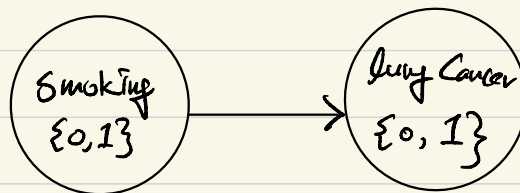


## \*Causal Structure Discovery.

- Does smoking cause lung cancer?

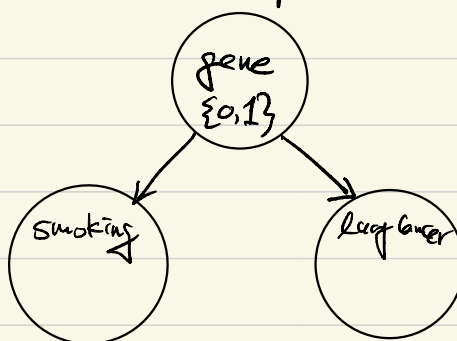
### Observational Data

		Lung Cancer	
		Yes	No
Smoking	Yes	15%	85%
	No	6%	94%



Correlation does not imply causality.

Alternative Explanation



people with specific gene is likely to smoke AND get lung cancer.

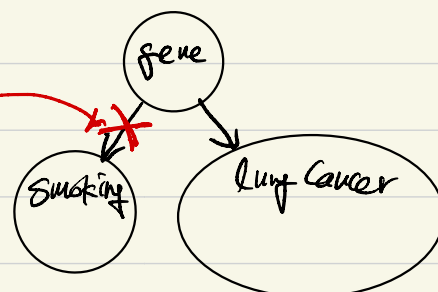
- these people would have gotten lung cancer even if they did not smoke.
- hence, smoking does not cause cancer.

### Interventional Data.

Randomized trials.

random 50% of population set smoking = 1

random 50% of population set smoking = 0.



you can identify causality by intervention, but it can be unethical and/or expensive.

\* All nodes are observed with observational data.

Recall: BN  $G=(V,E)$  is a DAG with  $P(x) = \prod_{i=1}^n P(x_i | \pi_i)$ .

Goal of Causal structure learning is to recover  $G$ .

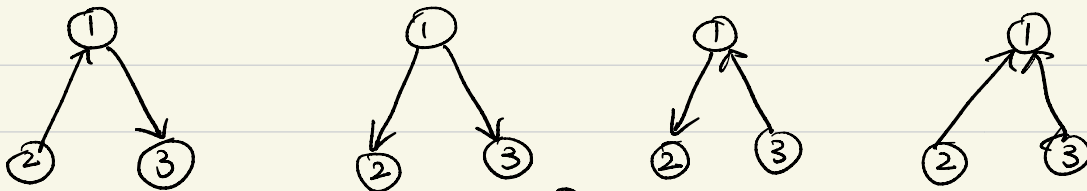
Def. **Markov Equivalence Class (MEC)**

$$G_1 \sim G_2 \iff I(G_1) = I(G_2)$$

$\left[ \begin{array}{l} \text{skeleton is the same} \\ \text{moral graph is the same} \end{array} \right.$

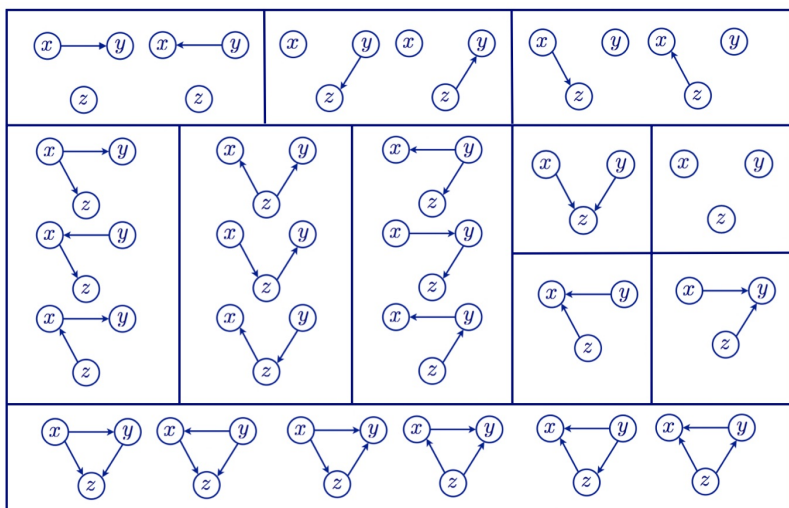
Claim: From observational data, we can only recover  $G$  up to its MEC.

