

Divide, Conquer, and Combine: a New Inference Strategy for Probabilistic Programs with Stochastic Support

Yuan Zhou, Hongseok Yang, Yee Whye Teh, Tom Rainforth Proceedings of the 37th International Conference on Machine Learning, PMLR 119:11534-11545, 2020.

Presented by: Bianca Cosma and Margot Pauëlsen | 25.09.2023



Introduction

- Stochastic variables lead to variable dimensionality
- Model declaration is easy, (automated) inference is hard
- **Goal of the paper:** provide stochastic support to wide range of models

```
1  z0 = sample(normal(0, 2))
2  y1 = 9
3  if z0 < 0:
4      z1 = sample(normal(-5, 2))
5      observe(normal(z1, 2), y1)
6      return (z0, z1)
7  else:
8      z2 = sample(normal(5, 2))
9      z3 = sample(normal(z2, 2))
10     observe(normal(z3, 2), y1)
11     return (z0, z2, z3)
```

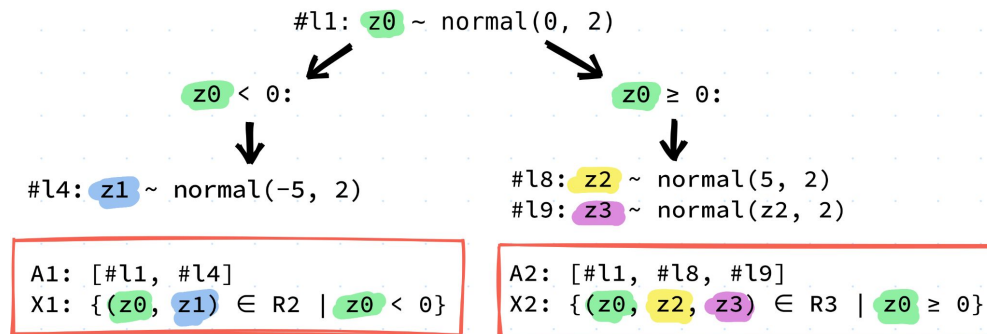


Figure adapted from Zhou et al.

Notation

Sample

The set of **random variables** $\{x_i\}_{i=1}^{n_x} = (x_1, \dots, x_{n_x})$,

where x_i has **address** a_i , **input (history)** η_i and **density (prior)** $f_{a_i}(x_i|\eta_i)$

Observe

The set of **observed values** $\{y_j\}_{j=1}^{n_y} = (y_1, \dots, y_{n_y})$,

where y_j has **address** b_j , **input** ϕ_j and **density (likelihood)** $g_{b_j}(y_j|\phi_j)$

Notation (continued)

The **posterior** $p(\mathbf{x} | \mathbf{y})$ can be expressed as $\pi(\mathbf{x}) = \gamma(\mathbf{x})/Z$,

where

$$\gamma(\mathbf{x}) := \prod_{i=1}^{n_x} \overset{\text{(prior)}}{f_{a_i}(x_i | \eta_i)} \prod_{j=1}^{n_y} \overset{\text{(likelihood)}}{g_{b_j}(y_j | \phi_j)},$$

