

# **Machine Learning 2**

**Graphical Models  
- Bayesian Networks**

Patrick Forré

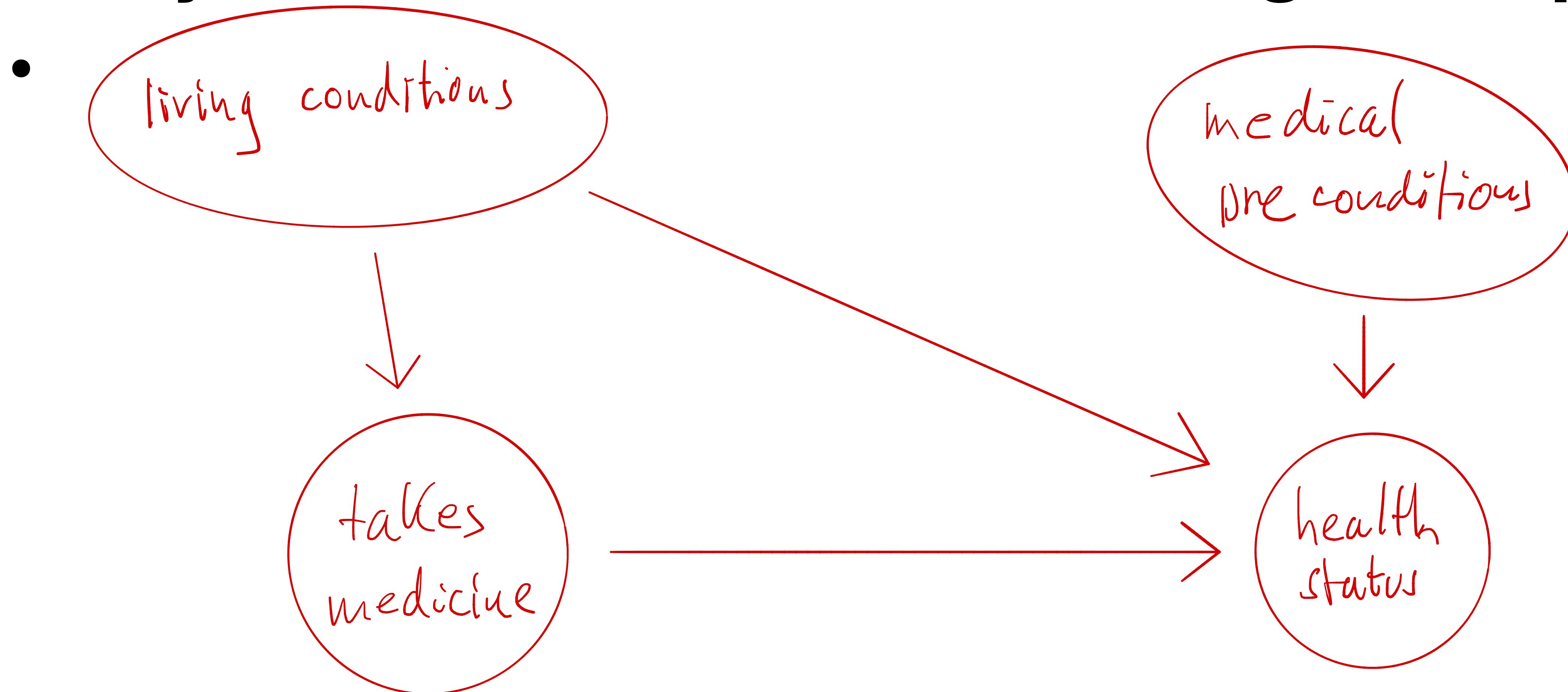
# **Machine Learning 2**

## **Graphical Models**

- Bayesian Networks**
  - Definitions**

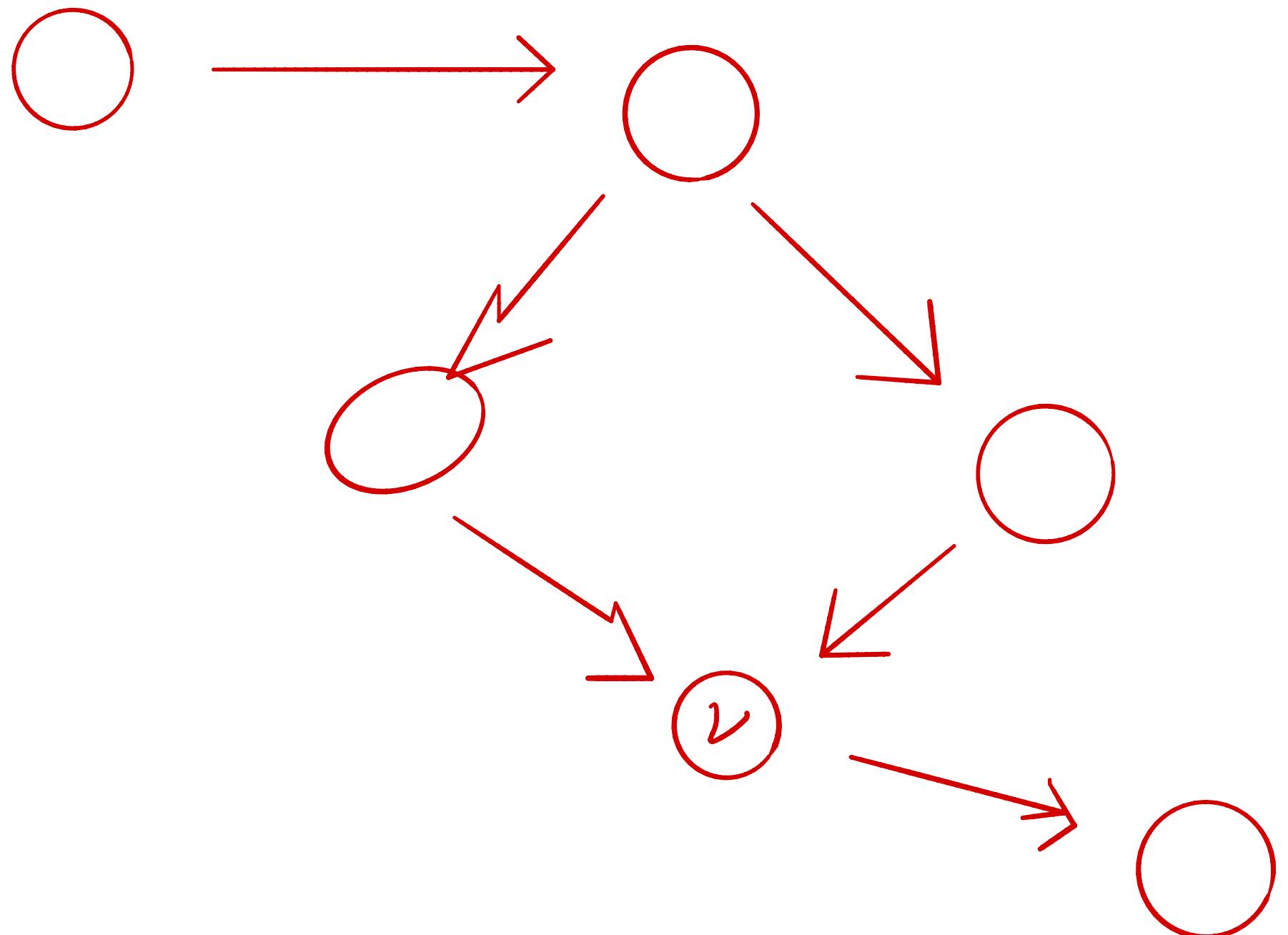
Patrick Forré

# Bayesian Networks - Motivating Example



# Directed Acyclic Graphs (DAGs) - Definitions

- A **Directed Acyclic Graph** (DAG) is a directed graph  $G = (V, E)$  with:
  - set of **nodes/vertices**  $V$ ,
  - set of **directed edges**  $E$
  - such that there is **no directed cycle**
- **Parents** of  $v$ : nodes pointing to  $v$ .
- **Children** of  $v$ : nodes that  $v$  points to.
- **Ancestors** of  $v$ : nodes with directed path ending in  $v$ .  
*(including  $v$  itself)*
- **Descendents** of  $v$ : nodes that have a directed path from  $v$  pointing to it.



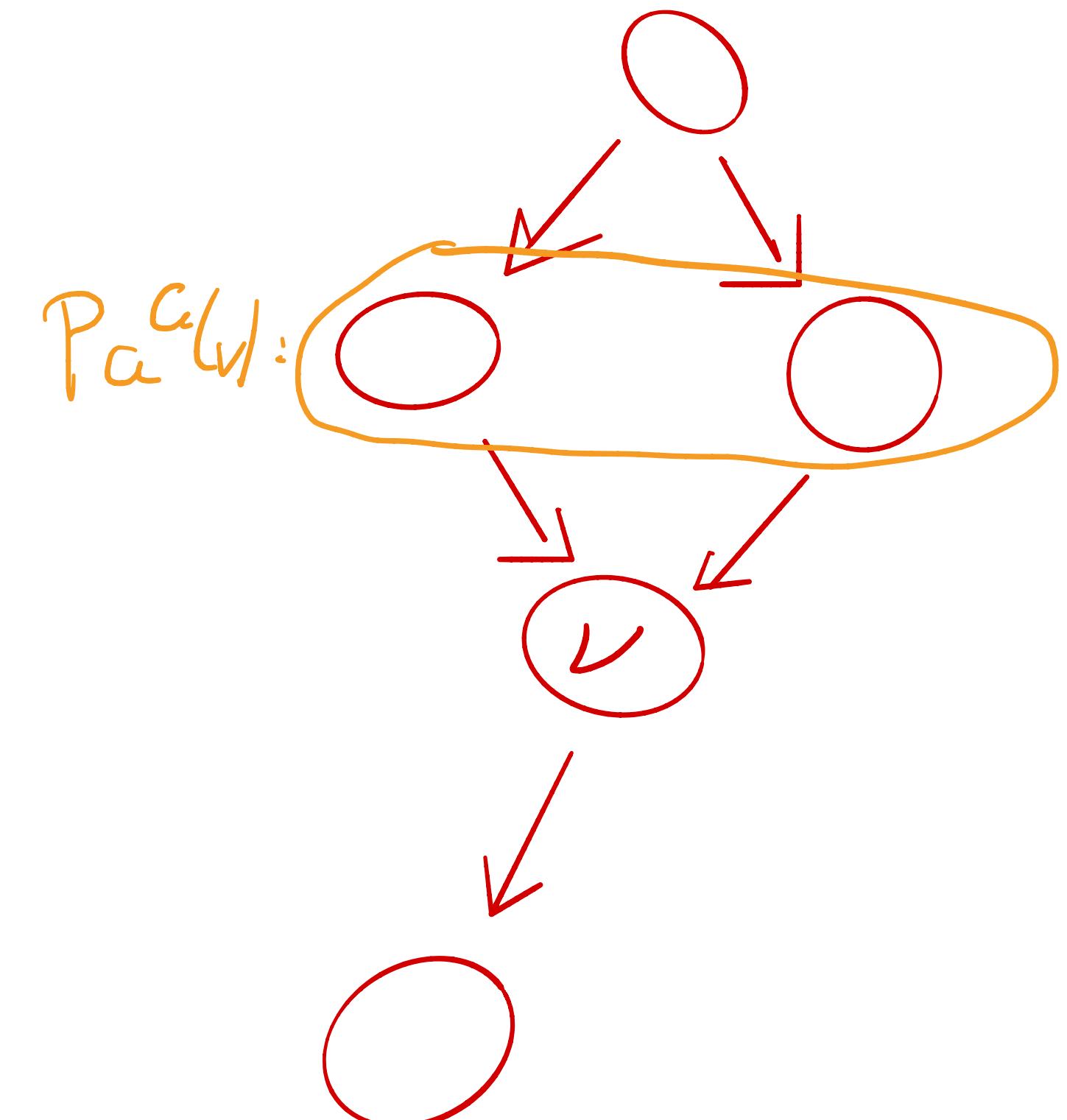
# Factorization Property

- Let  $V$  be an (index) set of random variables  $X_1, \dots, X_M$  with a joint distribution  $p(x_1, \dots, x_M)$ .
- Let  $G = (V, E)$  be a DAG.
- We say that  $p(x_1, \dots, x_M)$  factorizes over  $G$  if:

