

Adversarial Likelihood-Free Inference on Black-Box Generator

Dongjun Kim, Weonyoung Joo, Seungjae Shin, Kyungwoo Song, Il-Chul Moon
KAIST, Republic of Korea
{dongjoun57, es345, tmdwo0910, gtshs2, icmoon}@kaist.ac.kr

Abstract

Generative Adversarial Network (GAN) can be viewed as an implicit estimator of a data distribution, and this perspective motivates using the adversarial concept in the true input parameter estimation of black-box generators. While previous works on *likelihood-free inference* introduces an implicit proposal distribution on the generator input, this paper analyzes theoretic limitations of the proposal distribution approach. On top of that, we introduce a new algorithm, Adversarial Likelihood-Free Inference (ALFI), to mitigate the analyzed limitations, so ALFI is able to find the posterior distribution on the input parameter for black-box generative models. We experimented ALFI with diverse simulation models as well as pre-trained statistical models, and we identified that ALFI achieves the best parameter estimation accuracy with a limited simulation budget.

1 Introduction

Generative Adversarial Network (GAN) is highlighted recently for its success on the implicit estimation of the data distribution. In GAN, the generator is jointly learned with the discriminator, so the generator becomes a trainable and fine-tunable model. In contrast to training both the generator and the discriminator in GAN, there has been a line of work on applying the adversarial framework on the pre-trained and fixed generator to estimate the optimal input of the generator [1]. For example, a simulation model can be considered as a generator that is not successfully integrated into the adversarial concept. Researchers are interested in inferring the posterior distribution of a simulation input parameter with a snapshot of a validation observation from the real-world [2, 3, 1].

Before we move on, we define a *black-box generative model* to present our interested generator type clearly. The *black-box generative model* (g) indicates a generative model that has three properties. First, the model's internal structure is designed *before* the inference stage by domain experts or as another statistical model. Second, the internal coefficients (ω) do not change since the coefficients are obtained by the domain-specific knowledge or through a separate learning process. Third, the internal process contains the inherent stochasticity (u) that forms various sample paths for each of the model execution. For example, continuous, discrete, and agent-based simulation models can be black-box generators that produce the stochastic trajectories of modeled states from simulations. As another example, a pre-trained and fixed de-convolutional neural network can be another practical case of the black-box generators.

The research on the black-box generator emphasizes the regeneration of the observation, which is formulated as the inference on the posterior distribution $p(\theta|x_{obs})$, where $\theta \in \Theta \subseteq \mathbb{R}^d$ is a d -dimensional generator input parameter on a compact space Θ , and where $x_{obs} \in \mathbb{P}_r$ is the *single* instance of the real-world data to reconstruct from the real-world data distribution \mathbb{P}_r . The likelihood

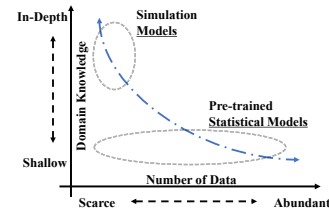


Figure 1: Black-box generative models

$p(x_{obs}|\theta)$ of the black-box generator requires the integration over the aforementioned generation stochasticity, or a nuisance variable $u \in \mathbb{R}^m$, because u determines the sample path of the generator. However, u is unknown in general, and the integration over u is likely to be intractable. Therefore, the Bayesian inference under black-box generators mainly focuses on estimating the intractable likelihood, called *likelihood-free inference*. This paper provides theoretic analysis on previous research, and suggests a new algorithm of *likelihood-free inference* under the adversarial setting.

2 Previous Research

In *likelihood-free inference* community, the summary statistics $s : \mathbb{R}^p \rightarrow \mathbb{R}^q$ extracts a set of statistics from either observation x_{obs} or generated fake data $g(\theta, u|\omega)$, and the discrepancy function $d : \mathbb{R}^q \times \mathbb{R}^q \rightarrow \mathbb{R}$ measures how $s(g(\theta, u|\omega))$ deviates from $s(x_{obs})$. Throughout the paper, we assume that the function s extracts the *sufficient* statistics [4–6] that assure the identity of the posteriors under the summary statistics as $p(\theta|x_{obs}) = p(\theta|s(x_{obs}))$. In addition, we assume the summary statistics to be embedded in the generator, so the generated data $g(\theta, u|\omega)$ or the observation x_{obs} becomes the extracted q -dimensional summary statistics of the generated raw data or the observed raw data, respectively. The likelihood, $p(x_{obs}|\theta) = \int p(x_{obs}|u, \theta)p_U(u)du = \int \delta(d(g(\theta, u|\omega), x_{obs}))p_U(u)du$, is intractable since the distribution $p_U(u)$ and the level set, $\{u | d(g(\theta, u|\omega), x_{obs}) = 0\}$, are generally not known in a black-box generator.

Approximate Bayesian Computation: Approximate Bayesian Computation (ABC) [7–11] estimates the likelihood through Monte-Carlo methods by approximating the singular distribution δ with mollifier kernels $\{K_\epsilon\}_{\epsilon>0}$, so that $p(x_{obs}|\theta) = \lim_{\epsilon \downarrow 0} \mathbb{E}_{u \sim p_U} [K_\epsilon(u; \theta)]$. A case of the ABC algorithm is the Rejection ABC [7] that uses the boxcar kernel [12]: $p(x_{obs}|\theta) = \lim_{\epsilon \downarrow 0} \frac{1}{|B_\epsilon(x_{obs})|} \mathbb{E}_u [1_{B_\epsilon(x_{obs})}(g(\theta, u|\omega))]$, where $B_\epsilon(x_{obs})$ is the ϵ -ball $\{x : d(x, x_{obs}) < \epsilon\}$.

