# Bayesian Causal Structure Learning

**Simon Rittel[1,2], Sebastian Tschiatschek[1]**

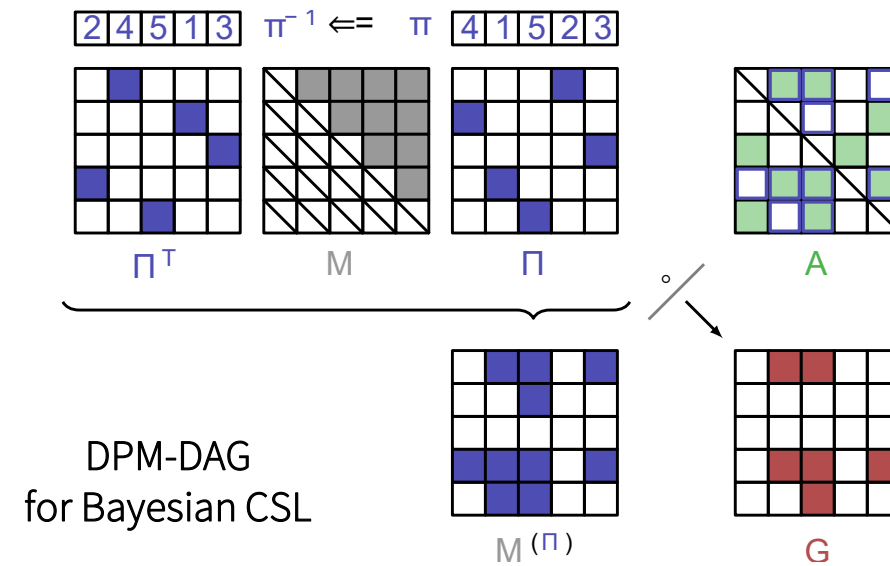[1]University of Vienna, Faculty of Computer Science, Vienna, Austria  [2]University of Vienna, UniVie Doctoral School for Computer Science, Vienna, Austria

Simon Rittel

Sebastian Tschiatschek
(supervison)

DPM-DAG
for Bayesian CSL

# Outline of this talk

- Intro to Causality

- Causal Structure Learning

- Bayesian Causal Structure Learning

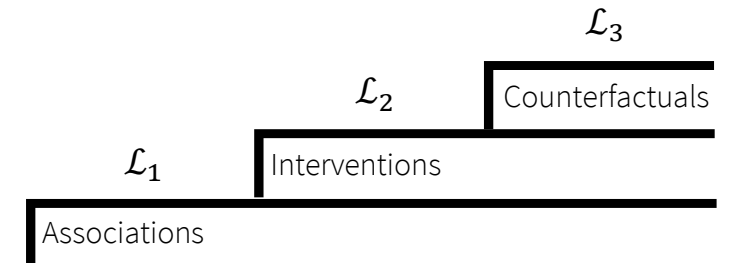- Differentiable Probabilistically Masked DAG (*DPM-DAG*)

# Causality for Machine Learning

- Fundamental differences

  Statistical relations ($\mathcal{L}_1$)

  $\longleftarrow$
  $\longrightarrow$

  Causal relations ($\mathcal{L}_{1:3}$)

  $\longrightarrow$

  $\mathcal{L}_3$

  Counterfactuals

  $\mathcal{L}_2$

  Interventions

  $\mathcal{L}_1$

  Associations

- Promising applications

  Robustness & generalization

  Interpretability & explainability

  Fairness

  Causal Insights

# Causal Model as Basis for Causal Inference

- Indiviudal effect of a treatment $T$ on an outcome $Y$: $\qquad$ $\text{ITE}_i := Y_i(T = 1) - Y_i(T = 0)$

- Average treatment effect: $\qquad$ $\text{ATE} := \mathbb{E}_i[\text{ITE}_i]$

- Identifiability: $\qquad$ Causal effect can be consistently estimated from observed data

| Causal Estimand | → | Statistical Estimand | → | Estimate |

Identification
using a causal model

Estimation
based on observed data

- Controlling/adjusting for a set of confounders $X$:

$$\text{CATE} = \mathbb{E}_X\big[\mathbb{E}_Y[Y|T = 1, X] - \mathbb{E}_Y[Y|T = 0, X]\big] = \frac{1}{|\mathcal{D}_{T=1}|} \sum_{i \in \mathcal{D}_{T=1}} \hat{\mu}_{Y|T=1, X}(X_i) - \frac{1}{|\mathcal{D}_{T=0}|} \sum_{j \in \mathcal{D}_{T=0}} \hat{\mu}_{Y|T=0, X}(X_j)$$

# Recap of Bayesian Networks

- Graphical model:
  One-to-one mapping between nodes $V_i \in V$ of a direct acyclic graph (DAG) $G = (V, E)$ and random variables $X_i \in \boldsymbol{X}$

- Local Markov Condition:
  Given the parents $\mathbf{pa}$ of a node $V_i$ in the DAG $G$, the corresponding random variable $X_i$ is independent of all its non-descendants $\mathbf{nd}$.