$$p(x_1, \dots, x_M) = \prod_{v \in V} p(x_v | x_{\text{Pa}^G(v)})$$

$$x_A := (x_v)_{v \in A}$$

$$A \subseteq V$$

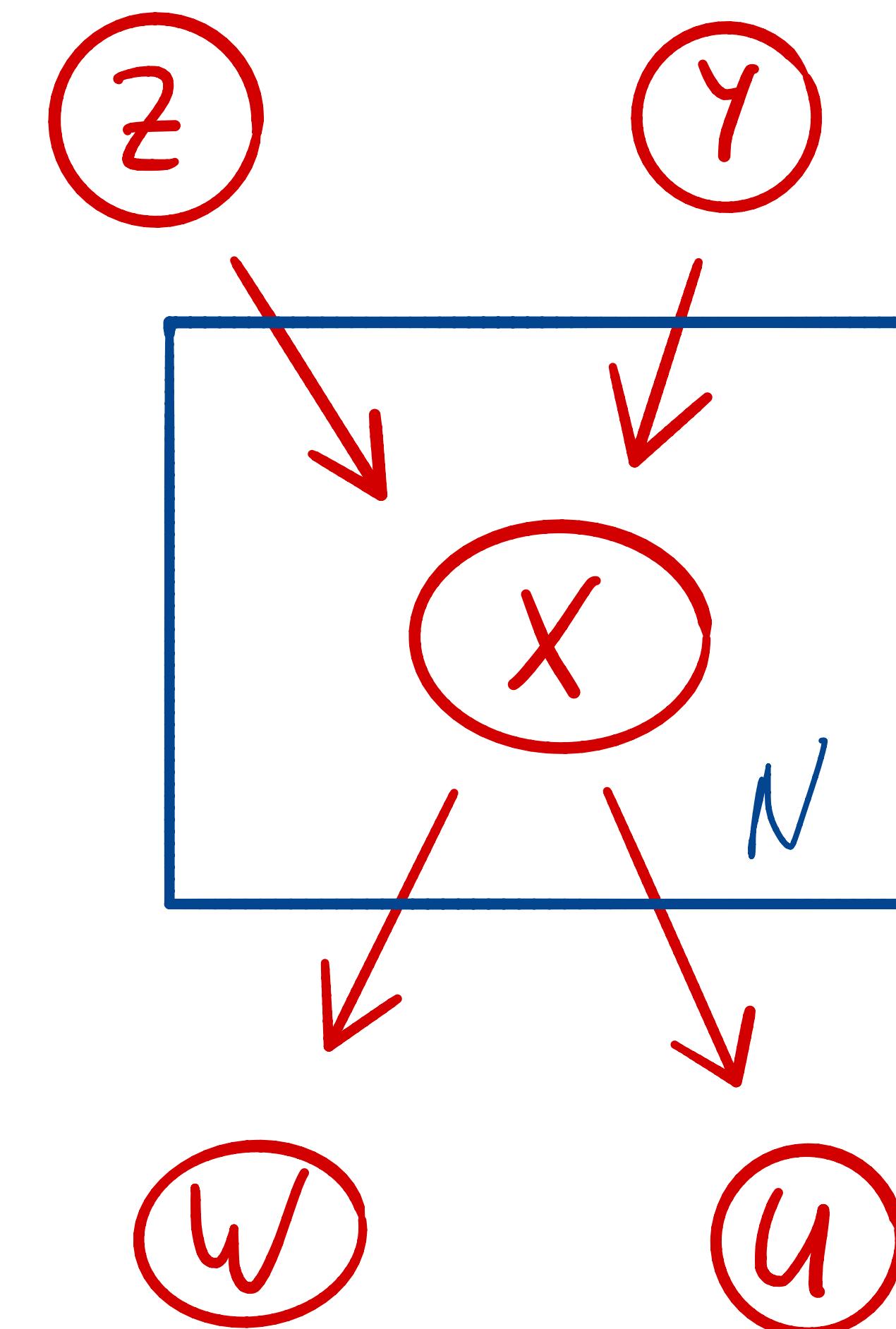
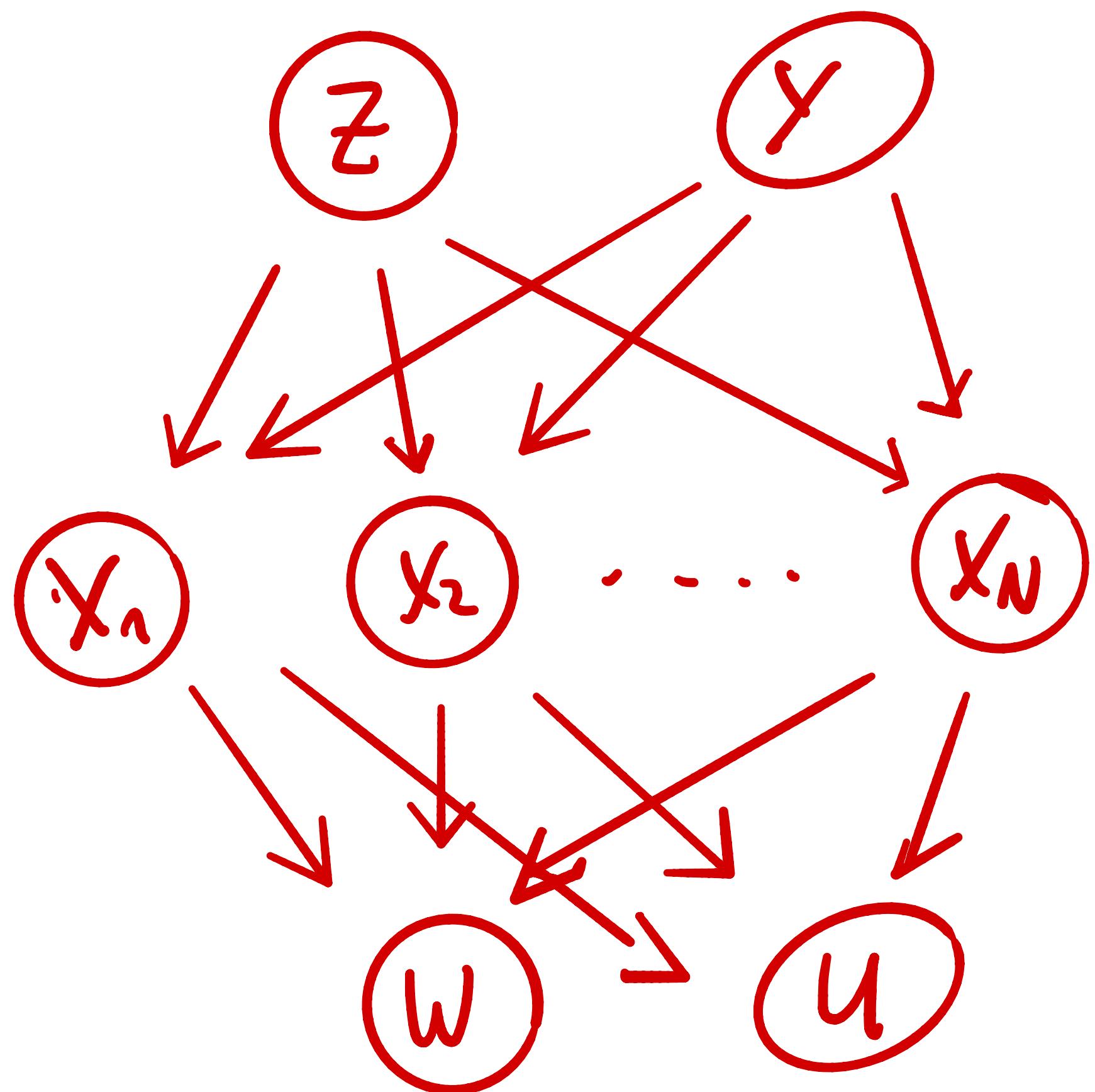


# Bayesian Network (BN) - Definition

- A Bayesian Network (BN)  $(G, p)$  by definition consists of:
  - an (index) set  $V$  of random variables  $X_1, \dots, X_M$  with a joint distribution  $p(x_1, \dots, x_M)$ ,
  - a DAG structure  $G = (V, E)$ , such that:
    - $p(x_1, \dots, x_M)$  factorizes over  $G$ .

