

---

# Adversarial Variational Optimization of Non-Differentiable Simulators

---

Gilles Louppe<sup>1</sup> Kyle Cranmer<sup>2</sup>

## Abstract

Complex computer simulators are increasingly used across fields of science as generative models tying parameters of an underlying theory to experimental observations. Inference in this setup is often difficult, as simulators rarely admit a tractable density or likelihood function. We introduce Adversarial Variational Optimization (AVO), a likelihood-free inference algorithm for fitting a non-differentiable generative model incorporating ideas from generative adversarial networks (GANs) and variational inference. We adapt the training procedure of (Wasserstein) GANs by replacing the differentiable generative network with a domain-specific simulator. We solve the resulting non-differentiable minimax problem by minimizing variational upper bounds of the two adversarial objectives. Effectively, the procedure results in learning a proposal distribution over simulator parameters, such that the corresponding marginal distribution of the generated data matches the observations. We present results of the method with simulators producing both discrete and continuous data.

## 1. Introduction

In many fields of science such as particle physics, epidemiology, and population genetics, computer simulators are used to describe complex data generation processes. These simulators relate observations  $\mathbf{x}$  to the parameters  $\theta$  of an underlying theory or mechanistic model. In most cases, these simulators are specified as procedural implementations of forward, stochastic processes involving latent variables  $\mathbf{z}$ . Rarely do these simulators admit a

tractable density (or likelihood)  $p(\mathbf{x}|\theta)$ . The prevalence and significance of this problem has motivated an active research effort in so-called *likelihood-free inference* algorithms such as Approximate Bayesian Computation (ABC) and density estimation-by-comparison algorithms (Beaumont et al., 2002; Marjoram et al., 2003; Sisson et al., 2007; Sisson & Fan, 2011; Marin et al., 2012; Cranmer et al., 2015).

In parallel, with the introduction of variational auto-encoders (Kingma & Welling, 2013) and generative adversarial networks (Goodfellow et al., 2014), there has been a vibrant research program around implicit generative models based on neural networks (Mohamed & Lakshminarayanan, 2016). While some of these models also do not admit a tractable density, they are all differentiable by construction. In addition, generative models based on neural networks are highly parametrized and the model parameters have no obvious interpretation. In contrast, scientific simulators can be thought of as highly regularized generative models as they typically have relatively few parameters and they are endowed with some level of interpretation. In this setting, inference on the model parameters  $\theta$  is often of more interest than the latent variables  $\mathbf{z}$ .

In this note, we develop a likelihood-free inference algorithm for non-differentiable, implicit generative models. We adapt the adversarial training procedure of generative adversarial networks (Goodfellow et al., 2014) by replacing the implicit generative network with a domain-based scientific simulator, and solve the resulting non-differentiable minimax problem by minimizing variational upper bounds (Wierstra et al., 2011; Staines & Barber, 2012) of the adversarial objectives. The objective of the algorithm is to match the marginal distribution of the generated data to the empirical distribution of the observations.

## 2. Problem statement

We consider a family of parametrized densities  $p(\mathbf{x}|\theta)$  defined implicitly through the simulation of a stochastic generative process, where  $\mathbf{x} \in \mathbb{R}^d$  is the data and  $\theta$  are the parameters of interest. The simulation may involve some complicated latent process where  $\mathbf{z} \in \mathcal{Z}$  is a latent variable providing an external source of randomness. Unlike

---

<sup>\*</sup>Equal contribution <sup>1</sup>University of Liège, Belgium <sup>2</sup>New York University, USA. Correspondence to: Gilles Louppe <g.louppe@uliege.be>, Kyle Cranmer <kyle.cranmer@nyu.edu>.

implicit generative models defined by neural networks, we do not assume  $\mathbf{z}$  to be a fixed-size vector with a simple density. Instead, the dimension of  $\mathbf{z}$  and the nature of its components (uniform, normal, discrete, continuous, etc.) are inherited from the control flow of the simulation code and may depend on  $\theta$  in some intricate way. Moreover, the dimension of  $\mathbf{z}$  may be much larger than the dimension of  $\mathbf{x}$ .

We assume that the stochastic generative process that defines  $p(\mathbf{x}|\theta)$  is specified through a non-differentiable deterministic function  $g(\cdot; \theta) : \mathcal{Z} \rightarrow \mathbb{R}^d$ . Operationally,

$$\mathbf{x} \sim p(\mathbf{x}|\theta) \triangleq \mathbf{z} \sim p(\mathbf{z}|\theta), \mathbf{x} = g(\mathbf{z}; \theta) \quad (1)$$

such that the density  $p(\mathbf{x}|\theta)$  can be written as

$$p(\mathbf{x}|\theta) = \int_{\{\mathbf{z}: g(\mathbf{z}; \theta) = \mathbf{x}\}} p(\mathbf{z}|\theta) \mu(d\mathbf{z}), \quad (2)$$

where  $\mu$  is a probability measure.

Given some observed data  $\{\mathbf{x}_i | i = 1, \dots, N\}$  drawn from the (unknown) true distribution  $p_r(\mathbf{x})$ , our goal is to estimate the parameters  $\theta^*$  that minimize the divergence or the distance  $\rho$  between  $p_r(\mathbf{x})$  and the implicit model  $p(\mathbf{x}|\theta)$ . That is,

$$\theta^* = \arg \min_{\theta} \rho(p_r(\mathbf{x}), p(\mathbf{x}|\theta)). \quad (3)$$

### 3. Background

#### 3.1. Generative adversarial networks

Generative adversarial networks (GANs) were first proposed by (Goodfellow et al., 2014) as a way to build an implicit generative model capable of producing samples from random noise  $\mathbf{z}$ . More specifically, a generative model  $g(\cdot; \theta)$  is pit against an adversarial classifier  $d(\cdot; \phi) : \mathbb{R}^d \rightarrow [0, 1]$  with parameters  $\phi$  and whose antagonistic objective is to recognize real data  $\mathbf{x}$  from generated data  $\tilde{\mathbf{x}} = g(\mathbf{z}; \theta)$ . Both models  $g$  and  $d$  are trained simultaneously, in such a way that  $g$  learns to fool its adversary  $d$  (which happens when  $g$  produces samples comparable to the observed data), while  $d$  continuously adapts to changes in  $g$ . When  $d$  is trained to optimality before each parameter update of the generator, it can be shown that the original adversarial learning procedure (Goodfellow et al., 2014) amounts to minimizing the Jensen-Shannon divergence  $\text{JSD}(p_r(\mathbf{x}) \parallel p(\mathbf{x}|\theta))$  between  $p_r(\mathbf{x})$  and  $p(\mathbf{x}|\theta)$ .

