

What?

We improve continuous Bayesian networks with normalizing flows.

How?

With a new conditioner that generalizes coupling and autoregressive conditioners.

Applications?

Learning the topology of Bayesian networks, Make normalizing flows more interpretable.



Code

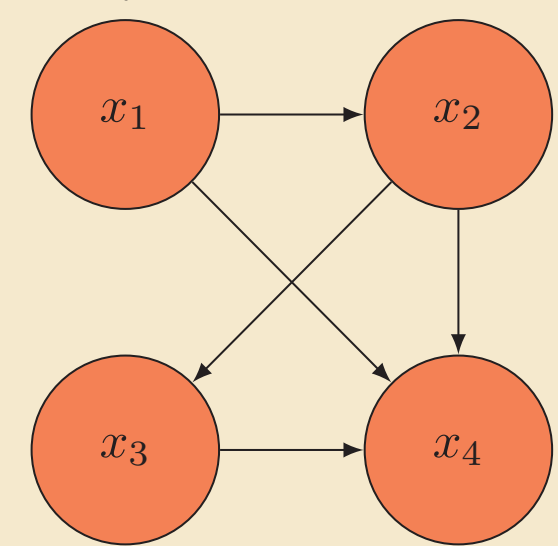


Arxiv

2006.02548

Bayesian Networks

A Bayesian network is a directed acyclic graph that factorizes the model distribution as $p(\mathbf{x}) = \prod_{i=1}^D p(x_i | \mathcal{P}_i)$.



$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_3|x_2)p(x_4|x_1, x_2, x_3)$$

Pros

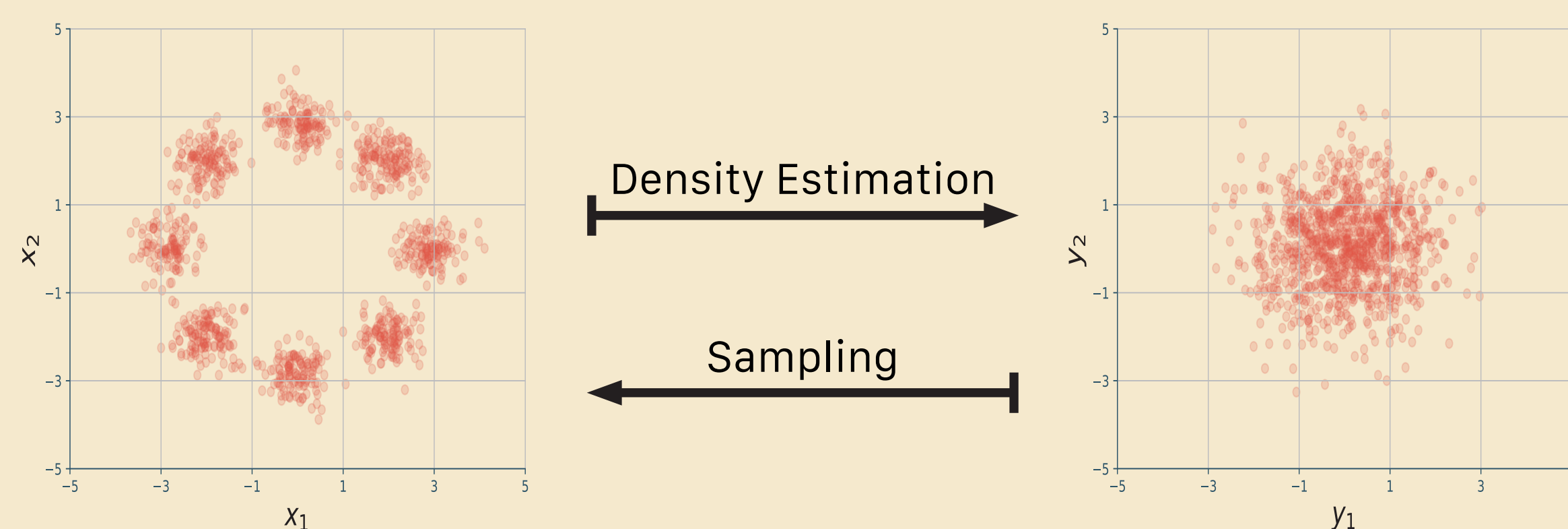
- Good for modeling independencies and checking their impact on the modeled density.
- Applications across science and technology.