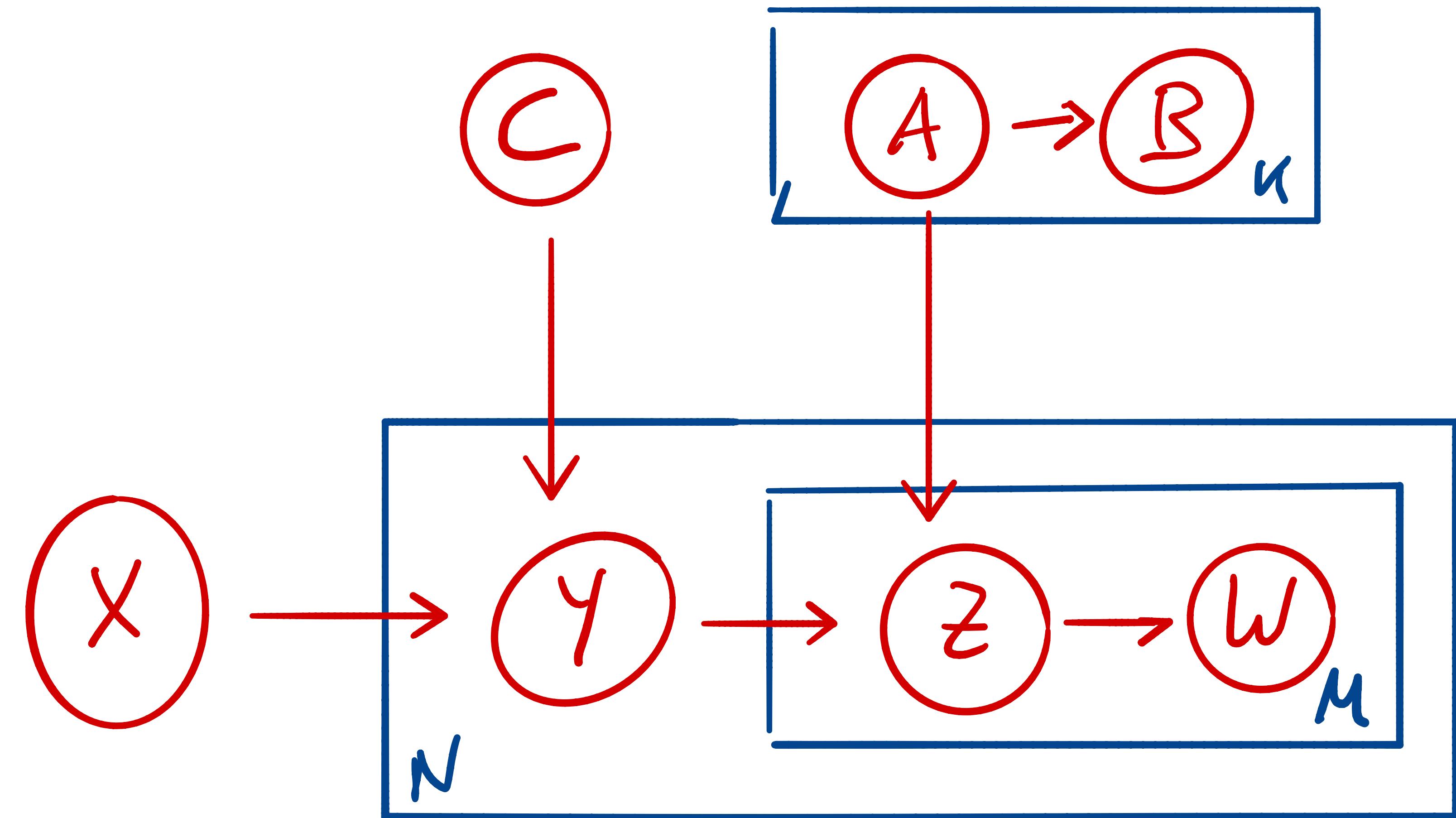
# Bayesian Networks - Plate Notation

.



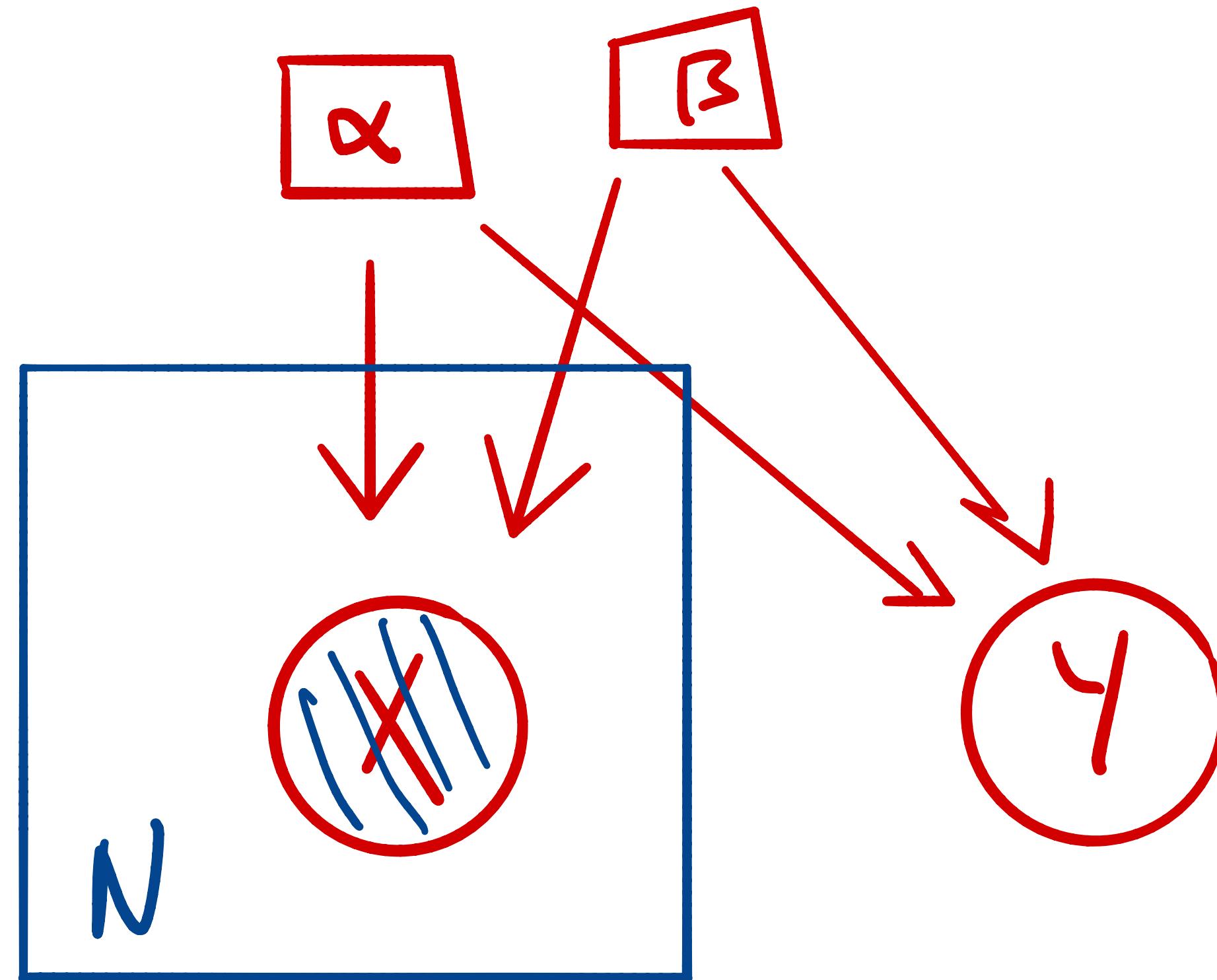
# Bayesian Networks - Nested Plate Notation

- 



# Bayesian Networks - Parameters & Observations

- 



# **Machine Learning 2**

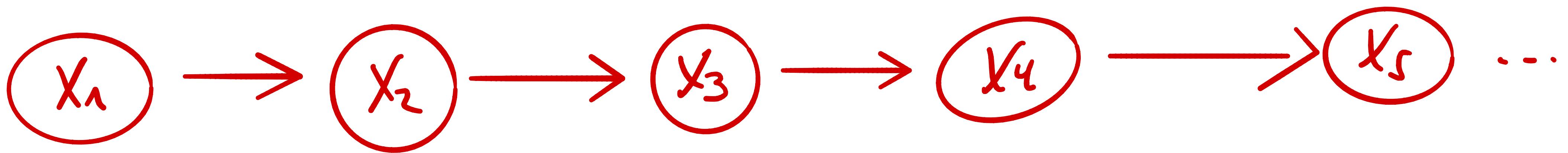
**Graphical Models**

- Bayesian Networks**
  - Examples**

Patrick Forré

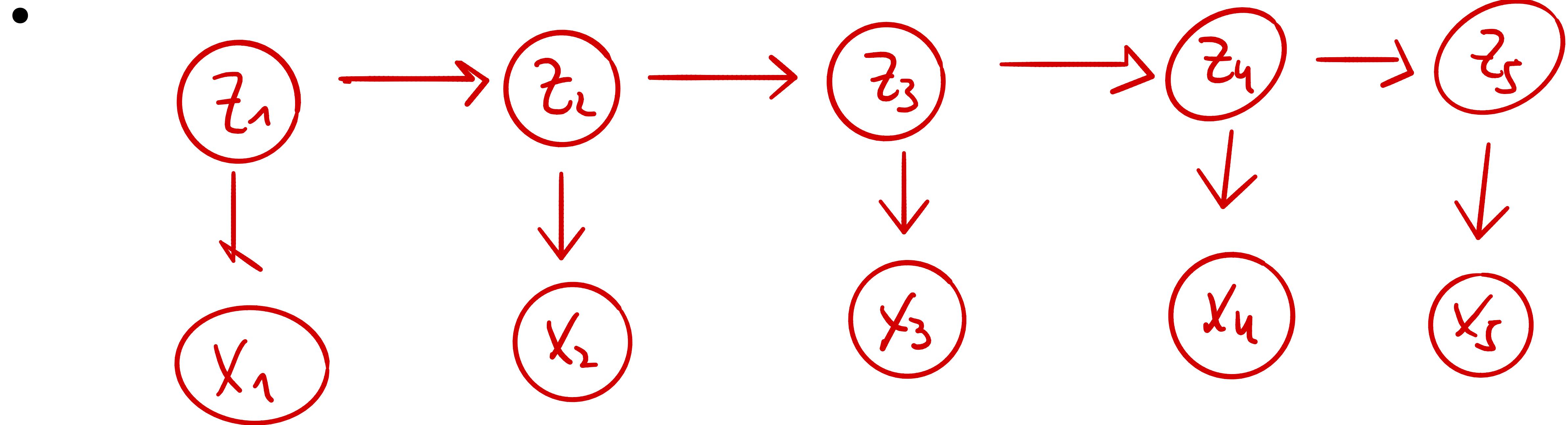
# Example: Markov Chain

.



$$\begin{aligned} p(x_v) &= \prod_{v \in V} p(x_v | x_{\text{Pa}(v)}) \\ &\stackrel{\text{Capital } V}{=} p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_2) \cdot p(x_4 | x_3) \cdot p(x_5 | x_4) \cdots \end{aligned}$$

