



subtractive mixture models

representation, learning & inference

antonio vergari (he/him)

 @tetraduzione

18th Feb 2025 - Flatiron Institute NYC

april

`april-tools.github.io`

april

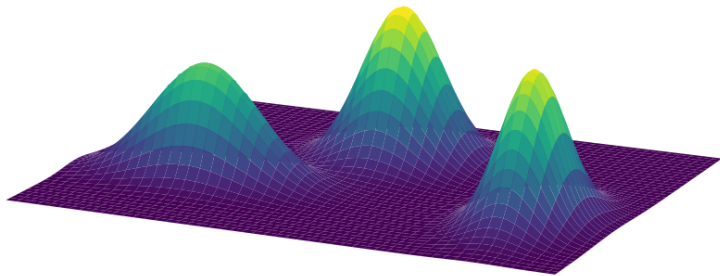
about
probabilities
integrals &
logic

april

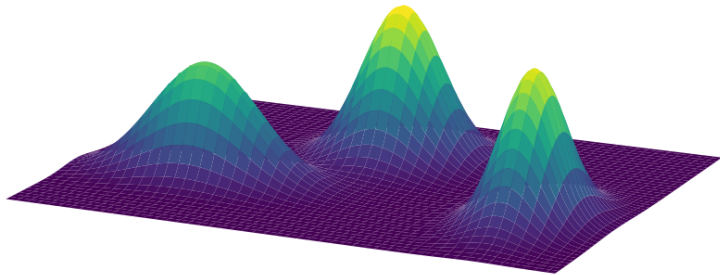
***autonomous &
provably
reliable
intelligent
learners***

april

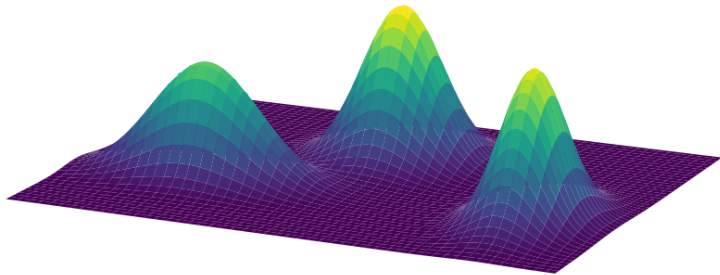
*april is
probably a
recursive
identifier of a
lab*



who knows mixture models?



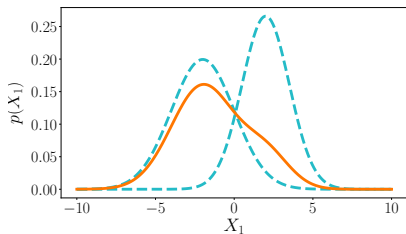
*who **loves** mixture models?*



$$c(\mathbf{X}) = \sum_{i=1}^K w_i c_i(\mathbf{X}), \quad \text{with } w_i \geq 0, \quad \sum_{i=1}^K w_i = 1$$

GMMs

as computational graphs

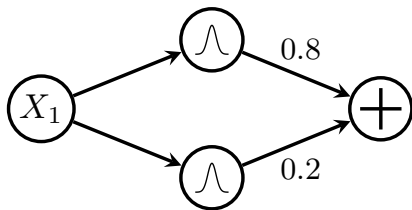


$$p(X) = w_1 \cdot p_1(X_1) + w_2 \cdot p_2(X_1)$$

⇒ *translating inference to data structures...*

GMMs

as computational graphs

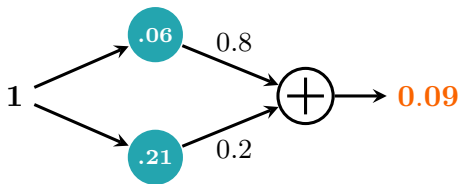


$$p(X_1) = 0.2 \cdot p_1(X_1) + 0.8 \cdot p_2(X_1)$$

\Rightarrow ...e.g., as a weighted sum unit over Gaussian input distributions

GMMs

as computational graphs

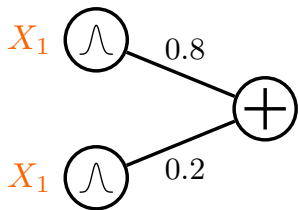


$$p(X = 1) = 0.2 \cdot p_1(X_1 = 1) + 0.8 \cdot p_2(X_1 = 1)$$

\Rightarrow inference = feedforward evaluation

GMMs

as computational graphs

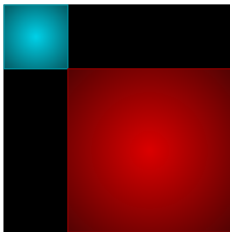


A simplified notation:

\Rightarrow **scopes** attached to inputs
 \Rightarrow edge directions omitted

GMMs

as computational graphs

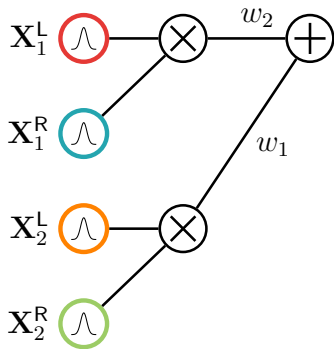


$$p(\mathbf{X}) = w_1 \cdot p_1(\mathbf{X}_1^L) \cdot p_1(\mathbf{X}_1^R) + \\ w_2 \cdot p_2(\mathbf{X}_2^L) \cdot p_2(\mathbf{X}_2^R)$$

\Rightarrow local factorizations...

GMMs

as computational graphs



$$p(\mathbf{X}) = w_1 \cdot p_1(\mathbf{X}_1^L) \cdot p_1(\mathbf{X}_1^R) + w_2 \cdot p_2(\mathbf{X}_2^L) \cdot p_2(\mathbf{X}_2^R)$$

\Rightarrow ...are product units

probabilistic circuits (PCs)

a grammar for tractable computational graphs

I. A simple tractable function is a circuit



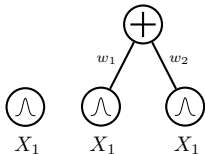
*e.g., a multivariate Gaussian or
orthonormal polynomial*



probabilistic circuits (PCs)

a grammar for tractable computational graphs

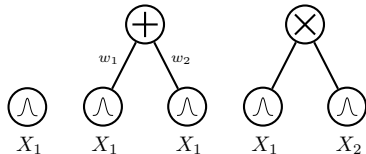
- I. *A simple tractable function is a circuit*
- II. *A weighted combination of circuits is a circuit*



probabilistic circuits (PCs)

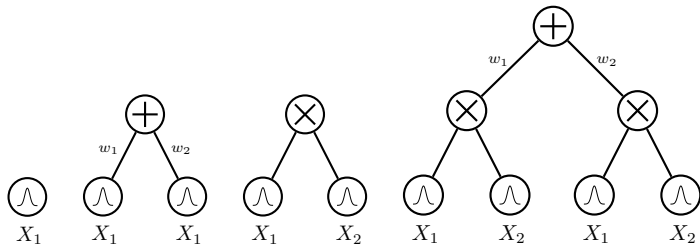
a grammar for tractable computational graphs

- I. *A simple tractable function is a circuit*
- II. *A weighted combination of circuits is a circuit*
- III. *A product of circuits is a circuit*