Bayesian Optimization Likelihood-Free Inference: Bayesian Optimization Likelihood-Free Inference (BOLFI) [3] infers the predictive distribution of the discrepancy $d(g(\theta, u|\omega), x_{obs})$, using a Gaussian process regression with a dataset $\mathcal{D}_{1:t} = \{(\theta_i, d(g(\theta_i, u_i|\omega), x_{obs}))\}_{i=1}^t$, and BOLFI samples a new parameter θ_{t+1} by the Bayesian optimization. The likelihood of the discrepancy being less than a threshold ϵ is proportional to $p(x_{obs}|\theta) \propto \Phi(\frac{\epsilon - \hat{\mu}(\theta)}{\hat{\sigma}(\theta)})$, where Φ is the cumulative distribution function of the standard Gaussian distribution, and where $\hat{\mu}$ and $\hat{\sigma}$ are the mean and the standard deviation of the predictive distribution of the Gaussian process regression.

3 Preliminary: Problems of Implicit Proposal Distribution

3.1 Motivation of Implicit Proposal Distribution

Figure 2 (a) illustrates the discrepancy (d) landscape of the Susceptible-Infectious-Recovered (SIR) simulation model [13], where the discrepancy is defined as the Euclidean measure. The choice of the summary statistics and the discrepancy measure are crucial in *likelihood-free inference*, since Figure 2 (b) illustrates the failure of inferring the posterior distribution if the discrepancy landscape is highly rugged, and if the landscape has a plateau near the true parameter θ^* .

This failure leads *likelihood-free inference* community to investigate the adaptive selection of summary statistics and a discrepancy measure. The *likelihood-free inference* under adversarial setting partially solves the selection problem by constructing discrepancy measure as a discriminator network. Besides, in some black-box generators, it is possible to put raw data into the discriminator network, without extracting the summary statistics. However, *likelihood-free inference* under the adversarial framework is rarely proposed since the gradient with respect to the input parameter θ is not backpropagated through an implicitly defined generator that has no closed-form solution. Recently, Louppe et al. [1] overcomes the backpropagation issue by suggesting the Adversarial Variational Optimization (AVO).

3.2 Adversarial Variational Optimization

The AVO in Figure 3 (a,c) introduces the implicit proposal distribution $p_\psi(\theta)$ for the black-box model input parameter, which enables the backpropagation through a non-differentiable black-box generator by switching the optimization target variable from θ to ψ . AVO has a fixed model internal coefficients

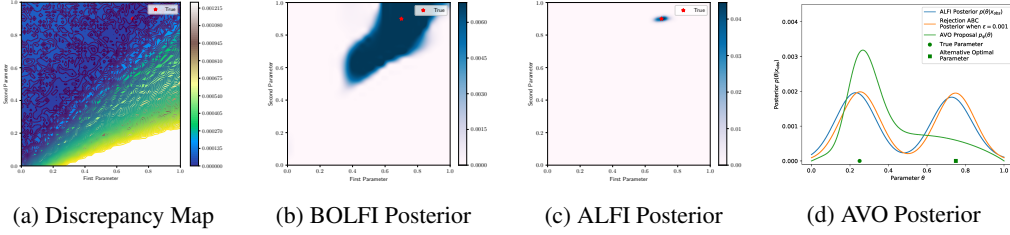


Figure 2: (a-c) The Susceptible-Infectious-Recovered simulation model result. (a) The value at (θ_1, θ_2) represents the discrepancy $\|x_{obs} - g((\theta_1, \theta_2), u)\|_2$. (b) BOLFI finds a huge area as the candidate region for the true parameter θ^* . (c) ALFI captures the true parameter θ^* within a tiny region. (d) The inferred posterior with simulation model $g(\theta, u) = (\theta - 0.5)^2 + u$, where $u \sim \mathcal{N}(0, 10^{-4})$.

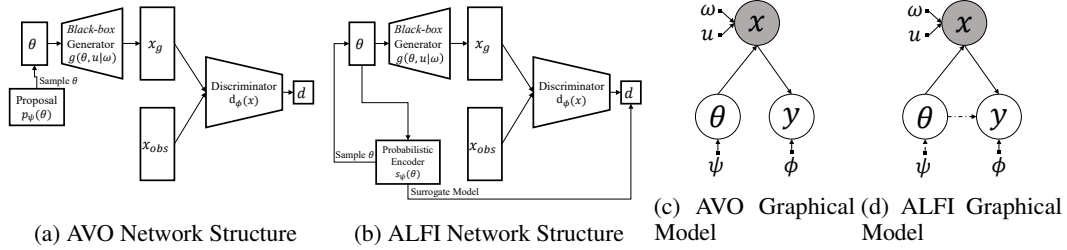


Figure 3: AVO and ALFI comparison

ω , so the generator distribution \mathbb{P}_g can only be adjusted by inferring the input distribution $p_\psi(\theta)$ in order to approximate the data distribution \mathbb{P}_r .

The minimax function of AVO is given by $V(\psi, \phi) = \mathbb{E}_{x \sim \mathbb{P}_r} [\log d_\phi(x)] + \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [\log(1 - d_\phi(\tilde{x}))]$, where a fake sample \tilde{x} from the generator distribution \mathbb{P}_g is the output of a black-box generator $g(\theta, u|\omega)$, with a sampled input $\theta \sim p_\psi$ and a sampled nuisance variable $u \sim p_U$. It should be noted that a discriminator $d_\phi(x)$ guides the proposal distribution to enforce the approximation of \mathbb{P}_g toward \mathbb{P}_r . The backpropagation through the generator is calculated by the REINFORCE algorithm, $\nabla_\psi V = \mathbb{E}_{\theta \sim p_\psi} [\nabla_\psi p_\psi(\theta) \mathbb{E}_{u \sim p_U} [\log(1 - d_\phi(g(\theta, u|\omega)))]]$.