As thoroughly explored in (Arjovsky & Bottou, 2017), GANs remain remarkably difficult to train because of vanishing gradients as  $d$  saturates, or because of unreliable updates when the training procedure is relaxed. As a remedy, Wasserstein GANs (Arjovsky et al., 2017) (WGANs) reformulate the adversarial setup in order to minimize the Wasserstein-1 distance  $W(p_r(\mathbf{x}), p(\mathbf{x}|\theta))$  by replacing the

adversarial classifier with a 1-Lipschitz adversarial critic  $d(\cdot; \phi) : \mathbb{R}^d \rightarrow \mathbb{R}$ , such that

$$\begin{aligned} W(p_r(\mathbf{x}), p(\mathbf{x}|\theta)) &= \sup_{\phi} \mathbb{E}_{\tilde{\mathbf{x}} \sim p(\mathbf{x}|\theta)} [d(\tilde{\mathbf{x}}; \phi)] - \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [d(\mathbf{x}; \phi)] \\ &\triangleq \sup_{\phi} \mathcal{L}_W. \end{aligned} \quad (4)$$

Under the WGAN-GP formulation of (Gulrajani et al., 2017) for stabilizing the optimization procedure, training  $d$  and  $g$  results in alternating gradient updates on  $\phi$  and  $\theta$  in order to respectively minimize

$$\mathcal{L}_d = \mathcal{L}_W + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim p(\tilde{\mathbf{x}})} [(\|\nabla_{\tilde{\mathbf{x}}} d(\tilde{\mathbf{x}}; \phi)\|_2 - 1)^2] \quad (5)$$

$$\mathcal{L}_g = -\mathcal{L}_W \quad (6)$$

where  $\tilde{\mathbf{x}} := \epsilon \mathbf{x} + (1 - \epsilon) \tilde{\mathbf{x}}$ , for  $\epsilon \sim U[0, 1]$ ,  $\mathbf{x} \sim p_r(\mathbf{x})$  and  $\tilde{\mathbf{x}} \sim p(\mathbf{x}|\theta)$ .

#### 3.2. Variational optimization

Variational optimization (Staines & Barber, 2012; Staines & Barber, 2013) and evolution strategies (Wierstra et al., 2011) are general optimization techniques that can be used to form a differentiable bound on the optima of a non-differentiable function. Given a function  $f$  to minimize, these techniques are based on the observation that

$$\min_{\theta} f(\theta) \leq \mathbb{E}_{\theta \sim q(\theta|\psi)} [f(\theta)] = U(\psi), \quad (7)$$

where  $q(\theta|\psi)$  is a proposal distribution with parameters  $\psi$  over input values  $\theta$ . That is, the minimum of a set of function values is always less than or equal to any of their average. Provided that the proposal distribution is flexible enough, the parameters  $\psi$  can be updated to place its mass arbitrarily tight around the optimum  $\theta^* = \min_{\theta \in \Theta} f(\theta)$ .

Under mild restrictions outlined in (Staines & Barber, 2012), the bound  $U(\psi)$  is differentiable with respect to  $\psi$ , and using the log-likelihood trick its gradient can be rewritten as:

$$\begin{aligned} \nabla_{\psi} U(\psi) &= \nabla_{\psi} \mathbb{E}_{\theta \sim q(\theta|\psi)} [f(\theta)] \\ &= \mathbb{E}_{\theta \sim q(\theta|\psi)} [f(\theta) \nabla_{\psi} \log q(\theta|\psi)] \end{aligned} \quad (8)$$

Effectively, this means that provided that the score function  $\nabla_{\psi} \log q(\theta|\psi)$  of the proposal is known and that one can evaluate  $f(\theta)$  for any  $\theta$ , then one can construct empirical estimates of Eqn. 8, which can in turn be used to minimize  $U(\psi)$  with stochastic gradient descent (or a variant thereof, robust to noise and parameter scaling).

### 4. Adversarial variational optimization

#### 4.1. Algorithm

The alternating stochastic gradient descent on  $\mathcal{L}_d$  and  $\mathcal{L}_g$  in WGANs (Section 3.1) inherently assumes that the gen-

**Algorithm 1** Adversarial variational optimization (AVO).

*Inputs:* observed data  $\{\mathbf{x}_i \sim p_r(\mathbf{x})\}_{i=1}^N$ , simulator  $g$ .

*Outputs:* proposal distribution  $q(\boldsymbol{\theta}|\boldsymbol{\psi})$ , such that  $q(\mathbf{x}|\boldsymbol{\psi}) \approx p_r(\mathbf{x})$ .