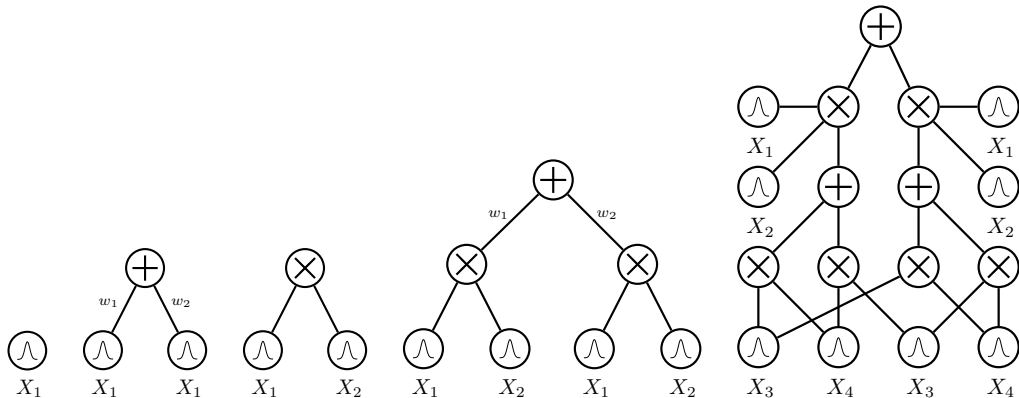
probabilistic circuits (PCs)

a grammar for tractable computational graphs

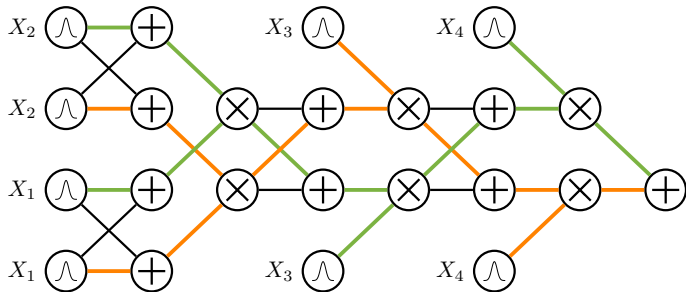


probabilistic circuits (PCs)

a grammar for tractable computational graphs

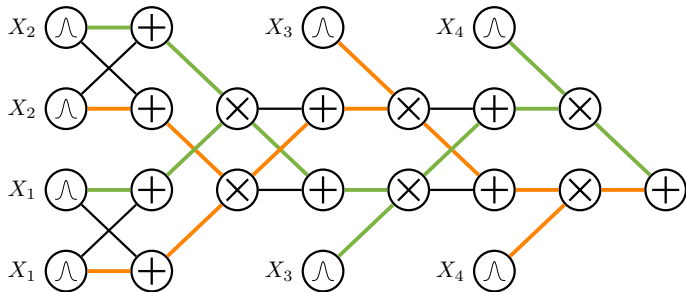


deep mixtures



$$p(\mathbf{x}) = \sum_{\mathcal{T}} \left(\prod_{w_j \in \mathbf{w}_{\mathcal{T}}} w_j \right) \prod_{l \in \text{leaves}(\mathcal{T})} p_l(\mathbf{x})$$

deep mixtures



an exponential number of mixture components!

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...

2. **Tractability** == **Structural Properties**!!!

Exact computations of reasoning tasks are certified by guaranteeing certain structural properties. #marginals #expectations #MAP, #product ...

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...

2. **Tractability** == **Structural Properties**!!!

Exact computations of reasoning tasks are certified by guaranteeing certain structural properties. #marginals #expectations #MAP, #product ...

3. **Reliable neuro-symbolic AI**

logical constraints as circuits, multiplied to probabilistic circuits

circuits

(and variants)

everywhere

Semantic Probabilistic Layers for Neuro-Symbolic Learning

Kareem Ahmed
CS Department
UCLA
ahmedk@cs.ucla.edu

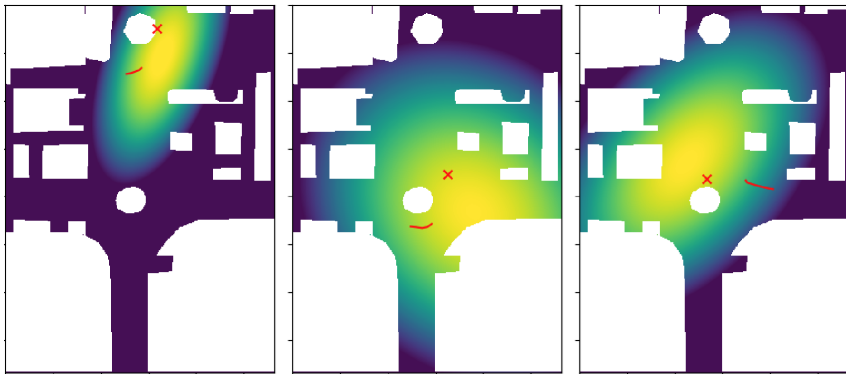
Stefano Teso
CIMEC and DISI
University of Trento
stefano.teso@unitn.it

Kai-Wei Chang
CS Department
UCLA
kwchang@cs.ucla.edu

Guy Van den Broeck
CS Department
UCLA
guyvdb@cs.ucla.edu

Antonio Vergari
School of Informatics
University of Edinburgh
avergari@ed.ac.uk

enforce constraints in neural networks at NeurIPS 2022



$$p(\mathbf{y} \mid \mathbf{x}) = \mathbf{q}_{\Theta}(\mathbf{y} \mid g(\mathbf{z})) \cdot \mathbf{c}_{\mathbf{K}}(\mathbf{x}, \mathbf{y}) / \mathbf{Z}(\mathbf{x})$$

SPL



Ground Truth



ResNet-18



Semantic Loss

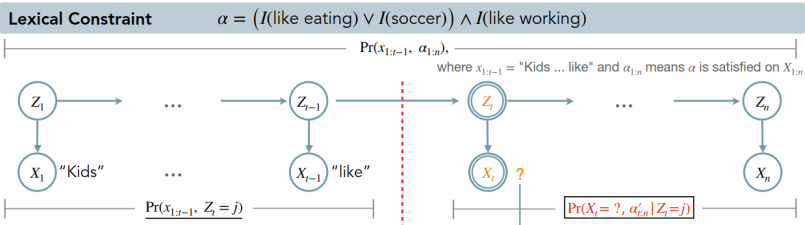


circuits

predictions guarantee a logical constraint 100% of the time!

