

Advanced Machine Learning

Lecture 5

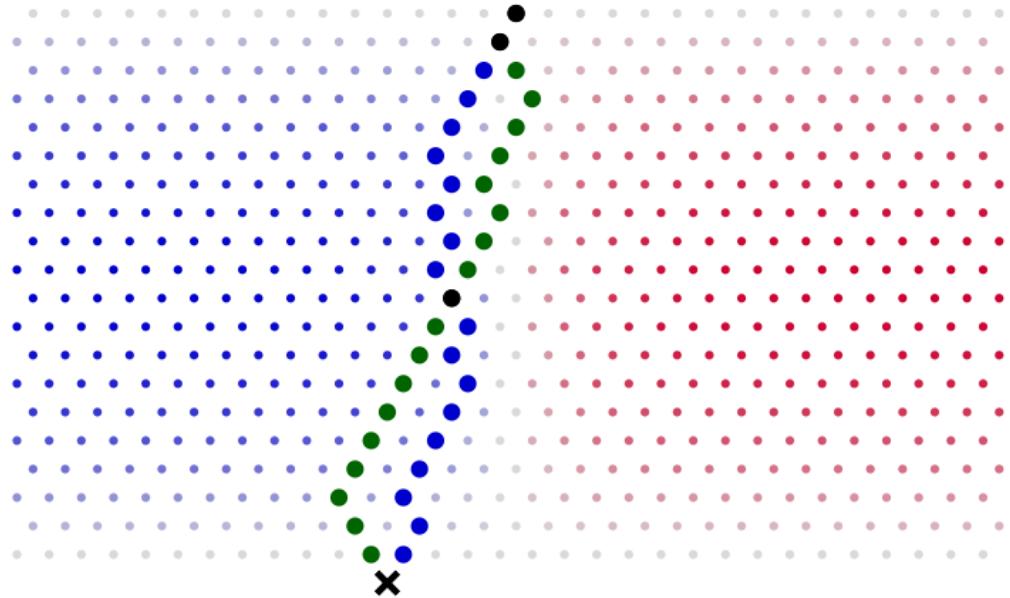
Paper: Gilles Louppe, Joeri Hermans and Kyle Cranmer,
Adversarial Variational Optimization of Non-Differentiable Simulators,
in International Conference on Artificial Intelligence and Statistics 2019.

Prof. Gilles Louppe
g.louppe@uliege.be

Likelihood-free inference



@physicsfun



The probability of ending in bin x corresponds to the total probability of all the paths z from start to x .

$$p(x|\theta) = \int p(x, z|\theta) dz = \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

Inference

Given a set of realizations $\mathbf{d} = \{x_i\}$ at the bins, inference consists in determining the value of θ that best describes these observations.

For example, following the principle of maximum likelihood estimation, we have

$$\hat{\theta} = \arg \max_{\theta} \prod_{x_i \in \mathbf{d}} p(x_i | \theta).$$

In general, when $p(x_i | \theta)$ can be evaluated, this problem can be solved either analytically or using optimization algorithms.

What if we shift or remove some of the pins?

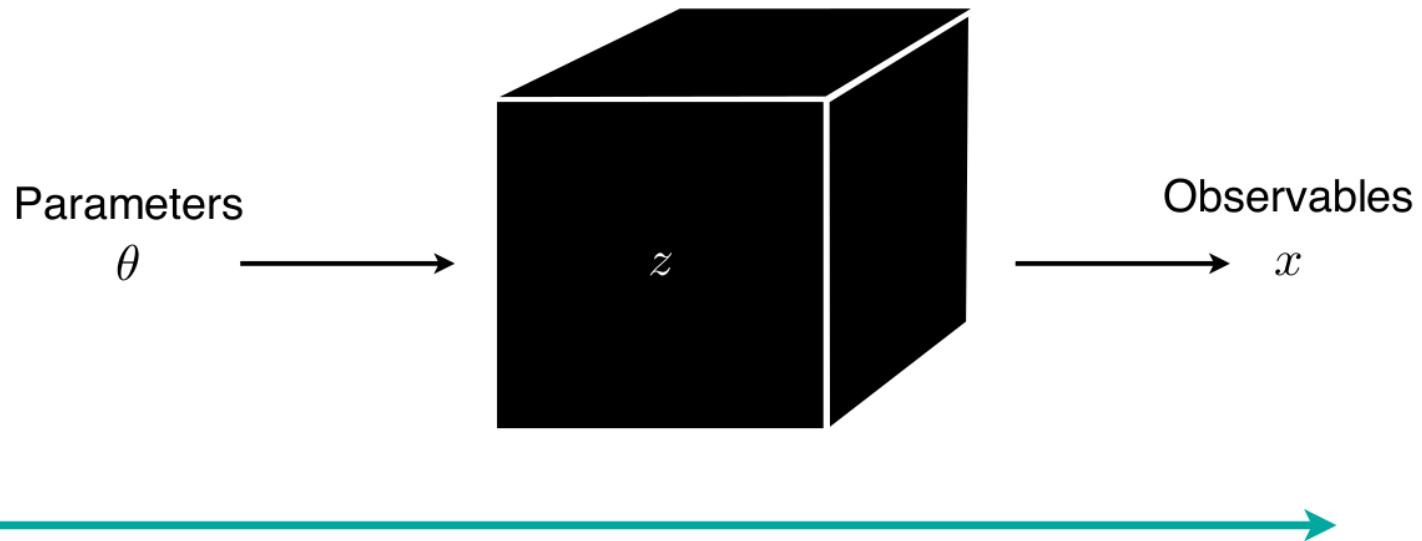
$$p(x|\theta) = \underbrace{\int p(x, z|\theta) dz}_{\text{intractable!}} \\ \neq \binom{n}{x} \theta^x (1-\theta)^{n-x}$$

Does this mean inference is no longer possible?

The Galton board is a **metaphor** of simulation-based science:

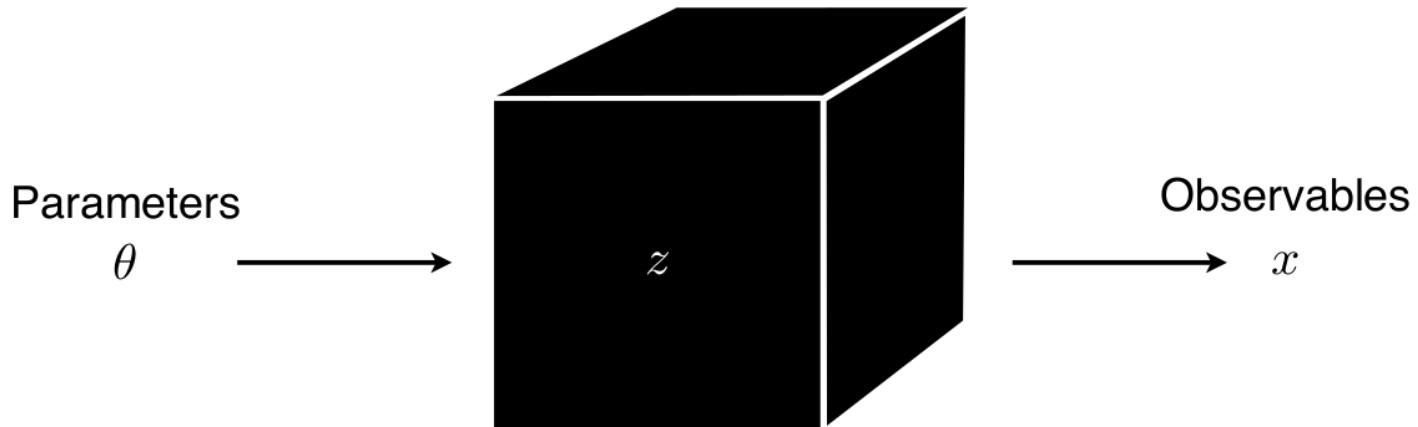
Galton board device	→	Computer simulation
Parameters θ	→	Model parameters θ
Buckets x	→	Observables x
Random paths z	→	Latent variables z (stochastic execution traces through simulator)

Inference in this context requires **likelihood-free algorithms**.