3.3 Gradient Vanishing Problem of Proposal Distribution Approach

The implicit proposal approach in AVO has two significant drawbacks. The first problem is the *gradient vanishing problem* [14]. Our analysis is different from the previous theory [14] on the gradient with respect to ω , the internal parameters of the black-box generator; since we analyze the gradient with respect to ψ , the proposal parameters of the black-box generator, see Appendix A.

Proposition 1 (Gradient Zero On Optimal Discriminator). *Let $a : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ and $b : [0, 1] \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ be strictly increasing and strictly decreasing functions, respectively. Assume the empirical value function is $V_{emp}(\psi, \phi) = \mathbb{E}_{x \in \mathbb{P}_{emp}} [a(d_\phi(x))] + \mathbb{E}_{\tilde{x} \in \mathbb{P}_g} [b(d_\phi(\tilde{x}))]$, where \mathbb{P}_{emp} is the empirical real-world data distribution with a single observation x_{obs} .*

- (i) *If a and b are upper bounded, the gradient $\nabla_\psi V_{emp}$ of the empirical value function V_{emp} with respect to ψ is always zero under the optimal discriminator.*
- (ii) *If neither a nor b is upper bounded, the empirical value function attains infinity under the optimal discriminator.*

Remark. *If $a(t) = \log t$ and $b(t) = \log(1 - t)$, the value function becomes the vanilla GAN [15] loss $V(\psi, \phi) = \mathbb{E}_{x \in \mathbb{P}_r} [\log d_\phi(x)] + \mathbb{E}_{\tilde{x} \in \mathbb{P}_g} [\log(1 - d_\phi(\tilde{x}))]$. If $a(t) = t$ and $b(t) = -t$, the value function becomes the Wasserstein GAN [16] loss $V(\psi, \phi) = \mathbb{E}_{x \in \mathbb{P}_r} [d_\phi(x)] - \mathbb{E}_{\tilde{x} \in \mathbb{P}_g} [d_\phi(\tilde{x})]$.*

Proposition 1 analyzes an extreme and special case that the dataset has a single instance x_{obs} , and this scenario occurs in the simulation calibration case (i.e. finding optimal input parameter) because the real-world observation with the same context happens only once. In this case, the *true* value function V is never approximated exactly, and the adversarial framework could gain a saddle point of the empirical value function V_{emp} , at best, not the *true* value function V .

Proposition 2. Let a and b be strictly increasing and strictly decreasing functions, respectively. For any ψ , there exists a constant M such that $\|\nabla_\psi V\|_2 = \|\nabla_\psi V_{emp}\|_2 \leq M|\Theta| \max(|b(\epsilon)|, |b(0)|)$, where $\|\cdot\|_2$ is the Euclidean norm, if the following conditions satisfy.

- (i) The proposal distribution, p_ψ , is differentiable with respect to ψ , and the derivative of $\nabla_\psi p_\psi$ is continuous with respect to θ .
- (ii) The space of the black-box model input parameter, Θ , is compact.
- (iii) The discriminator d_ϕ is ϵ -close to the optimal discriminator d^* of the empirical value function: $\|d_\phi - d^*\| < \epsilon$, where $\|d\| = \sup_x |d(x)|$.

Corollary 1 (Gradient Vanishing Near Optimal Discriminator). Let a and b be strictly increasing and strictly decreasing functions. Assume that 1) the conditions of Proposition 2 hold; 2) $\lim_{\epsilon \rightarrow 0} b(\epsilon) = 0$; and 3) $b(0) = 0$, then for any ψ , the limit of the gradients converge to zero as the discriminator d_ϕ converges to the optimal discriminator d^* : $\lim_{\|d - d^*\| \rightarrow 0} \nabla_\psi V = 0$ and $\lim_{\|d - d^*\| \rightarrow 0} \nabla_\psi V_{emp} = 0$.

Remark. Examples of value functions in Corollary 1 are vanilla GAN with $b(t) = \log(1 - t)$ and Wasserstein GAN with $b(t) = -t$. Note that AVO [1] uses Wasserstein GAN in their released code.

3.4 Implicit Relation Problem of Proposal Distribution Approach

The second problem is the implicit relation between the proposal distribution $p_\psi(\theta)$ and the posterior distribution $p(\theta|x_{obs})$. In the training time, the proposal distribution approximates \mathbb{P}_g to \mathbb{P}_{emp} instead of \mathbb{P}_r , at best. Proposition 3 analyzes marginal form when \mathbb{P}_g equals to \mathbb{P}_{emp} .

Proposition 3 (Implicit Relation). Assume that the generator distribution \mathbb{P}_g equals to the empirical data distribution \mathbb{P}_{emp} . Then, the marginal equivalence holds between the input parameter distribution $p_\psi(\theta)$ and the posterior distribution $p(\theta|x_{obs})$: $\int p(x|\theta)p_\psi(\theta)d\theta = \int p(x|\theta)p(\theta|x_{obs})d\theta$.

Even though the proposal approach succeeds on estimating \mathbb{P}_{emp} through \mathbb{P}_g , the marginal equivalence does not guarantee the equivalence between the optimal proposal distribution $p_{\psi^*}(\theta)$ and the posterior distribution $p(\theta|x_{obs})$. Figure 2 (d) illustrates an example of the non-equivalence between the implicit proposal $p_\psi(\theta)$ and the posterior $p(\theta|x_{obs})$ in AVO. This *implicit relation problem* can be mitigated by estimating the likelihood $p(x_{obs}|\theta)$ directly to infer the posterior distribution $p(\theta|x_{obs})$.

