# Lecture 4: Learning DAG models from Data

As we saw yesterday $\vec{K_n}$ is an I-map for any distribution $P$. So we want a minimal I-map.

A goal of the course is to develop knowledge of algorithms for learning representations from data.

Recall $P$ is <u>faithful</u> to DAG $G$ if $CI(P) = CI(G)$

Most algorithms assume faithfulness, "not such a strong assumption"

<u>Lemma:</u> The set of unfaithful distributions has Lebesgue measure zero

<u>Goal:</u> Assume $P$ faithful to true causal structure $G$ and recover $G$ from data $D$.

## The PC-algorithm

Assume $CI(P) = \{ 1 \not\perp 2, 1 \not\perp 4 | 3, 1 \not\perp 4 | \{2,3\}, 2 \not\perp 4 | 3, 2 \not\perp 4 | \{1,3\} \}$
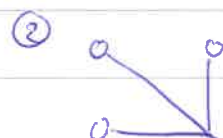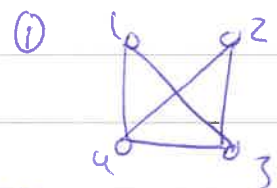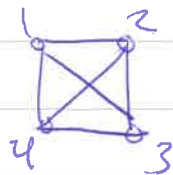
### Step 1 Learn skeleton

⓪ start with $K_m$, complete graph on $[n]$
① Eliminate edges $ij$ if $i \not\perp j | \emptyset$
② For each $ij$ still on edge & $k$ adjacent to either $i$ or $j$. (i.e. $k \in N(i) \cup N(j)$) eliminate $ij$ if $i \not\perp j | k$
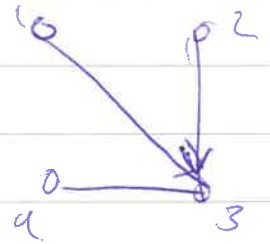③ Repeat ② with larger sets $\{k_1, k_2\} \subseteq N(i)$ or $\{k_1, k_2\} \subseteq N(j)$
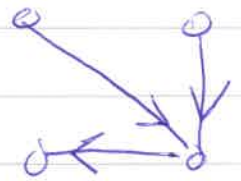④ Repeat ③ until we reached $|N(i)|$ $\forall i$
with larger neighbourhoods

# Step 2: Learn v-structures

⑤ For each triple ~~path~~ path $\langle i, j, k \rangle$
$i,k$ not adjacent orient path as $i \rightarrow j \leftarrow k$
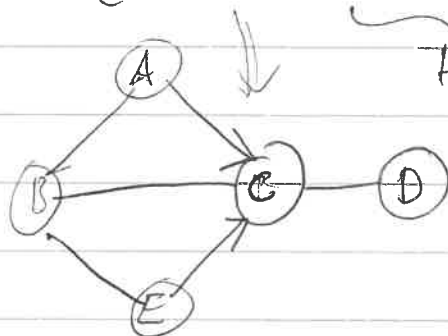if $j$ was not in the conditioning set $C$
in previous step

⑥ For each triple $i \rightarrow j - k$ orient $j \rightarrow k$ [$i,k$ not adjacent]
to avoid v-structure
Use also 4 configurations
~~⑦~~ from def. of strongly protected.

Example $6\mathcal{I} = \{A \perp\!\!\!\perp E \mid B, \ D \perp\!\!\!\perp E \mid C, \ D \perp\!\!\!\perp A \mid C, D \perp\!\!\!\perp B \mid C\}$

these give $\textcircled{D}$ only connected
to C

Must get:

$\textcircled{C} \longrightarrow \textcircled{D}$    Must counteract at least one edge in
the triangles

By fourth config in strongly protected it
becomes

Advantages
- poly time for sparse graphs.   $O(u^{\Delta H})$
- Recovers M.E.C if $\mathbb{P}$ faithful to some DAG
- software $\mathcal{I}$

Disadv.
- When it is wrong can be very wrong, not stable.

**theorem ⑥** Suppose the data generating distribution $P$ is faithful to $G$ and we are using a test for conditional independence that perfectly answers any query as $|D| \to \infty$. Then any acyclic orientation of the output of the PC algorithm that introduces no new v-structures will be ME to $G$.

