



# Assessing Inference Quality for Probabilistic Programs using Multivariate Simulation Based Calibration

Sharan Yalburgi\*, Jameson Quinn†, Veronica Weiner†§, Sam Witty‡, Vikash Mansinghka† and Cameron Freer†  
 \*BITS Pilani, India; †MIT, USA; ‡UMass Amherst, USA; §Probabilistic Computing Associates, USA



## Introduction

Problem:

- To have a principled approach that can test **arbitrary inference procedures** that are specified as code.

Solution:

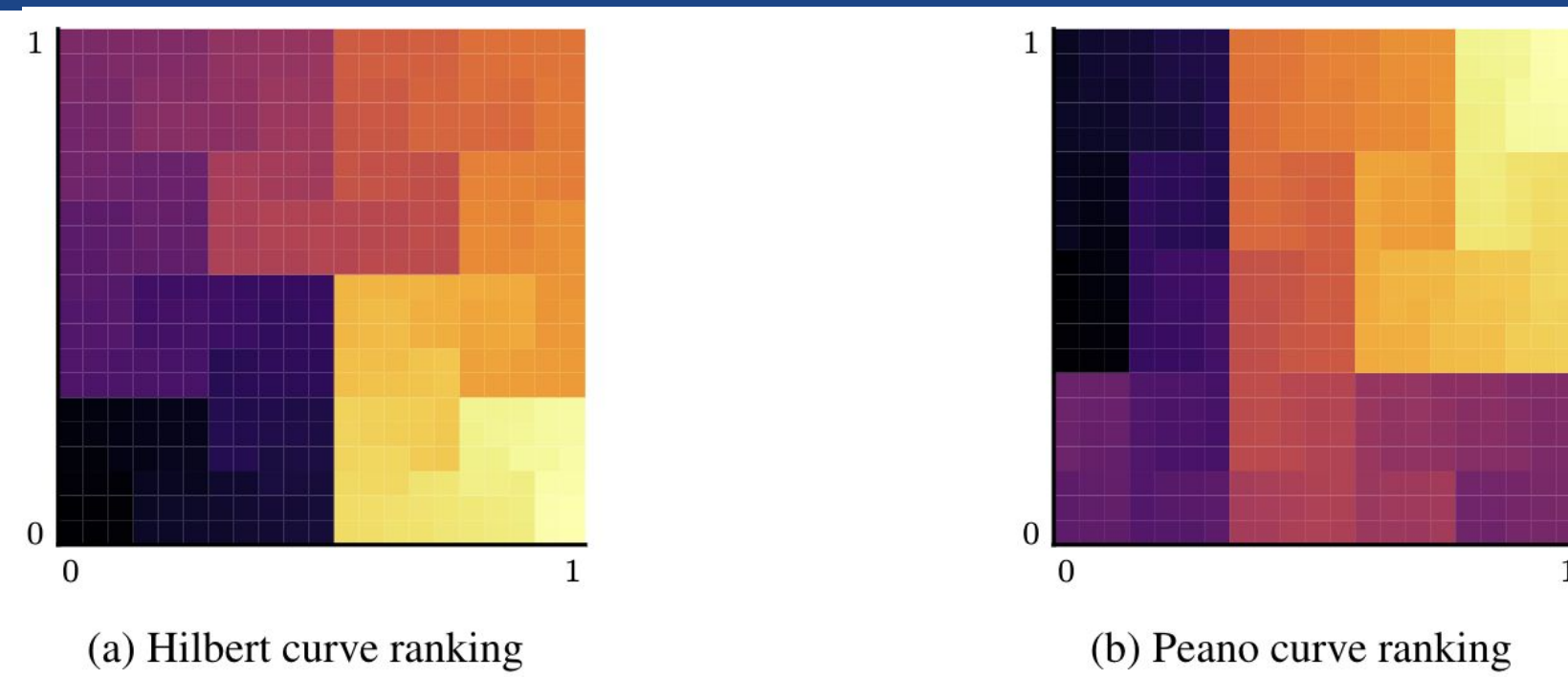
- This paper presents **Multivariate Simulation Based Calibration (mSBC)**, a multivariate extension to SBC [1], which assesses inference quality using the approximate posterior averaged across data simulated from the prior.
- For each iteration, mSBC only requires a bag of samples from the inferred parameter posterior and a sample from parameter prior.
- This paper addresses how to run mSBC on models with a mixture of continuous and discrete parameters; open-universe models where the number of variables is uncertain; and models where existence of one or more parameters in a given trace is uncertain.

