Causal Normalizing flow

January 08, 2024

Seoul National University

1 Contribution

- 2 Structural causal model(SCM)
- **3** Autoregressive normalizing flow (ANF)
- 4 Causal Autoregressive normalizing flow (CAF)
- 5 Intervention and Counterfactual

- Introduce causal autoregressive flows(CAFs) which reflect the causal structure of data.
- Through CAFs, enable sampling an intervened sample and counterfactual sample.

Contribution



Autoregressive normalizing flow (ANF)

4 Causal Autoregressive normalizing flow (CAF)

Intervention and Counterfactual

Structural causal model

A structural causal model (SCM) is a tuple M = (f̃, P_u) describing a data-generating process that transforms a set of d-dimensional latent variables, u ~ P_u, into a set of d-dimensional observed data, x, according to f̃ = (f̃₁, ..., f̃_d) : ℝ^d → ℝ^d.

• Specifically,
$$x = (x_1, \cdots, x_d)$$
 is computed as follows:

$$\mathbf{u} := (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_d) \sim P_{\mathbf{u}}, \quad \mathbf{x}_i = \tilde{f}_i \left(\mathbf{x}_{\mathrm{pa}_i}, \mathbf{u}_i \right), \quad \text{ for } i = 1, 2, \dots, d.$$

where x_{pa_i} is parents of x_i which directly cause x_i .

• Define the adjacency matrix of the causal graph as $A = (A_{ij}) \in \{0, 1\}^{d \times d}$,

$$A_{ij} = \mathbb{I}\left(\frac{\partial \tilde{f}_i(\mathsf{x}_{\mathrm{pa}_i},\mathsf{u}_i)}{\partial x_j} \neq 0\right)$$

where ${\ensuremath{\mathbb I}}$ is the indicator function.

- We assume that A is acyclic.
- Also assume that the x are sorted according to a causal ordering.

Contribution



3 Autoregressive normalizing flow (ANF)

- 4 Causal Autoregressive normalizing flow (CAF)
- 5 Intervention and Counterfactual

- Normalizing flows are generative models that express the probability density of data using the change-of-variables rule.
- Given an observed data $x \in \mathbb{R}^d$, a normalizing flow is a neural network with parameters θ that takes x as input, and outputs

$$T_{\theta}(\mathbf{x}) =: \mathbf{u} \sim P_{\mathbf{u}}$$

$$\log p(\mathbf{x}) = \log p(T_{\theta}(\mathbf{x})) + \log |\det (\nabla_{\mathbf{x}} T_{\theta}(\mathbf{x}))|$$

where P_u is a base distribution that is easy to sample.

• In normalizing flow, estimate θ by maximizing log p(x).

Autoregressive normalizing flow

- To sample data from u, T⁻¹_θ(u) should be exist and determinant of Jacobian matrix ∇_xT_θ(x) should be tractable.
- Autoregressive normalizing flows(ANFs) are models that satisfy the above conditions.
- ANFs are deep neural networks which is *i*-th output of each layer *l* denoted by z^l_i is computed as

$$\mathsf{z}_i' := au_i' \left(\mathsf{z}_i'^{-1}; \mathsf{h}_i'
ight), \quad ext{ where } \quad \mathsf{h}_i' := extsf{c}_i' \left(\mathsf{z}_{1:i-1}'^{-1}
ight)$$

and where τ_i and c_i are termed the transformer and the conditioner, respectively.

• Note that Jacobian matrices of ANFs are triangular matrices.





3 Autoregressive normalizing flow (ANF)

4 Causal Autoregressive normalizing flow (CAF)



Intervention and Counterfactual

- Want to express SCM as ANFs when A is known.
- One of the conditions that ANFs can be expressed as SCM is x_i is a function of x_{pai} which depend on u_{pai}, and u_i.

Theorem If a causal NF T_{θ} satisfies an condition in previous slide, then $\nabla_{x} T_{\theta}(x) \equiv I - A$ and $\nabla_{u} T_{\theta}^{-1}(u) \equiv I + \sum_{n=1}^{\text{diam}(A)} A^{n}$, where A is the causal adjacency matrix of \mathcal{M} .

where diam(A) = min{ $k : A^k = \tilde{0}$ } when $\tilde{0}$ is zero matrix.

