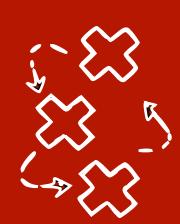


# Deep learning 2: Causality & DL

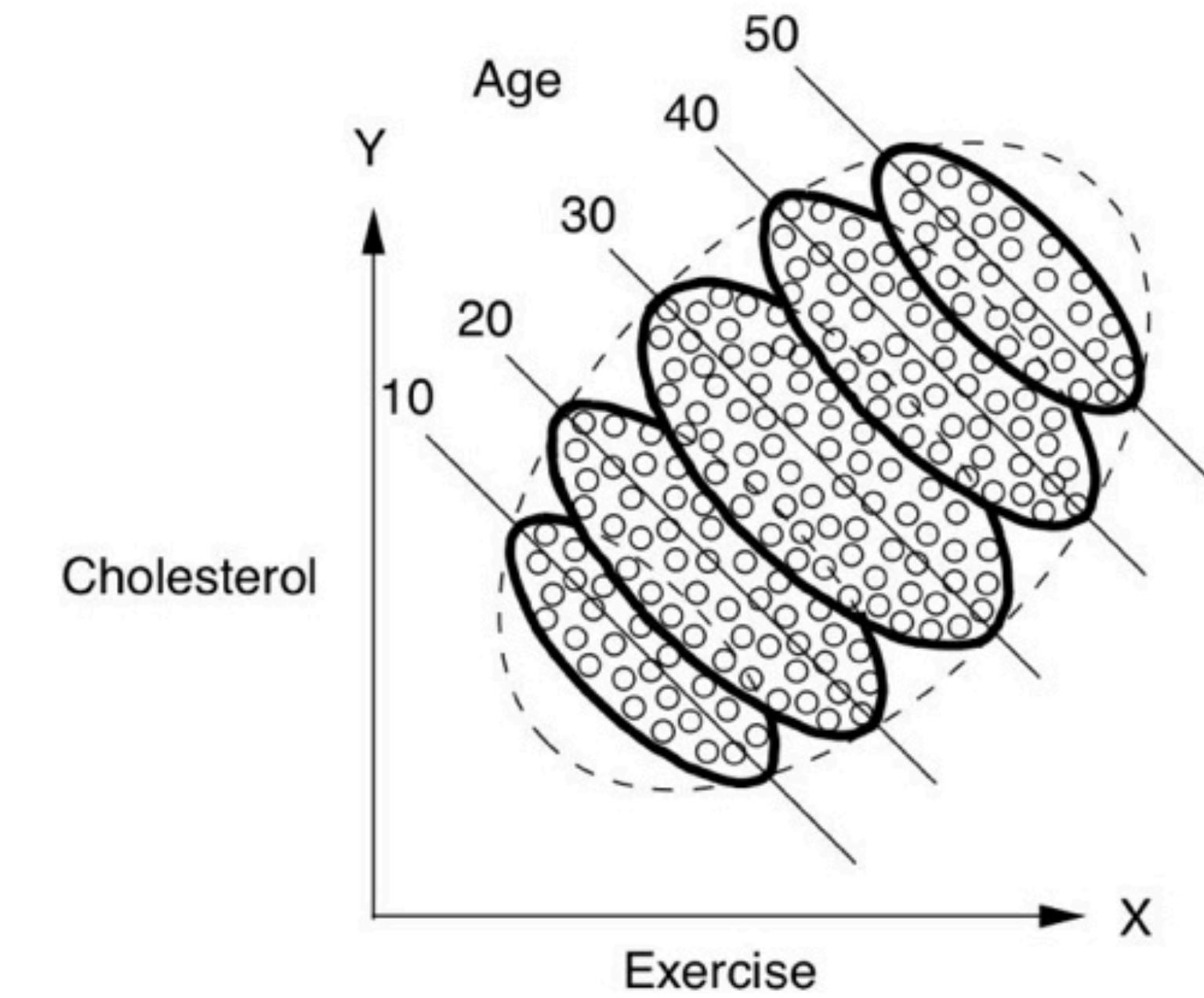
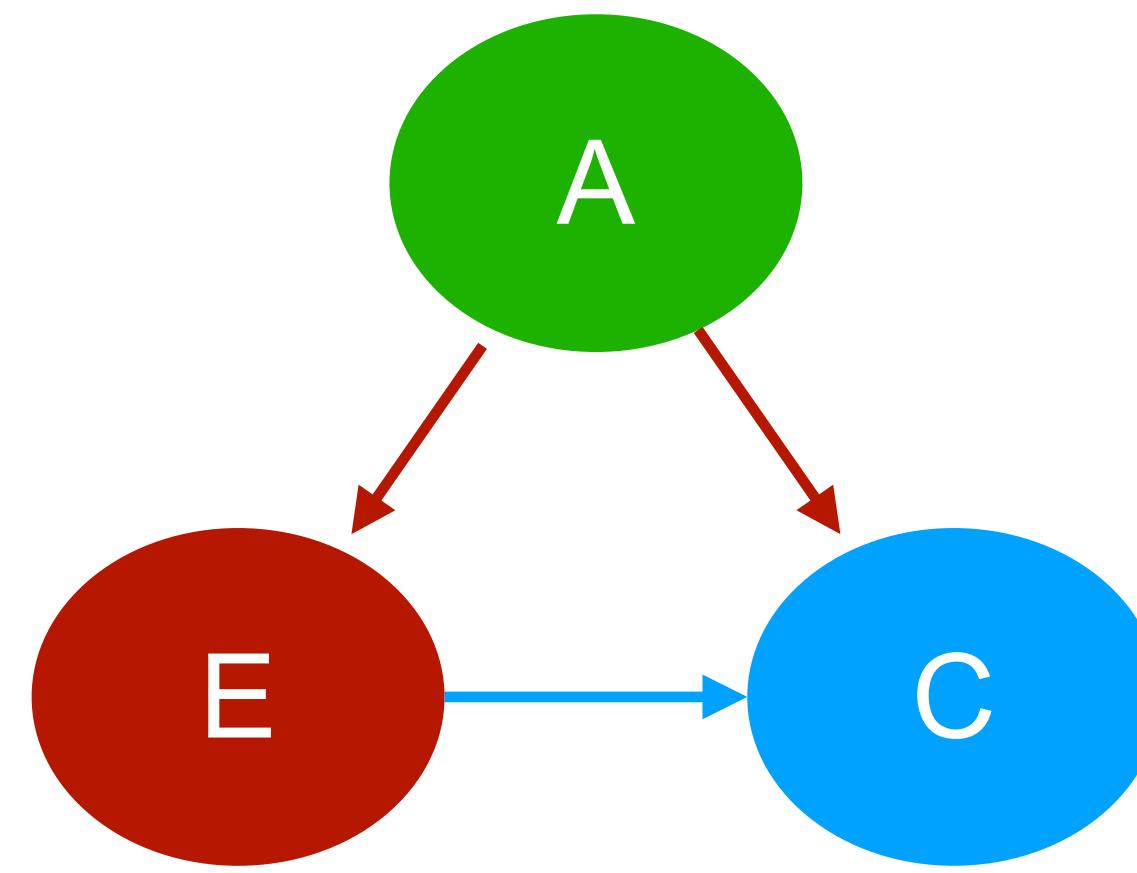
## 1.2: Graphical models

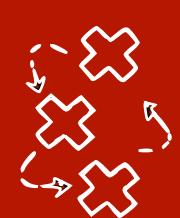
Lecturer: Sara Magliacane

UvA - Spring 2022



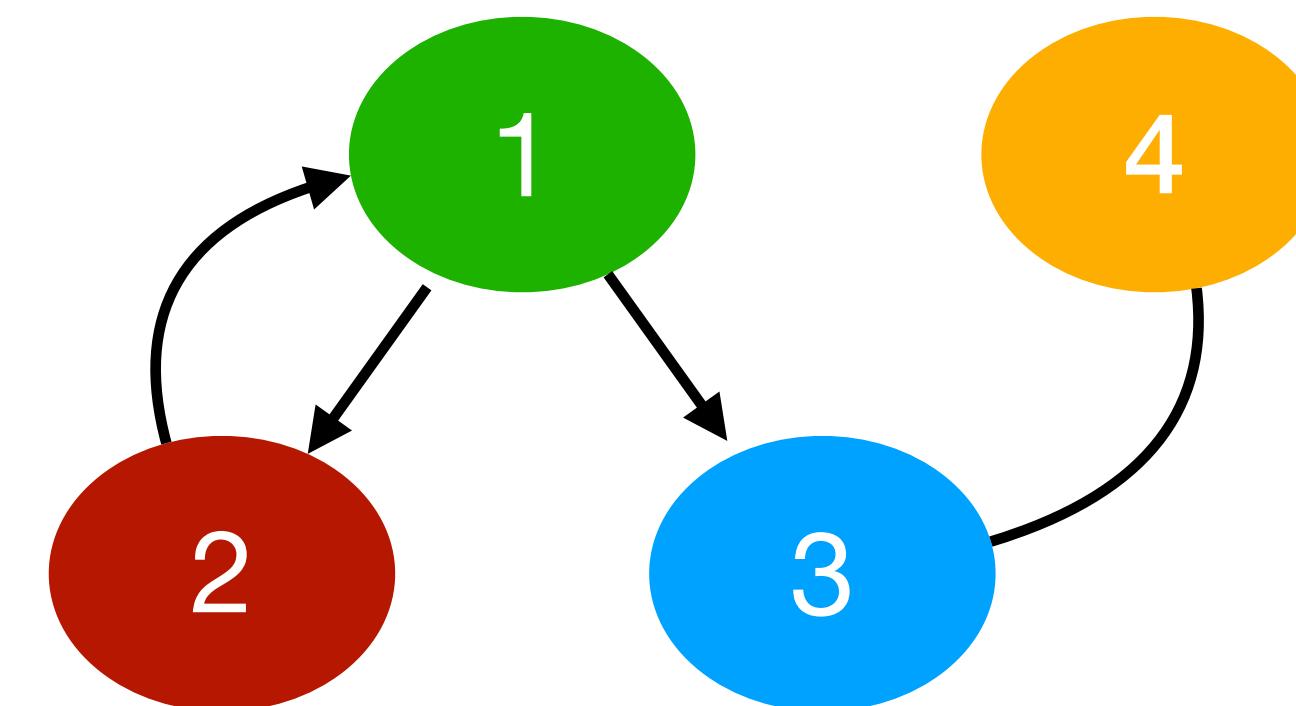
# Graphical models



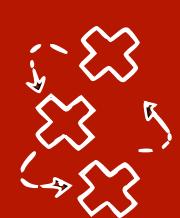


# Graph terminology

- A graph  $G$  is a tuple  $G = (V, E)$ :
  - $V$  is the set of **nodes** (vertices)
  - $E$  is the set of **edges** between two nodes, i.e.  $E \subseteq V \times V$

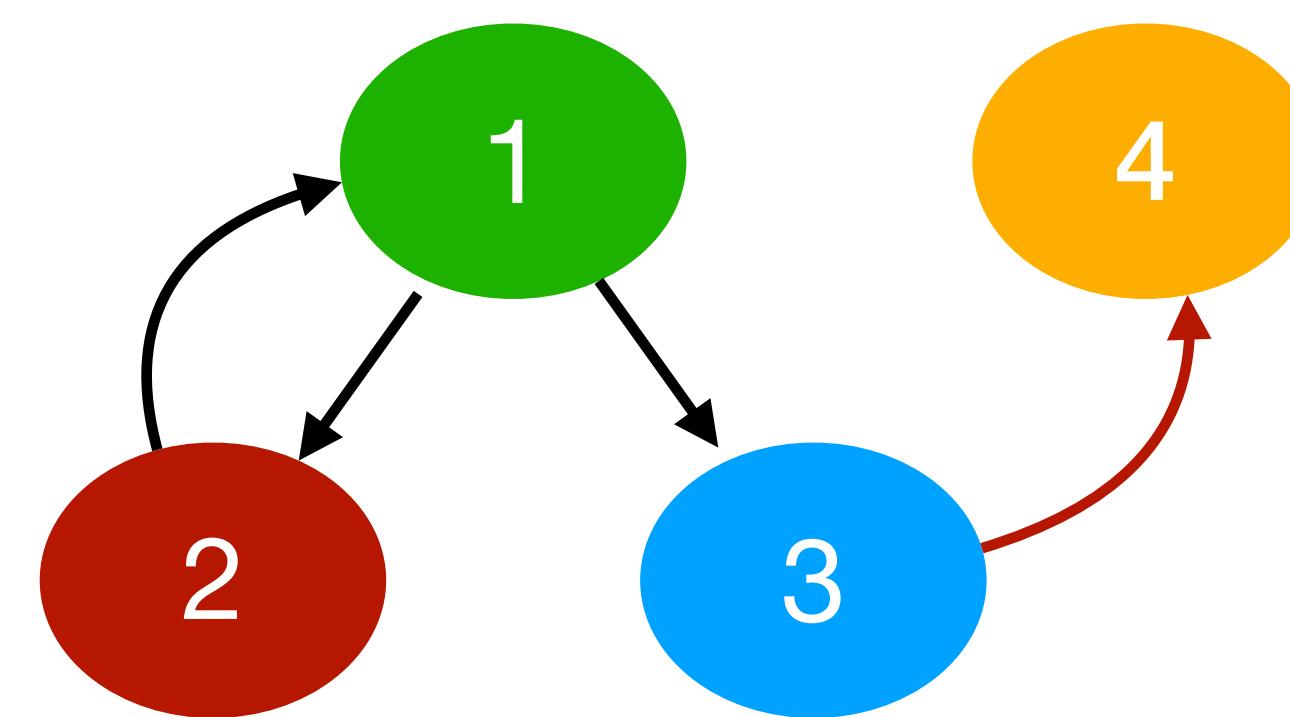


- Two nodes connected by an edge are **adjacent**

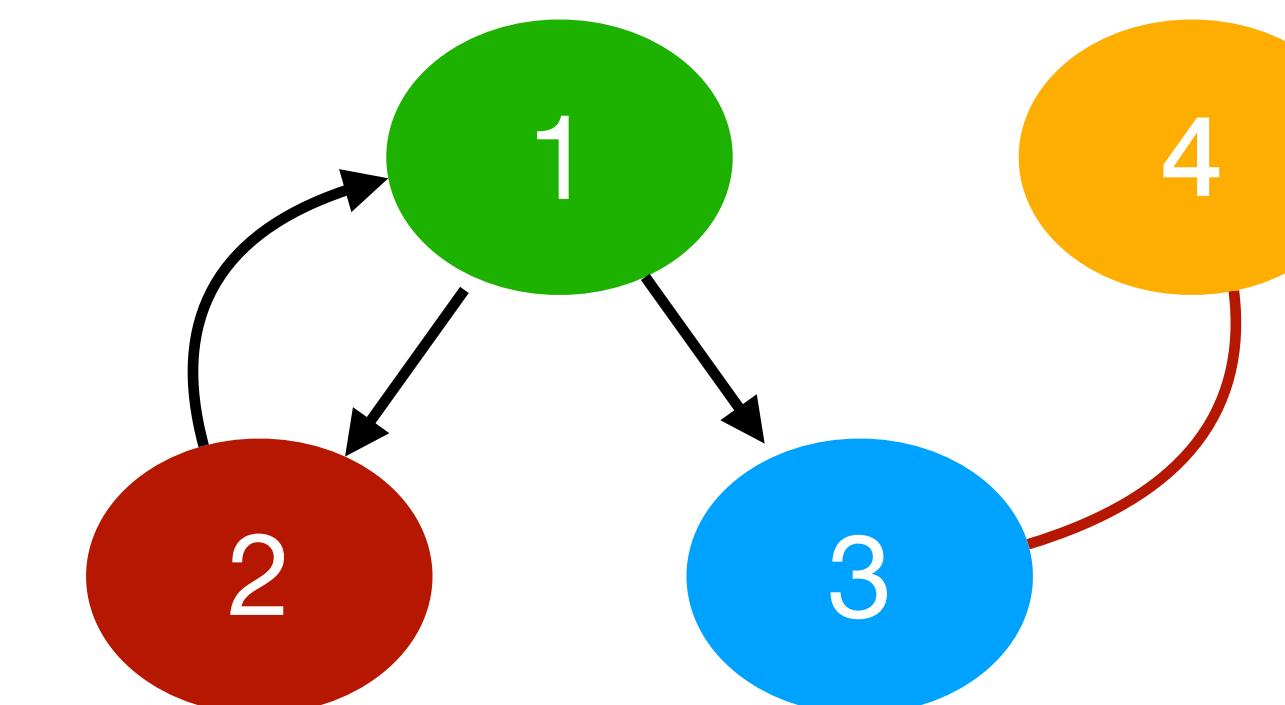


# Directed, mixed, undirected graphs

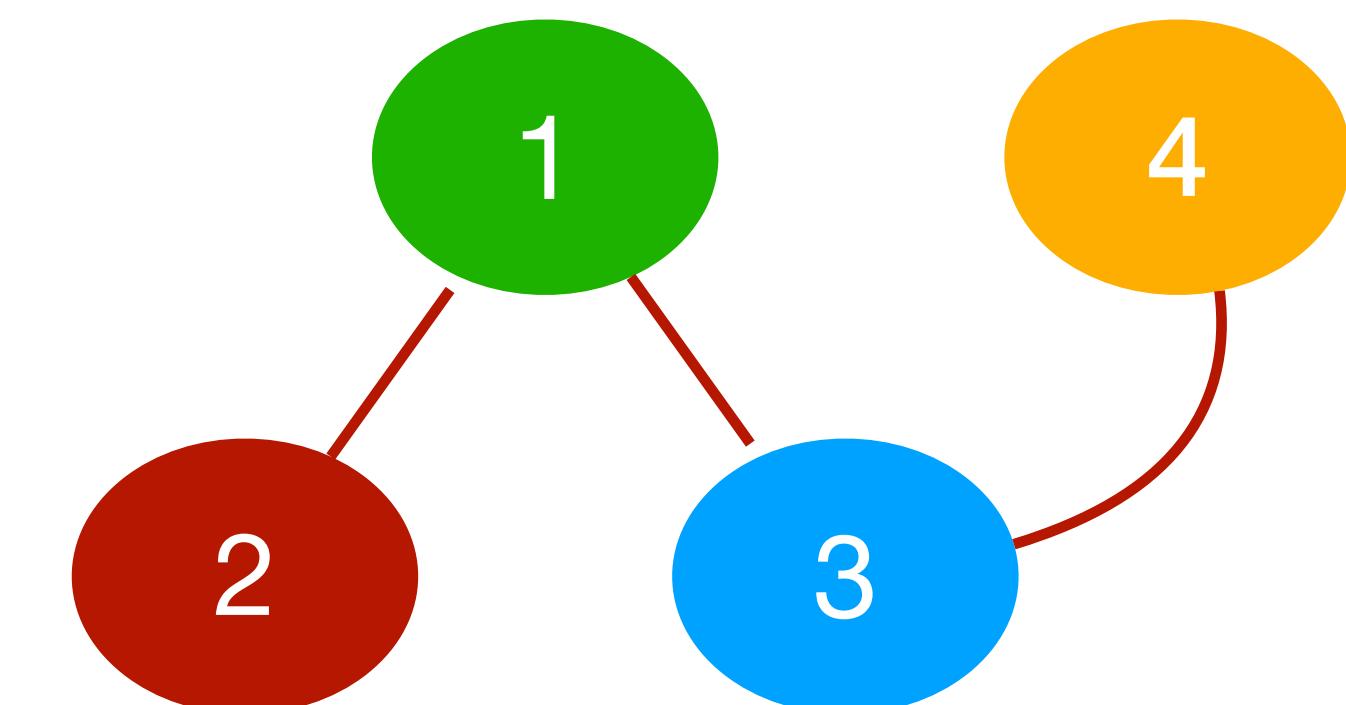
- A graph  $G$  is a tuple  $G = (\mathbf{V}, \mathbf{E})$ :
  - $\mathbf{V}$  is the set of **nodes** (vertices)
  - $\mathbf{E}$  is the set of **edges** between two nodes, i.e.  $\mathbf{E} \subseteq \mathbf{V} \times \mathbf{V}$
  - If all edges are **directed** → then the graph is **directed**



