

Learning Robust Statistics for Simulation-based Inference under Model Misspecification

Daolang Huang^{*1}, Ayush Bharti^{*1}, Amauri Souza¹, Luigi Acerbi², Samuel Kaski^{1,3}

Department of Computer Science, Aalto University¹

Department of Computer Science, University of Helsinki²

Department of Computer Science, University of Manchester³

December 1, 2023

Simulation-based inference

- Data $\mathbf{y} = \{y_i\}_{i=1}^n \subseteq \mathbb{R}^d$ denoted by empirical distribution \mathbb{Q}^n
- Simulator-based model $\mathcal{P}_\Theta = \{\mathbb{P}_\theta : \theta \in \Theta\}$
- \mathbb{P}_θ is intractable, but sampling $\mathbf{x} \sim \mathbb{P}_\theta$ is straightforward
- **Aim:** Estimate θ given data \mathbf{y}

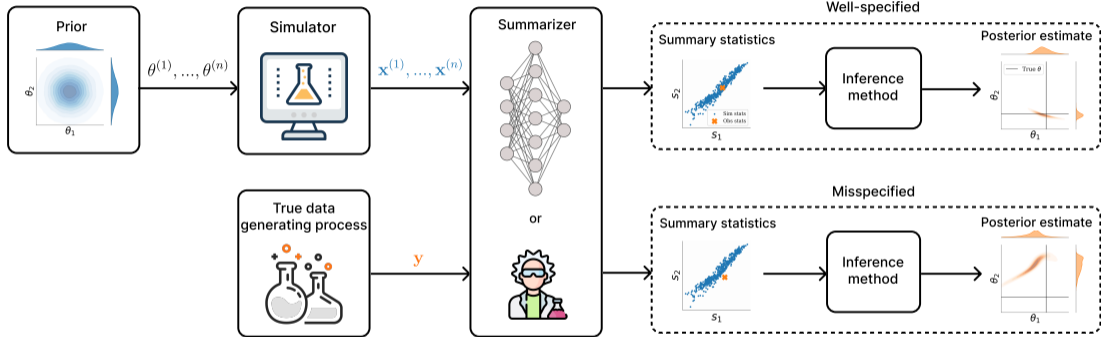
Simulation-based inference

- Data $\mathbf{y} = \{y_i\}_{i=1}^n \subseteq \mathbb{R}^d$ denoted by empirical distribution \mathbb{Q}^n
- Simulator-based model $\mathcal{P}_\Theta = \{\mathbb{P}_\theta : \theta \in \Theta\}$
- \mathbb{P}_θ is intractable, but sampling $\mathbf{x} \sim \mathbb{P}_\theta$ is straightforward
- **Aim:** Estimate θ given data \mathbf{y}
- **Solution:** methods based on distances like approximate Bayesian computation (ABC); methods based on deep neural networks like neural posterior estimation (NPE)

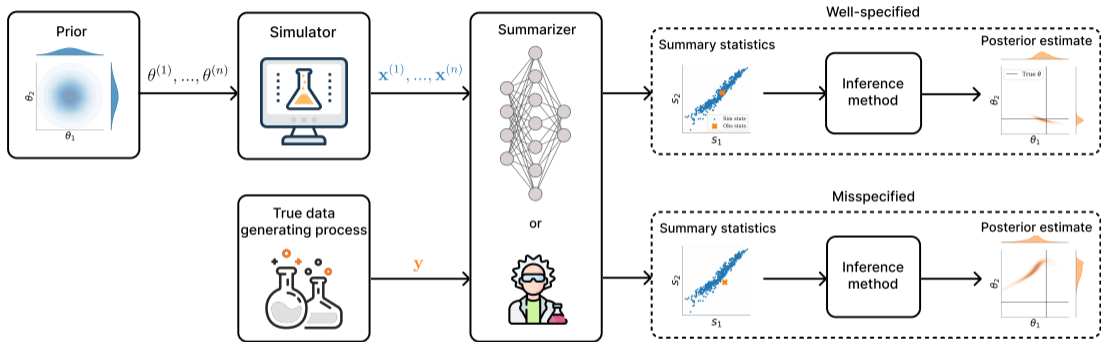
Simulation-based inference

- Data $\mathbf{y} = \{y_i\}_{i=1}^n \subseteq \mathbb{R}^d$ denoted by empirical distribution \mathbb{Q}^n
- Simulator-based model $\mathcal{P}_\Theta = \{\mathbb{P}_\theta : \theta \in \Theta\}$
- \mathbb{P}_θ is intractable, but sampling $\mathbf{x} \sim \mathbb{P}_\theta$ is straightforward
- **Aim:** Estimate θ given data \mathbf{y}
- **Assumption:** Model is “correct”, i.e., $\mathbb{Q}^n \in \mathcal{P}_\Theta$
- **Problem:** Model misspecification, i.e. $\mathbb{Q}^n \notin \mathcal{P}_\Theta \Rightarrow \nexists \theta \in \Theta$ s.t. $\mathbb{P}_\theta = \mathbb{Q}^n$
 - ▶ Stochasticity in data collection process (outliers, missing data, broken independence assumption, etc.)
 - ▶ “All models are wrong...”
- **Even more problem:** Inference is based on simulation from misspecified model!

