

# Rigorous bounds for posterior inference in universal probabilistic programming

Raven Beutner<sup>1</sup>   Luke Ong<sup>2</sup>   **Fabian Zaiser**<sup>2</sup>

<sup>1</sup>CISPA Helmholtz Center for Information Security

<sup>2</sup>University of Oxford

Languages for Inference @ POPL 2022

# A random walk as a probabilistic program

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

# A random walk as a probabilistic program

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

- ▶ continuous distributions
- ▶ unbounded loops
- ▶ unbounded number of **samples**

# Existing inference methods

**1. Approximate:** posterior  $\approx X$

- ▶ Monte Carlo (particle filter, MCMC)
- ▶ or optimization-based (variational inference)



# Existing inference methods

## 1. **Approximate:** posterior $\approx X$

- ▶ Monte Carlo (particle filter, MCMC)
- ▶ or optimization-based (variational inference)

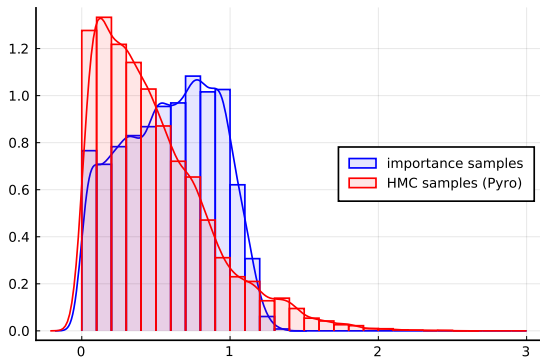


## 2. **Exact:** posterior = $X$

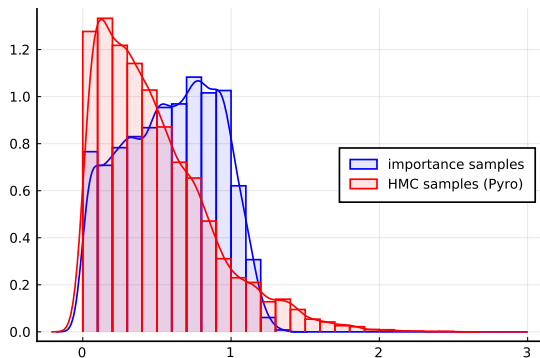
- ▶ symbolic expression



# Issues with existing methods



# Issues with existing methods



- ▶ **approximate methods:** convergence (e.g. for multimodal models)
- ▶ **exact methods:** restricted models (e.g. no recursion)

# Rigorous Bounds on the Posterior:

$$\text{posterior}(E) \in [a, b]$$



# Rigorous Bounds on the Posterior:

$$\text{posterior}(E) \in [a, b]$$

... for

- ▶ a universal PPL (including branching & recursion)
- ▶ with continuous distributions
- ▶ and conditioning (**observe**)

# Rigorous Bounds on the Posterior:

$$\text{posterior}(E) \in [a, b]$$

... for

- ▶ a universal PPL (including branching & recursion)
- ▶ with continuous distributions
- ▶ and conditioning (**observe**)

## Why?

- ▶ construct ground truth for inference problems
- ▶ to debug approximate inference

# Rigorous Bounds on the Posterior:

$$\text{posterior}(E) \in [a, b]$$

... for

- ▶ a universal PPL (including branching & recursion)
- ▶ with continuous distributions
- ▶ and conditioning (**observe**)

## Why?

- ▶ construct ground truth for inference problems
- ▶ to debug approximate inference

## How?

1. interval traces & interval arithmetic (basic idea)
2. interval type system (overapproximation)
3. symbolic execution (optimization of special case)

