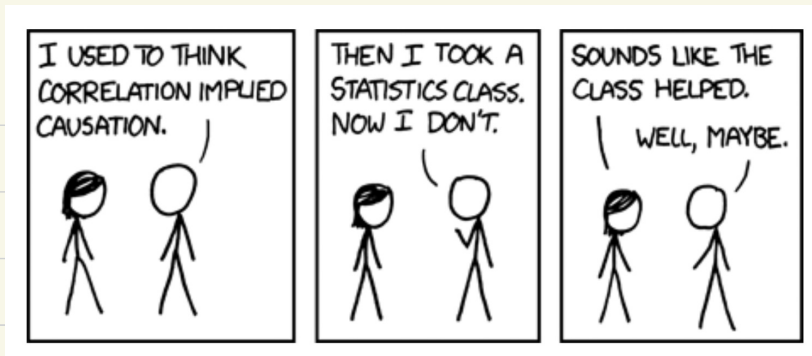


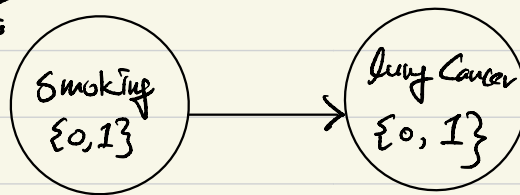
*Causal Structure Discovery.



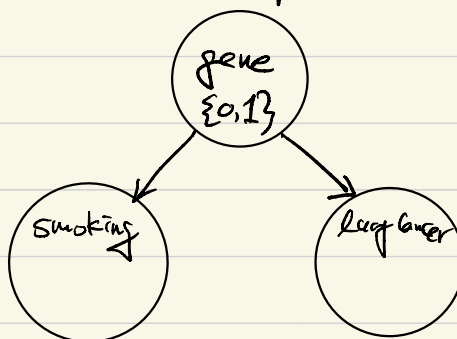
• Does smoking cause lung cancer?

Observational Data
Lung Cancer

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



Correlation does not imply causality.
Alternative Explanation



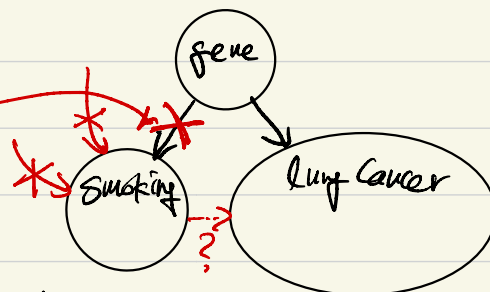
people with specific gene is likely to smoke AND get lung cancer.
 → those people would have gotten lung cancer even if they did not smoke.
 → hence, smoking does not cause cancer.

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.

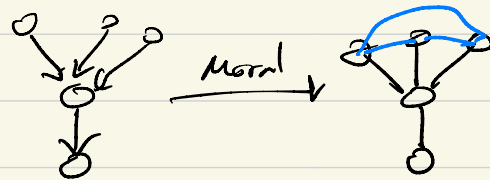
* All nodes are observed with Observational Data.

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

Def. Markov Equivalence Class (MEC)

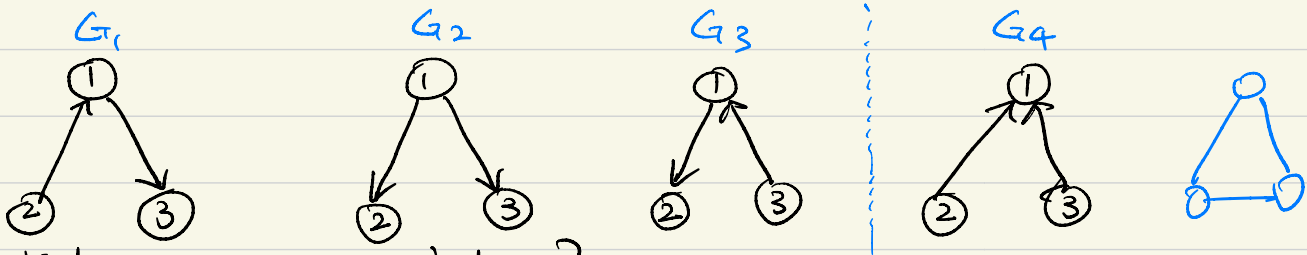
$$G_1 \sim G_2 \iff I(G_1) = I(G_2)$$

Skeleton is the same
Moral Graph is the same.