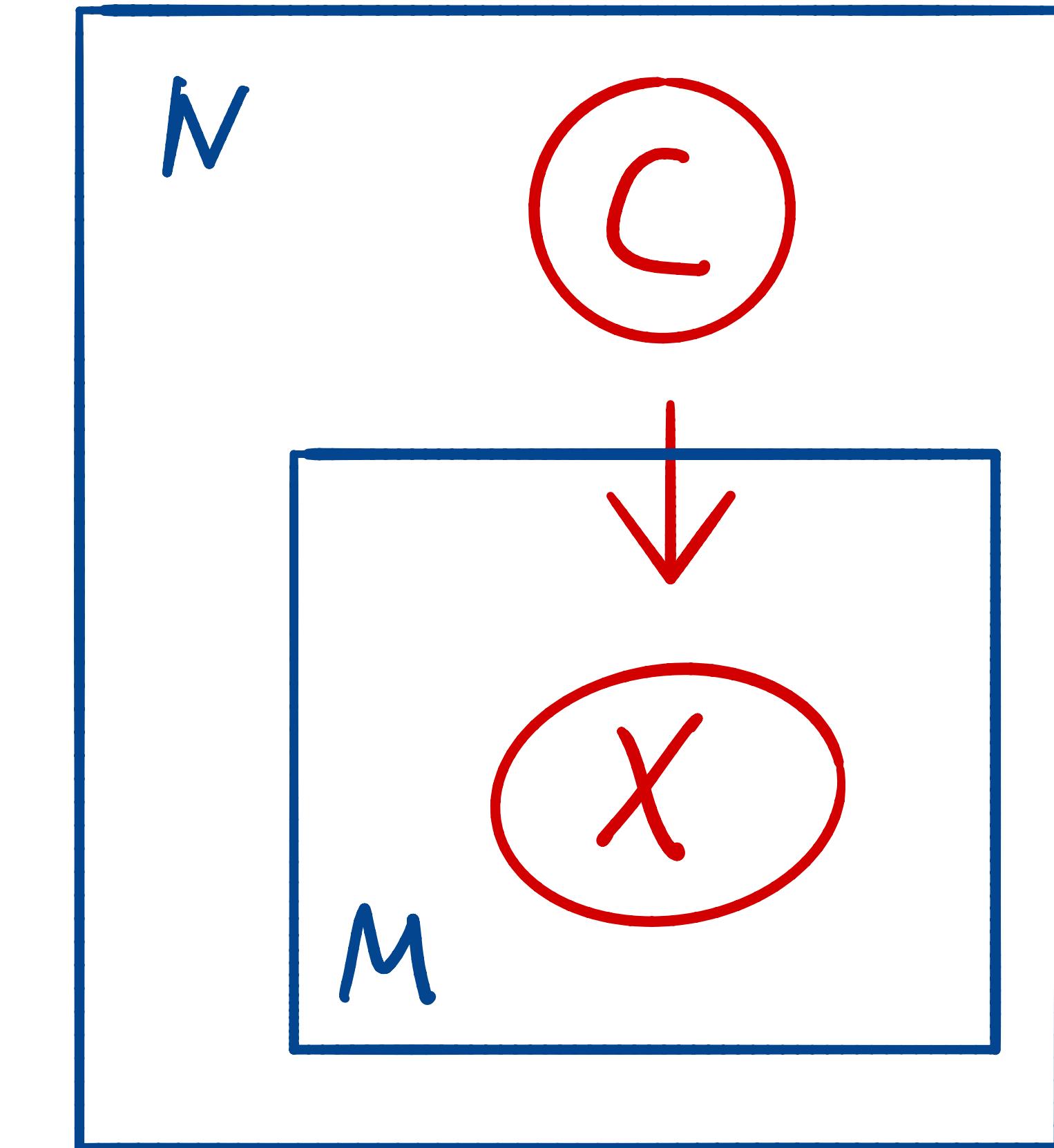
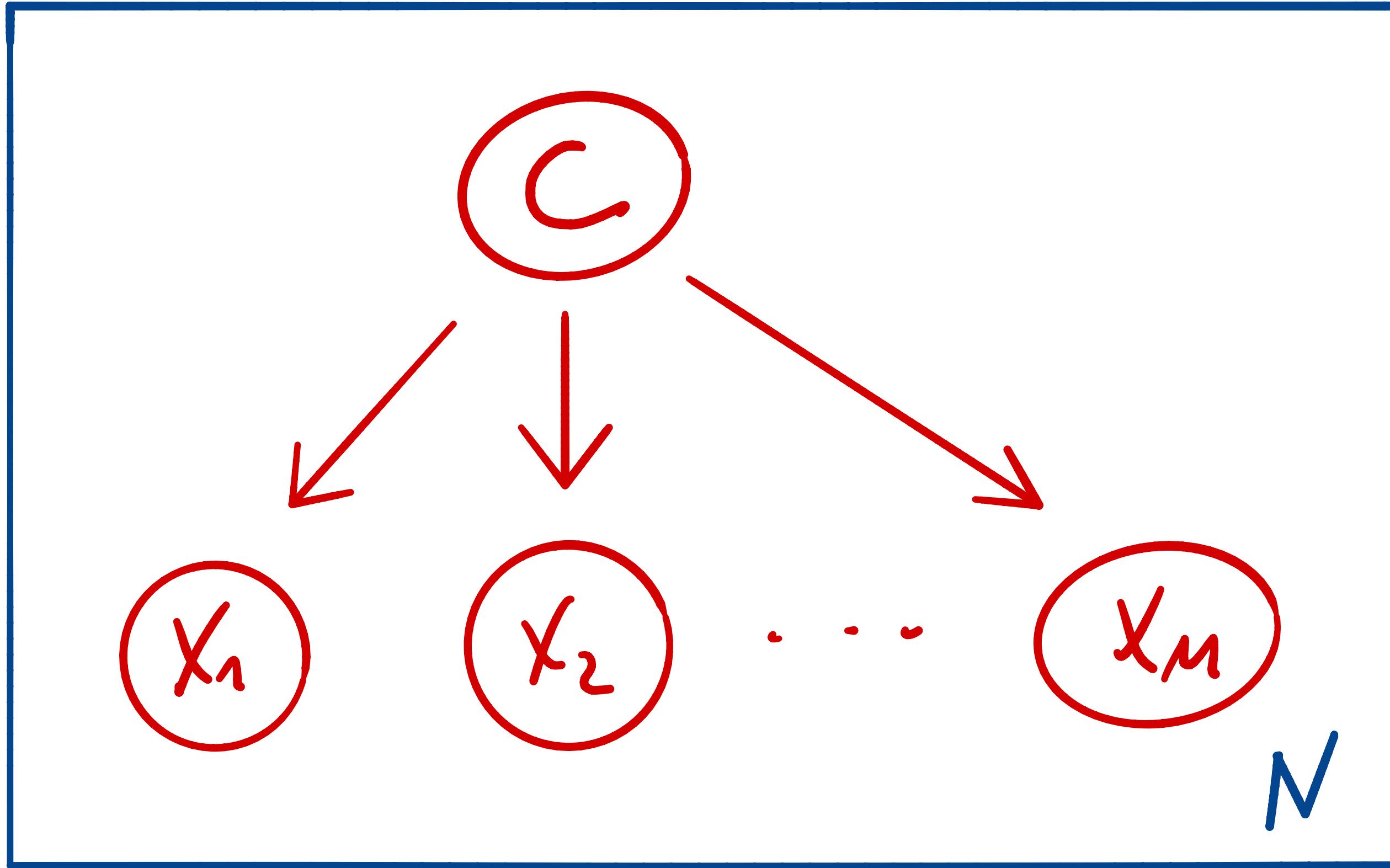
# Example: State-Space Model



$$p(x_t) = p(z_1) p(x_1|z_1) \cdot p(z_2|z_1) p(x_2|z_2) p(z_3|z_2) \cdots p(x_3|z_3) \cdots$$

*capital V*

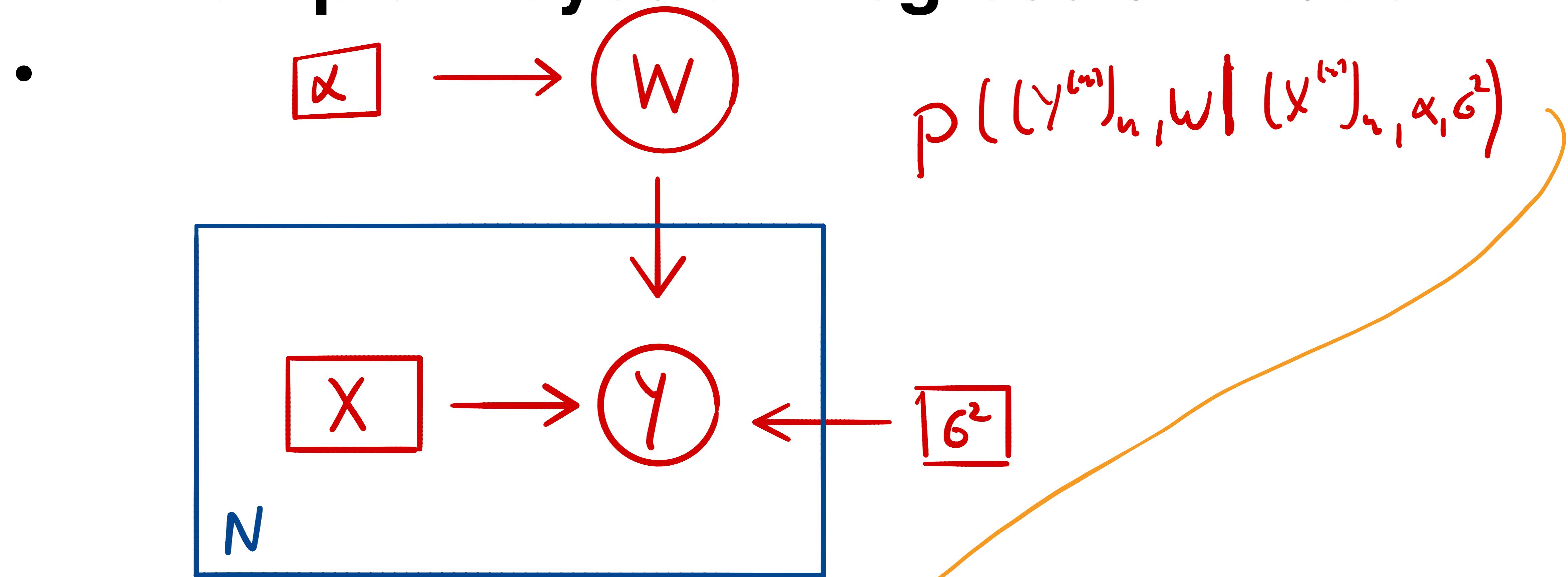
# Example: Naive Bayes Model



$$p(x_v) = \prod_{n=1}^N [p(C^{(n)}) \cdot p(x_1^{(n)} | C^{(n)}) \dots p(x_M^{(n)} | C^{(n)})]$$