Cons

- Often used with discrete or discretized data.
- Outdated with respect to the deep learning revolution.

Normalizing Flows

A normalizing flow is a sequence of K invertible transformation steps $g_k: \mathbb{R}^d \rightarrow \mathbb{R}^d$ composed together to create an expressive invertible mapping $g(\mathbf{x})$. Density estimation is performed via the change of variables theorem: $p(\mathbf{x}) = p_{\mathbf{y}}(g(\mathbf{x})) \left| \det J_{g(\mathbf{x})} \right|$.



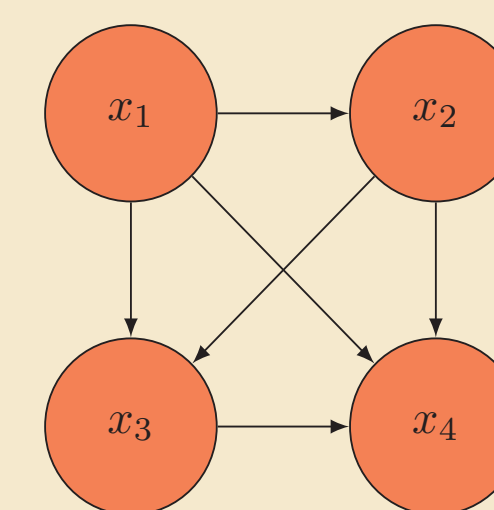
A transformation step can be written as: $g(\mathbf{x}) = [g^1(x_1; \mathbf{c}^1(\mathbf{x})) \dots g^d(x_d; \mathbf{c}^d(\mathbf{x}))]^T$.

Normalizer
(e.g. affine or monotonic)

Conditioner
(ensure a simple Jacobian)

Graphical Conditioners

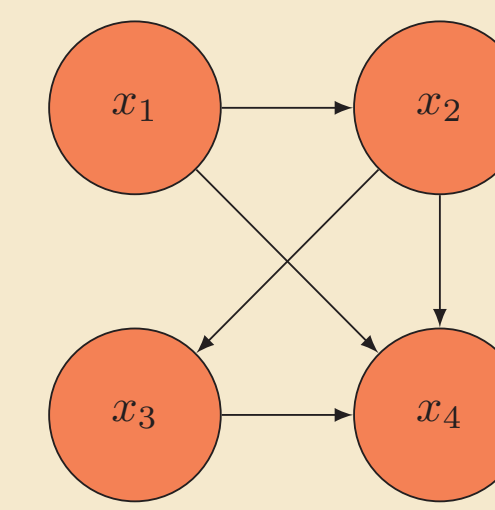
Autoregressive conditioner



$$\mathbf{c}^i(\mathbf{x}) = \mathbf{h}^i([x_1 \dots x_{i-1}]^T)$$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |
| 2 | 1 | 0 | 0 | 0 |
| 3 | 1 | 1 | 0 | 0 |
| 4 | 1 | 1 | 1 | 0 |

Graphical conditioner



$$\mathbf{c}^i(\mathbf{x}) = \mathbf{h}^i(\mathbf{x} \odot A_{i,:})$$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 |
| 2 | 1 | 0 | 0 | 0 |
| 3 | 0 | 1 | 0 | 0 |
| 4 | 1 | 1 | 1 | 0 |

Generalization

These conditioners can be generalized to any Bayesian network topology. Formally, let $A \in \{0, 1\}^{d \times d}$ be the adjacency of a Bayesian network, the graphical conditioner is defined as $\mathbf{c}^i(\mathbf{x}) = \mathbf{h}^i(\mathbf{x} \odot A_{i,:})$.

Topology learning

Choosing the network topology is not always easy!