- Bayesian Network Factorization:
  Given a joint probability distribution $P_X$ and a DAG $G$, $P_X$ factorizes according to $G$ if:

  $$P(\boldsymbol{X}) := P\left(\{X_i\}_{i=1}^{D}\right) = \prod_{i}^{D} P(x_i | \mathbf{pa}(X_i))$$
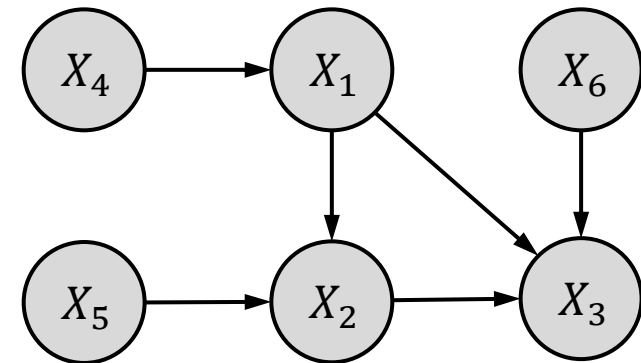
Fig. 1: DAG over six random variables
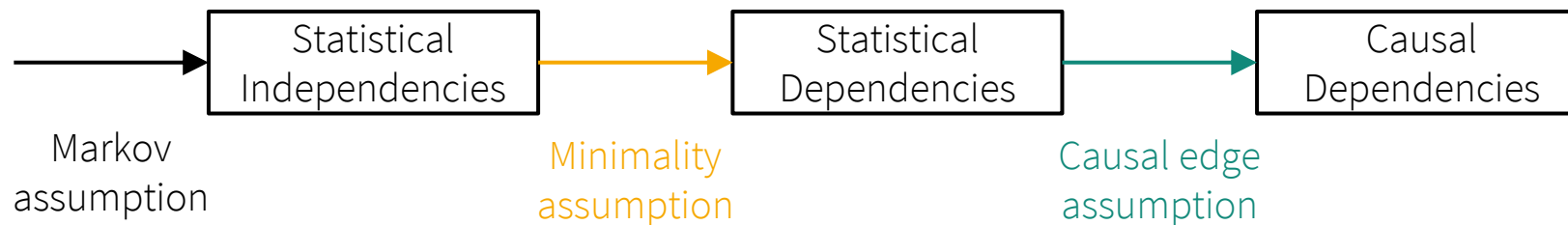
# Causal Graphs

- Minimality assumption:

    1. Local Markov condition
       (implies *d-separation* as *global Markov condition)*

       $$X_i \perp_P \mathrm{nd}(X_i) \mid \mathbf{pa}(X_i)$$

    2. Adjacent nodes in the DAG $G$ are dependent
       (no additional independences)

       $$X_i \sim X_j \; in \; G \implies X_i \perp_P X_j$$

- Strict causal edge assumption:
  Every parent is a direct cause of all its children, i.e. the children are affected by changes in their parents

# Causal Structure Learning (*CSL*)

- Functional Causal Model: indexed tuple of
  - endogenous variables $\boldsymbol{X}$,
  - exogeneous noise variables $\boldsymbol{\epsilon}$ with distribution $P_{\boldsymbol{\epsilon}}$ ,
  - deterministic functions $\boldsymbol{g}$, s. t. $X_i \coloneqq g_i(\mathbf{pa}_{\boldsymbol{G}}(X_i), \epsilon_i)$

- Assumptions:
  - Acyclic causal relations
    → Direct Acyclic Graph (*DAG*) $\boldsymbol{G}$
  - Causal sufficiency
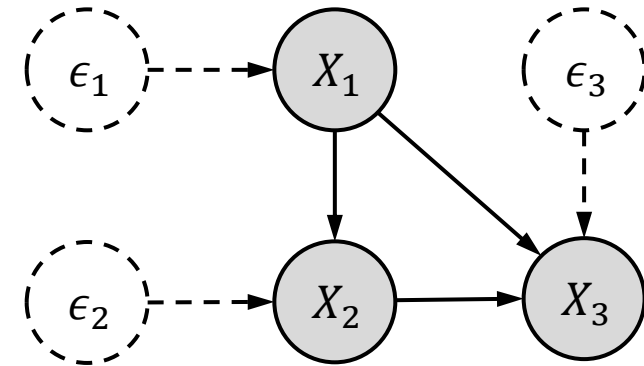    → no latent confounders and mutually independent noise $\boldsymbol{\epsilon}$



Fig. 2: Causal DAG induced by an acyclic FCM over three observed random variables

# Interventions by the do-Operator ($\mathcal{L}_1$)

- Definition
  Hard intervention $\mathbf{do}(X_i = x)$ replaces structural function $g_i$
  by the assignment $X_i = x$

- Truncated Factorization

$$P(\mathbf{X}|\mathrm{do}(X_i = x) \coloneqq \delta(X_i = x)\prod_{j \neq i} P(x_j|\mathbf{pa}(X_j))$$
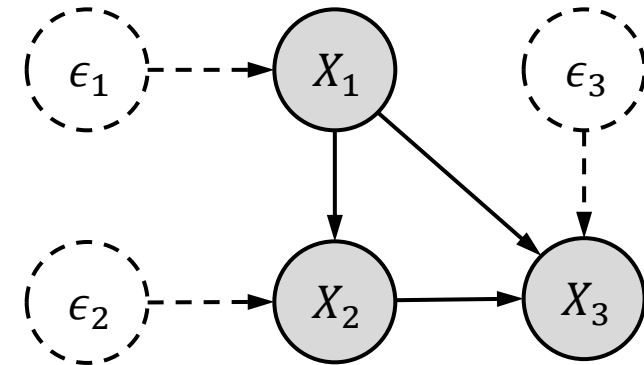


