

---

# Adversarial Bayesian Simulation

---

*Veronika Ročková and Yuexi Wang*



Booth School of Business  
University of Chicago

# Bayesian Inference with Intractable Likelihoods

**Our framework:** data  $X^{(n)} = \{X_j\}_{j=1}^n$  realized from  $P_0 = P_{\theta_0}$  indexed by  $\theta_0 \in \Theta$  with a prior  $\pi(\theta)$ .

We assume that  $P_\theta$ , for each  $\theta \in \Theta$ , admits a density  $p_\theta$ .

We want to draw from the posterior

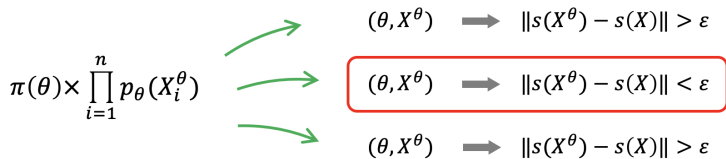
$$\pi_n(\theta | X^{(n)}) = \frac{p_\theta(X^{(n)})\pi(\theta)}{\int_{\Theta} p_\vartheta(X^{(n)}) d\Pi(\vartheta)}. \quad (1)$$

*Our focus is on situations where the likelihood  $p_\theta$  is too costly to evaluate ☹ but can be readily sampled from ☺.*

- ↪ Lotka-Volterra model: Population dynamics of animals in ecology
- ↪ Heston model: Stochastic volatility dynamics in finance
- ↪ Dynamic choice models: consumer dynamics in marketing

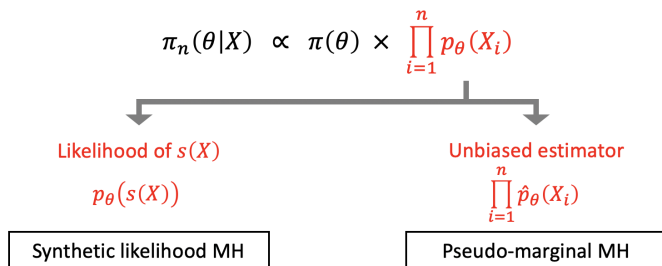
# Bayesian Inference with Intractable Likelihoods

**ABC:** A generator simulates fake data from  $p_\theta$  and reduces them to summary statistics



- ☹ *Reliance on summary statistics*
- ☹ *Only proper priors*
- ☹ *Not great for model comparisons*
- 😊 *Parallel computation feasible*

# Bayesian Inference with Intractable Likelihoods



**Bayesian Synthetic Likelihood:** Constructs a likelihood from summary statistics

☹ *Reliance on summary statistics*

**Pseudo-marginal MH:** Replace the likelihood in the MH ratio with an importance sampling estimate

☹ *Many simulation realizations*

☹ *May not yield the correct stationary distribution*

# Turning GANs into Posterior Simulators

*Generative Adversarial Networks* are a two-player minimax game

$$(g^*, d^*) = \arg \min_g \max_d [E_{X \sim P_0} \log d(X) + E_{Z \sim \pi_Z} \log(1 - d(g(Z)))] \quad (2)$$



$g^*$  minimizes the *Jensen-Shannon divergence*

$$JS(P_0, P_g), \quad \text{where } g(Z) \sim P_g \quad \text{and} \quad Z \sim \pi_Z.$$

*Wasserstein GANs* instead minimize

$$d_W(P_0, P_g) = \sup_{f \in \mathcal{F}_W} |E_{X \sim P_0} f(X) - E_{X \sim P_g} f(X)| \quad \text{where } \mathcal{F}_W = \{f : \|f\|_L \leq 1\}.$$

In practice, one replaces expectations with averages and restricts  $g$  and  $d$  to neural networks.

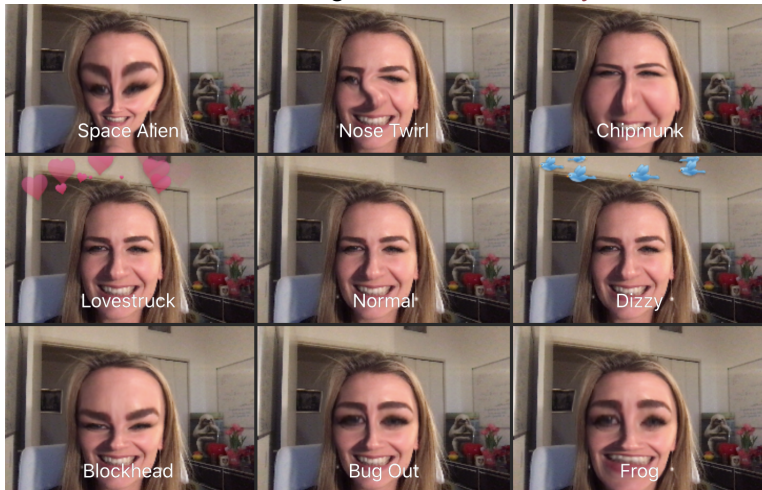
# Generative Adversarial Networks

Generator against the Adversary



# Generative Adversarial Networks

## Generator against the Adversary



# The Generator

Generate latent data  $Z \sim \pi_Z$  for some distribution  $\pi_Z$ .