# Method 1: Interval traces

standard semantics

---

**traces**

$\langle 0.2, 0.8 \rangle$

$$\mathbb{T} := \bigcup_{n \in \mathbb{N}} \mathbb{R}^n$$

---

**value**

$$\text{val}_P : \mathbb{T} \rightarrow \mathbb{R}$$

---

**weight**

$$\text{wt}_P : \mathbb{T} \rightarrow [0, \infty)$$

---

**posterior**

$$\llbracket P \rrbracket(E)$$

integral over  $\mathbb{T}$

# Method 1: Interval traces

**Idea:** summarize traces using intervals

	standard semantics	interval semantics
<b>traces</b>	$\langle 0.2, 0.8 \rangle$ $\mathbb{T} := \bigcup_{n \in \mathbb{N}} \mathbb{R}^n$	$\langle [0.2, 0.3], [0.7, 0.8] \rangle$ $\mathbb{T}_{\mathbb{I}} := \bigcup_{n \in \mathbb{N}} \mathbb{I}^n$
<b>value</b>	$\text{val}_P : \mathbb{T} \rightarrow \mathbb{R}$	$\text{val}_P^{\mathbb{I}} : \mathbb{T}_{\mathbb{I}} \rightarrow \mathbb{I}$
<b>weight</b>	$\text{wt}_P : \mathbb{T} \rightarrow [0, \infty)$	$\text{wt}_P^{\mathbb{I}} : \mathbb{T}_{\mathbb{I}} \rightarrow \mathbb{I}_{[0, \infty)}$
<b>posterior</b>	$\llbracket P \rrbracket(E)$ integral over $\mathbb{T}$	$[\text{lowerBd}_P^{\mathcal{T}}(E), \text{upperBd}_P^{\mathcal{T}}(E)]$ sum over partition $\mathcal{T} \subset \mathbb{T}_{\mathbb{I}}$

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position		
distance		
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.6	[0.5, 0.6]
distance	0.0	[0.0, 0.0]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		



# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.6	[0.5, 0.6]
distance	0.0	[0.0, 0.0]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.8	[0.6, 0.8]
distance	0.2	[0.1, 0.2]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.8	[0.6, 0.8]
distance	0.2	[0.1, 0.2]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.0	[-0.3, 0.0]
distance	1.0	[0.9, 1.1]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	1	[1, 1]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.0	[-0.3, 0.0]
distance	1.0	[0.9, 1.1]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	$\approx 2.4$	[0.53, 3.99]
<b>return value</b> $\text{val}(t)$		

# Interval trace semantics

```
start = sample uniform(0,3)
position = start; distance = 0
while position > 0:
    step = sample uniform(-1, 1)
    position += step
    distance += abs(step)
observe 1.1 from normal(distance, 0.12)
return start
```

	standard	interval semantics
start	0.6	[0.5, 0.6]
position	0.0	[-0.3, 0.0]
distance	1.0	[0.9, 1.1]
<b>trace</b> $t$	$\langle 0.6, 0.2, -0.8 \rangle$	$\langle [0.5, 0.6], [0.1, 0.2], [-0.9, -0.8] \rangle$
<b>weight</b> $\text{wt}(t)$	$\approx 2.4$	[0.53, 3.99]
<b>return value</b> $\text{val}(t)$	<b>0.6</b>	<b>[0.5, 0.6]</b>

## Lower bounds on the posterior

If  $\mathcal{T} \subset \mathbb{T}_{\mathbb{I}}$  is “non-overlapping”, then define for all intervals  $I$ :

$$\text{lowerBd}_P^{\mathcal{T}}(I) := \sum_{\mathbf{s}_{\mathbb{I}} \in \mathcal{T}} \text{vol}(\mathbf{s}_{\mathbb{I}}) \cdot (\min \text{wt}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}})) \cdot [\text{val}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}}) \subseteq I]$$

$$\text{vol}(\langle [a_1, b_1], \dots, [a_n, b_n] \rangle) := (b_1 - a_1) \times \dots \times (b_n - a_n)$$

## Lower bounds on the posterior

If  $\mathcal{T} \subset \mathbb{T}_{\mathbb{I}}$  is “non-overlapping”, then define for all intervals  $I$ :

$$\text{lowerBd}_P^{\mathcal{T}}(I) := \sum_{\mathbf{s}_{\mathbb{I}} \in \mathcal{T}} \text{vol}(\mathbf{s}_{\mathbb{I}}) \cdot (\min \text{wt}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}})) \cdot [\text{val}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}}) \subseteq I]$$

$$\text{vol}(\langle [a_1, b_1], \dots, [a_n, b_n] \rangle) := (b_1 - a_1) \times \dots \times (b_n - a_n)$$