## Simulation Based Calibration (SBC)

$$\underbrace{\pi(\theta)}_{\text{prior}} = \int \underbrace{\pi(\theta | \tilde{y})}_{\text{computed posterior}} \underbrace{\pi(\tilde{y} | \theta)}_{\text{data simulation}} \underbrace{\pi(\tilde{\theta})}_{\text{parameter sampled from prior}} d\tilde{y}d\tilde{\theta}$$

If parameters are sampled from the prior  
 AND  
 Data is simulated using the sampled parameters  
 THEN  
 The inferred parameter posteriors should in expectation (w.r.t. data) match the parameter priors.

## Space Filling Curves



## Multivariate Ranks

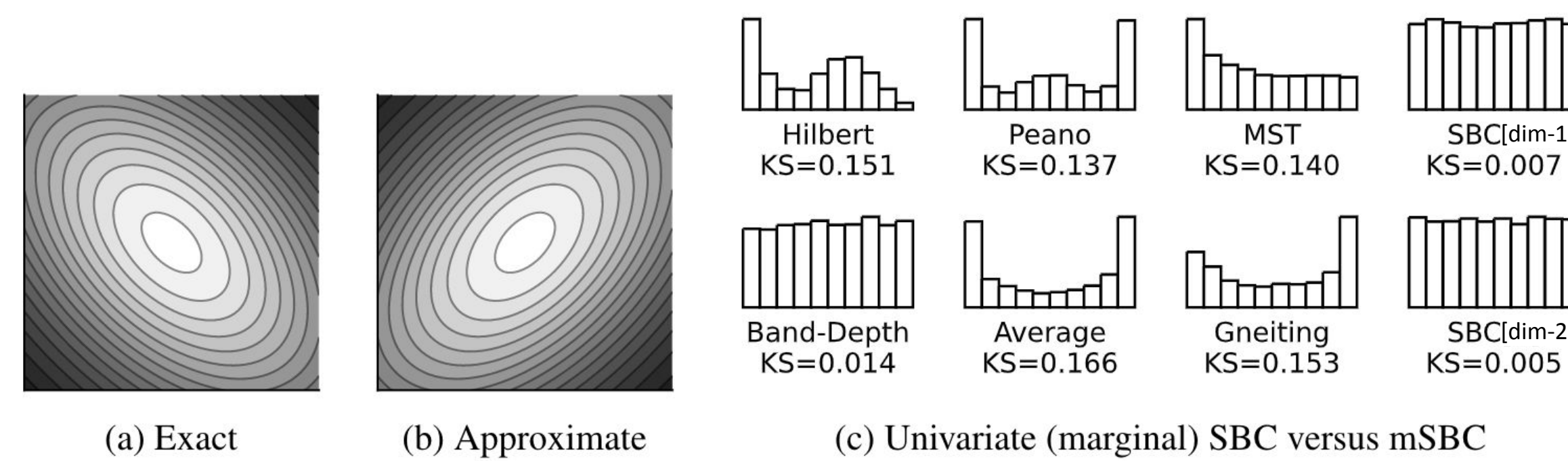
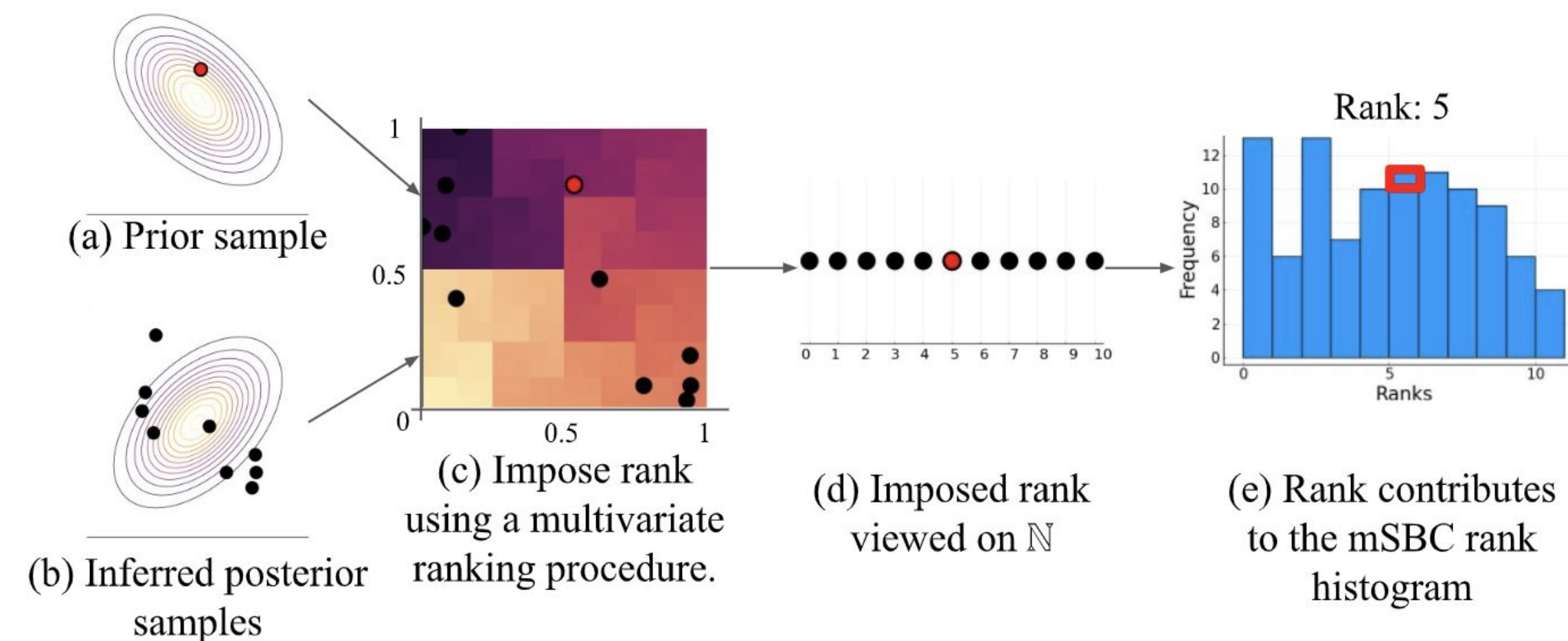
A *multivariate ranking function* is defined to be a function