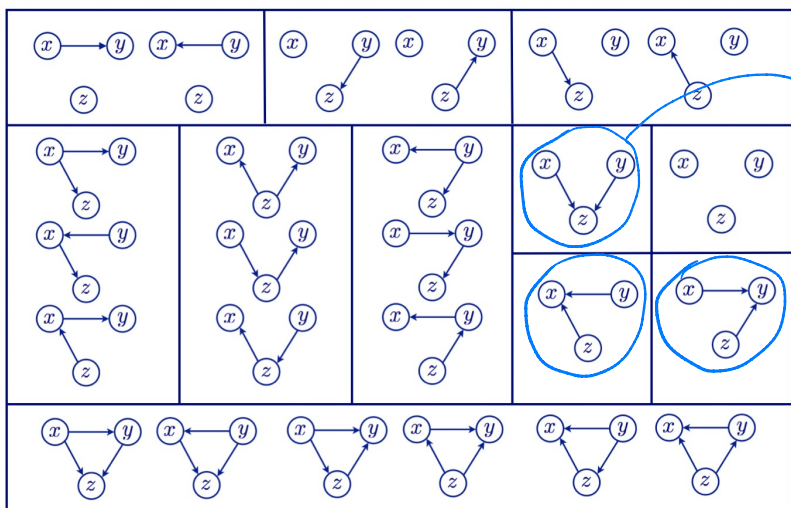
Claim: From observational data, we can only recover G up to its MEC.

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



Skeleton,
Moral.

$x \perp\!\!\!\perp y, x \not\perp\!\!\!\perp z, y \not\perp\!\!\!\perp z$
 $x \not\perp\!\!\!\perp y \mid z, y \not\perp\!\!\!\perp x \mid z, x \not\perp\!\!\!\perp z \mid y$

* Hence, G can only be partially identified

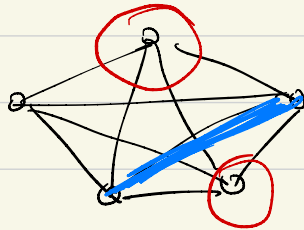
* To resolve the direction of edges within MEC, we need to use interventional data.

Constraint-based Algorithm [SGS-Algorithm]

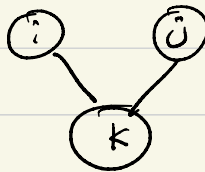
Sprites-Glymour-Scheines 2001.

Step 1. Start with Complete Graph $G=(V,E)$. undirected graph

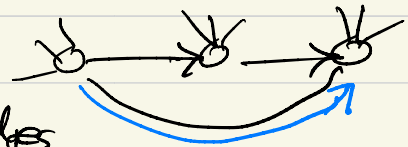
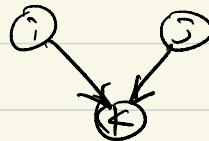
Step 2. Using observational data, for all $(i,j) \in V \times V$
remove (i,j) from E if $\exists S$ s.t. $X_i \perp\!\!\!\perp X_j \mid X_S$.



Step 3. for all triplets $(i,j,k) \in V \times V \times V$ s.t.



check if $X_i \perp\!\!\!\perp X_j \mid X_{rest \setminus \{i,j,k\}}$
If yes, direct edges as



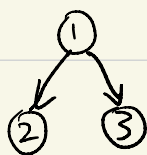
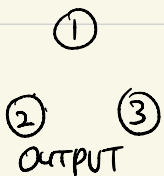
Step 4. Orient remaining undirected edges
by consistency, recursively do this

* How do we check $X_i \perp\!\!\!\perp X_j \mid X_S$?

Q. When does SGS-Algorithm fail
to recover MEC?

ex $\rightarrow ICP(X_1, X_2, X_3) = \{X_1 \perp\!\!\!\perp X_2 \mid X_3, X_2 \perp\!\!\!\perp X_3 \mid X_1, X_1 \perp\!\!\!\perp X_3 \mid X_2\}$

SGS algorithm

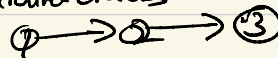


Possible Ground truth.

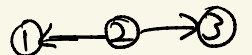
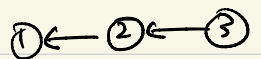
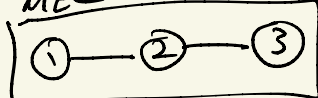
$$I(X_i \perp\!\!\!\perp X_j \mid X_S) = 0$$

$$\uparrow \text{all } I(X_i \perp\!\!\!\perp X_j \mid X_S)$$

Ground truths



MEC



Recall. Global Markov Property.