4 Adversarial Likelihood-Free Inference

The *implicit relation problem* causes AVO to be ill-posed, i.e. there could be many candidates for the optimal proposal distribution p_{ψ^*} that achieves \mathbb{P}_g to equal to \mathbb{P}_{emp} . Besides, the *gradient vanishing problem* forces p_ψ to remain to be pre-matured, so the convergence of p_ψ to an optimal proposal p_{ψ^*} is not guaranteed. Therefore, we introduce Adversarial Likelihood-Free Inference (ALFI), in place of proposal distribution approach, by breaking apart *likelihood-free inference* into the likelihood estimation problem and the sampling problem. ALFI estimates the likelihood by Theorem 1 with a surrogate model parametrized by ψ , and ALFI proposes the next inputs via a sampling algorithm.

ALFI consists of the three components: a non-optimizable black-box generator $g : \mathbb{R}^d \times \mathbb{R}^m \rightarrow \mathbb{R}^q$, a discriminator network $d_\phi : \mathbb{R}^q \rightarrow [0, 1]$, and a probabilistic encoder network $s_\psi : \mathbb{R}^d \rightarrow \mathbb{R}^K$, where the probabilistic encoder network forms a surrogate model to estimate the likelihood. The generator $g(\theta, u|\omega)$ is a function of $\theta \in \Theta \subseteq \mathbb{R}^d$ and $u \in \mathbb{R}^m$, conditioned on the fixed generator internal coefficients ω , throughout the inference stage, determined a-priori either by the domain experts or by another statistical model. Algorithm 1 presents three procedures of ALFI: 1) sampling procedure in line 3; 2) evaluation procedure in line 4; and 3) learning procedure in lines 6 and 7.

AVO suffers from *gradient vanishing problem* because AVO formulates *likelihood-free inference* as a saddle point problem with respect to ψ and ϕ of an adversarial value function V . On the other hand, ALFI detours *gradient vanishing problem* by formulating *likelihood-free inference* as a pair of maximization problems of separated losses $\mathcal{L}_d(\phi)$ and $\mathcal{L}_s(\psi)$ on the discriminator d_ϕ and the encoder s_ψ , respectively. In particular, while the discriminator maximizes the adversarially designed GAN loss $\mathcal{L}_d(\phi)$, the encoder updates parameter ψ via Maximum Likelihood Estimation of a loss $\mathcal{L}_s(\psi)$ for a surrogate model. Moreover, the introduction of probabilistic encoder s_ψ in ALFI mitigates *implicit relation problem* by estimating the likelihood directly.

ALFI chooses the Metropolis-Hastings algorithm as a sampling algorithm. The Metropolis-Hastings algorithm is a fast, parallelizable and mathematically well-developed sampler, and we prove Theorem 2 that guarantees the convergence of inhomogeneous Markov chain [17] to the posterior $p(\theta|x_{obs})$.

Algorithm 1: Adversarial Likelihood-Free Inference (ALFI)

Require: Discriminator network d_ϕ , Probabilistic encoder network b_ψ

```

1 for  $t$  steps do
2   for  $m$  steps do
3     Sample the next particles  $\{\theta'_i\}_{i=1}^n$  from the current particles  $\{\theta_i\}_{i=1}^n$ , using
       Metropolis-Hastings algorithm with acceptance ratio  $\bar{A}(\theta', \theta)$  given by Eq. 2
4   Execute the black-box generator with each particle in  $\{\theta'_i\}_{i=1}^n$ 
5   for  $l$  steps do
6     Update discriminator network parameters  $\phi$  from the gradient of Eq. 3
7     Update probabilistic encoder network parameters  $\psi$  from the gradient of Eq. 4

```

4.1 Likelihood Estimation

To estimate the intractable likelihood, we introduce Theorem 1, which states that the likelihood is a density of a 1-dimensional random variable Y_θ , see Appendix B.

Theorem 1. *The likelihood becomes $p(x_{obs}|\theta) = p_{Y_\theta}(d_\phi(x_{obs}))$, where $Y_\theta : \mathbb{R}^m \rightarrow \mathbb{R}$ is the random variable under the map $Y_\theta(u) = d_\phi(g(\theta, u|\omega))$, if $\|d_\phi - d^*\| < 0.5$ with the supremum norm $\|d\| = \sup_{x \in \mathbb{R}^q} |d(x)|$.*

Since the support of the random variable Y_θ is restricted to the unit interval, i.e. $d_\phi(\cdot) \in [0, 1]$, Corollary 2 expands the support to any of either a bounded interval, a semi-infinite interval, or a whole real line, where Corollary 2 is an application of change of variables in the probability densities.

Corollary 2. *Let $h : [0, 1] \rightarrow \mathbb{R}$ be strictly monotonic and continuous with a non-zero derivative at $d_\phi(x_{obs})$, then the likelihood becomes $p(x_{obs}|\theta) = p_{Z_\theta}(h(d_\phi(x_{obs})))|h'(d_\phi(x_{obs}))|$, where $Z_\theta = h(Y_\theta)$ is a transformed random variable of Y_θ under h .*

As the random variable Y_θ embeds the stochastic information of the stochastic nuisance variable u , $p_{Z_\theta}(z)dz$ is the probability of the random variable $h(d_\phi(g(\theta, u)))$ being observed in $[z, z + dz)$. Since the density on Z_θ is intractable because of the implicit nature of u , we use Corollary 3 to formulate the likelihood estimation problem as the shape parameter estimation problem by imposing an explicit parametric distribution on Z_θ .

Corollary 3. *If the random variable Z_θ follows a parametric probability distribution with shape parameters $\mathbf{s}_\theta = \{s_{\theta,k}\}_{k=1}^K$, the likelihood becomes*

$$p(x_{obs}|\theta) = f\left(h(d_\phi(x_{obs})); s_{\theta,1}, \dots, s_{\theta,K}\right) |h'(d_\phi(x_{obs}))|,$$

where $f(\cdot; s_{\theta,1}, \dots, s_{\theta,K})$ is the density of a parametric distribution with shape parameters $\{s_{\theta,k}\}_{k=1}^K$.