Inference is based on summary statistics



Inference is based on summary statistics



Insight 1: Under misspecification, observed statistic goes outside the set of simulated statistics
 \Rightarrow SBI methods have to generalize outside their training data

Insight 1: Under misspecification, observed statistic goes outside the set of simulated statistics
⇒ SBI methods have to generalize outside their training data

Insight 2: Even if model is misspecified ($\mathbb{Q}^n \notin \mathcal{P}_\Theta$), it may be well-specified w.r.t the statistics

- Example: Gaussian model, skewed data
- Misspecified if statistics are sample mean and sample skewness
- Well-specified if statistics are sample mean and sample variance
- If we pick statistics appropriately, we can be robust!

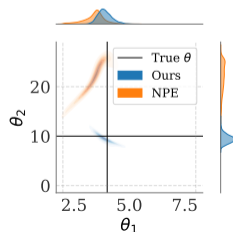
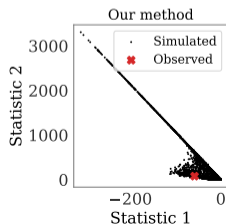
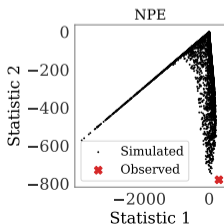
Learning robust statistics for SBI

proposed loss = usual loss + $\lambda \mathcal{D}(\text{simulated statistics}, \text{observed statistic})$

Learning robust statistics for SBI

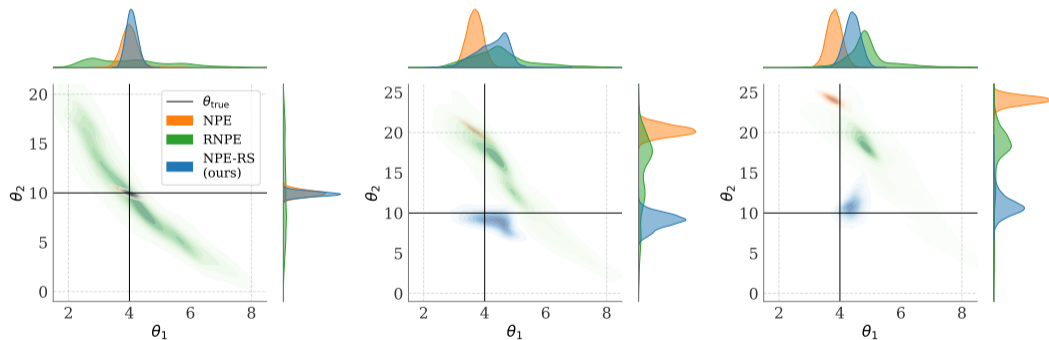
proposed loss = usual loss + $\lambda \mathcal{D}(\text{simulated statistics, observed statistic})$

- For ABC or other SBI methods, usual loss is autoencoder's reconstruction loss
- For NPE, statistics and posterior can be learned jointly
- We want \mathcal{D} to be outlier-robust. Hence, maximum mean discrepancy.
- Regularizer λ : encodes trade-off between accuracy and robustness



Results

- **Ricker model:** 2 parameters
- **Inference method:** Neural posterior estimation (NPE)
- **ϵ -contamination model:** $Q = (1 - \epsilon)\mathbb{P}_{\theta_{\text{true}}} + \epsilon\mathbb{P}_{\theta_c}$



(a) Well-specified ($\epsilon = 0$)

(b) Misspecified ($\epsilon = 10\%$)

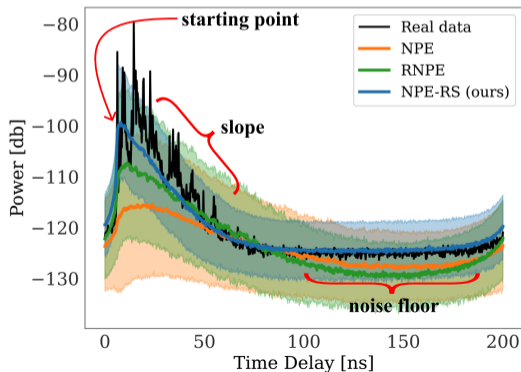
(c) Misspecified ($\epsilon = 20\%$)

Results

Application to real data

Radio propagation example

- 4 parameters
- Data dimension: 801
- Model misspecified due to broken iid assumption



- We propose a simple solution for tackling misspecification of simulator-based models.
- Our method can be applied to any SBI method that utilizes summary statistics.
- Our method only has one hyperparameter balancing efficiency and robustness.
- We show robustness under misspecified scenarios with both synthetic and real-world data.