*Hyper-parameters:* The number  $n_{\text{critic}}$  of training iterations of  $d$ ; the size  $M$  of a mini-batch; the gradient penalty coefficient  $\lambda$ ; the entropy penalty coefficient  $\gamma$ .

```

1:  $q(\boldsymbol{\theta}|\boldsymbol{\psi}) \leftarrow$  prior on  $\boldsymbol{\theta}$  (with differentiable and known density)
2: while  $\boldsymbol{\psi}$  has not converged do
3:   for  $i = 1$  to  $n_{\text{critic}}$  do ▷ Update  $d$ 
4:     Sample a mini-batch  $\{\mathbf{x}_m \sim p_r(\mathbf{x}), \boldsymbol{\theta}_m \sim q(\boldsymbol{\theta}|\boldsymbol{\psi}), \mathbf{z}_m \sim p(\mathbf{z}|\boldsymbol{\theta}_m), \epsilon_m \sim U[0, 1]\}_{m=1}^M$ .
5:     for  $m = 1$  to  $M$  do
6:        $\tilde{\mathbf{x}}_m \leftarrow g(\mathbf{z}_m; \boldsymbol{\theta}_m)$ 
7:        $\hat{\mathbf{x}}_m \leftarrow \epsilon_m \mathbf{x}_m + (1 - \epsilon_m) \tilde{\mathbf{x}}_m$ 
8:        $U_d^{(m)} \leftarrow d(\tilde{\mathbf{x}}_m; \boldsymbol{\phi}) - d(\mathbf{x}_m; \boldsymbol{\phi}) + \lambda(\|\nabla_{\tilde{\mathbf{x}}_m} d(\hat{\mathbf{x}}_m; \boldsymbol{\phi})\|_2 - 1)^2$ 
9:     end for
10:     $\boldsymbol{\phi} \leftarrow \text{Adam}(\nabla_{\boldsymbol{\phi}} \frac{1}{M} \sum_{m=1}^M U_d^{(m)})$ 
11:  end for
12:  Sample a mini-batch  $\{\boldsymbol{\theta}_m \sim q(\boldsymbol{\theta}|\boldsymbol{\psi}), \mathbf{z}_m \sim p(\mathbf{z}|\boldsymbol{\theta}_m)\}_{m=1}^M$ . ▷ Update  $q(\boldsymbol{\theta}|\boldsymbol{\psi})$ 
13:   $\nabla_{\boldsymbol{\psi}} U_g \leftarrow \frac{1}{M} \sum_{m=1}^M -d(g(\mathbf{z}_m; \boldsymbol{\theta}_m)) \nabla_{\boldsymbol{\psi}} \log q_{\boldsymbol{\psi}}(\boldsymbol{\theta}_m)$ 
14:   $\nabla_{\boldsymbol{\psi}} H(q_{\boldsymbol{\psi}}) \leftarrow \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\psi}} q_{\boldsymbol{\psi}}(\boldsymbol{\theta}_m) \log q_{\boldsymbol{\psi}}(\boldsymbol{\theta}_m)$ 
15:   $\boldsymbol{\psi} \leftarrow \text{Adam}(\nabla_{\boldsymbol{\psi}} U_g + \gamma \nabla_{\boldsymbol{\psi}} H(q_{\boldsymbol{\psi}}))$ 
16: end while
    
```

erator  $g$  is a differentiable function. In the setting where we are interested in estimating the parameters of a fixed non-differentiable simulator (Section 2), rather than in learning the generative model itself, gradients  $\nabla_{\boldsymbol{\theta}} g$  either do not exist or are inaccessible. As a result, gradients  $\nabla_{\boldsymbol{\theta}} \mathcal{L}_g$  cannot be constructed and the optimization procedure cannot be carried out.

In this work, we propose to rely on variational optimization to minimize  $\mathcal{L}_d$  and  $\mathcal{L}_g$ , thereby bypassing the non-differentiability of  $g$ . More specifically, we consider a proposal distribution  $q(\boldsymbol{\theta}|\boldsymbol{\psi})$  over the parameters of the simulator  $g$  and alternately minimize the variational upper bounds

$$U_d = \mathbb{E}_{\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\boldsymbol{\psi})}[\mathcal{L}_d] \quad (9)$$

$$U_g = \mathbb{E}_{\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\boldsymbol{\psi})}[\mathcal{L}_g] \quad (10)$$

respectively over  $\boldsymbol{\phi}$  and  $\boldsymbol{\psi}$ . When updating  $\boldsymbol{\phi}$ , unbiased estimates of  $\nabla_{\boldsymbol{\phi}} U_d$  can be obtained by directly evaluating the gradient of  $U_d$  over mini-batches of real and generated data. When updating  $\boldsymbol{\psi}$ ,  $\nabla_{\boldsymbol{\psi}} U_g$  can be estimated as described in the previous section. That is,

$$\nabla_{\boldsymbol{\psi}} U_g = \mathbb{E}_{\substack{\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\boldsymbol{\psi}), \\ \tilde{\mathbf{x}} \sim p(\mathbf{x}|\boldsymbol{\theta})}}[-d(\tilde{\mathbf{x}}; \boldsymbol{\phi}) \nabla_{\boldsymbol{\psi}} \log q(\boldsymbol{\theta}|\boldsymbol{\psi})], \quad (11)$$

which we can approximate with mini-batches of generated

data

$$\nabla_{\boldsymbol{\psi}} U_g \approx \frac{1}{M} \sum_{m=1}^M -d(g(\mathbf{z}_m; \boldsymbol{\theta}_m); \boldsymbol{\phi}) \nabla_{\boldsymbol{\psi}} \log q(\boldsymbol{\theta}_m|\boldsymbol{\psi}) \quad (12)$$

for  $\boldsymbol{\theta}_m \sim q(\boldsymbol{\theta}|\boldsymbol{\psi})$  and  $\mathbf{z}_m \sim p(\mathbf{z}|\boldsymbol{\theta}_m)$ . For completeness, Algorithm 1 outlines the proposed adversarial variational optimization (AVO) procedure, as built on top of WGAN-GP.

Algorithm 1 represents the simplest version of AVO; however, the variance of the noisy estimator of the gradients may be too large to be useful in many problems. We use Adam (Kingma & Ba, 2014), but note the opportunity to use instead the Natural Evolution Strategy algorithm (Wierstra et al., 2011) or variance reduction techniques as control variates (Grathwohl et al., 2017) or finite differences (Buesing et al., 2016).