If X_i and X_j are d -separated in G by S , then
 $X_i \perp\!\!\!\perp X_j \mid X_S$.

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

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

Claim. If $\left\{ \begin{array}{l} X^{(1)}, \dots, X^{(n)} \text{ iid } P(x) \\ P(x) \text{ is faithful to a graph } G^* \\ \text{all variables in } G \text{ are observed.} \end{array} \right.$

Then SGS-algorithm is consistent, i.e.

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

↑
MEC

· Constraint-based algorithm $\left\{ \begin{array}{l} \text{requires lots of samples} \\ \text{faithfulness assumption.} \end{array} \right.$
↑
Cond independencies as constraints.

* Score-based Algorithms

Recall: log-likelihood score of a DAG G . $X^{(1)} \dots X^{(N)}$

$$\text{SCORE}(G) = \underbrace{N \sum_{i=1}^n \log \hat{P}_G(X_i | X_{\pi_i})}_{\text{depends on } G} - \underbrace{N \sum_{i=1}^n H_{\hat{P}}(X_i)}_{\text{Does not depend on } G}$$

without further restrictions on G , complete DAG always achieves max SCORE .

Def. Bayesian Information Criteria (BIC) score.

$$\text{SCORE}_{\text{BIC}}(G) \triangleq \underbrace{\text{SCORE}(G)}_{\text{log-likelihood}} - \underbrace{\frac{\log N}{2} \cdot \text{dim}(G)}_{\text{how complex model.}}$$

how many variables needed to describe $P(x)$

$$\prod_{i=1}^n P_i(X_i | X_{\pi_i})$$

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

Minimum Description Lengths (MDL) principle.

- first term log-likelihood scales as N .
 - second term regularization scales as $\log N$.
- Samples \downarrow second dominates.
Samples \uparrow

$$\frac{N}{\log N} \approx \# \text{ variables in the model.}$$

* Properties:

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

② Consistency: If G^* is a perfect map for $P(x)$.
Then as $N \rightarrow \infty$, G^* is the unique maximizer.

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

\Rightarrow Greedy Equivalence Search (GES).

Algorithm [Greedy Equivalence Search].

Initialize $G^{(1)} = (V, E = \emptyset)$

Phase I: $t = 1, \dots, T \leftarrow$ time until no more gain
add an edge that maximizes $SCORE_{BIC}(G^{(t+1)})$

Phase II: $t = T+1, \dots$

remove an edge that maximizes $SCORE_{BIC}(G^{(t)})$

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

*How long can T be?

* Permutation-based Greedy Search Algorithm

Idea:

Table 1: Equivalence Class Counts

n	Equivalence classes	CI/ADG	CI ₁ /CI
1		1.00000	1.00000
2	00 00	0.66667	0.50000
3		0.44000	0.36364
4		0.34070	0.31892
5		0.29992	0.29788
6		0.28238	0.28667
7		0.27443	0.28068
8		0.27068	0.27754
9		0.26888	0.27590
10		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. $10^6 = 3,628,800$



we apply Greedy Search.

Greedy Search for Sparsest Permutation [GSP] Algorithm.

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

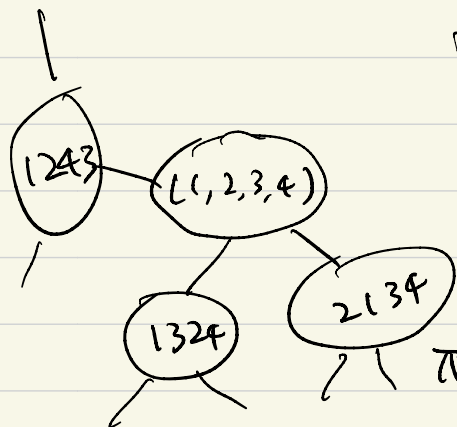
Repeat: $t=1, \dots$

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

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

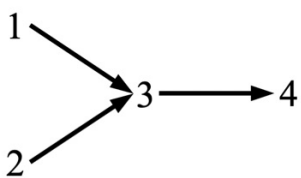
Evaluate $\text{SCORE}_{\text{BIG}}(G_\pi)$

$\pi^{(t+1)} \leftarrow$ 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:

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

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

