

# Causal Discovery

# Agenda

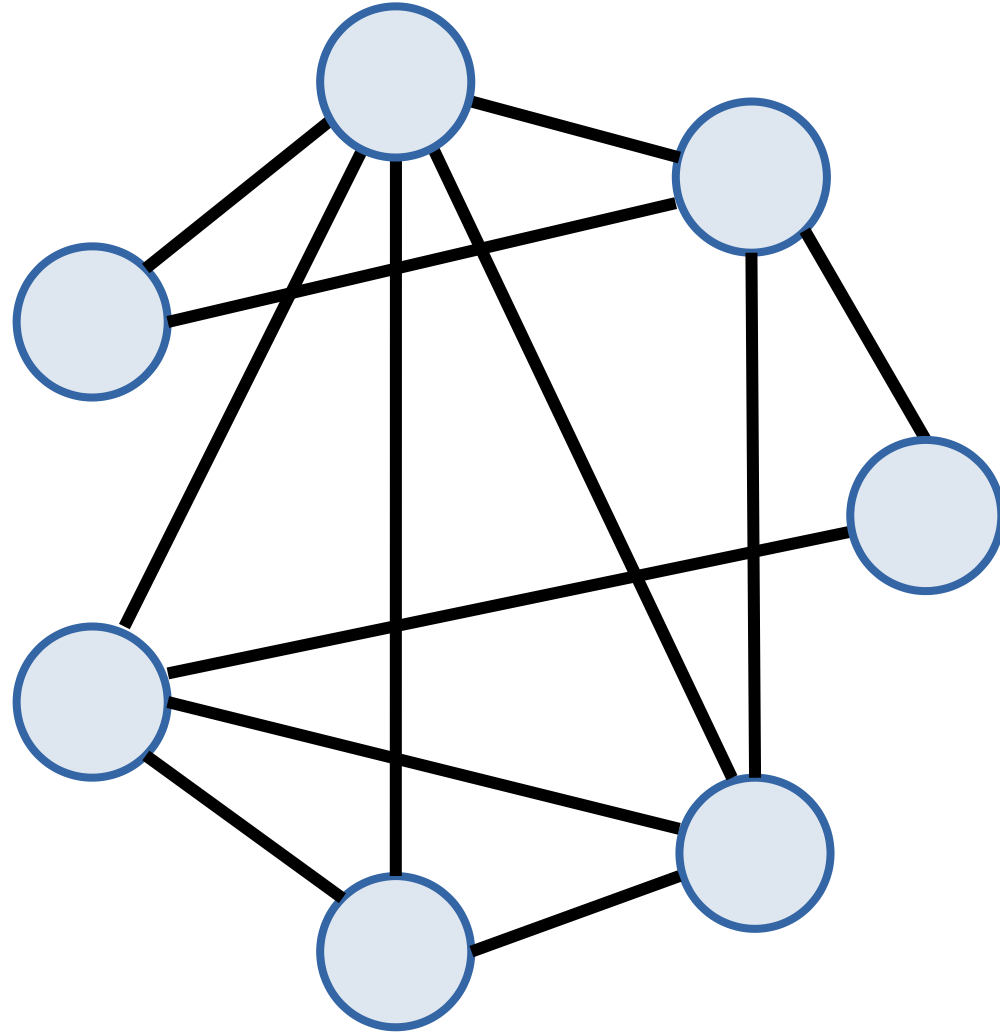
- Discovery vs inference
- Causal discovery algorithms
  - PC
  - GES
  - LiNGAM
  - NOTEARS
- Code and examples

**Discovery**

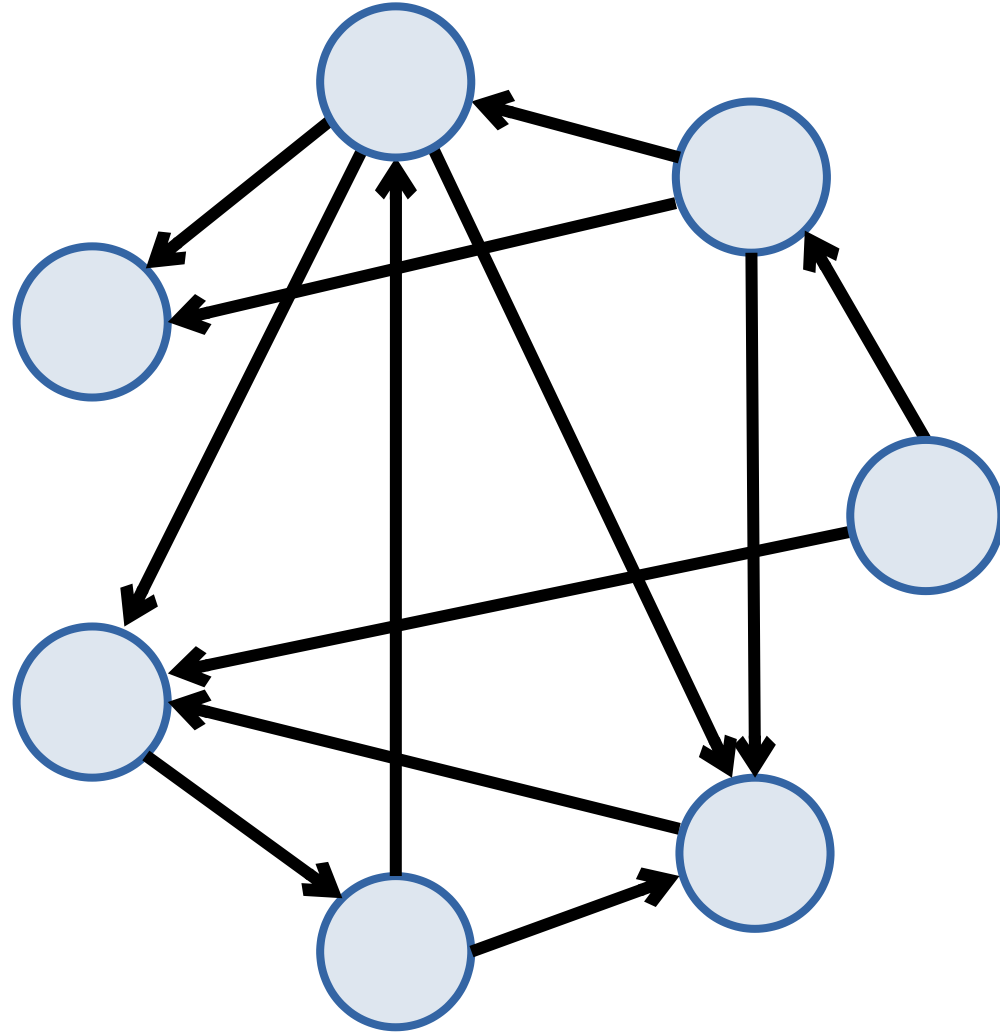
**vs**

**Inference**

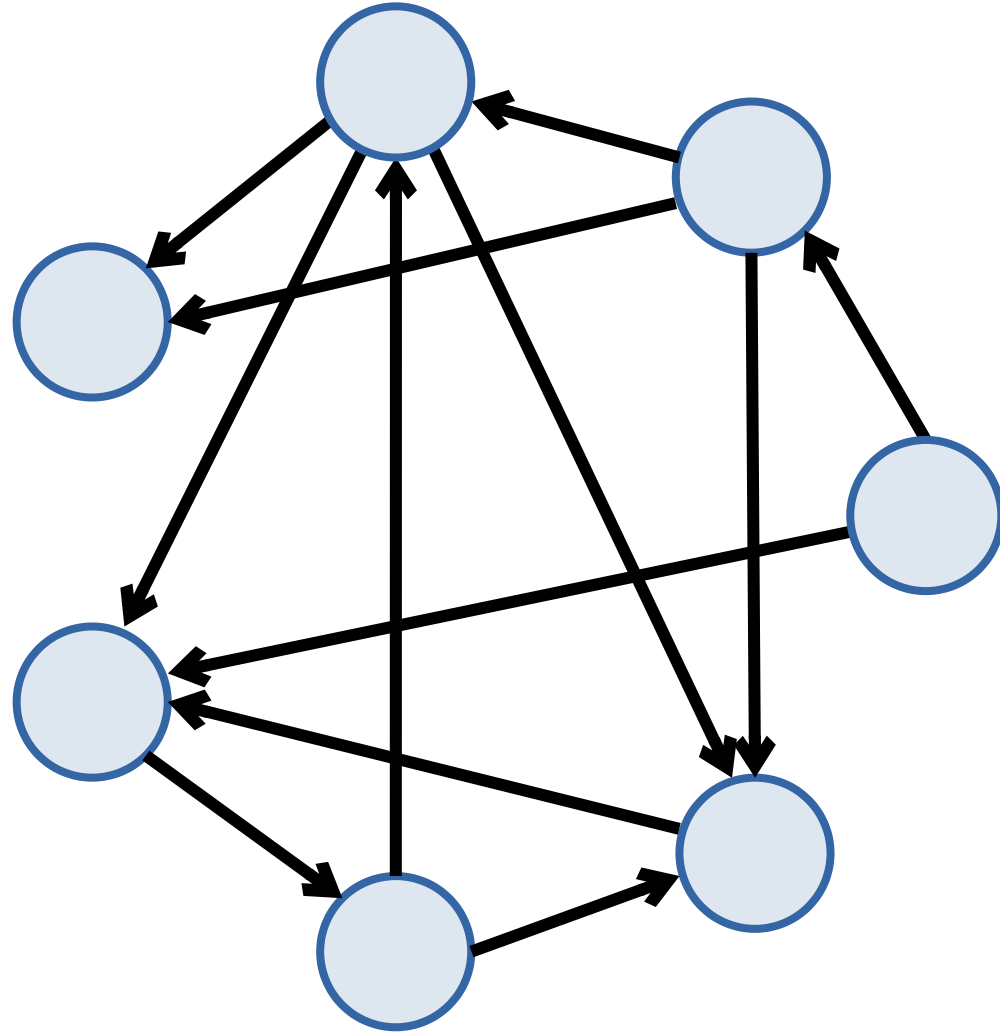
# Graph



# Directed graph



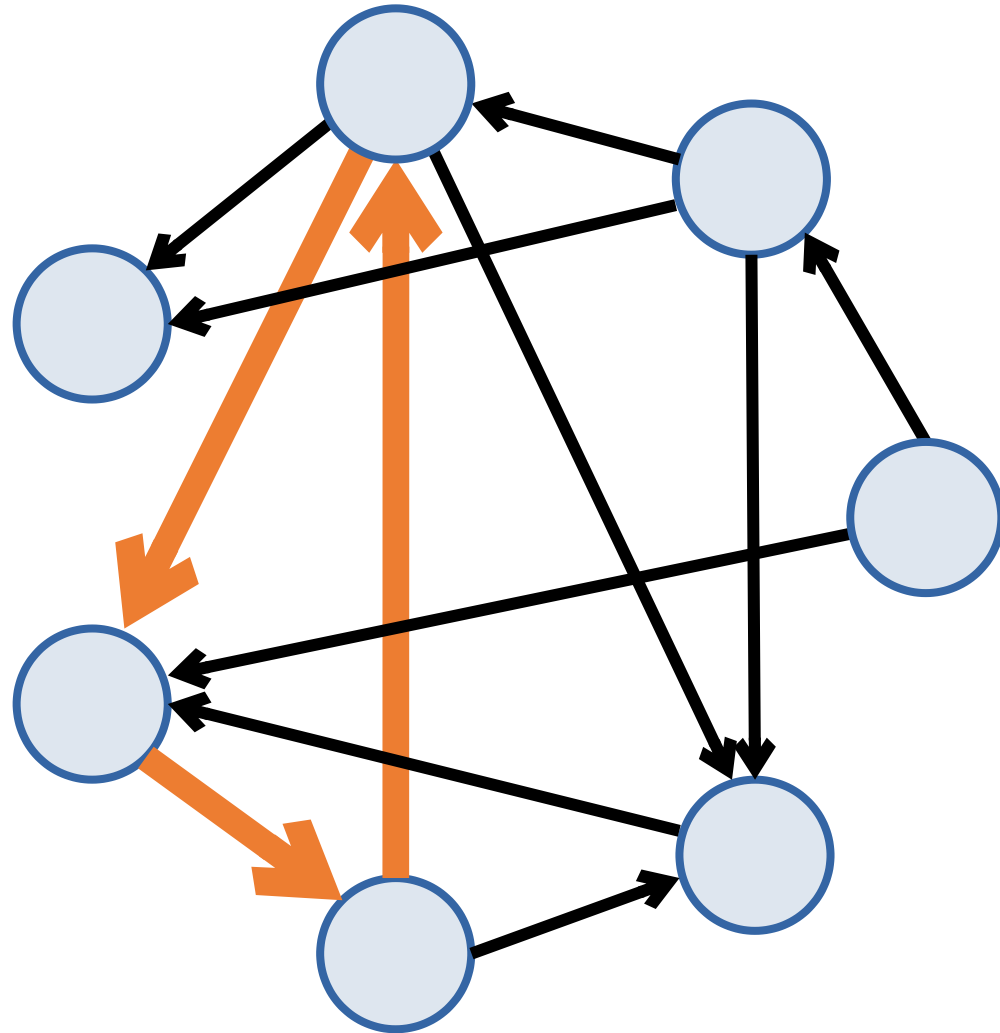
# Directed graph



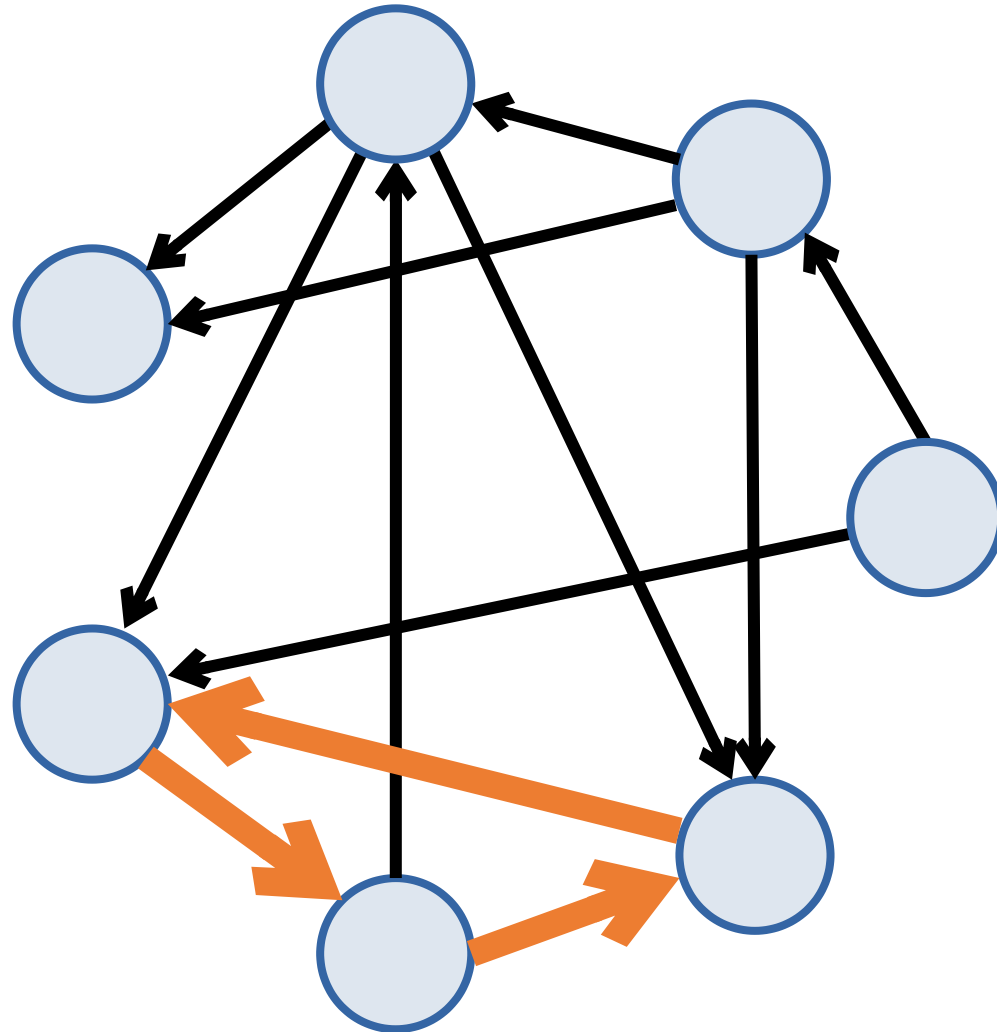
A directed graph  $(V, E)$  is:

- a set  $V$  (vertices, or nodes)
- and a set  $E \subset V \times V$  (edges).