Fig. 2: Causal DAG induced by an acyclic FCM
over three observed random variables

# Interventions by the do-Operator $(\mathcal{L}_1)$

- Definition
  Hard intervention $\mathbf{do}(X_i = x)$ replaces structural function $g_i$
  by the assignment $X_i = x$

- Truncated Factorization

$$P(\boldsymbol{X}|\mathrm{do}(X_i = x) \coloneqq \delta(X_i = x)\prod_{j \neq i} P(x_j|\mathbf{pa}(X_j))$$

- Example:

  ◦ $P(X_1, X_3) = \int P(X_3|X_1, X_2)P(X_2|X_1)P(X_1)\, dX_2$



Fig. 2: Causal DAG induced by an acyclic FCM
over three observed random variables

# Interventions by the do-Operator ($\mathcal{L}_1$)

- Definition
  Hard intervention $\mathbf{do}(X_i = x)$ replaces structural function $g_i$
  by the assignment $X_i = x$

- Truncated Factorization

$$P(\boldsymbol{X}|\mathrm{do}(X_i = x) \coloneqq \delta(X_i = x) \prod_{j \neq i} P(x_j|\mathbf{pa}(X_j))$$

- Example:

  ○ $P(X_1, X_3) = \int P(X_3|X_1, X_2)P(X_2|X_1)P(X_1)\, dX_2$

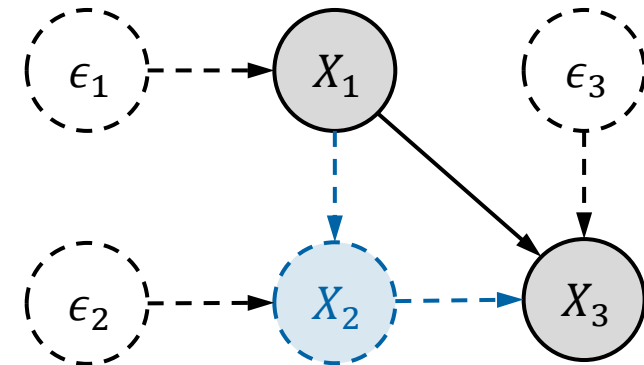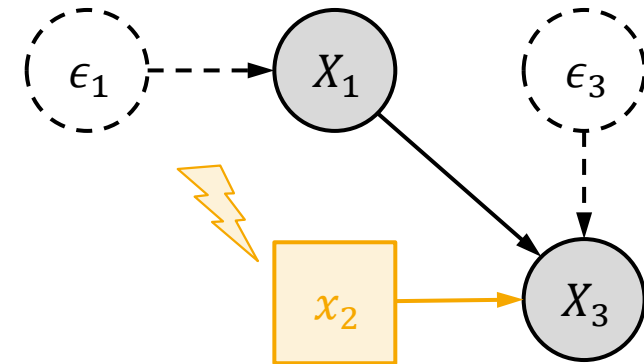  ○ $P(X_1, X_3|do(X_2 = x_2)) = P(X_3|X_1, X_2 = x_2)P(X_1)$

Fig. 2: Causal DAG induced by an acyclic FCM
over three observed random variables

# Interventions by the do-Operator ($\mathcal{L}_1$)

- Definition
  Hard intervention $\mathbf{do}(X_i = x)$ replaces structural function $g_i$
  by the assignment $X_i = x$

- Truncated Factorization

$$P(\boldsymbol{X}|\mathbf{do}(X_i = x) \coloneqq \delta(X_i = x)\prod_{j \neq i} P(X_j|\mathbf{pa}(X_j))$$

- Example:

  ○ $P(X_1, X_3) = \int P(X_3|X_1, X_2)P(X_2|X_1)P(X_1)\, dX_2$

  ○ $P(X_1, X_3|do(X_2 = x_2)) = P(X_3|X_1, X_2 = x_2)P(X_1)$

  ○ $P(X_1, X_3| X_2 = x_2) = \dfrac{P(X_3|X_1, X_2 = x_2)P(X_2 = x_2|X_1)P(X_1)}{P(X_2 = x_2)}$



Fig. 2: Causal DAG induced by an acyclic FCM over three observed random variables

# The Three Layer Causal Hierarchy by Pearl

| Level | Typical Quantity | Typical Activity | Typical Questions |
|---|---|---|---|
| 1. Association | $P(Y\|X=x)$ | Seeing | What is? How does observing X change my belief in Y? |
| 2. Intervention | $P(Y\|do(X=x))$ | Doing/Intervening | What if I do X? |
| 3. Counterfactuals | $P(Y_x\|do(X=x'),y')$ | Imagining, Retrospection | Why? Was it X that caused Y? |

# Typical Assumptions for Independence-based CSL

- Markov assumption:
$$X \perp_G Y \mid Z \implies X \perp_P Y \mid Z$$

- Faithfulness:
$$X \perp_G Y \mid Z \impliedby X \perp_P Y \mid Z$$



$$X_4 := \gamma X_2 + \delta X_3 = \underbrace{(\alpha\gamma + \beta\delta)}_{\neq 0} X_1$$

# Typical Assumptions for Independence-based CSL

- Markov assumption:
$$X \perp_G Y \mid Z \implies X \perp_P Y \mid Z$$

- Faithfulness:
$$X \perp_G Y \mid Z \impliedby X \perp_P Y \mid Z$$

- Causal sufficiency:
No unobserved confounders.