Prediction (simulation):

- Well-understood mechanistic model
- Simulator can generate samples



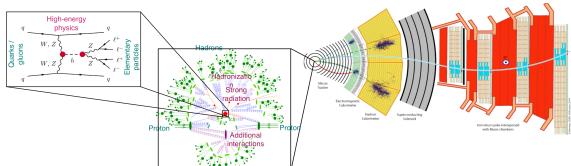
Prediction (simulation):

- Well-understood mechanistic model
- Simulator can generate samples

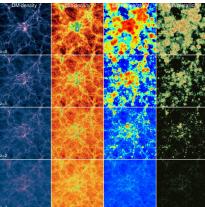
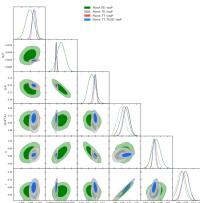
Inference:

- Likelihood function $p(x|\theta)$ is intractable
- Goal: estimator $\hat{p}(x|\theta)$

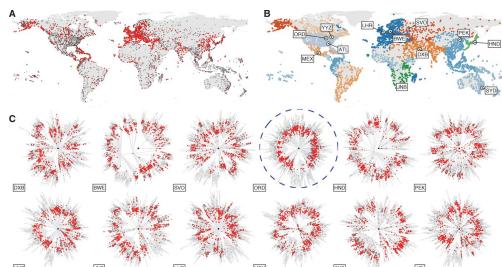
Applications



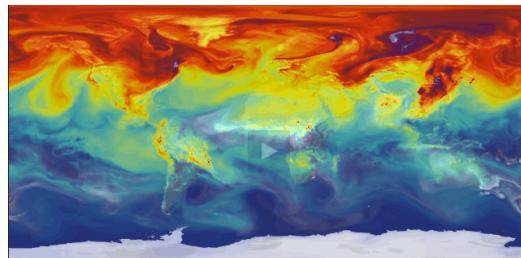
Particle physics



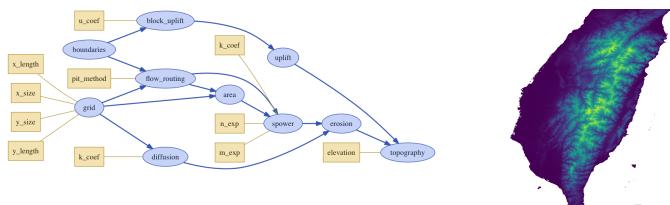
Cosmology



Epidemiology



Climatology

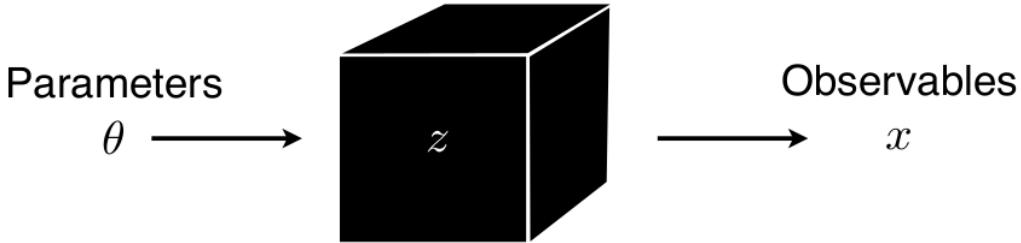


Computational topography



Astronomy

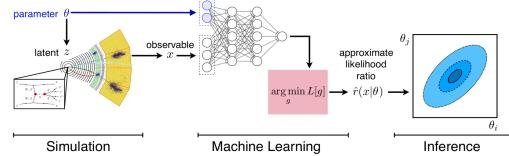
Particle physics



$$\begin{aligned}
\mathcal{L}_{SM} = & -\frac{1}{2}\partial_\nu g_\mu^a \partial_\nu g_\mu^a - g_s \epsilon^{abc} \partial_\mu g_\mu^a \partial_\mu g_\mu^c - \frac{1}{4}g_s^2 \epsilon^{abc} f^{ade} g_\mu^b g_\mu^d g_\mu^e - \partial_\mu W_\mu^+ \partial_\mu W_\mu^- \\
M^2 W_\mu^+ W_\mu^- = & -\frac{1}{2} \partial_\mu Z_\mu^0 Z_\mu^0 - \frac{1}{2} M^2 Z_\mu^0 Z_\mu^0 - \frac{1}{2} \partial_\mu A_\nu \partial_\mu A_\nu - ig s_w (\partial_\mu Z_\mu^0 W_\mu^+ W_\mu^- - \\
W_\mu^+ W_\mu^-) - & Z_\mu^0 (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) + Z_\mu^0 (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) - \\
ig s_w (\partial_\mu A_\nu (W_\mu^+ W_\mu^- - W_\mu^+ W_\mu^-) - A_\nu (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) + A_\nu (W_\mu^+ \partial_\mu W_\mu^- - \\
W_\mu^- \partial_\mu W_\mu^+)) - \frac{1}{2} g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + \frac{1}{2} g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + g^2 c_w^2 (Z_\mu^0 W_\mu^+ Z_\mu^0 W_\mu^- - \\
Z_\mu^0 Z_\mu^0 W_\mu^+ W_\mu^-) + g^2 s_w^2 (A_\mu W_\mu^+ A_\nu W_\nu^- + A_\mu A_\nu W_\mu^+ W_\nu^-) + g^2 s_w c_w (A_\mu Z_\mu^0 (W_\mu^+ W_\mu^- - \\
W_\mu^+ W_\mu^-) - 2 A_\mu Z_\mu^0 W_\mu^+ W_\mu^-) - \frac{1}{2} \partial_\mu H \partial_\mu H - 2 M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2} \partial_\mu \phi^0 \partial_\mu \phi^0 - \\
\beta_h \left(\frac{2M^2}{\sigma^2} + \frac{2M}{g} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2Ms}{\sigma^2} \alpha_h - \\
& g \alpha_h M (H^3 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\
\frac{1}{8} g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\
& g M W_\mu^+ W_\mu^- H - \frac{ig}{\sqrt{2}} Z_\mu^0 Z_\mu^0 H - \\
& \frac{1}{2} ig (W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - W_\mu^- (\phi^+ \partial_\mu \phi^+ - \phi^0 \partial_\mu \phi^0)) + \\
& \frac{1}{2} g (W_\mu^+ (H \partial_\mu \phi^- - \phi^- \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2} g \frac{1}{c_w} (Z_\mu^0 (H \partial_\mu \phi^0 - \phi^0 \partial_\mu H) + \\
M \frac{1}{c_w} Z_\mu^0 \partial_\mu \phi^0 + W_\mu^+ \partial_\mu \phi^- + W_\mu^- \partial_\mu \phi^+) - ig \frac{2}{c_w} M Z_\mu^0 (W_\mu^+ \phi^- - W_\mu^- \phi^+) + ig s_w M A_\mu (W_\mu^+ \phi^- - \\
W_\mu^- \phi^+) - ig \frac{1-2c_w}{2c_w} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\
& \frac{1}{4} g^2 W_\mu^+ W_\mu^- (H^2 + (\phi^0)^2 + 2\phi^+ \phi^-) - \frac{1}{8} g^2 \frac{s_w^2}{c_w^2} Z_\mu^0 Z_\mu^0 (H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+ \phi^-) - \\
& \frac{1}{2} g^2 \frac{s_w^2}{c_w^2} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2} g^2 \frac{s_w^2}{c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2} g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\
W_\mu^- \phi^+) + \frac{1}{2} g^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{2}{c_w} (2c_w^2 - 1) A_\mu \phi^+ \phi^- - \\
g^2 \frac{s_w^2}{c_w^2} A_\mu A_\nu \phi^+ \phi^- + \frac{1}{2} ig s_w \lambda_\mu^\alpha (q_\mu^\beta \gamma^\gamma q_\nu^\delta) g_\mu^\alpha - \bar{e}^\lambda (\gamma \partial + m_\lambda^2) \bar{e}^\lambda - \bar{\nu}^\lambda (\gamma \partial + m_\lambda^2) \bar{\nu}^\lambda - \bar{u}_j^\lambda (\gamma \partial + \\
m_\lambda^2) u_j^\lambda + \bar{d}_j^\lambda (\gamma \partial + m_\lambda^2) d_j^\lambda + ig s_w A_\mu (-\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3} (\bar{u}_j^\lambda \gamma^\mu u_j^\lambda) - \frac{1}{3} (\bar{d}_j^\lambda \gamma^\mu d_j^\lambda) + \\
\frac{ig}{4c_w} Z_\mu^0 ((\bar{p}^\nu \gamma^\mu (1 - \gamma^5) p^\lambda) + (\bar{e}^\lambda \gamma^\mu (1 - \gamma^5) e^\lambda) + (\bar{d}_j^\lambda \gamma^\mu (1 - \gamma^5) d_j^\lambda) + (\bar{u}_j^\lambda \gamma^\mu (1 - \gamma^5) u_j^\lambda) + \\
(\bar{u}_j^\lambda \gamma^\mu (1 - \frac{3}{2}s_w^2 + \gamma^5) u_j^\lambda)) + \frac{ig}{2\sqrt{2}} W_\mu^+ ((\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) l^\nu) e^\lambda + (\bar{l}^\nu C_{\lambda\kappa}^\mu (1 + \gamma^5) u_j^\lambda)) + \\
\frac{ig}{2\sqrt{2}} W_\mu^- ((\bar{e}^\nu U^{l\mu p}{}_{\lambda\kappa}^\lambda (1 + \gamma^5) \nu^\kappa) + (\bar{d}_j^\nu C_{\lambda\kappa}^\mu (1 + \gamma^5) u_j^\lambda)) + \\
\frac{ig}{2\sqrt{2}} \phi^+ (-m_\nu^\kappa (\bar{p}^\lambda U^{l\mu p}{}_{\lambda\kappa}^\lambda (1 - \gamma^5) e^\kappa) + m_\mu^\kappa (\bar{\nu}^\lambda U^{l\mu p}{}_{\lambda\kappa}^\lambda (1 - \gamma^5) \nu^\kappa) - \\
m_\lambda^\kappa (\bar{e}^\nu U^{l\mu p}{}_{\lambda\kappa}^\lambda (1 + \gamma^5) \nu^\kappa) - \frac{im_\lambda^\kappa}{2} H (\bar{p}^\lambda \bar{p}^\lambda) - \frac{im_\lambda^\kappa}{2} H (\bar{\nu}^\lambda \bar{\nu}^\lambda) - \\
\frac{ig}{2\sqrt{2}} H (\bar{e}^\lambda \lambda_\mu^\nu) + \frac{ig}{2} \frac{m_\lambda^2}{M} \phi^0 (\bar{p}^\lambda \gamma^\nu \bar{p}^\lambda) - \frac{ig}{2} \frac{m_\lambda^2}{M} \phi^0 (\bar{e}^\lambda \gamma^\nu \bar{e}^\lambda) - \frac{1}{4} \bar{\nu}_\kappa M_{\lambda\kappa}^R (1 - \gamma_5) \bar{\nu}_\kappa - \\
\frac{1}{4} \bar{\nu}_\lambda M_{\lambda\kappa}^R (1 - \gamma_5) \bar{\nu}_\kappa + \frac{ig}{2\sqrt{2}} \phi^+ (-m_d^\nu (\bar{u}_j^\lambda C_{\lambda\kappa}^\mu (1 - \gamma^5) d_j^\lambda) + m_u^\nu (\bar{u}_j^\lambda C_{\lambda\kappa}^\mu (1 + \gamma^5) d_j^\lambda) + \\
\frac{ig}{2\sqrt{2}} \phi^+ (m_d^\lambda (\bar{d}_j^\nu C_{\lambda\kappa}^\mu (1 + \gamma^5) u_j^\lambda) - m_u^\lambda (\bar{d}_j^\nu C_{\lambda\kappa}^\mu (1 - \gamma^5) u_j^\lambda) - \frac{im_\lambda^\kappa}{2} H (\bar{u}_j^\nu u_j^\lambda)) - \\
\frac{g}{2} \frac{m_\lambda^2}{M} H (\bar{d}_j^\nu d_j^\lambda) + \frac{ig}{2} \frac{m_\lambda^2}{M} \phi^0 (\bar{u}_j^\lambda \bar{u}_j^\nu) - \frac{ig}{2} \frac{m_\lambda^2}{M} \phi^0 (\bar{d}_j^\lambda \bar{d}_j^\nu) + G^\mu \partial^\nu G^\nu + g_s f^{abc} \partial_\mu G^\mu G^\nu g_\nu^a + \\
X^+ (\partial^2 - M^2) X^+ + X^- (\partial^2 - M^2) X^- + \bar{X}^0 (\partial^2 - \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^2 Y + ig c_w W_\mu^- (\partial_\mu \bar{X}^0 X^- - \\
\partial_\mu \bar{X}^+ X^+) + ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \partial_\mu \bar{X}^+ Y) + ig c_w W_\mu^- (\partial_\mu \bar{X}^- X^+ - \\
\partial_\mu \bar{X}^+ X^-) + ig s_w A_\mu (\partial_\mu \bar{X}^+ X^-) + \\
\partial_\mu X^- X^- - \frac{1}{2} g M \left(\bar{X}^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2} \bar{X}^0 X^0 H \right) + \frac{1-2c_w^2}{2c_w^2} ig M (\bar{X}^+ X^0 \phi^+ - \bar{X}^- X^0 \phi^-) + \\
\frac{1}{2c_w} ig M (\bar{X}^0 X^- \phi^+ - X^0 X^+ \phi^-) + ig M s_w (\bar{X}^0 X^- \phi^+ - X^0 X^+ \phi^-) + \\
\frac{1}{2} ig M (\bar{X}^- X^+ \phi^0 - \bar{X}^+ X^- \phi^0) .
\end{aligned}$$