# Directed graph

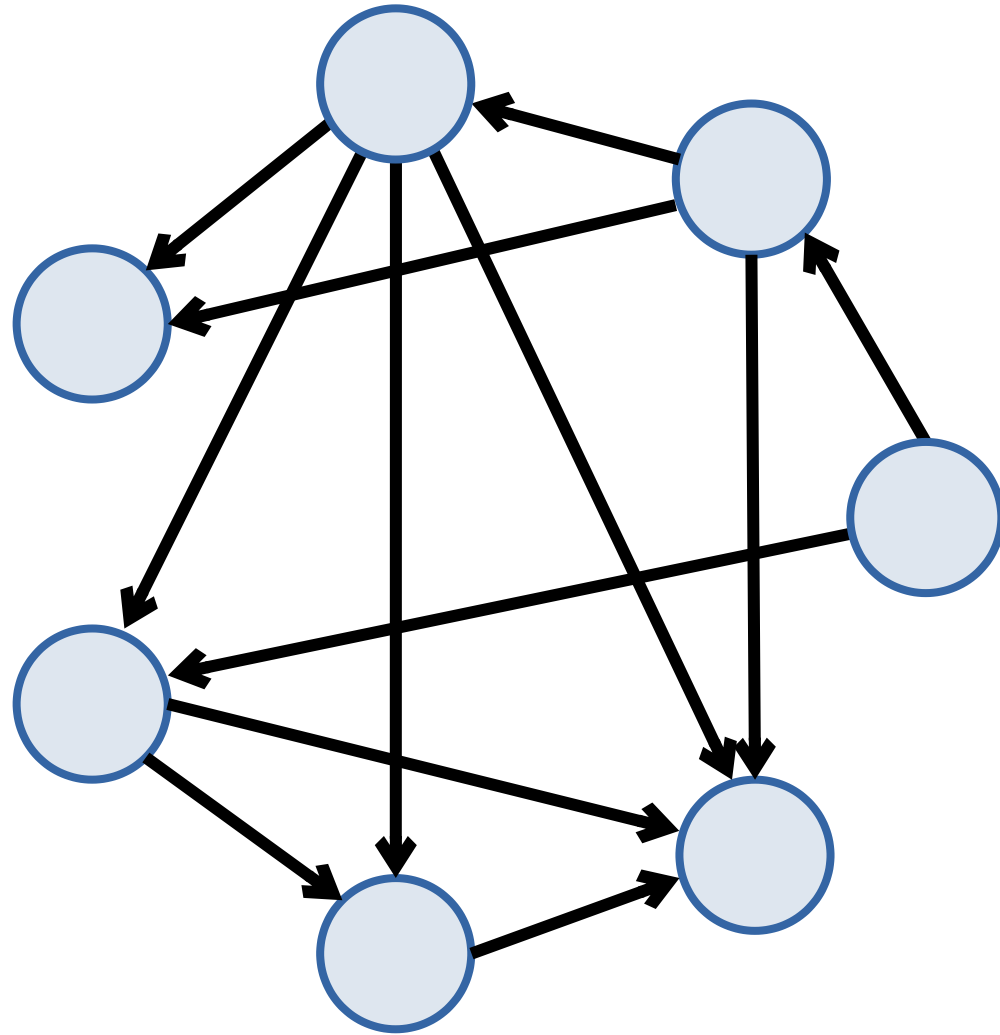


# Directed graph

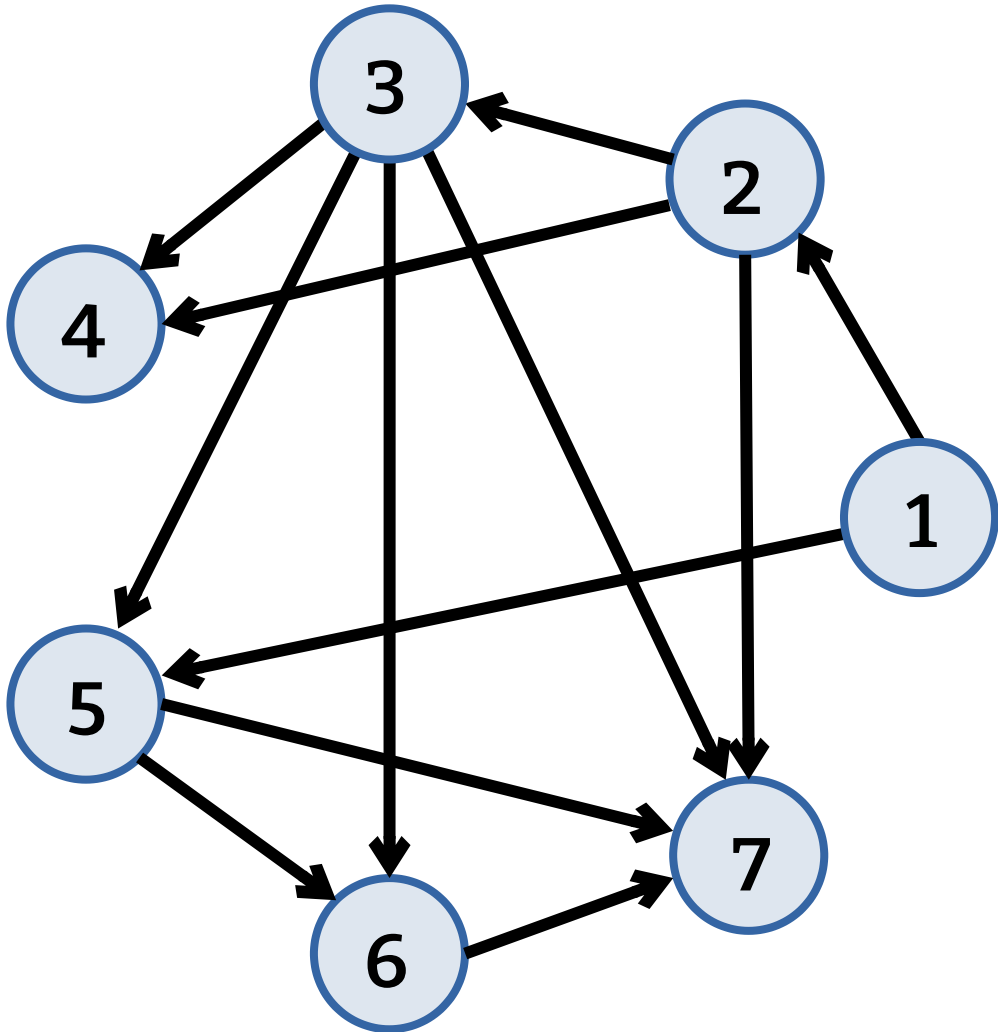




# Directed acyclic graph (DAG)

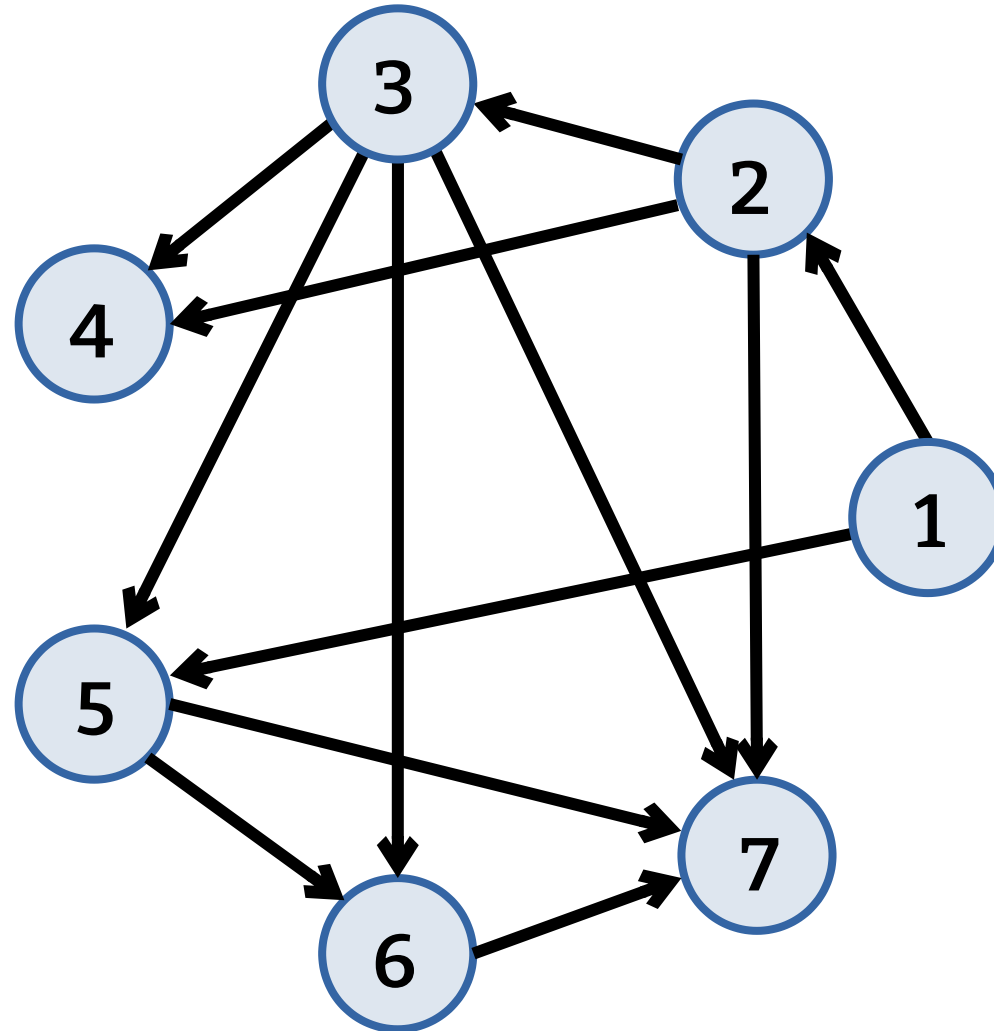


# Adjacency Matrix



1	0	1	0	0	0	0	0
2	0	0	1	0	0	0	1
3	0	0	0	1	1	1	1
4	0	0	0	0	0	0	0
5	0	0	0	0	0	1	1
6	0	0	0	0	0	0	1
7	0	0	0	0	0	0	0
	1	2	3	4	5	6	7

# Topological order



# Data Generation Process

$$X_1 = f_1(\varepsilon_1)$$

$$X_2 = f_2(X_1, \varepsilon_2)$$

$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = f_4(X_2, X_3, \varepsilon_4)$$

$$X_5 = f_5(X_4, \varepsilon_5)$$

# Structural Causal Model

$$X_1 = f_1(\varepsilon_1)$$

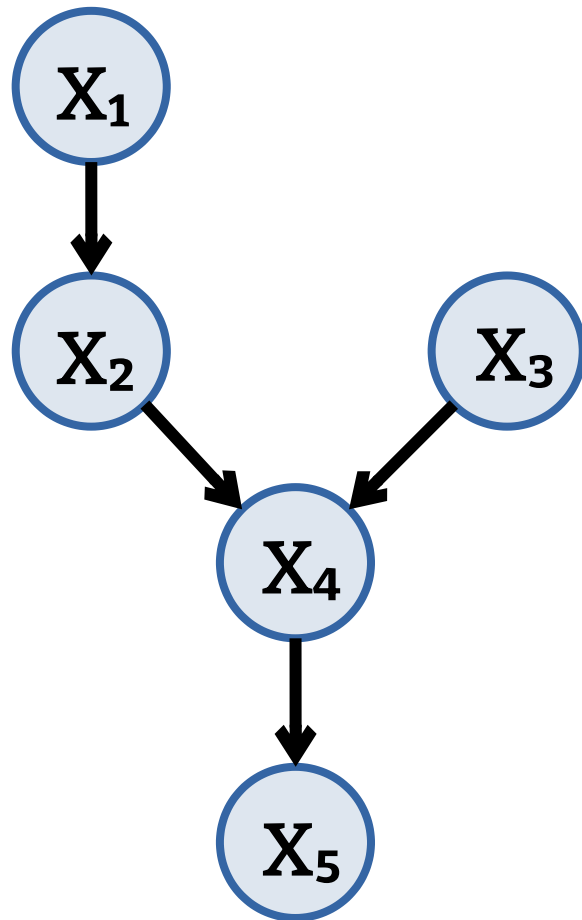
$$X_2 = f_2(X_1, \varepsilon_2)$$

$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = f_4(X_2, X_3, \varepsilon_4)$$

$$X_5 = f_5(X_4, \varepsilon_5)$$

# Causal Graph



# Structural Causal Model

$$X_1 = f_1(\varepsilon_1)$$

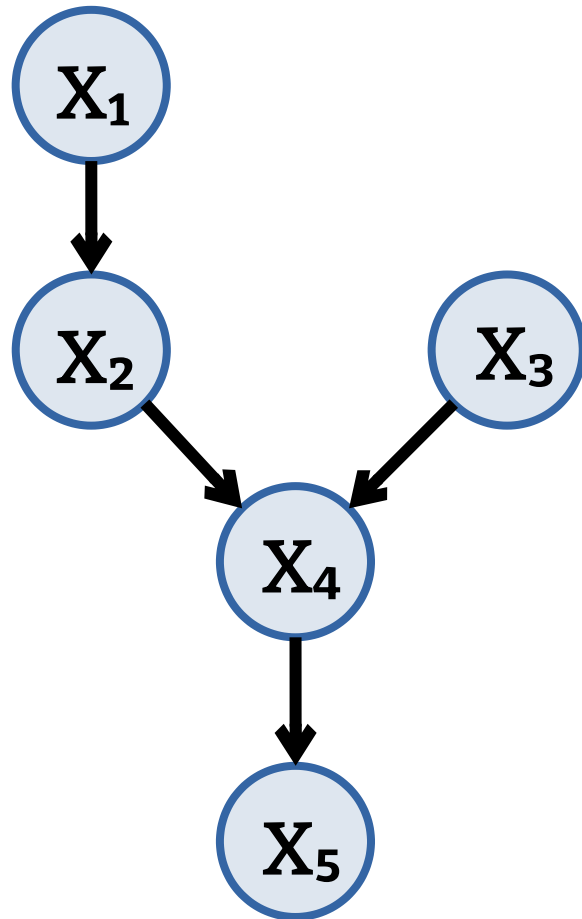
$$X_2 = f_2(X_1, \varepsilon_2)$$

$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = f_4(X_2, X_3, \varepsilon_4)$$

$$X_5 = f_5(X_4, \varepsilon_5)$$

# Causal Discovery



# Causal Inference

$$X_1 = f_1(\varepsilon_1)$$

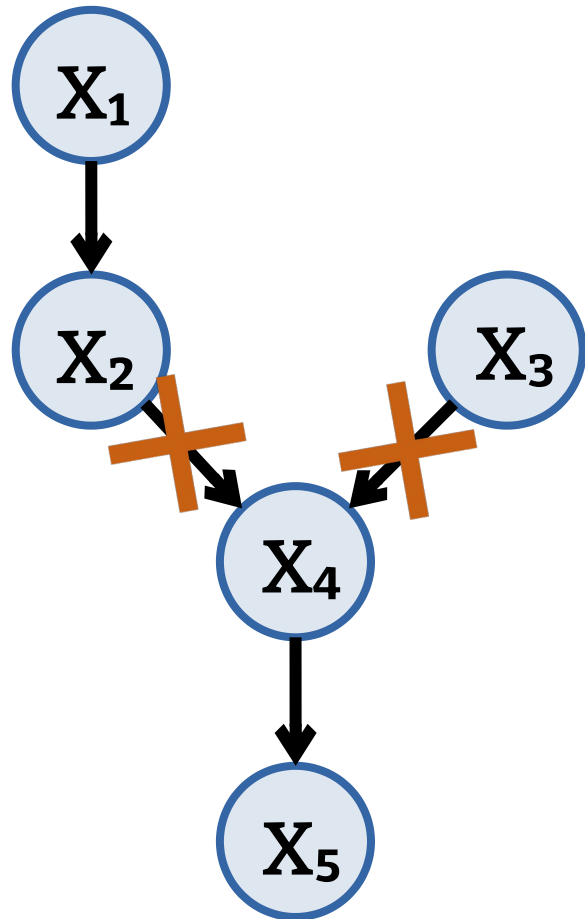
$$X_2 = f_2(X_1, \varepsilon_2)$$

$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = f_4(X_2, X_3, \varepsilon_4)$$

$$X_5 = f_5(X_4, \varepsilon_5)$$

# Intervention



$$X_1 = f_1(\varepsilon_1)$$

$$X_2 = f_2(X_1, \varepsilon_2)$$

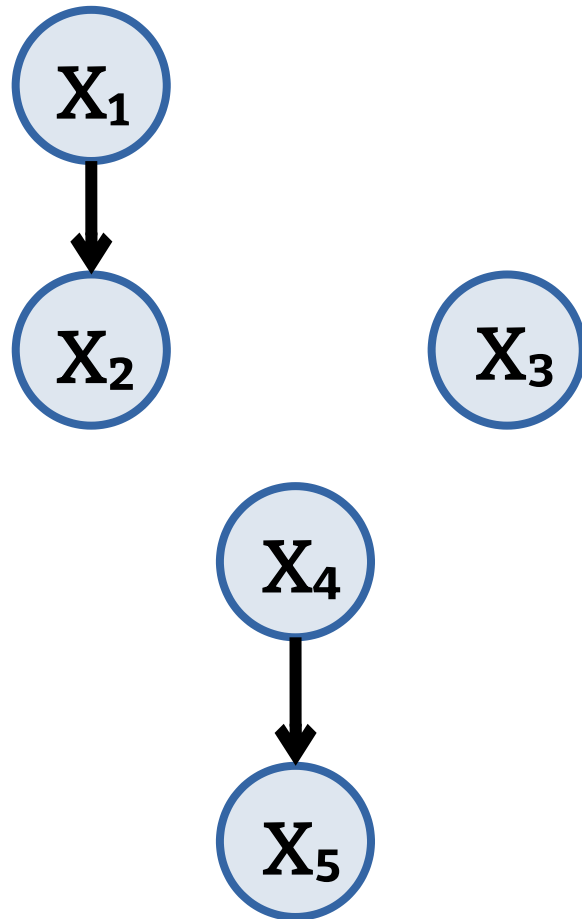
$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = X_4$$

$$X_5 = f_5(X_4, \varepsilon_5)$$



# Intervention



$$X_1 = f_1(\varepsilon_1)$$

$$X_2 = f_2(X_1, \varepsilon_2)$$

$$X_3 = f_3(\varepsilon_3)$$

$$X_4 = X_4$$

$$X_5 = f_5(X_4, \varepsilon_5)$$

# Discovery vs Inference

- **Causal discovery:** finding the causal graph from data
- **Causal inference:** using the structural causal model to answer “what if” questions

# Causality ladder

- **Association:**  $E[Y|T=t]$
- **Intervention:**  $E[Y|\text{do}(T=t)]$
- **Counterfactuals:**  $E[Y|\text{do}(T=t), T=t']$



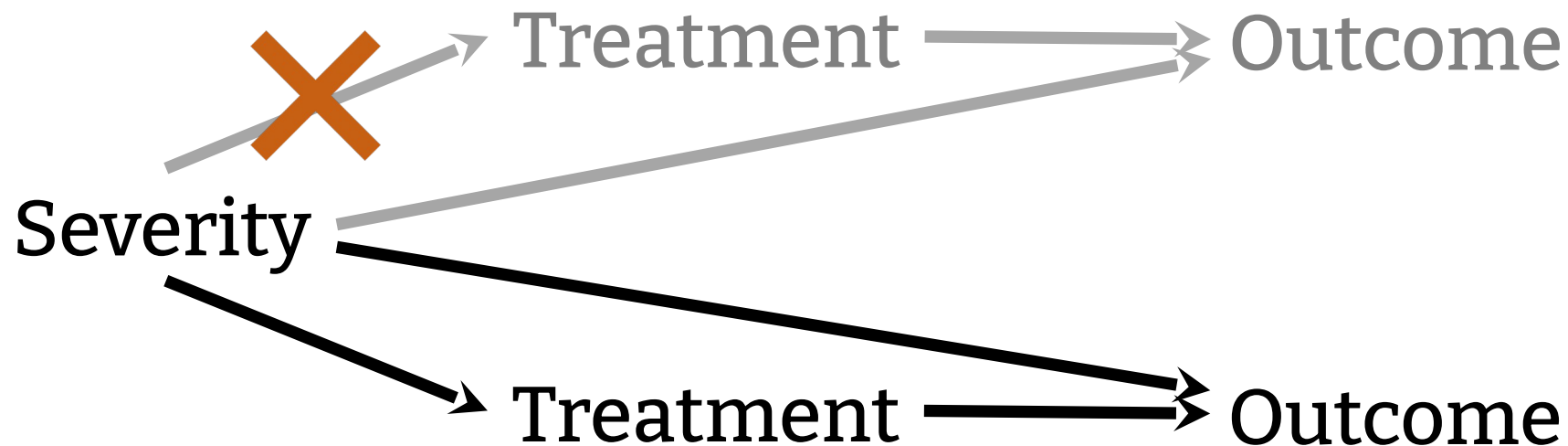
# Causality ladder

- Association:  $E[Y|T=t]$
- **Intervention**:  $E[Y|\text{do}(T=t)]$
- Counterfactuals:  $E[Y|\text{do}(T=t), T=t']$



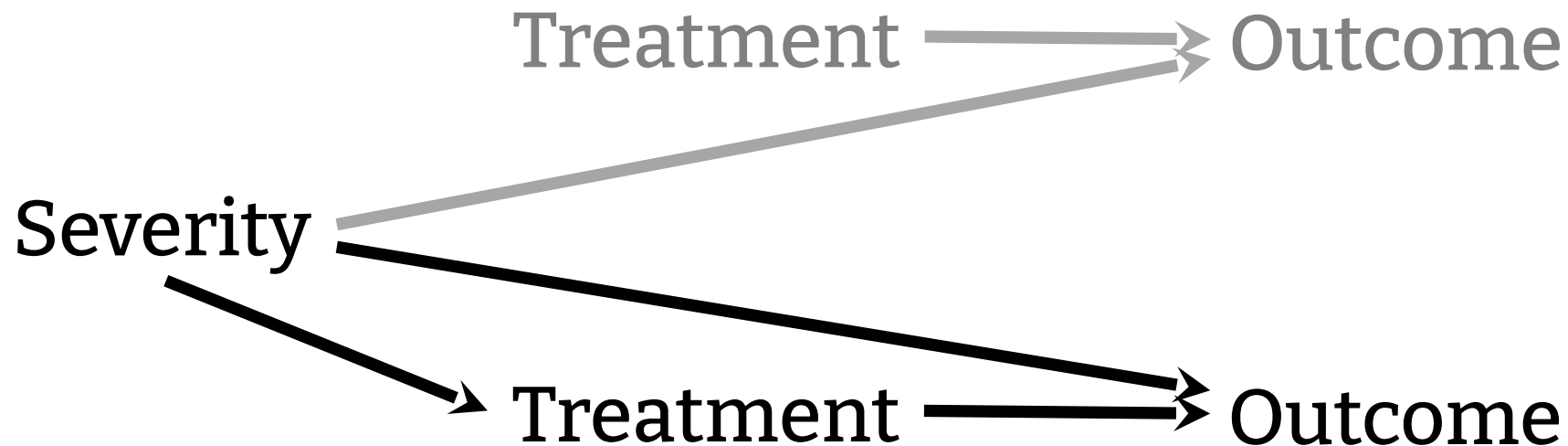
# Causality ladder

- Association:  $E[Y|T=t]$
- Intervention:  $E[Y|\text{do}(T=t)]$
- **Counterfactuals:  $E[Y|\text{do}(T=t), T=t']$**



# Causality ladder

- Association:  $E[Y|T=t]$
- Intervention:  $E[Y|\text{do}(T=t)]$
- **Counterfactuals:  $E[Y|\text{do}(T=t), T=t']$**



# Causality Ladder

- **Association**

If the patient received an aggressive treatment, it means his condition was already severe: the expected outcome is bad.

- **Intervention**

If we were to give all patients an aggressive treatment, the outcome would be good on average.

- **Counterfactual**

This specific patient received the non-aggressive treatment; this means his condition was mild; if we had given him the aggressive treatment, the outcome would have been good.

# Simpson's paradox

		Condition		
		Mild	Severe	Total
Treatment	A	15% (210 / 1400)	30% (30 / 100)	<b>16%</b> (240 / 1500)
	B	<b>10%</b> (5 / 50)	<b>20%</b> (100 / 500)	19% (105 / 550)



# Fundamental Problem

We often want to compute the average treatment effect,

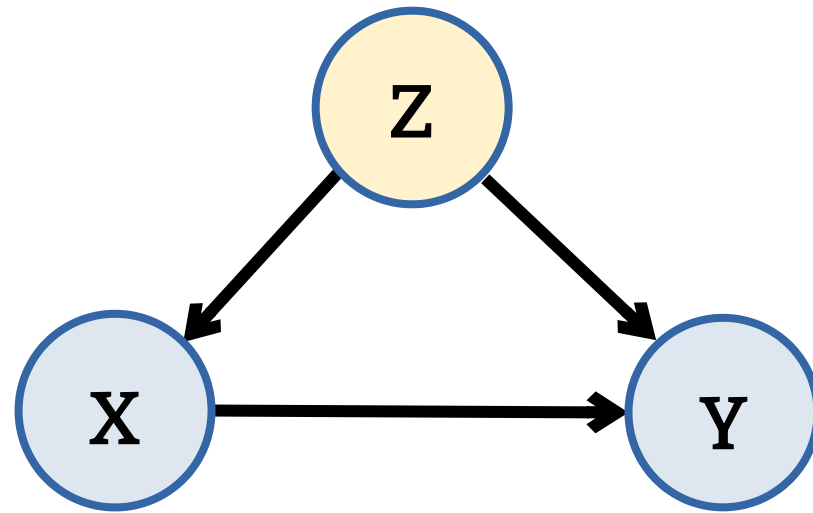
$$ATE = E[Y|do(T=1)] - E[Y|do(T=0)],$$

but, for each subject, we either have  $T=1$  or  $T=0$ , so we do not know the other.