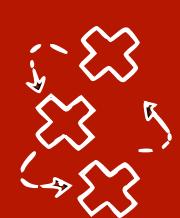
Directed graph



Mixed graph

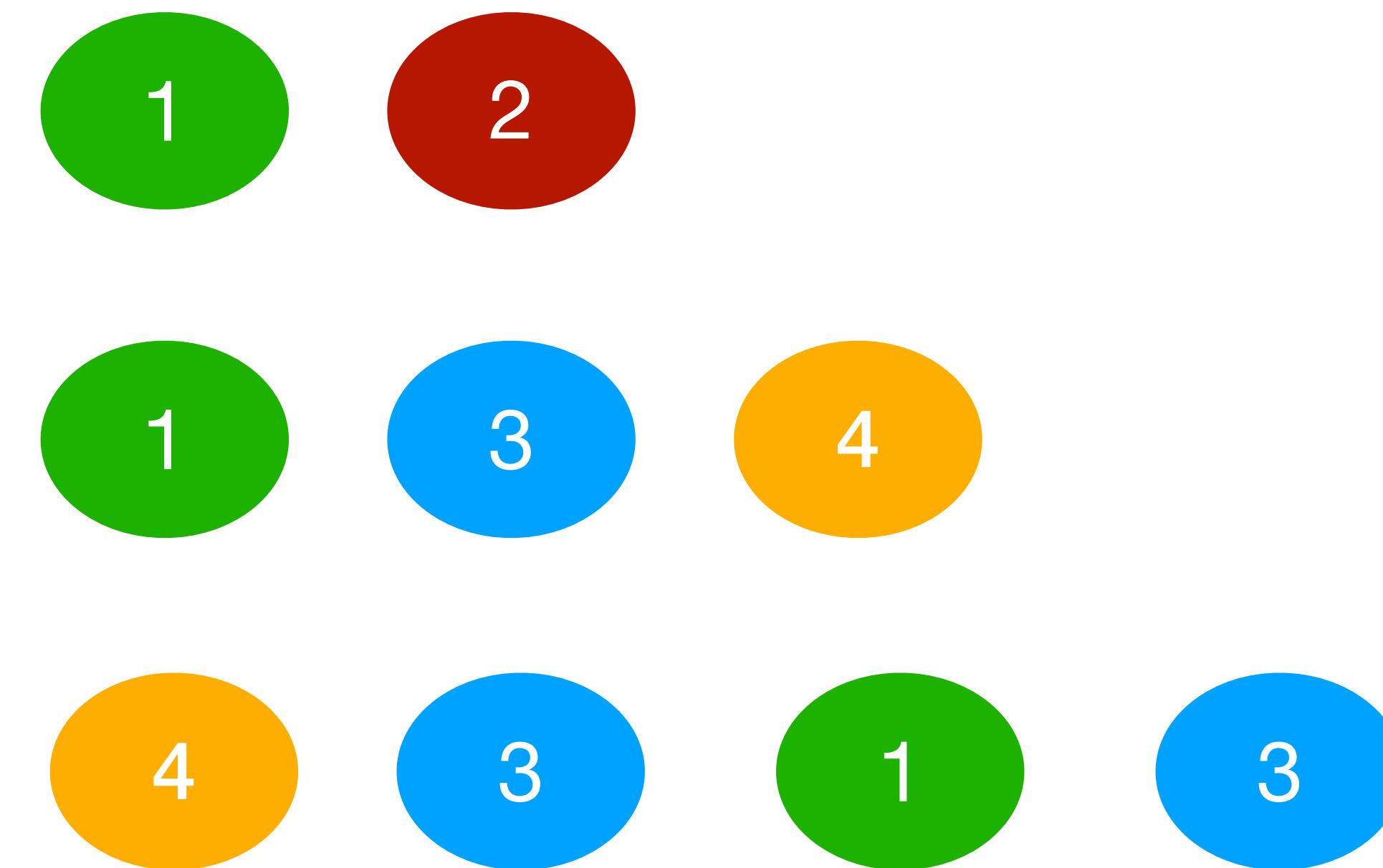
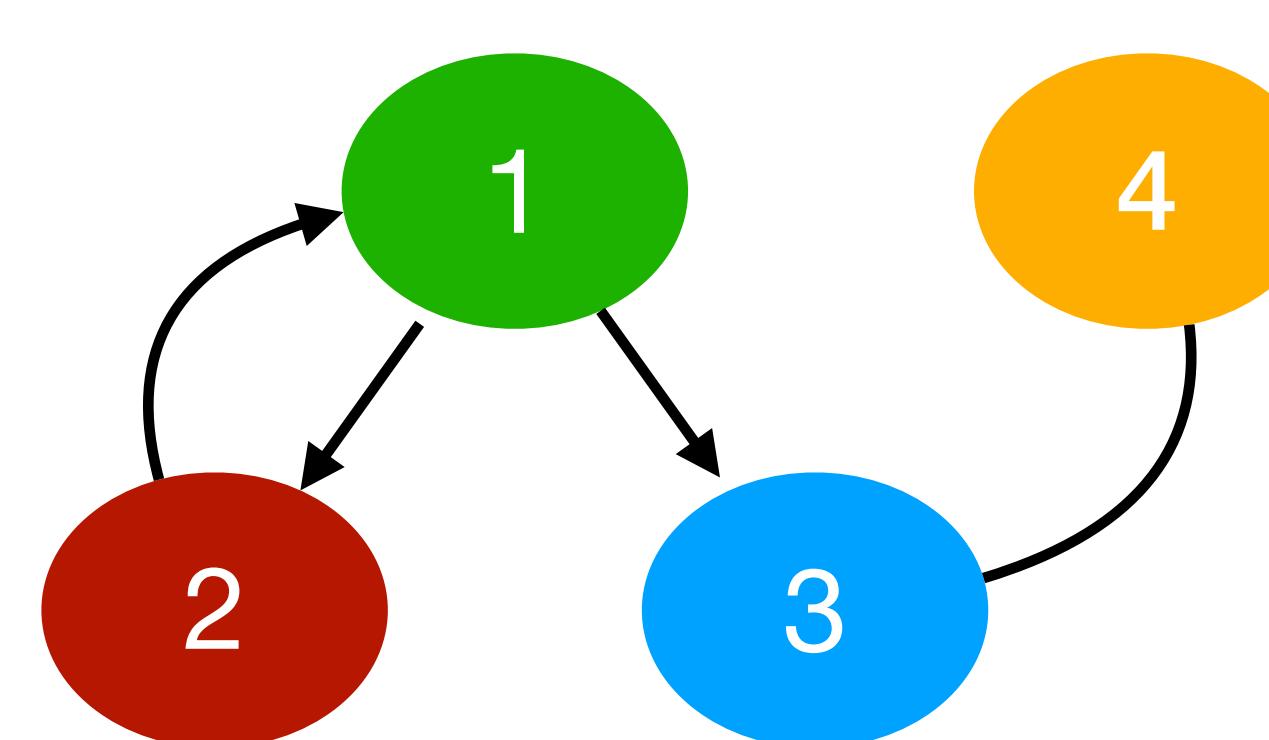


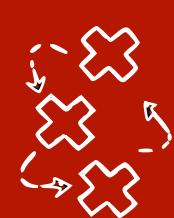
Undirected graph



# Graph terminology: paths

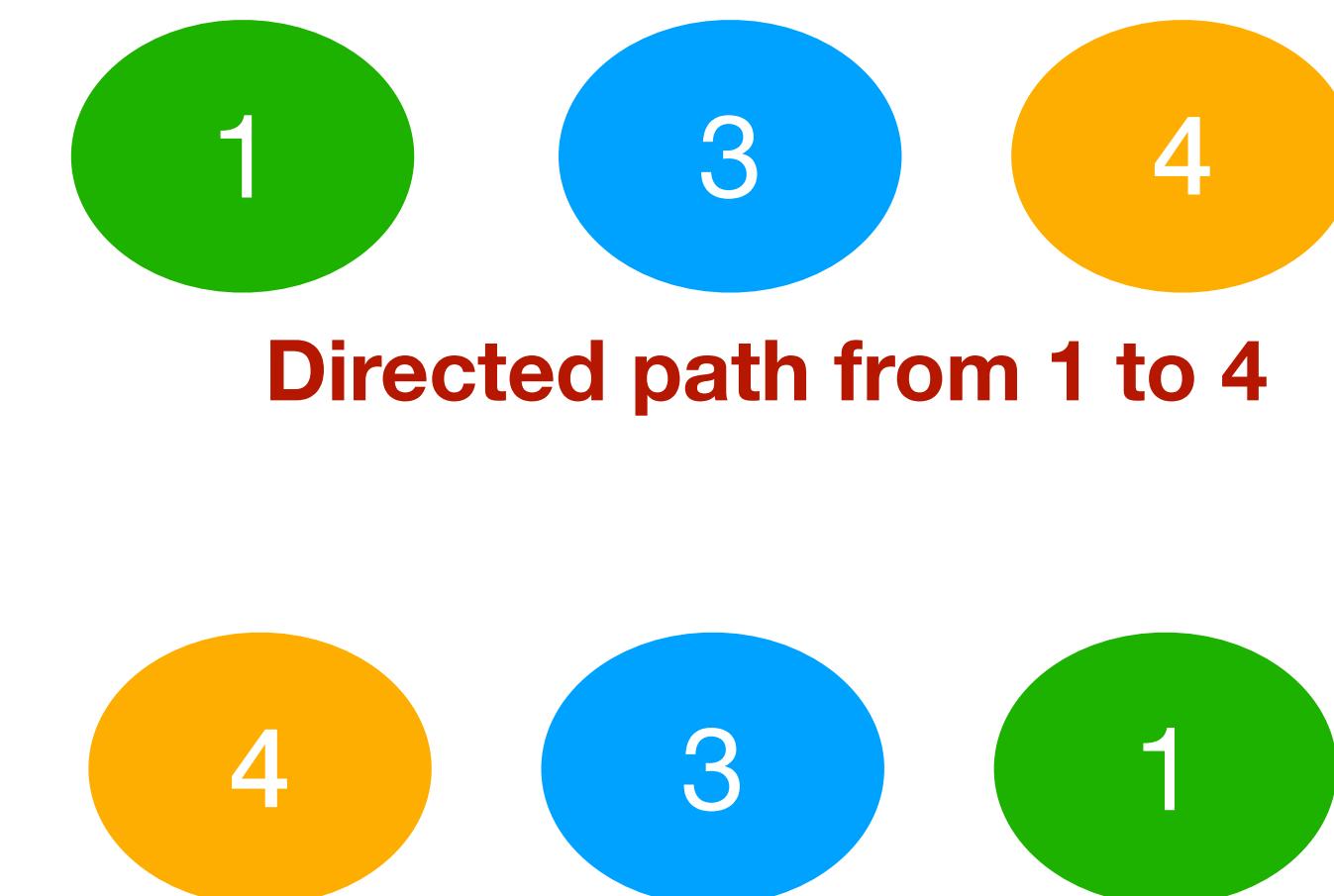
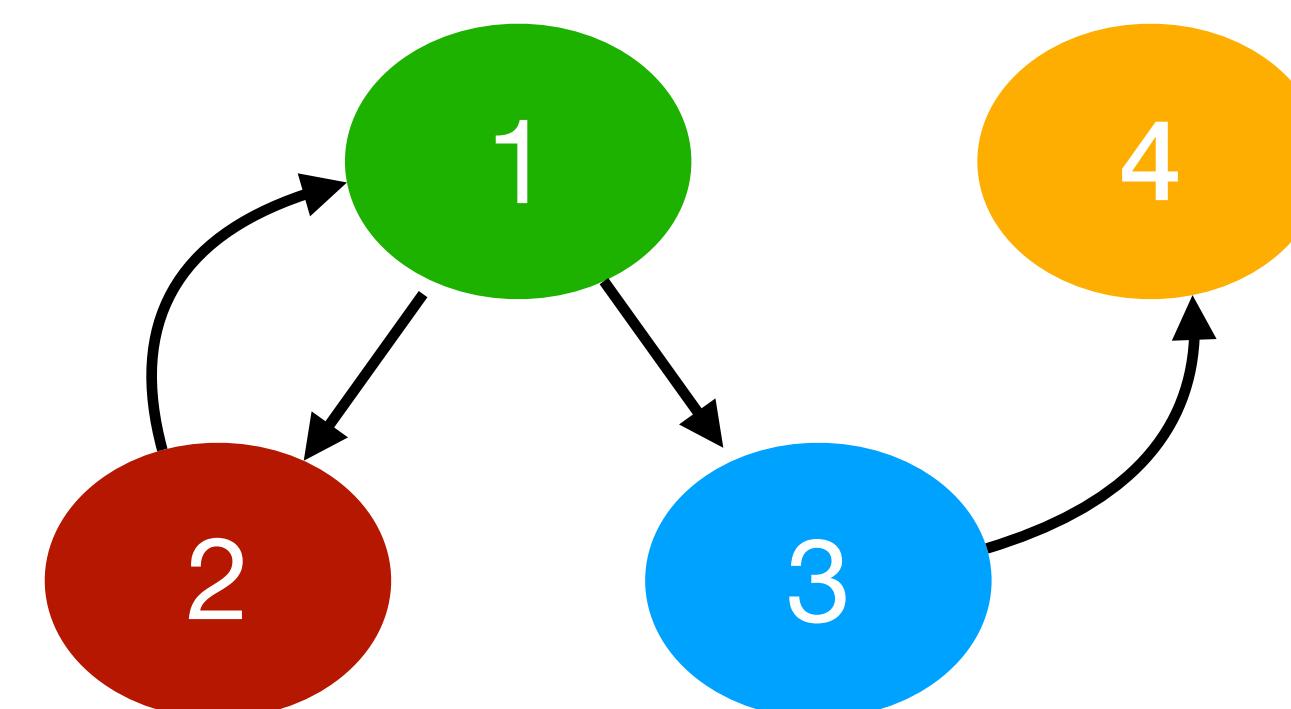
- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**



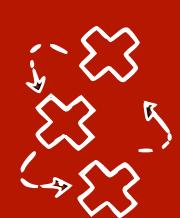


# Directed graphs: paths vs directed paths

- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**
- A **directed path** between **node i and node j** is a path where **all edges point towards j**, i.e.  $i \rightarrow \dots \rightarrow j$