$r_N: \mathbb{R}^{d \times N} \rightarrow S_N$  for some  $d, N \geq 1$ , which maps  $N$  points of  $d$ -dimensional real space to a permutation of  $[N]$ , with the natural symmetry condition.

Existing multivariate ranking functions.

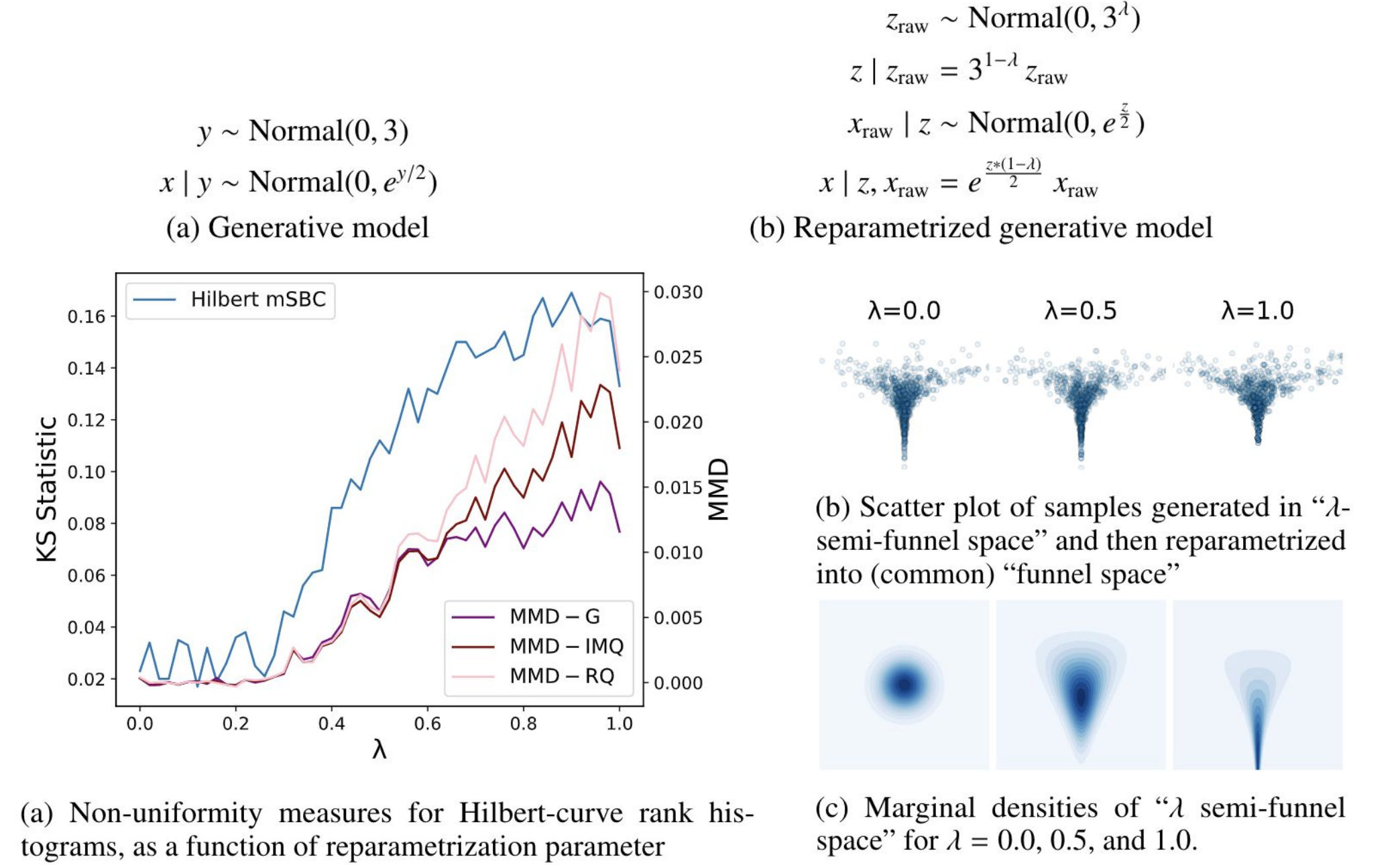
Name	Ranking function ( $r_N(x_0, \dots, x_{N-1})$ for $x_i \in \mathbb{R}^d$ )	Properties
MST [3]	$\arg \text{sort}_{i \in [N]} (\text{length of Euclidean MST on the set } x \setminus x_i)$	Center-outward
Band-Depth [3]	$\arg \text{sort}_{i \in [N]} \frac{1}{d} \sum_{1 \leq k \leq d, 1 \leq j_1 < j_2 \leq N} \mathbb{1}\{\min\{x_{j_1}^k, x_{j_2}^k\} \leq x_i^k \leq \max\{x_{j_1}^k, x_{j_2}^k\}\}$	Center-outward
Gneiting [2]	$\arg \text{sort}_{i \in [N]} \sum_{j=1}^N \mathbb{1}\{x_j^k \leq x_i^k\}$	Reduces to SBC for $d = 1$ .
Average [3]	$\arg \text{sort}_{i \in [N]} \frac{1}{d} \sum_{k=1}^d \sum_{j=1}^N \mathbb{1}\{x_j^k < x_i^k\}$	Reduces to SBC for $d = 1$ .