Tractable Control for Autoregressive Language Generation

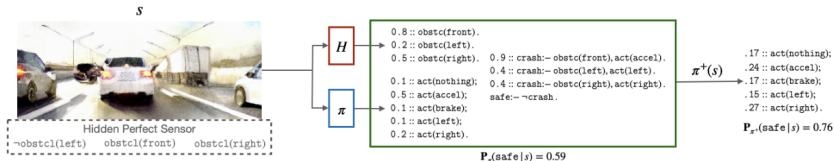
Honghua Zhang^{*1} Meihua Dang^{*1} Nanyun Peng¹ Guy Van den Broeck¹



constrained text generation with LLMs (ICML 2023)

Safe Reinforcement Learning via Probabilistic Logic Shields

Wen-Chi Yang¹, Giuseppe Marra¹, Gavin Rens and Luc De Raedt^{1,2}



reliable reinforcement learning (AAAI 23)

How to Turn Your Knowledge Graph Embeddings into Generative Models

Lorenzo Loconte
University of Edinburgh, UK
l.loconte@sms.ed.ac.uk































Nicola Di Mauro
University of Bari, Italy
nicola.dimauro@uniba.it

Robert Peharz
TU Graz, Austria
robert.peharz@tugraz.at

Antonio Vergari
University of Edinburgh, UK
avergari@ed.ac.uk

***enforce constraints in knowledge graph embeddings
oral at NeurIPS 2023***

Logically Consistent Language Models via Neuro-Symbolic Integration

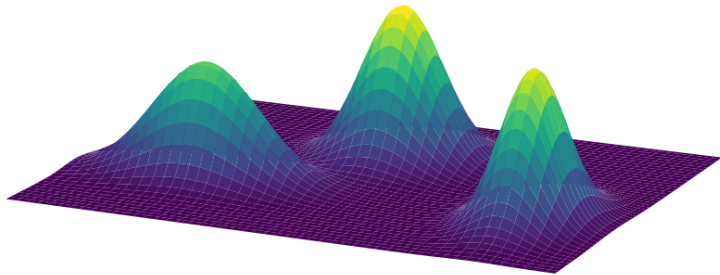
 = LLaMa 2  = LLaMa 2  = LLaMa 2	Forward Implication	Reverse Implication	Negation
	$A \rightarrow B$ A: (albatross, isA, bird) B: (albatross, isA, fish)	$\neg B \rightarrow \neg A$ B: (albatross, isNotA, organism) A: (albatross, isNotA, living thing)	$A \leftrightarrow A$ A: (computer, isA, airplane) A: (computer, isNotA, airplane)
	<div>Is an albatross a bird? </div> <div> Yes.</div> <div>Is an albatross a fish? </div> <div> Yes. Logical:  Factual: </div> <div> No. Logical:  Factual: </div>	<div>Is it true that an albatross is not an organism? </div> <div> No.</div> <div>Is it true that an albatross is not a living thing? </div> <div> Yes. Logical:  Factual: </div> <div> No. Logical:  Factual: </div>	<div>Is a computer a airplane? </div> <div> No.</div> <div>Is it true that a computer is not a airplane? </div> <div> No. Logical:  Factual: </div> <div> Yes. Logical:  Factual: </div>

improving logical (self-)consistency in LLMs at ICLR 2025



learning & reasoning with circuits in pytorch

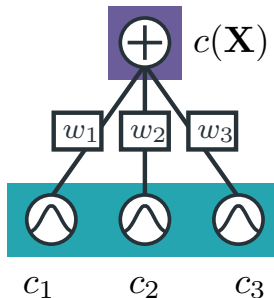
`github.com/april-tools/circuit`



$$c(\mathbf{X}) = \sum_{i=1}^K w_i c_i(\mathbf{X}), \quad \text{with } w_i \geq 0, \quad \sum_{i=1}^K w_i = 1$$

additive MMs

are so cool!



easily represented as shallow PCs

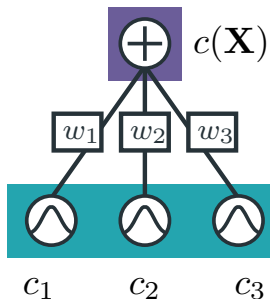
these are **monotonic** PCs

if marginals/conditionals are tractable for the components, they are tractable for the MM

universal approximators...

additive MMs

are so cool!



easily represented as shallow PCs

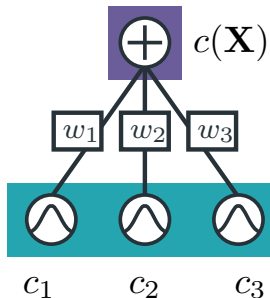
these are **monotonic** PCs

if marginals/conditionals are tractable for the components, they are tractable for the MM

universal approximators...

additive MMs

are so cool!



easily represented as shallow PCs

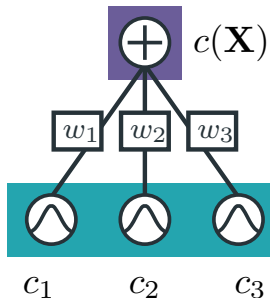
these are **monotonic** PCs

if marginals/conditionals are tractable for the components, they are tractable for the MM

universal approximators...

additive MMs

are so cool!



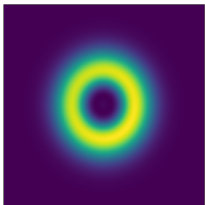
easily represented as shallow PCs

these are **monotonic** PCs

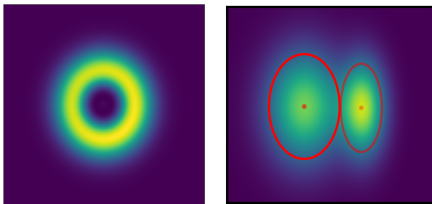
if marginals/conditionals are tractable for the components, they are tractable for the MM

universal approximators...

however...

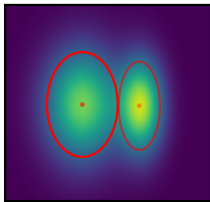
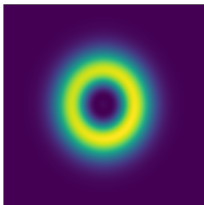


however...

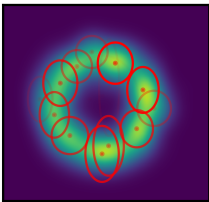


GMM ($K = 2$)

however...

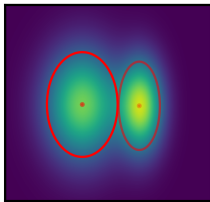
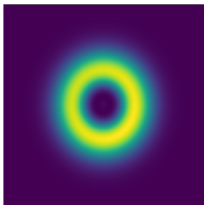


GMM ($K = 2$)

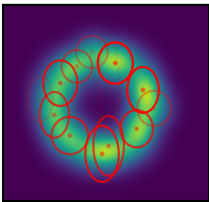


GMM ($K = 16$)

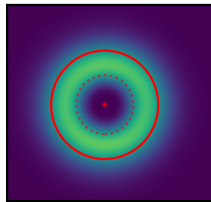
however...



GMM ($K = 2$)



GMM ($K = 16$)

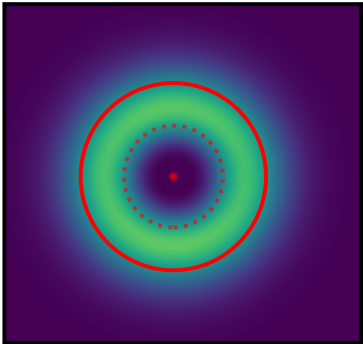


nGMM² ($K = 2$)

spoiler

**shallow mixtures
with negative parameters
can be *exponentially more compact* than
deep ones with positive parameters.**

subtractive MMs



also called negative/signed/**subtractive** MMs

⇒ or **non-monotonic** circuits,...