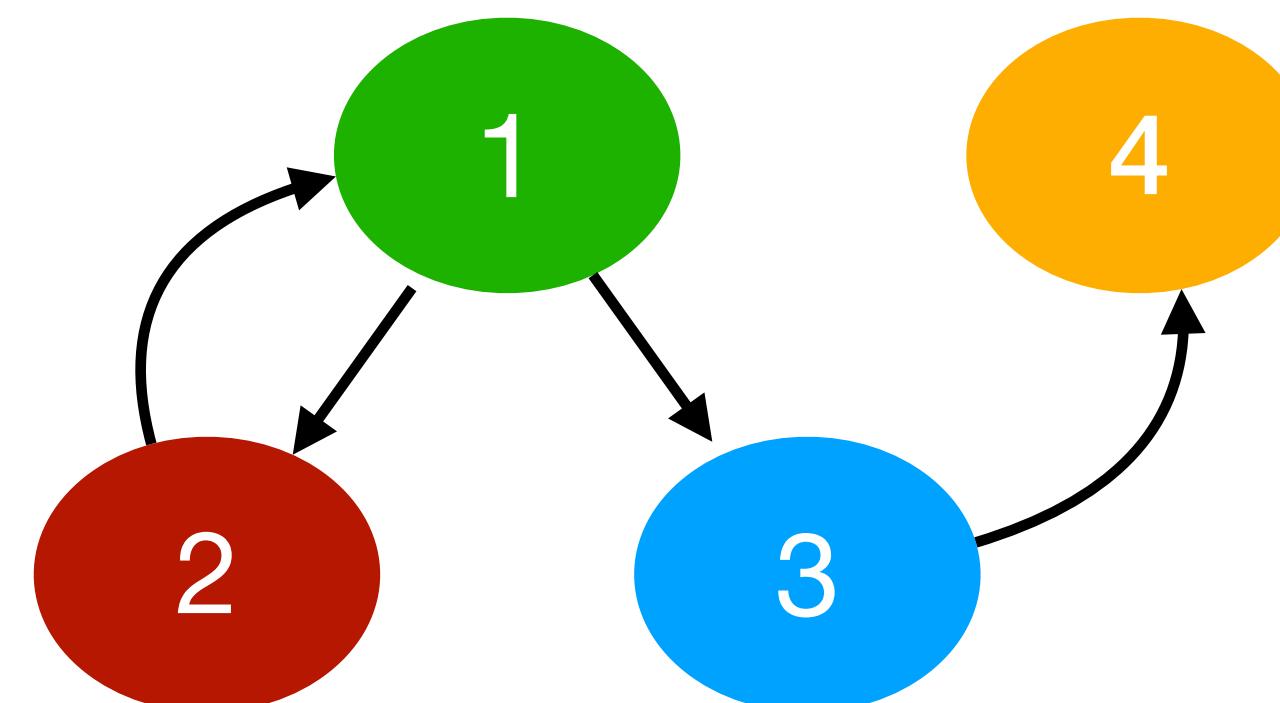


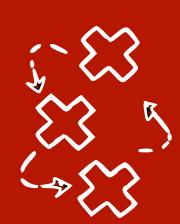
A path from 4 to 1, but not a directed path



# Directed graphs: paths vs directed paths

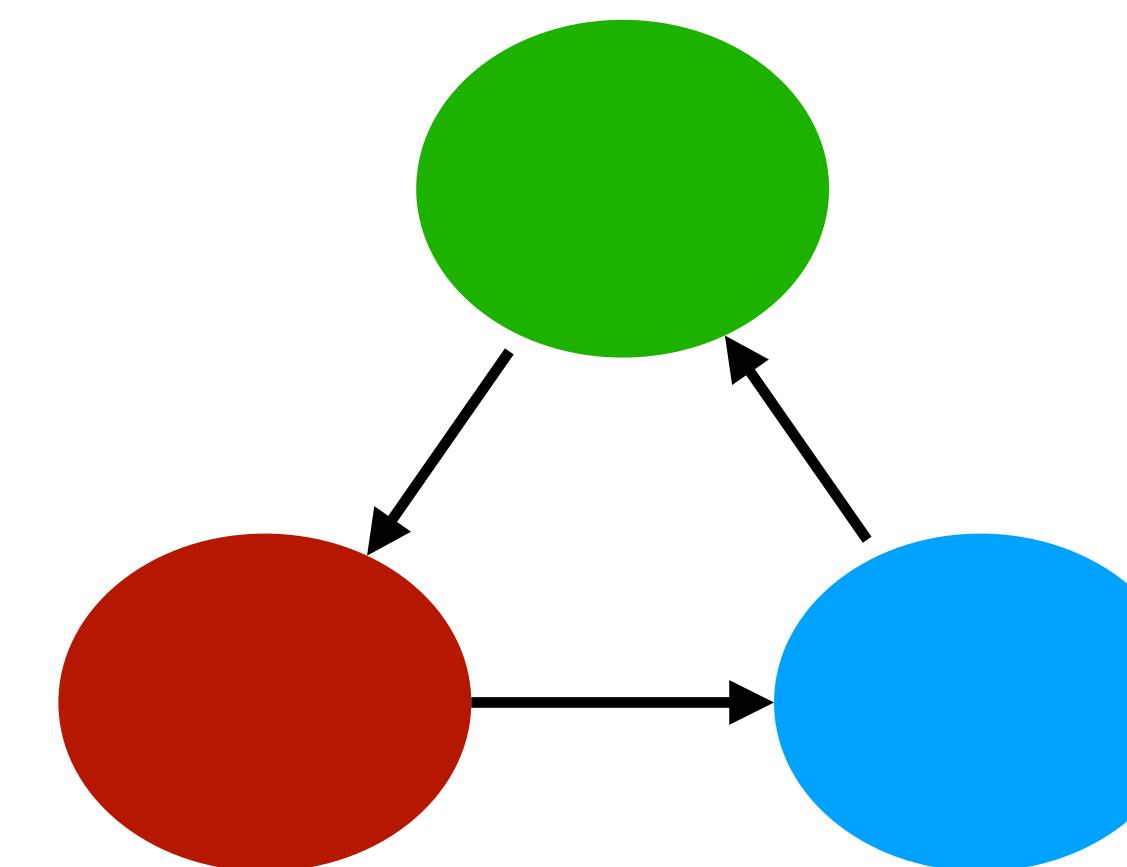
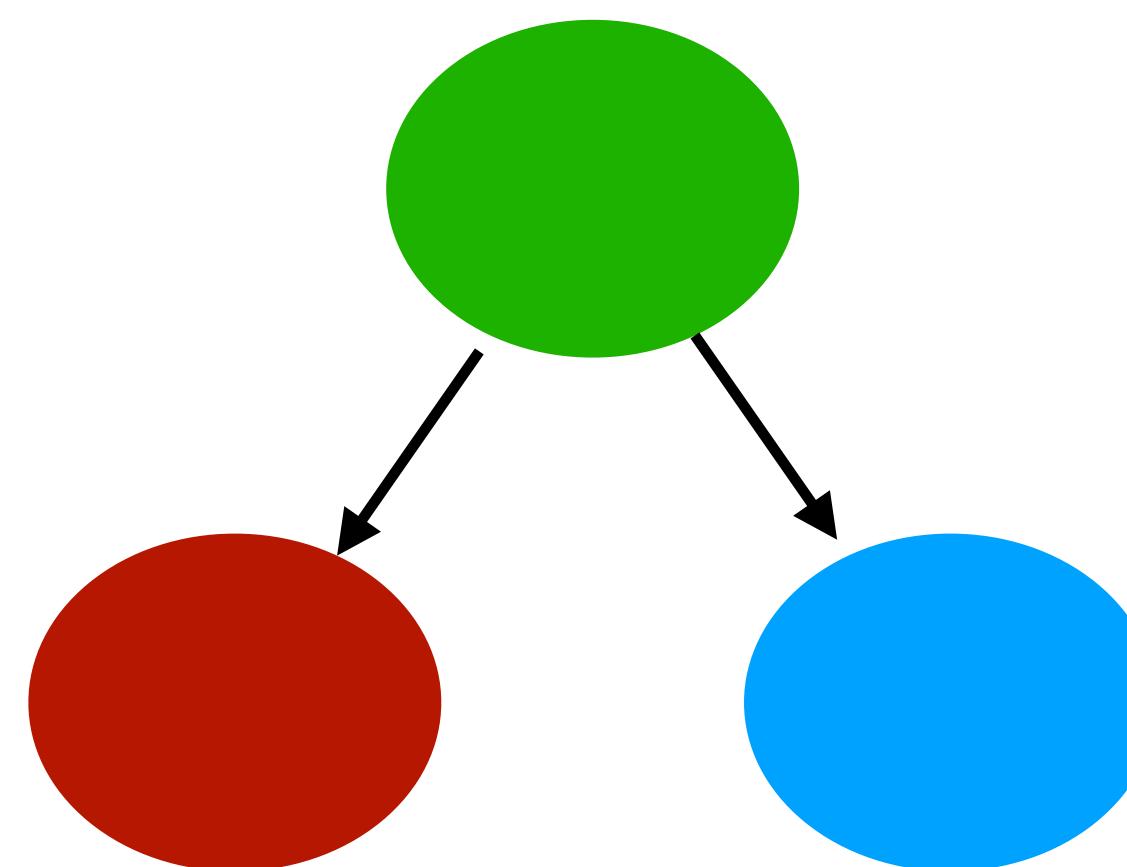
- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**
- A **directed path** between **node i and node j** is a path where **all edges point towards j**, i.e.  $i \rightarrow \dots \rightarrow j$
- A **directed cycle** is a directed path  $(i, \dots, j) + (j, i)$

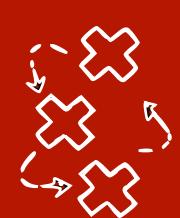




# Directed Acyclic Graphs (DAGs)

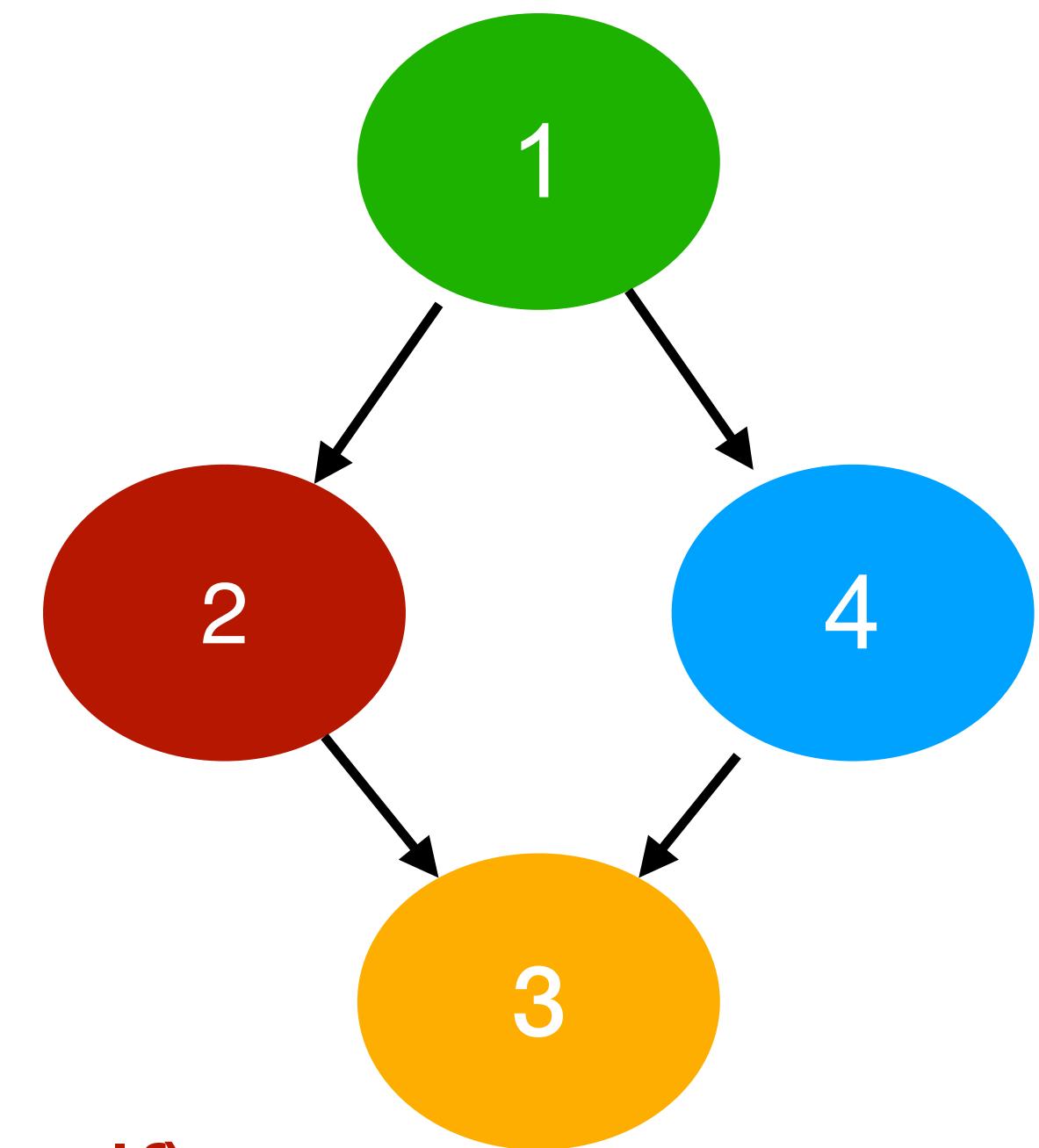
- A DAG is a directed graph  $G = (V, E)$ :
  - $V$  is the set of **nodes** (vertices)
  - $E$  is the set of **directed edges** between the nodes
  - There are **no directed cycles**

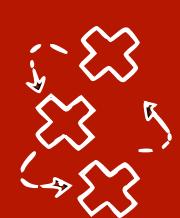




# Relationships between nodes in a DAG

- **Parents** of a node  $\text{Pa}_G(V)$ 
  - Nodes that have an edge pointing to  $V$
- **Children** of a node  $\text{Ch}_G(V)$ 
  - Nodes that have an edge pointing from  $V$
- **Ancestors** of a node  $\text{An}_G(V)$ 
  - Nodes that have a **directed path** to  $V$  (including  $V$  itself)
- **Descendants** of a node  $\text{Desc}_G(V)$ 
  - Nodes that are reached from  $V$  via **directed paths** (including  $V$  itself)

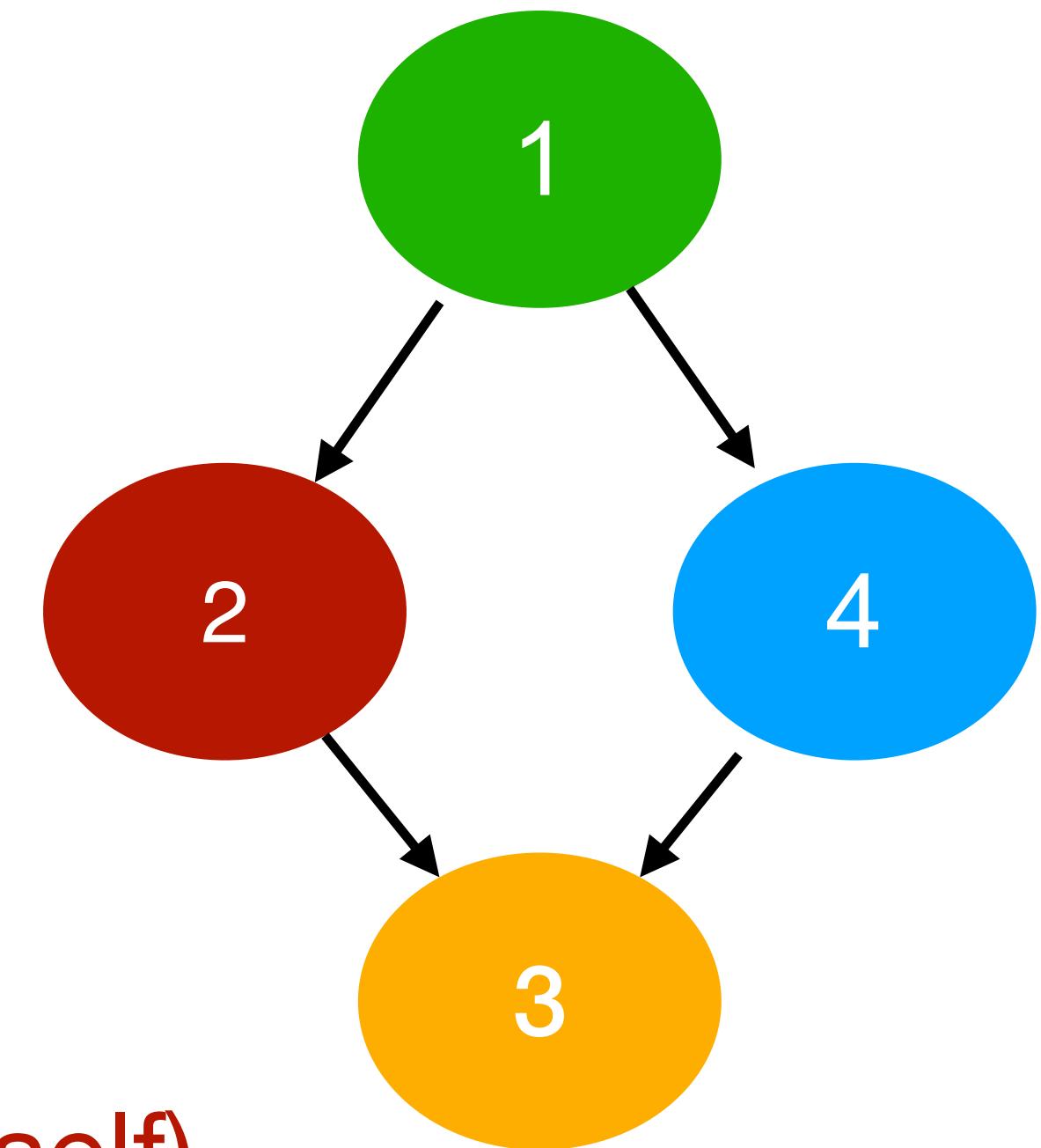


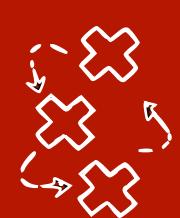


# Relationships between nodes in a DAG

[We omit  $G$ , when it is clear from the context]