**Proof:** By assumption, whenever we test $\mathbb{Z}_A \perp\!\!\!\perp \mathbb{Z}_b \mid \mathbb{Z}_C$, as $|D| \to \infty$, we will learn that the statement holds precisely when it holds in $P$.

So need only check that we test the right statements needed to recover a DAC ME to $G$. So need to check that we learn the correct skeleton ~~(problem 3 but have yesterday)~~ on v-structures.

_Skeleton:_ recall $i$ and $j$ adjacent in $G \iff i, j$ d-separated given either $pa_G(i)$ or $pa_G(j)$. In step 1, eventually the conditioning set is large enough that we test (wlog) $i \perp\!\!\!\perp j \mid pa_G(i)$, which holds in $P$, by faithfulness. $\Rightarrow$ learn correct skeleton.

_v-structures:_ Recall $i - j - k$ a v-structure $\iff^{①} \exists\, C \subset [m] \setminus \{i, j, k\}$ d-separating $i$ and $k$ in $G$ such that $j \notin C. \iff^{②}$ all sets $C \subset [m] \setminus \{i, k\}$ containing $j$ do not d-separate $i$ and $k$.

Let $i \perp\!\!\!\perp k \mid C$ be statement from step 1 that lead to removal of edge $i - k$. [If $j \notin C$, by faithfulness of $P$, $i \to j \leftarrow k$ v-structure in $G$ by ① If $j \in C$, by ② and faithfulness, no v-structure]

Similarly, by faithfulness, the orientation in step ⑦ will produce no additional v-structures.

So the output has exactly the skeleton & v-structures of $G$. ? Any such orientation of it will by ME to $G$ by VP.

Assuming faithfulness and a consistent test for conditional independence PC algorithm will learn (up to Markov equivalence) the true causal structure. However, we typically use hypothesis tests for testing CI, which are prone to error... Causes accuracy problems

PC algorithm is a constraint-based algorithm, relying on CI tests. To avoid problems with accuracy we can instead use greedy (score-based) algorithms...

# Greedy Equivalence Search (GES)

- Assign each DAG a score, such as the Bayesian Information Criterion (BIC):

$$BIC(G, D) := \log P(D | \hat{\Theta}, G) - \frac{d}{2} \log(|D|)$$

log-likelihood of $D$ given $G$      penalization term based on # of free parameters $d$ in $G$

**Lemma ④** BIC is _score equivalent_; i.e. if $G$ and $H$ are Markov equivalent then $BIC(G, D) = BIC(H, D)$ for all data sets

- GES relies on the following generalization of the characteriz[...] of Markov equivalence proven in problem set 3.

- Write $G \leq H$ if $CI(G) \supset CI(H)$ where
$$CI(G) := \{ A \perp\!\!\!\perp B \mid C : A, B \text{ d-separated given } C \text{ in } G \}.$$

**Theorem ⑦** (Chickering, 2002) Suppose $G \leq H$. Then $\exists$ a sequenc[...] of edge reversals and edge deletions such that

(1) all edges reversed are _covered_; $i \to j$ where $pa_G(j) = pa_G(i) \cup \{i\}$.

(2) After each edge reversal or addition get a DAG $\hat{G}$ s.t. $\hat{G} \leq H$

(3) After _all_ edge reversals and additions $\hat{G} = H$.

GES: Let $[G]$ denote MEC of $G$.

① Start at empty DAG $G := ([m], \emptyset)$.      (a DAG i[...]

② $E^+([G]) := \{$ MECs of DAGs produced by adding single edge to $[G]$[...] Pick $[H] \in E^+([G])$ with highest BIC and set $G := H$.

③ Repeat ② until no higher scoring MEC found.

④ $E^-([G]) := \{$ MECs produced by removing single edge from DAG in $[G]$[...] Pick $[H] \in E^-([G])$ with highest BIC and set $G := H$.

⑤ Repeat step ④ until no higher scoring MEC found.

⑥ Return $[G]$.

**Theorem ⑧** (Chickering, 2002) IF $P$ faithful to $G$, as $|D| \to \infty$, the output of GES will be the MEC of $G$.

- While GES is more accurate, PC has polynomial time performance guarantees (see HW).

- Hybrid Algorithms: Combine score-based and constraint-based approaches to strike balance.