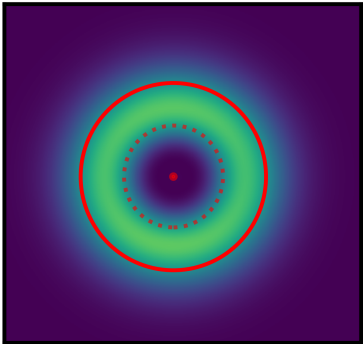
issue: how to preserve non-negative outputs?

well understood for simple parametric forms

e.g., Weibulls, Gaussians

⇒ constraints on variance, mean

subtractive MMs



also called negative/signed/**subtractive** MMs

⇒ or **non-monotonic** circuits,...

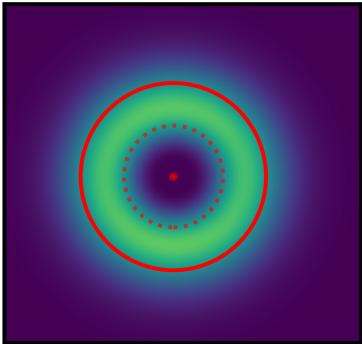
issue: how to preserve non-negative outputs?

well understood for simple parametric forms

e.g., Weibulls, Gaussians

⇒ constraints on variance, mean

subtractive MMs



also called negative/signed/**subtractive** MMs

⇒ or **non-monotonic** circuits,...

issue: how to preserve non-negative outputs?

well understood for simple parametric forms

e.g., Weibulls, Gaussians

⇒ *constraints on variance, mean*

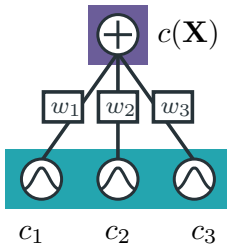
tl;dr

***“Understand when and how
we can use negative parameters
in deep **subtractive mixture models**”***

tl;dr

***“Understand when and how
we can use negative parameters
in deep **non-monotonic squared circuits**”***

subtractive MMs as circuits

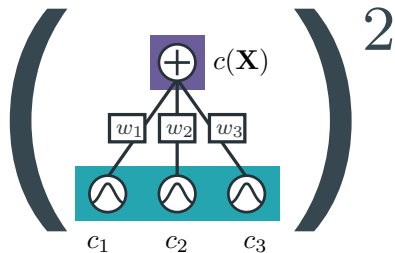


a **non-monotonic** smooth and (structured)
decomposable circuit

\Rightarrow possibly with negative outputs

$$c(\mathbf{X}) = \sum_{i=1}^K w_i c_i(\mathbf{X}), \quad w_i \in \mathbb{R},$$

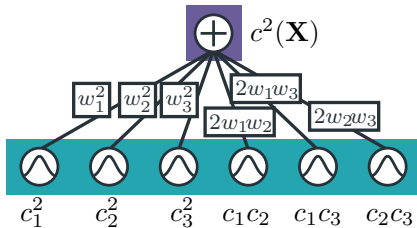
squaring shallow MMs



$$c^2(\mathbf{X}) = \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2$$

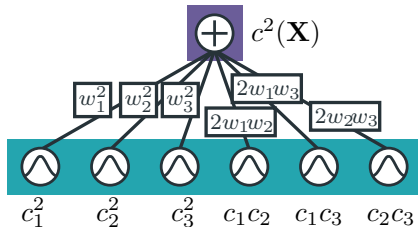
\Rightarrow ensure non-negative output

squaring shallow MMs



$$\begin{aligned} c^2(\mathbf{X}) &= \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2 \\ &= \sum_{i=1}^K \sum_{j=1}^K w_i w_j c_i(\mathbf{X}) c_j(\mathbf{X}) \end{aligned}$$

squaring shallow MMs

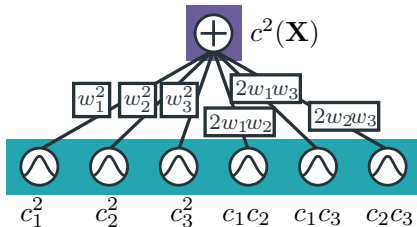


$$\begin{aligned} c^2(\mathbf{X}) &= \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2 \\ &= \sum_{i=1}^K \sum_{j=1}^K w_i w_j c_i(\mathbf{X}) c_j(\mathbf{X}) \end{aligned}$$

still a smooth and (str) decomposable PC with $\mathcal{O}(K^2)$ components!

\Rightarrow but still $\mathcal{O}(K)$ parameters

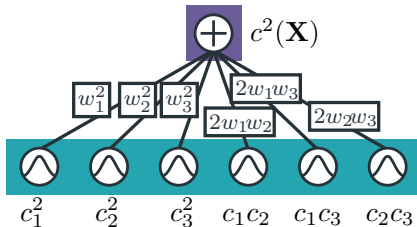
squaring shallow MMs



$$\begin{aligned} c^2(\mathbf{X}) &= \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2 \\ &= \sum_{i=1}^K \sum_{j=1}^K w_i w_j c_i(\mathbf{X}) c_j(\mathbf{X}) \end{aligned}$$

to **renormalize**, we have to compute $\sum_i \sum_j w_i w_j \int c_i(\mathbf{x}) c_j(\mathbf{x}) d\mathbf{x}$

squaring shallow MMs



$$\begin{aligned} c^2(\mathbf{X}) &= \left(\sum_{i=1}^K w_i c_i(\mathbf{X}) \right)^2 \\ &= \sum_{i=1}^K \sum_{j=1}^K w_i w_j c_i(\mathbf{X}) c_j(\mathbf{X}) \end{aligned}$$

to **renormalize**, we have to compute $\sum_i \sum_j w_i w_j \int c_i(\mathbf{x}) c_j(\mathbf{x}) d\mathbf{x}$
 \Rightarrow or we pick c_i, c_j to be **orthonormal**...!

EigenVI: score-based variational inference with orthogonal function expansions

Diana Cai

Flatiron Institute

dcai@flatironinstitute.org

Chirag Modi

Flatiron Institute

cmodi@flatironinstitute.org

Charles C. Margossian

Flatiron Institute

cmargossian@flatironinstitute.org

Robert M. Gower

Flatiron Institute

rgower@flatironinstitute.org

David M. Blei

Columbia University

david.blei@columbia.edu

Lawrence K. Saul

Flatiron Institute

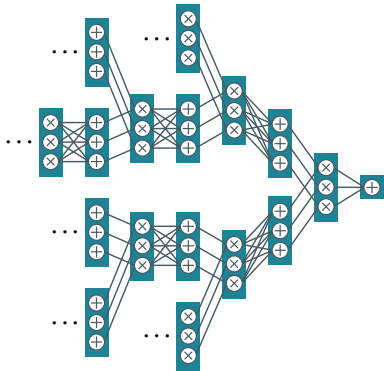
lsaul@flatironinstitute.org

orthonormal squared mixtures for VI

wait...