 **joint density $p(\mathbf{x}, \mathbf{y})$**

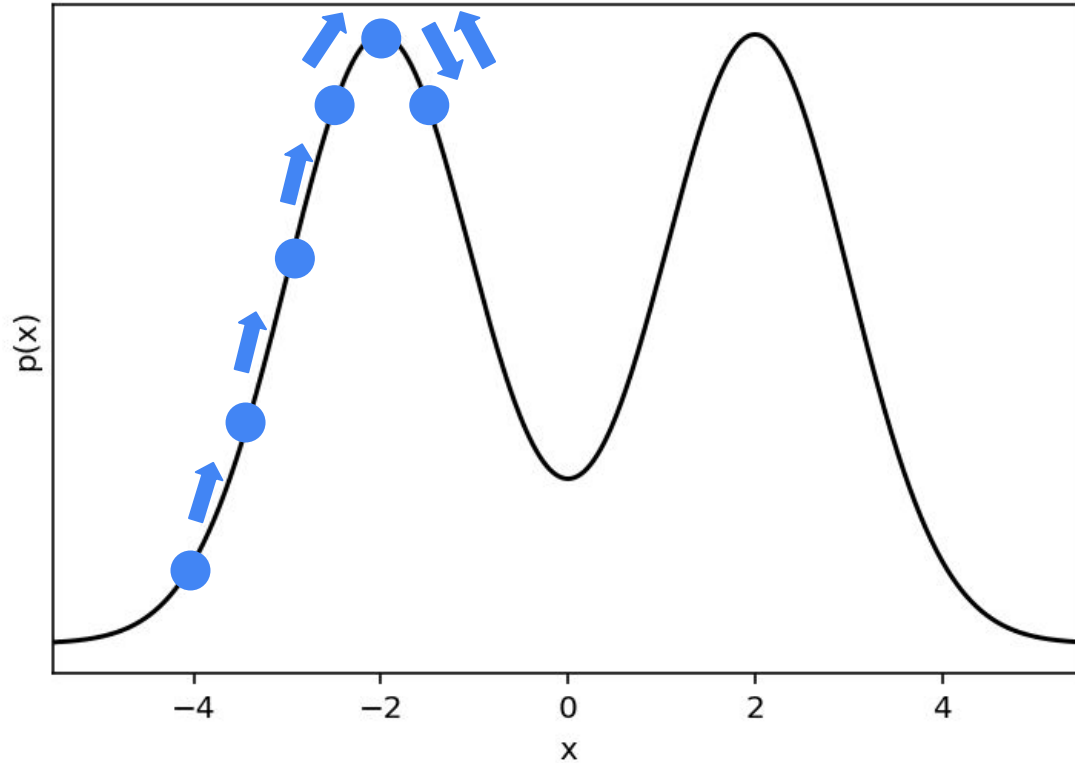
$$Z := \int \prod_{i=1}^{n_x} f_{a_i}(x_i | \eta_i) \prod_{j=1}^{n_y} g_{b_j}(y_j | \phi_j) dx_{1:n_x},$$

 **marginal likelihood $p(\mathbf{y})$**
(normalizing constant)

Problems with other inference methods

- Rejection and importance sampling suffer ‘curse of dimensionality’
- Variational inference
 - needs specific knowledge, so at odds with automation
 - model often approximated by model with fixed support
- MCMC needs kernel to switch between configurations
 - Posterior may vary across configurations
 - It is difficult to switch from high density region to new configuration

Problems with other inference methods

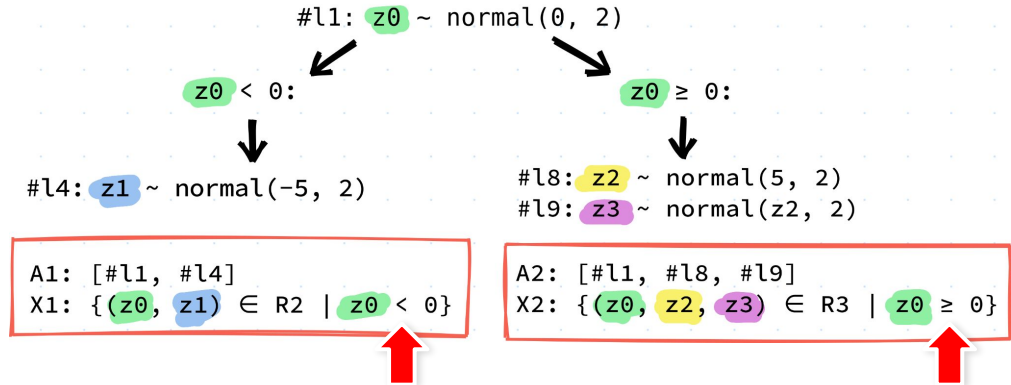


Multimodal distribution image:
<https://emcee.readthedocs.io/en/stable/tutorials/moves/>

