Simulation-based Bayesian Inference from Privacy Protected Data

¹Department of Statistics, Purdue University ²University of Chinese Academy of Sciences



Introduction

Many modern statistical analysis and machine learning applications require training models on sensitive user data. Differentially private algorithms inject calibrated noise into the confidential data or during the data analysis process to produce privacyprotected datasets or queries. This work proposes simulation-based inference methods from privacyprotected datasets.

Notations

- $\theta \in \Theta$: model parameter
- $\pi(\theta)$: prior of the model parameter
- $x = (x_1, \cdots, x_n) \in \mathbb{X}^n$: confidential database
- $f(x \mid \theta)$: intractable likelihood function as a 'simulator'
- s_{dp} : differentially private queries for the confidential database x

Likelihood-free Inference

Given observed dataset x^{o} , how to learn the parameter posterior $\pi(\theta \mid x^o) \propto \pi(\theta) f(x^o \mid \theta)$?

- Approximate Bayesian Computation (ABC) 1: Circumvent likelihood evaluations with simulations.
- Accept $\theta^{(i)}$ if $d(x^{(i)}, x^o) < \epsilon$, where $x^{(i)} \sim f(x \mid \theta^{(i)})$
- Depends on distance function $d(\cdot, \cdot)$; Requires multiple simulations
- SMC-ABC improves ABC by specifying a sequence of intermediate target distributions
- Neural density estimation (NDE) [2]: approximate $\pi(\theta \mid x^o)$ with neural density $\{q_{\phi}(\theta \mid x)\}_{\phi}$ by an architecture q and weight parameters ϕ .

$$\hat{\phi} = \arg\min_{\phi} \mathbb{E}_{p(\theta, x)} \left[-\log q_{\phi}(\theta \mid x) \right]$$

• Training relies on Monte Carlo approximations

Yifei Xiong¹, Nianqiao Phyllis Ju¹, Sanguo Zhang²



Differential Privacy

Let $\eta(s_{dp} \mid x)$ be the conditional density of the private output s_{dp} given the confidential data x. We say η satisfies ϵ -DP if, for all possible values of s_{dp} and for neighboring datasets x, x',

 $\frac{\int_A \eta(s_{dp} \mid x) \, \mathrm{d}s_{dp}}{1} \leq \exp(\epsilon), \quad \forall A \subseteq \operatorname{Range}(\eta).$ $\int_A \eta(s_{\mathrm{dp}} \mid x') \, \mathrm{d}s_{\mathrm{dp}}$

Inference with Privatized Queries

Our goal is to approximate the posterior distribution of θ given privatized queries s_{dp} .

$$egin{aligned} \pi(heta \mid s_{\mathrm{dp}}) &\propto \pi(heta) f(s_{\mathrm{dp}} \mid heta) \ &\propto \pi(heta) \int_{\mathbb{X}^n} f(x \mid heta) \eta(s_{\mathrm{dp}} \mid x) \mathrm{d}x \end{aligned}$$

- Intractable $f(x \mid \theta)$ and integration over \mathbb{X}^n
- Data-augmentation MCMC methods [3] work when $f(x \mid \theta)$ is tractable

Our Methods

We present two complementary approaches: **Private data likelihood estimation (PLE)**: approximate the private data marginal likelihood $f(s_{dp} \mid \theta)$. Minimizing

 $\mathbb{E}_{\pi(\theta)} \left[\mathcal{D}_{\mathrm{KL}} \left(f(s_{\mathrm{dp}} \mid \theta) \| q_{\phi}(s_{\mathrm{dp}} \mid \theta) \right) \right]$ is equivalent to minimizing $\ell_{\text{PLE}}(\phi) = \mathbb{E}_{p(\theta, x)} \left| - \int_{\mathbb{S}} \eta(s_{\text{dp}} \mid x) \log q_{\phi}(s_{\text{dp}} \mid \theta) \mathrm{d}s_{\text{dp}} \right|$ up to a constant independent of ϕ . The posterior approximation can be $\hat{\pi}_{\text{PLE}}(\theta) \propto \pi(\theta) q_{\hat{\phi}}(s_{\text{dp}} \mid \theta)$. **Private data posterior estimation (PPE)**: approximate $\pi(\theta \mid s_{dp})$ directly. Minimizing $\mathbb{E}_{\pi(x)} \left[\mathcal{D}_{\mathrm{KL}} \left(\pi(\theta \mid s_{\mathrm{dp}}) \| q_{\phi}(\theta \mid s_{\mathrm{dp}}) \right) \right]$ is equivalent to minimizing

 $\ell_{\text{PPE}}(\phi) = \mathbb{E}_{p(\theta, x)} \left| - \int_{\mathbb{S}} \eta(s_{\text{dp}} \mid x) \log q_{\phi}(\theta \mid s_{\text{dp}}) \mathrm{d}s_{\text{dp}} \right|$



• Randomized quasi-Monte Carlo (RQMC) can generate correlated, low-discrepancy sequences to reduce the RMSE to $\mathcal{O}(M^{-1+\delta})$ under mild conditions.

HSW 10⁻² 10⁻³



estimator q_{ϕ} into the proposal distribution of the next training round • Sequential training procedures can gradually move q_{ϕ} towards high-density regions of the private data posterior, and thus achieve good accuracy with fewer samples from the simulator

Published in Transactions on Machine Learning Research (TMLR)

Nested RQMC Estimators





Sequential Neural Estimations

• We incorporate the current neural density



Experiments

- PMLR, 2021.

arXiv:2310.12781

Figure: Inference on the SIR model with 1,000 simulations per round. Top: Convergence of posterior estimations; Bottom: Approximation accuracy by SPPE (orange) and SPLE (red) against the number of rounds.

Figure: Posterior comparison on the Bayesian linear regression. Grey: Ground truth posterior; orange: SPPE; red: SPLE.

References

[1] Mark A Beaumont, Jean-Marie Cornuet, Jean-Michel Marin, and Christian P Robert. Adaptive approximate Bayesian computation. Biometrika, 96(4):983–990, 2009. [2] Jan-Matthis Lueckmann, Jan Boelts, David Greenberg, Pedro Goncalves, and Jakob Macke. Benchmarking simulation-based inference. In International Conference on Artificial Intelligence and Statistics, pages 343–351.

[3] Nianqiao Ju, Jordan Awan, Ruobin Gong, and Vinayak Rao. Data augmentation MCMC for Bayesian inference from privatized data. Advances in Neural Information Processing Systems, 35:12732–12743, 2022.