## Upper bounds on the posterior

If  $\mathcal{T} \subset \mathbb{T}_{\mathbb{I}}$  is “exhaustive” (covers every trace), then define for all intervals  $I$ :

$$\text{upperBd}_P^{\mathcal{T}}(I) := \sum_{\mathbf{s}_{\mathbb{I}} \in \mathcal{T}} \text{vol}(\mathbf{s}_{\mathbb{I}}) \cdot (\max \text{wt}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}})) \cdot [\text{val}_P^{\mathbb{I}}(\mathbf{s}_{\mathbb{I}}) \cap I \neq \emptyset]$$



# Soundness

$$\text{lowerBd}_P^{\mathcal{T}} \leq \llbracket P \rrbracket \leq \text{upperBd}_P^{\mathcal{T}}.$$

## Soundness

$$\text{lowerBd}_P^{\mathcal{T}} \leq \llbracket P \rrbracket \leq \text{upperBd}_P^{\mathcal{T}}.$$

## Completeness

For all intervals  $I$  and  $\epsilon > 0$ , there is a countable set  $\mathcal{T} \subseteq \mathbb{T}_{\mathbb{I}}$  s.t.

$$\text{upperBd}_P^{\mathcal{T}}(I) - \epsilon \leq \llbracket P \rrbracket(I) \leq \text{lowerBd}_P^{\mathcal{T}}(I) + \epsilon$$

under the assumptions:

- ▶ the primitive functions are continuous\*
- ▶ each **sampled** value is used at most once in each condition, **observe** statement, and in the return value.

## Soundness

$$\text{lowerBd}_P^{\mathcal{T}} \leq \llbracket P \rrbracket \leq \text{upperBd}_P^{\mathcal{T}}.$$

## Completeness

For all intervals  $I$  and  $\epsilon > 0$ , there is a **countable** set  $\mathcal{T} \subseteq \mathbb{T}_{\mathbb{I}}$  s.t.

$$\text{upperBd}_P^{\mathcal{T}}(I) - \epsilon \leq \llbracket P \rrbracket(I) \leq \text{lowerBd}_P^{\mathcal{T}}(I) + \epsilon$$

under the assumptions:

- ▶ the primitive functions are continuous\*
- ▶ each **sampled** value is used at most once in each condition, **observe** statement, and in the return value.

## Soundness

$$\text{lowerBd}_P^{\mathcal{T}} \leq \llbracket P \rrbracket \leq \text{upperBd}_P^{\mathcal{T}}.$$

## Completeness

For all intervals  $I$  and  $\epsilon > 0$ , there is a **finite** set  $\mathcal{T} \subseteq \mathbb{T}_I$  s.t.

$$\llbracket P \rrbracket(I) \leq \text{lowerBd}_P^{\mathcal{T}}(I) + \epsilon$$

under the assumptions:

- ▶ the primitive functions are continuous\*
- ▶ each **sampled** value is used at most once in each condition, **observe** statement, and in the return value.

## Method 2: Interval type system

→ to **overapproximate** recursion and conditionals (not resolvable by intervals)

## Method 2: Interval type system

→ to **overapproximate** recursion and conditionals (not resolvable by intervals)

- ▶ types keep track of the value and weight interval
- ▶  $\vdash P : \left\{ \begin{array}{l} [v, v'] \\ [w, w'] \end{array} \right\}$  means  $\text{val}_P(\mathbf{s}) \in [v, v']$  and  $\text{wt}_P(\mathbf{s}) \in [w, w']$ .
- ▶ efficient type inference
- ▶ uses interval arithmetic & **widening** to approximate fixpoints