***“do negative parameters
really boost expressiveness?
and...always?”***

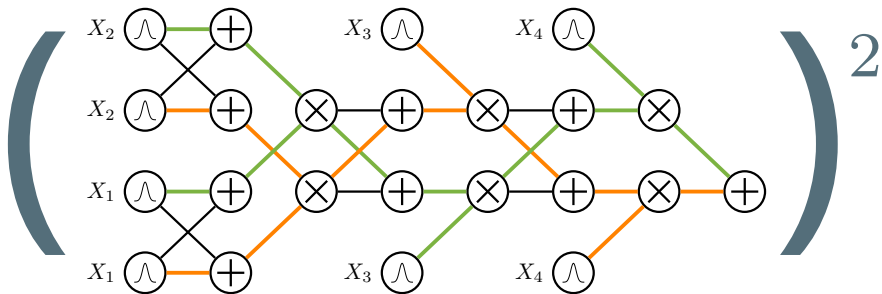
theorem



$$\left(\begin{matrix} \text{node} \\ \text{node} \\ \vdots \\ \text{node} \end{matrix} \right)^2$$

...but compact

squared non-monotonic circuits



how to efficiently square (and **renormalize**) a deep PC?

compositional inference /



```
1 from cirq.symbolic.functional import integrate, multiply
2
3 #
4 # create a deep circuit
5 c = build_symbolic_circuit('quad-tree-4')
6
7 #
8 # compute the partition function of  $c^2$ 
9 def renormalize(c):
10     c2 = multiply(c, c)
11     return integrate(c2)
```

probabilistic circuits (PCs)

the unit-wise definition

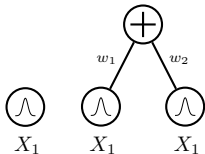
1. A simple tractable function is a circuit



probabilistic circuits (PCs)

the unit-wise definition

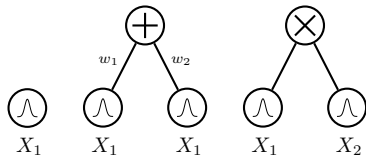
- I. A simple tractable function is a circuit
- II. A weighted combination of circuits is a circuit



probabilistic circuits (PCs)

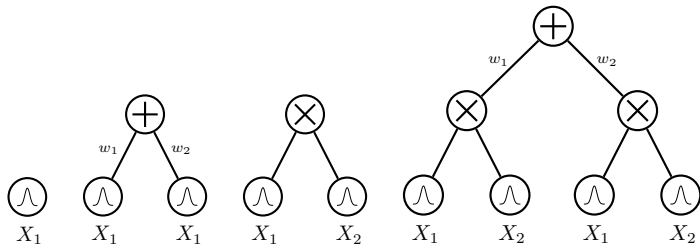
the unit-wise definition

- I. A simple tractable function is a circuit
- II. A weighted combination of circuits is a circuit
- III. A product of circuits is a circuit



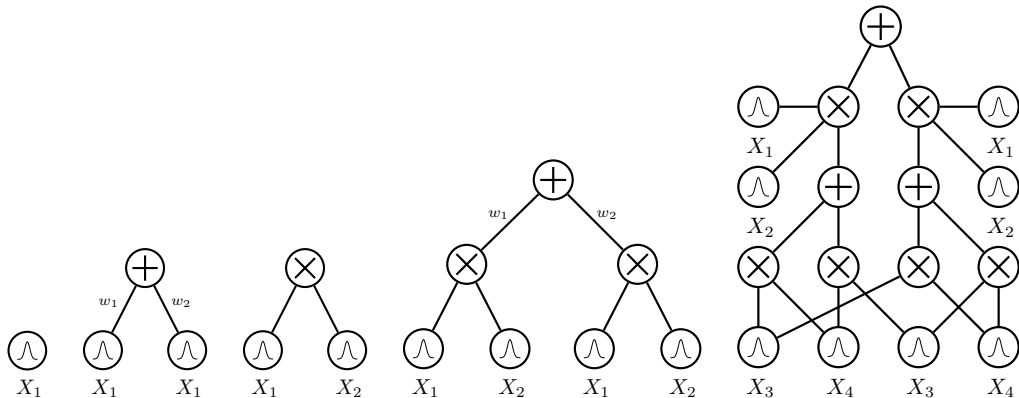
probabilistic circuits (PCs)

the unit-wise definition



probabilistic circuits (PCs)

the unit-wise definition



probabilistic circuits (PCs)

a tensorized definition

I. *A set of tractable functions is a circuit layer*



probabilistic circuits (PCs)

a tensorized definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer

$$c(\mathbf{x}) = \mathbf{W}l(\mathbf{x})$$



probabilistic circuits (PCs)

a tensorized definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer

$$c(\mathbf{x}) = \mathbf{W}l(\mathbf{x})$$



probabilistic circuits (PCs)

a tensorized definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \mathbf{l}(\mathbf{x}) \odot \mathbf{r}(\mathbf{x}) \quad // \text{Hadamard}$$



probabilistic circuits (PCs)

a tensorized definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \mathbf{l}(\mathbf{x}) \odot \mathbf{r}(\mathbf{x}) \quad // \text{Hadamard}$$

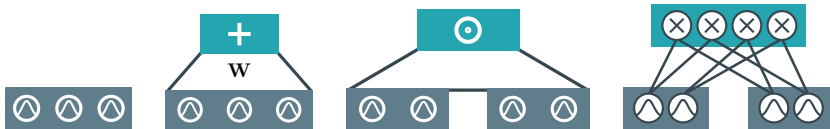


probabilistic circuits (PCs)

a tensorized definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \text{vec}(\mathbf{l}(\mathbf{x})\mathbf{r}(\mathbf{x})^\top) \quad // \text{Kronecker}$$



probabilistic circuits (PCs)

a tensorized definition

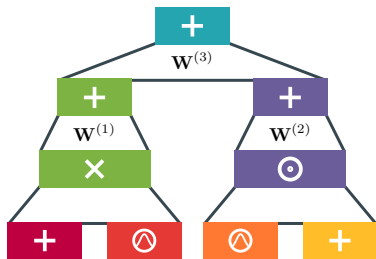
- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \text{vec}(\mathbf{l}(\mathbf{x})\mathbf{r}(\mathbf{x})^\top) \quad // \text{Kronecker}$$



probabilistic circuits (PCs)

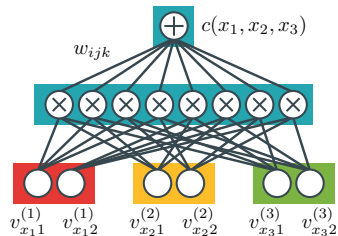
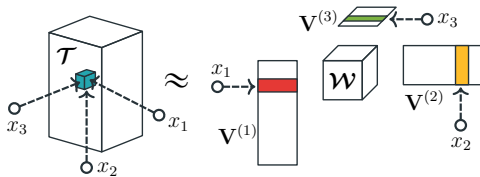
a tensorized definition



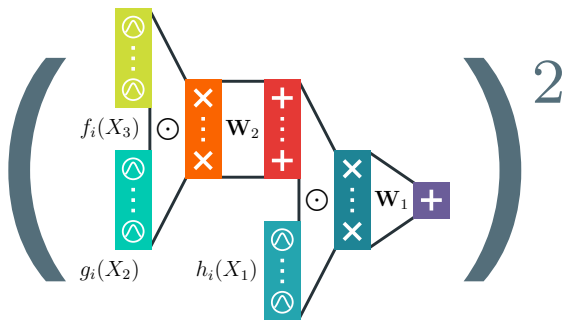
- I. A set of tractable functions is a circuit layer
 - II. A linear projection of a layer is a circuit layer
 - III. The product of two layers is a circuit layer
- stack layers to build a deep circuit!***

circuits layers

as tensor factorizations



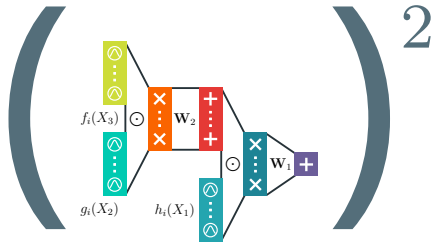
Loconte et al., "What is the Relationship between Tensor Factorizations and Circuits (and How Can We Exploit it)?", TMLR, 2025



how to efficiently square (and **renormalize**) a deep PC?