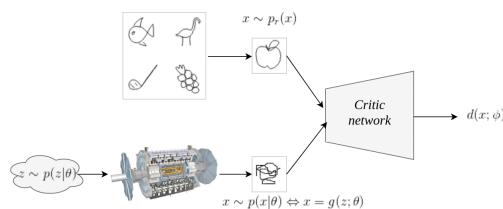
Algorithms

Learn a proxy for inference



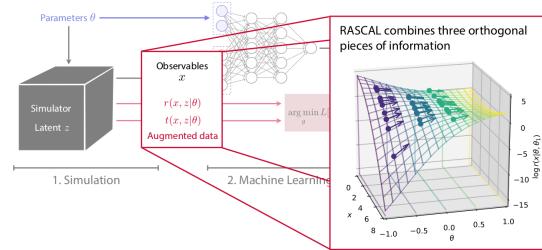
Histograms of observables
Neural density (ratio)
estimation

Learn to control the simulator

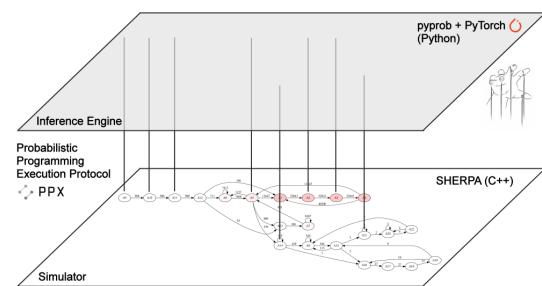


Adversarial variational
optimization

Make use of the inner structure



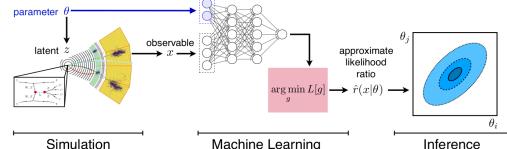
Mining gold from implicit models



Probabilistic programming

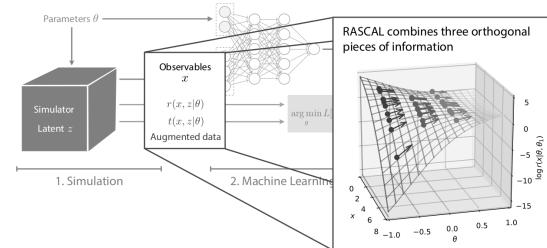
Treat the simulator as a black box

**Learn a proxy
for inference**



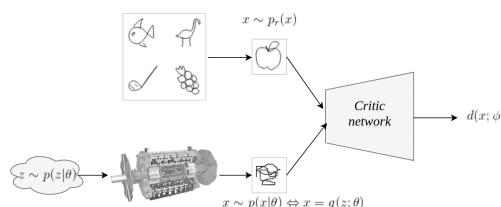
Histograms of observables
Neural density (ratio)
estimation

Make use of the inner structure

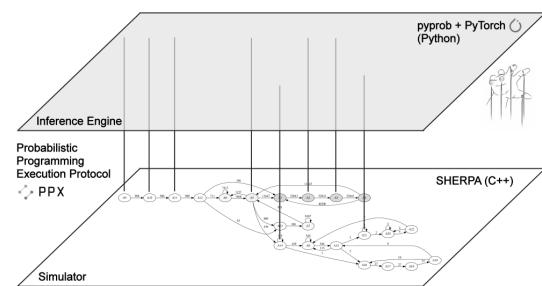


Mining gold from implicit models

**Learn to
control the
simulator**



Adversarial variational
optimization



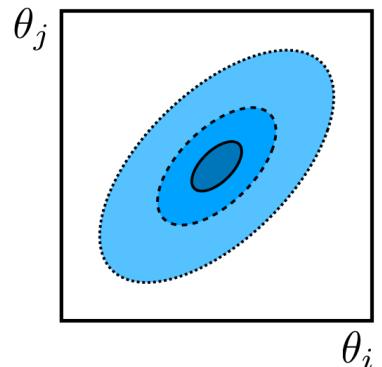
Probabilistic programming

The physicist's way

The Neyman-Pearson lemma states that the likelihood ratio

$$r(x|\theta_0, \theta_1) = \frac{p(x|\theta_0)}{p(x|\theta_1)}$$

is the most powerful test statistic to discriminate between a null hypothesis θ_0 and an alternative θ_1 .



IX. On the Problem of the most Efficient Tests of Statistical Hypotheses.

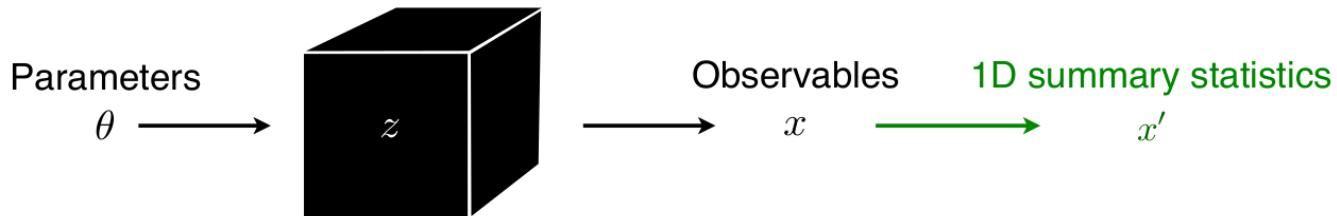
By J. NEYMAN, Nencki Institute, Soc. Sci. Lit. Varsoviensis, and Lecturer at the Central College of Agriculture, Warsaw, and E. S. PEARSON, Department of Applied Statistics, University College, London.

