\!\perp X_B | X_S \text{ in } P \Rightarrow A \text{ and } B \text{ are d-separated by } S \text{ in } G$$

Example



$$X \perp Y|Z$$

e.g. fire \rightarrow smoke \rightarrow alarm



$$X \perp Y|Z$$

e.g. shoe size \leftarrow age of child \rightarrow reading skills



$$X \not\perp Y|Z$$

e.g. talent \rightarrow celebrity \leftarrow beauty

DAG models

- A **DAG model** or **Bayesian network** is a combination (G, P) , where G is a DAG and P is a distribution that factorizes according to G
- DAG models can be used for various purposes:
 - Estimating the joint density from low order conditional densities
 - Reading off conditional independencies from the DAG
 - **Probabilistic reasoning (expert systems)**
 - Causal inference

Probabilistic reasoning

- Conditional probabilities are rather counterintuitive for many people
- DAGs allow us to obtain conditional probabilities efficiently, using a “message passing” algorithm
 - See R script `02_graphical_models.R`
 - We won't discuss the details behind these algorithms

Discussion

Any comments or questions?

We may not always find an answer, and since we're not very familiar with causality, we will need to dedicate more time to this topic.