## Method 3: Symbolic execution

→ optimization for a common special case

## Method 3: Symbolic execution

→ optimization for a common special case

For each program path,

- ▶  $\alpha_k$ : the  $k$ -th **sample**
- ▶  $\mathcal{V}$ : result value, e.g.  $\alpha_1 + 2\alpha_2$
- ▶  $\Delta$ : guards, e.g.  $\{\alpha_1 \leq 0, \alpha_1 + \alpha_2 > 1\}$
- ▶  $\Xi$ : weights, e.g.  $\{\text{pdf}_{\text{Normal}(0,1)}(\alpha_1 - \alpha_2), \text{pdf}_{\text{Normal}(1,2)}(\alpha_3)\}$



## Method 3: Symbolic execution

→ optimization for a common special case

For each program path,

- ▶  $\alpha_k$ : the  $k$ -th **sample**
- ▶  $\mathcal{V}$ : result value, e.g.  $\alpha_1 + 2\alpha_2$
- ▶  $\Delta$ : guards, e.g.  $\{\alpha_1 \leq 0, \alpha_1 + \alpha_2 > 1\}$
- ▶  $\Xi$ : weights, e.g.  $\{\text{pdf}_{\text{Normal}(0,1)}(\alpha_1 - \alpha_2), \text{pdf}_{\text{Normal}(1,2)}(\alpha_3)\}$

$$\llbracket P \rrbracket(I) = \sum_{\text{paths}} \int_{\Delta \cup \{\mathcal{V} \in I\}} \left( \prod \Xi \right) d\alpha$$

## Method 3: Symbolic execution

→ optimization for a common special case

For each program path,

- ▶  $\alpha_k$ : the  $k$ -th **sample**
- ▶  $\mathcal{V}$ : result value, e.g.  $\alpha_1 + 2\alpha_2$
- ▶  $\Delta$ : guards, e.g.  $\{\alpha_1 \leq 0, \alpha_1 + \alpha_2 > 1\}$
- ▶  $\Xi$ : weights, e.g.  $\{\text{pdf}_{\text{Normal}(0,1)}(\alpha_1 - \alpha_2), \text{pdf}_{\text{Normal}(1,2)}(\alpha_3)\}$

$$\llbracket P \rrbracket(I) = \sum_{\text{paths}} \int_{\Delta \cup \{\mathcal{V} \in I\}} \left( \prod \Xi \right) d\alpha \leq \sum_{\text{paths}} \text{vol}(\Delta \cup \{\mathcal{V} \in I\}) \prod_{\mathcal{W} \in \Xi} \max_{\alpha} \mathcal{W}$$

## Method 3: Symbolic execution

→ optimization for a common special case

For each program path,

- ▶  $\alpha_k$ : the  $k$ -th **sample**
- ▶  $\mathcal{V}$ : result value, e.g.  $\alpha_1 + 2\alpha_2$
- ▶  $\Delta$ : guards, e.g.  $\{\alpha_1 \leq 0, \alpha_1 + \alpha_2 > 1\}$
- ▶  $\Xi$ : weights, e.g.  $\{\text{pdf}_{\text{Normal}(0,1)}(\alpha_1 - \alpha_2), \text{pdf}_{\text{Normal}(1,2)}(\alpha_3)\}$

$$\llbracket P \rrbracket(I) = \sum_{\text{paths}} \int_{\Delta \cup \{\mathcal{V} \in I\}} \left( \prod \Xi \right) d\alpha \leq \sum_{\text{paths}} \text{vol}(\Delta \cup \{\mathcal{V} \in I\}) \prod_{\mathcal{W} \in \Xi} \max_{\alpha} \mathcal{W}$$

If  $\Delta$  and  $\mathcal{V}$  are affine then use

- ▶ polytope volume computation (→ Vinci tool)

## Method 3: Symbolic execution

→ optimization for a common special case

For each program path,

- ▶  $\alpha_k$ : the  $k$ -th **sample**
- ▶  $\mathcal{V}$ : result value, e.g.  $\alpha_1 + 2\alpha_2$
- ▶  $\Delta$ : guards, e.g.  $\{\alpha_1 \leq 0, \alpha_1 + \alpha_2 > 1\}$
- ▶  $\Xi$ : weights, e.g.  $\{\text{pdf}_{\text{Normal}(0,1)}(\alpha_1 - \alpha_2), \text{pdf}_{\text{Normal}(1,2)}(\alpha_3)\}$