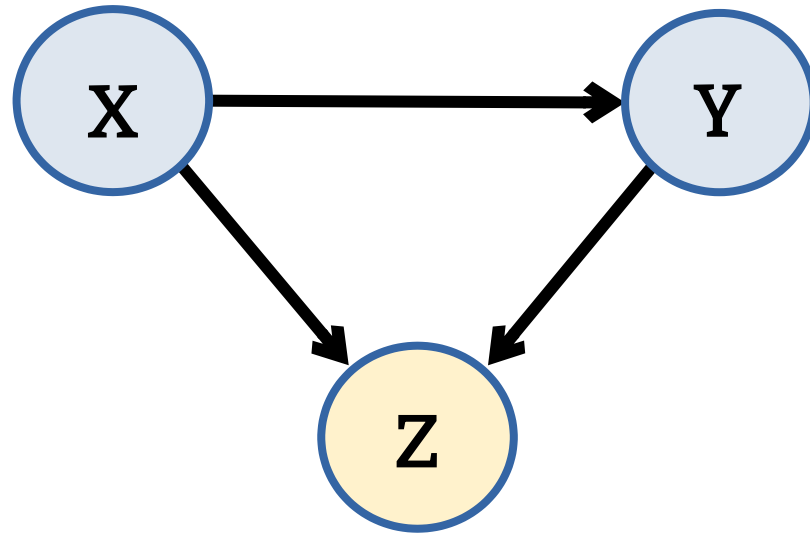
# Do Calculus

**Do Calculus is a set of rules to compute, when possible, the effect interventions would have from observational data alone, given the causal graph.**

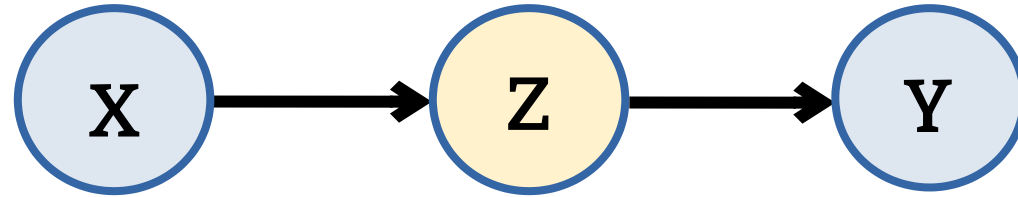
# Confounder



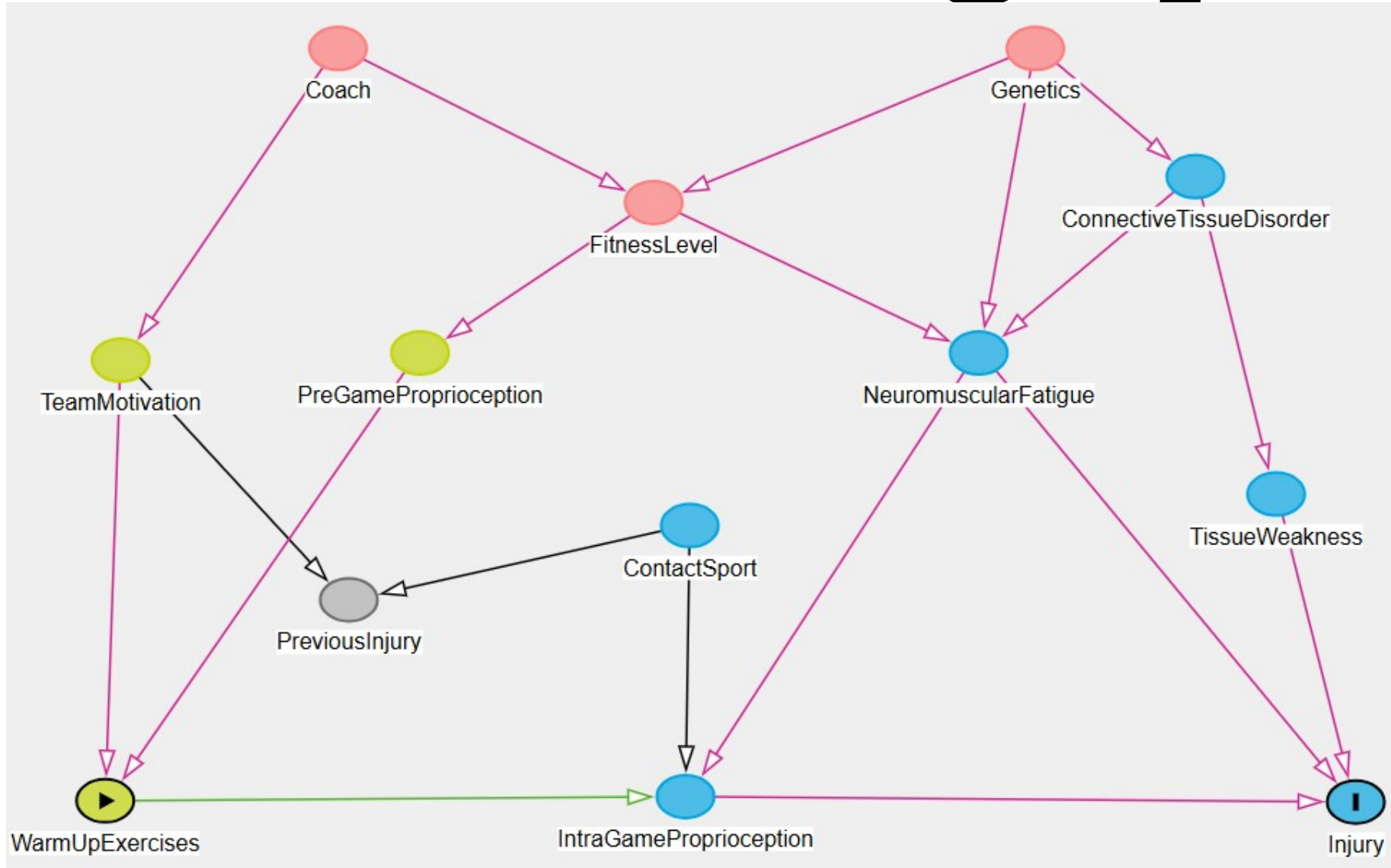
# Collider



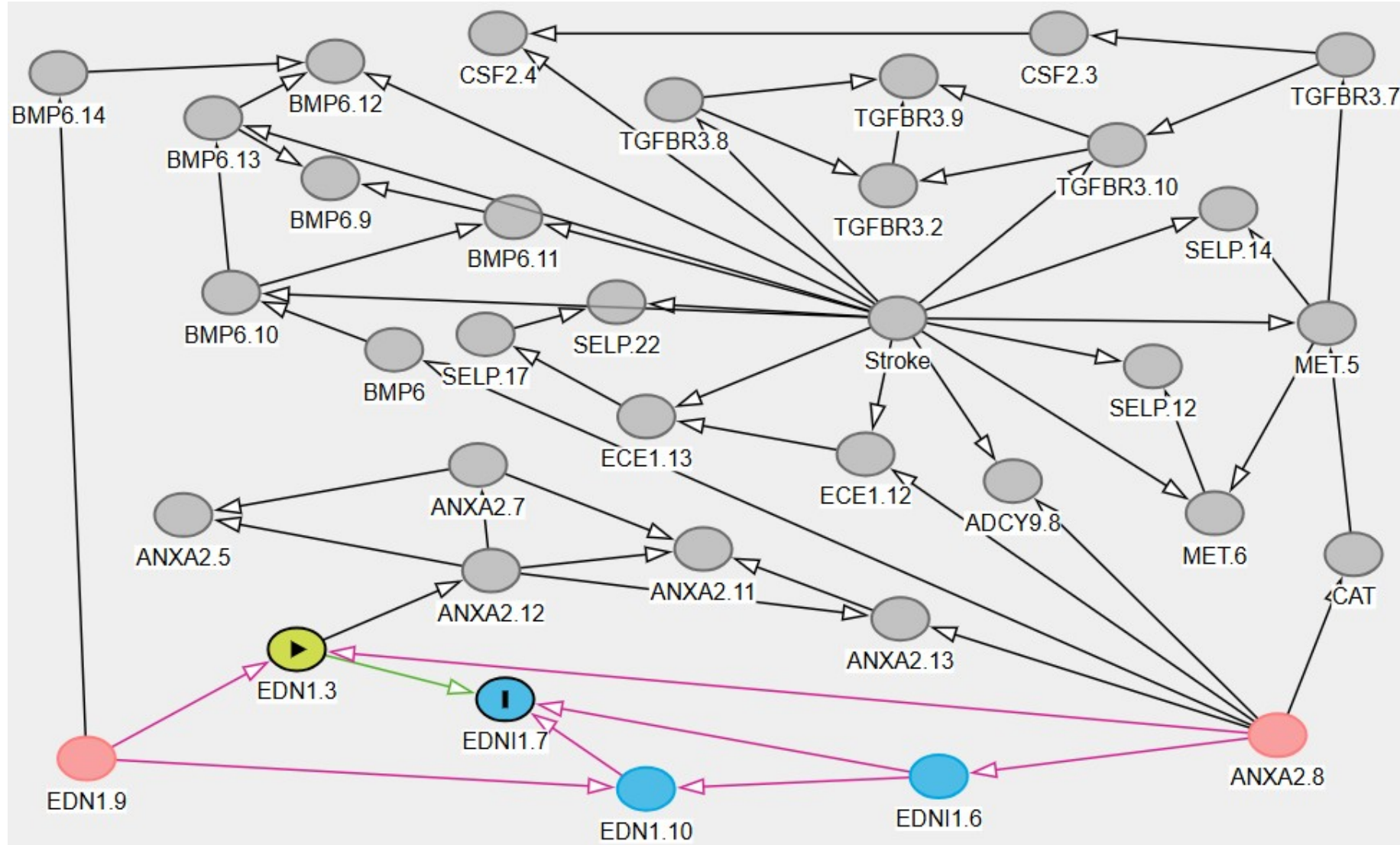
# Mediator



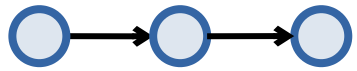
# Real-world causal graphs



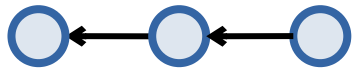
# Real-world causal graphs



# Open and closed paths



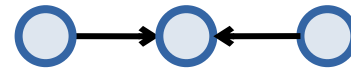
open



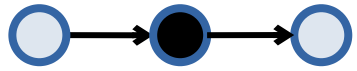
open



open



closed



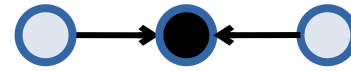
closed



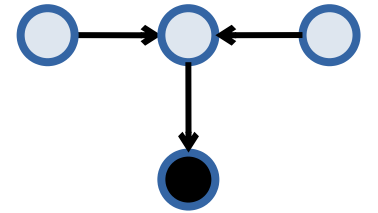
closed



closed



open



open

○ Not in the conditioning set

● In the conditioning set



# Causal inference

To assess the strength of the causal relation  $X \rightarrow Y$ :

- List all the (undirected) paths from X to Y
- All the non-causal paths should be blocked; if not, condition on one or more nodes to block them.
- All the causal paths should be open; if not, adjust the conditioning set to unblock them.

# **Causal Discovery Algorithms**

# Causal Discovery Algorithms

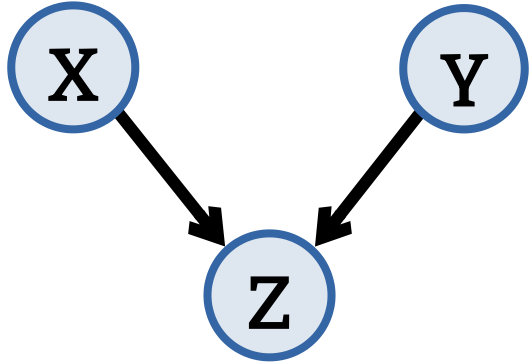
- **PC:** Conditional independence tests
- **GES:** Scores
- **LiNGAM:** Independent component analysis
- **NOTEARS:** Optimization

# PC Algorithm

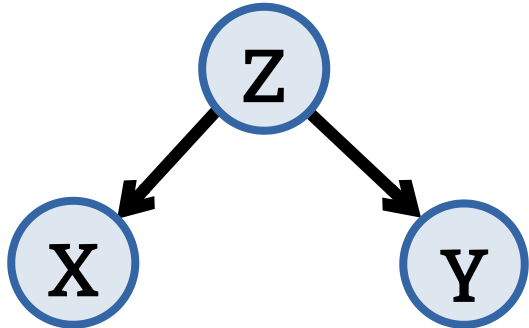
# Conditional Independence



$$X \not\perp Y$$

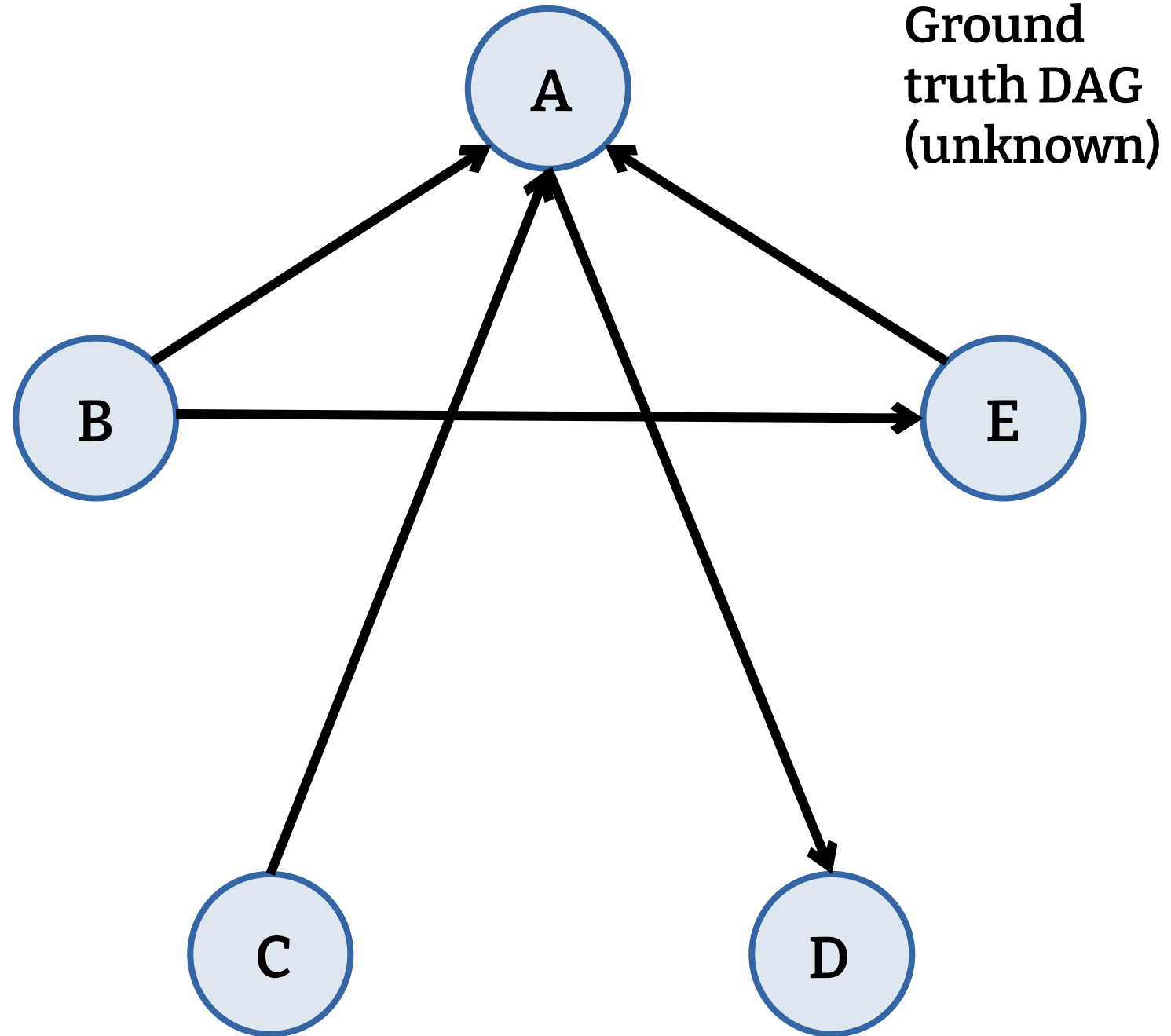


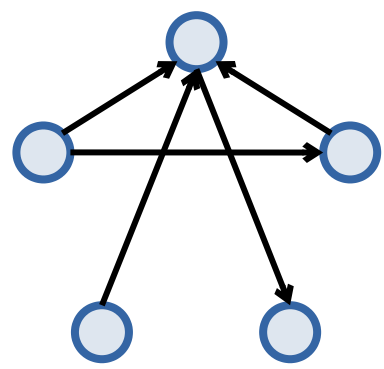
$$X \perp Y$$
$$X \not\perp Y | Z$$



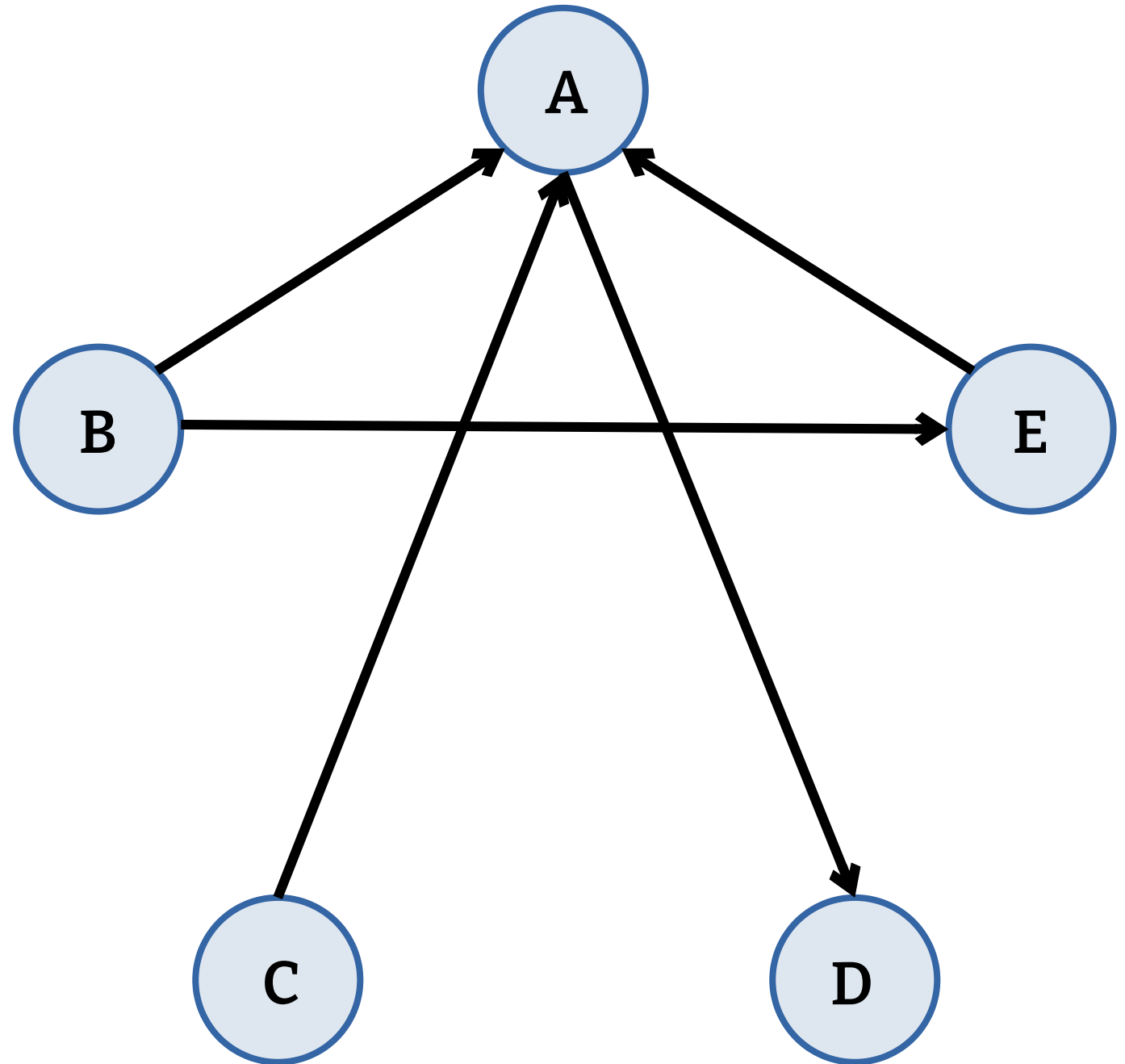
$$X \not\perp Y$$
$$X \perp Y | Z$$

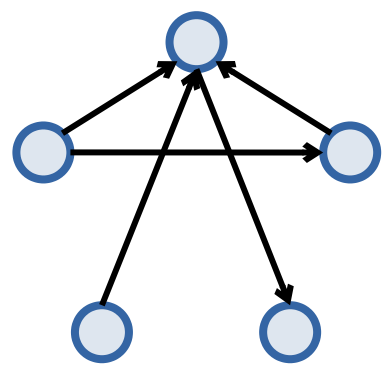
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