## 4.2. Parameter Point Estimates

The variational objectives 9-10 have the effect of replacing the modeled data distribution of Eqn. 1 with the parametrized marginal distribution of the generated data

$$q(\mathbf{x}|\boldsymbol{\psi}) = \int p(\mathbf{x}|\boldsymbol{\theta}) q(\boldsymbol{\theta}|\boldsymbol{\psi}) d\boldsymbol{\theta}. \quad (13)$$

We can think of  $q(\mathbf{x}|\boldsymbol{\psi})$  as a *variational program* as described in (Ranganath et al., 2016), though more complicated than a simple reparametrization of normally dis-

tributed noise  $\mathbf{z}$  through a differentiable function. In our case, the variational program is a marginalized, non-differentiable simulator. Its density is intractable; nevertheless, it can generate samples for  $\mathbf{x}$  whose expectations are differentiable with respect to  $\psi$ . Operationally, we sample from this marginal model via

$$\mathbf{x} \sim q(\mathbf{x}|\psi) \triangleq \boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\psi), \mathbf{z} \sim p(\mathbf{z}|\boldsymbol{\theta}), \mathbf{x} = g(\mathbf{z}; \boldsymbol{\theta}). \quad (14)$$

We can view the optimization of  $q(\mathbf{x}|\psi)$  with respect to  $\psi$  through the lens of empirical Bayes, where the data are used to optimize a prior within the family  $q(\boldsymbol{\theta}|\psi)$ . If  $\rho$  is the KL distance,  $\psi^*$  would correspond to the maximum marginal likelihood estimator advocated by Rubin (Rubin, 1984). When  $\rho$  is the Wasserstein distance,  $\psi^*$  is referred to as the minimum Wasserstein estimator (MWE). When the model is well specified, the MWE coincides with the true data-generating parameter; however, if the model is misspecified, the MWE is typically different from the maximum likelihood estimator (MLE). Thus, if the simulator  $p(\mathbf{x}|\boldsymbol{\theta})$  is misspecified,  $q(\boldsymbol{\theta}|\psi)$  will attempt to smear it so that the marginal model  $q(\mathbf{x}|\psi)$  is closer to  $p_r(\mathbf{x})$ . However, if the simulator is well specified, then  $q(\boldsymbol{\theta}|\psi)$  will concentrate its mass around the true data-generating parameter.

In order to more effectively target point estimates  $\boldsymbol{\theta}^*$ , we augment Eqn. 10 with an entropic regularization term  $H(q(\boldsymbol{\theta}|\psi))$ , that is

$$U_g = \mathbb{E}_{\boldsymbol{\theta} \sim q(\boldsymbol{\theta}|\psi)}[\mathcal{L}_g] + \gamma H(q(\boldsymbol{\theta}|\psi)) \quad (15)$$

where  $\gamma \in \mathbb{R}^+$  is a hyper-parameter controlling the trade-off between the generator objective and the tightness of the proposal distribution. For small values of  $\gamma$ , proposal distributions with large entropy are not penalized, which results in learning a smeared variation of the original simulator. On the other hand, for large values of  $\gamma$ , the procedure is encouraged to fit a proposal distribution with low entropy, which has the effect of concentrating its density tightly around one or a few  $\boldsymbol{\theta}$  values.

Finally, we note that very large penalties may eventually make the optimization unstable, as the variance of  $\nabla_{\psi} \log q(\boldsymbol{\theta}_m|\psi)$  typically increases as the entropy of the proposal decreases.

## 5. Experiments

### 5.1. Univariate discrete data

As a first illustrative experiment, we evaluate inference for a discrete Poisson distribution with unknown mean  $\lambda$ . We artificially consider the distribution as a parametrized simulator, from which we can only generate data.

The observed data is sampled from a Poisson with mean

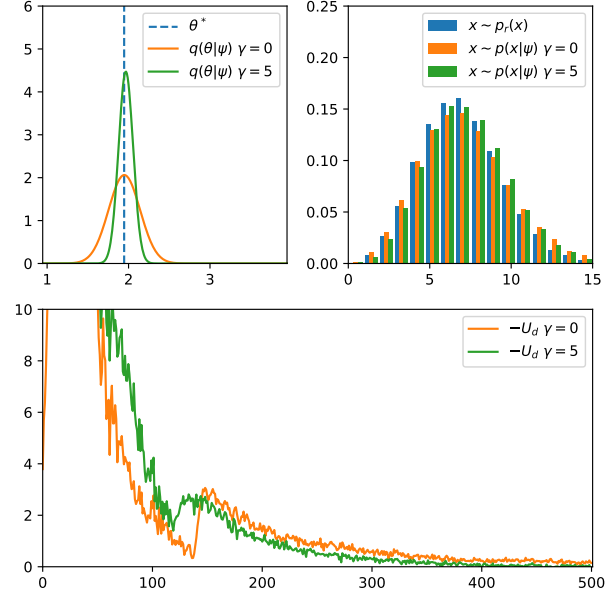


Figure 1. Discrete Poisson model with unknown mean. (Top left) Proposal distributions  $q(\boldsymbol{\theta}|\psi)$  after training. For both  $\gamma = 0$  and  $\gamma = 5$ , the distributions correctly concentrate their density around the true value  $\log(\lambda^*)$ . Entropic regularization ( $\gamma = 5$ ) results in a tighter density. (Top right) Model distributions  $q(\mathbf{x}|\psi)$  after training. This plot shows that the resulting parametrizations of the simulator closely reproduce the true distribution. (Bottom) Empirical estimates of the variational upper bound  $U_d$  as optimization progresses.

$\lambda^* = 7$ . Algorithm 1 is run for 500 epochs with mini-batches of size  $M = 64$  and the following configuration. For the critic  $d$ , we use a 3-layer MLP with 10 hidden nodes per layer and ReLU activations. At each epoch, Adam is run for  $n_{\text{critic}} = 10$  iterations with a step size  $\alpha = 0.005$  and  $\beta_1 = \beta_2 = 0.05$ . For estimating  $\lambda^*$ , we parametrize  $\boldsymbol{\theta}$  as  $\log(\lambda)$  and use a univariate Gaussian proposal distribution  $q(\boldsymbol{\theta}|\psi)$  initialized with a mean at  $\log(10)$  and unit variance. At each epoch, parameters  $\psi$  are updated by taking one Adam step and the same hyper-parameters as for the critic. The gradient penalty coefficient is set to  $\lambda = 0.1$ , and the entropy penalty is evaluated at both  $\gamma = 0$  and  $\gamma = 5$ .