• (Abductive model) Causal NF from x to u

$$\mathbf{z}_{i}^{\prime}= au_{i}\left(\mathbf{z}_{i}^{\prime-1};\mathbf{h}_{i}^{\prime}
ight), \hspace{0.5cm} ext{where} \hspace{0.5cm} \mathbf{h}_{i}^{\prime}=c_{i}\left(\mathbf{z}_{ ext{pa}_{i}}^{\prime-1}
ight)$$

 To penalize spurious correlations from x to T_θ(x), add penalized term : minimize E_x [-log p (T_θ(x)) + ||∇_xT_θ(x) ⊙ (ĩ - A)||₂] where ĩ is a metric of energy

where $\tilde{1}$ is a matrix of ones.

Contribution

- 2 Structural causal model(SCM)
- **3** Autoregressive normalizing flow (ANF)
- 4 Causal Autoregressive normalizing flow (CAF)
- **5** Intervention and Counterfactual

Intervention

 The do-operator, denoted as do (x_i = α), is a mathematical operator that fixes the observational value x_i = α, and thus removes any causal dependency on x_i.

Algorithm 1 Algorithm to sample from the interventional distribution, $P(\mathbf{x} | do(\mathbf{x}_i = \alpha))$.

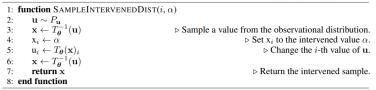


Figure 1: Do-operation

Algorithm 2 Algorithm to sample from the counterfactual distribution, $P(\mathbf{x}^{cf} | do(\mathbf{x}_i = \alpha), \mathbf{x}^{f})$.

1: **function** GETCOUNTERFACTUAL($\mathbf{x}^{f}, i, \alpha$) 2: $\mathbf{u} \leftarrow T_{\boldsymbol{\theta}}(\mathbf{x}^{\mathrm{f}})$ \triangleright Get u from the factual sample. $\mathbf{x}_i^{\mathrm{f}} \leftarrow \alpha$ \triangleright Set \mathbf{x}_i to the intervened value α . 3: 4: $\mathbf{u}_i \leftarrow T_{\boldsymbol{\theta}}(\mathbf{x}^{\mathrm{f}})_i$ \triangleright Change the *i*-th value of **u**. $\mathbf{x}^{cf} \leftarrow T_{\boldsymbol{\theta}}^{-1}(\mathbf{u})$ 5: return x^{cf} 6: Return the counterfactual value. 7: end function

Figure 2: Counterfactual sample

Contribution

- 2 Structural causal model(SCM)
- **3** Autoregressive normalizing flow (ANF)
- 4 Causal Autoregressive normalizing flow (CAF)
 - 5 Intervention and Counterfactual



3-CHAIN_{LIN}:

$$\begin{split} \tilde{f}_1(\mathbf{u}_1) &= \mathbf{u}_1 \\ \tilde{f}_2(\mathbf{x}_1, \mathbf{u}_2) &= 10 \cdot \mathbf{x}_1 - \mathbf{u}_2 \\ \tilde{f}_3(\mathbf{x}_2, \mathbf{u}_3) &= 0.25 \cdot \mathbf{x}_2 + 2 \cdot \mathbf{u}_3 \end{split}$$

3-CHAIN_{NLIN}:

$$\begin{split} \tilde{f}_1(\mathbf{u}_1) &= \mathbf{u}_1 \\ \tilde{f}_2(\mathbf{x}_1, \mathbf{u}_2) &= e^{\mathbf{x}_1/2} + \mathbf{u}_2/4 \\ \tilde{f}_3(\mathbf{x}_2, \mathbf{u}_3) &= \frac{(\mathbf{x}_2 - 5)^3}{15} + \mathbf{u}_3 \end{split}$$

4-CHAIN_{LIN}:

$$\begin{split} \tilde{f}_1(\mathbf{u}_1) &= \mathbf{u}_1 \\ \tilde{f}_2(\mathbf{x}_1, \mathbf{u}_2) &= 5 \cdot \mathbf{x}_1 - \mathbf{u}_2 \\ \tilde{f}_3(\mathbf{x}_2, \mathbf{u}_3) &= -0.5 \cdot \mathbf{x}_2 - 1.5 \cdot \mathbf{u}_3 \\ \tilde{f}_4(\mathbf{x}_3, \mathbf{u}_4) &= \mathbf{x}_3 + \mathbf{u}_4 \end{split}$$