- **Parents** of a node  $\text{Pa}(V)$ 
  - Nodes that have an edge pointing to  $V$
- **Children** of a node  $\text{Ch}(V)$ 
  - Nodes that have an edge pointing from  $V$
- **Ancestors** of a node  $\text{An}(V)$ 
  - Nodes that have a **directed path** to  $V$  (including  $V$  itself)
- **Descendants** of a node  $\text{Desc}(V)$ 
  - Nodes that are reached from  $V$  via **directed paths** (including  $V$  itself)

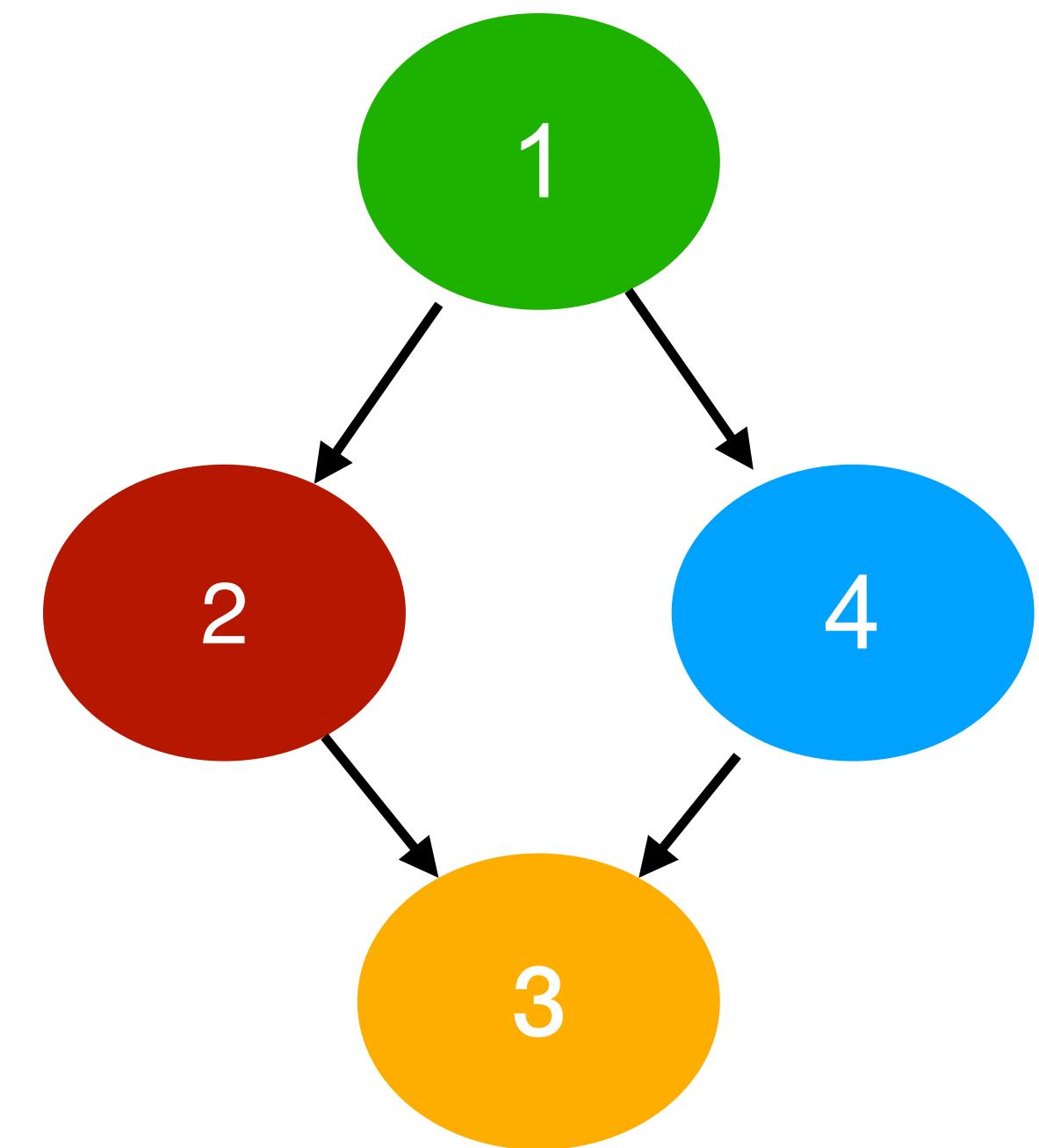




# Kinship relationships for sets of nodes

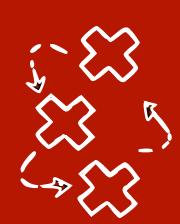
- We will use **bold** for sets (including sets of nodes)
- **Parents** of a set of nodes  $\mathbf{A} \subseteq V$ :

$$\text{Pa}(\mathbf{A}) := \bigcup_{V \in \mathbf{A}} \text{Pa}(V)$$



$$\text{Pa}(\{2,3\}) = \text{Pa}(2) \cup \text{Pa}(3) = \{1\} \cup \{2,4\} = \{1,2,4\}$$

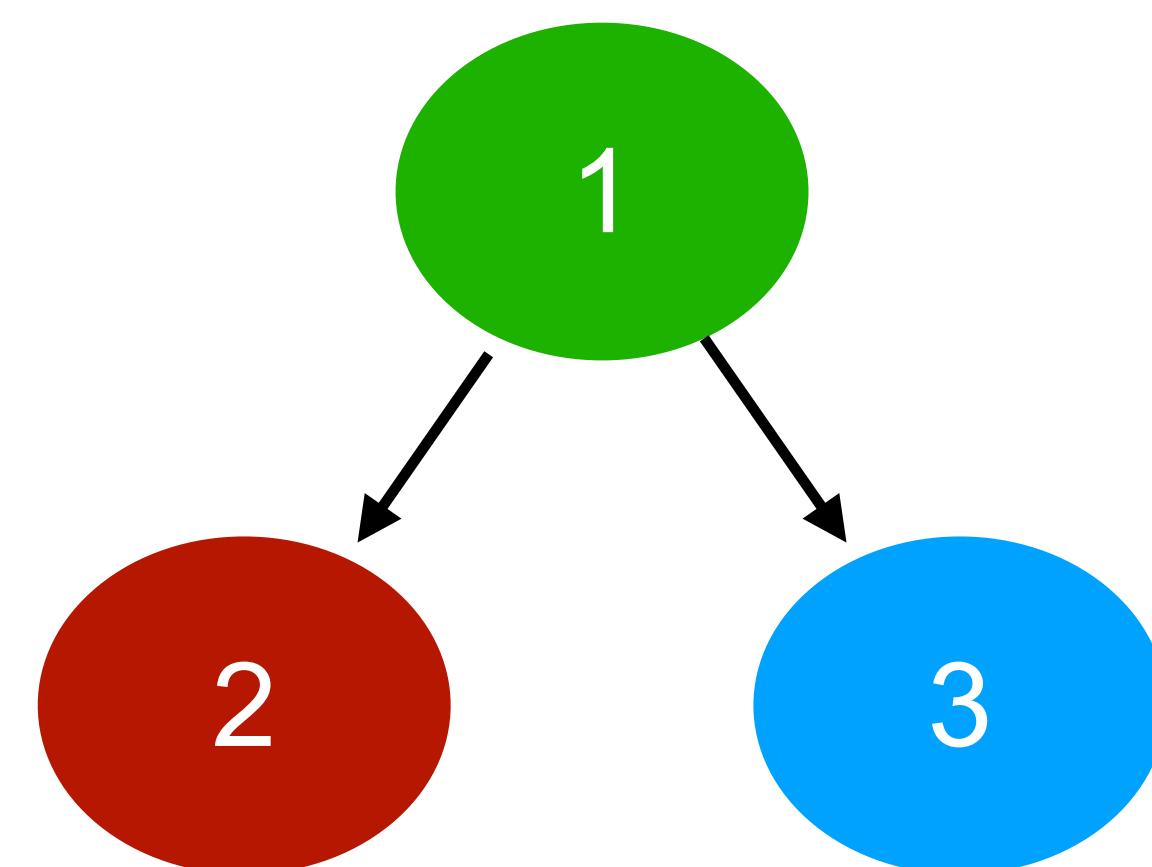
- Similarly for children, ancestors and descendants

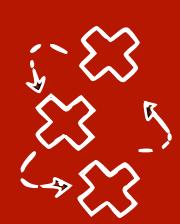


# DAGs and random variables

- We can represent a factorisation of a joint distribution  $p(X_1, \dots, X_p)$  as a **DAG**
- **Each node  $i \in V$**  represents a **random variable  $X_i$** 
  - For  $A \subseteq V$ , we can define  $X_A := \{X_i : i \in A\}$
- **Edges** are related to the factors

$$p(X_1, X_2, X_3)$$

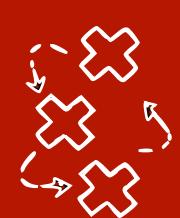




# Factorizing joint distributions

- A joint distribution can always be factorized in several ways by iterating the **chain rule**

$$p(X_1, X_2) = p(X_1)p(X_2 | X_1) = p(X_2)p(X_1 | X_2)$$



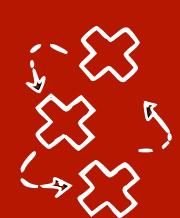
# Factorizing joint distributions

- A joint distribution can always be factorized in several ways by iterating the **chain rule**

$$p(X_1, X_2) = p(X_1)p(X_2 | X_1) = p(X_2)p(X_1 | X_2)$$

- In general, given any **ordering** of the variables  $(X_1, \dots, X_p)$ , we can write:

$$p(X_1, \dots, X_p) = p(X_1)p(X_2 | X_1) \dots p(X_p | X_1, \dots, X_{p-1})$$



# Factorizing joint distributions

- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

$$p(X_1, \dots, X_p) = p(X_1)p(X_2 | X_1)\dots p(X_p | X_1, \dots, X_{p-1})$$

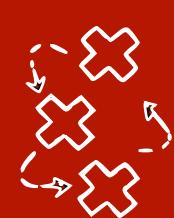
- For example  $p(X, Y, Z)$  can be equivalently factorized as:

- $p(X, Y, Z) =$

- $p(X, Z, Y) =$

- $p(Z, Y, X) =$

- $\dots$



# Exploiting conditional independences

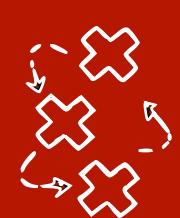
- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

$$p(X_1, X_2) = p(X_1)p(X_2 | X_1) = p(X_2)p(X_1 | X_2)$$

- We can **simplify** the factorisation by using **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies p(X_i | X_j, X_{\mathbf{Z}}) = p(X_i | X_{\mathbf{Z}})$$

(special case  $X_i \perp\!\!\!\perp X_j \implies p(X_i | X_j) = p(X_i)$ )



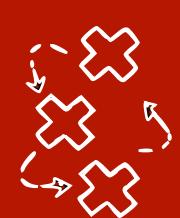
## Quick recap: Independent random variables (discrete case)

- **Definition:** Two discrete random variables  $X$  and  $Y$  are **independent** iff:

$$\forall x, y : P(X = x, Y = y) = P(X = x)P(Y = y)$$

- We then write  $X \perp\!\!\!\perp Y$  (equiv.  $Y \perp\!\!\!\perp X$ ), otherwise  $X \not\perp\!\!\!\perp Y$
- This is equivalent to  $P(X = x | Y = y) = P(X = x)$  (and vice versa for  $Y$ )
- Intuitively, this means that knowing the value of  $Y$  **will not tell us anything** about the distribution of  $X$ .





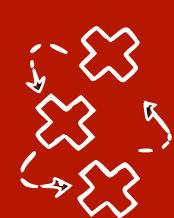
## Quick recap: conditional independence (discrete case)

- $X$  is independent of  $Y$  **conditioned/given**  $Z$  (possibly a set) iff

$$\forall x, y, z : P(X = x, Y = y | Z = z) = P(X = x | Z = z)P(Y = y | Z = z)$$

(for  $P(Z = z) > 0$ )

- We then write  $X \perp\!\!\!\perp Y | Z$  (equiv.  $Y \perp\!\!\!\perp X | Z$ ), otherwise  $X \perp\!\!\!\perp Y | Z$
- This is equivalent to  $P(X = x | Y = y, Z = z) = P(X = x | Z = z)$
- Intuitively this means that **Y does not add any information** to predict  $X$  that isn't already offered by  $Z$



# Exploiting conditional independences

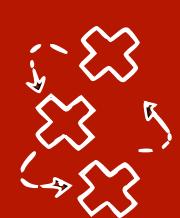
- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

$$p(X_1, X_2) = p(X_1)p(X_2 | X_1) = p(X_2)p(X_1 | X_2)$$

- We can **simplify** the factorisation by using **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies p(X_i | X_j, X_{\mathbf{Z}}) = p(X_i | X_{\mathbf{Z}})$$

(special case  $X_i \perp\!\!\!\perp X_j \implies p(X_i | X_j) = p(X_i)$ )



# Factorizing joint distributions

- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

$$p(X_1, \dots, X_p) = p(X_1)p(X_2 | X_1)\dots p(X_p | X_1, \dots, X_{p-1})$$

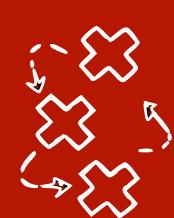
- For example if  $X \perp\!\!\!\perp Y | Z$ :

- $p(X, Y, Z) =$

- $p(X, Z, Y) =$

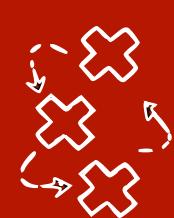
- $p(Z, Y, X) =$

- ...



# Bayesian networks

- We have a set of random variables  $X_1, \dots, X_p$  with joint  $p(X_1, \dots, X_p)$
- We have a DAG  $G$ , s.t. **each random variable  $X_i$  is represented by node  $i$**
- We then say  $p(X_1, \dots, X_p)$  **factorizes over  $G$**  if
$$p(X_1, \dots, X_p) = \prod_{i \in V} p(X_i | \mathbf{X}_{\text{pa}(i)})$$
- A **Bayesian network** (BN) is the tuple  $(G, p)$  s.t.  $p$  **factorizes over  $G$**

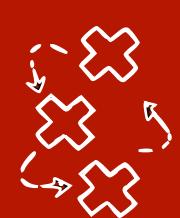


# Multiple BN can represent the same distribution

- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

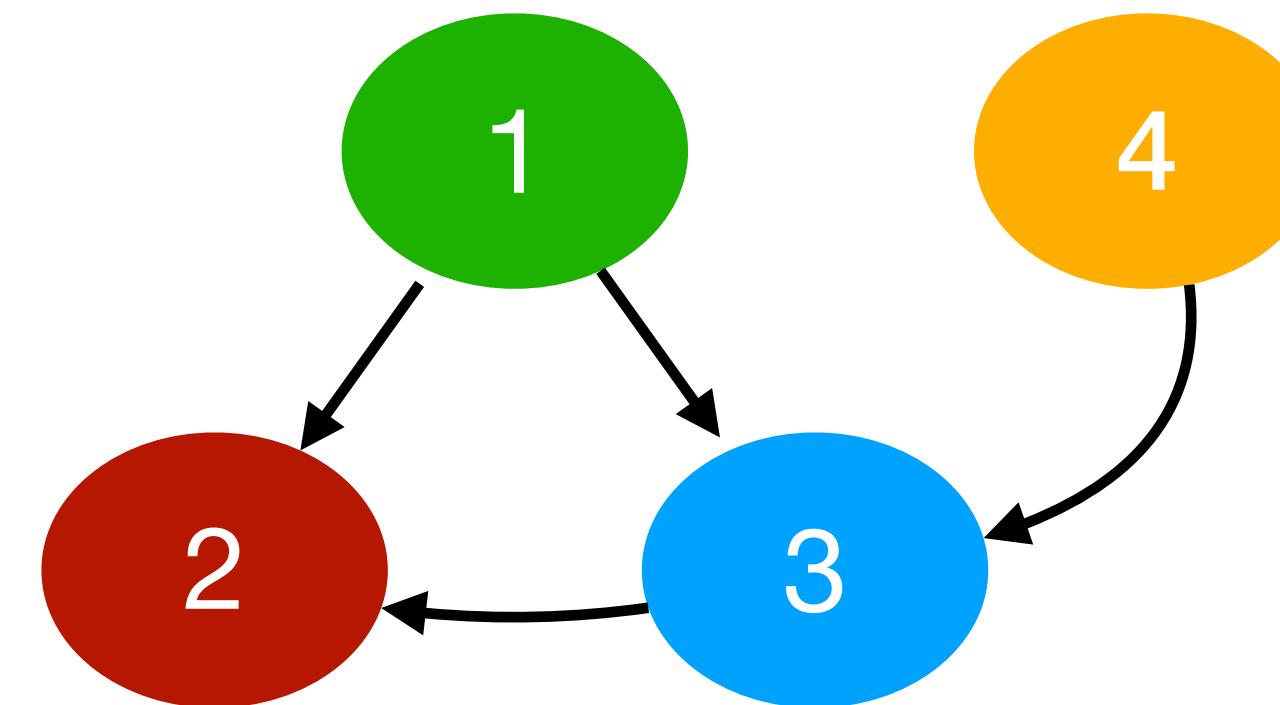
$$p(X_1, \dots, X_p) = p(X_1)p(X_2 | X_1)\dots p(X_p | X_1, \dots, X_{p-1})$$

- For example if  $X \perp\!\!\!\perp Y | Z$ : [Each factorisation can be represented with a DAG]
  - $p(X, Y, Z) = p(X)p(Y)p(Z | X, Y)$
  - $p(X, Z, Y) = p(X)p(Z | X)p(Y | Z)$
  - $p(Z, Y, X) = p(Z)p(Y | Z)p(X | Z)$
  - ...