C  
Capital V

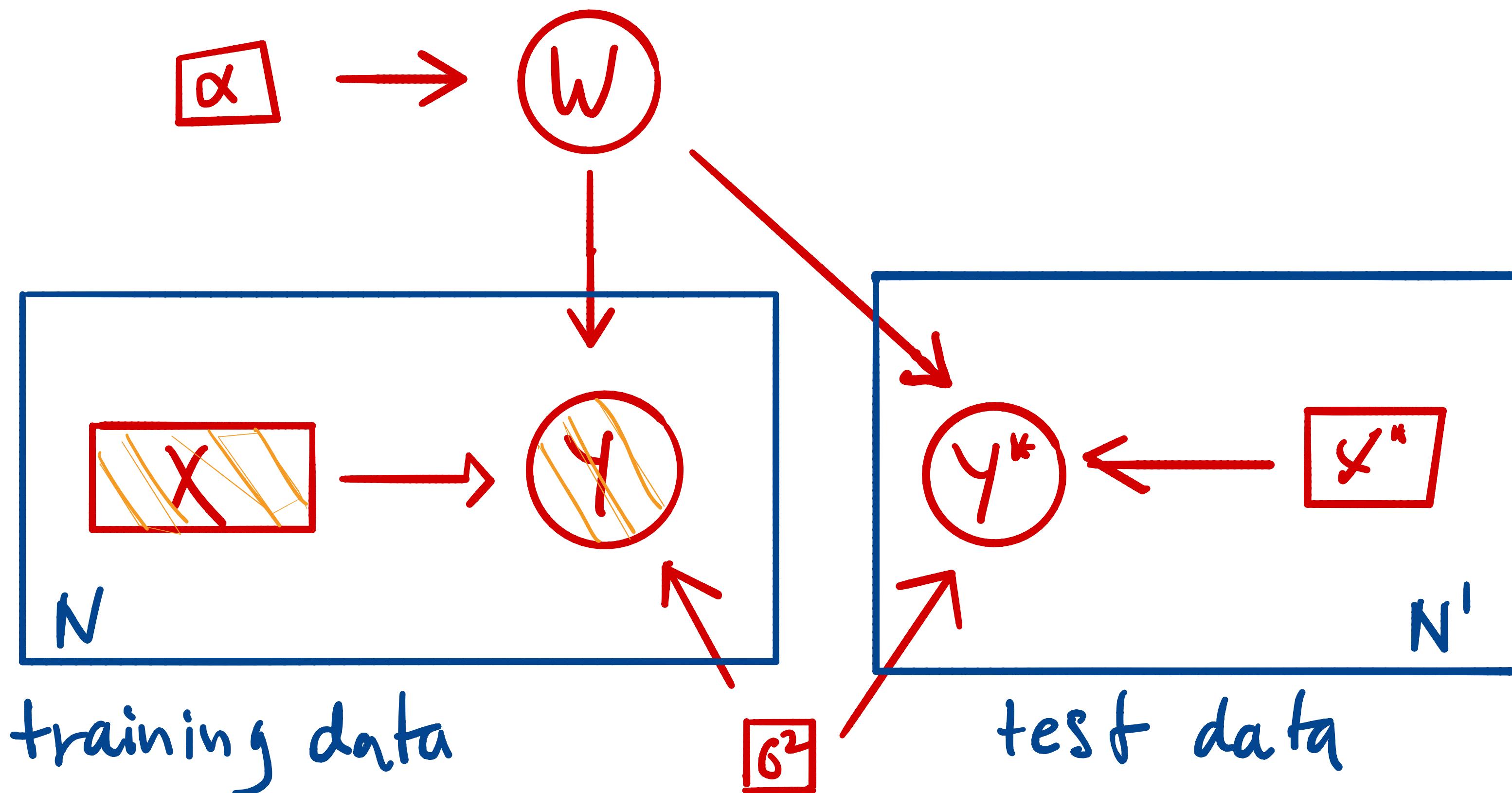
# Example: Bayesian Regression Model



$$= p(w|\alpha) \cdot \prod_{n=1}^N p(y^{(n)}|x^{(n)}, w, \sigma^2)$$

# Example: Prediction in Bayesian Regression Model

- 



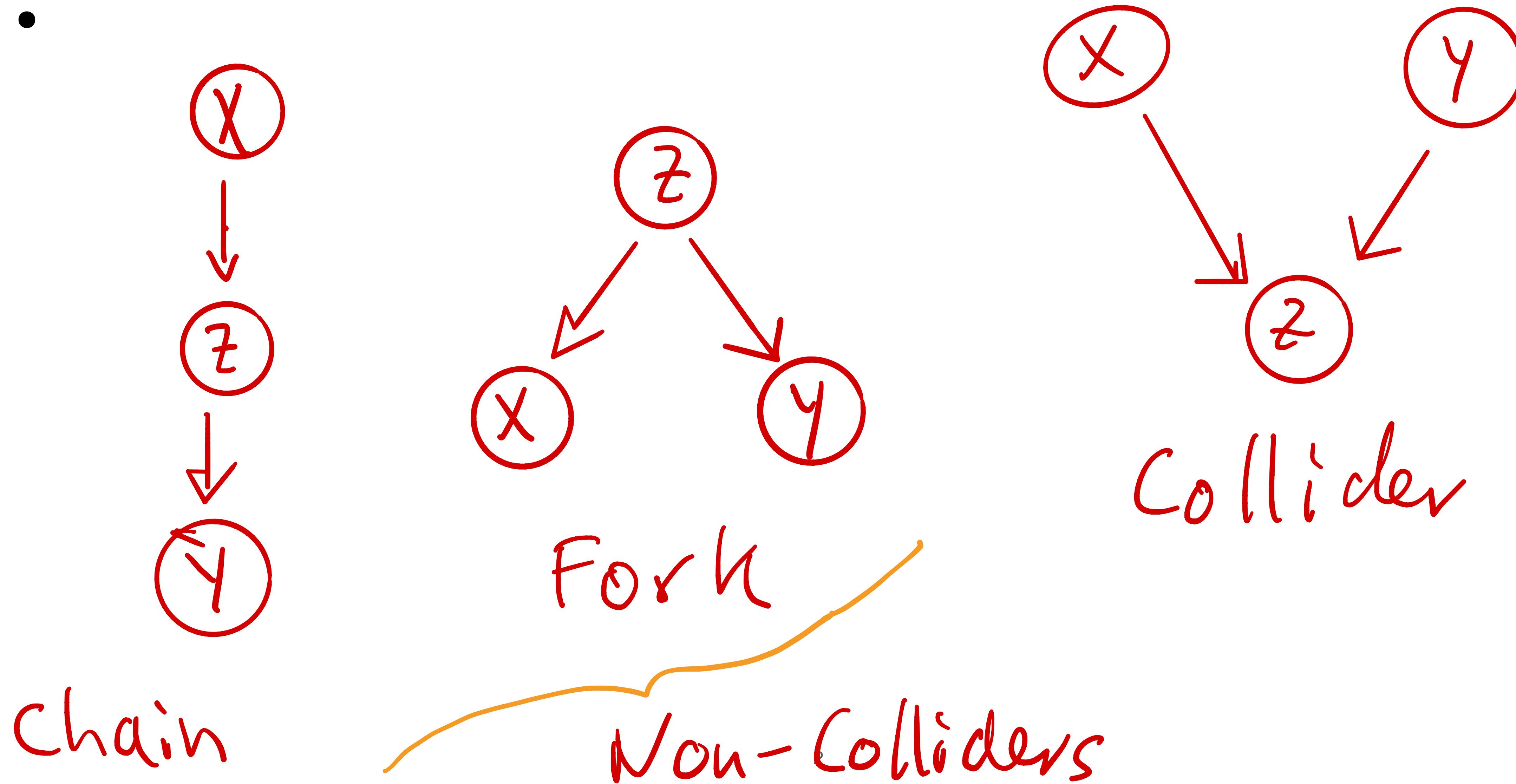
# **Machine Learning 2**

## **Graphical Models**

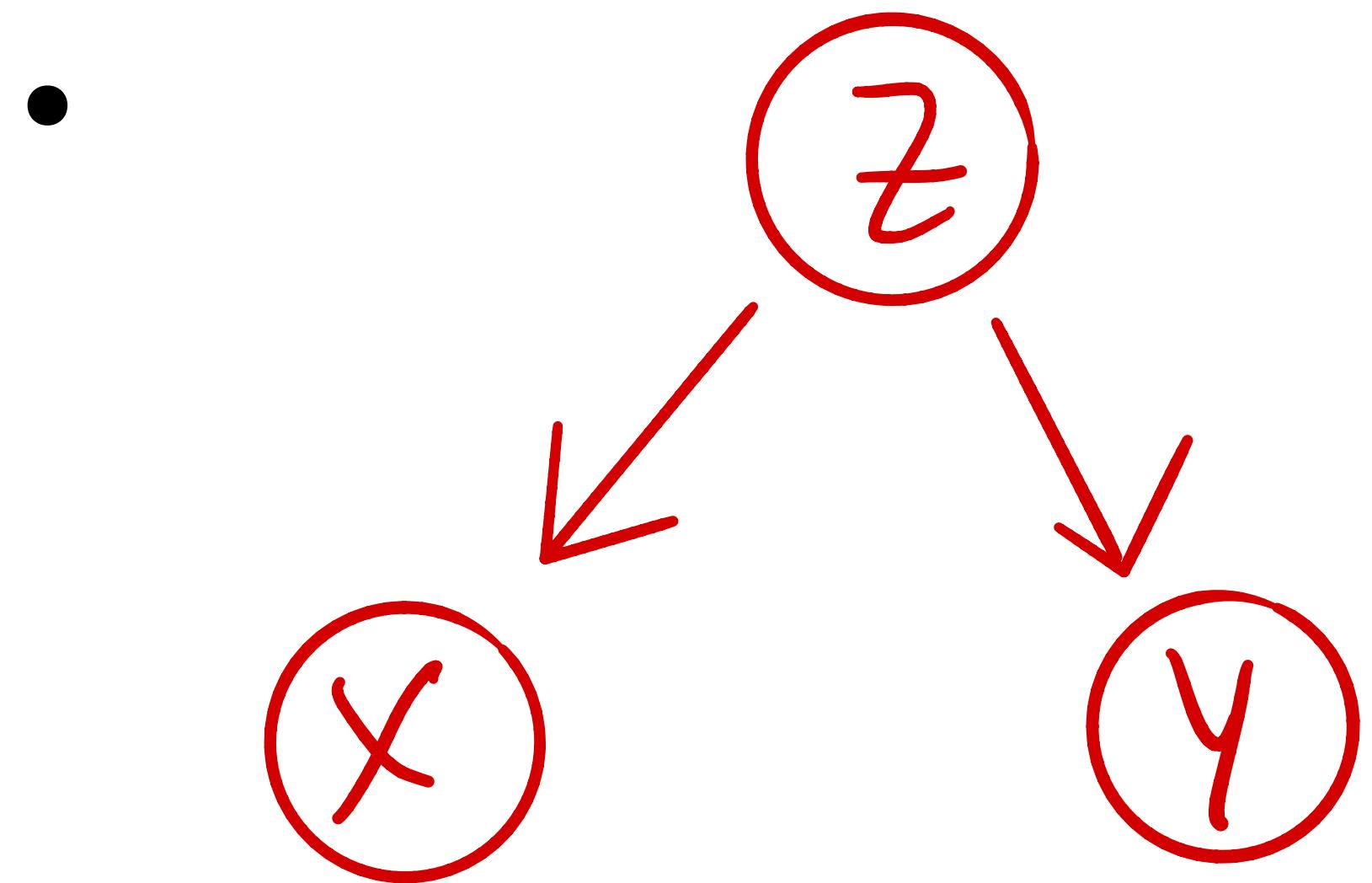
- Bayesian Networks**
- the 3 building blocks**