But learning a good topology can be cast as a continuous optimization problem:

$$\max_{A \in \mathbb{R}^{d \times d}} F(A) \quad \Longleftrightarrow \quad \max_{A \in \mathbb{R}^{d \times d}} F(A) \quad \text{where} \quad F(A) = \sum_{j=1}^N \log(p(\mathbf{x}^j)) + \lambda_{\ell_1} \|A\|_1$$

s.t. $\mathcal{G}(A) \in \text{DAGs}$ s.t. $w(A) = 0$

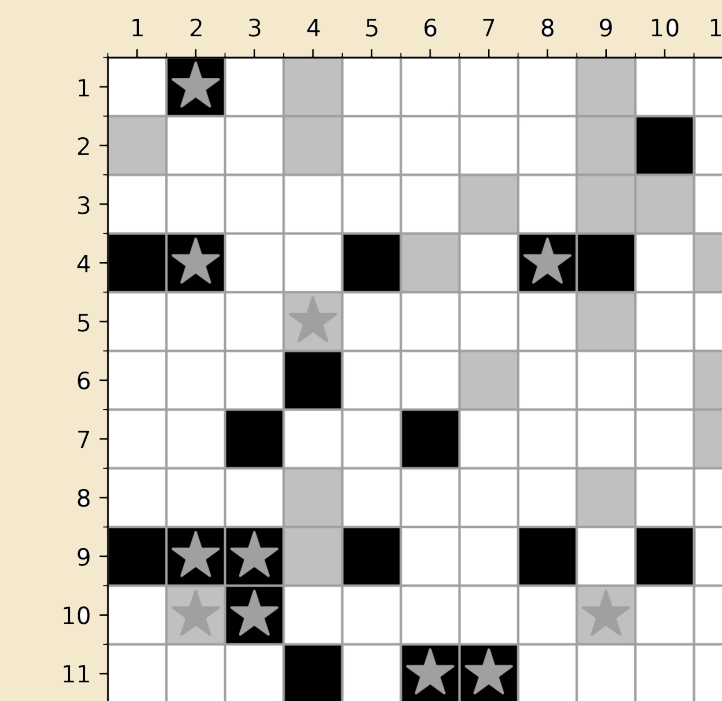
The log-density is evaluated with a graphical normalizing flow, introducing the neural networks parameters it can be written as $p(\mathbf{x}; \theta) = p_{\mathbf{z}}(g(\mathbf{x}; \theta)) \prod_{i=1}^d \left| \frac{\partial g^i(x_i; \mathbf{h}^i(\mathbf{x} \odot A_{i,:}), \theta)}{\partial x_i} \right|$.

The acyclicity constraint is expressed as $w(A) := \text{tr}((I + \alpha A)^d) - d \propto \text{tr}\left(\sum_{k=1}^d \alpha^k A^k\right)$.

Lagrangian formulation: $\max_A \mathbb{E}_{\gamma_1, \gamma_2} [F(A)] - \lambda_t w(A) - \frac{\mu_t}{2} w(A)^2$.