## Multivariate Simulation Based Calibration (mSBC)



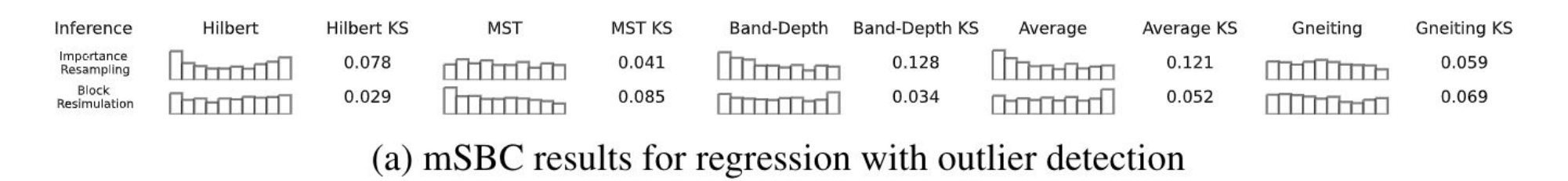
## Case Studies

### Neal's Funnel



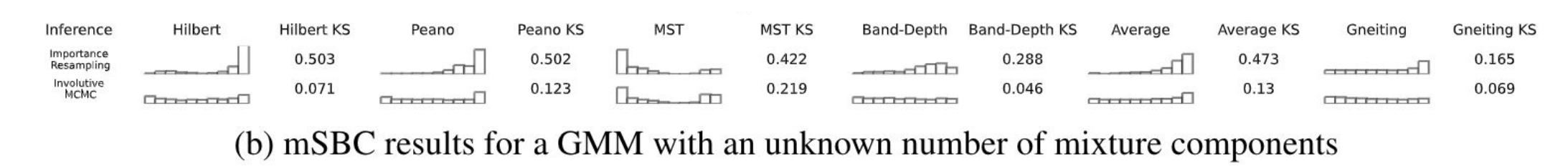
### Linear regression with outliers

```
function block_resimulation_kernel(tr)
    line_params = Gen.select(:noise, :slope, :
        ↪ intercept)
    (tr, _) = mh(tr, line_params)
    (xs, _) = get_args(tr)
    n = length(xs)
    for i=1:n
        (tr, _) = mh(tr, Gen.select((:
            ↪ is_outlier, i)))
    end
    (tr, _) = mh(tr, Gen.select(:prob_outlier))
    tr
end
```



### Gaussian mixture model with an unknown number of mixture components

```
function imcmc_kernel(tr)
    tr, = mh(tr, gibbs_update_w, ())
    tr, = mh(tr, gibbs_update_xi, ())
    tr, = mh(tr, gibbs_update_means, ())
    tr, = mh(tr, gibbs_update_vars, ())
    tr, = mh(tr, gibbs_update_allocations,
        ↪ ())
    tr, = split_merge(tr)
    tr
end
```



## Validity of mSBC

**Theorem 1** Suppose  $(R_n)_{n \in \mathbb{N}}$  is a ranking procedure. Sample  $\theta \sim \pi(\theta)$  and  $y | \theta \sim \pi(y | \theta)$ . Let  $L \geq 1$ , and for each  $\ell = 1, \dots, L$ , sample  $\theta_\ell | y \sim \pi(\theta | y)$ . Assume that  $\pi(\theta | y)$  is a continuous measure with probability 1. Define  $\sigma = R_{L+1}(\theta, \theta_1, \dots, \theta_L)$ , so that  $\sigma(0)$  is the rank of  $\theta$  with respect to  $\{\theta_1, \dots, \theta_L\}$ . Then  $\sigma(0)$  is distributed uniformly on  $[L + 1]$ .

Non-uniform rank histogram  $\Rightarrow$  Incorrect inference

## References

- Talts, Sean, Michael Betancourt, Daniel Simpson, Aki Vehtari, and Andrew Gelman. "Validating Bayesian inference algorithms with simulation-based calibration." arXiv preprint arXiv:1804.06788 (2018).
- Gneiting, Tilmann, Larissa I. Stanberry, Eric P. Grit, Leonhard Held, and Nicholas A. Johnson. "Assessing probabilistic forecasts of multivariate quantities, with an application to ensemble predictions of surface winds." Test 17, no. 2 (2008): 211-235.
- Thorarindottir, Thordis L., Michael Scheuerer, and Christopher Heinz. "Assessing the calibration of high-dimensional ensemble forecasts using rank histograms." Journal of computational and graphical statistics 25, no. 1 (2016): 105-122.