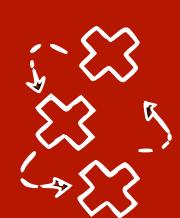


# Graph terminology: collider on a path

- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**



- A **collider**  $k$  on a **path**  $\pi = (i, \dots, j)$  is a non-endpoint node ( $k \neq i, j$ ) s.t.  $\pi$  contains  $\rightarrow k \leftarrow$  (other non-endpoint nodes are **non-colliders**)

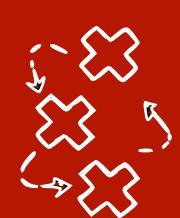


# Graph terminology: collider on a path

- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**

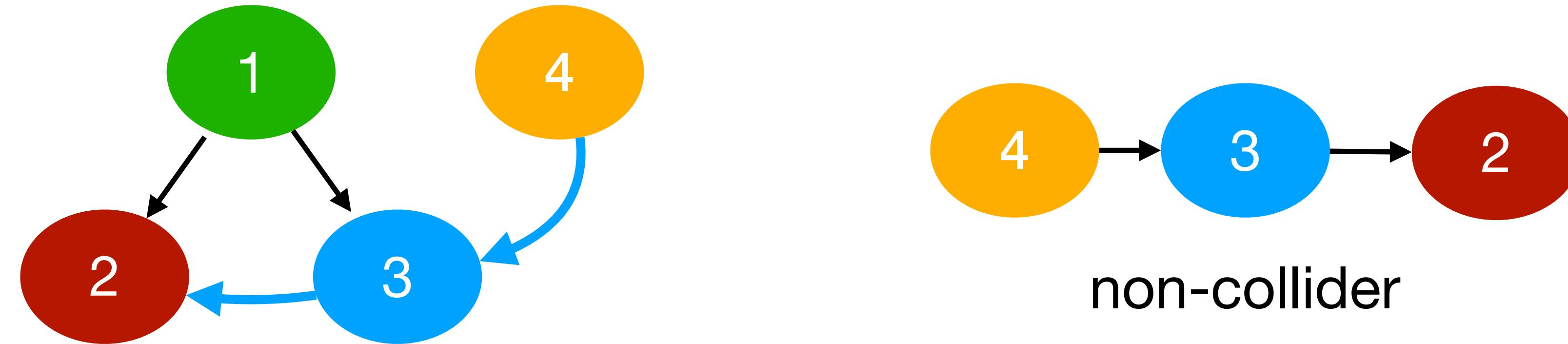


- A **collider**  $k$  on a **path**  $\pi = (i, \dots, j)$  is a non-endpoint node ( $k \neq i, j$ ) s.t.  $\pi$  contains  $\rightarrow k \leftarrow$  (other non-endpoint nodes are **non-colliders**)

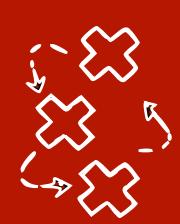


# Graph terminology: collider on a path

- A **path** between **node i and node j** is a sequence of **distinct nodes**  $(i, \dots, j)$  such that each two **consecutive nodes** are **adjacent**

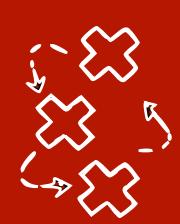


- A **collider**  $k$  on a **path**  $\pi = (i, \dots, j)$  is a non-endpoint node ( $k \neq i, j$ ) s.t.  $\pi$  contains  $\rightarrow k \leftarrow$  (other non-endpoint nodes are **non-colliders**)



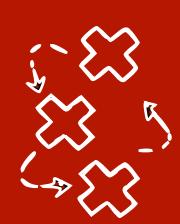
# d-separation: blocked paths

- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  at least one holds:
  - There is a **non-collider** on the path that is in  $A$ , or
  - There is a **collider  $k$**  on the path, but  $\text{Desc}(k) \cap A = \emptyset$



# d-separation: blocked paths

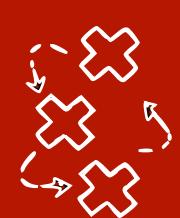
- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  at least one holds:
  - There is a non-collider on the path that is in  $A$ , or
  - There is a collider  $k$  on the path, but  $\text{Desc}(k) \cap A = \emptyset$
- Otherwise it is **active**



# d-separation: blocked paths

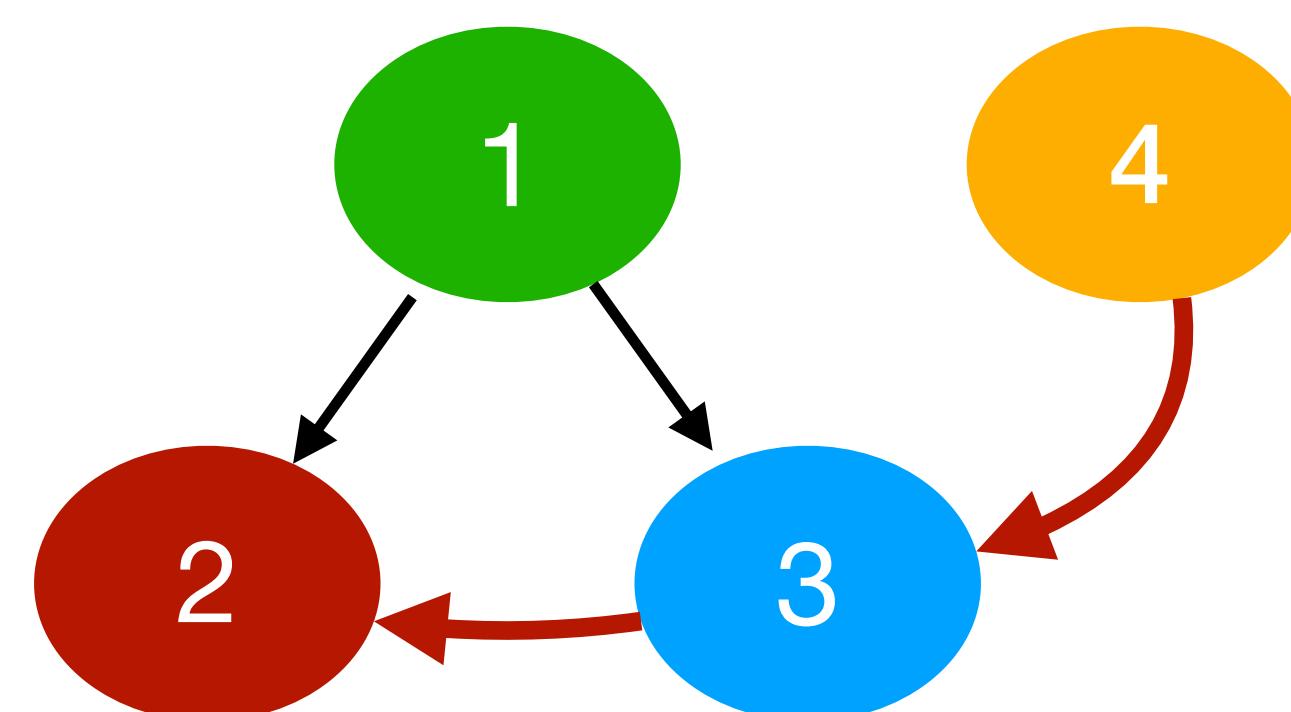
- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  at least one holds:
  - There is a non-collider on the path that is in  $A$ , or
  - There is a collider  $k$  on the path, but  $\text{Desc}(k) \cap A = \emptyset$
- Otherwise it is **active**

**Note:** descendants w.r.t. the **whole graph**

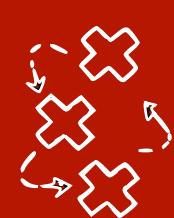


# d-separation: blocked paths - example 1

- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  at least one holds:
  - There is a non-collider on the path that is in  $A$ , or
  - There is a collider  $k$  on the path, but  $\text{Desc}(k) \cap A = \emptyset$
- Otherwise it is **active**

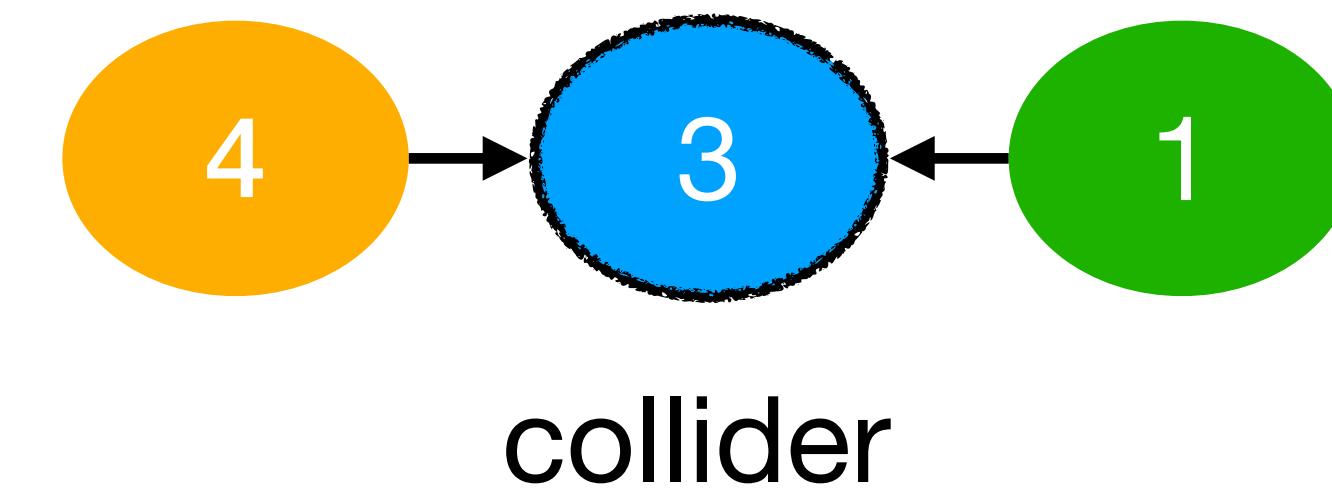
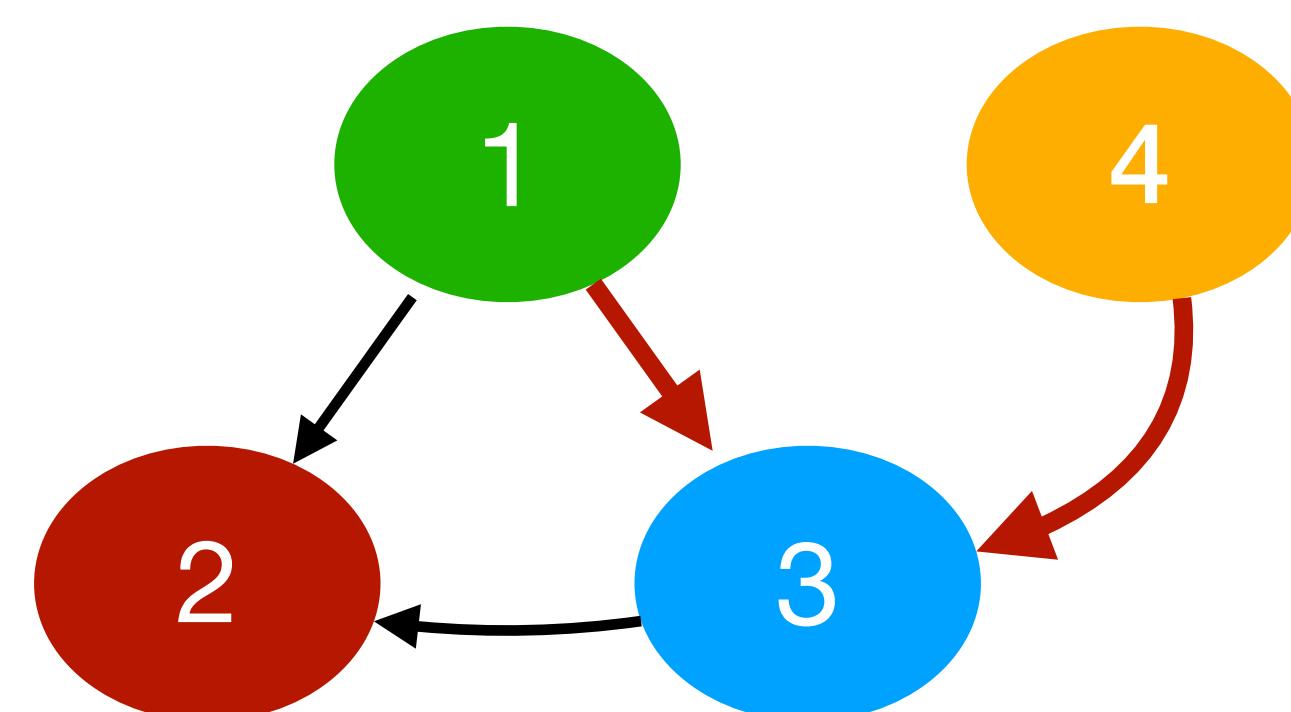


If  $3 \in A$ , the path is **blocked**,  
otherwise it is **active**

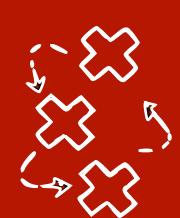


# d-separation: blocked paths - example 2

- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  at least one holds:
  - There is a non-collider on the path that is in  $A$ , or
  - There is a collider  $k$  on the path, but  $\text{Desc}(k) \cap A = \emptyset$
- Otherwise it is **active**

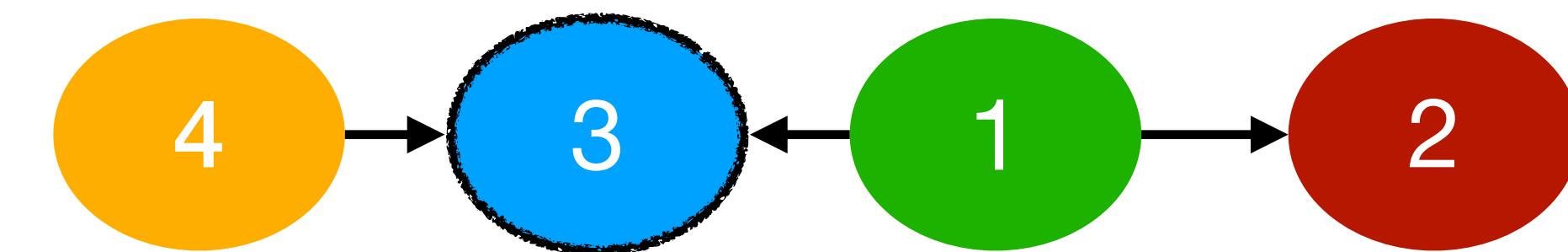
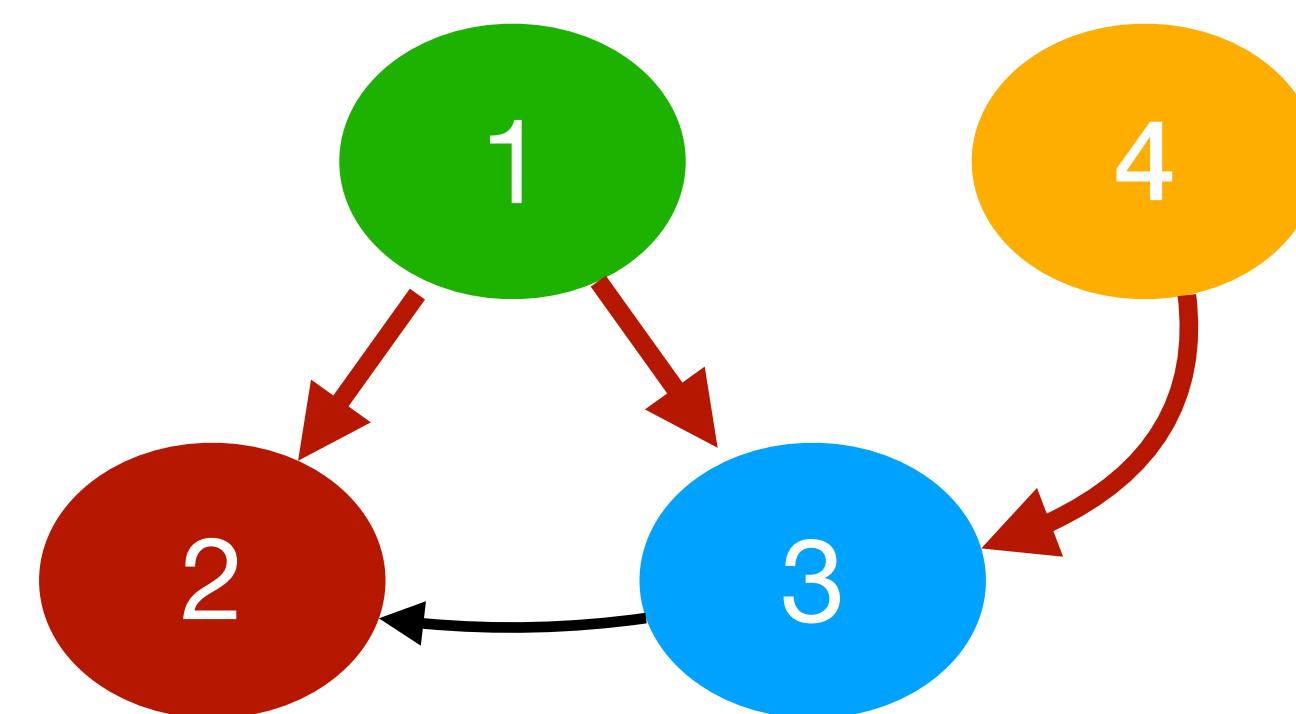