The top left plot in Figure 1 illustrates the resulting proposal distributions  $q(\boldsymbol{\theta}|\psi)$  after AVO. For both  $\gamma = 0$  and  $\gamma = 5$ , the proposal distributions correctly concentrate their density around the true parameter value  $\log(\lambda^*) = 1.94$ . Under the effect of entropic regularization, the proposal distribution for  $\gamma = 5$  concentrates its mass tightly, yielding in this case precise inference. The top right plot compares the model distributions to the true distribution. As theoretically expected from adversarial training, we see that the resulting distributions align with the true distribution, with in this case visually slightly better results for the

penalized model. The bottom plot of Figure 1 shows empirical estimates of  $-U_d$  with respect to the epoch number. For both  $\gamma = 0$  and  $\gamma = 5$ , the curves fall towards 0, which indicates that  $\mathbb{E}_{\tilde{\mathbf{x}} \sim p(\mathbf{x}|\theta)}[d(\tilde{\mathbf{x}}; \phi)] \approx \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})}[d(\mathbf{x}; \phi)]$  and that the critic cannot distinguish between true and model data. This confirms that adversarial variational optimization works despite the discreteness of the data and lack of access to the density  $p(\mathbf{x}|\theta)$  or its gradient.

## 5.2. Multidimensional continuous data

As a second toy example, we consider a generator producing 5-dimensional continuous data, as originally specified in Section 4.2 of (Cranmer et al., 2015). Again, AVO does not have access to the density or its gradient, only samples from the generative model. We consider observed data generated at the nominal values  $\theta^* = (\alpha^* = 1, \beta^* = -1)$ . The simulator parameters are modeled with a factored (mean field) Gaussian proposal distribution  $q(\theta|\psi) = q(\alpha|\psi)q(\beta|\psi)$ , where each component was initialized with zero mean and unit variance. Hyperparameters are set to  $M = 64$ ,  $n_{\text{critic}} = 10$ ,  $\lambda = 0.1$ ,  $\gamma = 1$  and Adam configured with a step size  $\alpha = 0.005$  and  $\beta_1 = \beta_2 = 0.05$ . The architecture of the critic is the same as in the previous example.

Starting with a proposal distribution  $q(\theta|\psi)$  largely spread over the parameter space, as illustrated in the left plot of Figure 2, AVO converges towards a proposal distribution whose density concentrates around the nominal values  $\theta^*$ , as shown in the right plot of Figure 2. Overall, this example further illustrates and confirms the ability of adversarial variational optimization for inference with multiple parameters and multidimensional data, where reliable approximations of  $p(\mathbf{x}|\theta)$  in a traditional MLE setting would otherwise be difficult to construct.

## 5.3. Benchmarks

[GL: Compare against ABC and others.] [GL: Sections 6.1 and 6.4 of Bernton.]

## 6. Related work

This work sits at the intersection of several lines of research related to likelihood-free inference, approximate Bayesian computation (ABC), implicit generative models, and variational inference. Viewed from the literature around implicit generative models based on neural networks, the proposed method can be considered as a direct adaptation of generative adversarial networks (Goodfellow et al., 2014; Arjovsky & Bottou, 2017) to non-differentiable simulators using variational optimization (Staines & Barber, 2012).

From the point of view of likelihood-free inference, where non-differentiable simulators are the norm, our contribu-

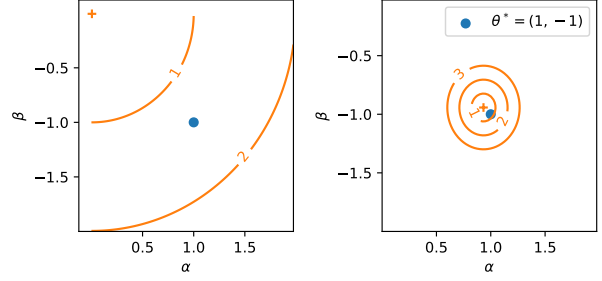


Figure 2. Multidimensional continuous data. (Left) Density  $q(\theta|\psi)$  at the beginning of the procedure, for a proposal distribution initialized with zero mean and unit variance. Contours correspond to parameters  $\theta$  within 1-2-3 Mahalanobis distance units from the mean of  $q(\theta|\psi)$ . (Right) Density  $q(\theta|\psi)$  after adversarial variational optimization ( $\gamma = 10$ ). The proposal density correctly converges towards a distribution whose density concentrates around  $\theta^* = (1, -1)$ .

tions are threefold. First is the process of lifting the expectation with respect to the non-differentiable simulator  $\mathbb{E}_{\tilde{\mathbf{x}} \sim p(\mathbf{x}|\theta)}$  to a differentiable expectation with respect to the variational program  $\mathbb{E}_{\tilde{\mathbf{x}} \sim q(\mathbf{x}|\psi)}$ . Secondly, is the introduction of a novel form of variational inference that works in a likelihood-free setting. Thirdly, the AVO algorithm can be viewed as a form of empirical Bayes where the prior is optimized based on the data.

Perhaps the closest to our work is (Bernton et al., 2017), which uses the Wasserstein distance both to find point estimates  $\theta^*$  and as a part of a rejection sampler in an ABC-like setup (as opposed to variational inference). They emphasize that this approach allows one to eliminate the summary statistics typically used in ABC and calculate the Wasserstein distance explicitly, without making use of the Kantorovich-Rubinstein duality and of a critic  $d$ . For high-dimensional data, they note that this is computationally expensive and introduce an approximation based on projection of the data onto Hilbert space-filling curves. Their Remark 5.1 points out that (Montavon et al., 2016) proposed an approximation of the gradient of an entropy-regularized Wasserstein distance, which uses a similar duality. They note that “unfortunately, it is not applicable in the setting of purely generative models, as it involves point-wise evaluations of the derivative of the log-likelihood.” Thus, our contribution is to provide gradients of an approximate MWE by taking expectations with the variational program  $q(\mathbf{x}|\psi)$ . This paired with the critic Kantorovich-Rubinstein dual formulation of the Wasserstein distance allows us to work in high dimensions without summary statistics and to optimize  $\psi$  with stochastic gradient descent. Additionally, our procedure avoids the inefficiencies of their Wasserstein rejection sampler incurred from using the ABC-likelihood. In the case of small  $N$ , we note however that fast al-

gorithms for calculating the exact (entropically regularized) Wasserstein distance on empirical distributions (Cuturi, 2013; Genevay et al., 2016; Montavon et al., 2016) are alternatives to adversarial learning that are worth considering.