Remark. *An example of a parametric distribution is the beta distribution with shape parameters α and β , where h is the identity function. The other example could be the Gaussian distribution with shape parameters μ and σ , where $h(y) = h_0 \circ \tilde{h}(y)$ with $h_0(t) = \frac{-2 \sin 2\pi t}{1 - \cos 2\pi t}$ and $\tilde{h}(y) = 1/(1 + e^{-(y - d_\phi(x_{obs}))})$, or $h(y) = \tilde{h}^{-1}(y)$, see Appendix D for the detailed explanation on h .*

4.2 Acceptance Ratio

Corollary 4. *If the prior distribution on θ is uniform and the random variables Z_θ and $Z_{\theta'}$ follow a parametric probability distribution with shape parameters $\mathbf{s}_\theta = \{s_{\theta,k}\}_{k=1}^K$ and $\mathbf{s}_{\theta'} = \{s_{\theta',k}\}_{k=1}^K$, respectively, the acceptance ratio of the Metropolis-Hastings algorithm of jumping to θ' from θ is*

$$A(\theta', \theta) = \min \left(1, \frac{f(h(d_\phi(x_{obs})); s_{\theta',1}, \dots, s_{\theta',K})}{f(h(d_\phi(x_{obs})); s_{\theta,1}, \dots, s_{\theta,K})} \right). \quad (1)$$

If we define a stochastic process $\{Z_\theta | \theta \in \Theta\}$ to be the collection of the random variables Z_θ , the optimal shape parameters \mathbf{s}_θ become diverse by the random variables Z_θ with different parameters θ . The probabilistic encoder network, $s_\psi(\theta) = \hat{\mathbf{s}}_\theta$, is a surrogate model that estimates the K -dimensional optimal shape parameters \mathbf{s}_θ of the probability distribution for Z_θ .

4.3 Algorithm

At the t -th iteration, the Metropolis-Hastings algorithm samples the next n independent set of parameters $\{\theta_{t+1}^{(i)}\}_{i=1}^n$. The i -th particle $\theta_t^{(i)}$ searches the neighborhood of $\theta_t^{(i)}$ by suggesting an intermediate particle $\tilde{\theta}_{t+1}^{(i)}$ from the symmetric proposal distribution. The intermediate particle $\tilde{\theta}_{t+1}^{(i)}$ will be accepted to be the next parameter $\theta_{t+1}^{(i)}$ by the below probability:

$$\bar{A}(\tilde{\theta}_{t+1}^{(i)}, \theta_t^{(i)}) = \min \left(1, \frac{f(h(d_\phi(x_{obs})); s_\psi(\tilde{\theta}_{t+1}^{(i)}))}{f(h(d_\phi(x_{obs})); s_\psi(\theta_t^{(i)}))} \right). \quad (2)$$

The ratio in Eq. 2 equals to the true acceptance ratio in Eq. 1 when the probabilistic encoder network s_ψ estimates the exact shape parameters s_θ .

After sampling the n independent particles, ALFI puts each particle into the black-box generator to evaluate. Then, the discriminator classifies the n generated fake data with the real data. To maximize $d_\phi(x)$ for $x \in \mathbb{P}_r$ and minimize $d_\phi(\tilde{x})$ for $\tilde{x} \in \mathbb{P}_g$, we use the Wasserstein loss [16]

$$\mathcal{L}_d(\phi) = -\mathbb{E}_{x \sim \mathbb{P}_r} [d_\phi(x)] + \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [d_\phi(\tilde{x})]. \quad (3)$$

We calculate the above expectation $\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g}$ through the Monte-Carlo estimation with sampled fake data $\{g(\theta_{t+1}^{(i)}, u_{t+1}^{(i)}|\omega)\}_{i=1}^n$, where $\{\theta_{t+1}^{(i)}\}_{i=1}^n$ are selected from the Metropolis-Hastings algorithm and $\{u_{t+1}^{(i)}\}_{i=1}^n$ are the sampled nuisance variables that are determined for each execution.

The probabilistic encoder estimates the shape parameters by minimizing the negative log-likelihood of $h(d_\phi(g(\theta, u|\omega)))$ being observed from the parametric distribution with shape parameters $s_\psi(\theta)$. The below expectation $\mathbb{E}_{\theta, u}$ equals to the expectation $\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g}$.

$$\mathcal{L}_s(\psi) = -\mathbb{E}_{\theta, u} \left[\log f \left(h(d_\phi(g(\theta, u))) \right); s_\psi(\theta) \right]. \quad (4)$$

4.4 Convergence of Inhomogeneous Markov Chain

After t iterations of learning, the sampling procedure from the Metropolis-Hastings algorithm is equivalent to sample from the transition kernel, $P_t(\theta'|\theta) = q(\theta'|\theta) \min \left\{ 1, \frac{p_t(\theta'|x_{obs})}{p_t(\theta|x_{obs})} \right\} + \delta(\theta' - \theta) \int q(\tilde{\theta}|\theta) \left(1 - \min \left\{ 1, \frac{p_t(\tilde{\theta}|x_{obs})}{p_t(\theta|x_{obs})} \right\} \right) d\tilde{\theta}$. Here, $q(\theta'|\theta)$ is a symmetric proposal distribution of the Metropolis-Hastings algorithm and $p_t(\theta|x_{obs})$ is the approximate posterior at t -th iteration, where the approximate likelihood, $p_t(x_{obs}|\theta) = f(h(d_{\phi_t}(x_{obs})); s_{\psi_t}(\theta)) |h'(d_{\phi_t}(x_{obs}))|$, is estimated by Corollary 3. The parameter update of the discriminator and the encoder networks causes the transition kernel to be adjusted for every iteration. Therefore, the standard theory on Markov chain with fixed transition kernel [18] is not applicable, which means that the Markov chain is no longer guaranteed to asymptotically follow the posterior distribution $p(\theta|x_{obs})$.