Patrick Forré

# The Good, the Bad and the Ugly



# The Fork



$$(*) \quad p(x,y,z) = p(x|z) \cdot p(y|z) \cdot p(z)$$

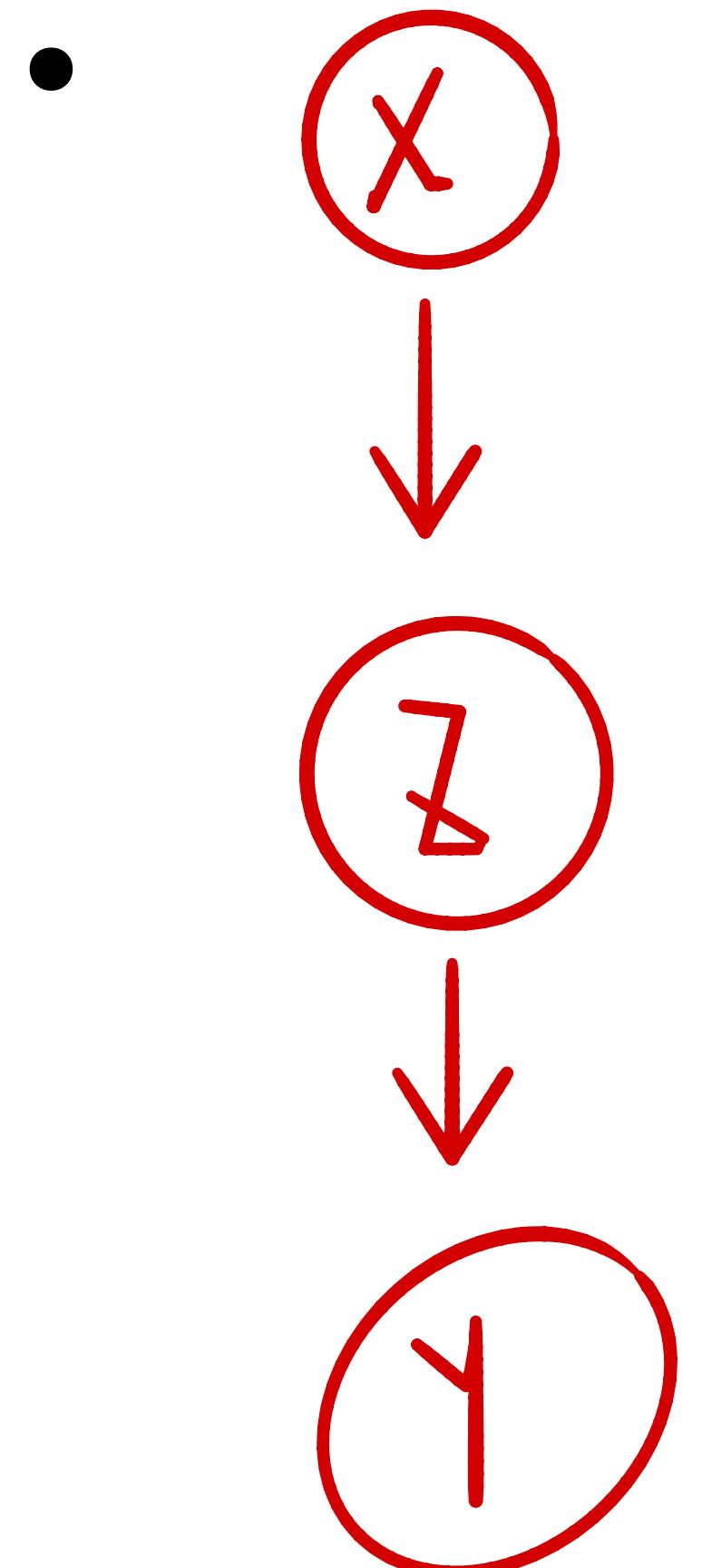
generally :  $x \perp\!\!\!\perp y$

always :  $x \perp\!\!\!\perp y \mid z$

$$(*) \Rightarrow p(x,y|z) = p(x|z) \cdot p(y|z)$$

$$\Rightarrow x \perp\!\!\!\perp y \mid z$$

# The Chain



$$p(x,y,z) = p(y|z) \cdot p(z|x) \cdot p(x)$$

generally:  $x \not\perp\!\!\!\perp y$

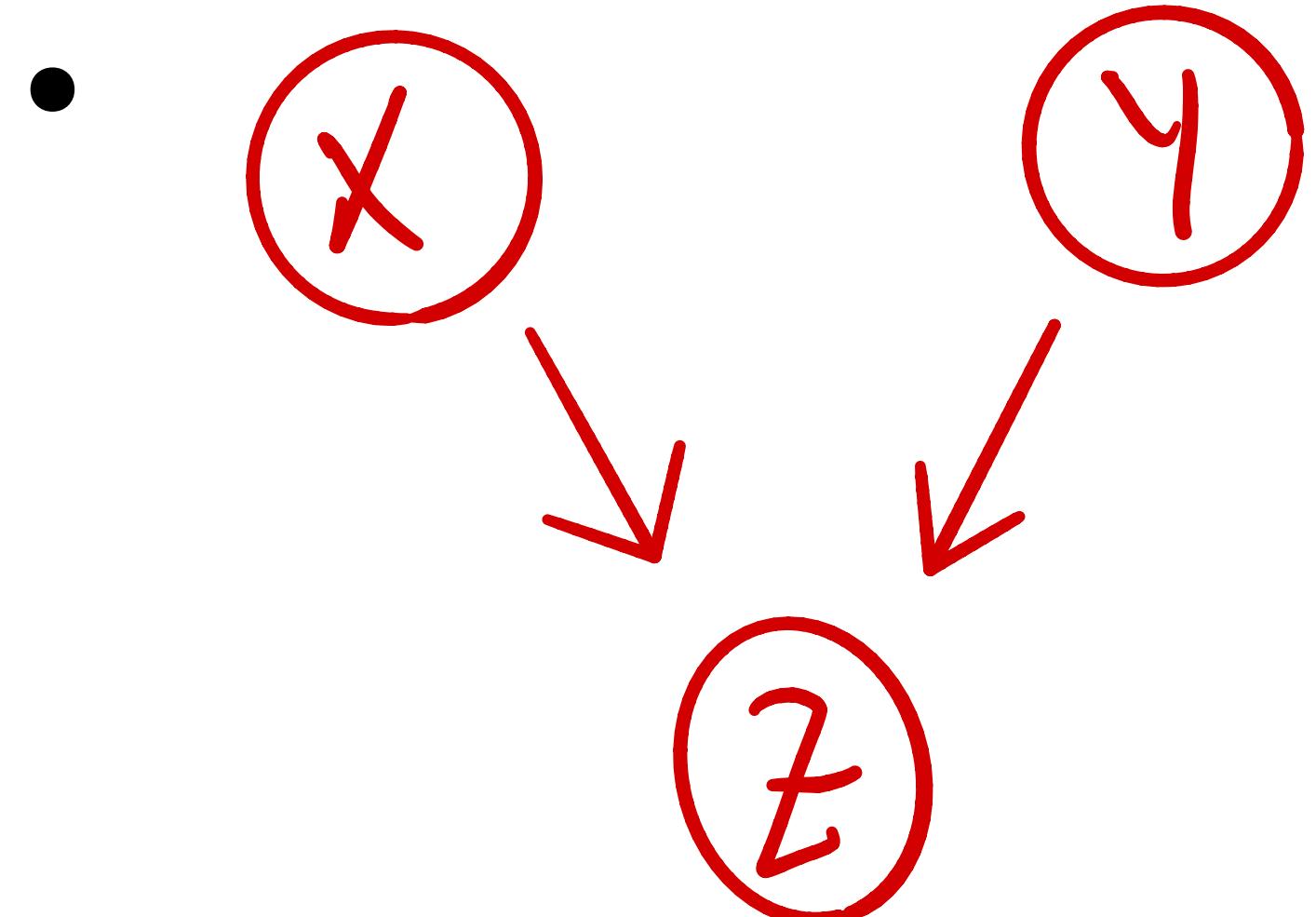
always:  $x \perp\!\!\!\perp y | z$

$$\begin{aligned} (*) \Rightarrow p(x,y,z) &= p(y|z) \cdot p(z|x) \\ &= p(y|z) \cdot p(x|z) \cdot p(z) \end{aligned}$$

$$\Rightarrow p(x,y|z) = p(x|z) \cdot p(y|z)$$

$$\Rightarrow x \perp\!\!\!\perp y | z$$

# The Collider



$$P(x,y,z) = P(z|x,y) \cdot P(x) \cdot P(y)$$
$$\stackrel{!-dz}{\Rightarrow} P(x,y) = P(x) \cdot P(y)$$

∴ always:  $x \perp\!\!\!\perp y$   
generally:  $x \not\perp\!\!\!\perp y \mid z$

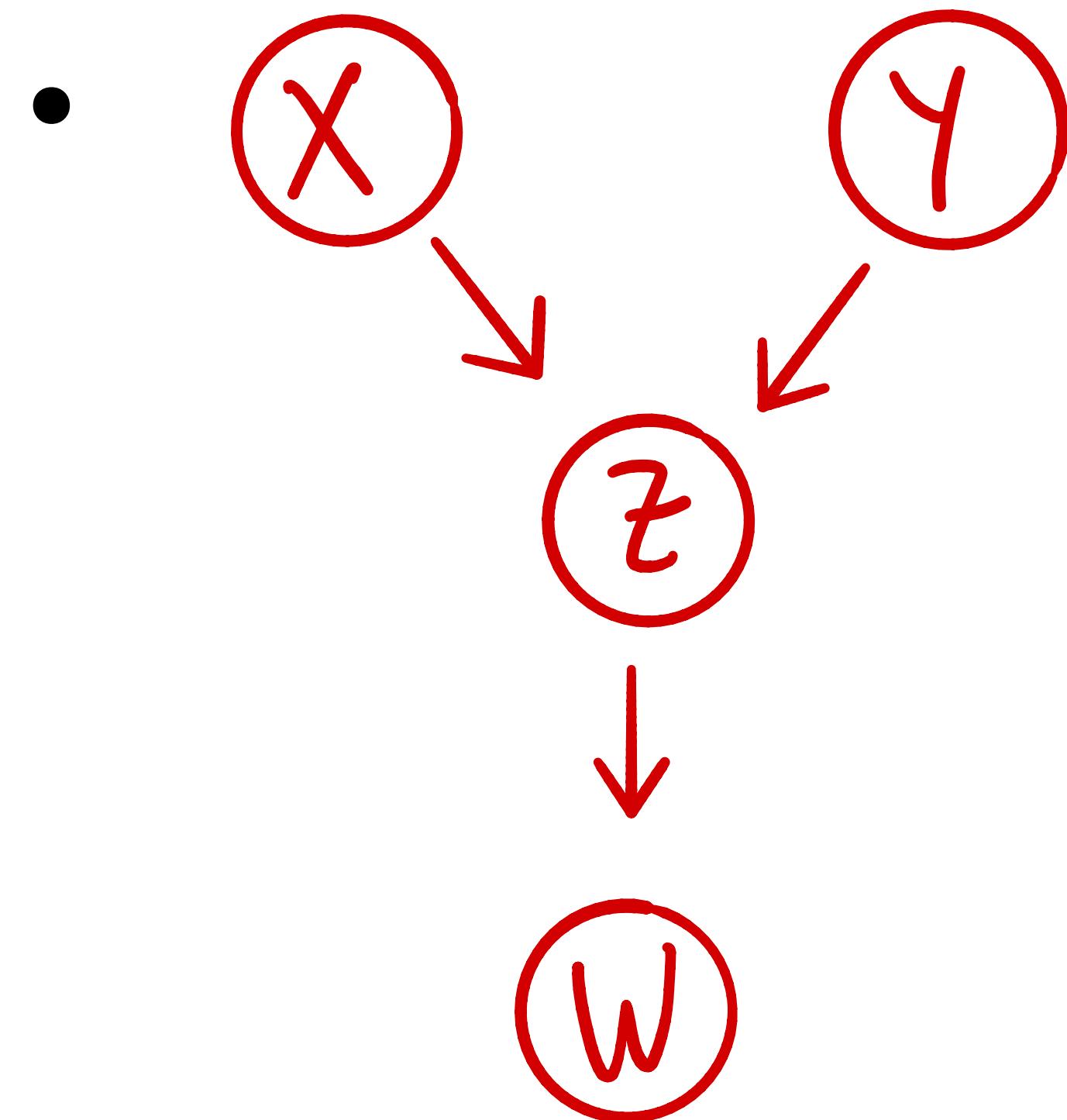
Example: Coin flips  $X, Y \in \{0, 1\}$

$$z = x + y$$

$$z=0 \Rightarrow X=0 \wedge Y=0$$

$$z=1 : X=0 \rightsquigarrow Y=1$$

# The Collider



$$p(x,y,z,w) = p(w|z) \cdot p(z|x,y) \cdot p(x) \cdot p(y)$$

generally:  $X \not\perp\!\!\!\perp Y | W$

always:  $X \perp\!\!\!\perp Y$

Example: as before  
 $w := z$

## Summary

- 

Fork :