proof:  $G_1 \sim G_2$  implies  $\forall P(x)$  that factorizes as  $G_1$  also factorizes as  $G_2$ .



Which ones are equivalent?

MEC on 3node graphs



\* Hence,  $G$  can only be partially identified

\* To resolve the direction of edges within MEC, we need to use interventional data.

# Constraint-based Algorithm

checks conditional independencies

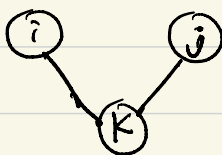
[SGS-Algorithm]

Spirites-Glymour-Scheines 2001

Step 1. start with a complete undirected  $G = (V, E)$

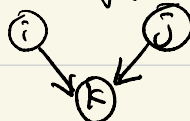
Step 2. using given observational data, for all  $(i, j) \in V \times V$   
remove  $(i, j)$  from  $E$  if  $\exists S \subseteq V$  s.t.  $X_i \perp\!\!\!\perp X_j \mid X_S$ .

Step 3. for all triplet  $(i, j, k)$  s.t.



Check if  $X_i \perp\!\!\!\perp X_j \mid X_{rest \setminus \{k\}}$

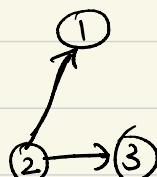
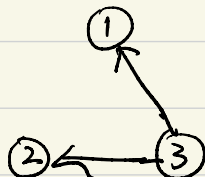
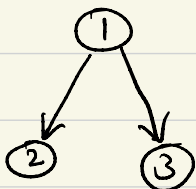
If yes, direct edges as



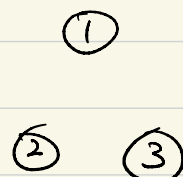
Step 4. Orient remaining undirected edges by consistency and recursively do this until no more can be oriented.

Q. When does SGS algorithm fail to recover the MEC?

$$\text{ex) } I(P(X_1, X_2, X_3)) = \{ X_1 \perp\!\!\!\perp X_2 \mid X_3, X_2 \perp\!\!\!\perp X_3 \mid X_1, X_1 \perp\!\!\!\perp X_3 \mid X_2 \}$$



possible ground truths,



SGS output

Recall. Global Markov Property

If  $X_i$  and  $X_j$  are d-separated by  $S$ , then

$$X_i \perp\!\!\!\perp X_j \mid X_S$$

Def.  $P(x)$  is faithful w.r.t  $G$  if

$$X_i \perp\!\!\!\perp X_j \mid X_S \text{ for } S \implies (i, j) \notin E.$$

This justifies step 2 of SGS.

Claim. If  $X^{(1)}, \dots, X^{(N)}$  iid  $P(x)$ ,

$P(x)$  is faithful to a graph  $G$ ,  
all variables in  $G$  are observed.

Then SGS is consistent, i.e.,

$$\lim_{N \rightarrow \infty} \mathbb{P}(\hat{G}_{\text{SGS}} \neq G) = 0$$

Constraint-based algorithms require a lot of samples  
↳ faithfulness assumption

## \* Score-based Algorithms

Recall: log-likelihood score of a DAG

$$\text{SCORE}(G) = N \sum_{i=1}^n \log \hat{p}(X_i; X_{\pi_i}) - N \sum_{i=1}^n H_{\hat{p}}(X_i)$$

and without further assumption on  $G$ , the complete DAG has the highest SCORE.