$X_1 \perp X_2 \mid W$      True graph

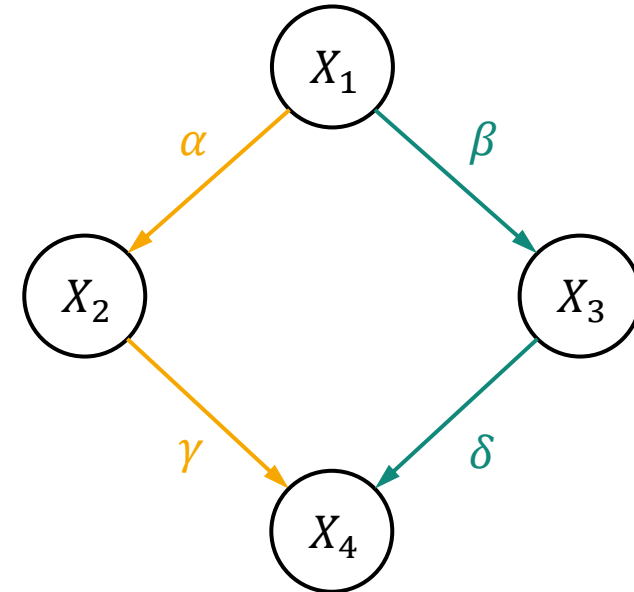$X_1 \not\perp X_2$      Projected graph

# Typical Assumptions for Independence-based CSL

- Markov assumption:
$$X \perp_G Y \mid Z \implies X \perp_P Y \mid Z$$

- Faithfulness:
$$X \perp_G Y \mid Z \impliedby X \perp_P Y \mid Z$$

- Causal sufficiency:
No unobserved common causes

- No selection bias:
No conditioning on unobserved colliders

$X_1 \perp X_2$

$X_1 \not\perp X_2 \mid C$

True graph

Projected graph

# Review of Unshielded 3-node Structures

Causal trail

$$X \longrightarrow Z \longrightarrow Y$$

Evidential trail

$$X \longleftarrow Z \longleftarrow Y$$

Common cause

$$X \longleftarrow Z \longrightarrow Y$$

$X \not\perp Y$
$X \perp Y \mid Z$

$$X \text{—} Z \text{—} Y$$

Common effect

$$X \longrightarrow Z \longleftarrow Y$$

$X \perp Y$
$X \not\perp Y \mid Z$

# Sketch of the PC-algorithm

1) Start with a complete undirected graph

2) Eliminate edges between variables that are (conditionally) independent

3) Add arrow marks at colliders in identified v-structure

4) Propagate arrows such that no additional v-structures are formed that were not detected

True graph

PCDAG :
skeleton + v-structures

# Known Identifiable Causal Models

- Linear Gaussian model with equal or known variance $\qquad Y := aX + \epsilon \quad$ with $\quad \epsilon \sim \mathcal{N}(\mu, \sigma)$
  (Loh & Bühlmann 2014)

- Linear non-Gaussian model (LiNGAM) $\qquad\qquad Y := aX + \epsilon \quad$ with $\quad \epsilon \nsim \mathcal{N}(\mu, \sigma)$
  (Shimizu et al.. 2006)

- Nonlinear additive noise model (ANL) $\qquad\qquad Y := f(X) + \epsilon \quad$ where $f$ is nonlinear
  (Hoyer et al., 2008)

- Post-nonlinear causal model (PNL) $\qquad\qquad Y := g(f(X) + \epsilon) \quad$ where $g$ is nonlinear & invertible
  (Zhang & Hyvärinen, 2009)

# Greedy Equivalence Search ( *GES* )

1) Initialization by an empty graph

2) Forward equivalence search
   Add the edge that most increases the score and maps
   the resulting graph then to its MEC

3) Backward equivalence search
   Remove the  edge that will most increase the score
   until no further edges can be removed

- Score equivalence:
  Graphs of the same MEC are assigned the same score

- Locally consistent scoring criterion:
  Score prefers edge additions that remove incorrect
  dependencies and edge deletions that remove
  incorrect dependencies

- Decomposable score function:

$$S(\boldsymbol{G}, \boldsymbol{X}) = \sum_{d=1}^{D} S(X_i, \boldsymbol{pa_G}(X_i))$$

# Continuous Relaxation of the Discrete Graph Structure & Acyclicity

- Converting the combinatorical optimization problem into a continuous program

$$\min_{G \in \{0,1\}^{D \times D}} S(G) \qquad \Leftrightarrow \qquad \min_{G \in [0,1]^{D \times D}} S(G)$$

$$\text{subject to } G \in G_{acyclic} \qquad \qquad \text{subject to } h(G) = 0$$

- Differentiable DAG-Constraint $\quad h(G_{\text{acyclic}}) = 0 \, , \quad h(G_{\text{cyclic}}) > 0$

  - $h_1(G) = \text{tr}(e^{G \circ G}) - D$          (Zheng et al. , 2018)