(Communicated by K. PEARSON, F.R.S.)

(Received August 31, 1932.—Read November 10, 1932.)

CONTENTS.

	PAGE.
I. Introductory	289
II. Outline of General Theory	293
III. Simple Hypotheses	



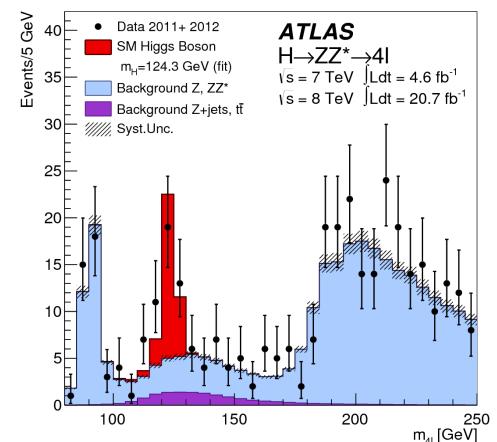
Define a projection function $s : \mathcal{X} \rightarrow \mathbb{R}$ mapping observables x to a summary statistics $x' = s(x)$.

Then, approximate the likelihood $p(x|\theta)$ as

$$p(x|\theta) \approx \hat{p}(x|\theta) = p(x'|\theta).$$

From this it comes

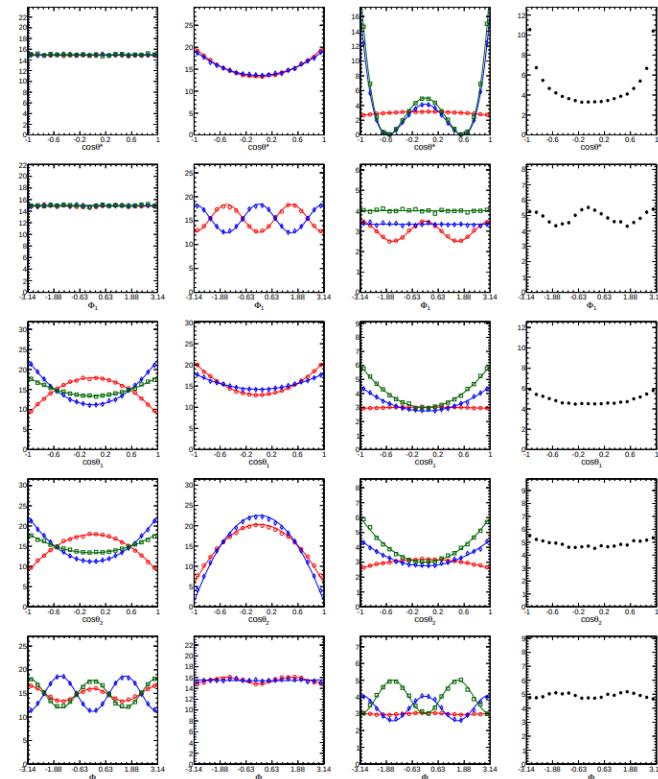
$$\frac{p(x|\theta_0)}{p(x|\theta_1)} \approx \frac{\hat{p}(x|\theta_0)}{\hat{p}(x|\theta_1)} = \hat{r}(x|\theta_0, \theta_1).$$



This methodology has worked great for physicists for the last 20-30 years, but ...

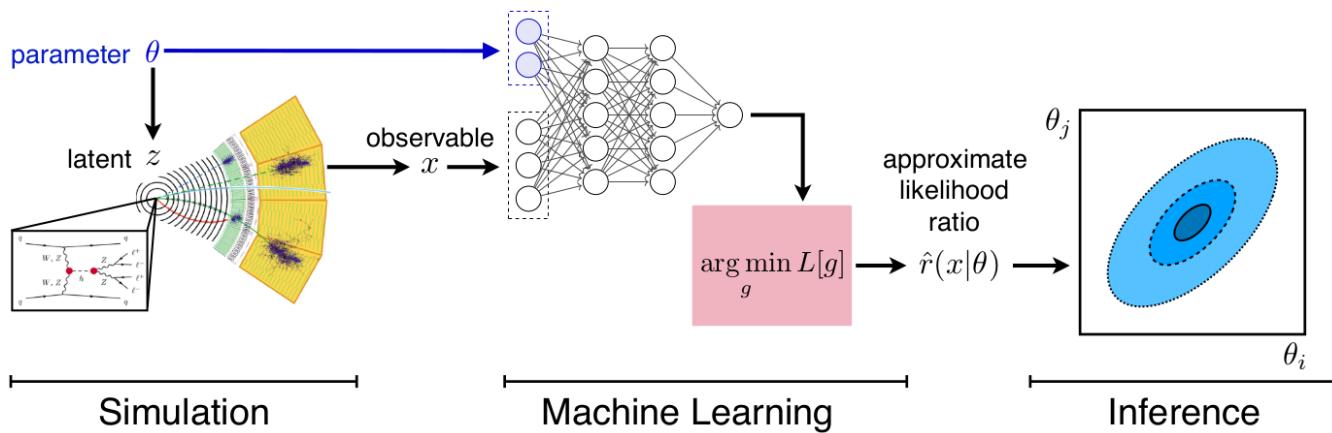
- Choosing the projection s is difficult and problem-dependent.
- Often there is no single good variable: compressing to any x' loses information.
- Ideally: analyse high-dimensional x' , including all correlations.

Unfortunately, filling high-dimensional histograms is **not tractable**.



Who you gonna call? **Machine learning!**

CARL



Key insights

- The likelihood ratio is **sufficient** for inference.
- Evaluating the likelihood ratio **does not** require evaluating the individual likelihoods.
- Supervised learning indirectly estimates likelihood ratios.

Supervised learning provides a way to automatically construct \hat{s} :

- Let us consider a binary classifier \hat{s} (e.g., a neural network) trained to distinguish $x \sim p(x|\theta_0)$ from $x \sim p(x|\theta_1)$.
- \hat{s} is trained by minimizing the cross-entropy loss

$$L_{XE}[\hat{s}] = -\mathbb{E}_{p(x|\theta)\pi(\theta)}[1(\theta = \theta_0) \log \hat{s}(x) + 1(\theta = \theta_1) \log(1 - \hat{s}(x))]$$

The solution \hat{s} found after training approximates the optimal classifier

$$\hat{s}(x) \approx s^*(x) = \frac{p(x|\theta_1)}{p(x|\theta_0) + p(x|\theta_1)}.$$

Therefore,

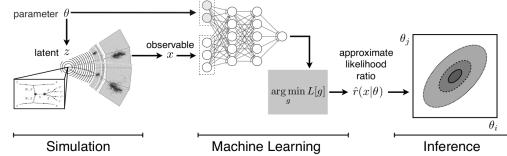
$$r(x|\theta_0, \theta_1) \approx \hat{r}(x|\theta_0, \theta_1) = \frac{1 - \hat{s}(x)}{\hat{s}(x)}$$

That is, **supervised classification is equivalent to likelihood ratio estimation** and can therefore be used for MLE inference.

Adversarial Variational Optimization

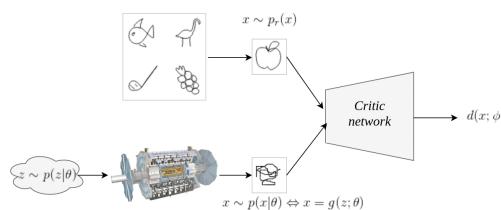
Learn a proxy for inference

Treat the simulator as a black box



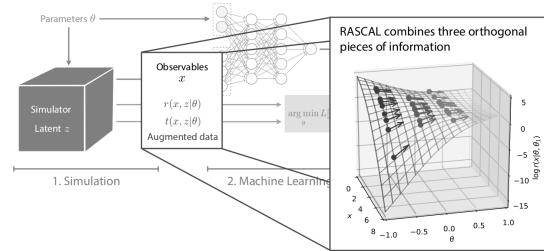
Histograms of observables
Neural density (ratio)
estimation

Learn to control the simulator

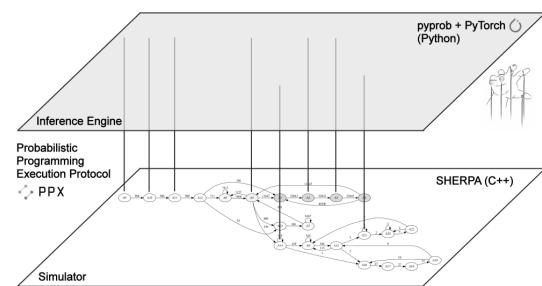


Adversarial variational
optimization

Make use of the inner structure



Mining gold from implicit models



Probabilistic programming

Adversarial Variational Optimization of Non-Differentiable Simulators

Gilles Louppe
University of Liège

Joeri Hermans
University of Liège

Kyle Cranmer
New York University

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, variational optimization and empirical Bayes. We adapt the training procedure of generative adversarial networks 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 JS divergence between the marginal distribution of the synthetic data and the empirical distribution of observed data is minimized. We evaluate and compare the method with simulators producing both discrete and continuous data.

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 and Fan, 2011](#); [Marin et al., 2012](#); [Cranmer et al., 2015](#)).