If  $3 \notin A$  and  $2 \notin A$ , the path is **blocked**,  
otherwise it is **active**



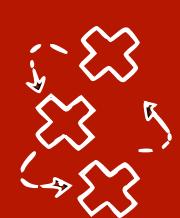
# d-separation: blocked paths - example 3

- A **path** between **i and j** is **blocked by  $A \subseteq V \setminus \{i, j\}$**  if at least one of the following holds:
  - There is a non-collider on the path that is in  $A$ , or
  - There is a collider  $k$  on the path, but  $\text{Desc}(k) \cap A = \emptyset$
- Otherwise it is **active**



If  $1 \in A$ , the path is **blocked**  
OR

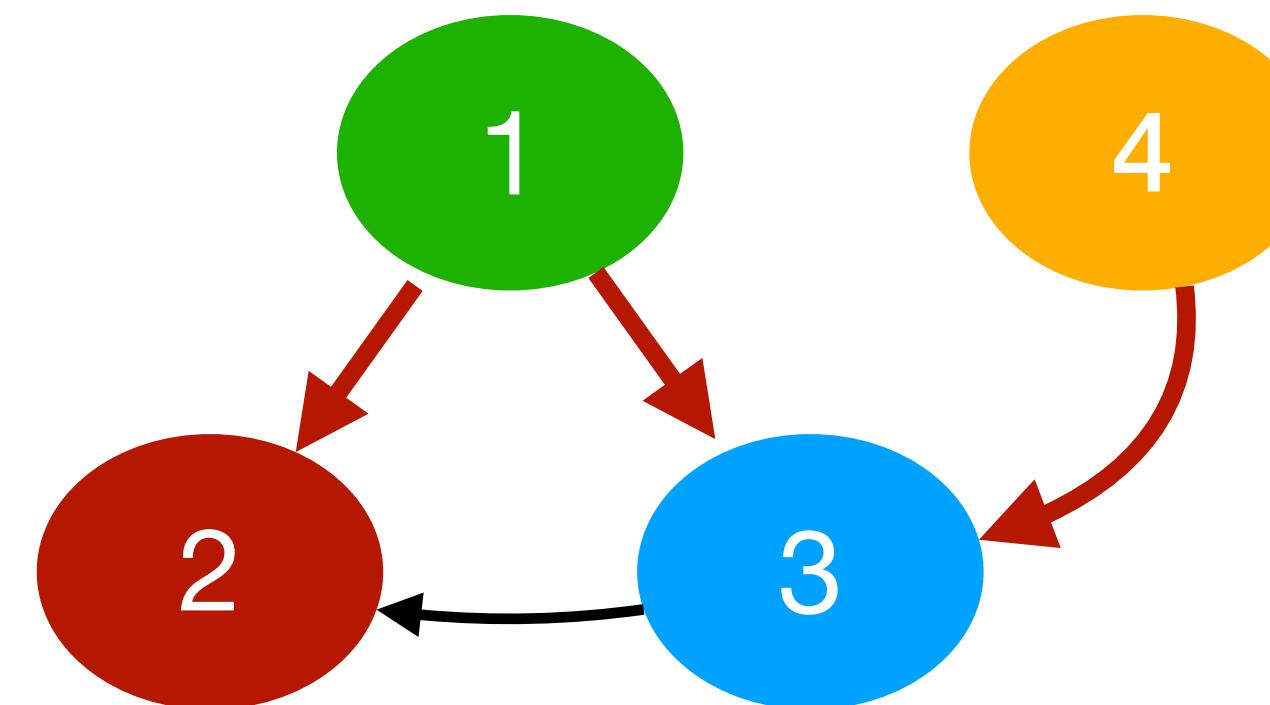
If  $3 \notin A$  and  $2 \notin A$ , the path is **blocked**

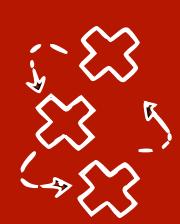


# d-separation

- **i and j are d-separated by  $A \subseteq V \setminus \{i, j\}$**  if all paths between  $i, j$  are **blocked by A**
  - We denote d-separation as  $i \perp j | A$
- Otherwise we say they are **d-connected**  $i \not\perp j | A$

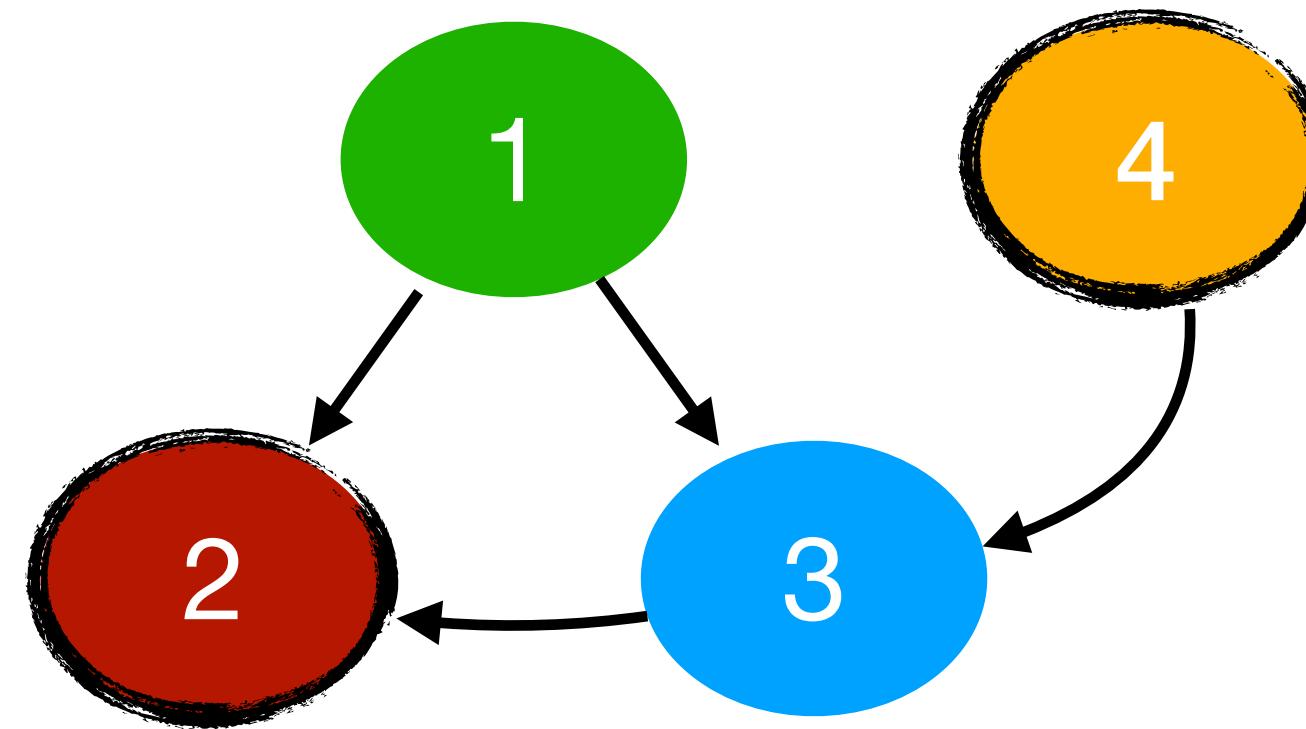
**Note:** d-separation is symmetric

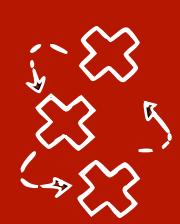




# d-separation - complete example

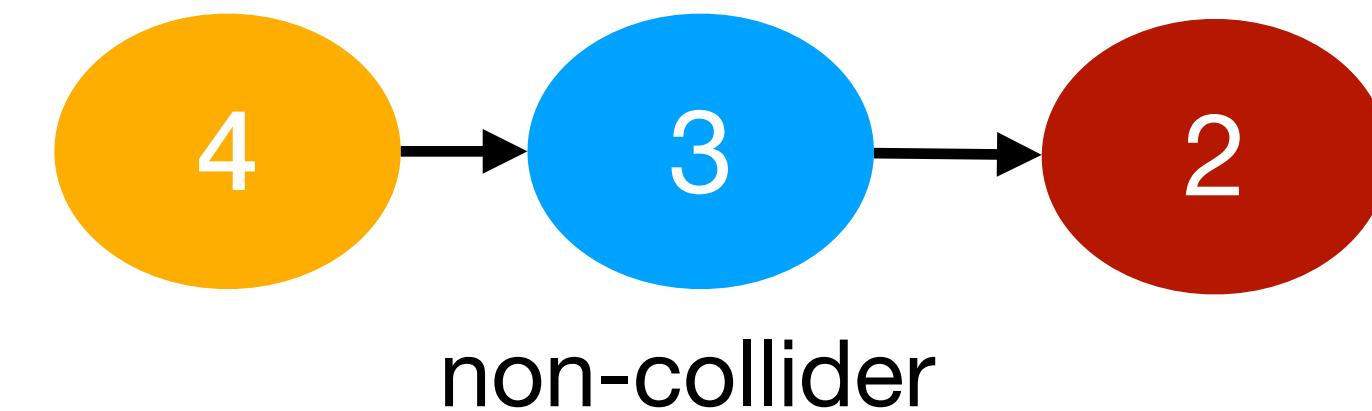
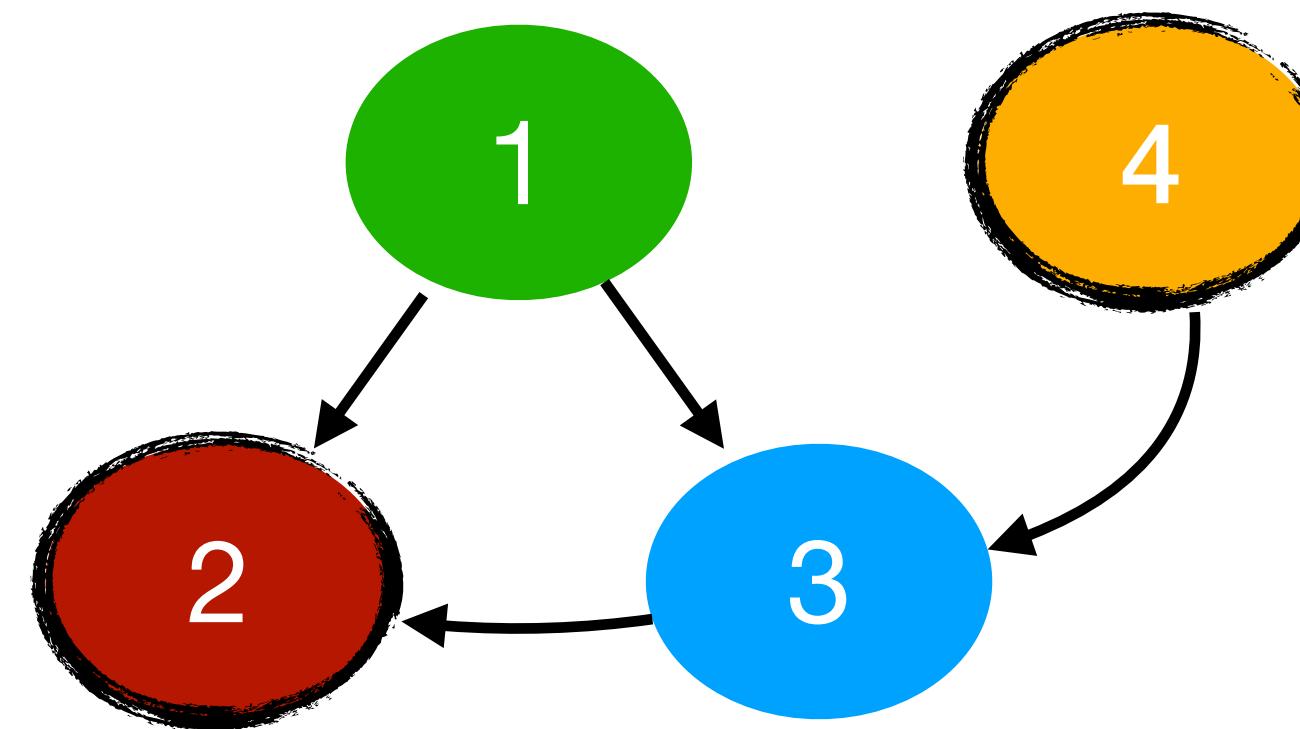
- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?



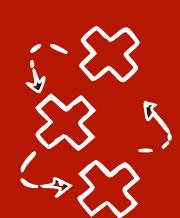


# d-separation - complete example

- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?

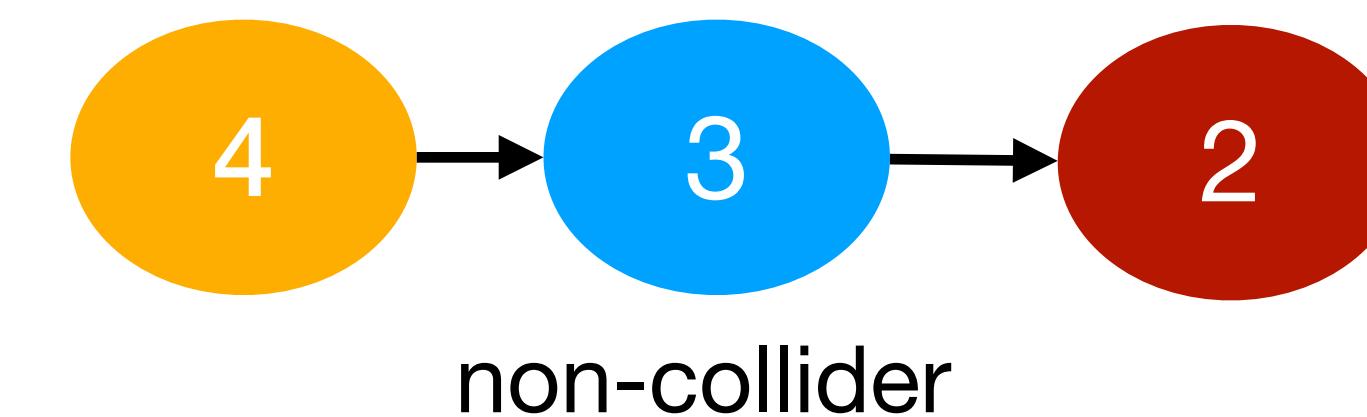
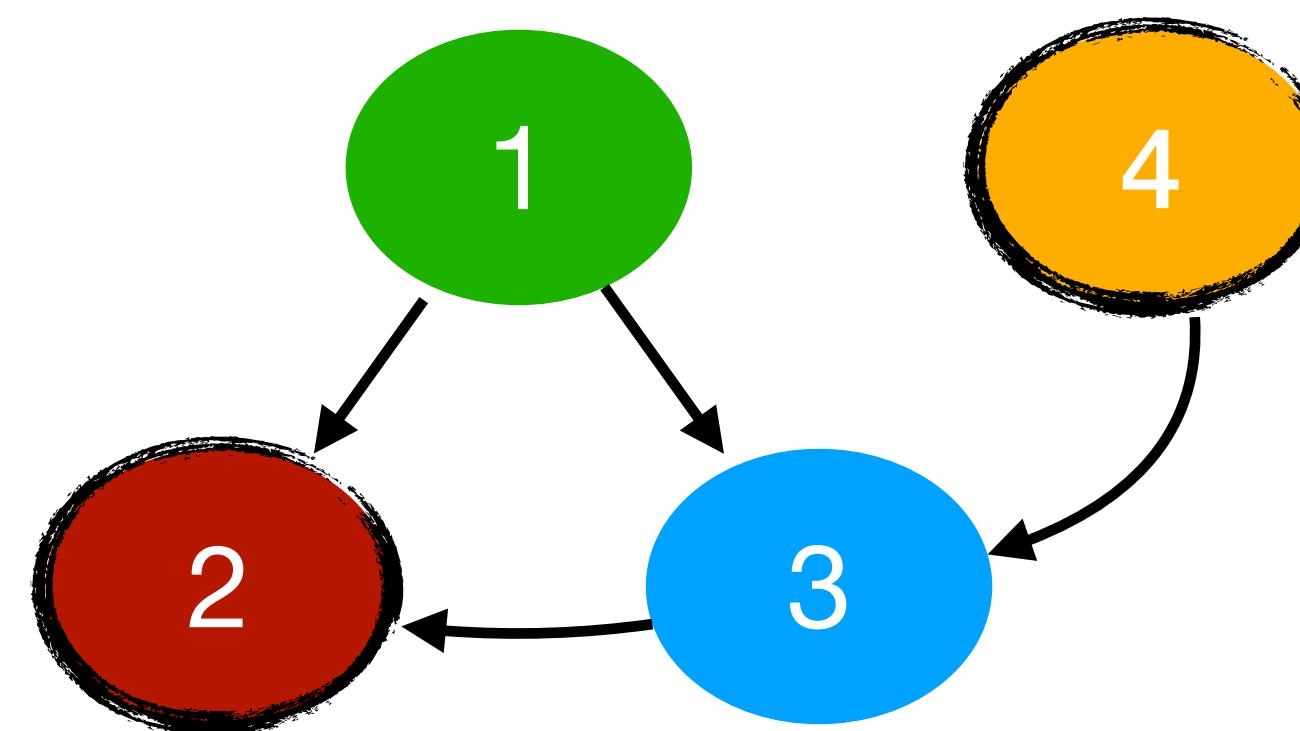


If  $3 \in A$ , the path is **blocked**

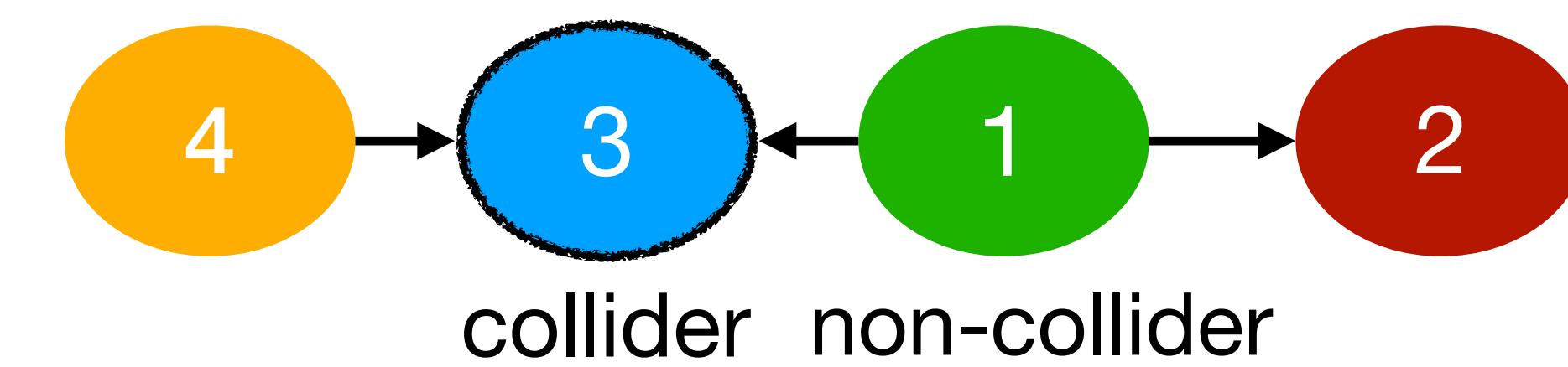


# d-separation - complete example

- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?