  - $h_2(G) = \text{tr}\left(\left(I + \frac{1}{D}(G \circ G)\right)^D\right)$     (Yu et al., 2019)

  - $h_3(G) = \log \det(sI - G \circ G) + D \log s$    (Bello et al. , 2023)

# Research Areas in Causal Structure Learning

- Relaxing assumptions
  - No assumed causal sufficiency :     FCI algorithm  (Spirtes et al., 2001)
  - No assumed acyclicity                  CCD algorithm  (Richardson, 1996)
  - Neither of both:                          SAT-based causal discovery  (Hyttinen et. al., 2013)

# Research Areas in Causal Structure Learning

- Relaxing assumptions

- Improving computational scalability
  - Limiting the number of potential parents:         PNS-algorithm (Bühlmann et al., 2014)
  - Omitting some CI test:                             RFCI-algorithm (Colombo et al, 2012)
  - Considering only one edge change at a time:       GES-algorithm (Chickering, 2002)
  - Continuous relaxation of the binary adjacency matrix:    NOTEARS-algorithm (Zheng et al., 2018)

# Research Areas in Causal Structure Learning

- Relaxing assumptions

- Improving computational scalability

- Increasing robustness:
  - Additional CI-tests:    Order-independent PC/FCI (Colombo & Maathuis, 2014)

# Research Areas in Causal Structure Learning (non-exhaustive)

- Relaxing assumptions

- Improving computational scalability

- Increasing robustness

- Identifiable functional models

- Focus only on local structure relevant for downstream task

- Modeling uncertainty in the prediction

- Combining with interventional data

# Independence-based CSL

- Based on Conditional Independence (CI) tests

- Additional assumption of faithfulness

- Iterative restriction of CI test to avoid all pairwise combinations

- Point estimate as output

- Sound in the large sample limit

# Bayesian CSL

$$p(\boldsymbol{G}, \boldsymbol{\Theta}|\boldsymbol{X}) \propto p(\boldsymbol{G})p(\boldsymbol{\Theta}|\boldsymbol{G})p(\boldsymbol{X}|\boldsymbol{G}, \boldsymbol{\Theta})$$

- Quantifying the uncertainty in the posterior

- Incorporation of probabilistic domain knowledge via prior

- Sound in the large sample limit

Fig. 2: Generative model

# Generative Model

$$p(\boldsymbol{G}, \boldsymbol{\Theta}, \boldsymbol{X}) = p(\boldsymbol{G})p(\boldsymbol{\Theta}|\boldsymbol{G})p(\boldsymbol{X}|\boldsymbol{G}, \boldsymbol{\Theta})$$

$$p(\boldsymbol{X}|\boldsymbol{G}, \boldsymbol{\Theta}) = \prod_{n=1}^{N}\prod_{d=1}^{D} p\left(X_d^{(n)}\Big|\mathrm{pa}_{\boldsymbol{G}}\left(X_d^{(n)}\right), \boldsymbol{\Theta}\right)$$



Fig. 2: Generative model

# Marginalized Generative Model

$$\mathrm{p}(\boldsymbol{G}, \boldsymbol{X}) = p(\boldsymbol{G}) \int p(\boldsymbol{\Theta}|\boldsymbol{G}) p(\boldsymbol{X}|\boldsymbol{G}, \boldsymbol{\Theta}) d\boldsymbol{\Theta}$$

$$\leq p(\boldsymbol{G}) \boldsymbol{p_{\Theta^*}}(\boldsymbol{X}|\boldsymbol{G})$$

$$\text{where } \boldsymbol{\Theta}^* \coloneqq \arg\max_{\boldsymbol{\Theta}} p(\boldsymbol{X}|\boldsymbol{G}, \boldsymbol{\Theta})$$



Fig. 2: Generative model

# Graph Posterior

$$p_{\boldsymbol{\Theta}^*}(\boldsymbol{G}|\boldsymbol{X}) = \frac{p_{\boldsymbol{\Theta}^*}(\boldsymbol{G},\boldsymbol{X})}{p(\boldsymbol{X})} \propto p_{\boldsymbol{\Theta}^*}(\boldsymbol{G},\boldsymbol{X})$$

$$\text{where } \boldsymbol{\Theta}^* := \arg\max_{\boldsymbol{\Theta}} p(\boldsymbol{X}|\boldsymbol{G},\boldsymbol{\Theta})$$



Fig. 2: Generative model

# Probabilistic Graph: REINFORCE estimator [1]

- Independent Bernoulli distributed RV models each edge

$$G_{ij} \sim \mathrm{Bern}(\phi_{ij})$$

- Score function gradient estimator for its parameters

$$\eta = \boldsymbol{\nabla}_{\boldsymbol{\phi}}\, \mathbb{E}_{p_{\Theta^*(X)}}[f(\boldsymbol{X})] = \boldsymbol{\nabla}_{\boldsymbol{\phi}} \int p_{\Theta^*(X)} f(\boldsymbol{X}) \mathbf{d}\boldsymbol{X} = \int f(\boldsymbol{X}) \boldsymbol{\nabla}_{\boldsymbol{\phi}} p_{\Theta^*(X)} \mathbf{d}\boldsymbol{X} = \int f(\boldsymbol{X}) p_{\Theta^*(X)}\, \boldsymbol{\nabla}_{\boldsymbol{\phi}} \log p_{\Theta^*(X)}\, \mathbf{d}\boldsymbol{X} =$$