In parallel, with the introduction of variational auto-encoders ([Kingma and 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 and 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 parameterized 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 θ is often of more interest than the latent variables \mathbf{z} .

In this work, we introduce Adversarial Variational Optimization (AVO), a likelihood-free inference algorithm for non-differentiable, implicit generative models. We adapt the adversarial training procedure of generative adversarial networks by replacing the implicit generative network with a domain-based scientific simula-

Generative adversarial networks

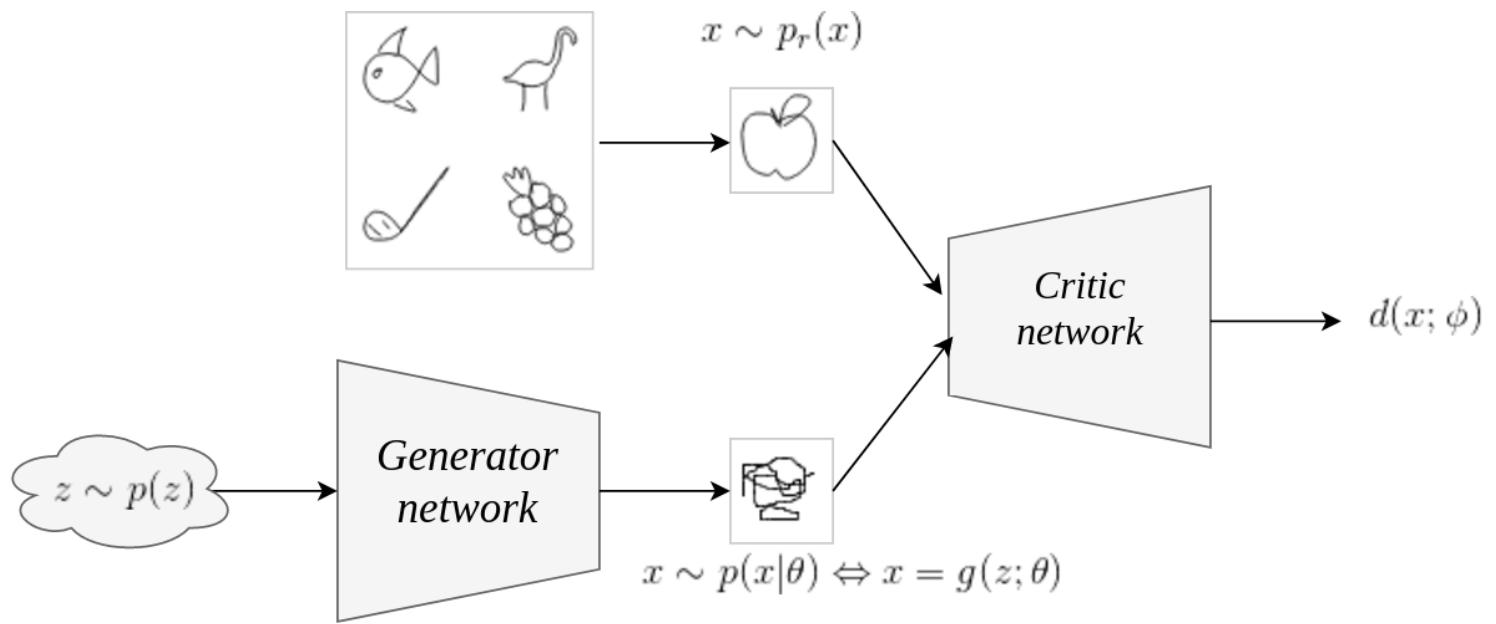
Express the task of learning a generative model as a two-player zero-sum game between two networks:

- The first network is a generator $g(\cdot; \theta) : \mathcal{Z} \rightarrow \mathcal{X}$, mapping a latent space equipped with a prior distribution $p(\mathbf{z})$ to the data space, thereby inducing a distribution

$$\mathbf{x} \sim p(\mathbf{x}; \theta) \Leftrightarrow \mathbf{z} \sim p(\mathbf{z}), \mathbf{x} = g(\mathbf{z}; \theta).$$

- The second network $d(\cdot; \phi) : \mathcal{X} \rightarrow [0, 1]$ is a classifier trained to distinguish between true samples $\mathbf{x} \sim p_r(\mathbf{x})$ and generated samples $\mathbf{x} \sim p(\mathbf{x}; \theta)$.

The central mechanism will be to use supervised learning to guide the learning of the generative model.



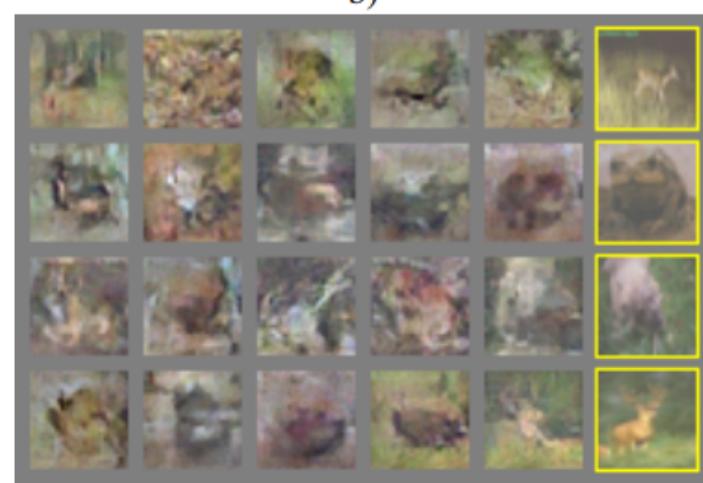
$$\arg \min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log d(\mathbf{x}; \phi)] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - d(g(\mathbf{z}; \theta); \phi))]$$

In practice, the min-max solution is approximated using **alternating stochastic gradient descent** on the losses

$$\begin{aligned}\mathcal{L}_d(\phi) &= \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [-\log(d(\mathbf{x}; \phi))] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [-\log(1 - d(g(\mathbf{z}; \theta); \phi))] \\ \mathcal{L}_g(\theta) &= \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - d(g(\mathbf{z}; \theta); \phi))],\end{aligned}$$

for which unbiased gradient estimates can be computed with Monte Carlo integration.

- For one step on θ , we can optionally take k steps on ϕ , since we need the classifier to remain near optimal.
- Note that to compute $\nabla_{\theta} \mathcal{L}_g$, it is necessary to backprop all the way through d before computing the partial derivatives with respect to g 's internals.



Goodfellow et al, 2014.



Odena et al
2016



Miyato et al
2017



Zhang et al
2018



Brock et al
2018



Figure 2. Uncurated set of images produced by our style-based generator (config F) with the FFHQ dataset. Here we used a variation of the truncation trick [5, 29] with $\psi = 0.7$ for resolutions $4^2 - 32^2$. Please see the accompanying video for more results.

Karras et al, 2018.

Game analysis

Consider a generator \mathbf{g} fixed at θ . Given a set of observations

$$\mathbf{x}_i \sim p_r(\mathbf{x}), i = 1, \dots, N,$$

we can generate a two-class dataset

$$\mathbf{d} = \{(\mathbf{x}_1, 1), \dots, (\mathbf{x}_N, 1), (g(\mathbf{z}_1; \theta), 0), \dots, (g(\mathbf{z}_N; \theta), 0)\}.$$

The best classifier d is obtained by minimizing the cross-entropy

$$\begin{aligned}\mathcal{L}_d(\phi) &= -\frac{1}{2N} \left(\sum_{i=1}^N [\log d(\mathbf{x}_i; \phi)] + \sum_{i=1}^N [\log(1 - d(g(\mathbf{z}_i; \theta); \phi))] \right) \\ &\approx -\frac{1}{2} \left(\mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log d(\mathbf{x}; \phi)] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - d(g(\mathbf{z}; \theta); \phi))] \right)\end{aligned}$$

with respect to ϕ .

Let us define the **value function**

$$V(\phi, \theta) = \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} [\log d(\mathbf{x}; \phi)] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - d(g(\mathbf{z}; \theta); \phi))].$$

We have,

- $V(\phi, \theta)$ is high if d is good at recognizing true from generated samples.
- If d is the best classifier given g , and if V is high, then this implies that the generator is bad at reproducing the data distribution.
- Conversely, g will be a good generative model if V is low when d is a perfect opponent.

Therefore, the ultimate goal is

$$\theta^* = \arg \min_{\theta} \max_{\phi} V(\phi, \theta).$$

For a generator \mathbf{g} fixed at θ , the classifier d with parameters ϕ_θ^* is optimal if and only if

$$\forall \mathbf{x}, d(\mathbf{x}; \phi_\theta^*) = \frac{p_r(\mathbf{x})}{p(\mathbf{x}; \theta) + p_r(\mathbf{x})}.$$

Therefore,

$$\begin{aligned} \min_{\theta} \max_{\phi} V(\phi, \theta) &= \min_{\theta} V(\phi_\theta^*, \theta) \\ &= \min_{\theta} \mathbb{E}_{\mathbf{x} \sim p_r(\mathbf{x})} \left[\log \frac{p_r(\mathbf{x})}{p(\mathbf{x}; \theta) + p_r(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x}; \theta)} \left[\log \frac{p(\mathbf{x}; \theta)}{p(\mathbf{x}; \theta) + p_r(\mathbf{x})} \right] \\ &= \min_{\theta} \text{KL} \left(p_r(\mathbf{x}) \parallel \frac{p_r(\mathbf{x}) + p(\mathbf{x}; \theta)}{2} \right) \\ &\quad + \text{KL} \left(p(\mathbf{x}; \theta) \parallel \frac{p_r(\mathbf{x}) + p(\mathbf{x}; \theta)}{2} \right) - \log 4 \\ &= \min_{\theta} 2 \text{JSD}(p_r(\mathbf{x}) \parallel p(\mathbf{x}; \theta)) - \log 4 \end{aligned}$$

where **JSD** is the Jensen-Shannon divergence.