Once the Markov chain does not follow the posterior distribution, the parameter learning may not succeed to estimate the exact likelihood, since the area near the true parameter θ^* could not have been visited in the process of learning. Therefore, we provide a theoretic analysis on the limit behavior of inhomogeneous Markov chain with updating transition kernel, which guarantees the success on *likelihood-free inference* through ALFI structure that integrates the Metropolis-Hastings algorithm with the likelihood estimation networks. Theorem 2 ensures the convergence of distribution for the Markov chain with trainable transition kernel P_t to the posterior distribution, $p(\theta|x_{obs})$.

Theorem 2. Assume that Z_θ follows either beta or Gaussian distribution and the probabilistic encoder network asymptotically estimate the true shape parameters. With the minorization condition [19], sufficiently large m and continuously differentiable black-box generator with respect to θ , the distribution of the inhomogeneous Markov chain uniformly converges to the posterior distribution

$$\lim_{N \rightarrow \infty} \|P_1^m \dots P_t^m - P^\infty\| = 0, \quad (5)$$

where P is the transition kernel that has the posterior, $p(\theta|x_{obs})$, as the unique stationary distribution; and where the operator norm, $\|P_1^m \dots P_t^m - P^\infty\| = \sup_{\theta_0} \|\delta_{\theta_0} P_1^m \dots P_t^m - p(\cdot|x_{obs})\|_{TV}$, is the supremum of the total variation norm, $\|q\|_{TV} = \int |q(\theta)| d\theta$.

Table 1: The performance of *likelihood-free inference* algorithms. The boldface indicates the highest performance among algorithms. Rejection ABC takes simulation budget 10 times more than ALFI.

	TUMOR [20]	SIR [13]	POISSON [21]	STOKES [22]	NPA [23]	MA(2) [24]	M/G/1 [25]	WEALTH [26]
REJECTION ABC (REFERENCE) [8]	5.0 ±1.3	4.0±1.4	1.6±0.4	1.3±0.2	1.9±0.7	2.3±1.0	2.4±0.5	2.9 ±0.6
MCMC ABC [9]	2.2±0.7	2.3±1.0	1.7±0.6	1.3±0.7	1.5±0.6	2.6±1.1	1.6±0.5	2.2±1.0
SMC ABC [10]	4.1±0.7	4.8 ±2.5	1.6±0.8	1.3±0.5	2.1±1.1	2.4±0.6	3.2 ±1.1	2.3±0.4
BOLFI [3]	0.9±0.7	1.4±1.0	0.5±0.3	0.9±0.4	0.7±0.6	1.1±0.9	1.0±0.6	1.1±0.9
ROMC [27]	1.8±0.6	2.3±0.3	1.8±0.3	0.5±0.3	0.5±0.3	1.9±1.1	1.3±0.5	1.3±0.6
AVO (GAUSSIAN) [1]	0.7±0.6	0.8±0.5	0.1±0.5	0.1±0.5	0.3±0.4	0.7±0.5	0.3±0.4	0.6±0.4
AVO (IMPLICIT) [1]	1.7±0.8	1.6±0.8	1.1±0.5	0.4±0.2	1.2±0.4	1.5±0.5	1.0±0.1	0.9±0.3
ALFI-BETA	4.9±1.1	3.9±1.5	2.8 ±1.0	2.2 ±0.8	2.4 ±0.5	3.0±0.5	2.7±0.7	2.5±1.4
ALFI-GAUSSIAN	3.7±0.5	3.1±0.9	2.4±0.4	1.9±0.6	2.3±0.5	3.3 ±0.9	2.6±0.4	2.0±0.8

Table 2: Computational complexity and the wall clock time of *likelihood-free inference* algorithms, see Appendix F for the further discussion. Rejection ABC, MCMC ABC, SMC ABC and BOLFI present (simulation time/sampling time) of the Poisson simulation model for the wall clock time. Other algorithms present (simulation time/sampling time/optimization time).

	REJECTION ABC	MCMC ABC	SMC ABC	BOLFI	ROMC	AVO	ALFI
SAMPLING COMPLEXITY	$O(d)$	$O(d)$	$O(dn^2)$	$O(dt^2)$ [28, 29]	$O(dM^2)$	$O(P)$	$O(P)$
OPTIMIZATION COMPLEXITY	—	—	—	—	$O(dLM^2)$	$O(P+Q)$	$O(P+Q)$
WALL CLOCK TIME (WEALTH)*	29h/2s	3h/15s	3h/1h	24m/14h	25h/14h/18h	20h/48s/9m	3h/15s/78s

* h: hours, m: minutes, s: seconds

Remark. See Appendix C for the general version of Theorem 2 and the proofs. The distribution $\delta_{\theta_0} P_1^m \dots P_t^m$ is the distribution of the inhomogeneous Markov chain after t iterations of learning, starting at θ_0 . The uniform convergence of Theorem 2 states that the convergence speed of the Markov chain to the posterior distribution is uniform, i.e. the mixing time of the Markov chain to the posterior distribution $p(\theta|x_{obs})$ is uniform with respect to the initial point θ_0 .

5 Experiments

5.1 Simulations as Black-box Generative Models

Simulation models with in-depth domain knowledge are the examples of the black-box generative models. Table 1 presents the performance of *likelihood-free inference* algorithms on eight simulation models (see Appendix E), where the performance is the negative log Euclidean distance, $\mathbb{E}_{\theta^*} [-\log(\|\theta^* - \hat{\theta}\|_2)]$, between the true parameter θ^* and the estimated posterior mode $\hat{\theta}$. The observation $x_{obs} = \frac{1}{100} \sum_{j=1}^{100} g(\theta^*, u_j)$ is the average of 100 simulation executions. The algorithms are replicated for 10 times with a different set of true parameters $\{\theta_k^*\}_{k=1}^{10}$ to calculate the performance statistics.