Def. Bayesian Information Criterion (BIC) score:

$$\text{SCORE}_{\text{BIC}}(G) = \underbrace{\text{SCORE}(G)}_{\text{log-likelihood}} - \frac{\log N}{2} \underbrace{\text{dim}(G)}_{\substack{\text{how many bits required to describe } \hat{P}_G \\ \text{Minimum Description Length (MDL) \\ principle}}}$$

$$\text{where } \text{dim}(G) = \sum_{i=1}^n (|\mathcal{A}| - 1) |\mathcal{X}_i|^{|\pi_i|}$$

- $\text{SCORE} = O(N)$ ,  $\text{DL} = O(\log N)$ ,  $\rightarrow$  Second term dominates when there is not enough samples.

properties:

① Score equivalence:  $G_1 \sim G_2 \iff \text{SCORE}_{\text{BIC}}(G_1) = \text{SCORE}_{\text{BIC}}(G_2)$

② Consistency: If  $G^*$  is a perfect map for  $P(X)$ , then as  $N \rightarrow \infty$ ,  $G^*$  is the unique maximizer of  $\text{SCORE}_{\text{BIC}}(G)$ .

③ Decomposability:  $\text{SCORE}_{\text{BIC}}(G) = \sum_{i=1}^n \text{SCORE}(X_i, X_{\pi_i})$

$\Rightarrow$  Greedy Equivalence Search (GES).

## Algorithm. [Greedy Equivalence Search]

Initialize  $G^{(0)} = (V, E = \emptyset)$

phase 1:  $t=1, \dots, T$

add an edge that maximizes  $\text{SCORE}_{\text{BIC}}(G^{(t+1)})$ .

phase 2:  $t=T+1, \dots$

remove an edge that maximizes  $\text{SCORE}_{\text{BIC}}(G^{(t+1)})$

Claim: As  $N \rightarrow \infty$ , GES greedily finds MEC under faithfulness

## \* Permutation-based Greedy Search Algorithm

Idea:

Table 1: Equivalence Class Counts

$n$	Equivalence classes	CI/ADG	CI <sub>1</sub> /CI
1	1	1.00000	1.00000
2	2	0.66667	0.50000
3	11	0.44000	0.36364
4	185	0.34070	0.31892
5	8782	0.29992	0.29788
6	1067825	0.28238	0.28667
7	312510571	0.27443	0.28068
8	212133402500	0.27068	0.27754
9	326266056291213	0.26888	0.27590
10	1118902054495975141	0.26799	0.27507

(Gillispie & Perlman, 2001)

the # of MECs for  $n$ -node graph  
explodes.



we instead search over all  
permutations (and skeletons)

→ # MEC  $\approx 10^{18}$  vs.  $10! = 3,628,800$



we apply Greedy Search.

# Greedy Search for Sparsest Permutation [GSP] Algorithm.

Initialize:  $\pi^{(1)}$  as arbitrary ordering.

Repeat:  $t=1, \dots$

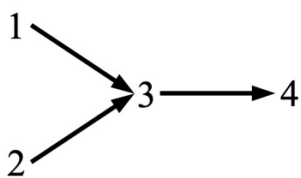
for each permutation/ordering  $\pi$  in the neighborhood of  $\pi^{(t)}$   
construct a DAG  $G_\pi$  by

$$(\pi_i, \pi_j) \in E_\pi \iff x_{\pi_i} \not\leq x_{\pi_j} \mid x_{\pi_1 \dots \pi_{i-1}, \pi_{i+1} \dots \pi_{j-1}}$$

Evaluate  $\text{SCORE}_{\text{BIG}}(G_\pi)$

$\pi^{(t+1)} \leftarrow$  the best scoring candidate permutation.

- two permutations are neighboring if they differ only in two adjacent positions  
 e.g.  $(2, 5, 3, 1, 4)$   
 $(2, 3, 5, 1, 4)$
- Claim: GSP is consistent under strictly weaker condition than faithfulness



**CI relations:**  $1 \perp\!\!\!\perp 2, 1 \perp\!\!\!\perp 4 \mid 3, 1 \perp\!\!\!\perp 4 \mid \{2, 3\}$   
 $2 \perp\!\!\!\perp 4 \mid 3, 2 \perp\!\!\!\perp 4 \mid \{1, 3\}$