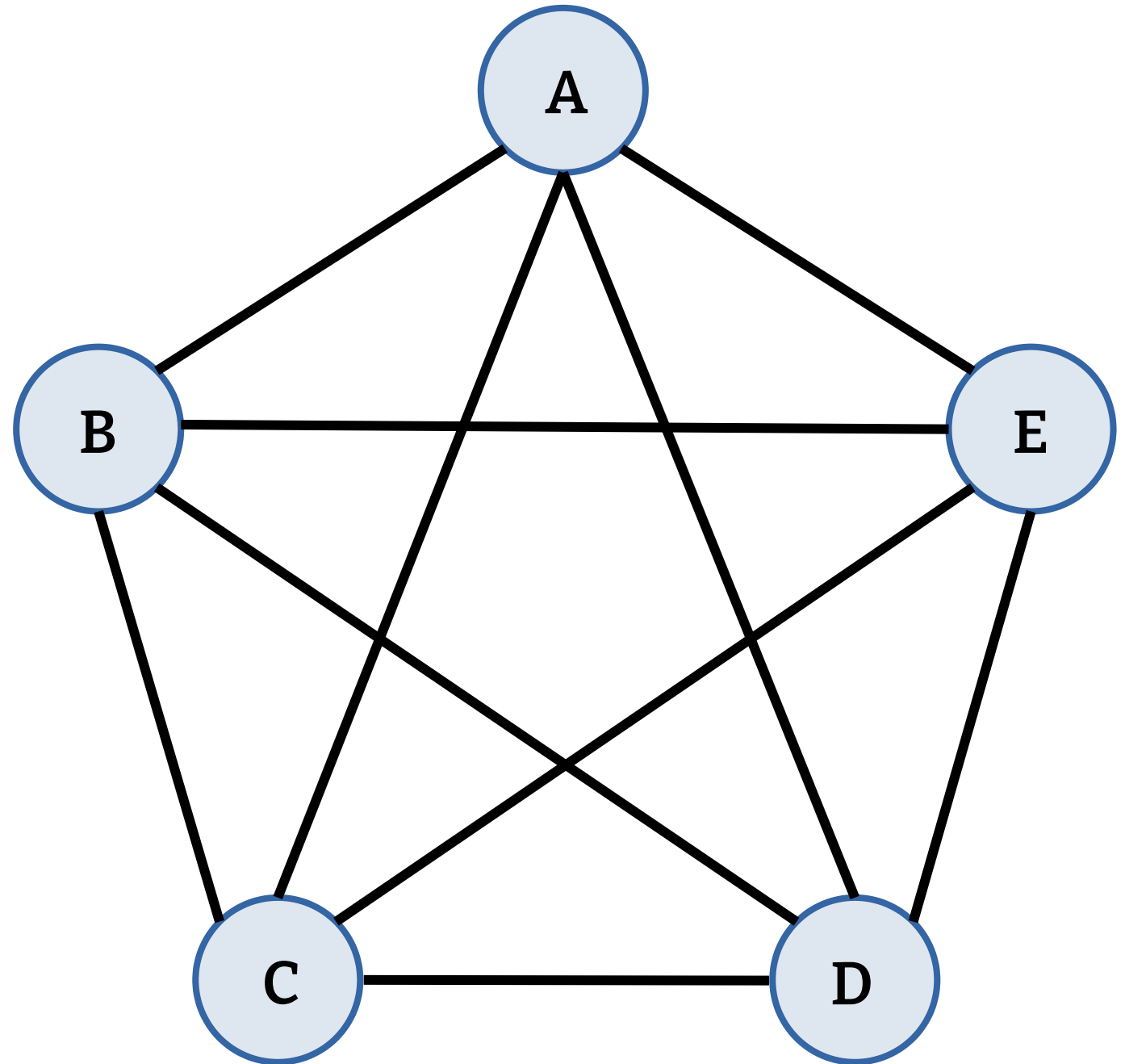


- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders

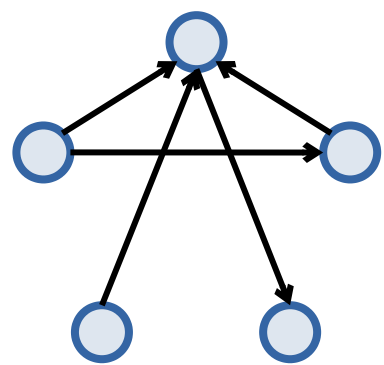




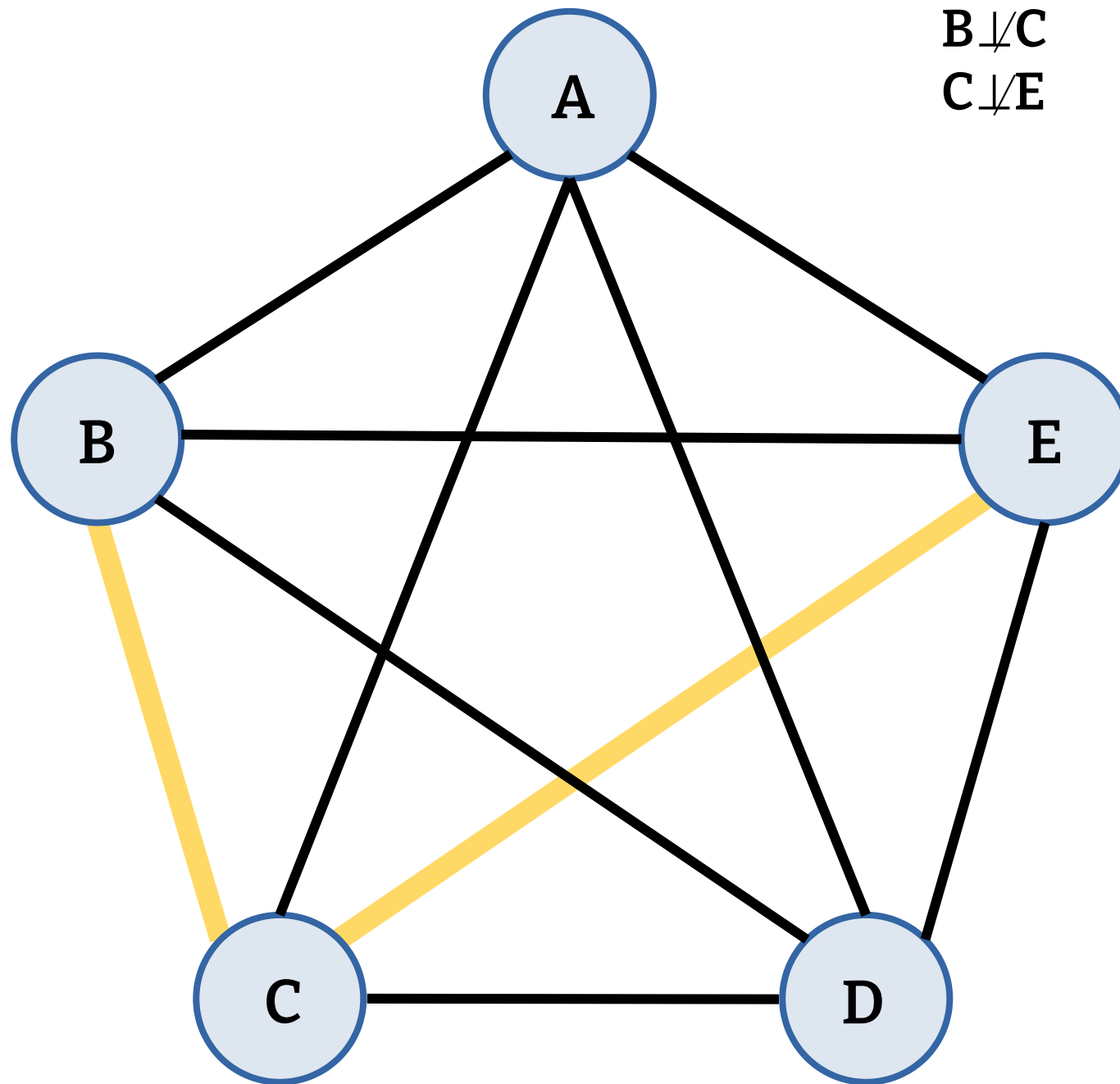
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders

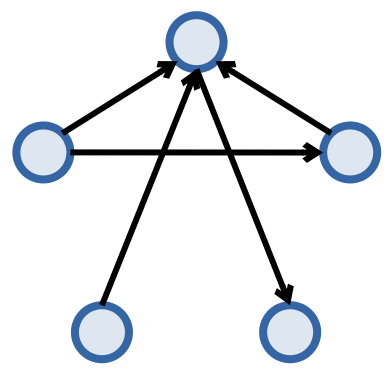




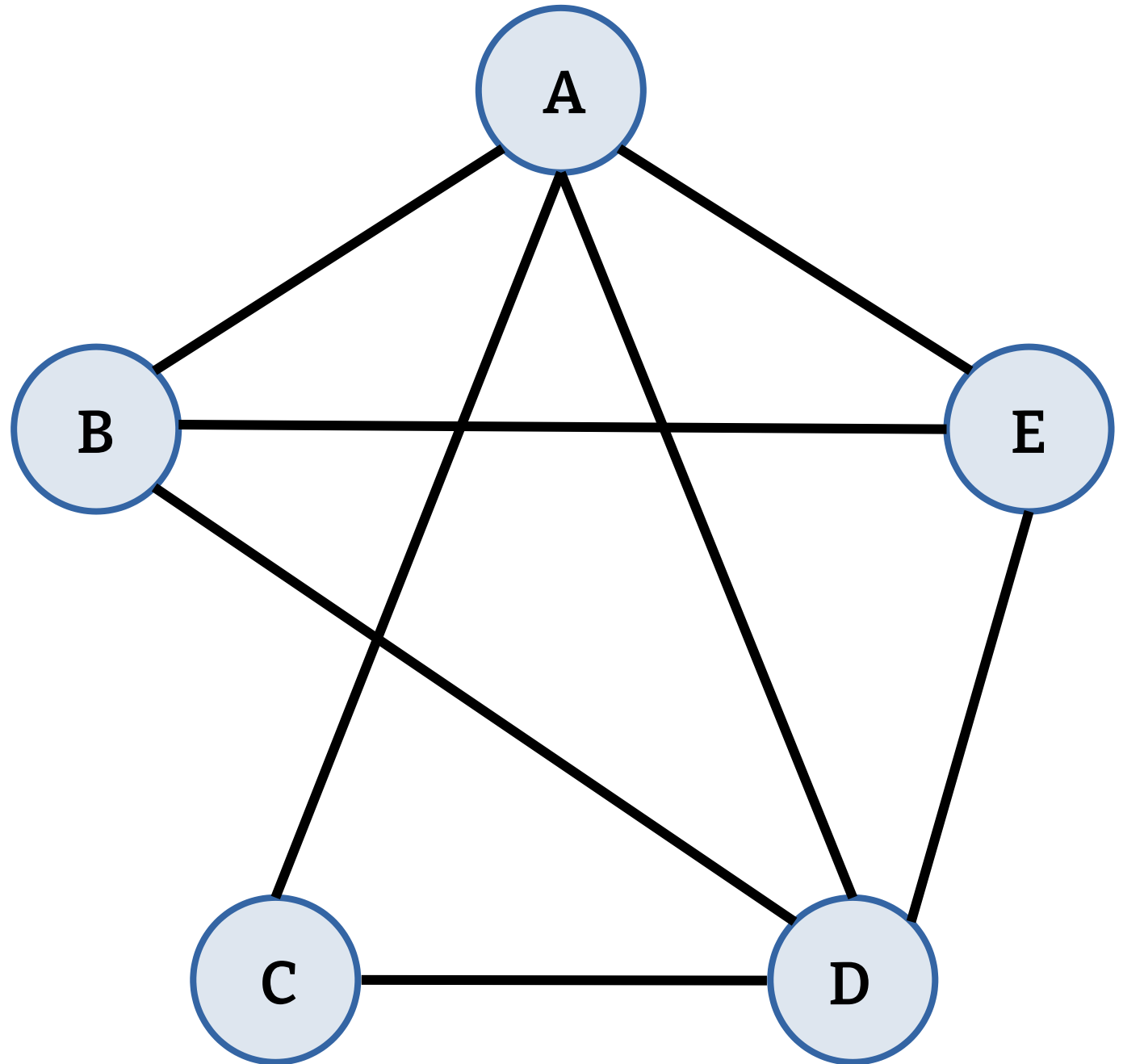


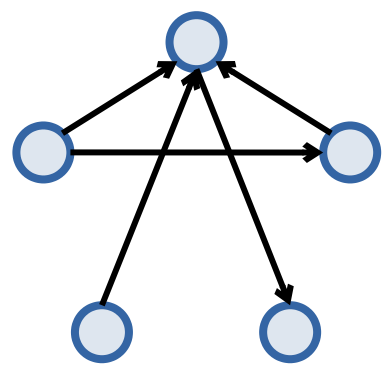
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



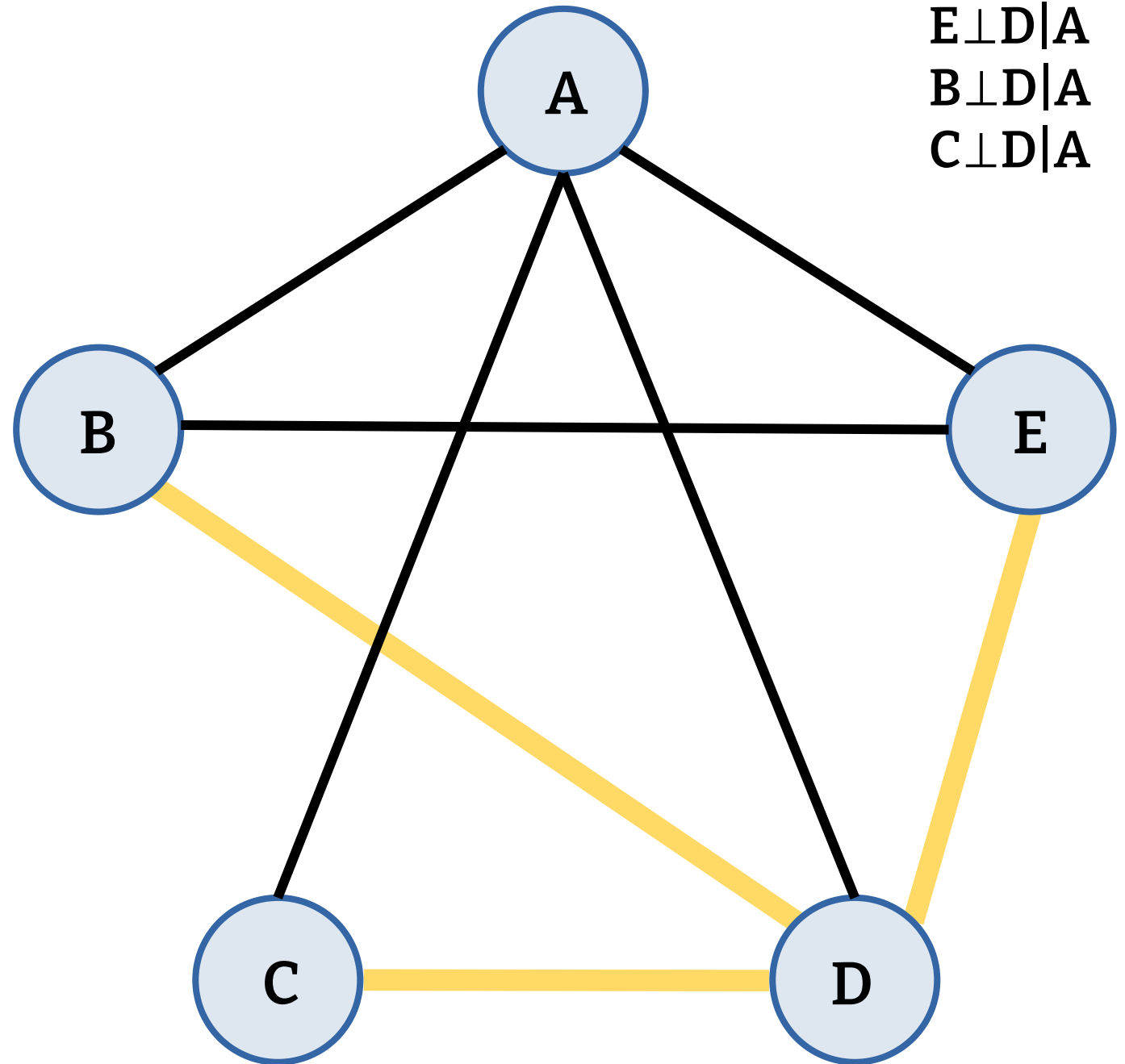


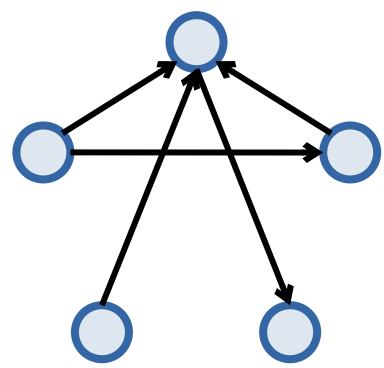
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



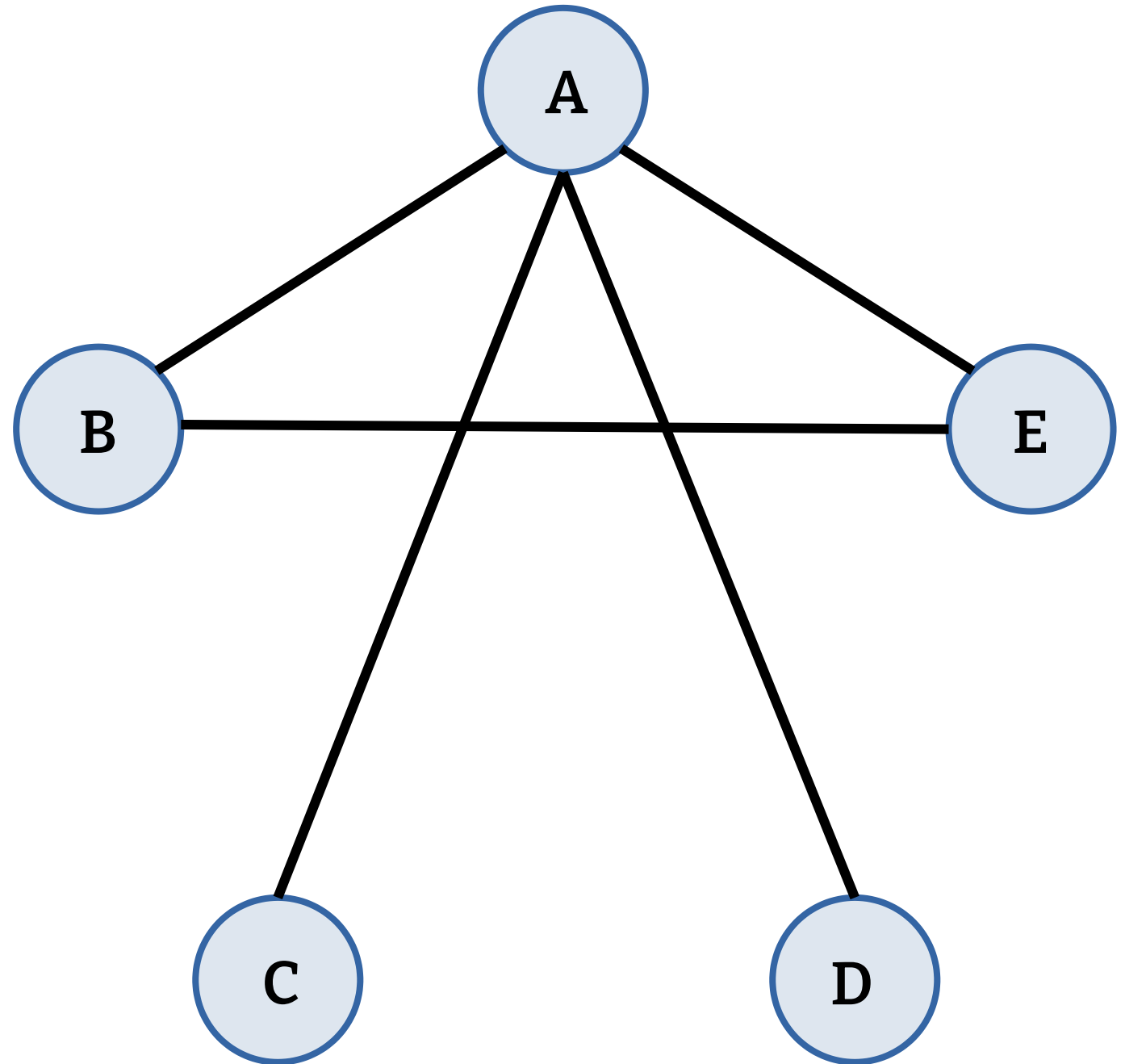


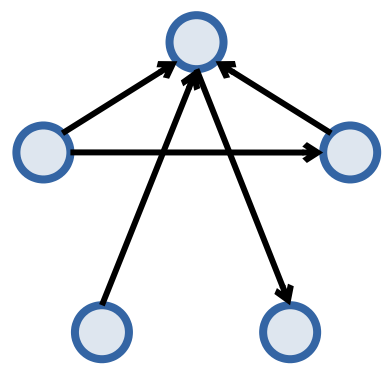
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



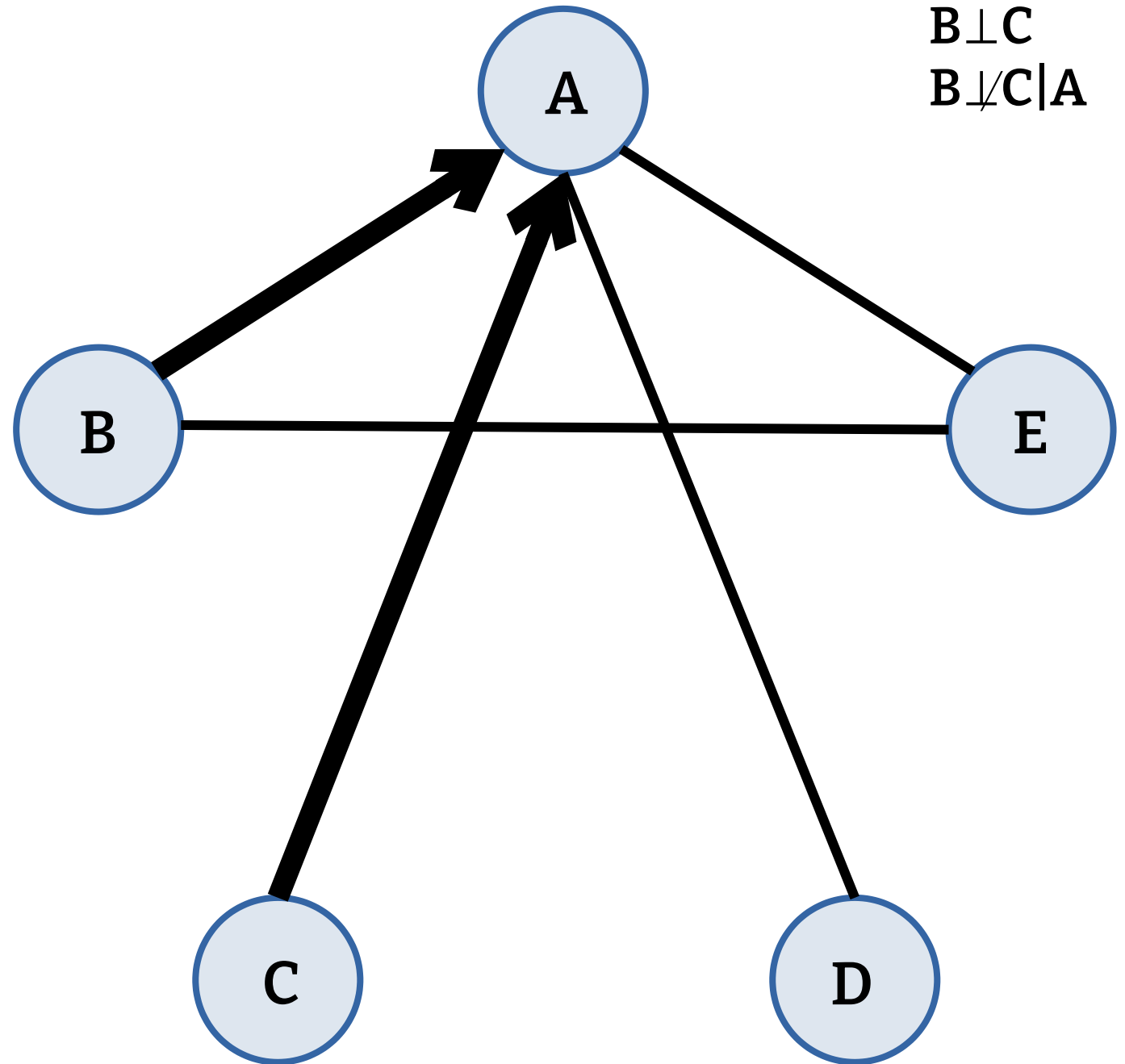


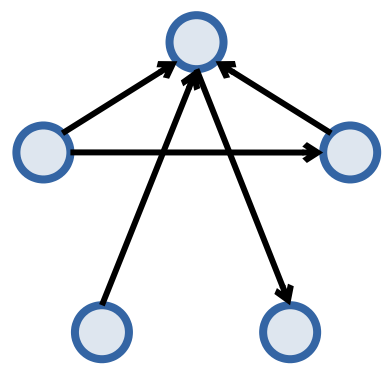
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