More generally, likelihood-free inference is intimately tied to a class of algorithms that can be framed as density estimation-by-comparison, as reviewed in (Mohamed & Lakshminarayanan, 2016). In most cases, these inference algorithms are formulated as an iterative two-step process where the model distribution is first compared to the true data distribution and then updated to make it more comparable to the latter. Relevant work in this direction includes those that rely on a classifier to estimate the discrepancy between the observed data and the model distributions (Gutmann & Hyvärinen, 2012; Cranmer et al., 2015; 2016; Dutta et al., 2016; Gutmann et al., 2017; Rosca et al., 2017). Also of direct relevance in the likelihood-free setup, Hamiltonian ABC (Meeds et al., 2015) estimate gradients with respect to  $\theta$  through finite differences from multiple forward passes of the simulator with variance reduction strategies based on controlling the source of randomness used for the latent variable  $\mathbf{z}$ .

Likewise, AVO closely relates to recent extensions of GANs, such as Adversarial Learned Inference (ALI) (Dumoulin et al., 2016), Bidirectional GANs (BiGANs) (Donahue et al., 2016),  $\alpha$ -GAN (Rosca et al., 2017), Adversarial Variational Bayes (AVB) (Mescheder et al., 2017), and the PC-Adv algorithm of (Huszár, 2017), which add an inference network to the generative model. Each of these assume a tractable density  $p(\mathbf{x}|\theta)$  that is differentiable with respect to  $\theta$ , which is not satisfied in the likelihood-free setting. Our lifting of the non-differentiable simulator  $p(\mathbf{x}|\theta)$  to the variational program  $q(\mathbf{x}|\psi)$  provides the ability to differentiate expectations with respect to  $\psi$  as in Eqn 8; however, the density  $q(\mathbf{x}|\psi)$  is still intractable. Moreover, we do not attempt to define a recognition model  $q(\mathbf{z}, \theta|\psi)$  as the latent space  $\mathcal{Z}$  of many real-world simulators is complicated and not amenable to a neural recognition model.

This work has also many connections to work on variational inference, in which the goal is to optimize the recognition model  $q(\mathbf{z}, \theta|\psi)$  so that it is close to the true posterior  $p(\mathbf{z}, \theta|\mathbf{x})$ . There have been efforts to extend variational inference to intractable likelihoods; however, many require restrictive assumptions. In (Tran et al., 2017), the authors consider Variational Bayes with an Intractable Likelihood (VBIL). In that approach “the only requirement is that the intractable likelihood can be estimated unbiasedly.” In the case of simulators, they propose to use the ABC-likelihood with an  $\epsilon$ -kernel. The ABC likelihood is only unbiased as  $\epsilon \rightarrow 0$ , thus this method inherits the drawbacks of the ABC-likelihood including the choice of summary statistics

and the inefficiency in evaluating the ABC likelihood for high-dimensional data and small  $\epsilon$ . More recently, (Tran et al., 2017) adapted variational inference to hierarchical implicit models defined on simulators. In this work, the authors step around the intractable likelihoods by reformulating the optimization of the ELBO in terms of a neural and differentiable approximation  $r$  of the log-likelihood ratio  $\log \frac{p}{q}$ , thereby effectively using the same core principle as used in GANs (Mohamed & Lakshminarayanan, 2016).

Lastly, we make a connection to Operator Variational Inference (OPVI) (Ranganath et al., 2016), which is a generalization of variational inference formulated as the following optimization problem:

$$\lambda^* = \inf_{\lambda} \sup_{\phi} \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\lambda)} [(O^{p,q} f_{\phi})]. \quad (16)$$

In traditional VI with the KL distance, this corresponds to  $(O^{p,q} f)(\mathbf{z}) = \log q(\mathbf{z}) - \log p(\mathbf{z}|\mathbf{x}) \forall f \in \mathcal{F}$ . AVO can be cast into a similar form with expectations over  $\mathbf{x}$  instead of  $\mathbf{z}$  and

$$\begin{aligned} \psi^* &= \inf_{\psi} \sup_{\phi} \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x}|\psi)} [(O^{p_r,q\psi} d_{\phi})] \\ &= \inf_{\psi} \sup_{\phi} \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [d(\mathbf{x}; \phi)] - \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x}|\psi)} [d(\mathbf{x}; \phi)]. \end{aligned} \quad (17)$$

Rewriting Eqn. 17 as above is possible through importance sampling, corresponding to an implicit form of the operator

$$(O^{p_r,q\psi} d_{\phi})(\mathbf{x}) = \left( \frac{p_r(\mathbf{x})}{q(\mathbf{x}|\psi)} - 1 \right) d(\mathbf{x}; \phi), \quad (18)$$

which reinforces the link to density ratio estimation and inference with implicit models.

## 7. Summary

In this note, we develop a likelihood-free inference algorithm for non-differentiable, implicit generative models. The algorithm combines adversarial training with variational optimization to minimize variational upper bounds on the otherwise non-differentiable adversarial objectives. The AVO algorithm enables empirical Bayes through variational inference in the likelihood-free setting. This approach does not incur the inefficiencies of an ABC-like rejection sampler. When the model is well-specified, the AVO algorithm provides point estimates for the generative model, which asymptotically corresponds to the data generating parameters. The algorithm works on continuous or discrete data.