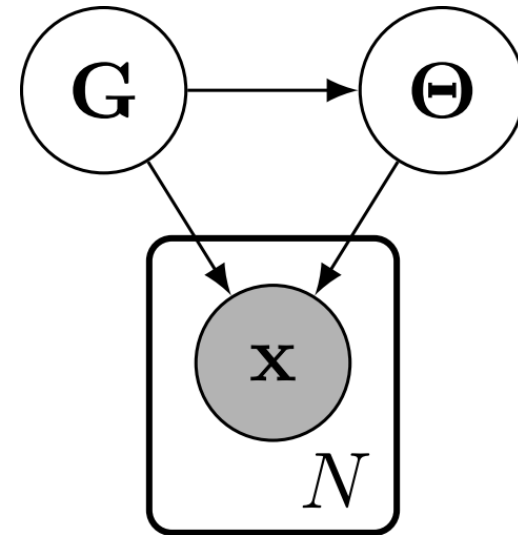$$= \mathbb{E}_{p_{\Theta^*(X)}}\big[f(\boldsymbol{X}) \boldsymbol{\nabla}_{\boldsymbol{\phi}} \log p_{\Theta^*(X)}\big]$$

$$\hat{\eta}_N = \frac{1}{N}\sum_{n=1}^{N} f\big(\widehat{\boldsymbol{X}}^{(\boldsymbol{n})}\big) \boldsymbol{\nabla}_{\boldsymbol{\phi}} \log p_{\Theta^*(\widehat{X}^{(n)})} \qquad \text{where} \quad \widehat{\boldsymbol{X}}^{(\boldsymbol{n})} \sim p_{\Theta^*}(\boldsymbol{X})$$

[1] Rezende et al. , 'MC Gradient Estimation in Machine Learning', in the Journal of Machine Learning Research, (2020)

# Probabilistic Graph: Pathwise Gradient Estimator [2]

- Independent perturbed Gumbel distributed RV models each edge

$$G_i \sim \text{Gumbel}(0, 1) \,, \qquad \phi_i + G_i \sim \text{Gumbel}(\phi_i, 1)$$

- Perturbed Gumbel-Softmax samples

$$\arg\max_{i \in \mathbb{I}}(\phi_i + G_i) \sim \frac{\exp(\phi_i)}{\sum_{j \in \mathbb{I}} \exp(\phi_j)}$$

- **Softmax** as continuous, differentiable relaxation of the **arg max** operator (equivalence for $\boldsymbol{\tau \to 0}$)

$$Z_i = \frac{\exp\big((\phi_i + G_i)/\tau\big)}{\sum_{j \in \mathbb{I}} \exp\big((\phi_i + G_i)/\tau\big)}$$

[2] Maddison et al. , 'The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables', in Proceedings of the International Conference on Learning Representations, (2017)

# Probabilistic Graph: Pathwise Gradient Estimator

- Straight-through estimator
  discrete samples ($\mathbf{arg\ max}$) in the forward pass and continuous samples ($\boldsymbol{softmax}$) in the backward pass
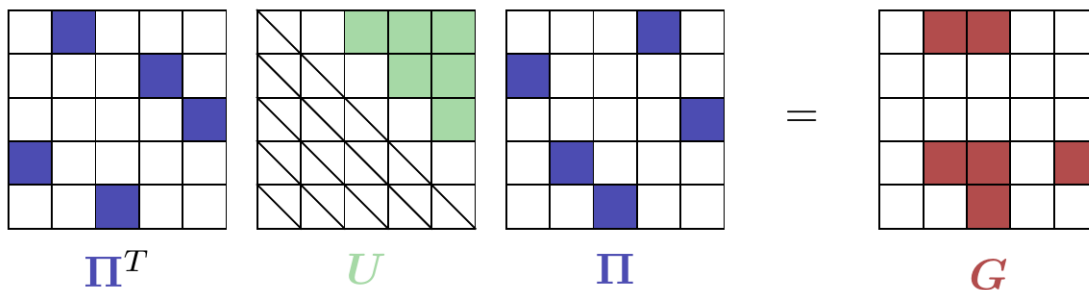
- Logistic samples for binary RV
  $$[G_1 + \phi_1 > G_0 + \phi_0] = [\ \underbrace{G_1 - G_0}_{\doteq L} + \underbrace{\phi_1 - \phi_0}_{=:\ \phi} > 0\ ]\ , \qquad \text{where } L \sim \text{Logistic}(0,1)$$

- Sigmoid as 2-dim version of Softmax

$$Z_i = \left( 1 + \exp\left( -\frac{L_i + \phi_i}{\tau} \right) \right)^{-1}$$

# Enforcing Acyclicity

1)  Permuted upper triangular matrix [3]



$$\mathbf{\Pi}^T \qquad U \qquad \mathbf{\Pi} \qquad = \qquad G$$

$$p(G) = \sum_{\mathbf{\Pi} \in \mathcal{P}_D(G)} p(G, \mathbf{\Pi})$$

Mean-field approximation: $\quad p(G, \mathbf{\Pi}) = p(U)p(\mathbf{\Pi})$

2)  Differentiable acyclicity constraint [4]

$$h(G_{\text{cyclic}}) > 0 \quad , \quad h(G_{\text{acyclic}}) = 0$$

$$p(G) \propto e^{-\lambda h(G)}$$

$$p(G_{\text{cyclic}}) \xrightarrow[\lambda \to \infty]{} 0 \quad , \quad p(G_{\text{acyclic}}) \xrightarrow[\lambda \to \infty]{} \frac{1}{|\mathbb{G}_{acyclic}|}$$