In summary, solving the min-max problem

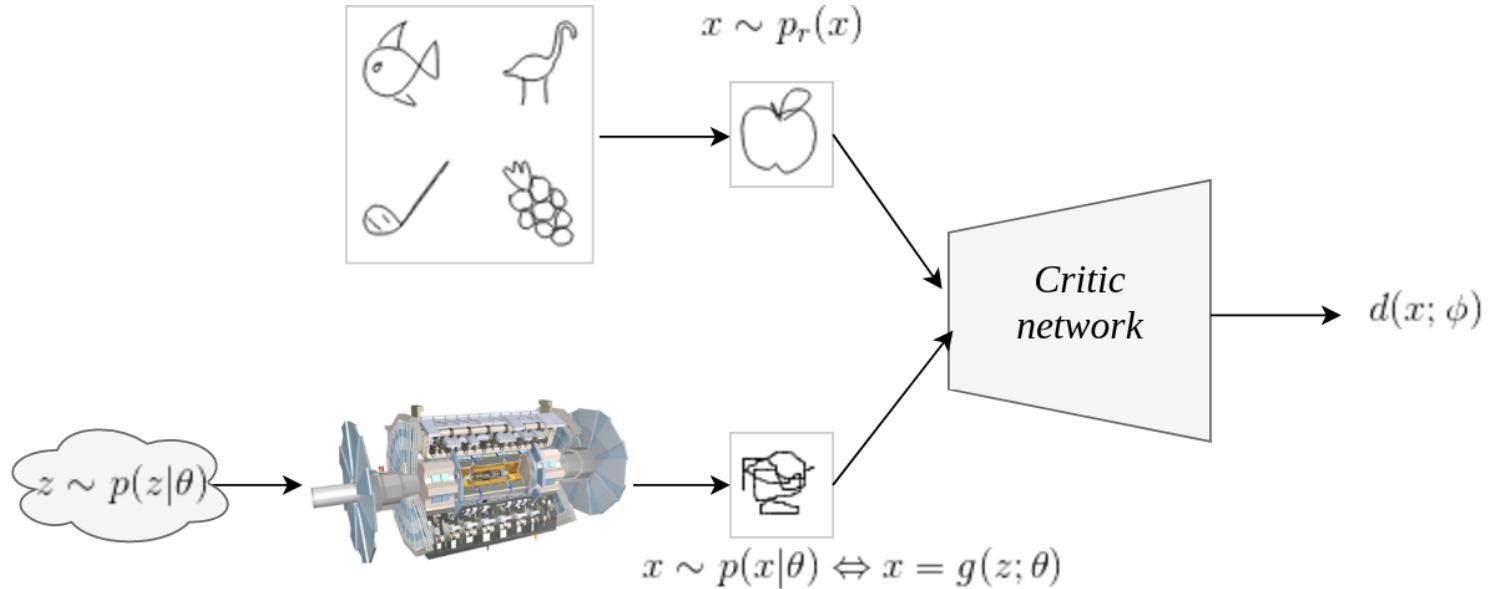
$$\theta^* = \arg \min_{\theta} \max_{\phi} V(\phi, \theta)$$

is equivalent to

$$\theta^* = \arg \min_{\theta} \text{JSD}(p_r(\mathbf{x}) || p(\mathbf{x}; \theta)).$$

Since $\text{JSD}(p_r(\mathbf{x}) || p(\mathbf{x}; \theta))$ is minimum if and only if $p_r(\mathbf{x}) = p(\mathbf{x}; \theta)$, this proves that the min-max solution corresponds to a generative model that perfectly reproduces the true data distribution.

AVO



Replace g with an actual scientific simulator!

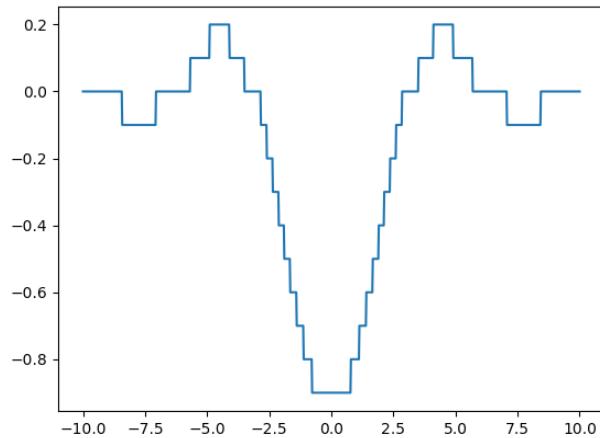
Key insights

- Replace the generative network with a non-differentiable forward simulator $g(\mathbf{z}; \theta)$.
- Let the neural network critic figure out how to adjust the simulator parameters.
- Bypass the non-differentiability with variational optimization.

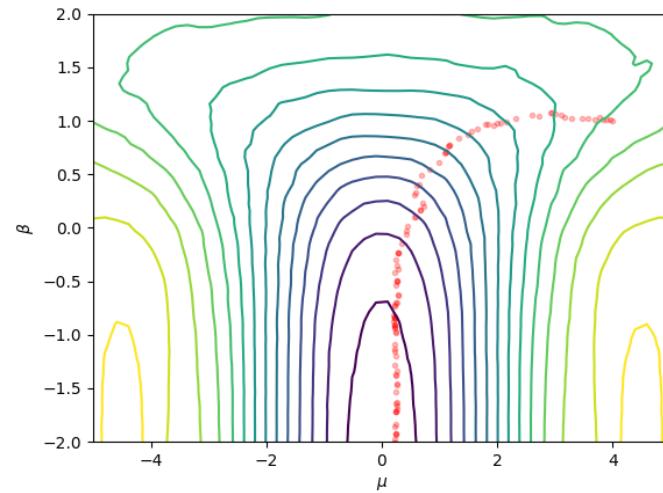
Variational optimization

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

$$\nabla_{\psi} U(\psi) = \mathbb{E}_{\theta \sim q(\theta|\psi)} [f(\theta) \nabla_{\psi} \log q(\theta|\psi)]$$



Piecewise constant — $\frac{\sin(x)}{x}$



$$q(\theta|\psi = (\mu, \beta)) = \mathcal{N}(\mu, e^{\beta})$$

(Similar to REINFORCE gradient estimates)

Variational optimization can be used to bypass the non-differentiability of the adversarial game by optimizing upper bounds of the adversarial objectives

$$U_d(\phi) = \mathbb{E}_{\theta \sim q(\theta; \psi)} [\mathcal{L}_d(\phi)]$$
$$U_g(\psi) = \mathbb{E}_{\theta \sim q(\theta; \psi)} [\mathcal{L}_g(\theta)]$$

respectively over ϕ and ψ .

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 k of training iterations of the discriminator d (default: $k = 1$),
The size M of a mini-batch (default: $M = 32$),
The R_1 regularization coefficient λ (default: $\lambda = 10$),
The entropy penalty coefficient γ (default: $\gamma = 0$).
The baseline strategy b in REINFORCE estimates (default: Eqn. 13).