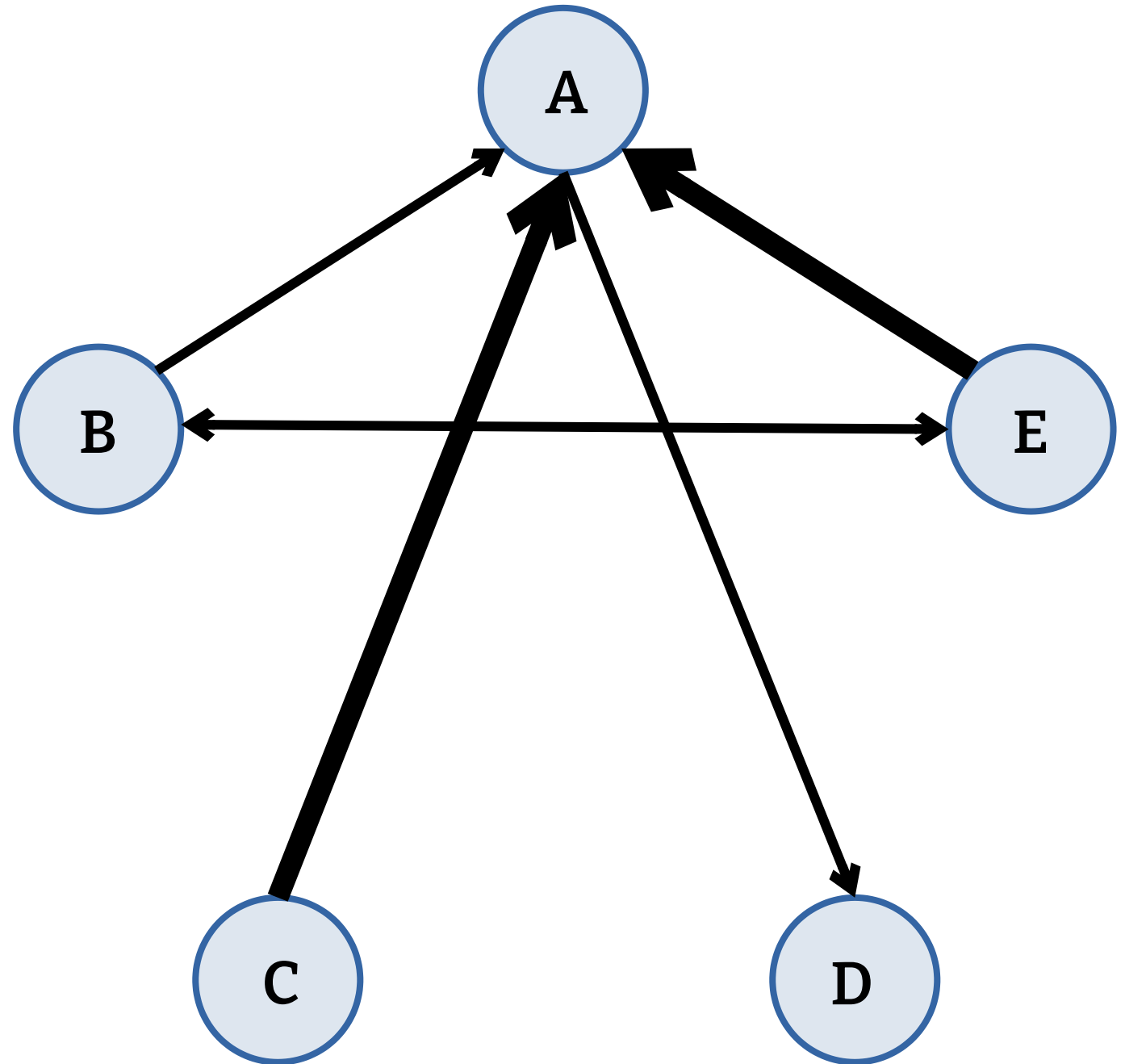


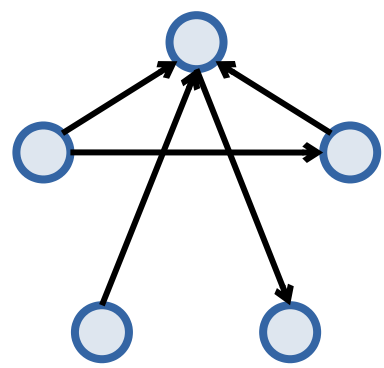
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



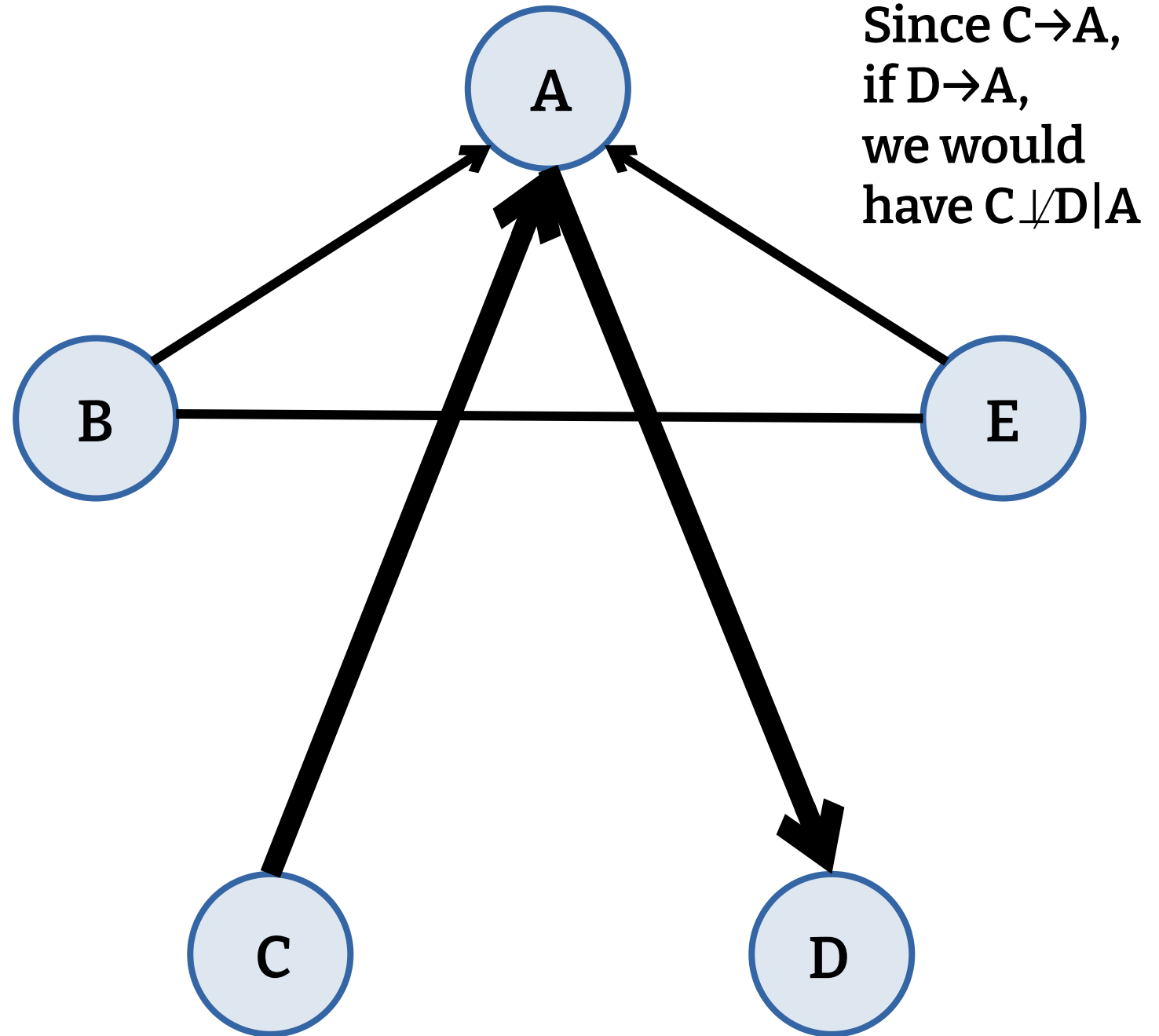


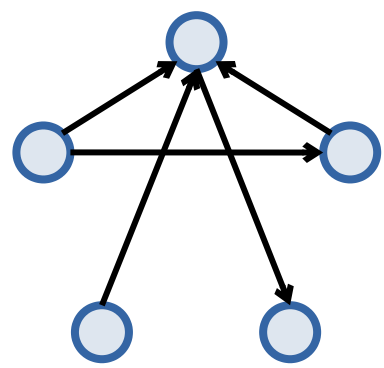
- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



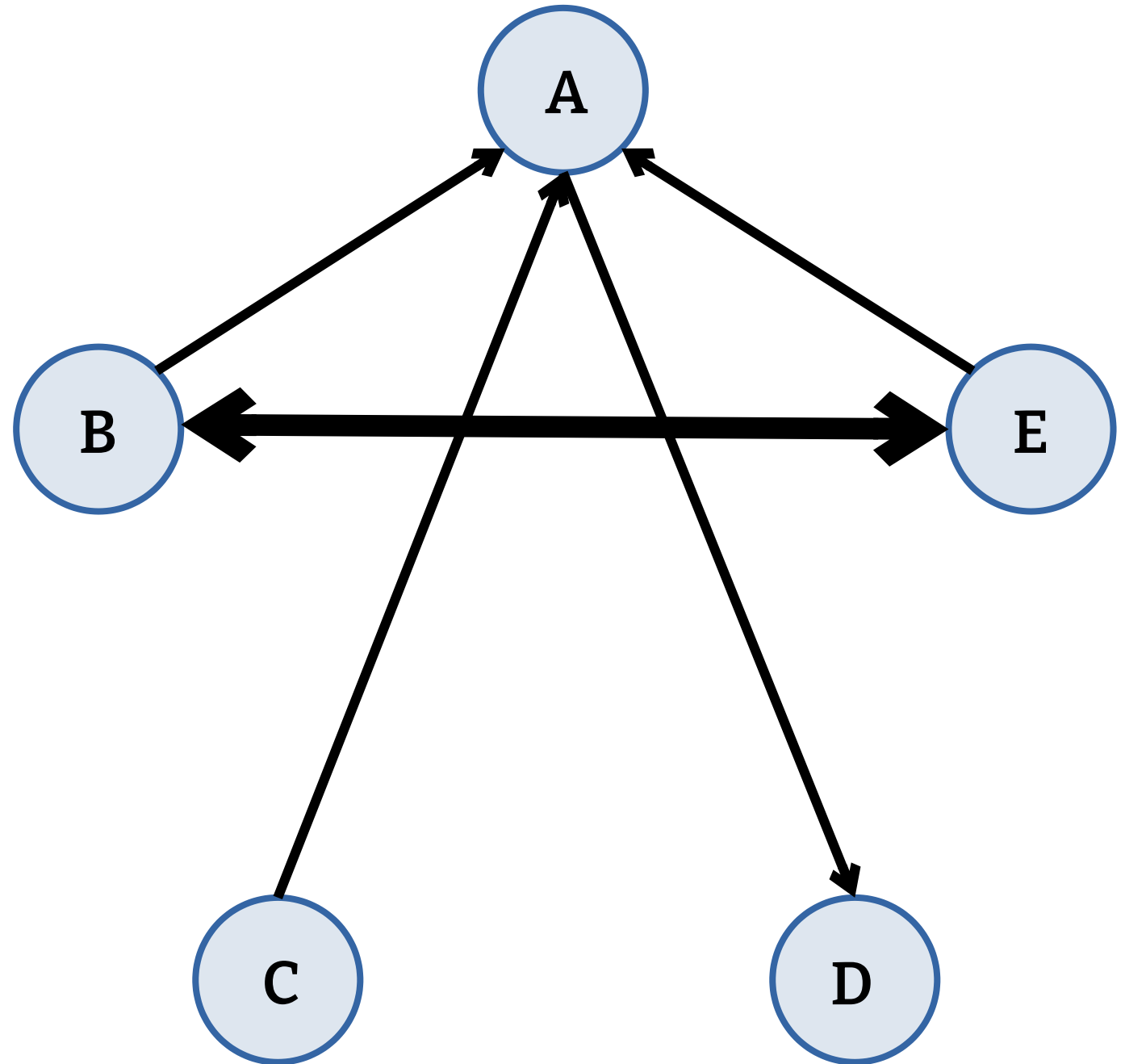


- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders



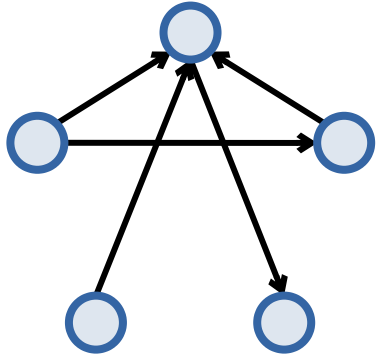


- Start with a complete (undirected) graph
- Remove the edge  $X-Y$  if  $X \perp Y | Z$  for some (possibly empty) set of nodes  $Z$
- For all  $X-Z-Y$ , if  $X \perp Y$  and  $X \not\perp Y | Z$ , we have a collider  $X \rightarrow Z \leftarrow Y$
- Propagate the orientation, assuming we have found all the colliders

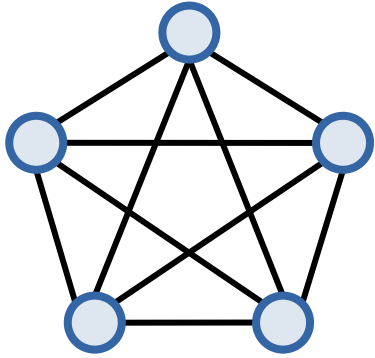




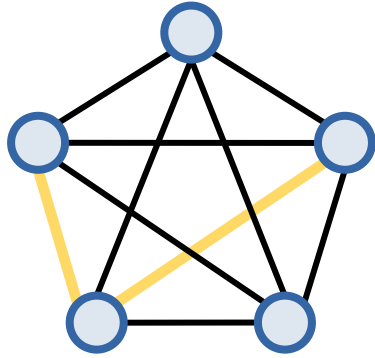
# PC algorithm



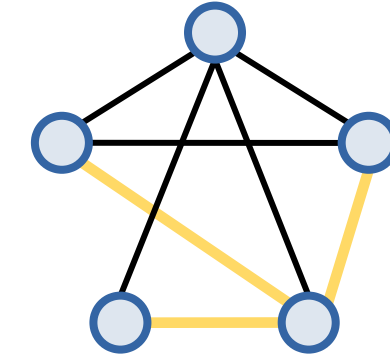
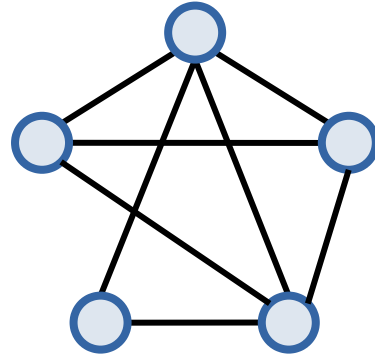
Ground truth



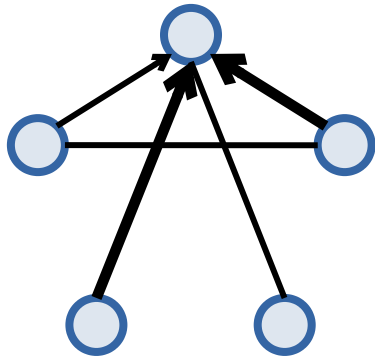
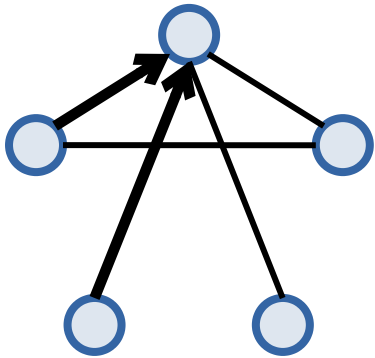
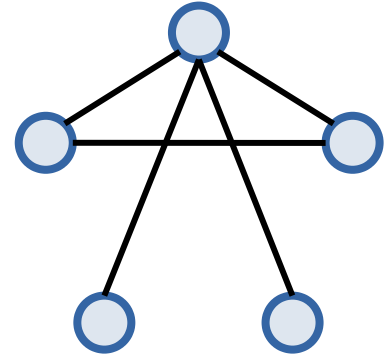
Complete graph



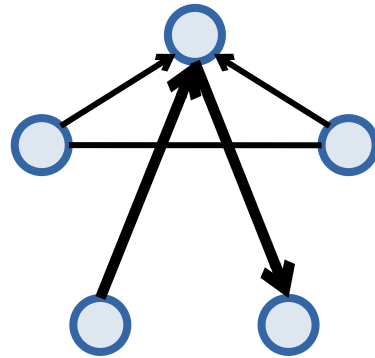
Independence tests



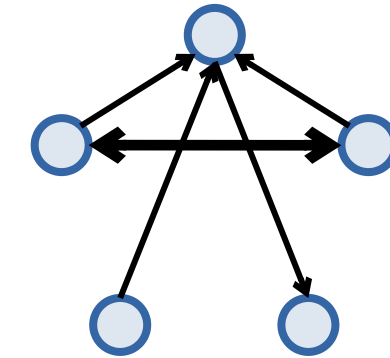
Conditional independence tests



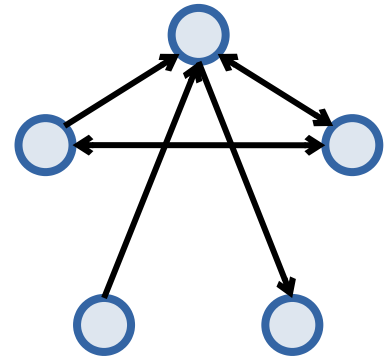
Colliders



Non-collider



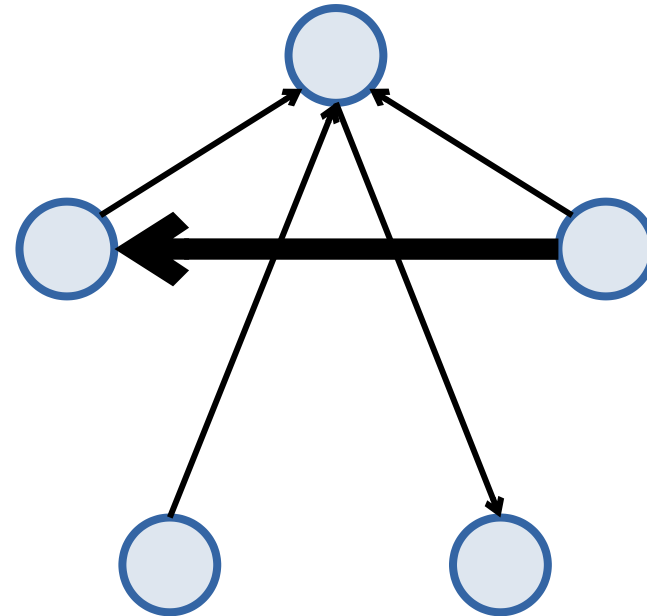
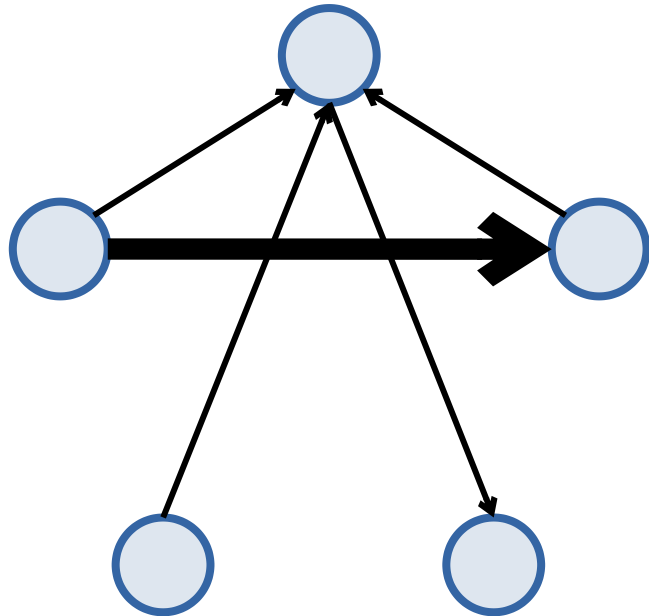
Undirected edges



Final result

# Markov equivalence class

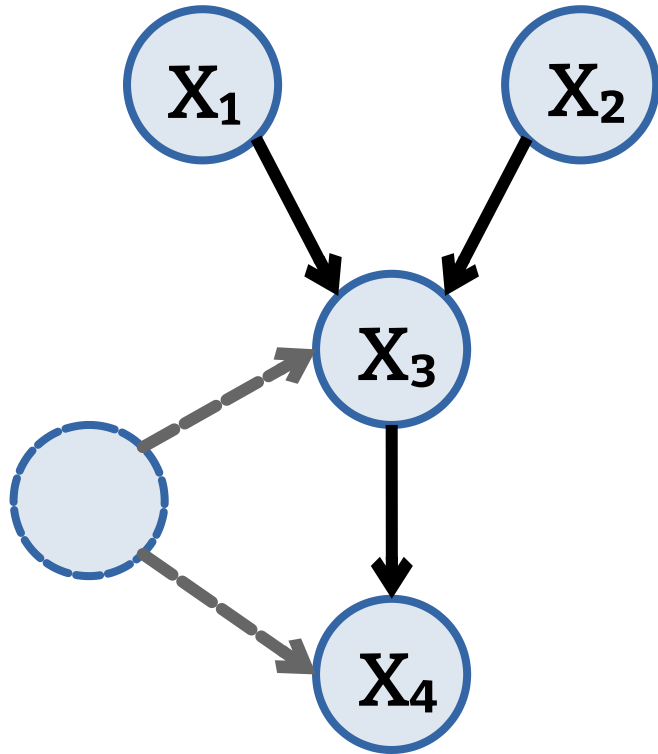
Conditional independence relations cannot separate all DAGs: they can only recover the Markov equivalence class (MEC) of the causal model.



# Independence tests

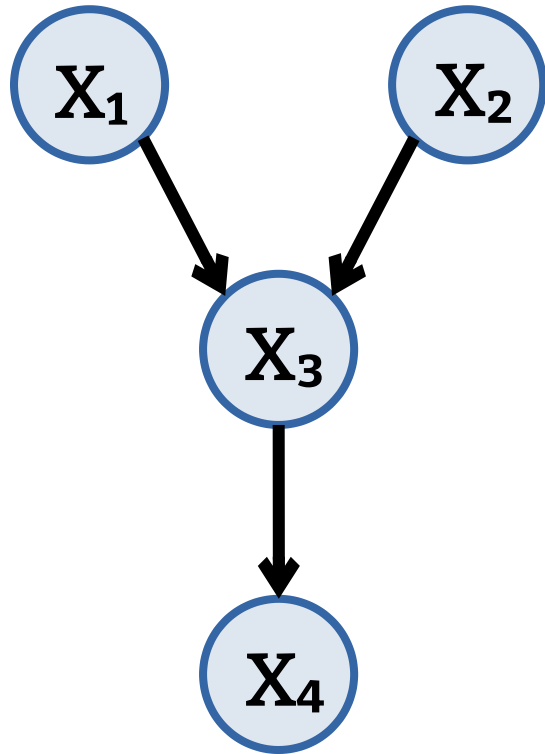
- Partial correlation (Fisher Z)
- Kernel-based tests
- $\chi^2$  test (for discrete data)

# Unobserved confounder



$$X_1 \not\perp\!\!\!\perp X_4 \mid X_3$$

# Unobserved confounder



$$X_1 \perp X_4 \mid X_3$$

**GES**

# **GES (Greedy Equivalent Search)**

- **Start with an empty graph**
- **Greedily add the edges whose addition increases the score the most**
- **Greedily remove the edges whose removal increase the score the most**
- **Score: BIC, when forecasting a variable from its parents**

**LiNGAM**



# LinGAM

The structural causal model (SCM)

$$X_1 = \mu_1 + \varepsilon_1$$

$$X_2 = \mu_2 + a_{21} X_1 + \varepsilon_2$$

$$\vdots$$

$$X_k = \mu_k + a_{k1} + \dots + a_{kk-1} X_{k-1} + \varepsilon_k$$

can be written  $X = \mu + AX + \varepsilon$ , or  $\varepsilon = (I - A)X - \mu$ , with  $\forall i \neq j \ \varepsilon_i \perp \varepsilon_j$

If the  $\varepsilon_i$ 's are non-Gaussian, ICA (independent component analysis) can recover the linear transformation  $(I - A)X$  by looking for the directions in which the data is the least Gaussian.

