

# Learning Robust Statistics for Simulation-based Inference under Model Misspecification

Daolang Huang<sup>\*1</sup>, Ayush Bharti<sup>\*1</sup>, Amauri Souza<sup>1</sup>, Luigi Acerbi<sup>2</sup>, Samuel Kaski<sup>1,3</sup>

Department of Computer Science, Aalto University<sup>1</sup>

Department of Computer Science, University of Helsinki<sup>2</sup>

Department of Computer Science, University of Manchester<sup>3</sup>

December 1, 2023

# Simulation-based inference

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

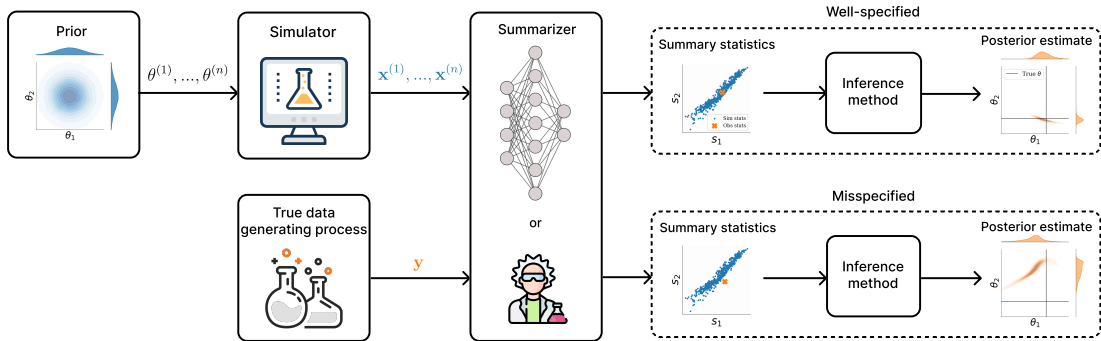
# Simulation-based inference

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

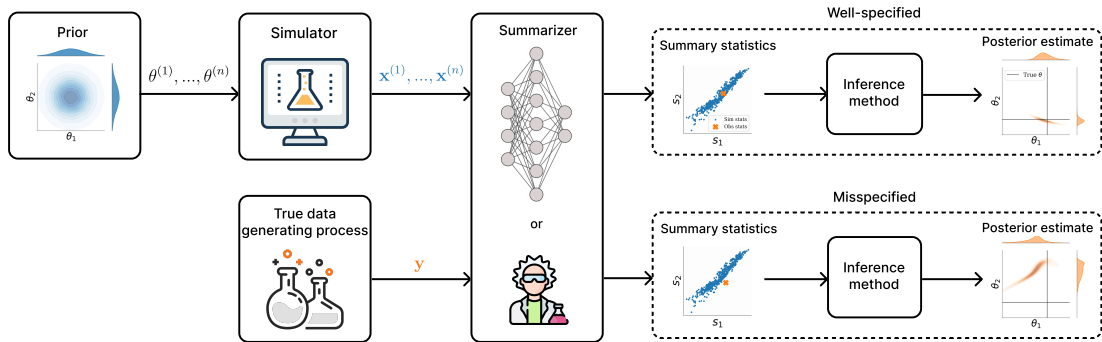
# Simulation-based inference

- Data  $\mathbf{y} = (y_i)_{i=1}^n \in \mathbb{R}^d$  denoted by empirical distribution  $Q^n$
- Simulator-based model  $P_\Theta = \{P_\theta : \theta \in \Theta\}$
- $P_\theta$  is intractable, but sampling  $\mathbf{x} \sim P_\theta$  is straightforward
- **Aim:** Estimate  $g$  given data  $\mathbf{y}$
- **Assumption:** Model is “correct”, i.e.,  $Q^n \in P_\Theta$
- **Problem:** Model misspecification, i.e.  $Q^n \notin P_\Theta$  ) @  $\theta \in \Theta$  s.t.  $P_\theta = Q^n$ 
  - ▶ 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

) 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 ( $Q^n \not\approx P_\Theta$ ), 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

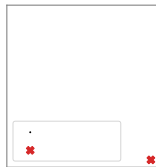
proposed loss = usual loss +  $D(\text{simulated statistics, observed statistic})$



# Learning robust statistics for SBI

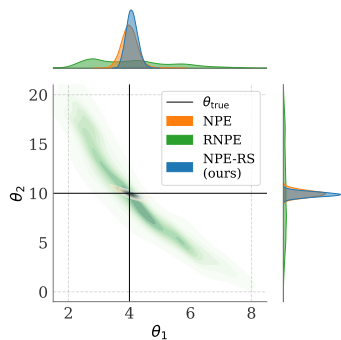
proposed loss = usual loss +  $D(\text{simulated statistics, observed statistic})$

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

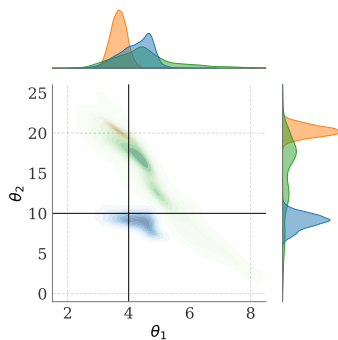


# Results

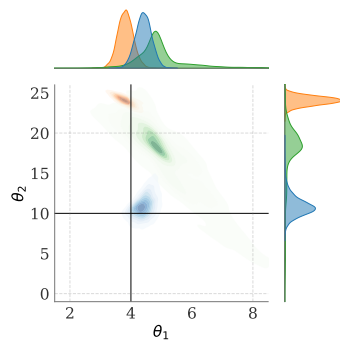
- **Ricker model:** 2 parameters
- **Inference method:** Neural posterior estimation (NPE)
- **-contamination model:**  $Q = (1 - \epsilon)P_{\theta_{\text{true}}} + \epsilon P_{\theta_c}$



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



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



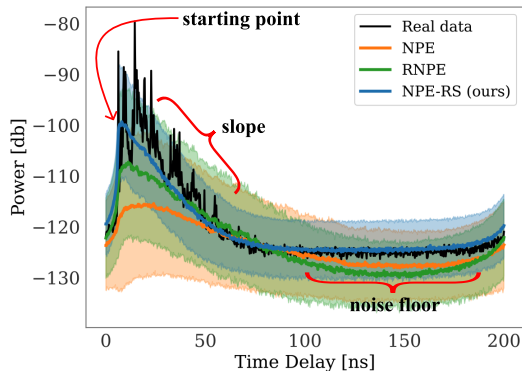
(c) Misspecified ( $\epsilon = 20\%$ )

# 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.