Filter  $Z$  through a **deterministic mapping**  $g_\beta$  such that

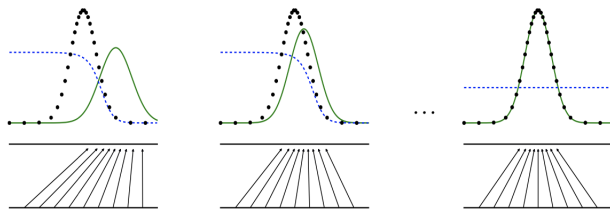
$$g_\beta(Z) \sim P_\theta.$$

Based on the feedback from the Classifier,  $\beta$  is updated so that  $P_\theta$  is closer and closer to  $P_0$ , where the game solution satisfies

$$g^*(Z) \sim P_0.$$

You can think of  $\tilde{X}_i^\theta = g_\beta(Z_i) \sim P_\theta$  as the **'fake' data**.

*Classifier against the Generator*





# The Classifier

The classification problem defined through

$$\max_{d \in \mathcal{D}} \left[ \frac{1}{n} \sum_{i=1}^n \log d(X_i) + \frac{1}{m} \sum_{i=1}^m \log(1 - d(\tilde{X}_i^\theta)) \right], \quad (3)$$

where  $d : \mathcal{X} \rightarrow (0, 1)$  (1 for 'real' and 0 for 'fake' data)

Oracle discriminator

$$d_\theta^*(X) := \frac{p_{\theta_0}(X)}{p_{\theta_0}(X) + p_\theta(X)}.$$

The optimal Generator leaves  $d_\theta^*$  maximally **confused** (assigning score 1/2), which occurs when  $p_\theta = p_{\theta_0}$ .



# Contrastive Learning for Bayesian Simulation

For iid data ( $p_\theta^{(n)} = \prod_i p_\theta(X_i)$ ), we can rewrite the likelihood as

$$p_\theta^{(n)} = p_0^{(n)} \times \exp\left(\sum_{i=1}^n \log \frac{1 - d_\theta^*(X_i)}{d_\theta^*(X_i)}\right).$$

---

**(1) Likelihood estimator?** Deploy  $\widehat{d}_{n,m}(\cdot)$  (e.g. neural network)

$$\widehat{p}_\theta^{(n)} = p_0^{(n)} \times \exp\left(\sum_{i=1}^n \log \frac{1 - \widehat{d}_{n,m}^\theta(X_i)}{\widehat{d}_{n,m}^\theta(X_i)}\right).$$

Kaji and Rockova (2021): *Metropolis Hastings via Classification*

---

**(2) KL estimator?**

$$\widehat{K}(\mathbf{X}, \tilde{\mathbf{X}}^\theta) = P_n \log \frac{\widehat{d}_{n,m}^\theta}{1 - \widehat{d}_{n,m}^\theta} = \frac{1}{n} \sum_{i=1}^n \log \frac{\widehat{d}_{n,m}^\theta(X_i)}{1 - \widehat{d}_{n,m}^\theta(X_i)}. \quad (4)$$

Wang, Kaji and Rockova (2022): *ABC via Classification*

---

☹ *Both require iid data and to run classification at every iteration!*

# Conditional GANs

GANs can be trained to simulate from *conditional distributions*.

Consider a two-player minimax game

$$(g^*, d^*) = \arg \min_{g \in \mathcal{G}} \max_{d \in \mathcal{D}} D(g, d)$$

prescribed by

$$D(g, d) = E_{(X, \theta) \sim \pi(X, \theta)} \log d(X, \theta) + E_{X \sim \pi(X), Z \sim \pi_Z} \log[1 - d(X, g(Z, X))].$$

Uniformly on  $\mathcal{X}$  and  $\Theta$  (for ‘flexible’  $\mathcal{G}$  and  $\mathcal{D}$ ), the solution  $(g^*, d^*)$  satisfies

$$\pi_{g^*}(\theta | X) = \frac{\pi(X, \theta)}{\pi(X)} = \pi(\theta | X)$$

and

$$d_g^*(X, \theta) = \frac{\pi(X, \theta)}{\pi(X, \theta) + \pi_g(\theta | X)\pi(X)} \quad \text{for any } g \in \mathcal{G}.$$

# Bayesian GANs

The **B-GAN** Algorithm:

**Simulate** the ABC reference table  $\{(\theta_i, X_i)\}_{i=1}^T$  from

$$\pi(\theta, X) = \pi(\theta) p_{\theta}^{(n)}(X^{(n)}) \quad \text{and} \quad \{Z_j\}_{j=1}^T \stackrel{\text{iid}}{\sim} \pi_Z(\cdot).$$

**Choose** function classes  $\mathcal{G}$  and  $\mathcal{D}$  (e.g. neural networks parametrized by  $\beta$  and  $\omega$ ).