[3] Charpentier  et al. , 'Differentiable DAG sampling', in Proceedings of the International Conference of
    Learning Representations, (2022)
[4] Lorch et al. , 'DiBS: Differentiable Bayesian structure learning', in Advances in Neural Information
    Processing Systems, volume 34, pp. 24111-2413, (2021)

# Variational Posterior Not Constrained to DAGs

$$D_{\mathrm{KL}}\big(p_{\boldsymbol{\phi}}(\boldsymbol{G})\big|p_{\lambda}(\boldsymbol{G})\big)$$

KL-Divergence

Posterior [5]　　　　　　　　Prior

$$p_{\boldsymbol{\phi}}(\boldsymbol{G}) = \prod_{i \neq j} p_{\phi_{ij}}(\boldsymbol{G}_{ij})$$

$$p_{\lambda}(\boldsymbol{G}) \propto \exp(-\lambda h(\boldsymbol{G}))$$

DAG $\boldsymbol{G}_1$　　$p_1$

DAG $\boldsymbol{G}_2$　　$p_2$

Cyclic $\boldsymbol{G}_3$　　$p_3$

[5] Geffner et al. , 'Deep end-to-end causal inference', in NeurIPS Workshop on Causality for Real-world Impact, (2022)

# Incorporating Probabilistic Knowledge in a Gibbs Prior

- Number of expected causes for every node [3]

  ○ Erdös-Renyi graphs
  $$p(\boldsymbol{G}) \propto p^{\|\boldsymbol{G}\|_1}(1-p)^{E-\|\boldsymbol{G}\|_1}$$

  ○ Scale-free graphs
  $$p(\boldsymbol{G}) \propto \prod_{i=1}^{D} \left(1 + \|G_i^{\mathrm{T}}\|_1\right)^{-3}$$

- Additional sparsity regularization [4]
  $$p(\boldsymbol{G}) \propto \beta\|\boldsymbol{G}\|_F^2$$

- Prior over a single edge $\boldsymbol{p_{ij}}$

$$p(\boldsymbol{G}) \propto \left(q_{ij}\boldsymbol{G}_{ij} + (1-q_{ij})(1-\boldsymbol{G}_{ij})\right)$$

$$p(\boldsymbol{G} \in \mathbb{G}_{ij}) = \frac{q_{ij}}{p_{ij} + (1-p_{ij})} = q_{ij} := \frac{p_{ij}}{|\mathbb{G}_{ij}|}$$

[4] Lorch et al. , 'DiBS: Differentiable Bayesian structure learning', in Advances in Neural Information Processing Systems, volume 34, pp. 24111-2413, (2021)
[5] Geffner et al. , 'Deep end-to-end causal inference', in NeurIPS Workshop on Causality for Real-world Impact, (2022)

# Differentiable Probabilistic DAG (*DP-DAG*) [3]



1) $\mathbf{\Pi} \sim p_{\psi}(\mathbf{\Pi})$      Gumbel-Softsort

2) $\mathbf{U} \sim p_{\phi}(\mathbf{U})$      Gumbel-Softmax

3) $\mathbf{G} = \mathbf{U}^{(\mathbf{\Pi})} = \mathbf{\Pi}^T \mathbf{U}\, \mathbf{\Pi}$

$$p_{\psi,\phi}(\boldsymbol{G}, \boldsymbol{\Pi}) = p_{\psi}(\mathbf{\Pi}) \prod_{i \neq j} \boldsymbol{p}_{\phi_{ij}}\left(\mathbf{U}_{ij}^{(\mathbf{\Pi})}\right) \left[\mathbf{U}_{ij}^{(\mathbf{\Pi})} = \mathbf{G}_{ij}\right]$$

[3] Charpentier et al. , 'Differentiable DAG sampling', in in Proceedings of the *International Conference of Learning Representations*, (2022)

# Differentiable Probabilistically Masked DAG (*DPM-DAG*) [6]



$$2\ 4\ 5\ 1\ 3 \quad \pi^{-1} \Longleftarrow \quad \pi \quad 4\ 1\ 5\ 2\ 3$$

$\mathbf{\Pi}^T \qquad M \qquad \mathbf{\Pi} \qquad\qquad A$

$\circ$

$\mathbf{M}^{(\mathbf{\Pi})} \qquad\qquad \mathbf{G}$

1) $\quad \mathbf{\Pi} \sim p_{\boldsymbol{\psi}}(\boldsymbol{\Pi}) \qquad$ Gumbel-Softsort

2) $\quad \mathbf{M}^{(\mathbf{\Pi})} = \mathbf{\Pi}^T M\, \mathbf{\Pi}$

3) $\quad \mathbf{A} \sim p_{\boldsymbol{\phi}}(\boldsymbol{A}) \qquad$ Gumbel-Softmax