Step 1: Divide

- Divide program in straight-line sub-programs (SLPs)
- **The individual programs are disjoint (independent)**

```
1  z0 = sample(normal(0, 2))
2  y1 = 9
3  if z0 < 0:
4      z1 = sample(normal(-5, 2))
5      observe(normal(z1, 2), y1)
6      return (z0, z1)
7  else:
8      z2 = sample(normal(5, 2))
9      z3 = sample(normal(z2, 2))
10     observe(normal(z3, 2), y1)
11     return (z0, z2, z3)
```



Step 2: Conquer

For each SLP k , we estimate, using any conventional inference approach:

The posterior:

$$\pi_k(x)$$

$p(x|y)$

The marginal likelihood:
(evidence)

$$Z_k$$

$p(y)$

Step 3: Combine

joint distributions $p(x, y)$ for all SLPs 1:K

$$\pi(x) = \frac{\sum_{k=1}^K \gamma_k(x)}{\sum_{k=1}^K Z_k}$$

we want the posterior $p(x|y)$ for the entire program

marginal likelihoods $p(y)$ for all SLPs 1:K

We can sum up disjoint events to get the joint probability

Step 3: Combine (continued)

Why can we just sum the independent SLPs?

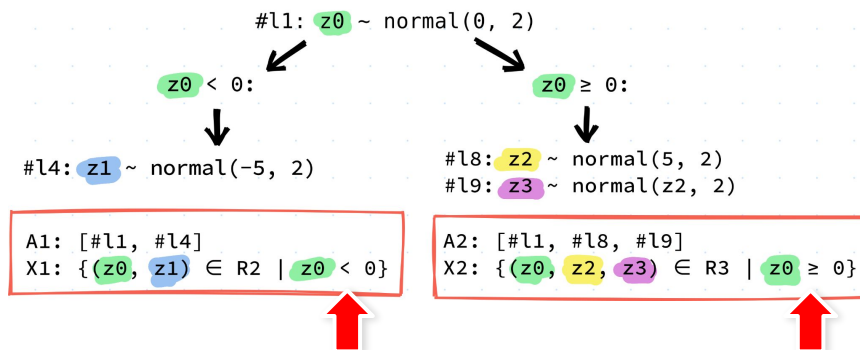
for the joint distribution:

$$\gamma(x) = \sum_{k=1}^K \gamma_k(x)$$

for the marginal likelihood:

$$Z = \sum_{k=1}^K Z_k$$

```
1 z0 = sample(normal(0, 2))
2 y1 = 9
3 if z0 < 0:
4     z1 = sample(normal(-5, 2))
5     observe(normal(z1, 2), y1)
6     return (z0, z1)
7 else:
8     z2 = sample(normal(5, 2))
9     z3 = sample(normal(z2, 2))
10    observe(normal(z3, 2), y1)
11    return (z0, z2, z3)
```



Because the program supports are disjoint (and we can sum up disjoint events to get the joint probability) 10

How do we find SLPs?