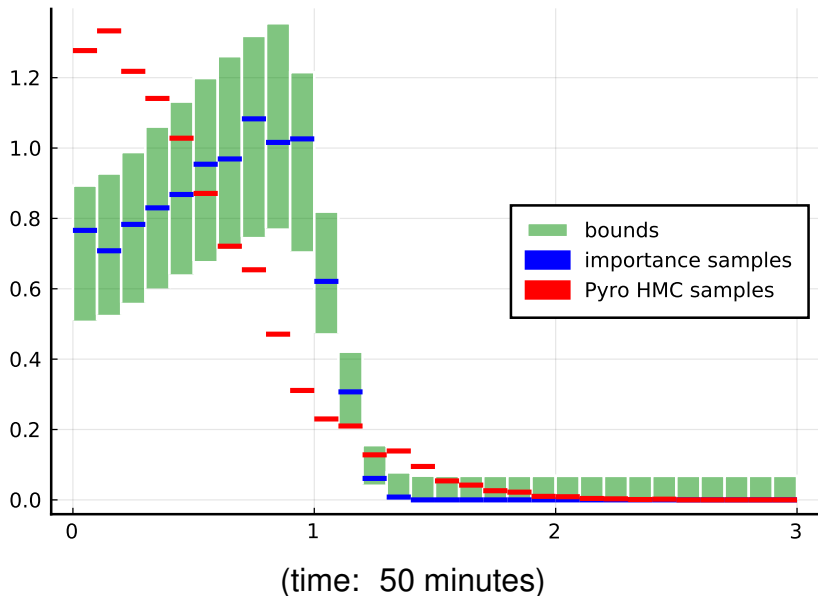
$$\llbracket P \rrbracket(I) = \sum_{\text{paths}} \int_{\Delta \cup \{\mathcal{V} \in I\}} \left( \prod \Xi \right) d\alpha \leq \sum_{\text{paths}} \text{vol}(\Delta \cup \{\mathcal{V} \in I\}) \prod_{\mathcal{W} \in \Xi} \max_{\alpha} \mathcal{W}$$

If  $\Delta$  and  $\mathcal{V}$  are affine then use

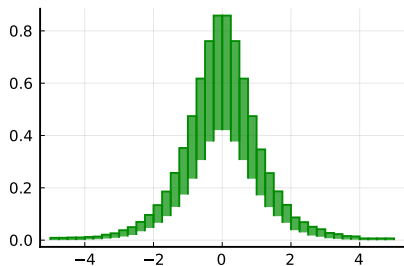
- ▶ polytope volume computation (→ Vinci tool)
- ▶ linear optimization & interval arithmetic

# Empirical evaluation

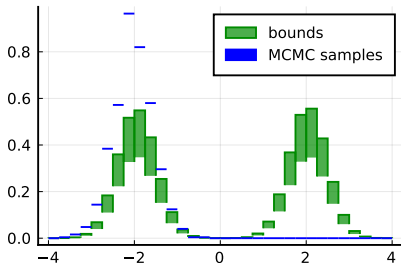
# Empirical evaluation



# Examples that are hard for MCMC



(a) Neal's funnel (5 seconds)



(b) Binary Gaussian mixture model (90 seconds)

# Comparison with previous work

## Sankaranarayanan et al. (PLDI13)

- ▶ bounding probabilities (but no **observe**)
- ▶ ours is usually slower, but often better bounds



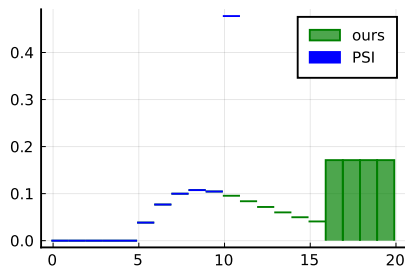
# Comparison with previous work

## Sankaranarayanan et al. (PLDI13)

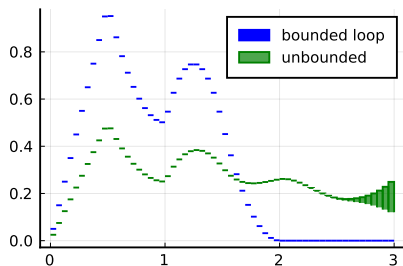
- ▶ bounding probabilities (but no **observe**)
- ▶ ours is usually slower, but often better bounds