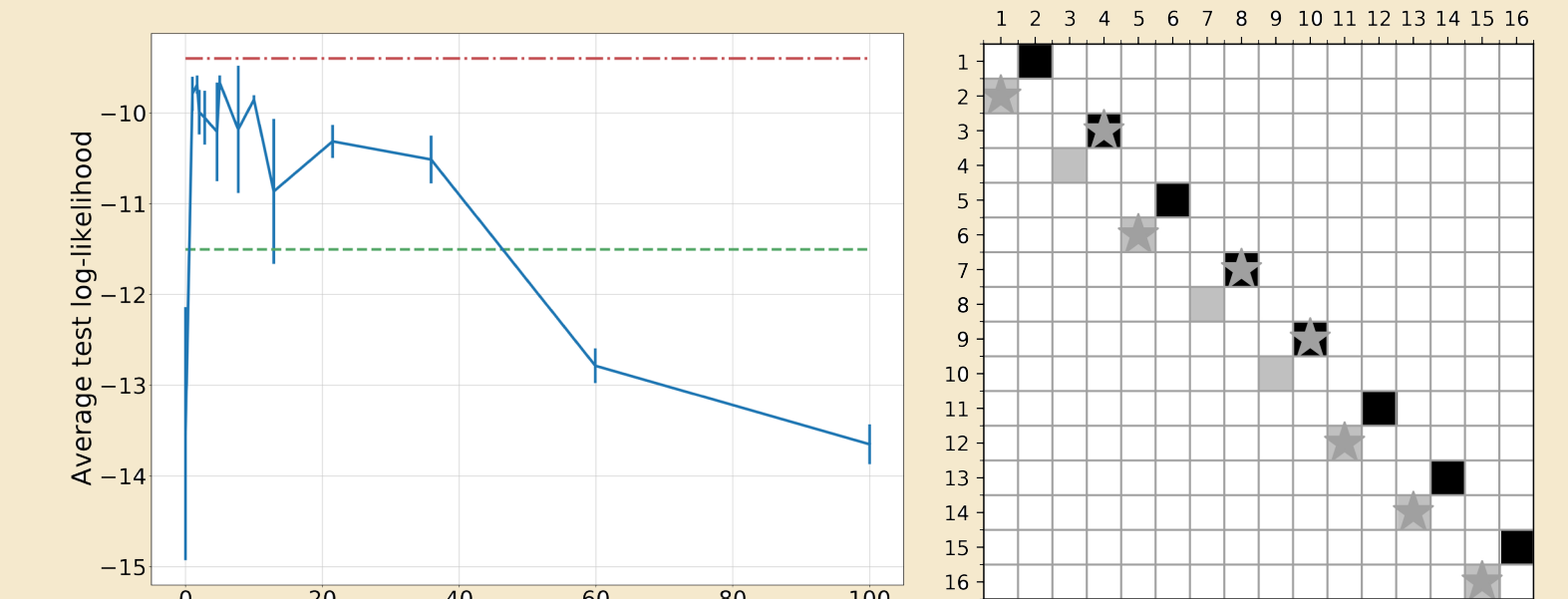
The importance of graph topology

Protein interaction network



On real data, the optimization recovers (stars) a graph that is similar to the one designed by experts.

8 pairs of 2D toy problems



(left) Test log-likelihood as a function of ℓ_1 -penalization. Red: Correct topology. Green: Autoregressive conditioner. (right) Recovered topology. Learning the right topology leads to better results than autoregressive conditioners.

Density estimation

| Dataset | d | V | Train | Test |
|--------------------|----|--------------|-----------|---------|
| Arithmetic Circuit | 8 | 8 | 10,000 | 5,000 |
| 8 Pairs | 16 | 8 | 10,000 | 5,000 |
| Tree | 7 | 8 | 10,000 | 5,000 |
| Protein | 11 | 20 | 6,000 | 1,466 |
| POWER | 6 | ≤ 15 | 1,659,917 | 204,928 |
| GAS | 8 | ≤ 28 | 852,174 | 105,206 |
| HEPMAS | 21 | ≤ 210 | 315,123 | 174,987 |
| MINIBOONE | 43 | ≤ 903 | 29,556 | 3,648 |
| BSDS300 | 63 | $\leq 1,953$ | 1,000,000 | 250,000 |

Prescribed topology

| Conditioner | Graphical | Autoreg. |
|--------------------|-------------|--------------|
| Arithmetic Circuit | 3.99 ± .16 | 3.06 ± .38 |
| 8 Pairs | -9.40 ± .06 | -11.50 ± .27 |
| Tree | -6.85 ± .02 | -6.96 ± .05 |
| Protein | 6.46 ± .08 | 7.52 ± .10 |

Learned topology

| Dataset | POWER | GAS | HEPMAS | MINIBOONE | BSDS300 |
|-----------|-------------|-------------|--------------|--------------|--------------|
| Coup. | -5.60 ± .00 | -3.05 ± .01 | -25.74 ± .01 | -38.34 ± .02 | 57.33 ± .00 |
| (a) Auto. | -3.55 ± .00 | -0.34 ± .01 | -21.66 ± .01 | -16.70 ± .05 | 63.74 ± .00 |
| Graph. | -2.80 ± .01 | 1.99 ± .02 | -21.18 ± .07 | -19.67 ± .06 | 62.85 ± .07 |
| Coup. | 0.25 ± .00 | 5.12 ± .03 | -20.55 ± .04 | -32.04 ± .12 | 107.17 ± .46 |
| (b) Auto. | 0.58 ± .00 | 9.79 ± .04 | -14.52 ± .16 | -11.66 ± .02 | 151.29 ± .31 |
| Graph. | 0.62 ± .04 | 10.15 ± .15 | -14.17 ± .13 | -16.23 ± .52 | 155.22 ± .11 |

We tested density estimation performance on many diverse datasets. Graphical conditioners outperform coupling and autoregressive flows.

Take home messages

- Continuous Bayesian networks can be combined with deep generative models.
- A correct prescribed topology improves the performance of normalizing flows.
- It is possible to discover relevant Bayesian network topology with graphical normalizing flows.