**NOTEARS**

# NOTES

Find a matrix  $W$  such that  $WX \approx X$  and the corresponding graph is acyclic.

The acyclicity condition can be written

$$\text{trace exp } |W| = n$$

(where  $|\cdot|$  is the elementwise absolute value)

$$\exp A = I + A + \frac{1}{2!} A^2 + \frac{1}{3!} A^3 + \dots$$

$\text{diag } A^k$  : number of cycles of length  $k$

# Code

# gCastle

```
import castle.algorithms  
import networkx as nx
```

```
model = castle.algorithms.PC()  
model.learn(X)
```

```
A = model.causal_matrix  
A = pd.DataFrame( A, columns = X.columns, index = X.columns )  
g = nx.from_pandas_adjacency( A, create_using = nx.DiGraph )
```

# causal-learn

```
from causallearn.search.ConstraintBased.PC import pc
from causallearn.utils.cit import fisherz, kci, chisq, gsq
```

```
# Computation
g = pc(X.values)
```

```
# Extract the adjacency matrix
A = pd.DataFrame( g.G.graph )
A = ( A == -1 ).astype(int)
A.columns = A.index = X.columns.copy()
```

# causal-learn

```
from causallearn.search.ScoreBased.GES import ges
```

```
r = ges(X.values)
```

```
A = pd.DataFrame( r['G'].graph )
```

```
A = ( A == -1 ).astype(int)
```

```
A.columns = A.index = X.columns.copy()
```

# causal-learn

```
from causallearn.search.FCMBased.lingam import ICALiNGAM
```

```
model = ICALiNGAM()
```

```
model.fit(X)
```

```
A = model.adjacency_matrix_
```

```
A = pd.DataFrame( A != 0, index = X.columns, columns = X.columns )
```



# cdt

```
import cdt                                # Causal discovery toolbox
import networkx as nx

pc = cdt.causality.graph.PC()             # Continuous, Gaussian variables
g = pc.predict(d)

A = nx.adjacency_matrix(g).todense()
A = pd.DataFrame( X, index = g.nodes, columns = g.nodes )
```

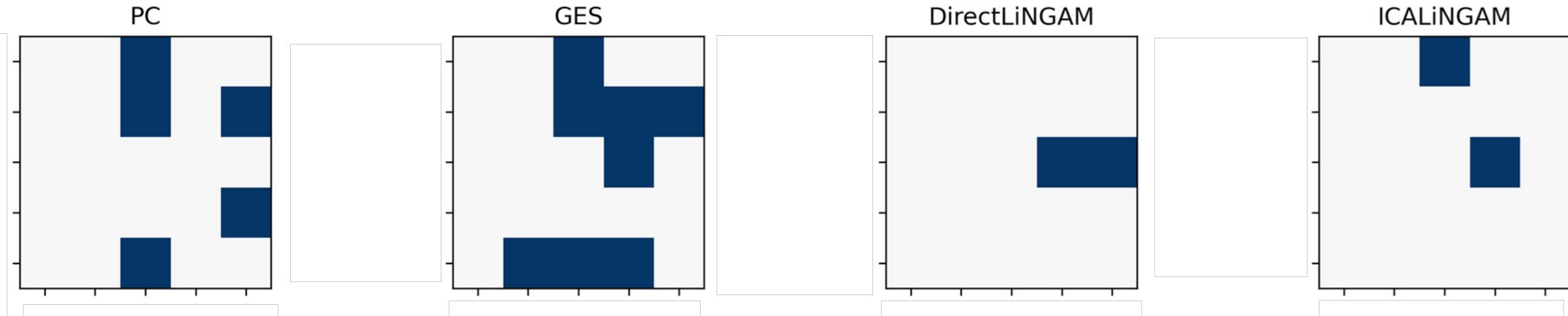
# dowhy

```
from dowhy import CausalModel
model = CausalModel(
    data = X,
    treatment = "T",
    outcome = "Y",
    graph = ' '.join( nx.generate_gml(g) ),
)
estimand = model.identify_effect()
estimate = model.estimate_effect(
    estimand,
    method_name = "backdoor.linear_regression",
)
model.refute_estimate(
    estimand, estimate,
    method_name = "random_common_cause",
)
```

# Examples

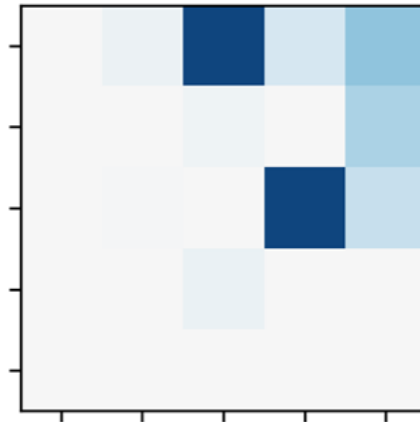
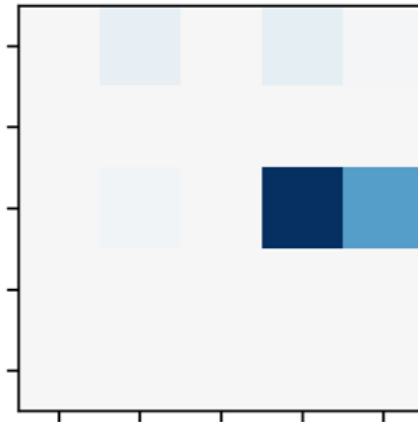
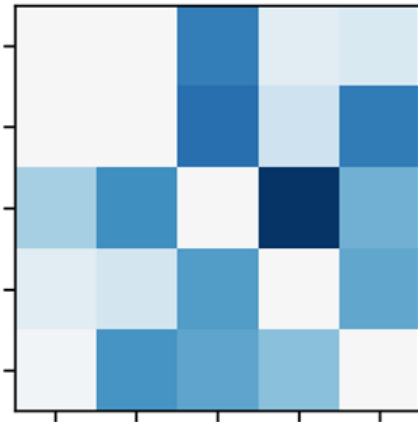
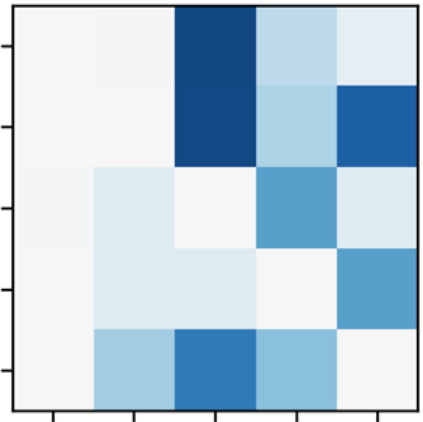
# Examples

- Different algorithms give very different results



# Examples

- Different bootstrap samples give very different results



# Conclusion

# Summary

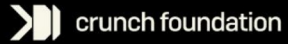
- **PC**: conditional independence test
- **GES**: goodness-of-fit of models predicting  $y$  from its parents
- **LiNGAM**: independent component analysis
- **NOTEARS**: optimization problem, with acyclicity constraint
- **Software**: gCastle, causal-learn, cdt

# Conclusion

- Do not start from data, but from a domain knowledge causal graph
- Do not use a single causal discovery algorithm, but several
- Do not run them on just the data, but also on bootstrap samples
- Do not only look at the output, look at the ingredients of the algorithms



# Causality competition



ADIA : Lab

COMPETITION GRANTS ABOUT CRUNCH TESTIMONIALS

JOIN NOW

GRANTS

\$100,000

## HOW **CAUSAL DISCOVERY** CAN HELP **RESOLVE CRITICAL CHALLENGES** IN AI RESEARCH?

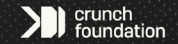
Join the latest ADIA Lab x Crunch Foundation competition to build a groundbreaking AI framework. Get a chance to win up to \$100,000.

JOIN THE CRUNCH →



ADIA : Lab

**"TRUTH IS  
COMPLETELY OUT  
OF CONTROL IF  
YOU DON'T  
UNDERSTAND  
WHY"**



ADIA : Lab

**"WHY IS  
WHAT  
SEPARATE  
YOU FROM  
THEM"**



ADIA : Lab

**"LUCKY  
IS WHO HAS  
BEEN ABLE TO  
UNDERSTAND  
THE CAUSE OF  
THINGS"**

# References

*Introduction to Causal Inference* (B. Neal, 2020)

*A survey on causal discovery: theory and practice* (A. Zanga and F. Stella, 2023)

<https://github.com/huawei-noah/trustworthyAI/tree/master/gcastle>

<https://causal-learn.readthedocs.io/en/latest/>

<https://cran.r-project.org/web/views/CausalInference.html#dag>

# **Extra Slides**

# Kernel methods

- Many machine learning algorithms do not really require coordinates, but just the Gram matrix  $K_{ij} = \langle x_i, x_j \rangle$ .
- Increasing the dimension,  $K_{ij} = \langle \phi(x_i), \phi(x_i) \rangle$  does not change the size of the Gram matrix.
- We do not even need to compute  $\phi(x_i)$ : we just need a (positive definite) kernel,  $\kappa(x, y) = \langle \phi(x), \phi(y) \rangle$ .

# Local BIC score of $X \rightarrow y$

$$H = \log( \Sigma[y,y] - \Sigma[y,X] \Sigma[X,X]^{-1} \Sigma[X,y] )$$

$$-\text{BIC} = n \cdot H + \lambda \cdot k \cdot \log(n)$$

$n$  = number of observations

$k$  = number of variables in  $X$

$\Sigma$  = variance matrix of the data

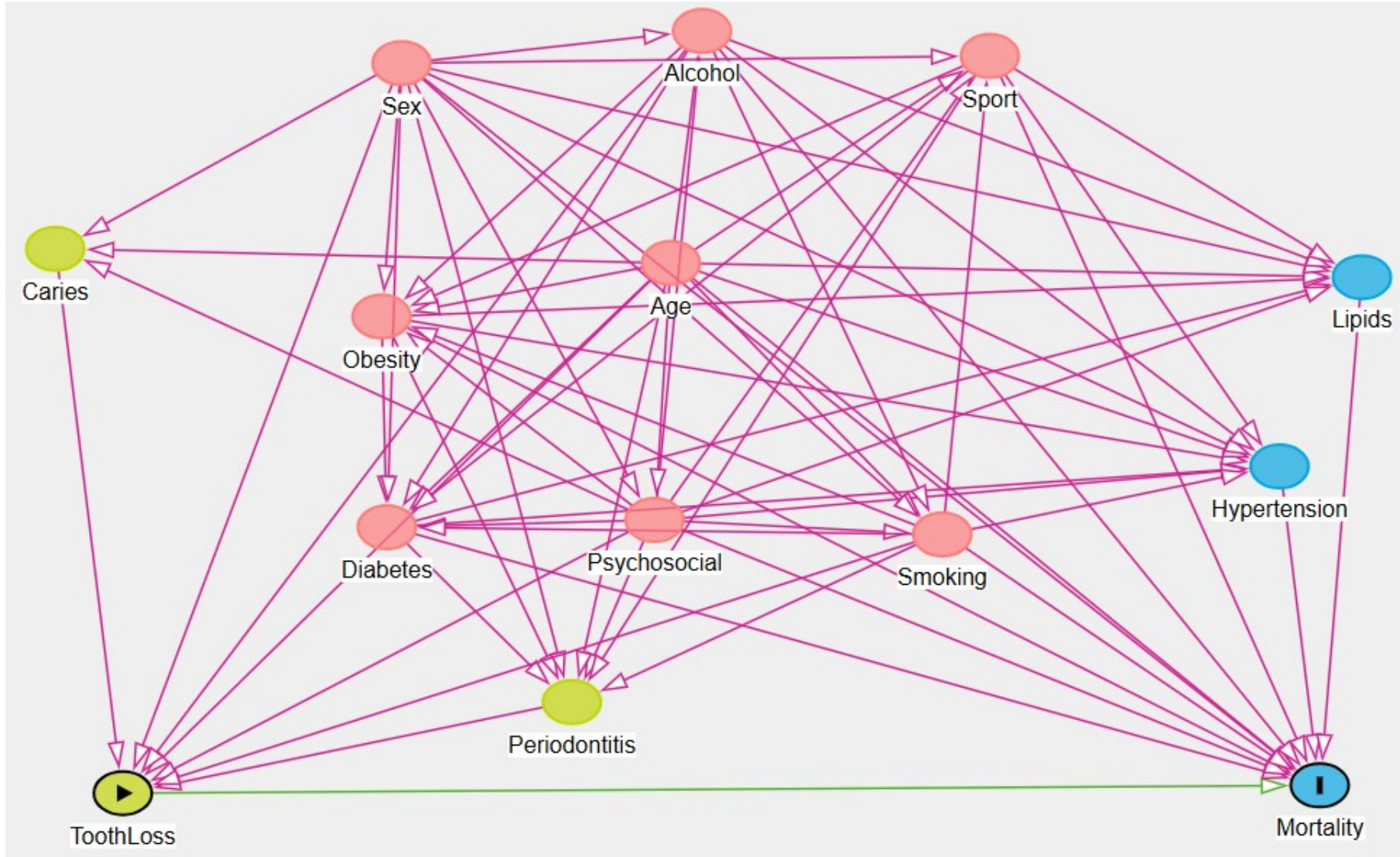
# Markov equivalence class

A **CPDAG** (complete partially directed acyclic graph) is a PDAG where

- All undirected edges are **reversible** (you can choose either direction, that does not change the MEC)
- All directed edges are **compelled** (if you flip them, the graph moves to another MEC)

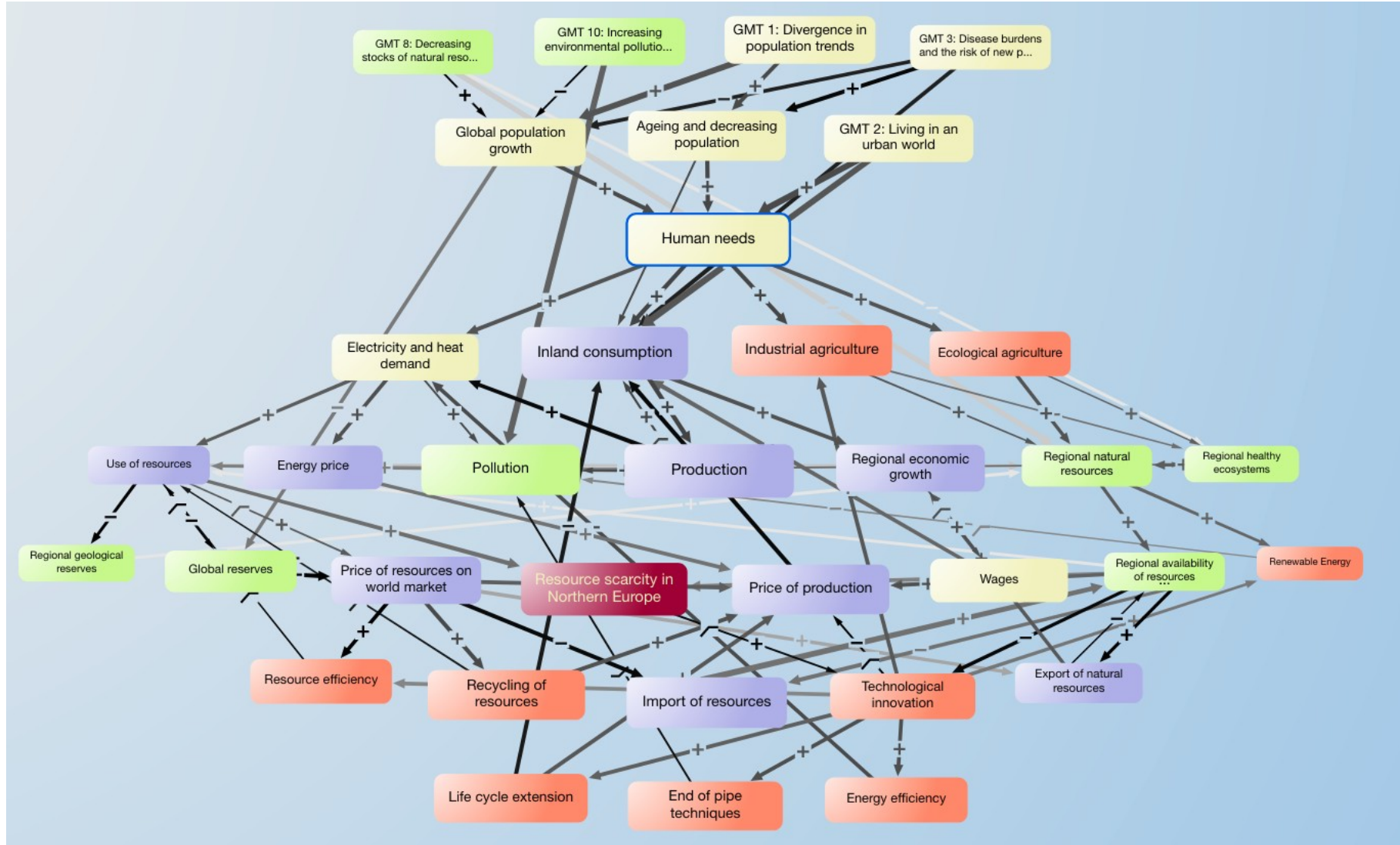
# Unused Slides

# Real-world causal graphs

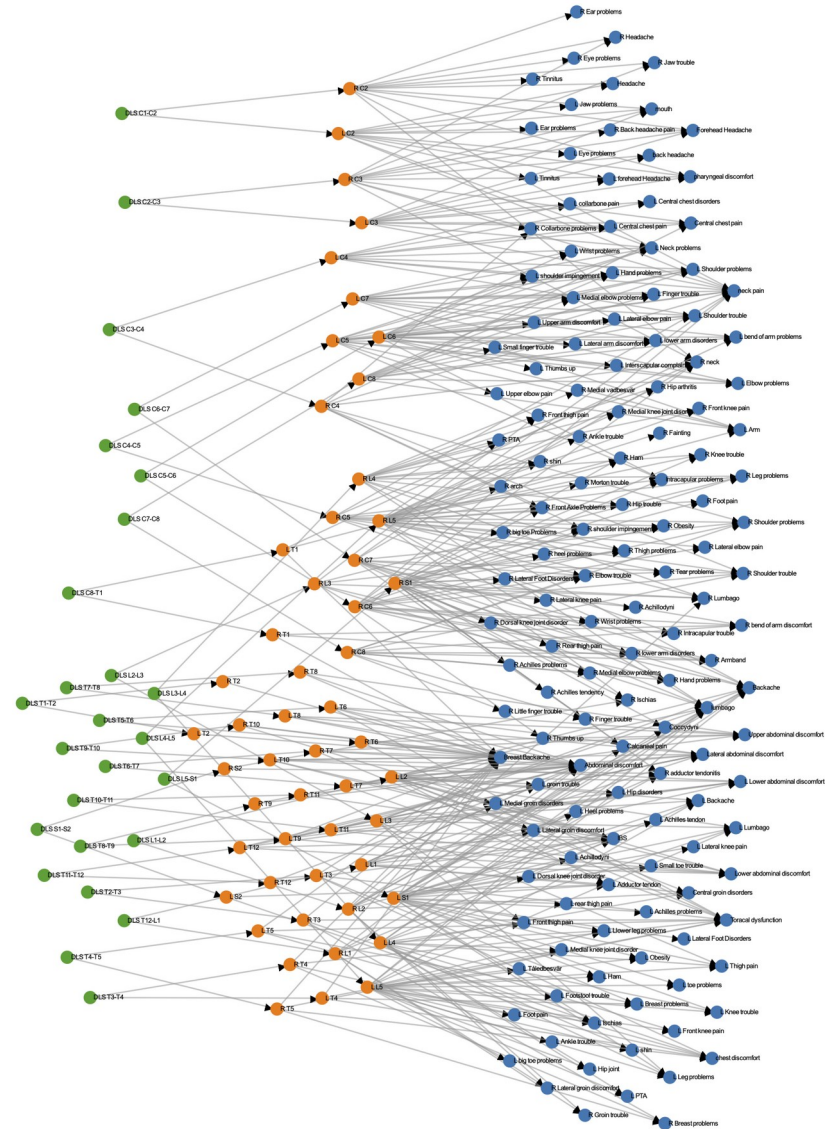




# Real-world causal graphs



# Real-world causal graphs

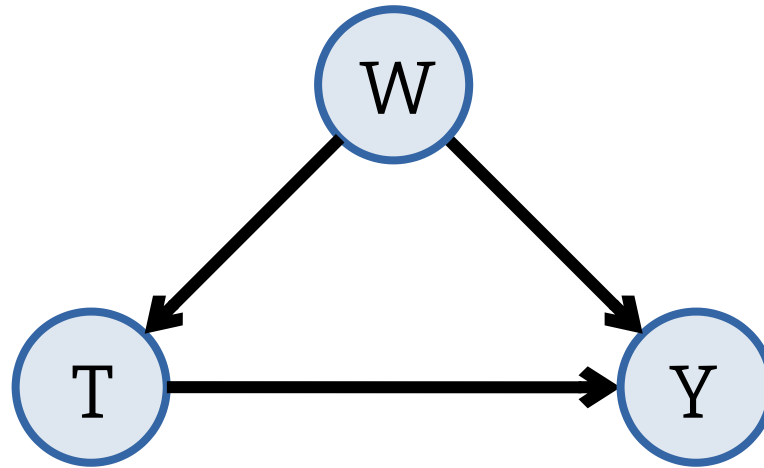


<https://arxiv.org/abs/1906.01732>

# **Variants of those algorithms**

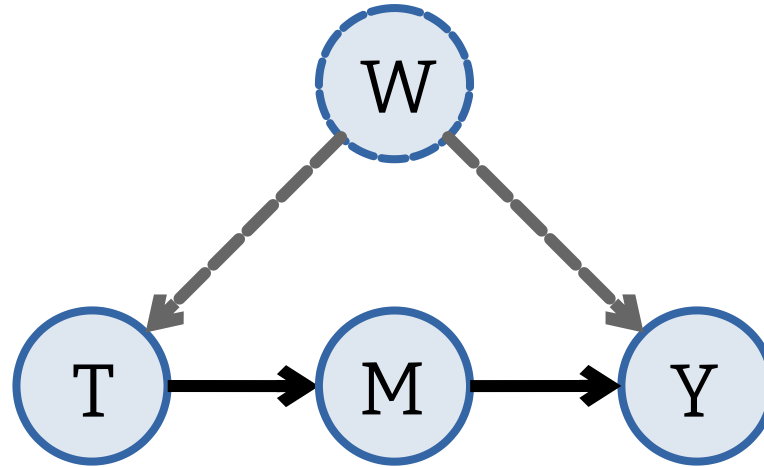
- **FCI: PC variant, allowing for unobserved confounders**
- **FGES: faster implementation of GES**
- **ARGES: another GES variant, for high-dimensional data**
- **GFCI: GES variant allowing for unobserved confounders (FCI on the FGES skeleton)**
- **CCD: PC/FCI with feedback (cycles)**
- **LiNG: LiNGAM with feedback**
- **Other LiNGAM variants: non-linear, post-non-linear**
- **CD-NOD**