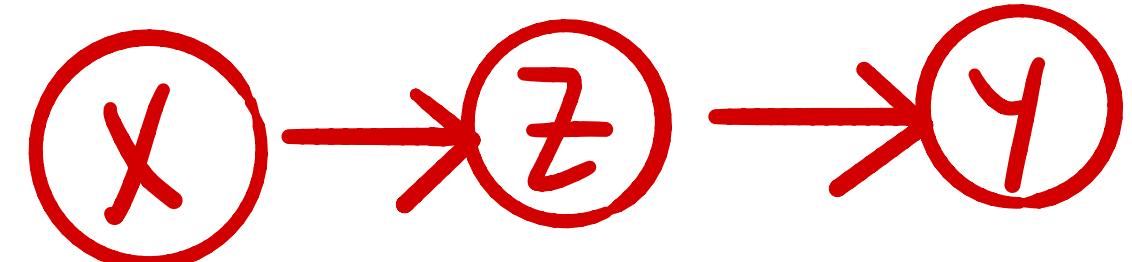
always

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

generally

$X \not\perp\!\!\!\perp Y$

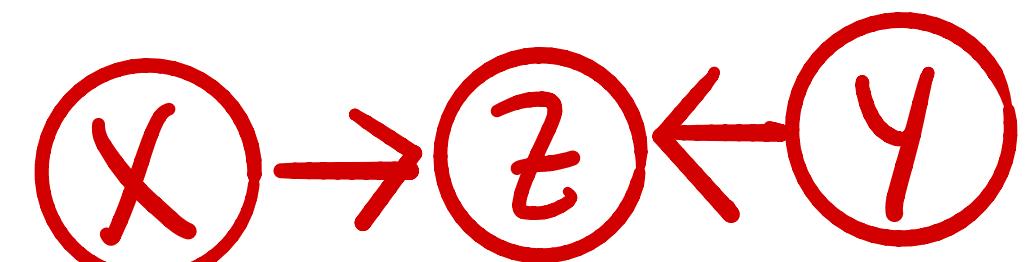
Chain :



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

$X \not\perp\!\!\!\perp Y$

Collider :



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

$X \not\perp\!\!\!\perp Y \mid Z$

# **Machine Learning 2**

**Graphical Models**

- Bayesian Networks**
- Global Markov Property**

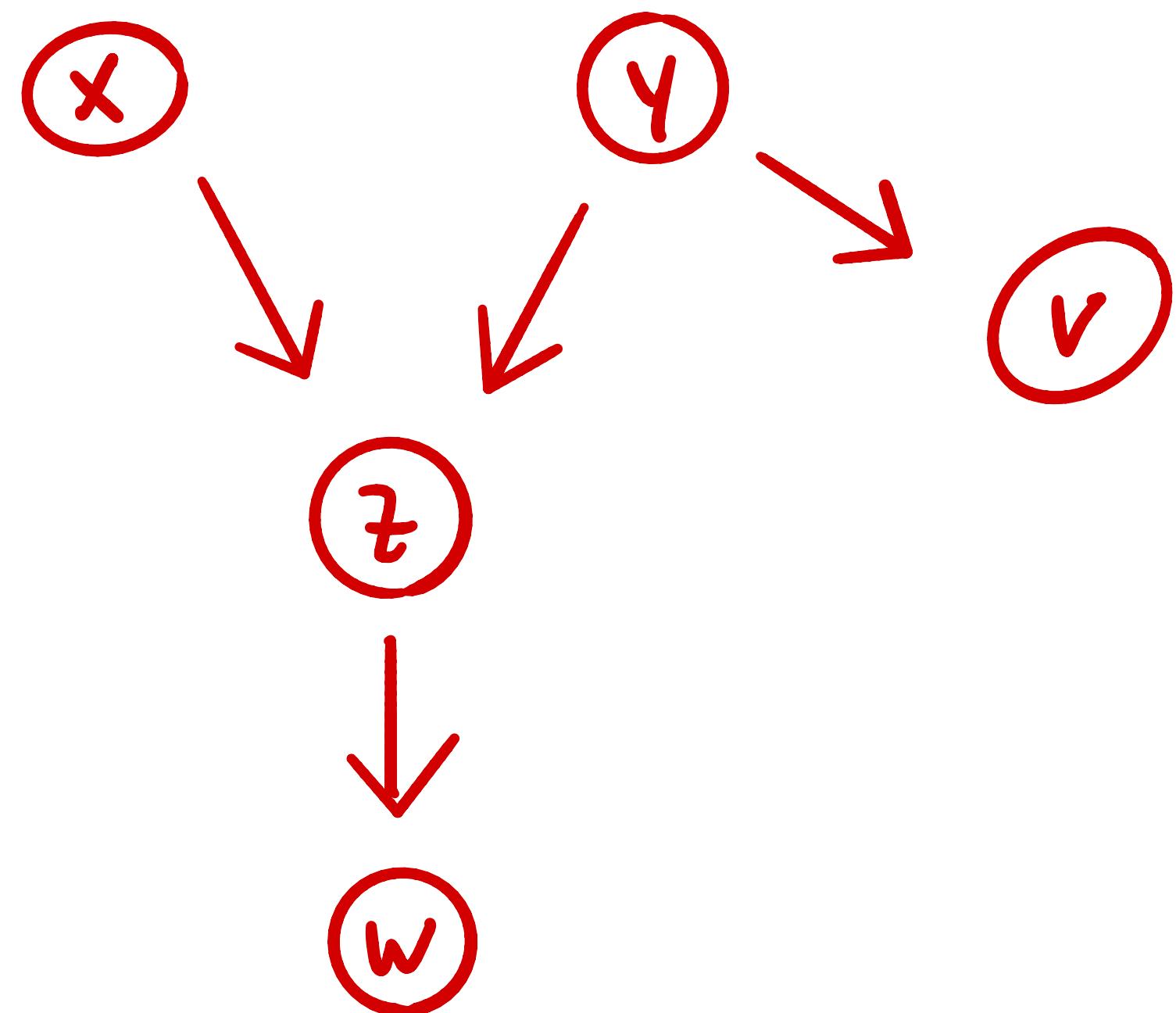
Patrick Forré

# Blocked and open paths in DAGs - Definition

- Let  $G = (V, E)$  be a DAG and  $C \subseteq V$  a subset,  $v, w \in V$  nodes.
- A path  $v = v_0 \rightarrow \dots \leftarrow \dots \rightarrow \dots v_n = w$  in  $G$  (any direction of arrow heads allowed) is called blocked by  $C$  if at least one of the following cases holds:
  - at least one of the end-nodes  $v$  or  $w \in C$ ,
  - there is a non-collider ( $\leftarrow v_i \rightarrow / \rightarrow v_i \rightarrow / \leftarrow v_i \leftarrow$ ) on the path with  $v_i \in C$ ,
  - there is a collider ( $\rightarrow v_i \leftarrow$ ) on the path with  $v_i \notin \text{Anc}^G(C)$ .
- Otherwise the path is called  $C$ -open or not blocked by  $C$ .
- If  $C = \emptyset$  then a path is blocked iff it contains at least one collider.

# Examples: blocked and open paths in DAGs

- 



- $x \rightarrow z \leftarrow y$  :  $\phi$ -blocked  
 $\{z\}$ -open,  $\{w\}$  open
- $x \rightarrow z \leftarrow y \rightarrow v$  :  $\phi$ -Blocked  
 $\{z\}$ -open,  $\{w\}$  open  
 $\{z, w\}$ -open  
 $\{z, w, y\}$  - Blocked  
 $\{y\}$ -Blocked

# d-Separation in DAGs - Definition

- Let  $G = (V, E)$  be a DAG and  $A, B, C \subseteq V$  subset of nodes.
- We say that  $A$  and  $B$  are d-separated by  $C$  in  $G$  if every path from a node  $v \in A$  to a node  $w \in B$  in  $G$  is blocked by  $C$ .
- In symbols:

$$A \perp B \mid C \quad \text{or} \quad A \xrightarrow[G]{d} B \mid C$$

$$A \perp B \Leftrightarrow A \perp B \mid \emptyset$$

# Exercise: Separoid Axioms for d-Separation

- Let  $G = (V, E)$  be a DAG and  $A, B, C, F \subseteq V$  be subset of nodes.

- Redundancy:  $A \perp B | A$  always holds.

- Symmetry:  $A \perp B | C \iff B \perp A | C$

- Decomposition:  $A \perp B \cup F | C \implies A \perp B | C$

- Weak Union:  $A \perp B \cup F | C \implies A \perp F | B \cup C$

- Contraction:  $A \perp F | B \cup C \wedge A \perp B | C \implies A \perp B \cup F | C$

- So we have the equivalence:

$$A \perp B \cup F | C \iff A \perp F | B \cup C \wedge A \perp B | C$$

# Theorem - Global Markov Property for BN

- Let  $(G, p)$  be a Bayesian Network (BN) with DAG  $G = (V, E)$  and random variables  $X_1, \dots, X_M$  with joint distribution  $p(x_V)$ . For a set of nodes  $F \subseteq V$  we write  $X_F = (X_v)_{v \in F}$ .
- For any three subsets  $A, B, C \subseteq V$  of nodes we have the implication:

$$A \perp B | C \implies X_A \perp\!\!\!\perp X_B | X_C.$$

i.e. d-separation implies corresponding conditional independence.