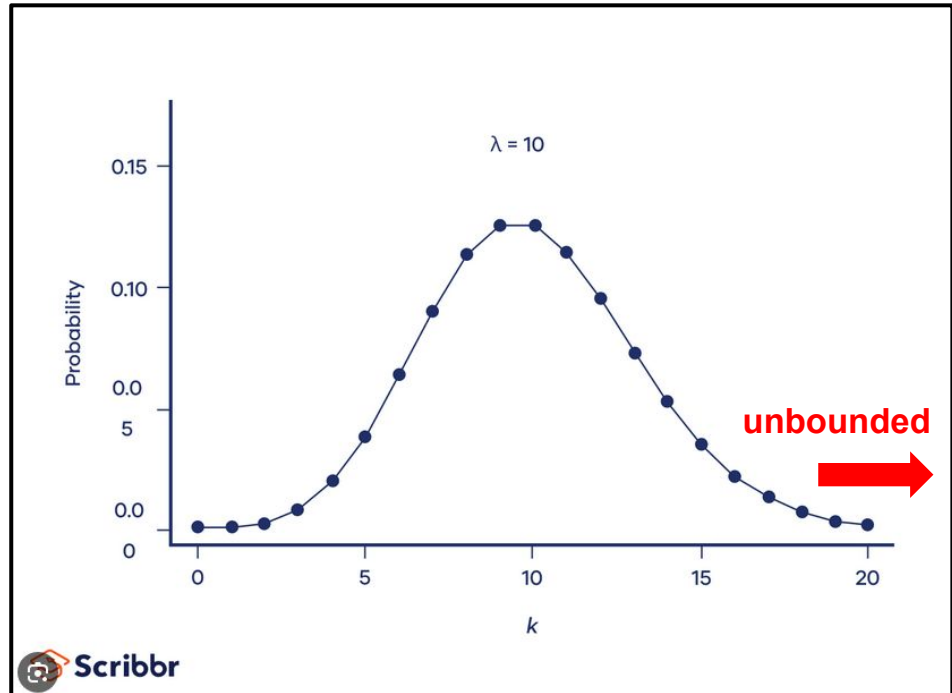
Static analysis doesn't always work: sometimes there's no upper bound on the number of code paths

$K \sim \text{Poisson}(10)$ \longrightarrow

for $k \in K$:

$$w \sim N(\mu_k, \Sigma_k)$$

...



Implementation details: finding SLPs at inference time

1. Run the program a number of times to generate an initial set of SLPs

Initial set: $S = \{[\#11, \#14]\}$

2. At each iteration:

- a. Perform local inference on an SLP (for example, MCMC)
- b. If we find another path, we update our set

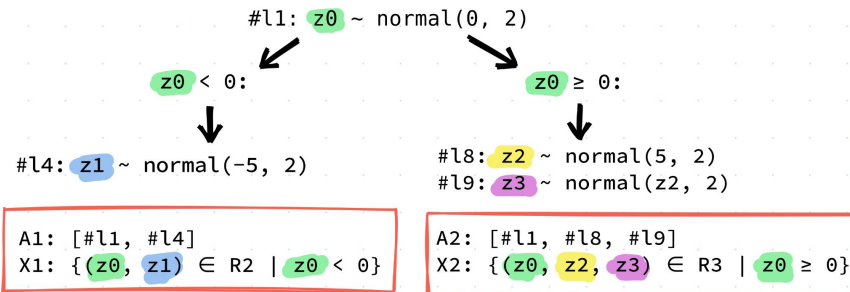
Assume on #11 we sample $z_0 = 0.4 > 0$:

- so we jump to #18
- we found another path
- update:
 $S = \{[\#11, \#14], [\#11, \#18, \#19]\}$

which one?



```
1  z0 = sample(normal(0, 2))
2  y1 = 9
3  if z0 < 0:
4      z1 = sample(normal(-5, 2))
5      observe(normal(z1, 2), y1)
6      return (z0, z1)
7  else:
8      z2 = sample(normal(5, 2))
9      z3 = sample(normal(z2, 2))
10     observe(normal(z3, 2), y1)
11     return (z0, z2, z3)
```



Implementation details: resource allocation

At each iteration, how do we choose an SLP to do local inference on?