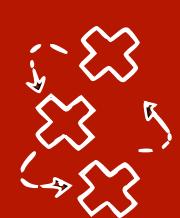


If  $3 \in A$ , the path is **blocked**



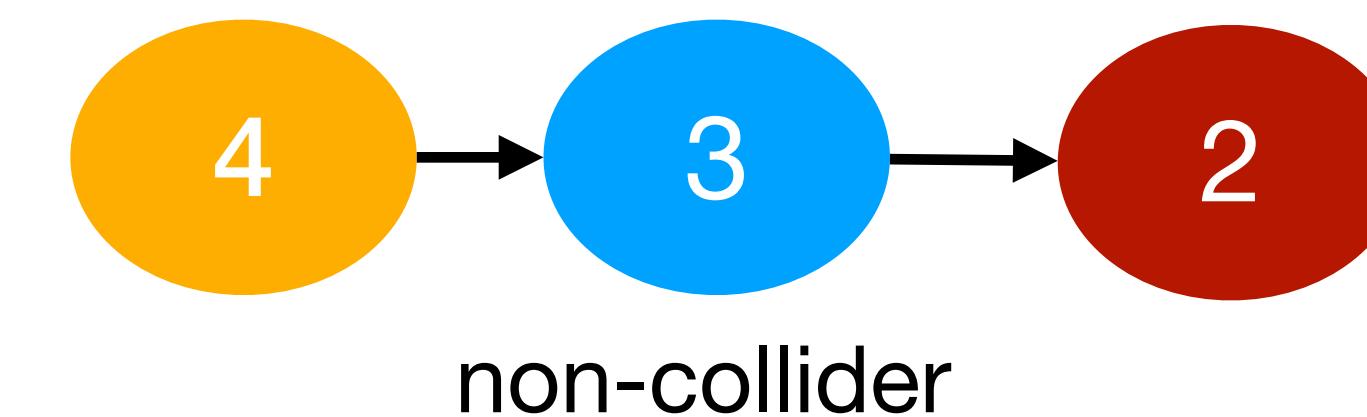
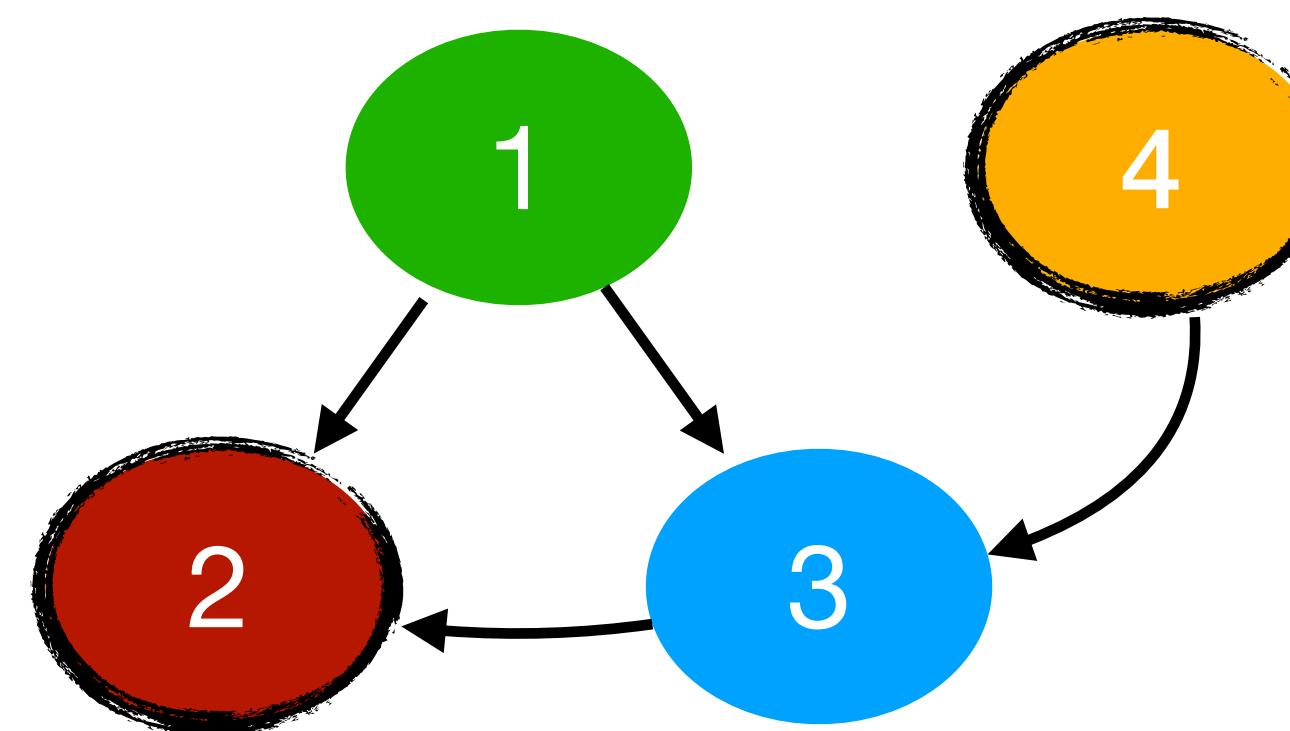
If  $1 \in A$ , the path is **blocked**, OR

If  $3 \notin A$ , the path is **blocked**

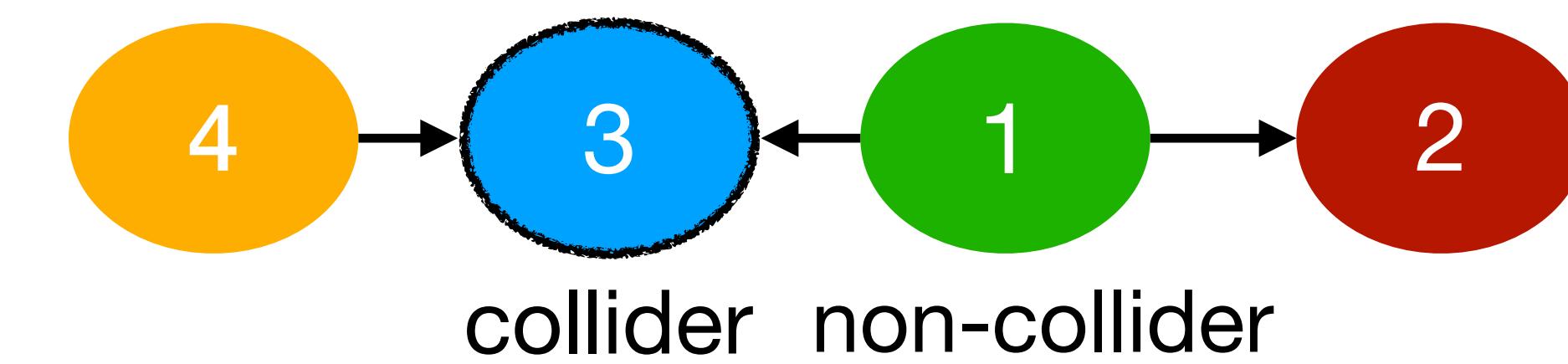


# d-separation - complete example

- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?

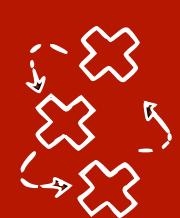


If  $3 \in A$ , the path is **blocked**



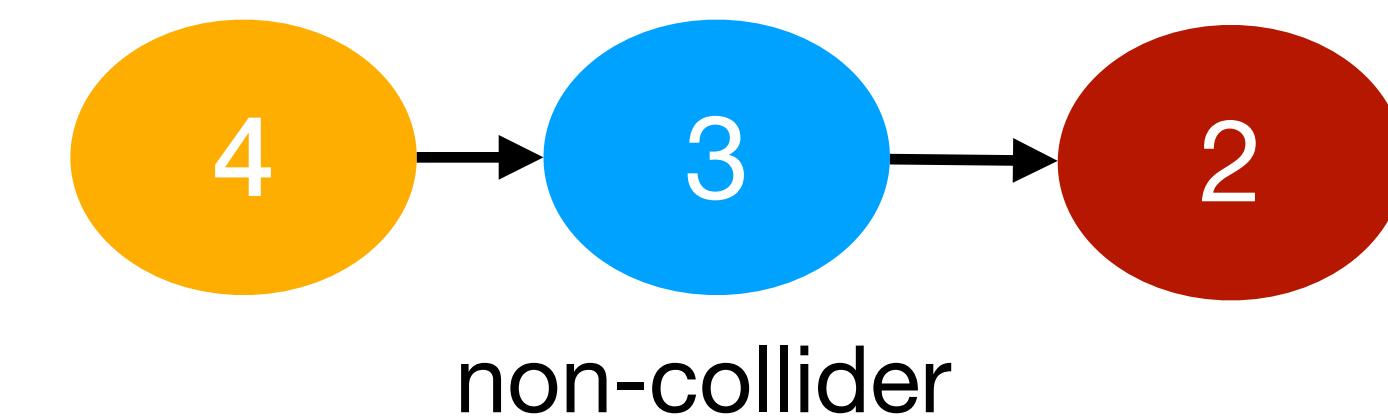
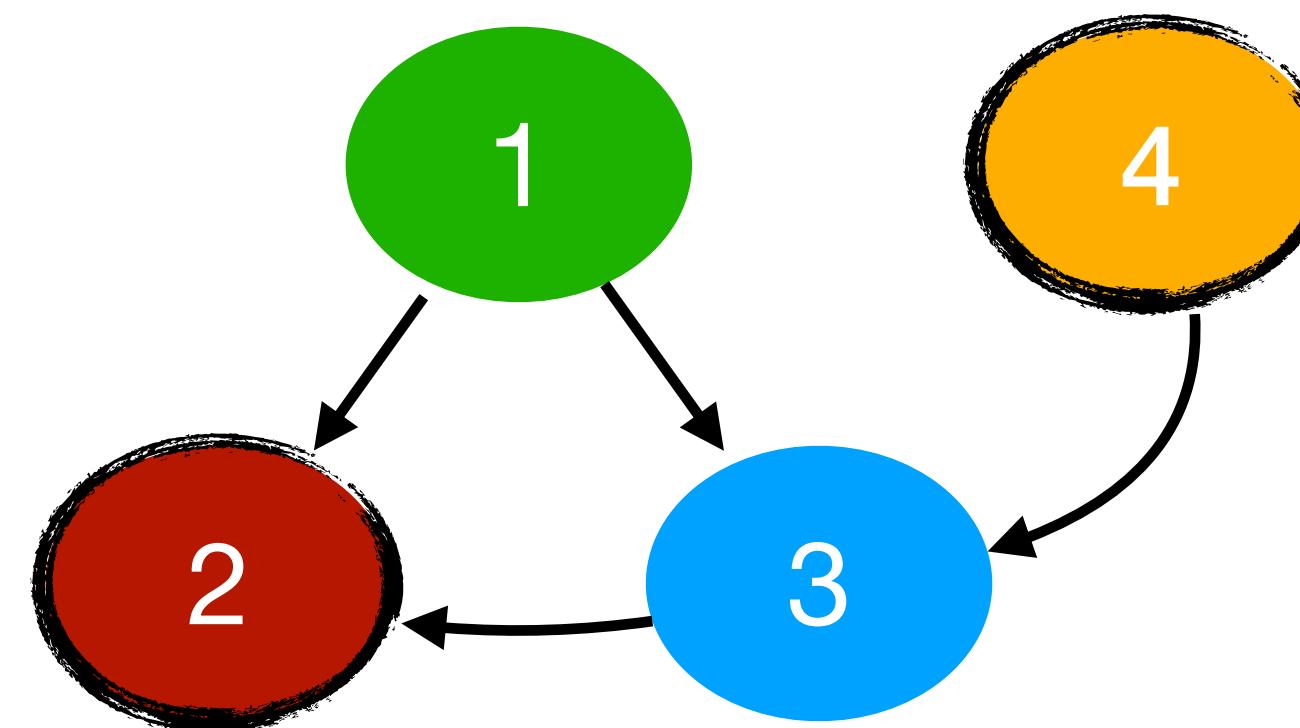
If  $1 \in A$ , the path is **blocked**, OR

~~If  $3 \notin A$ , the path is **blocked**~~



# d-separation - complete example

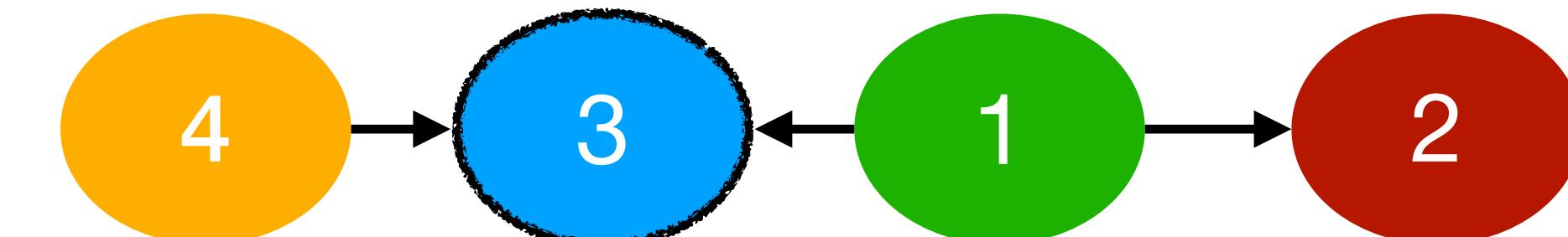
- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?