Figure 3: Counterfactual sample

		Performance			Time Evaluation (µs)			
Dataset	Model	\mathbf{KL}	ATE _{RMSE}	CF _{RMSE}	Training	Evaluation	Sampling	
3-CHAIN LIN [36]	Causal NF CAREFL [†] VACA	$\begin{array}{c} 0.00_{0.00} \\ 0.00_{0.00} \\ 4.44_{1.03} \end{array}$	0.05 _{0.01} 0.20 _{0.13} 5.76 _{0.07}	$\begin{array}{c} \textbf{0.04_{0.01}} \\ 0.20_{0.09} \\ 4.98_{0.10} \end{array}$	$\begin{array}{c} \mathbf{0.41_{0.06}} \\ 0.68_{0.24} \\ 36.19_{1.54} \end{array}$	0.48_{0.10} 0.97 _{0.33} 28.33 _{0.72}	$\begin{array}{r} \textbf{0.76}_{0.06} \\ 1.94_{0.77} \\ 75.34_{4.58} \end{array}$	
3-CHAIN NLIN [36]	Causal NF CAREFL [†] VACA	0.00 _{0.00} 0.00 _{0.00} 12.82 _{1.00}	$\begin{array}{c} \mathbf{0.03_{0.01}} \\ \mathbf{0.05_{0.02}} \\ 1.54_{0.03} \end{array}$	$\begin{array}{c} \mathbf{0.02_{0.01}} \\ 0.04_{0.02} \\ 1.32_{0.02} \end{array}$	$\begin{array}{c} \mathbf{0.52_{0.06}}\\ \mathbf{0.60_{0.22}}\\ \mathbf{39.45_{4.12}} \end{array}$	$\begin{array}{c} \mathbf{0.56_{0.03}} \\ 0.84_{0.22} \\ 30.93_{2.30} \end{array}$	$\frac{1.02_{0.05}}{1.66_{0.41}}$ $84.36_{9.60}$	
4-CHAIN LIN	Causal NF CAREFL [†] VACA	0.00 _{0.00} 0.00 _{0.00} 13.14 _{0.73}	0.07 _{0.02} 0.16 _{0.07} 3.82 _{0.01}	$\begin{array}{c} \textbf{0.04_{0.01}} \\ 0.14_{0.04} \\ 3.72_{0.05} \end{array}$	$\begin{array}{c} \mathbf{0.56_{0.08}} \\ \mathbf{0.70_{0.28}} \\ \mathbf{61.85_{5.06}} \end{array}$	$\begin{array}{c} \mathbf{0.62_{0.15}} \\ 0.99_{0.20} \\ 49.31_{4.11} \end{array}$	$\frac{1.54_{0.40}}{2.85_{0.54}}$ 92.067.93	

Figure 4: Result 1

German data set when a sensitive variable is sex.

Table 3: Accuracy, F1-score, and counterfactual unfairness of the audited classifiers. Causal NFs enable both fair classifiers and accurate unfairness metrics. Results are averaged on five runs.

	Logistic classifier				SVM classifier			
	full	unaware	fair \mathbf{x}	fair \mathbf{u}	full	unaware	fair \mathbf{x}	fair \mathbf{u}
f1	$72.28_{6.16}$	$72.37_{4.90}$	$59.66_{8.57}$	$73.08_{4.38}$	$76.04_{2.86}$	$76.80_{5.82}$	$68.28_{5.74}$	$77.39_{1.52}$
accuracy	$67.00_{3.83}$	$66.75_{2.63}$	$54.75_{5.91}$	$66.50_{3.70}$	$69.50_{3.11}$	$71.00_{3.83}$	$59.25_{2.99}$	$69.75_{1.26}$
unfairness	$5.84_{2.93}$	$2.81_{0.72}$	$0.00_{0.00}$	$0.00_{0.00}$	$6.65_{2.45}$	$2.78_{0.40}$	$0.00_{0.00}$	$0.00_{0.00}$

Figure 5: Result 2