Pick the one with the highest utility:

inversely proportional to the number of times we have already performed inference

$$U_k := \frac{1}{S_k} \left(\frac{(1 - \delta) \hat{\tau}_k}{\max_k \{ \hat{\tau}_k \}} + \frac{\delta \hat{p}_k}{\max_k \{ \hat{p}_k \}} + \frac{\beta \log \sum_k S_k}{\sqrt{S_k}} \right)$$

utility function for an SLP k

exploitation

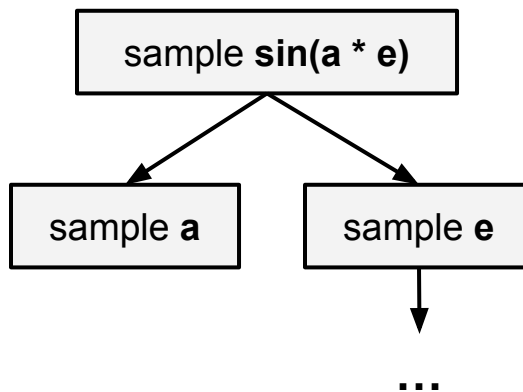
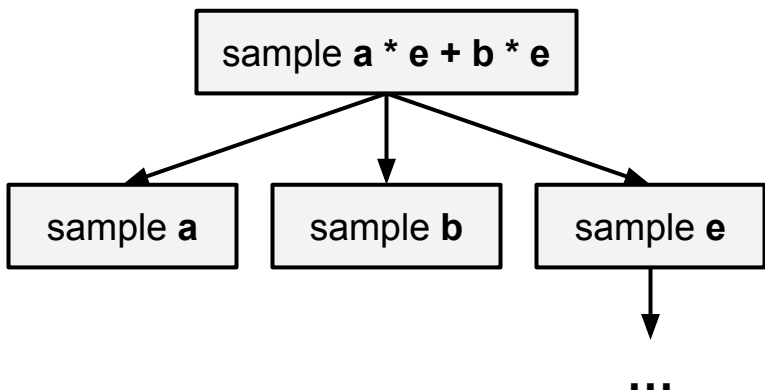
exploration

proportional to Z_k
(likelihood of the SLP)

Experiments

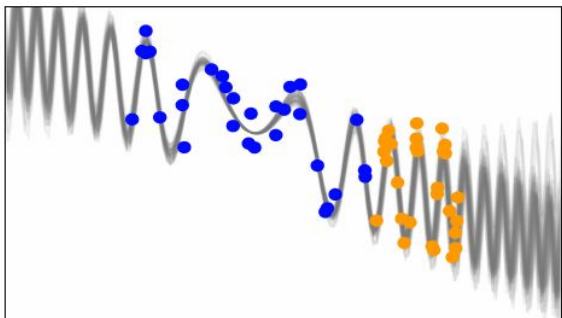
Function induction: we want to estimate the function that generated our data

- the set of rules: $e \rightarrow \{x \mid x^2 \mid \sin(a * e) \mid a * e + b * e\}$
- why does it make sense to apply “Divide, Conquer and Combine”?



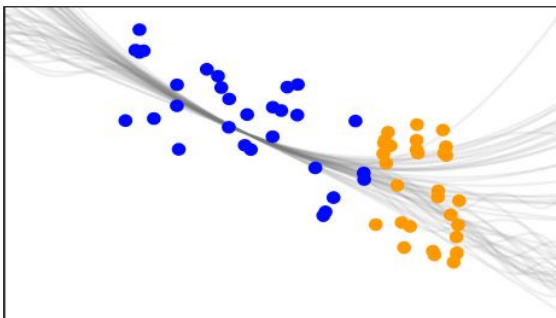
Experiments

Function induction: we want to estimate the function that generated our data



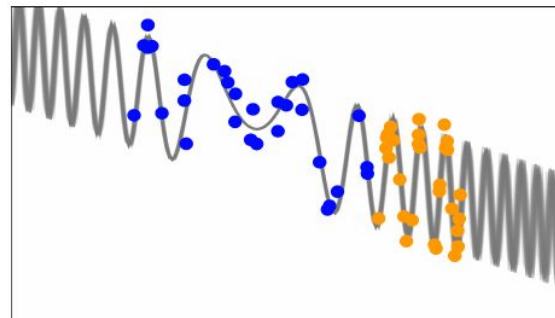
(a) DCC (ours)

(Divide, Conquer and Combine)



(b) IS

(importance sampling)



(c) RMH

(MCMC: random-walk Metropolis Hastings)