- 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 k **do** ▷ Update d
 - 4: Sample true data $\{\mathbf{x}_m \sim p_r(\mathbf{x}), y_m = 1\}_{m=1}^{M/2}$.
 - 5: Sample synthetic data $\{\boldsymbol{\theta}_m \sim q(\boldsymbol{\theta}|\boldsymbol{\psi}), \mathbf{z}_m \sim p(\mathbf{z}|\boldsymbol{\theta}_m), \tilde{\mathbf{x}}_m = g(\mathbf{z}_m; \boldsymbol{\theta}_m), y_m = 0\}_{m=M/2+1}^M$.
 - 6: $\nabla_{\boldsymbol{\phi}} U_d \leftarrow \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\phi}} [-y_m \log(d(\mathbf{x}_m; \boldsymbol{\phi})) - (1 - y_m) \log(1 - d(\mathbf{x}_m; \boldsymbol{\phi}))]$
 - 7: $\nabla_{\boldsymbol{\phi}} R_1 \leftarrow \frac{1}{M/2} \sum_{m=1}^{M/2} \nabla_{\boldsymbol{\phi}} [\|\nabla_{\boldsymbol{\phi}} d(\mathbf{x}_m; \boldsymbol{\phi})\|^2]$
 - 8: $\boldsymbol{\phi} \leftarrow \text{RMSPROP}(\nabla_{\boldsymbol{\phi}} U_d + \lambda \nabla_{\boldsymbol{\phi}} R_1)$
 - 9: **end for**
 - 10: Sample synthetic data $\{\boldsymbol{\theta}_m \sim q(\boldsymbol{\theta}|\boldsymbol{\psi}), \mathbf{z}_m \sim p(\mathbf{z}|\boldsymbol{\theta}_m), \tilde{\mathbf{x}}_m = g(\mathbf{z}_m; \boldsymbol{\theta}_m)\}_{m=1}^M$. ▷ Update $q(\boldsymbol{\theta}|\boldsymbol{\psi})$
 - 11: $\nabla_{\boldsymbol{\psi}} U_g \leftarrow \frac{1}{M} \sum_{m=1}^M [\nabla_{\boldsymbol{\psi}} \log q(\boldsymbol{\theta}_m|\boldsymbol{\psi})(\log(1 - d(\tilde{\mathbf{x}}_m; \boldsymbol{\phi})) - b)]$
 - 12: $\nabla_{\boldsymbol{\psi}} H \leftarrow \frac{1}{M} \sum_{m=1}^M \nabla_{\boldsymbol{\psi}} [-q(\boldsymbol{\theta}_m|\boldsymbol{\psi}) \log q(\boldsymbol{\theta}_m|\boldsymbol{\psi})]$
 - 13: $\boldsymbol{\psi} \leftarrow \text{RMSPROP}(\nabla_{\boldsymbol{\psi}} U_g + \gamma \nabla_{\boldsymbol{\psi}} H)$
 - 14: **end while**
-

Operationally,

$$\mathbf{x} \sim q(\mathbf{x}|\psi) \Leftrightarrow \theta \sim q(\theta|\psi), \mathbf{z} \sim p(\mathbf{z}|\theta), \mathbf{x} = g(\mathbf{z}; \theta)$$

Therefore, $q(\mathbf{x}|\psi)$ is the marginal $\int p(\mathbf{x}|\theta)q(\theta|\psi)d\theta$ and the AVO procedure results in solving

$$\psi^* = \arg \min_{\psi} \text{JSD}(p_r(\mathbf{x}) || q(\mathbf{x}|\psi)).$$

- If the simulator $p(\mathbf{x}|\theta)$ is misspecified, $q(\mathbf{x}|\psi)$ will attempt to smear the simulator to approach $p_r(\mathbf{x})$.
- If not, $q(\theta|\psi)$ will concentrate its mass around the true data-generating parameters.

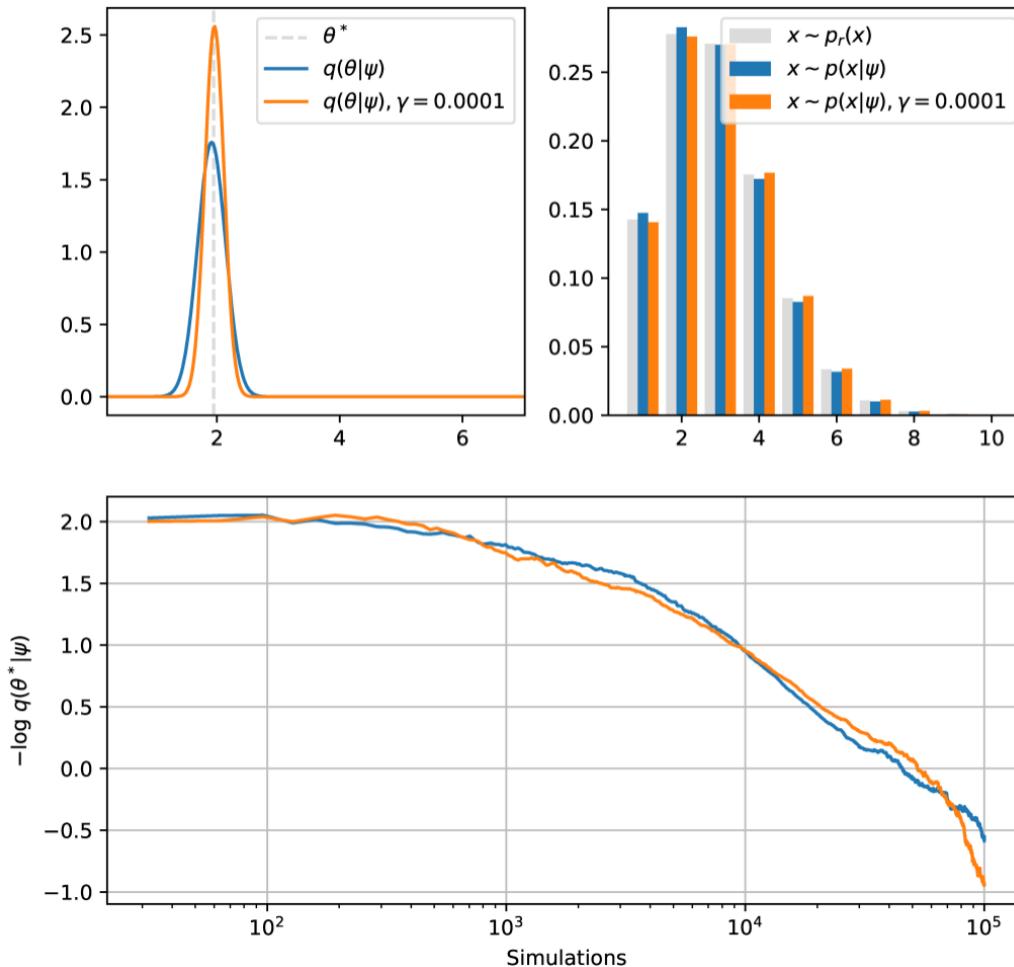
Weaknesses

AVO can been seen through the lens of empirical Bayes, where the data are used to optimize a prior within the family $q(\theta|\psi)$.

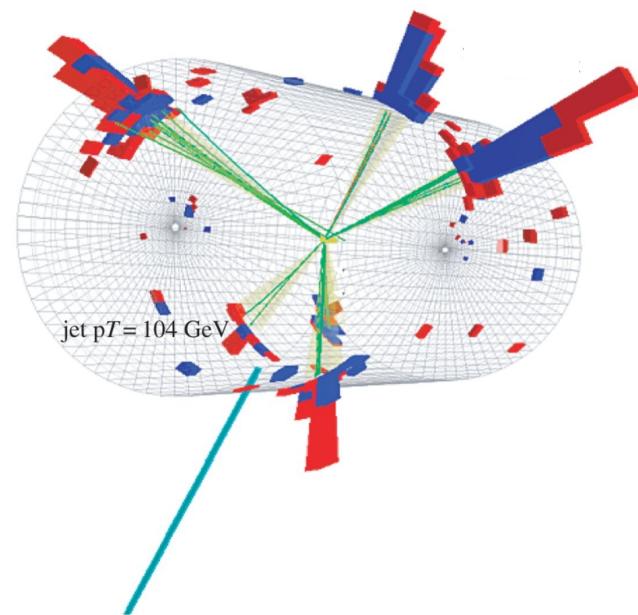
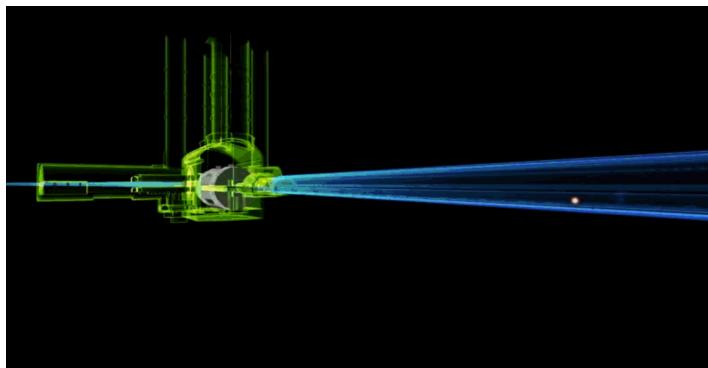
- Therefore, if the simulator is well specified, AVO can **only** be used to determine a point estimate of the MLE. It does not enable likelihood-free posterior inference.
- The solution may change from one run to the next in case of local minima.