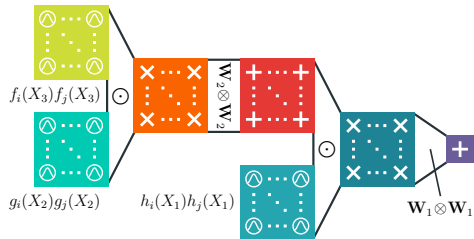
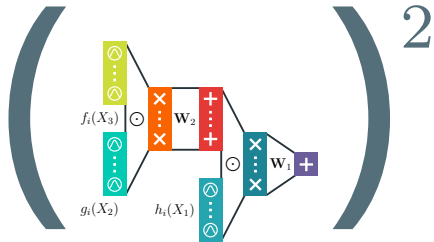
squaring deep PCs

the tensorized way



squaring deep PCs

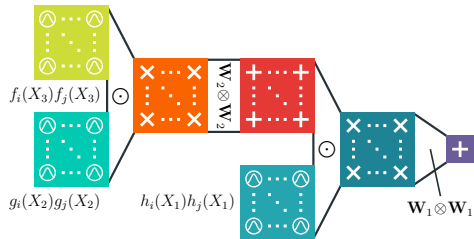
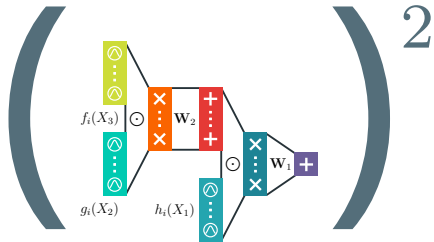
the tensorized way



squaring a circuit = squaring layers

squaring deep PCs

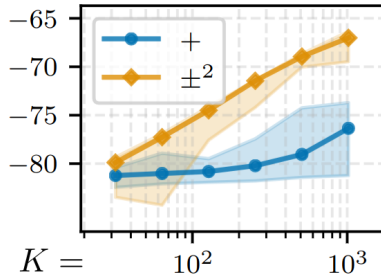
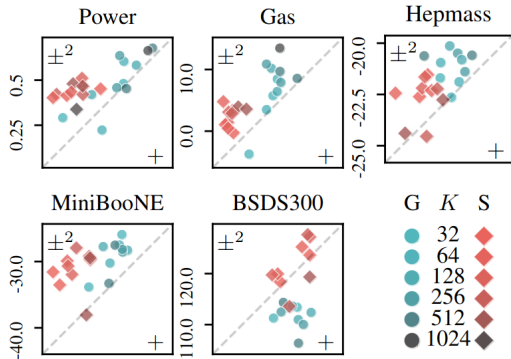
the tensorized way



exactly compute $\int \mathbf{c}(\mathbf{x}) \mathbf{c}(\mathbf{x}) d\mathbf{X}$ in time $O(LK^2)$

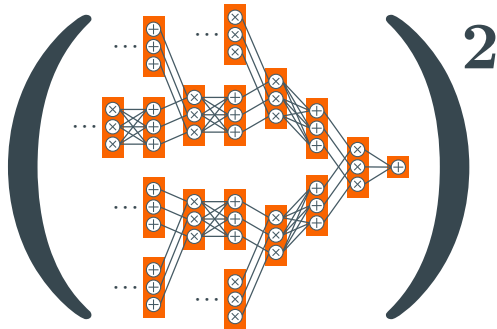
how more expressive?

for the ML crowd



theorem

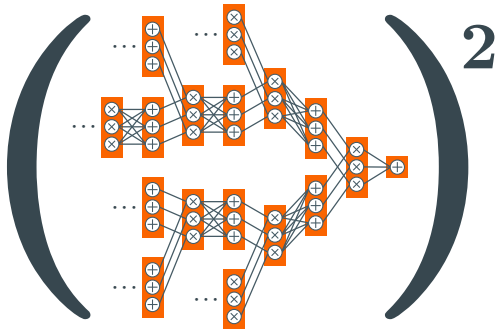
$\exists p$ requiring exponentially large
squared non-mono circuits...

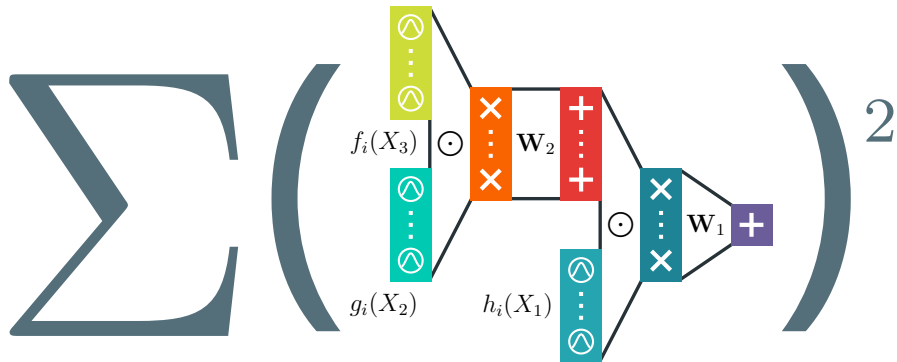


theorem



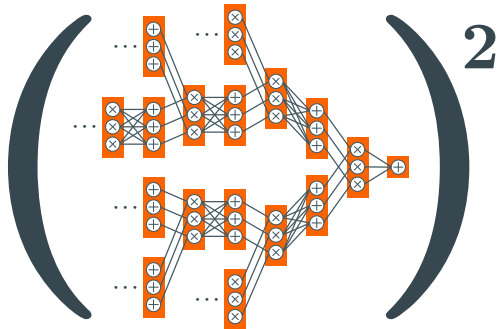
...but compact
monotonic circuits...!