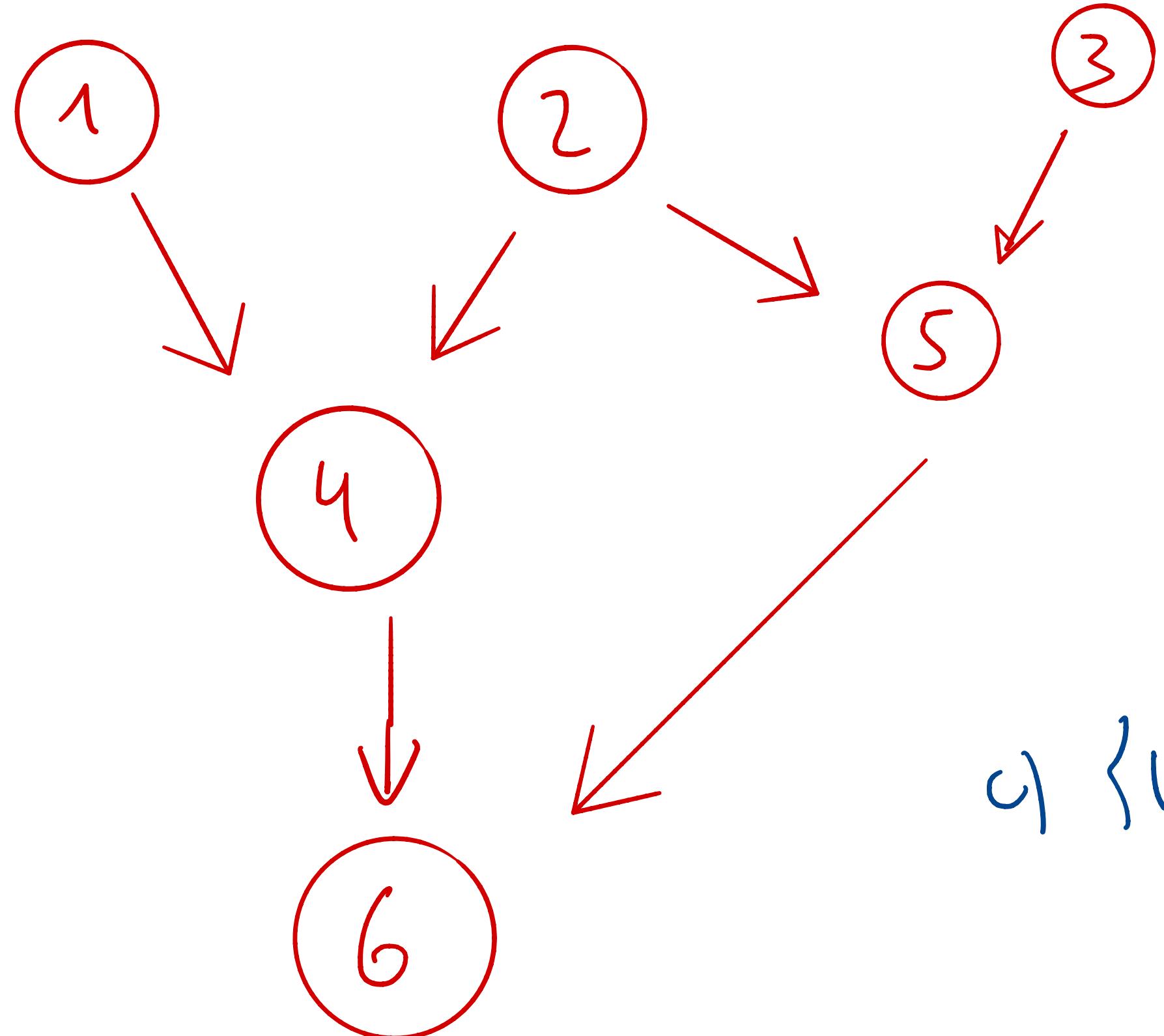
- The reverse implication is NOT true in general.

# Exercise - Prove Global Markov Property - Sketch

- Induction by the number of nodes.
- Check that removing a **child-less** node  $w \in V$  (exists because of acyclicity) doesn't hurt factorization property for remaining variables.
- Check (\*):  $X_w \perp\!\!\!\perp X_{-w} | X_{\text{Pa}^G(w)}$ .
- Check that w.l.o.g.  $A, B, C$  can be assumed to be pairwise disjoint.
- Check 4 cases: 0)  $w \notin A \cup B \cup C$ , 1)  $w \in A$ , 2)  $w \in B$ , 3)  $w \in C$ .
- For those use induction step, relation (\*) and the separoid axioms.
- For case 3) show that  $A \perp B | C' \cup \{w\}$  implies  $A \cup \{w\} \perp B | C'$  or  $A \perp B \cup \{w\} | C'$  and then use case 1) and 2).

# Examples: d-Separation and Conditional Independence

- 



$$a) \{v_4\} \perp \{v_5\} | \{v_2\} \Rightarrow X_4 \perp\!\!\!\perp X_5 | X_2$$

$$b) \{v_1, v_2, v_3\} \perp \{v_6\} | \{v_4, v_5\} \\ \Rightarrow (X_1, X_2, X_3) \perp\!\!\!\perp X_6 | (X_4, X_5)$$

$$c) \{v_3\} \perp \{v_6\} | \{v_5, v_4\} \Rightarrow X_3 \perp\!\!\!\perp X_6 | X_5, X_4$$

# **Machine Learning 2**

## **Graphical Models**

- Bayesian Networks**
- Construction of BNs**

Patrick Forré

# Bayesian Networks - What are they good for?

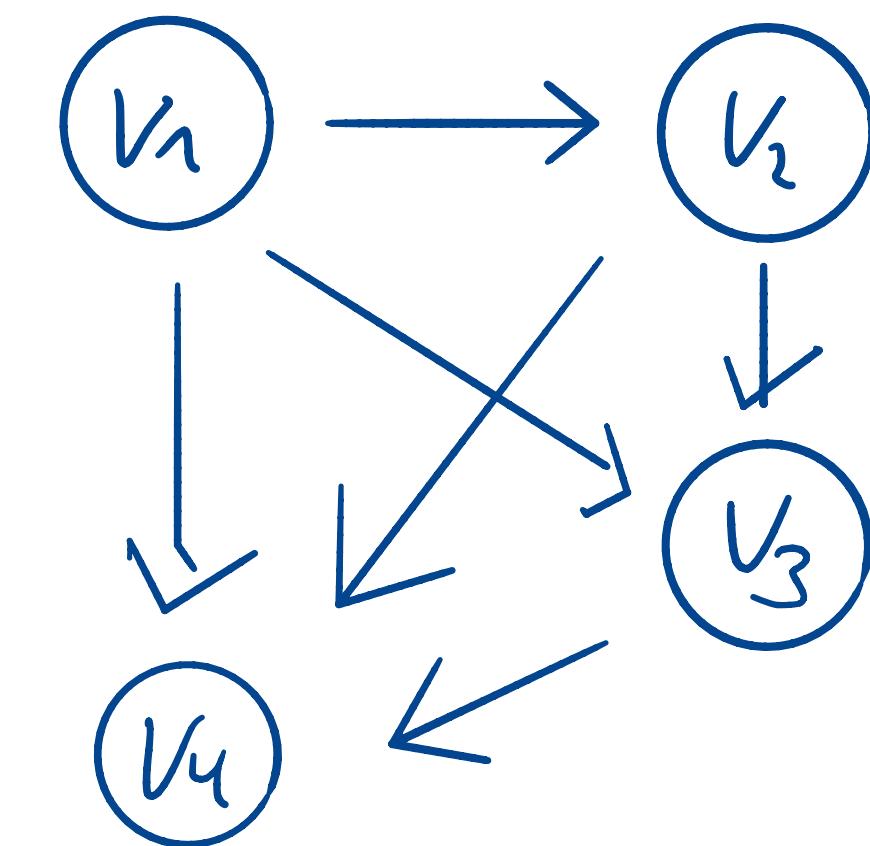
- Let  $p(x_1, \dots, x_M)$  be the joint distribution of discrete random variables  $X_1, \dots, X_M$ , where each  $X_m$  can take on  $K$  values.
- A naive table representation of  $p(x_V)$  needs  $\sim K^M$  entries.
- If  $(G, p)$  is a 'sparse' BN, where every node has at most  $L \ll M$  parents.
- Representing  $p(x_v | x_{\text{Pa}^G(v)})$  in table form each takes  $\sim K^{L+1}$  entries.
- So  $p(x_V) = \prod_{v \in V} p(x_v | x_{\text{Pa}^G(v)})$  takes  $\sim M \cdot K^{L+1} \ll K^M$  entries.
- $\implies$  BNs model independence relations and allow for sparse representations.

# Construction of bad Bayesian Networks

- Let  $p(x_1, \dots, x_M)$  be the joint distribution and  $V = \{1, \dots, M\}$ .
- We want to construct a DAG  $G = (V, E)$  such that  $(G, p)$  is a BN.
- **Product rule:**

$$p(x_1, \dots, x_M) = \prod_{m=1}^M p(x_m | \underbrace{x_1, \dots, x_{m-1}}_{\text{Pa}_m^G(V_m)})$$

$p(x_m | x_1, \dots, x_{m-1}) \sim k^M$



- So **fully connected** DAG always turns  $(G, p)$  into BN.
- Note, different ordering of  $X_1, \dots, X_M$  results in different DAG.

# Construction of better Bayesian Networks

- Inductively, take a **minimal subset**  $P_m \subseteq \{1, \dots, m-1\}$  such that:
  - $X_m \perp\!\!\!\perp (X_1, \dots, X_{m-1}) | X_{P_m}$ .
- Then:  $p(x_m | x_1, \dots, x_{m-1}) = p(x_m | X_{P_m})$
- Include in  $G$  edges from nodes  $w \in P_m$  to  $v_m$ , then:
  - $G$  is a DAG,
  - $\text{Pa}^G(v_m) = P_m$ ,
  - $(G, p)$  is a BN with **minimal number of parents** (given the total ordering).
  - Note, different ordering of  $X_1, \dots, X_M$  results in different DAG and BN.

$$p(x_V) = \prod_{m=1}^M p(x_m | X_{P_m}) = \prod_{m=1}^M p(x_m | X_{P_m})$$

# **Machine Learning 2**

**Graphical Models**

- Bayesian Networks**
  - Summary**

Patrick Forré

# Bayesian Networks - Summary

- BN consists of a DAG  $G = (V, E)$  and distribution  $p(x_V) = \prod_{v \in V} p(x_v | x_{\text{Pa}_G(v)})$ .
- Introduced Plate notation with parameters and observations.
- d-Separation:  $A \perp B | C$  iff every path from  $A$  to  $B$  in  $G$  either contains a non-collider in  $C$  or a collider NOT in  $\text{Anc}^G(C)$ .
- Global Markov Property of BN:  $A \perp B | C \implies X_A \perp\!\!\!\perp X_B | X_C$ .
- Construction of BN for  $p(x_V)$ : Inductively, take a minimal subset  $\text{Pa}_m \subseteq \{1, \dots, m-1\}$  s.t.  $X_m \perp\!\!\!\perp X_{[1:m]} | X_{\text{Pa}_m}$  as the parents of  $v_m$ .
- BNs model independence relations and allow for sparse representations.