Table 2 presents the computational complexities; and the wall clock time of 1) simulation, 2) sampling and 3) optimization procedures for *likelihood-free inference* algorithms. The Bayesian optimization-based algorithms, such as BOLFI and ROMC, take most of the computation time at the sampling procedure, whereas the sampling and the optimization time in ALFI are ignorable compared to the simulation time. Although ALFI allows parallelization by a cheap sampler, the sampling from parallel Bayesian optimization [30] in BOLFI and ROMC is prohibitive due to the heavy computations in practice. Consequently, we take 100 simultaneous simulation executions in an iteration of ALFI, yet BOLFI takes a single simulation evaluation for an iteration. This property is important if a simulation is expensive.

The second column of Figure 4 illustrates that the estimated beta distribution converges to the nonparametrically estimated density of $\{d_\phi(g(\theta^*, u_j)) | j = 1, \dots, 100\}$. The last column illustrates the contour map of the likelihood estimation, i.e. $Beta(d_\phi(x_{obs}); s_\psi(\theta))$, which concentrates to the true parameter θ^* after iterations. Figure 5 compares the Poisson simulation results among ALFI and implicit AVO; and Figure 5 concludes that ALFI is more identical to the observation than AVO.

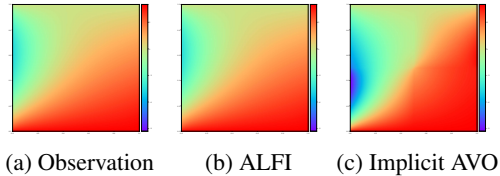


Figure 5: Simulation results

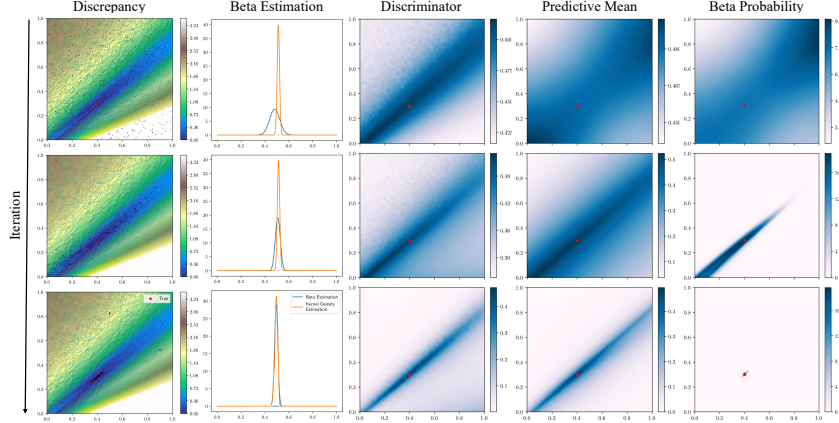


Figure 4: SIR model ALFI result (1st column) The Euclidean discrepancy landscape. The Metropolis-Hastings samples plotted as black dots (2nd column) The estimated beta distribution converges to the nonparametric kernel density estimation of $\{d_\phi(\theta^*, u_j) | j = 1, \dots, 100\}$ (3rd column) The discriminator vanish except $\{\theta | \|x_{obs} - g(\theta, u)\|_2 < \epsilon\}$ for $\epsilon \ll 1$ (4th column) The mean of estimated beta distribution (5th column) The likelihood estimation concentrates to θ^*

5.2 Pre-trained Statistical Models as Black-box Generative Models

5.2.1 Estimation of Spectral Density Mixture

We estimate the mixture parameters of the spectral density [31] modeled by the mixture of Gaussian. The badly selected initial parameters make the optimization stuck at a local optimum with gradient descent. Considering ALFI as a gradient-free optimization algorithm, Figure 6 compares the interpolation/extrapolation results between the inference by ALFI and the gradient learning of the Adam optimizer. Figure 6 illustrates that ALFI optimizes the spectral mixture parameters better than Adam optimizer.

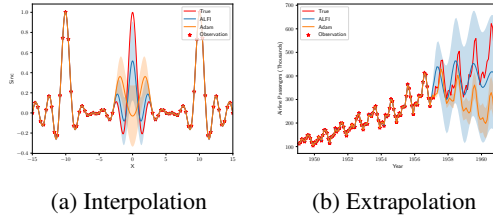


Figure 6: ALFI interpolates/extrapolates better than Adam optimizer after 100 iterations of ALFI learning.

5.2.2 Corrupted Image Inpainting

This experiment assumes a pre-trained DCGAN [32] generator (ω is pre-trained) as the black-box model, and this experiment finds the best regeneration of the observation x_{obs} by estimating the nearest latent embedding of x_{obs} . The observation x_{obs} is a masked MNIST image. Figure 7 concludes that ALFI can generate an image more similar to the true image than implicit AVO can generate.

6 Conclusions

The contribution of this paper is four-fold. First, this paper analyzes the *gradient vanishing problem* and the *implicit relation problem* of the previous research on *likelihood-free inference*. Second, this paper provides a formula of the intractable likelihood as a 1-dimensional density of a random variable in Theorem 1. Third, this paper suggests a new *likelihood-free inference* that uses the adversarial framework. Fourth, this paper suggests Theorem 2 that proves the convergence of the distributions of Markov chain to the posterior distribution, where the transition kernel of inhomogeneous Markov chain in the Metropolis-Hastings algorithm is dynamically updated.