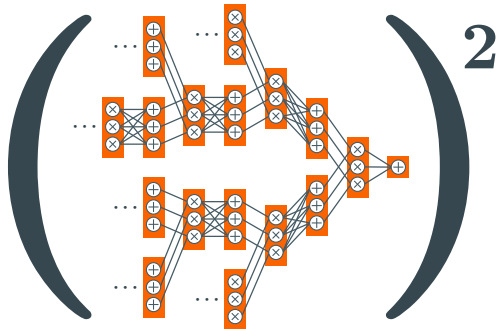
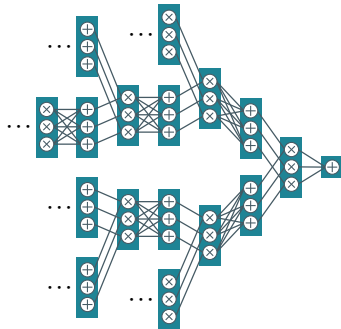
what if we use more than one square?

theorem



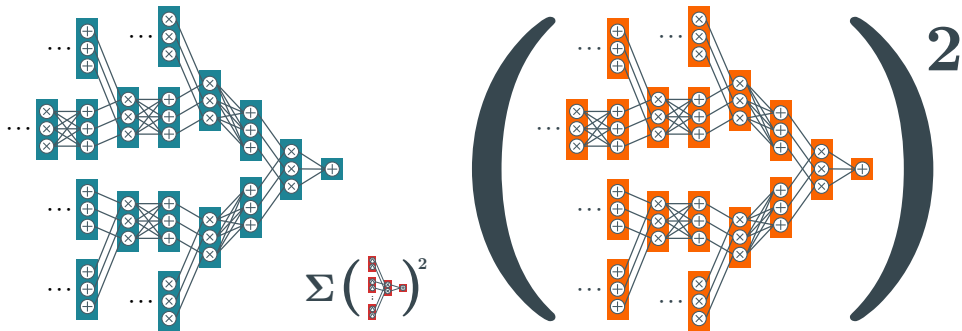
$\exists p$ requiring exponentially large **squared non-mono circuits**...

theorem

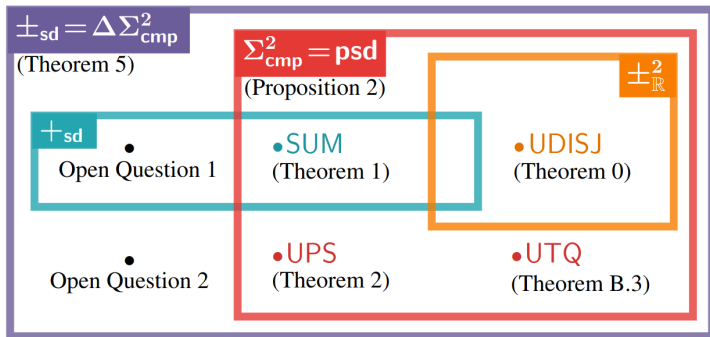


...exponentially large **monotonic circuits**...

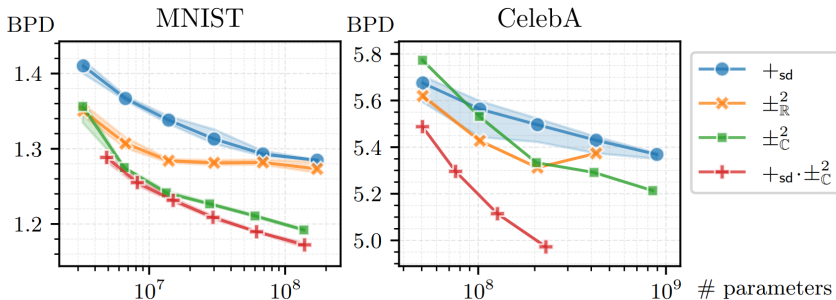
theorem



...but compact **SOS circuits...**



a hierarchy of subtractive mixtures



complex circuits are SOS (and scale better!)

compositional inference I



```
1 from cirq.symbolic.functional import integrate, multiply,  
   ↪ conjugate  
2  
3 # create a deep circuit with complex parameters  
4 c = build_symbolic_complex_circuit('quad-tree-4')  
5  
6 # compute the partition function of  $c^2$   
7 def renormalize(c):  
8     c1 = conjugate(c)  
9     c2 = multiply(c, c1)  
10    return integrate(c2)
```

EigenVI: score-based variational inference with orthogonal function expansions

Diana Cai

Flatiron Institute

dcai@flatironinstitute.org

Chirag Modi

Flatiron Institute

cmodi@flatironinstitute.org

Charles C. Margossian

Flatiron Institute

cmargossian@flatironinstitute.org

Robert M. Gower

Flatiron Institute

rgower@flatironinstitute.org

David M. Blei

Columbia University

david.blei@columbia.edu

Lawrence K. Saul

Flatiron Institute

lsaul@flatironinstitute.org

*what about **deep** orthonormal mixtures
and **arbitrary** marginals?*

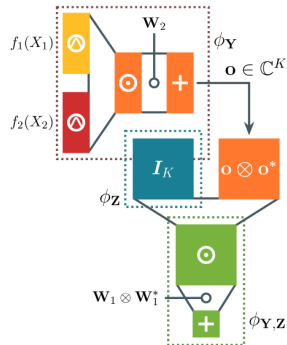
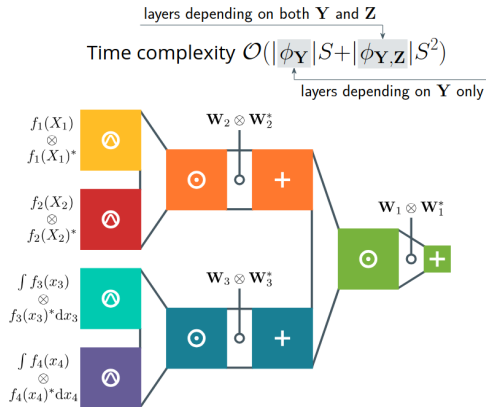
On Faster Marginalization with Squared Circuits via Orthonormalization

Lorenzo Loconte¹ **Antonio Vergari¹**

¹ School of Informatics, University of Edinburgh, UK

l.loconte@sms.ed.ac.uk, avergari@ed.ac.uk

*it suffices to **orthonormalize** each layer!*



faster marginalization of arbitrary subsets of features

approximate inference

e.g., via sampling

Can we use a subtractive mixture model to approximate expectations?

$$\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [f(\mathbf{x})] \approx \frac{1}{S} \sum_{i=1}^S f(\mathbf{x}^{(i)}) \quad \text{with} \quad \mathbf{x}^{(i)} \sim q(\mathbf{x})$$

\Rightarrow *but how to sample from q ?*

approximate inference

e.g., via sampling

Can we use a subtractive mixture model to approximate expectations?

$$\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} [f(\mathbf{x})] \approx \frac{1}{S} \sum_{i=1}^S f(\mathbf{x}^{(i)}) \quad \text{with} \quad \mathbf{x}^{(i)} \sim q(\mathbf{x})$$

\Rightarrow but how to sample from q ?

use **autoregressive inverse transform sampling**:

$$x_1 \sim q(x_1), \quad x_i \sim q(x_i | \mathbf{x}_{<i}) \quad \text{for } i \in \{2, \dots, d\}$$