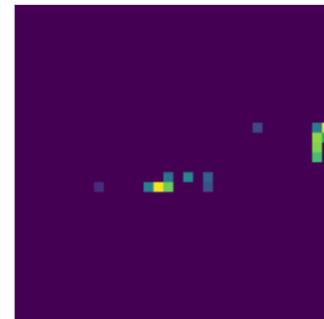
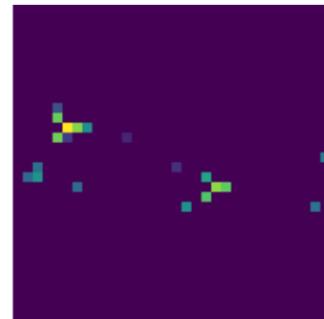
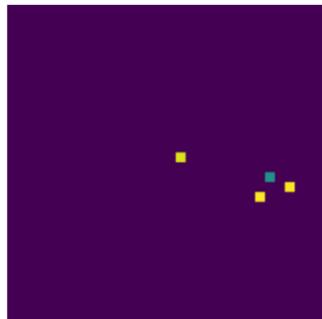
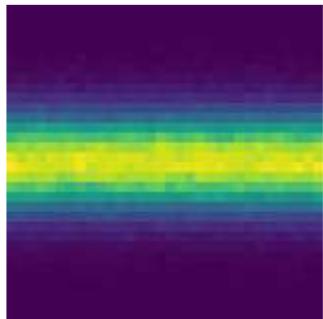
Experiments

Fitting a discrete Poisson distribution



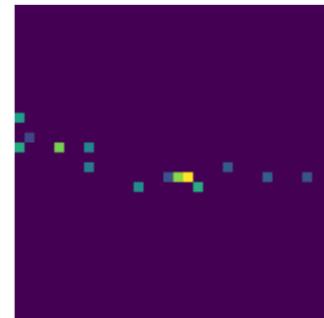
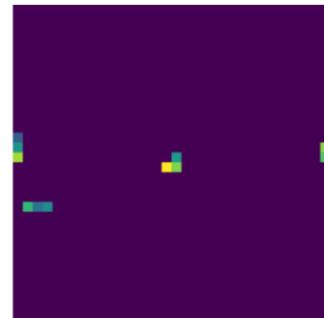
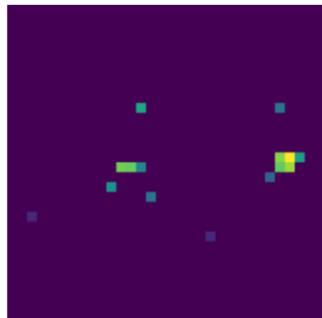
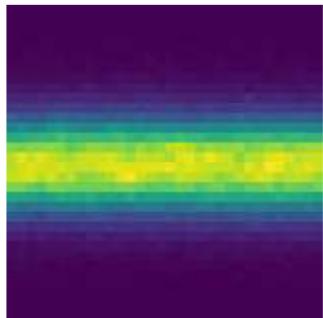
High-energy particle collisions



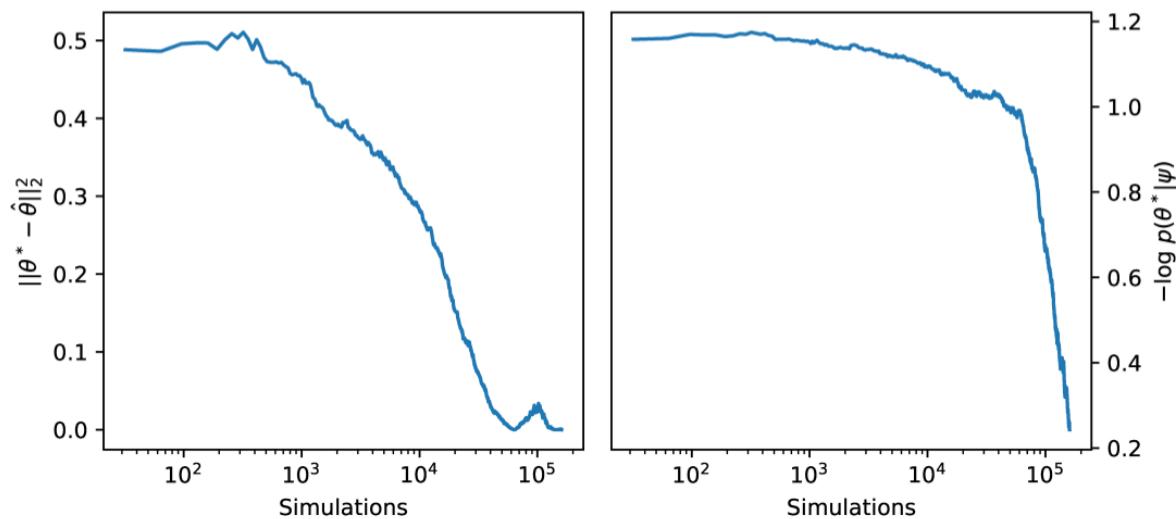
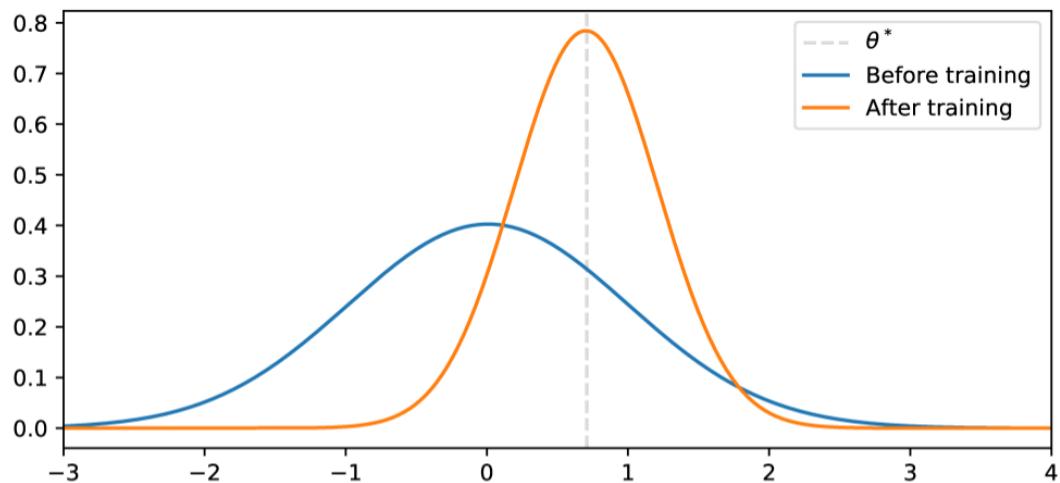


$z = 0$

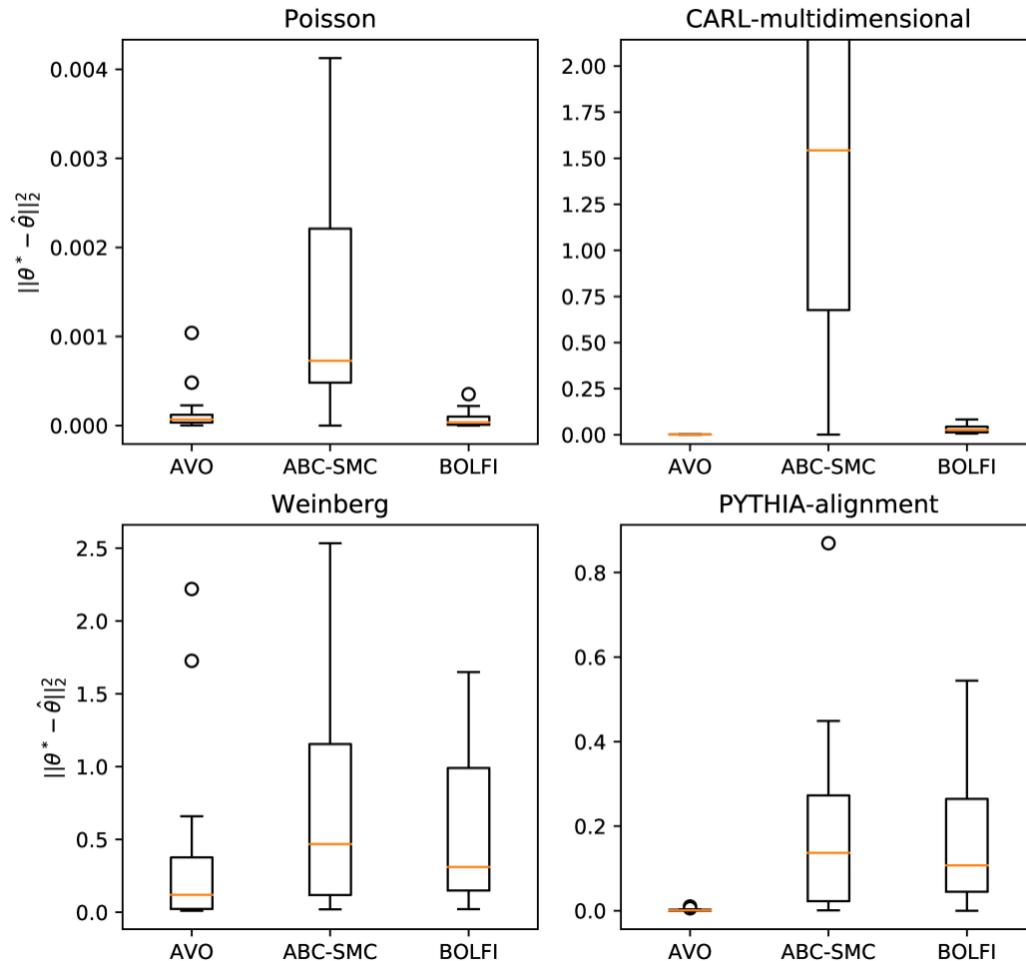
vs



$z = 1$



Benchmarks



Summary

In this work, 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 nor the disadvantages of likelihood-free inference algorithms that rely on ad hoc summary statistics. When the model is well-specified, the AVO algorithm provides point estimates for the generative model, which asymptotically corresponds to the data generating parameters. Experimental results highlight the good performance of AVO in comparison to the well-established ABC-SMC and BOLFI algorithms.

The end.