Broader Impact

We believe that ALFI is particularly useful in calibrating a simulation model to be realistic, and such simulations provide foundations for policymaking by what-if simulations. In this aspect, ALFI could aid policymakers in the decision-making process, by optimizing simulation models more congruent to the real-world system of interest. Additionally, ALFI allows us to infer unknown quantities in diverse domains, such as the diffusivity of a porous material; and this ability provides the profound efficiency in scientific simulations. However, we emphasize that ALFI is just an approximation of the posterior distribution, so the calibration using ALFI should not be fully trusted. Once ALFI fails to calibrate a simulation model and a practitioner depends exclusively on the simulation model in making policy, one might suggest a policy that could lead to catastrophic results.

References

- [1] Gilles Louppe, Joeri Hermans, and Kyle Cranmer. Adversarial variational optimization of non-differentiable simulators. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1438–1447, 2019.
- [2] Ted Meeds and Max Welling. Optimization monte carlo: Efficient and embarrassingly parallel likelihood-free inference. In *Advances in Neural Information Processing Systems*, pages 2080–2088, 2015.
- [3] Michael U Gutmann and Jukka Corander. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *The Journal of Machine Learning Research*, 17(1): 4256–4302, 2016.
- [4] Mark A Beaumont. Approximate bayesian computation in evolution and ecology. *Annual review of ecology, evolution, and systematics*, 41:379–406, 2010.
- [5] George Casella and Roger L Berger. *Statistical inference*, volume 2. Duxbury Pacific Grove, CA, 2002.
- [6] Andrei Nikolaevitch Kolmogorov. Determination of the centre of dispersion and degree of accuracy for a limited number of observation. *Izv. Akad. Nauk, USSR Ser. Mat.*, 6:3–32, 1942.
- [7] Simon Tavaré, David J Balding, Robert C Griffiths, and Peter Donnelly. Inferring coalescence times from dna sequence data. *Genetics*, 145(2):505–518, 1997.
- [8] Yun-Xin Fu and Wen-Hsiung Li. Estimating the age of the common ancestor of a sample of dna sequences. *Molecular biology and evolution*, 14(2):195–199, 1997.
- [9] Paul Marjoram, John Molitor, Vincent Plagnol, and Simon Tavaré. Markov chain monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328, 2003.
- [10] Scott A Sisson, Yanan Fan, and Mark M Tanaka. Sequential monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.
- [11] Mark A Beaumont, Jean-Marie Cornuet, Jean-Michel Marin, and Christian P Robert. Adaptive approximate bayesian computation. *Biometrika*, 96(4):983–990, 2009.
- [12] Elias M Stein and Rami Shakarchi. *Functional Analysis: Introduction to Further Topics in Analysis*, volume 4. Princeton University Press, 2011.
- [13] Odo Diekmann and Johan Andre Peter Heesterbeek. *Mathematical epidemiology of infectious diseases: model building, analysis and interpretation*, volume 5. John Wiley & Sons, 2000.
- [14] Martin Arjovsky and Léon Bottou. Towards principled methods for training generative adversarial networks, 2017.
- [15] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

- [16] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.
- [17] David A Levin and Yuval Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [18] Dani Gamerman and Hedibert F Lopes. *Markov chain Monte Carlo: stochastic simulation for Bayesian inference*. CRC Press, 2006.
- [19] Kirill Neklyudov, Evgenii Egorov, and Dmitry P Vetrov. The implicit metropolis-hastings algorithm. In *Advances in Neural Information Processing Systems*, pages 13932–13942, 2019.
- [20] Pranav Unni and Padmanabhan Seshaiyer. Mathematical modeling, analysis, and simulation of tumor dynamics with drug interventions. *Computational and mathematical methods in medicine*, 2019, 2019.
- [21] Lawrence C Evans. *Partial differential equations*, volume 19. American Mathematical Soc., 2010.
- [22] Roger Temam. *Navier-Stokes equations: theory and numerical analysis*, volume 343. American Mathematical Soc., 2001.
- [23] John D Sterman. System dynamics modeling: tools for learning in a complex world. *California management review*, 43(4):8–25, 2001.
- [24] Ngai Hang Chan. Autoregressive moving average models. *Time Series: Applications to Finance with R and S-Plus, Second Edition*, John Wiley & Sons, Inc., Hoboken, NJ, 2010.
- [25] C Newell. *Applications of Queueing Theory*, volume 4. Springer Science & Business Media, 2013.
- [26] Uri Wilensky. Netlogo wealth distribution model. *Center for Connected Learning and Computer-Based Modeling, Northwestern University, Evanston, IL.*[En línea] Disponible en: <http://ccl.northwestern.edu/NetLogo/models/WealthDistribution>, 1998.
- [27] Borislav Ikonov and Michael U Gutmann. Robust optimisation monte carlo. *arXiv preprint arXiv:1904.00670*, 2019.
- [28] Ching-An Cheng and Byron Boots. Variational inference for gaussian process models with linear complexity. In *Advances in Neural Information Processing Systems*, pages 5184–5194, 2017.
- [29] Dewi Retno Sari Saputro and Purnami Widyaningsih. Limited memory broyden-fletcher-goldfarb-shanno (l-bfgs) method for the parameter estimation on geographically weighted ordinal logistic regression model (gwolr). In *AIP Conference Proceedings*, volume 1868, page 040009. AIP Publishing LLC, 2017.
- [30] Jian Wu and Peter Frazier. The parallel knowledge gradient method for batch bayesian optimization. In *Advances in Neural Information Processing Systems*, pages 3126–3134, 2016.
- [31] Andrew Wilson and Ryan Adams. Gaussian process kernels for pattern discovery and extrapolation. In *International conference on machine learning*, pages 1067–1075, 2013.
- [32] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.