\Rightarrow can be slow for large dimensions, requires **inverting the CDF**

Loconte et al., "What is the Relationship between Tensor Factorizations and Circuits (and How Can We Exploit it)?", TMLR, 2025

approximate inference

difference of expectation estimator

Idea: represent q as a difference of two additive mixtures

$$q(\mathbf{x}) = Z_+ \cdot q_+(\mathbf{x}) - Z_- \cdot q_-(\mathbf{x})$$

\Rightarrow *expectations will break down in two “parts”*

approximate inference

difference of expectation estimator

Idea: represent q as a difference of two additive mixtures

$$q(\mathbf{x}) = Z_+ \cdot q_+(\mathbf{x}) - Z_- \cdot q_-(\mathbf{x})$$

\Rightarrow expectations will break down in two “parts”

$$\frac{Z_+}{S_+} \sum_{s=1}^{S_+} f(\mathbf{x}_+^{(s)}) - \frac{Z_-}{S_-} \sum_{s=1}^{S_-} f(\mathbf{x}_-^{(s)}), \text{ where } \begin{matrix} \mathbf{x}_+^{(s)} \sim q_+(\mathbf{x}_+) \\ \mathbf{x}_-^{(s)} \sim q_-(\mathbf{x}_-) \end{matrix}, \quad (1)$$

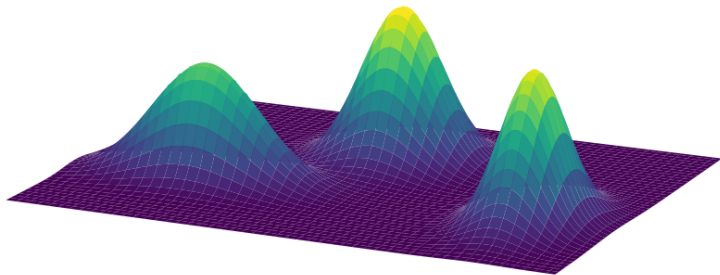
approximate inference

difference of expectation estimator

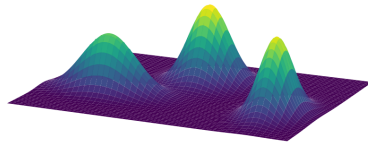
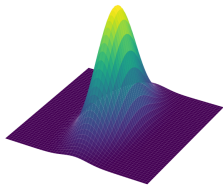
		Number of components (K)					
		2		4		6	
Method	d	$\log(\hat{I} - I)$	Time (s)	$\log(\hat{I} - I)$	Time (s)	$\log(\hat{I} - I)$	Time (s)
ΔExS	16	-19.507 ± 1.025	0.293 ± 0.004	-19.062 ± 0.823	1.049 ± 0.077	-19.497 ± 1.974	2.302 ± 0.159
ARITS	16	-19.111 ± 1.103	7.525 ± 0.038	-19.299 ± 1.611	7.52 ± 0.023	-18.739 ± 1.024	7.746 ± 0.032
ΔExS	32	-48.411 ± 1.265	0.325 ± 0.012	-48.046 ± 0.972	1.027 ± 0.107	-48.34 ± 0.814	2.213 ± 0.177
ARITS	32	-47.897 ± 1.165	15.196 ± 0.059	-47.349 ± 0.839	15.535 ± 0.059	-47.3 ± 0.978	17.371 ± 0.06
ΔExS	64	-108.095 ± 1.094	0.38 ± 0.034	-107.56 ± 0.616	0.9 ± 0.14	-107.653 ± 0.945	1.512 ± 0.383
ARITS	64	-107.898 ± 1.129	30.459 ± 0.098	-107.33 ± 0.929	33.892 ± 0.119	-107.374 ± 1.138	52.02 ± 0.127

faster than autoregressive sampling

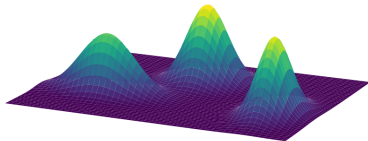
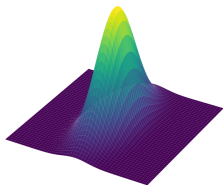
*Zellinger et al., "Scalable Expectation Estimation with Subtractive Mixture Models",
Under submission, 2025*



oh mixtures, you're so fine you blow my mind!



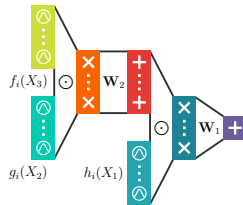
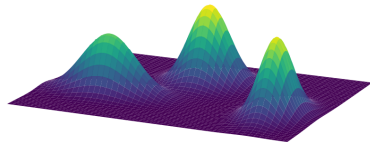
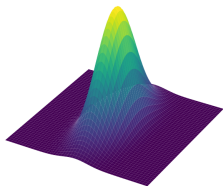
$$p(\mathbf{X}) \quad \longrightarrow \quad \sum_{i=1}^K w_i p_i(\mathbf{X}) \quad w_i > 0$$



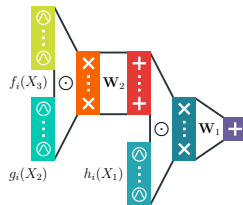
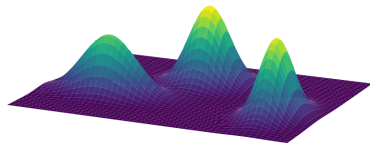
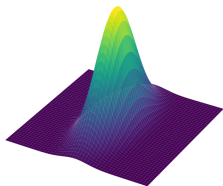
$$p(\mathbf{X}) \quad \longrightarrow \quad \sum_{i=1}^K w_i p_i(\mathbf{X}) \quad w_i > 0$$

*“if someone publishes a paper on **model A**, there will be a paper about **mixtures of A** soon, with high probability”*

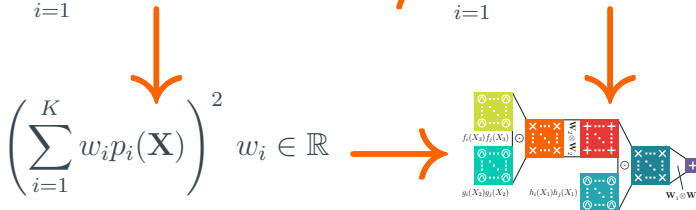
A. Vergari



$$p(\mathbf{X}) \longrightarrow \sum_{i=1}^K w_i p_i(\mathbf{X}) \quad w_i > 0 \longrightarrow \sum_{i=1}^{2^D} w_i p_i(\mathbf{X}) = \text{PC}(\mathbf{X})$$



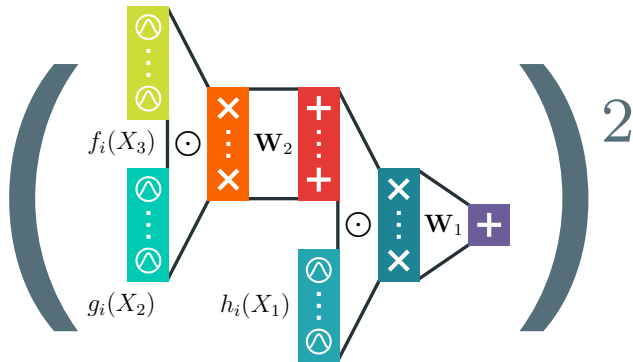
$$p(\mathbf{X}) \longrightarrow \sum_{i=1}^K w_i p_i(\mathbf{X}) \quad w_i > 0 \longrightarrow \sum_{i=1}^{2^D} w_i p_i(\mathbf{X}) = \text{PC}(\mathbf{X})$$





learning & reasoning with circuits in pytorch

`github.com/april-tools/circuit`



questions?