## PSI solver

- ▶ consistency check: benchmarks from the PSI repository
- ▶ we can handle unbounded loops, contrary to PSI



(a)  $\approx$  2 minutes



(b)  $\approx$  20 seconds

## Limitations

- ▶ lots of branching
- ▶ high-dimensional models (many **samples**)

## Limitations

- ▶ lots of branching
- ▶ high-dimensional models (many **samples**)

## Future work

- ▶ better heuristics for finding a “good”  $\mathcal{T}$
- ▶ can this be refined into an approximate inference algorithm?

## Rigorous Bounds on the Posterior:

$$\text{posterior}(E) \in [a, b]$$

... for a **universal** PPL with **continuous** distributions and **conditioning** (**observe**)

### Why?

- ▶ construct ground truth
- ▶ debug approximate inference

### How?

1. interval trace semantics
2. interval type system
3. symbolic execution

# Backup slides

# Semantics of a probabilistic program

Let  $P : \mathbb{R}$  a probabilistic program.

- ▶ trace space:  $\mathbb{T} := \bigcup_{n \in \mathbb{N}} \mathbb{R}^n$
- ▶ value function:  $\text{val}_P : \mathbb{T} \rightarrow \mathbb{R}$
- ▶ weight function:  $\text{wt}_P : \mathbb{T} \rightarrow [0, \infty)$

Unnormalized posterior:

$$\llbracket P \rrbracket(E) := \int_{\text{val}_P^{-1}(E)} \text{wt}_P(t) \mu_{\mathbb{T}}(dt) \quad \text{for event } E \subseteq \mathbb{R}.$$

Normalized posterior (probability distribution):

$$\text{posterior}_P(E) := \frac{1}{Z} \llbracket P \rrbracket(E) \quad \text{where } Z := \llbracket P \rrbracket(\mathbb{R}).$$

# Trace partitioning heuristics

**Option 1:** split equidistantly in each dimension

**Option 2:**

- ▶ start with the full interval trace  $\langle [-\infty, \infty], \dots \rangle$
- ▶ pick the next interval  $s_{\mathbb{I}}$  trace or, depending on the input program, select it with a mix of the following criteria
  - ▶ high weight  $\text{wt}_P^{\mathbb{I}}(s_{\mathbb{I}})$
  - ▶ wide value interval  $\text{val}_P^{\mathbb{I}}(s_{\mathbb{I}})$
  - ▶ large volume  $\text{vol}(s_{\mathbb{I}})$
- ▶ split that box in half along the dimension that reduces the width of the interval of the posterior expected value the most
- ▶ repeat.

# Interval type system

Types:

► unweighted:  $\sigma ::= [v, v'] \mid \sigma \rightarrow \mathcal{A}$

► weighted:  $\mathcal{A} ::= \left\{ \begin{array}{c} \sigma \\ [w, w'] \end{array} \right\}$

Selected typing rules:

$$\frac{\Gamma; \varphi : \sigma \rightarrow \mathcal{A}; x : \sigma \vdash M : \mathcal{A}}{\Gamma \vdash \mu_x^\varphi.M : \left\{ \begin{array}{c} \sigma \rightarrow \mathcal{A} \\ [1, 1] \end{array} \right\}}$$
$$\frac{\Gamma \vdash M : \left\{ \begin{array}{c} \sigma_1 \rightarrow \left\{ \begin{array}{c} \sigma_2 \\ [e, f] \end{array} \right\} \\ [a, b] \end{array} \right\} \quad \Gamma \vdash N : \left\{ \begin{array}{c} \sigma_1 \\ [c, d] \end{array} \right\}}{\Gamma \vdash MN : \left\{ \begin{array}{c} \sigma_2 \\ [a, b] \times^{\mathbb{I}} [c, d] \times^{\mathbb{I}} [e, f] \end{array} \right\}}$$