# Back-door adjustment



$$P[Y|\text{do}(T = t)] = \sum_w P[Y, t, w]P[w]$$

# Front-door adjustment

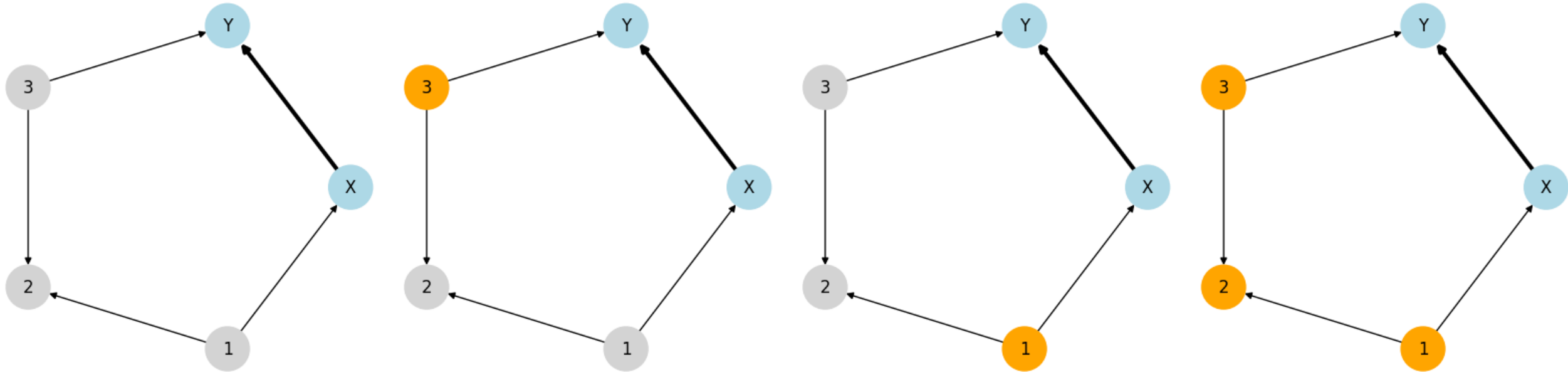


$$P[Y|\text{do}(t)] = \sum_m P[m|t] \sum_{t'} P[y|m, t'] P[t']$$

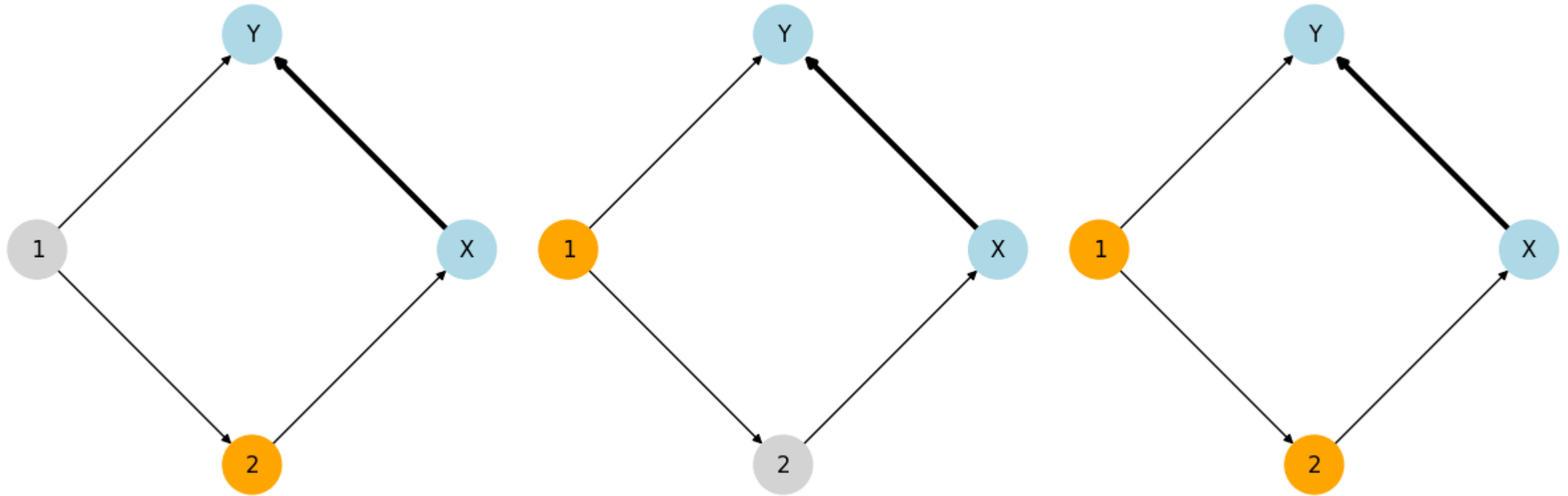
# Causal inference

- The sufficient conditioning set is not unique.
- The parents of  $X$  form a sufficient conditioning set, but it may be needlessly large.
- More generally, a sufficient conditioning set is a set of nodes blocking all the non-causal paths from  $X$  to  $Y$ , and leaving all the causal paths open.
- If some variables are not observed, things get more complicated, but do-calculus can tell you if the effect can be estimated from observational data alone.

# Sufficient Conditioning Sets

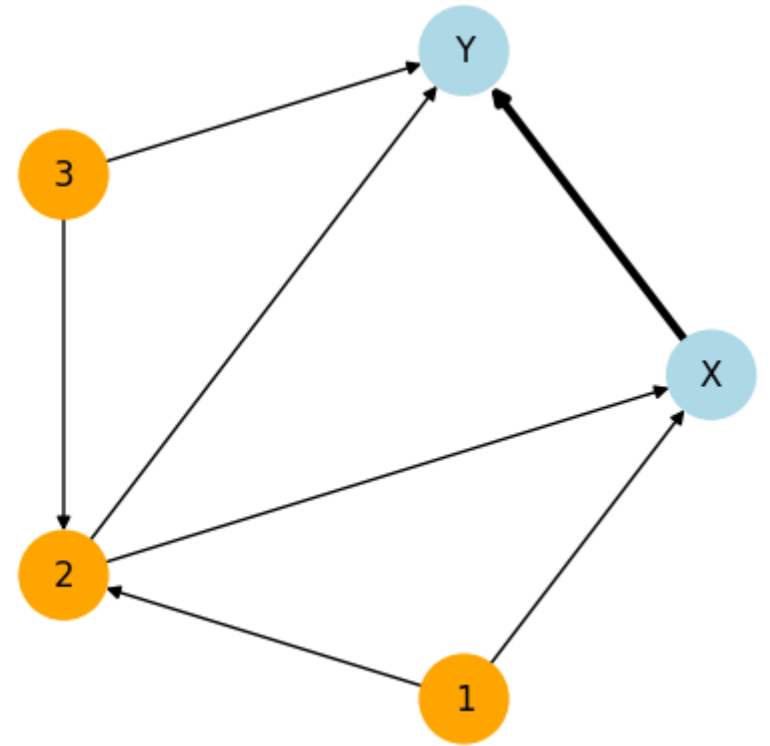
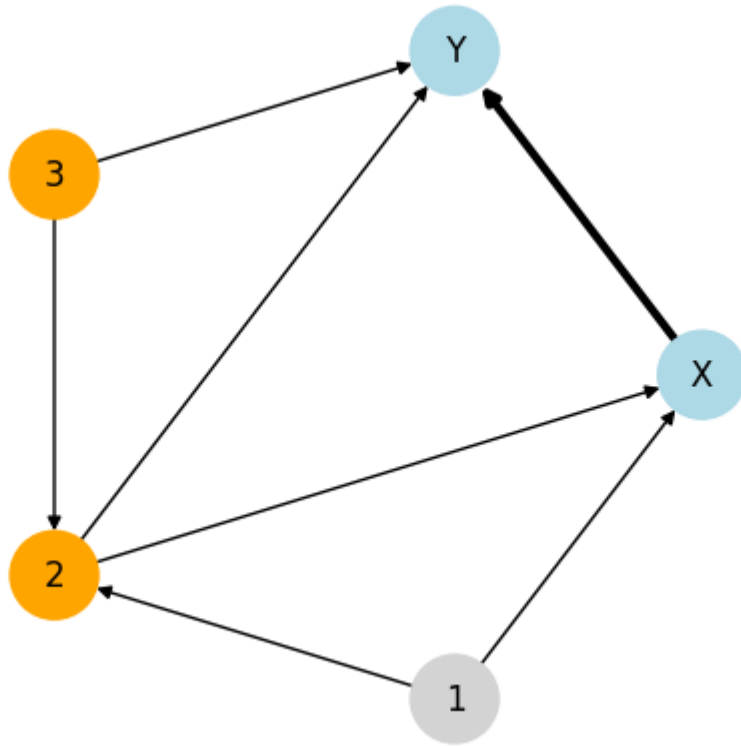
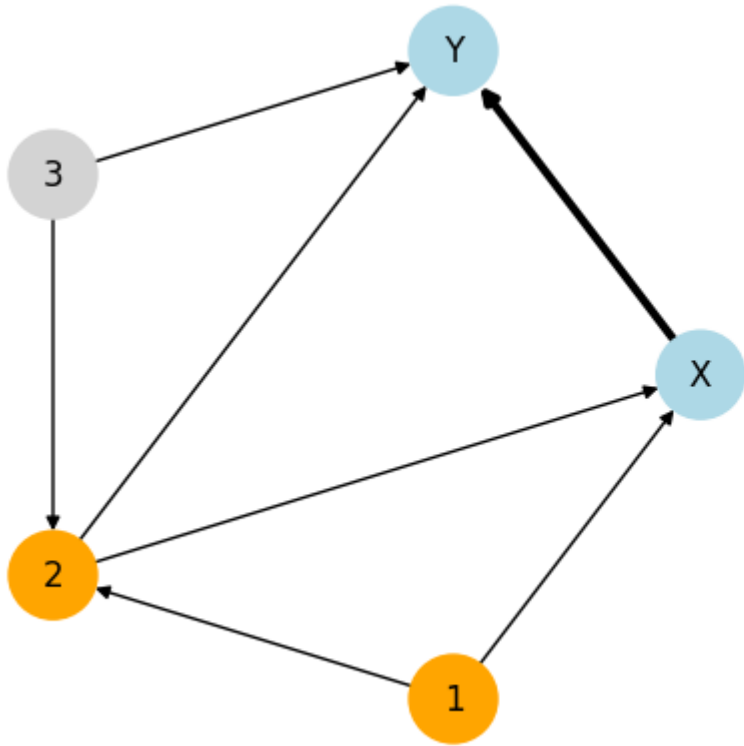