4) $\quad \mathbf{G} = \mathbf{M}^{(\mathbf{\Pi})} \circ \mathbf{A}$

$$p_{\boldsymbol{\psi},\boldsymbol{\phi}}(\boldsymbol{G}, \boldsymbol{\Pi}) = p_{\boldsymbol{\psi}}(\mathbf{\Pi}) \prod_{i \prec j \text{ in } \mathbf{\Pi}} p_{\boldsymbol{\phi}}(\boldsymbol{A}_{ij} = \boldsymbol{G}_{ij}) \prod_{j \prec i \text{ in } \mathbf{\Pi}} [0 = \boldsymbol{G}_{ij}]$$

# Prior specification

- Gumbel-SoftSort is equal in distribution to the Plackett-Luce distribution

$$\underset{i \in \mathbb{I} \setminus \mathbb{S}}{\arg \max}(\psi_i + g_i) \sim p\left(\frac{\exp(\psi_i)}{\sum_{j \in \mathbb{I} \setminus \mathbb{S}} \exp(\psi_j)}\right) \quad \Longrightarrow \quad p^{(\mathrm{PL})}(i \prec j) = p(\mathbf{M}_{ij}^{(\mathbf{\Pi})} = 1) = \frac{w_i}{w_i + w_j}$$

- Prior over permutation

$$D_{\mathrm{KL}}\big(p_{\boldsymbol{\psi}}(\mathbf{\Pi}) \big| p_{\boldsymbol{\omega}}(\mathbf{\Pi})\big) \approx \sum_i^D w_i (\log w_i - \log \omega_i)$$

- Prior over unmasked part of $A$

$$D_{\mathrm{KL}}\big(p_{\boldsymbol{\psi},\boldsymbol{\phi}}(\boldsymbol{G}|\mathbf{\Pi}) \big| p_{\boldsymbol{\gamma}}(\boldsymbol{G}|\mathbf{\Pi})\big) = \sum_{\mathbf{\Pi}} \sum_{i \prec j \ \mathrm{in} \ \mathbf{\Pi}} a_{ij} \frac{\log a_{ij}}{\log \gamma_{ij}} + (1 - a_{ij}) \frac{\log(1 - a_{ij})}{\log(1 - \gamma_{ij})}$$

# Variational Loss for Bayesian CSL



$$\psi \qquad \phi$$

Fig. 3: Generative model of DPM-DAG

- Maximizing evidence lower bound (*ELBO*) $\quad \max_{\boldsymbol{\psi},\boldsymbol{\phi},\boldsymbol{\Theta}} \mathcal{L}$

- DP-DAG $\quad \mathcal{L} = \mathbb{E}_{\mathbf{G}\sim p_{\boldsymbol{\psi},\boldsymbol{\phi}}(\boldsymbol{G})}[\log p_{\boldsymbol{\Theta}}(\boldsymbol{x}|\boldsymbol{G})] - \beta \underbrace{D_{\mathrm{KL}}\big(p_{\boldsymbol{\phi}}(\boldsymbol{A})\big|\big(p_{\boldsymbol{\phi}}(\mathbf{A})\big|p(\mathbf{A})\big)\big)}_{\Pi_{i\neq j} D_{\mathrm{KL}}\big(p_{\boldsymbol{\phi}_{ij}}(\boldsymbol{A}_{ij})\big|p\big)}$

- DPM-DAG $\quad \mathcal{L} = \mathbb{E}_{\mathbf{G}\sim p_{\boldsymbol{\psi},\boldsymbol{\phi}}(\boldsymbol{G})}[\log p_{\boldsymbol{\Theta}}(\boldsymbol{X}|\boldsymbol{G})] - D_{\mathrm{KL}}\big(p_{\boldsymbol{\phi}}(\boldsymbol{G}|\boldsymbol{\Pi})\big|p_{\boldsymbol{\gamma}}(\boldsymbol{G}|\boldsymbol{\Pi})\big) - D_{\mathrm{KL}}\big(p_{\boldsymbol{\psi}}(\boldsymbol{\Pi})\big|p_{\boldsymbol{\omega}}(\boldsymbol{\Pi})\big)$

# Influence of the prior over unmasked edges $p_\gamma$ on AUROC (↑) & AUCPR (↑)



Probabilistic knowledge of true causal graph $\boldsymbol{G}^*$

- For $G^*_{ij} = 1$:
$$p(A_{ij} = 1) \coloneqq a_{ij}$$

- For $G^*_{ij} = 0$:
$$p(A_{ij} = 1) \coloneqq 1 - a_{ij}$$

# Influence of the prior over the order $p_\omega$ on AUROC ($\uparrow$) & AUCPR ($\uparrow$)



- Favorable order: Decreasing permutation weights $\{\boldsymbol{w}_i\}_1^D$ according to a total order admitting $\boldsymbol{G}^*$

- Uninformative order: same permutation weight $\boldsymbol{w}_i$ for each $\boldsymbol{X}_i$

- Adverse order: reversed favorable order

# Conclusion

- Introduction to CSL and Bayesian models for it

- Probability distribution over DAGs that enables differentiable sampling (DPM-DAG)

- Edge-wise priors in Bayesian CSL can speed up convergence w.r.t sample efficiency

- Using DPM-DAG for both models allows to reuse the posterior as the next prior

# Thank you very much for your Attention & Interest

Invited Talk on Bayesian Causal Structure Learning