If  $3 \in A$ , the path is **blocked**



$$A = \{1, 3\}$$

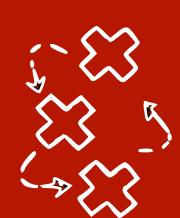


collider non-collider

If  $1 \in A$ , the path is **blocked**, OR

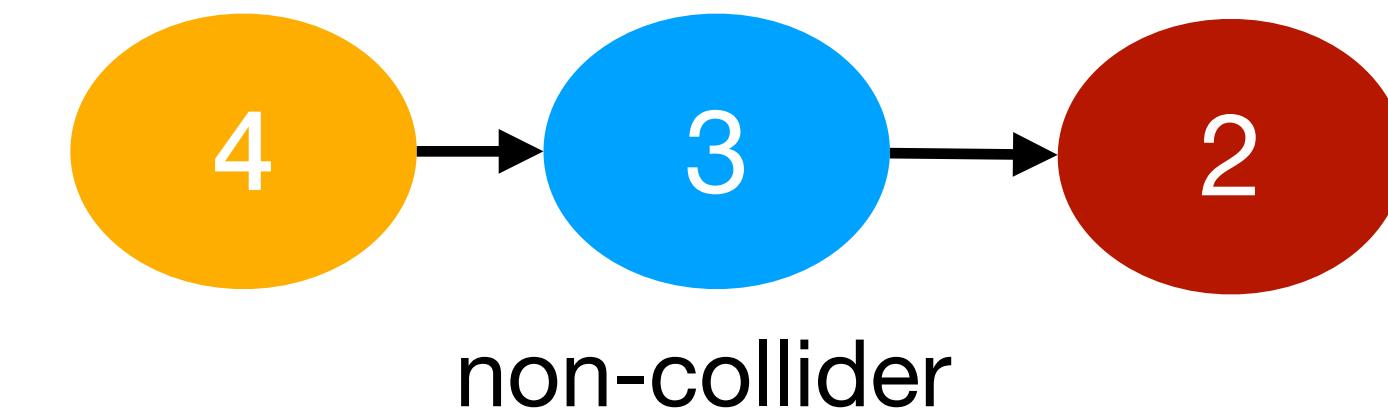
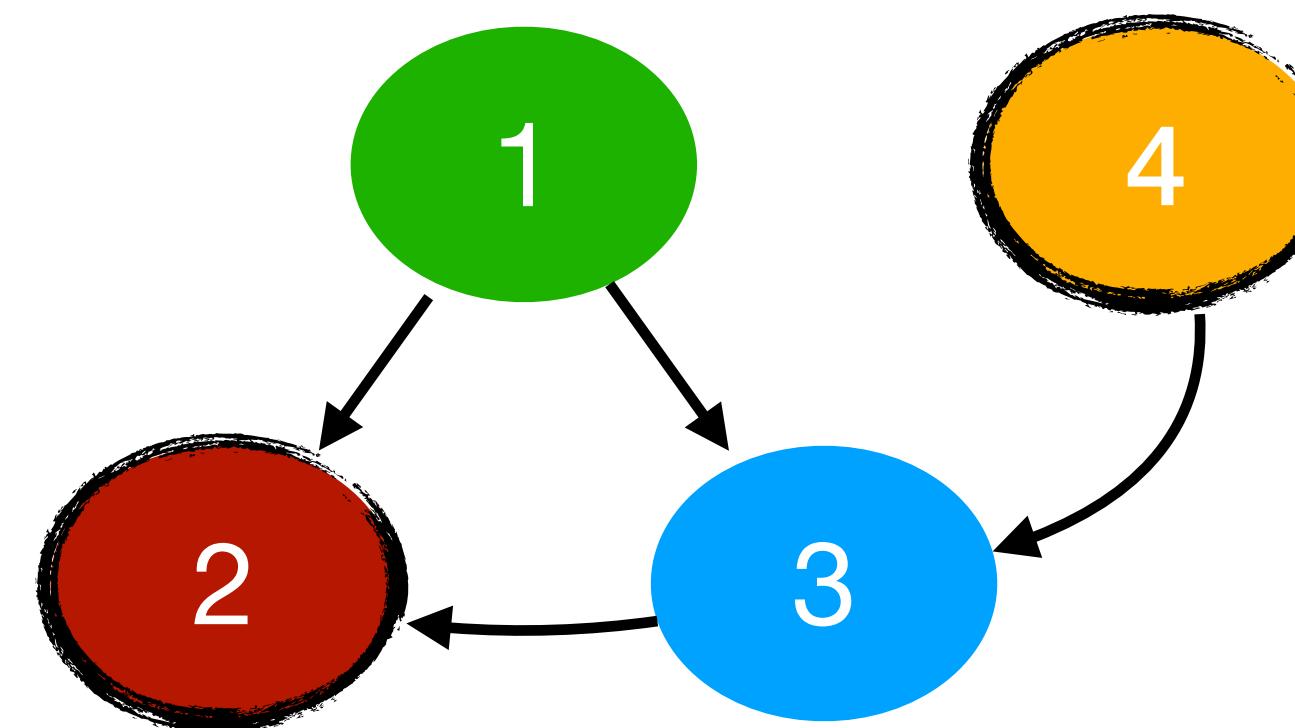


~~If  $3 \notin A$ , the path is **blocked**~~



# d-separation - complete example

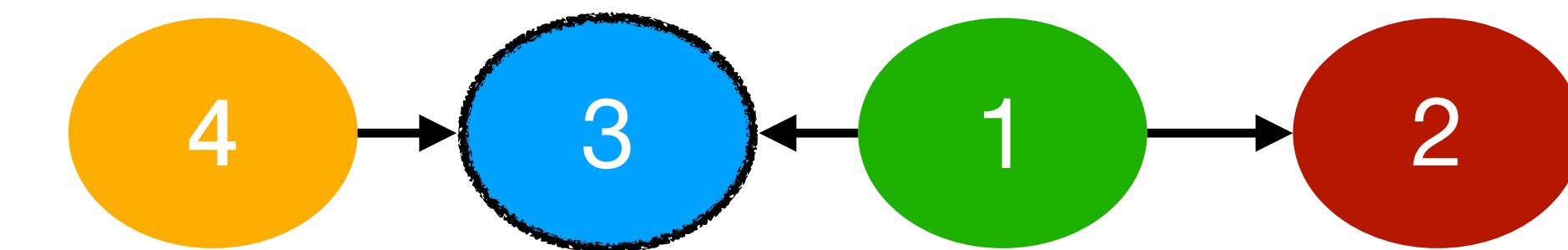
- Which  $A \subseteq V \setminus \{2,4\}$  d-separate 4 and 2?



If  $3 \in A$ , the path is **blocked**



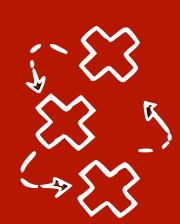
$$A = \{1, 3\}$$



$$2 \perp 4 \mid \{1, 3\}$$

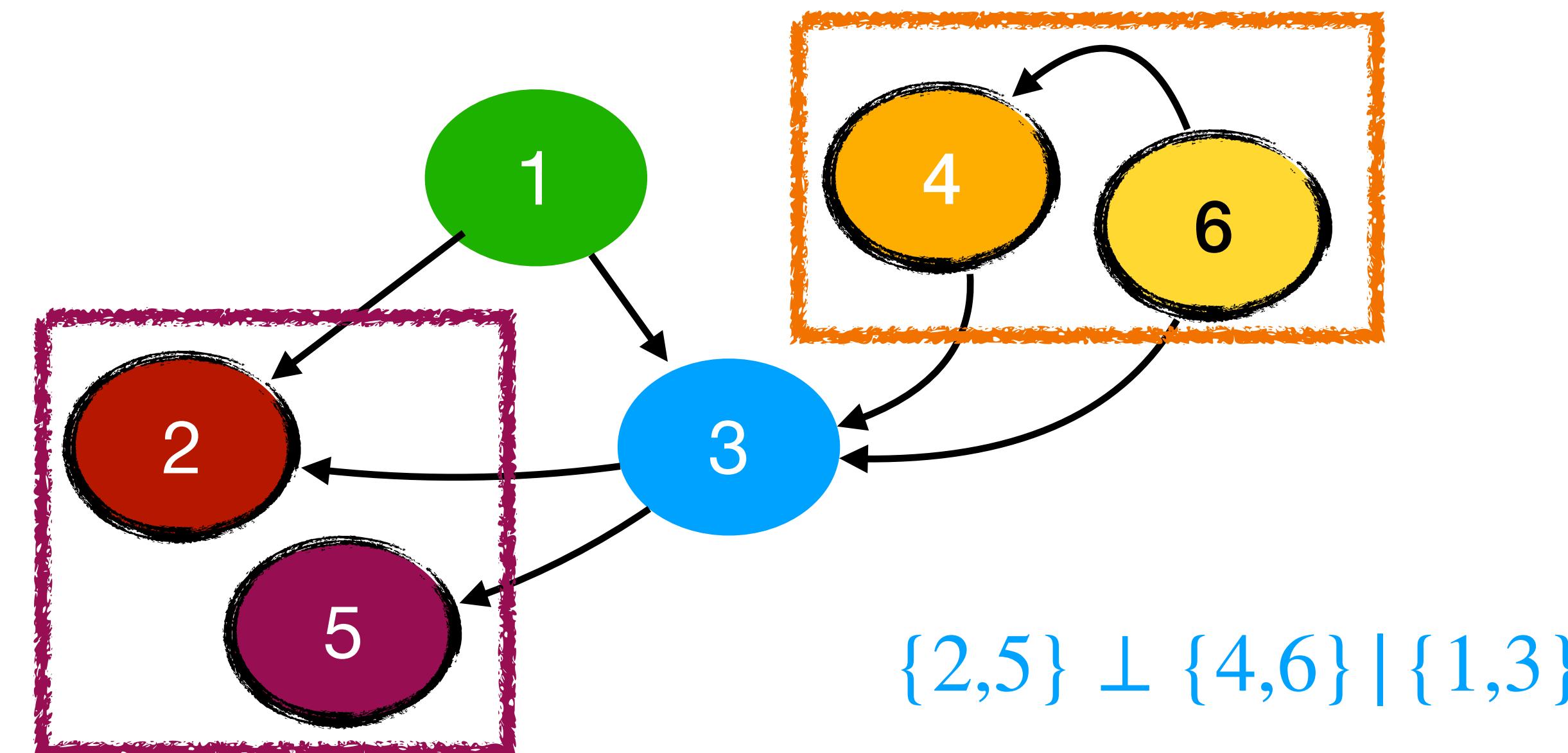
If  $1 \in A$ , the path is **blocked**

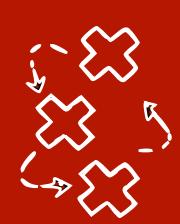




# d-separation for sets of nodes

- Sets of nodes  $\mathbf{I}$  and  $\mathbf{J}$  are **d-separated by  $A \subseteq V \setminus I \cup J$** , if all paths between  $i \in I$  and  $j \in J$  are **blocked by  $A$**
- We denote d-separation as  $I \perp J | A$

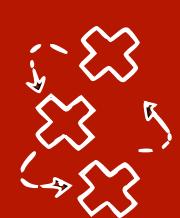




# Global Markov Property and faithfulness

- If  $(G, p)$  is a Bayesian network with a DAG  $G = (\mathbf{V}, \mathbf{E})$ , i.e.  $p$  **factorizes according to  $\mathbf{G}$** , then for any disjoint  $A, B, C \subseteq \mathbf{V}$ :

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

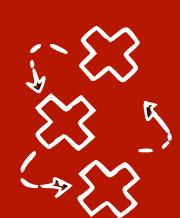


# Global Markov Property and faithfulness

- If  $(G, p)$  is a Bayesian network with a DAG  $G = (\mathbf{V}, \mathbf{E})$ , i.e.  $p$  **factorizes according to  $\mathbf{G}$** , then for any disjoint  $A, B, C \subseteq \mathbf{V}$ :

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

- **d-separations** that can be read purely from a graph imply **conditional independences** in the random variables and data generated by the graph



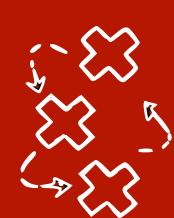
# Global Markov Property and faithfulness

- If  $(G, p)$  is a Bayesian network with a DAG  $G = (\mathbf{V}, \mathbf{E})$ , i.e.  $p$  **factorizes according to  $\mathbf{G}$** , then for any disjoint  $A, B, C \subseteq \mathbf{V}$ :

$$A \perp_G B | C \implies X_A \perp_p X_B | X_C$$

- **d-separations** that can be read purely from a graph imply **conditional independences** in the random variables and data generated by the graph
- The reverse is not true in general, but if it is we say  $p$  **is faithful to  $\mathbf{G}$**  and:

$$A \perp_G B | C \iff X_A \perp_p X_B | X_C$$



# Why should we care about Bayesian networks?

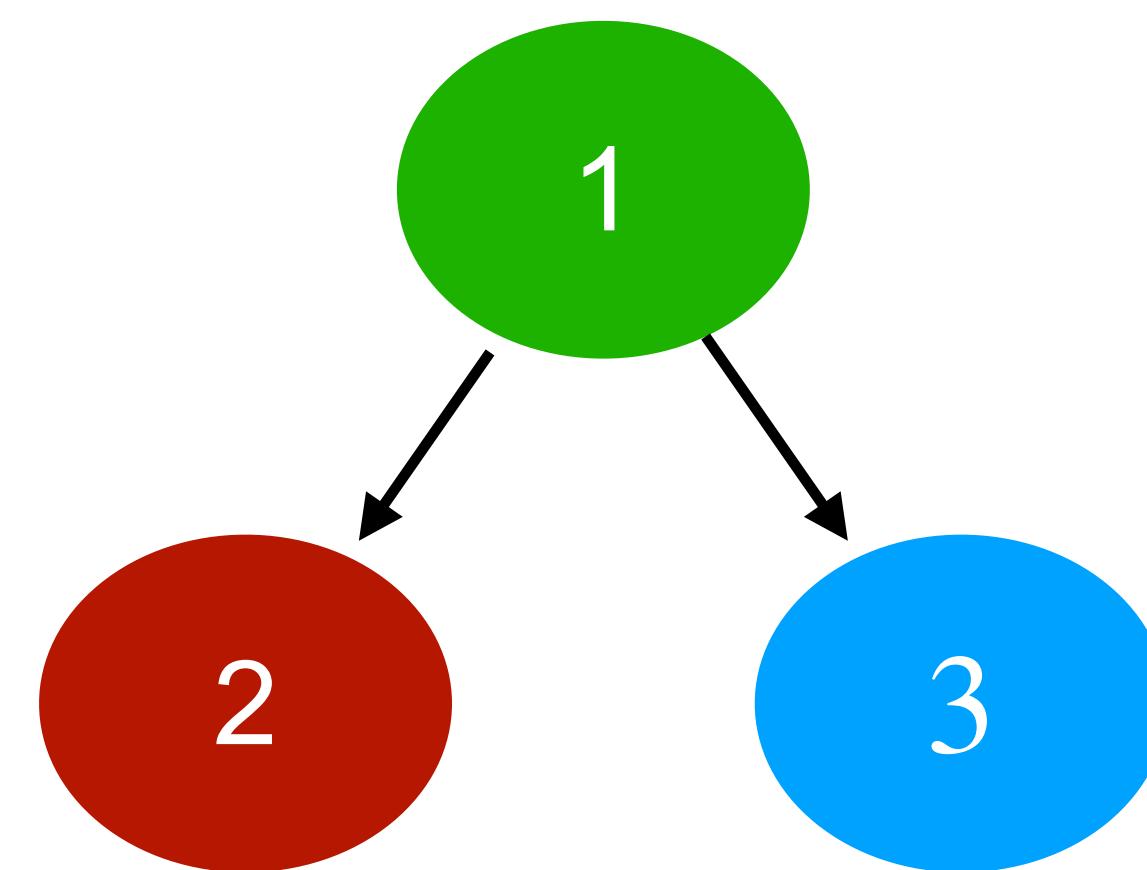
- We have a set of random variables  $X_1, \dots, X_p$  with joint  $p(X_1, \dots, X_p)$
- We have a DAG  $G$ , s.t. **each random variable  $X_i$  is represented by node  $i$**
- We then say  $p(X_1, \dots, X_p)$  **factorizes over  $G$**  if

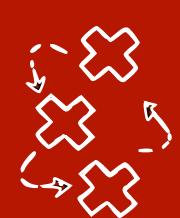
$$p(X_1, \dots, X_p) = \prod_{i \in V} p(X_i | \mathbf{X}_{\text{pa}(i)})$$

They can help simplify the factorisation

We can easily read conditional independences (d-separation)

They can represent causal models

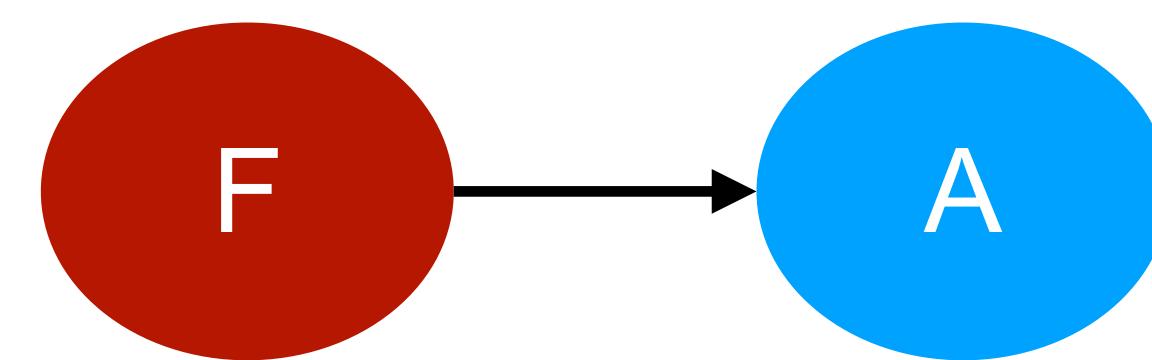




# BNs vs causal BNs - example

- Fire (F) and Alarm (A) with  $p(F, A)$  and  $A \not\perp\!\!\!\perp F$  can be factorized as:

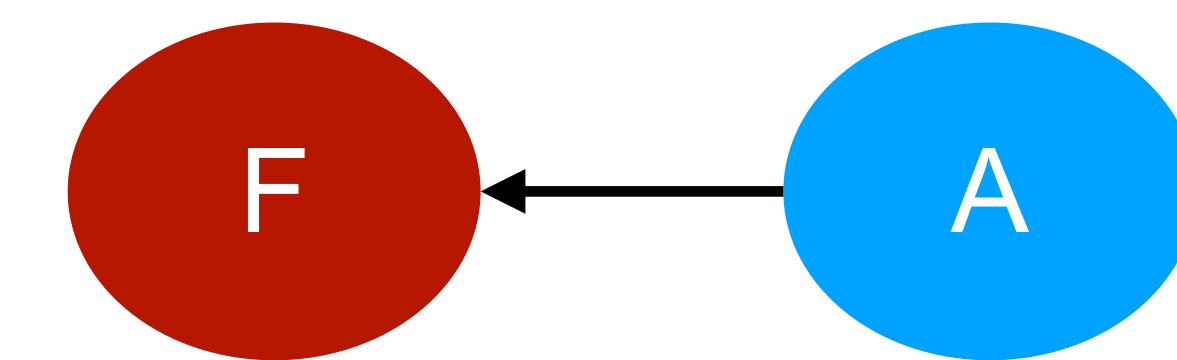
$$p(F, A) = p(F) p(A|F)$$



CAUSAL

(lighting a fire triggers alarm)

$$p(F, A) = p(A) p(F|A)$$



NOT-CAUSAL

(triggering alarm does not light a fire)