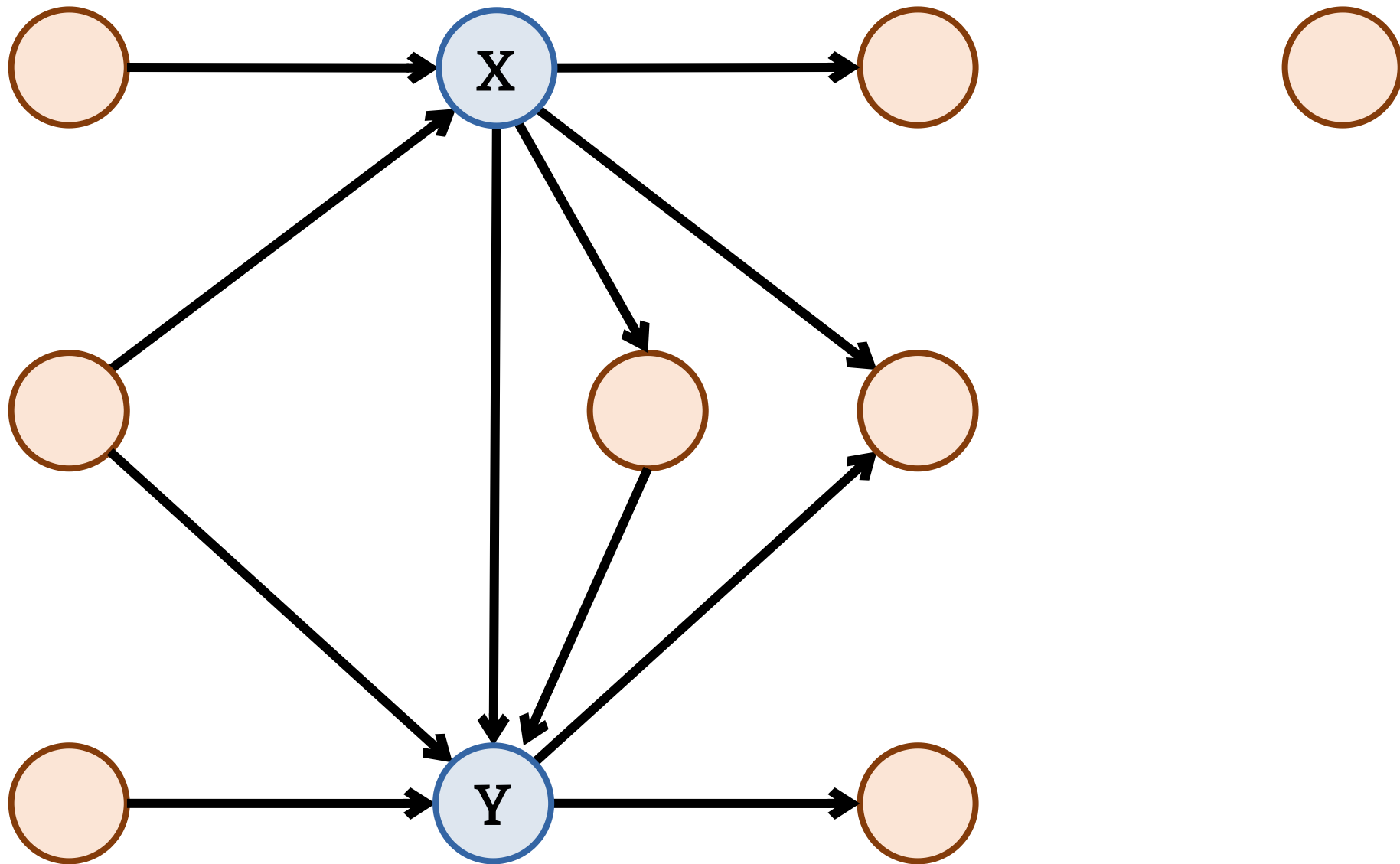
# Sufficient Conditioning Sets



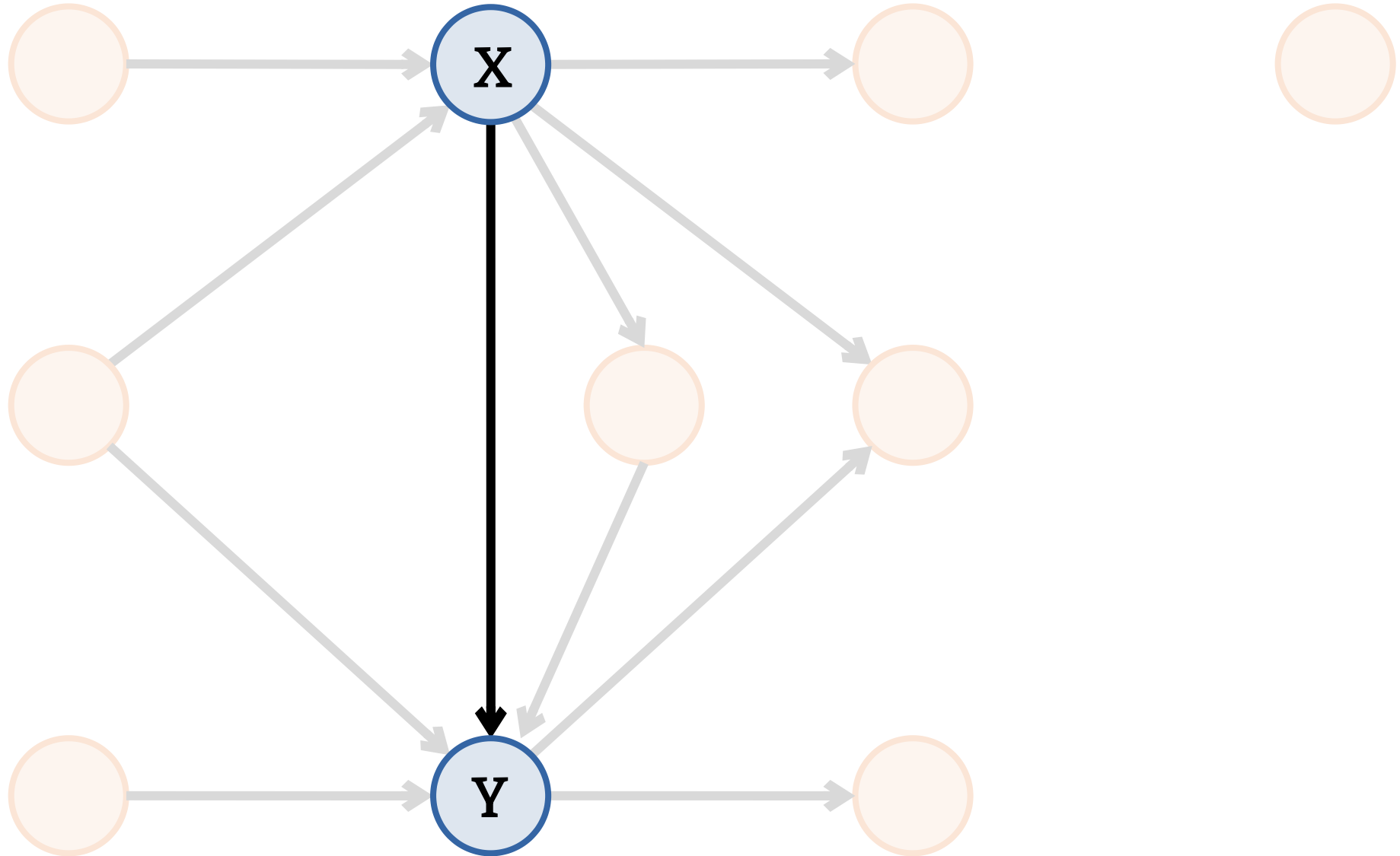


# Sufficient Conditioning Sets

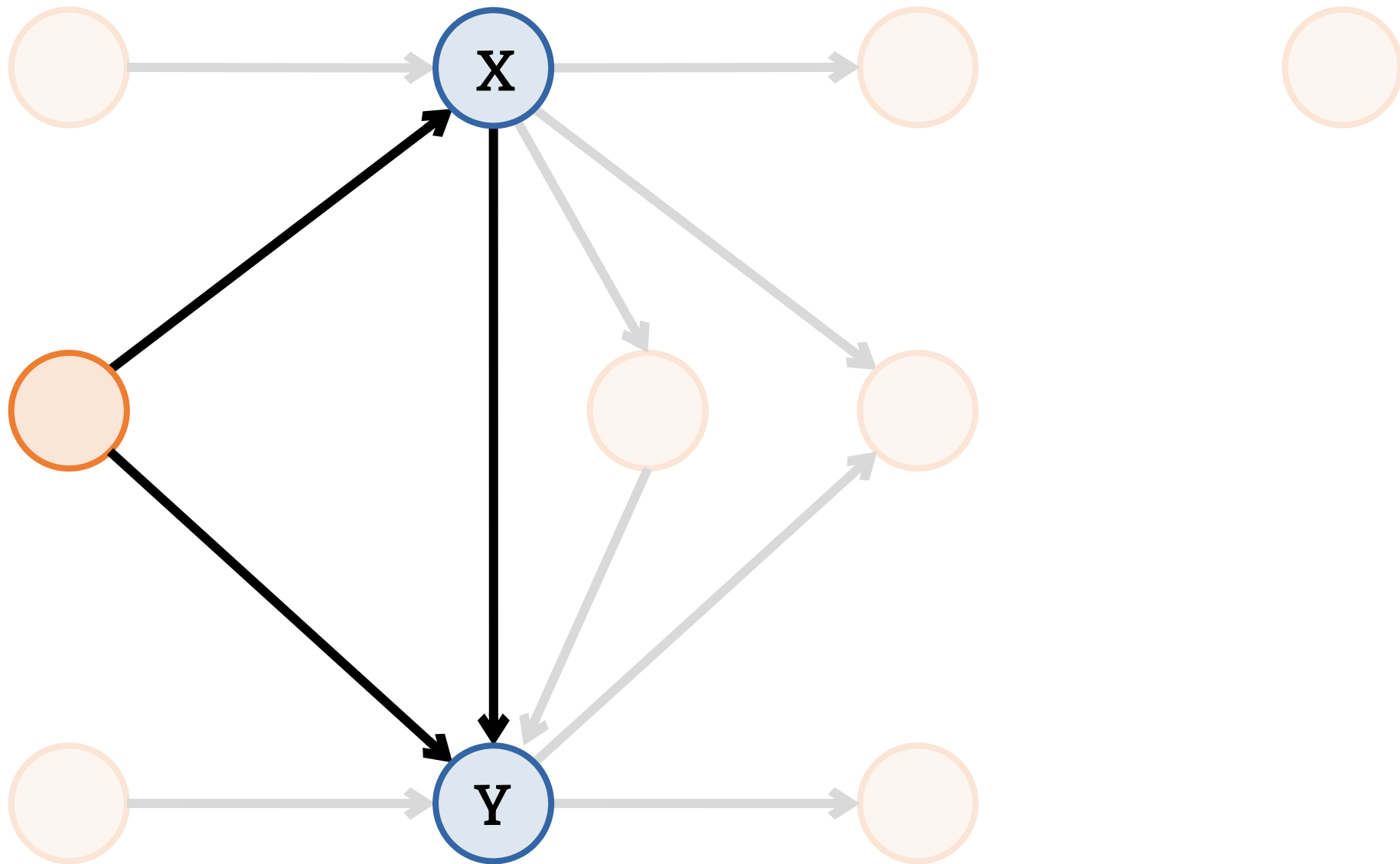




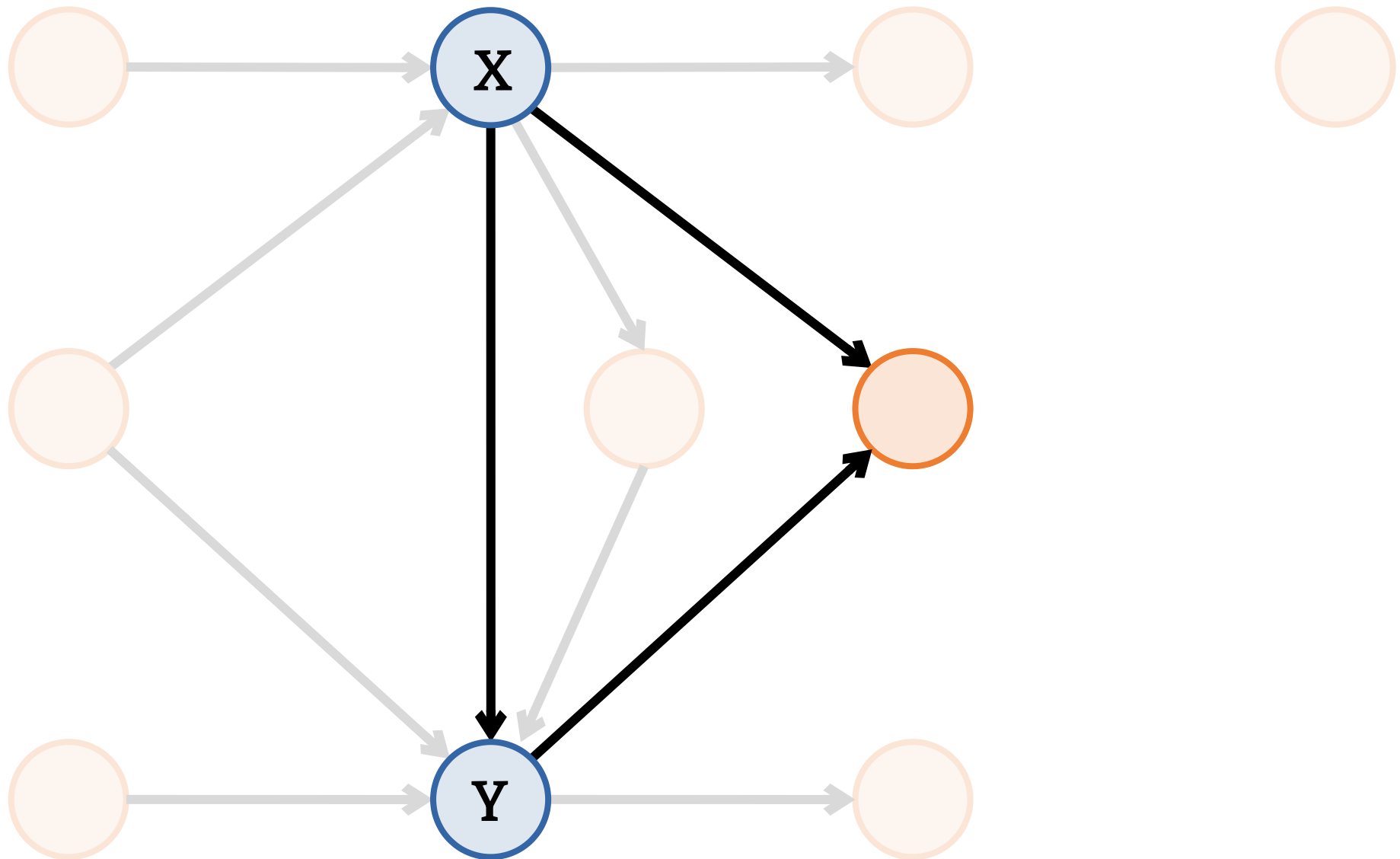
# Relation of interest



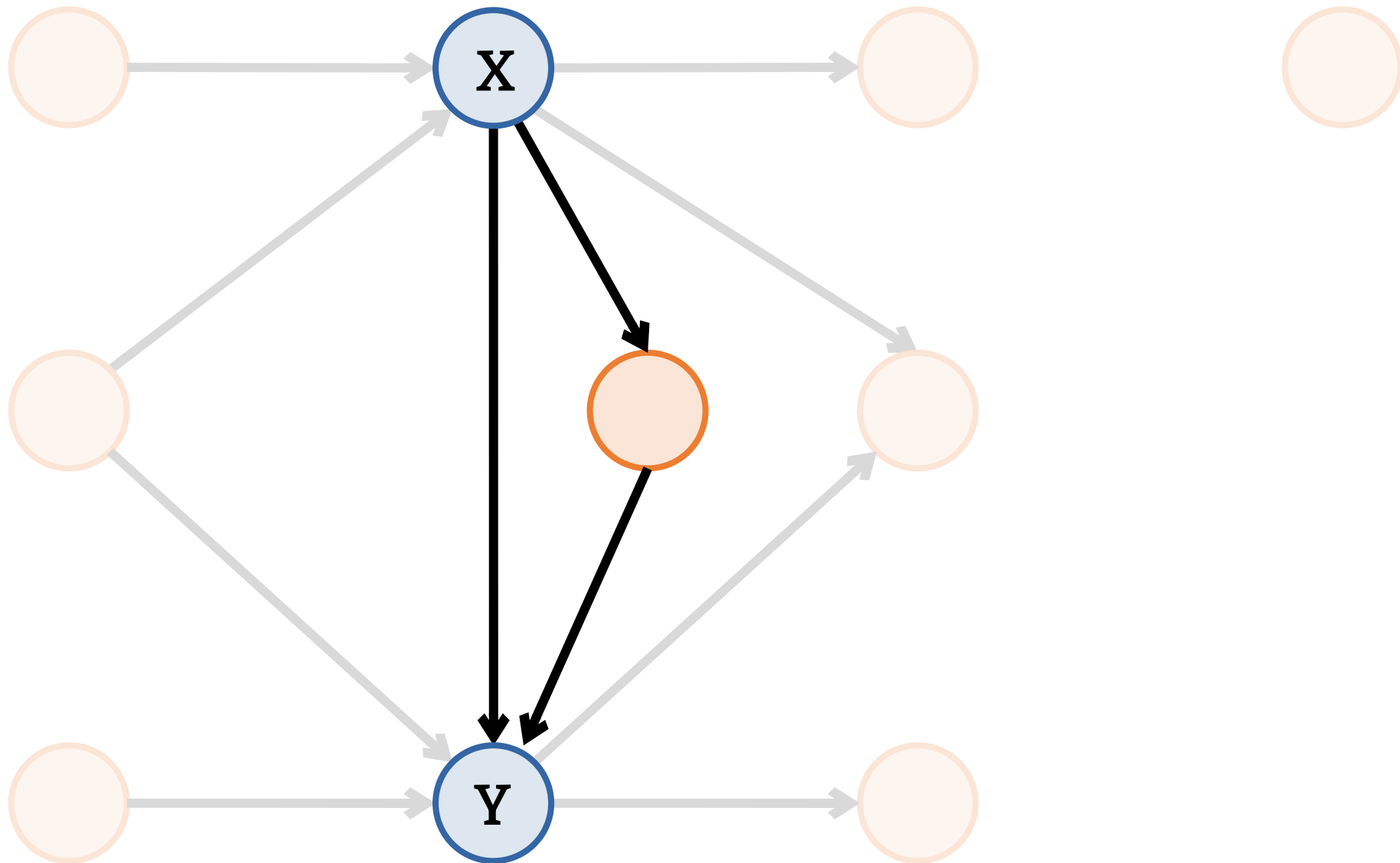
# Confounder



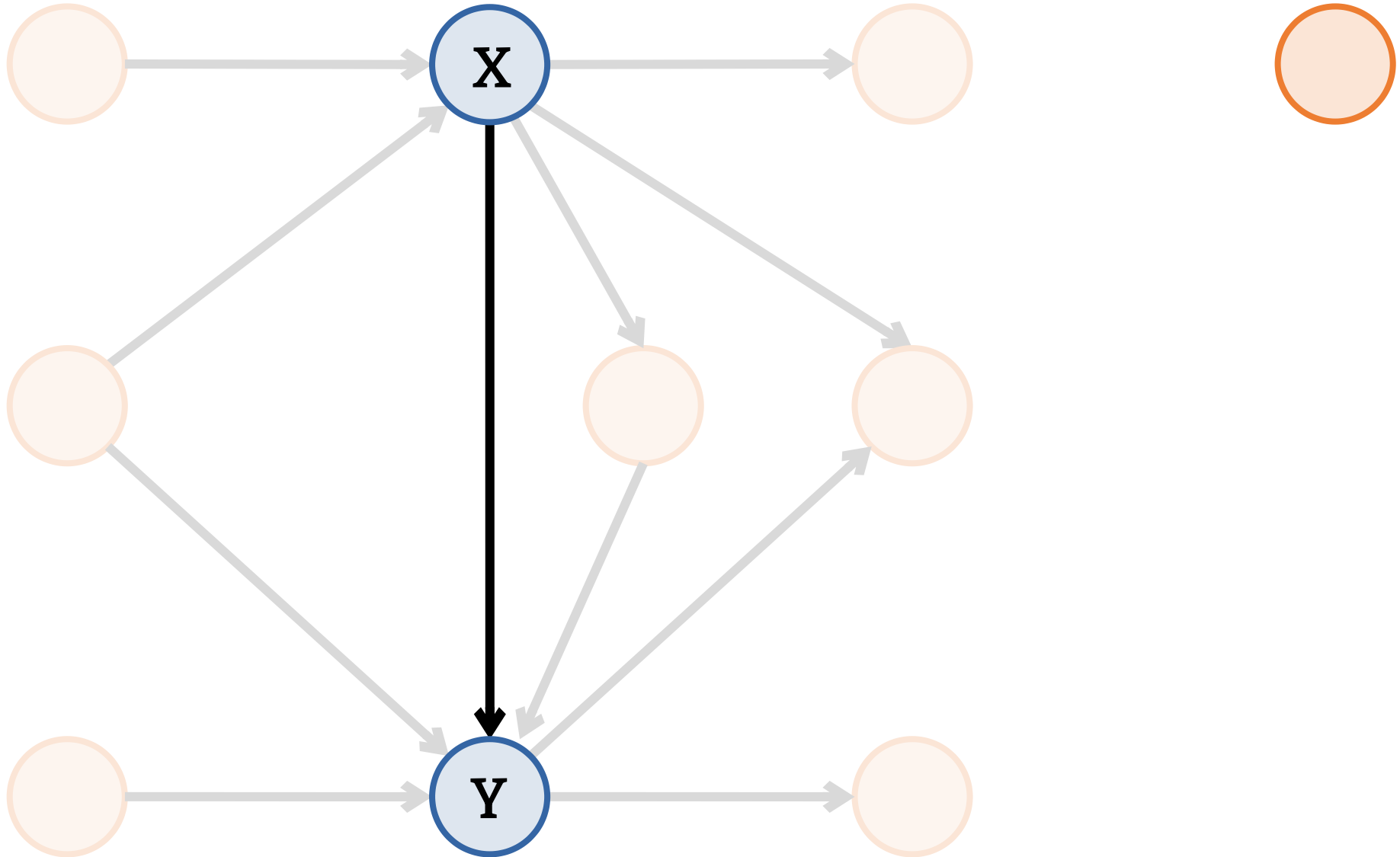
# Collider



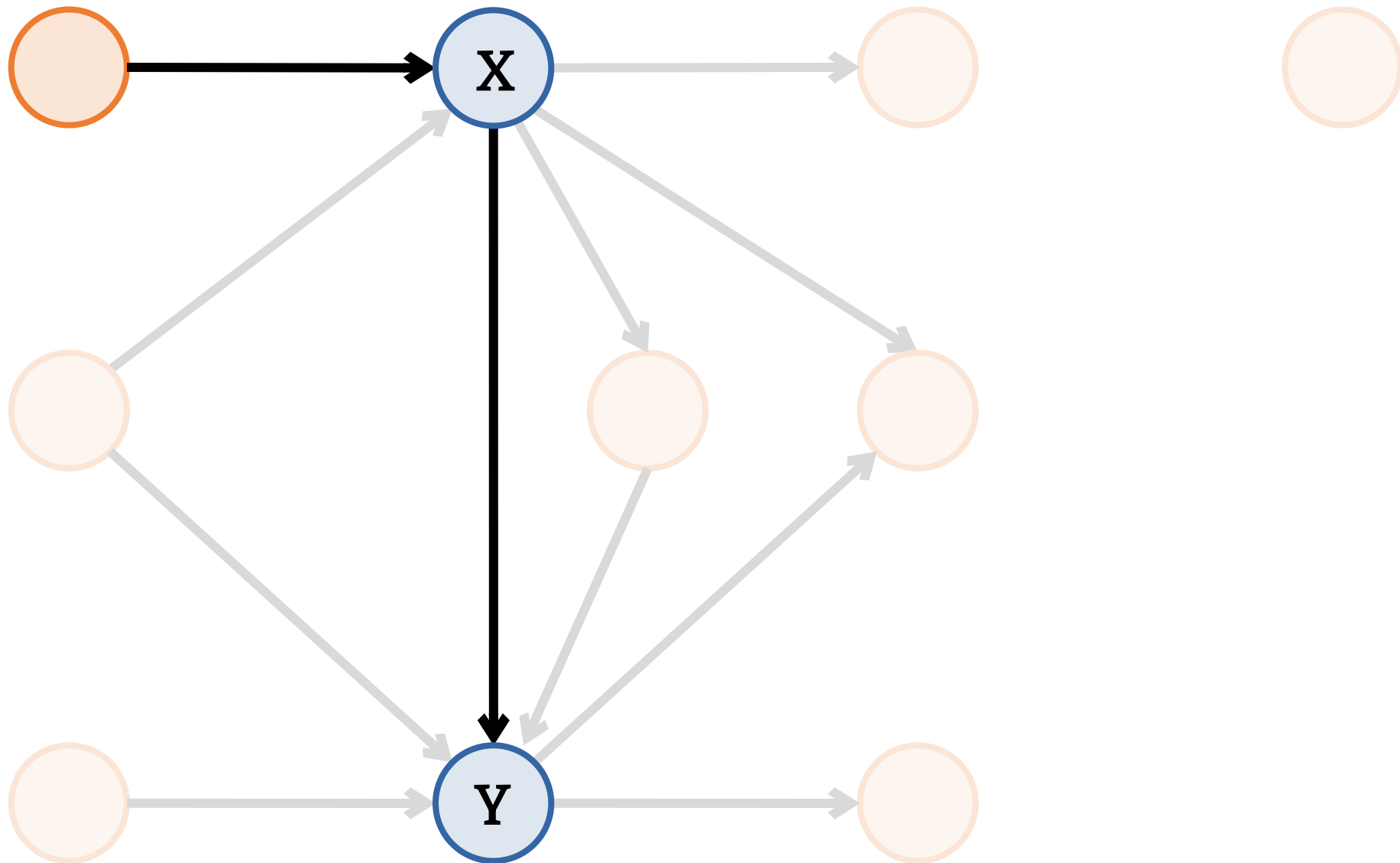
# Mediator



# Independent

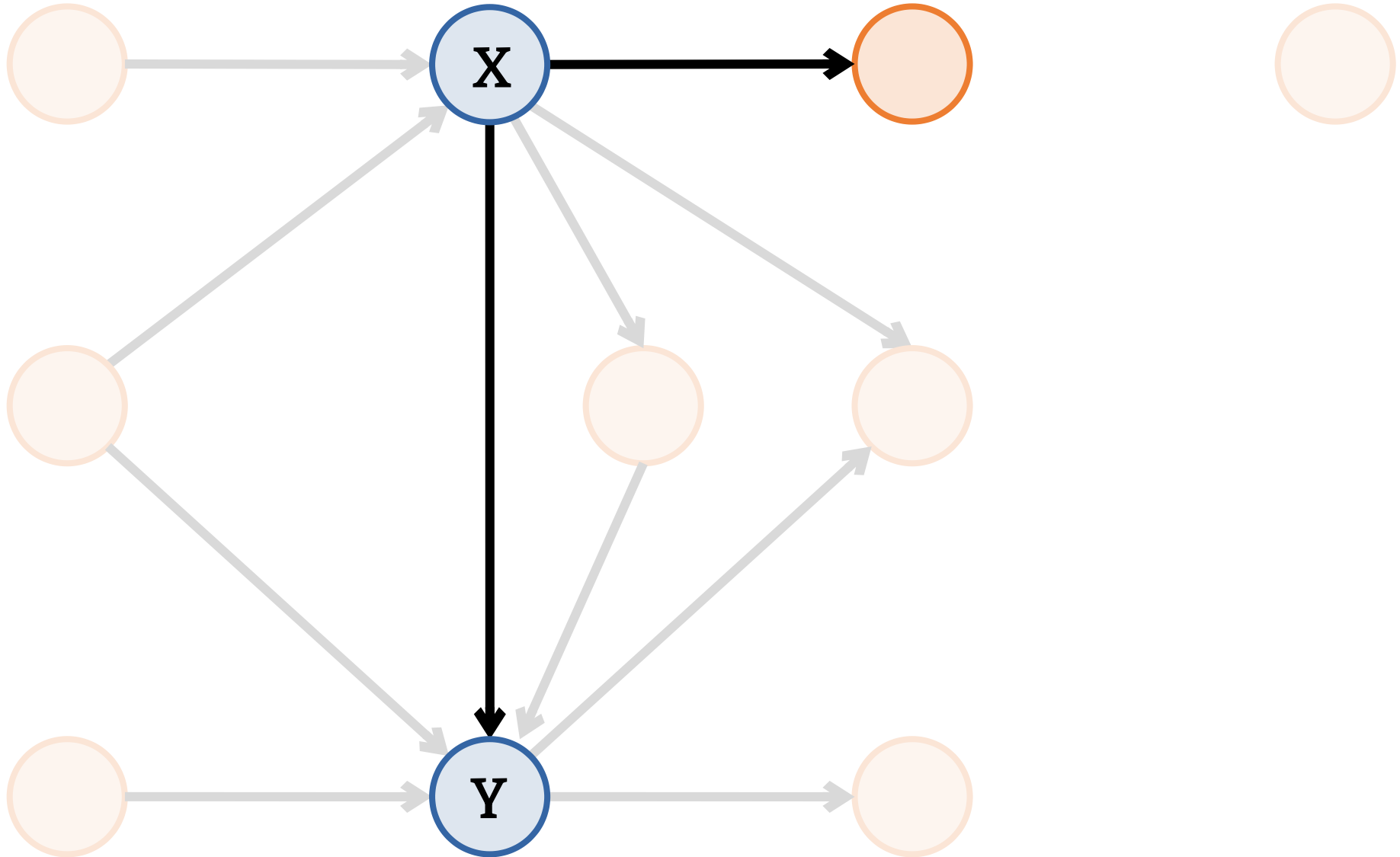


# Cause of X

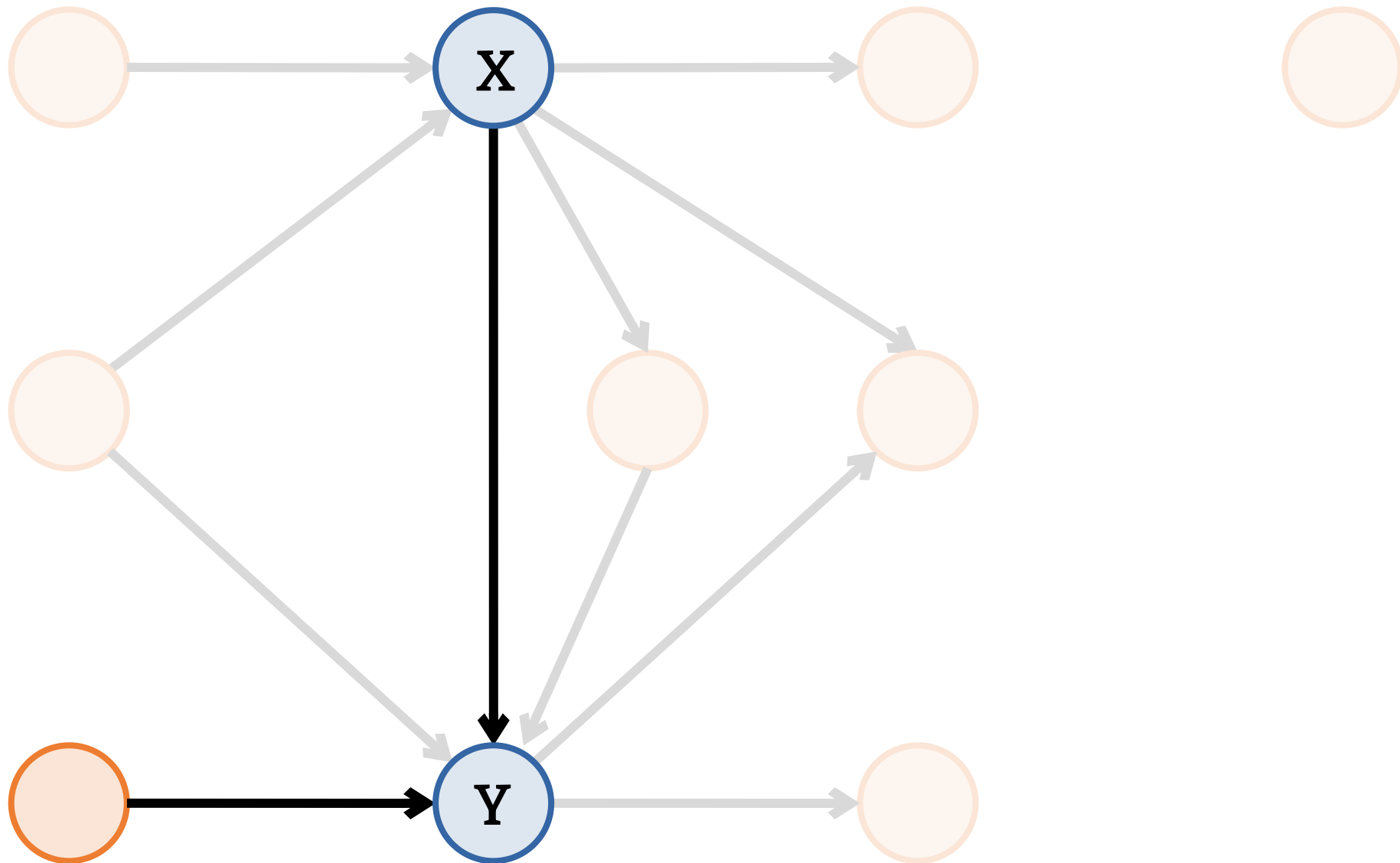




# Consequence of X



# Cause of Y



# Consequence of Y