Preliminary results on toy problems with discrete and continuous data validate the proposed method. While the obtained results are encouraging, the complete validation of the method remains to be carried out in real conditions on a full fledged scientific simulator – which we plan to

achieve for a next version of this work [GL: to update]. In terms of method, several components need further investigation. First, we need to better study the interplay between the entropic regularization and the adversarial objectives. Second, we should better understand the dynamics of the optimization procedure, in particular when combined with momentum-based optimizers like Adam. Third, we need to consider whether less noisy estimates of the gradients  $\nabla_{\psi} U_g$  can be computed.

## References

- Agostinelli, S. et al. GEANT4: A Simulation toolkit. *Nucl. Instrum. Meth.*, A506:250–303, 2003. doi: 10.1016/S0168-9002(03)01368-8.
- Alwall, Johan, Herquet, Michel, Maltoni, Fabio, Mattelaer, Olivier, and Stelzer, Tim. MadGraph 5 : Going Beyond. *JHEP*, 06:128, 2011. doi: 10.1007/JHEP06(2011)128.
- Arjovsky, M. and Bottou, L. Towards Principled Methods for Training Generative Adversarial Networks. *ArXiv e-prints*, January 2017.
- Arjovsky, M., Chintala, S., and Bottou, L. Wasserstein GAN. *ArXiv e-prints*, January 2017.
- Beaumont, Mark A, Zhang, Wenyang, and Balding, David J. Approximate bayesian computation in population genetics. *Genetics*, 162(4):2025–2035, 2002.
- Bernton, Espen, Jacob, Pierre E, Gerber, Mathieu, and Robert, Christian P. Inference in generative models using the wasserstein distance. *arXiv preprint arXiv:1701.05146*, 2017.
- Buesing, Lars, Weber, Theophane, and Mohamed, Shakir. Stochastic gradient estimation with finite differences. 2016.
- Cranmer, K, Pavez, J, Louppe, G, and Brooks, WK. Experiments using machine learning to approximate likelihood ratios for mixture models. In *Journal of Physics: Conference Series*, volume 762, pp. 012034. IOP Publishing, 2016.
- Cranmer, Kyle, Pavez, Juan, and Louppe, Gilles. Approximating likelihood ratios with calibrated discriminative classifiers. *arXiv preprint arXiv:1506.02169*, 2015.
- Cuturi, Marco. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pp. 2292–2300, 2013.
- Donahue, Jeff, Krähenbühl, Philipp, and Darrell, Trevor. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- Dumoulin, Vincent, Belghazi, Ishmael, Poole, Ben, Lamb, Alex, Arjovsky, Martin, Mastropietro, Olivier, and Courville, Aaron. Adversarially learned inference. *arXiv preprint arXiv:1606.00704*, 2016.
- Dutta, R., Corander, J., Kaski, S., and Gutmann, M. U. Likelihood-free inference by ratio estimation. *ArXiv e-prints*, November 2016.
- Genevay, Aude, Cuturi, Marco, Peyré, Gabriel, and Bach, Francis. Stochastic optimization for large-scale optimal transport. In *Advances in Neural Information Processing Systems*, pp. 3440–3448, 2016.
- Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Grathwohl, W., Choi, D., Wu, Y., Roeder, G., and Duvenaud, D. Backpropagation through the Void: Optimizing control variates for black-box gradient estimation. *ArXiv e-prints*, October 2017.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. Improved Training of Wasserstein GANs. *ArXiv e-prints*, March 2017.
- Gutmann, Michael U and Hyvärinen, Aapo. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of Machine Learning Research*, 13(Feb):307–361, 2012.
- Gutmann, Michael U, Dutta, Ritabrata, Kaski, Samuel, and Corander, Jukka. Likelihood-free inference via classification. *Statistics and Computing*, pp. 1–15, 2017.
- Huszár, F. Variational Inference using Implicit Distributions. *ArXiv e-prints*, February 2017.
- Kingma, D. P. and Ba, J. Adam: A Method for Stochastic Optimization. *ArXiv e-prints*, December 2014.
- Kingma, Diederik P. and Welling, Max. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013. URL <http://arxiv.org/abs/1312.6114>.
- Marin, Jean-Michel, Pudlo, Pierre, Robert, Christian P, and Ryder, Robin J. Approximate bayesian computational methods. *Statistics and Computing*, pp. 1–14, 2012.
- Marjoram, Paul, Molitor, John, Plagnol, Vincent, and Tavaré, Simon. Markov chain monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328, 2003.



Meeds, Edward, Leenders, Robert, and Welling, Max. Hamiltonian abc. *arXiv preprint arXiv:1503.01916*, 2015.

Mescheder, Lars M., Nowozin, Sebastian, and Geiger, Andreas. Adversarial variational bayes: Unifying variational autoencoders and generative adversarial networks. *CoRR*, abs/1701.04722, 2017. URL <http://arxiv.org/abs/1701.04722>.

Mohamed, S. and Lakshminarayanan, B. Learning in Implicit Generative Models. *ArXiv e-prints*, October 2016.

Montavon, Grégoire, Müller, Klaus-Robert, and Cuturi, Marco. Wasserstein training of restricted boltzmann machines. In *Advances in Neural Information Processing Systems*, pp. 3718–3726, 2016.

Ranganath, R., Altosaar, J., Tran, D., and Blei, D. M. Operator Variational Inference. *ArXiv e-prints*, October 2016.

Rosca, Mihaela, Lakshminarayanan, Balaji, Warde-Farley, David, and Mohamed, Shakir. Variational approaches for auto-encoding generative adversarial networks. *arXiv preprint arXiv:1706.04987*, 2017.

Rubin, Donald B. Bayesianly justifiable and relevant frequency calculations for the applied statistician. *Ann. Statist.*, 12(4):1151–1172, 12 1984. doi: 10.1214/aos/1176346785. URL <http://dx.doi.org/10.1214/aos/1176346785>.

Sisson, Scott A and Fan, Yanan. *Likelihood-free MCMC*. Chapman & Hall/CRC, New York.[839], 2011.

Sisson, Scott A, Fan, Yanan, and Tanaka, Mark M. Sequential monte carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 104(6):1760–1765, 2007.

Staines, J. and Barber, D. Variational Optimization. *ArXiv e-prints*, December 2012.

Staines, J and Barber, D. Optimization by variational bounding. In *ESANN 2013 proceedings, 21st European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pp. 473–478, 2013.

Tran, D., Ranganath, R., and Blei, D. M. Hierarchical Implicit Models and Likelihood-Free Variational Inference. *ArXiv e-prints*, February 2017.

Tran, Minh-Ngoc, Nott, David J, and Kohn, Robert. Variational bayes with intractable likelihood. *Journal of Computational and Graphical Statistics*, (just-accepted), 2017.

Wierstra, D., Schaul, T., Glasmachers, T., Sun, Y., and Schmidhuber, J. Natural Evolution Strategies. *ArXiv e-prints*, June 2011.

## A. Further experiments

### A.1. Electron–positron annihilation

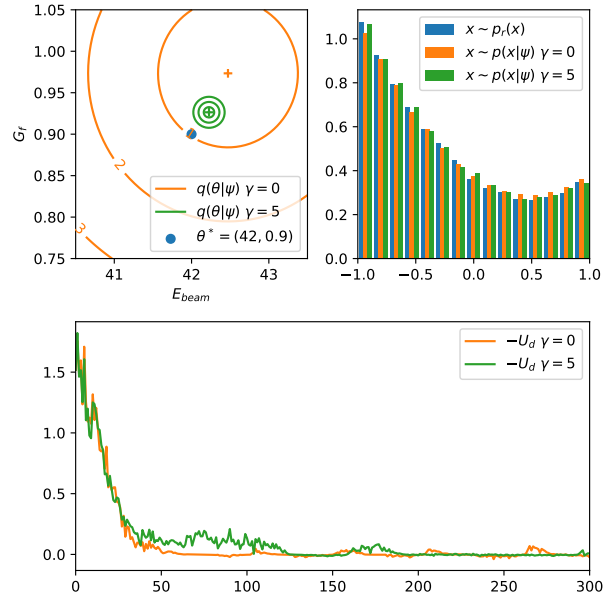


Figure 3. Electron–positron annihilation. (Top left) Proposal distributions  $q(\theta|\psi)$  after adversarial variational optimization. Contours correspond to parameters  $\theta$  within 1-2-3 Mahalanobis distance units from the mean of  $q(\theta|\psi)$ . The density of the penalized distribution ( $\gamma = 5$ ) is highly concentrated, resulting in the green mass near  $\theta^*$ . (Top right) Model distributions  $q(\mathbf{x}|\psi)$  after training. Despite the differences between their proposal distributions, both models closely match the observed data. (Bottom) Empirical estimates of the variational upper bound  $U_d$  as optimization progresses.

As a more realistic example, we here consider a (simplified) simulator from particle physics for electron–positron collisions resulting in muon–antimuon pairs ( $e^+e^- \rightarrow \mu^+\mu^-$ ). The simulator approximates the distribution of observed measurements  $\mathbf{x} = \cos(A) \in [-1, 1]$ , where  $A$  is the polar angle of the outgoing muon with respect to the originally incoming electron. Neglecting measurement uncertainty induced from the particle detectors, this random variable is approximately distributed as

$$p(\mathbf{x}|E_{\text{beam}}, G_f) = \frac{1}{Z} [(1 + \mathbf{x}^2) + c(E_{\text{beam}}, G_f)\mathbf{x}] \quad (19)$$

where  $Z$  is a known normalization constant and  $c$  is an asymmetry coefficient function. Due to the linear term in



the expression, the density  $p(\mathbf{x}|E_{\text{beam}}, G_f)$  exhibits a so-called *forward-backward* asymmetry. Its size depends on the values of the parameters  $E_{\text{beam}}$  (the beam energy) and  $G_f$  (the Fermi constant) through the coefficient function  $c$ .

A typical physics simulator for this process includes a more precise treatment of the quantum mechanical  $e^+e^- \rightarrow \mu^+\mu^-$  scattering using MadGraph (Alwall et al., 2011), ionization of matter in the detector due to the passage of the out-going  $\mu^+\mu^-$  particles using GEANT4 (Agostinelli et al., 2003), electronic noise and other details of the sensors that measure the ionization signal, and the deterministic algorithms that estimate the polar angle  $A$  based on the sensor readouts. The simulation of this process is highly non-trivial as is the space of latent variables  $\mathcal{Z}$ .

In this example, we consider observed data generated with the simplified generator of Eqn. 19 using  $\theta^* = (E_{\text{beam}}^* = 42, G_f^* = 0.9)$ . The simulator parameters are modeled with a factored (mean field) Gaussian proposal distribution  $q(\theta|\psi) = q(E_{\text{beam}}|\psi)q(G_f|\psi)$ , where each component is respectively initialized with mean 45 and 1 and variance 1 and 0.01. Hyper-parameters are set to  $M = 64$ ,  $n_{\text{critic}} = 100$ ,  $\lambda = 0.0025$  and Adam configured with  $\alpha = 0.01$ ,  $\beta_1 = 0.9$  and  $\beta_2 = 0.999$ . As with the first example, we compare entropy penalties  $\gamma = 0$  and  $\gamma = 5$ .

The top left plot in Figure 3 illustrates the resulting proposal distributions  $q(\theta|\psi)$  for  $\gamma = 0$  and  $\gamma = 5$  after AVO. We see that the distributions both arrive in the neighborhood of  $\theta^*$ , with a density more highly concentrated for  $\gamma = 5$  than for  $\gamma = 0$ . Despite these differences and the relative distance with  $\theta^*$ , both models closely match the observed data, as shown in the top right plot of Figure 3, with again slightly better results when using entropic regularization. This suggests either a relatively flat landscape around  $\theta^*$  or that the observed data can in this case also be reproduced with the predictive distribution  $q(\mathbf{x}|\psi)$ . Finally, the bottom plot of Figure 3 shows that for both  $\gamma = 0$  and  $\gamma = 5$  the variational upper bound  $-U_d$  quickly fall towards 0, which indicates convergence towards a distribution that the critic cannot distinguish from the true distribution.