**Train** the empirical version of Wasserstein conditional GANs

$$\hat{\beta}_T = \arg \min_{\beta: g_{\beta} \in \mathcal{G}} \left[ \max_{\omega: f_{\omega} \in \mathcal{F}_W} \left| \sum_{j=1}^T f_{\omega}(X_j, g_{\beta}(Z_j, X_j)) - \sum_{j=1}^T f_{\omega}(X_j, \theta_j) \right| \right]. \quad (5)$$

**Simulate**

$$\tilde{\theta}_j = g_{\hat{\beta}_T}(Z_j, X_0) \quad \text{for} \quad Z_j \stackrel{\text{iid}}{\sim} \pi_Z. \quad (6)$$

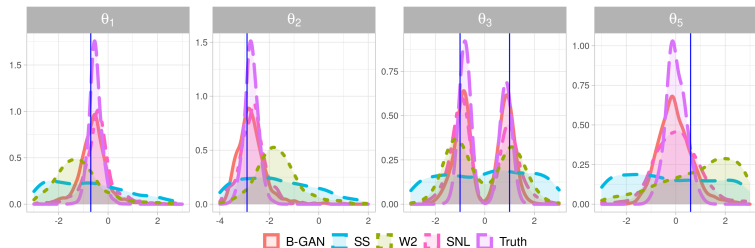
Observed data  $X_0$  at evaluation stage, not training stage!

# Toy Example

$X = (x_1, x_2, x_3, x_4)'$  consists of  $n = 4$  two-dimensional Gaussian observations with  $x_j \sim \mathcal{N}(\mu_\theta, \Sigma_\theta)$  parametrized by  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4, \theta_5)'$ , where

$$\mu_\theta = (\theta_1, \theta_2)' \quad \text{and} \quad \Sigma_\theta = \begin{pmatrix} s_1^2 & \rho s_1 s_2 \\ \rho s_1 s_2 & s_2^2 \end{pmatrix}$$

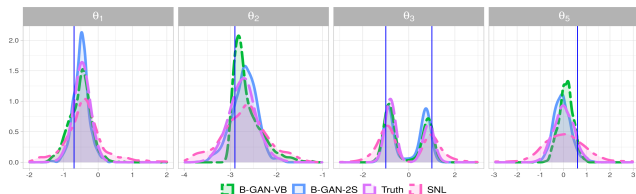
with  $s_1 = \theta_3^2$ ,  $s_2 = \theta_4^2$  and  $\rho = \tanh(\theta_5)$ .



$T = 100K$ , batchsize for SGD is 6400

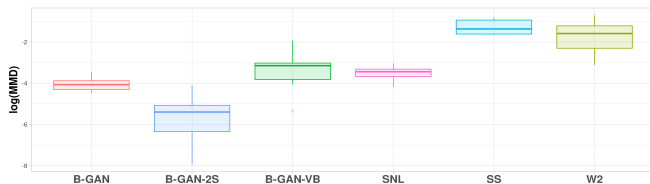
# Sequential Refinement

- ☹ **B-GAN is not trained on observed data!**
- ☹ We want  $\pi(\theta | X)$  *at observed data*  $X = X_0$ , not at any  $X$ .
- ☹ The ABC reference table  $\{(\theta_j, X_j)\}_{j=1}^T$  may not contain enough data points  $X_j$  in the vicinity of  $X_0$  to train the simulator when the prior  $\pi(\theta)$  is too vague.
- ☺ Use pilot simulator  $g_{\beta_T}(Z, X_0)$  in (6) obtained under the original prior  $\pi(\theta)$  as a proposal for the next round
- ☺ The ‘wrong’ prior can be corrected for by importance re-weighting with weights  $r(\theta) = \pi(\theta)/\tilde{\pi}(\theta)$ .



# Toy Example

## Performance summary



**Figure:** Maximum Mean Discrepancies (MMD, log scale) between the true posteriors and the approximated posteriors. The box-plots are computed from 10 repetitions.

# TV Bounds: The Three Terms

(1) The ability of the *critic* to tell the true model apart from the approximating model

$$\mathcal{A}_1(\mathcal{F}, \widehat{\beta}_T) \equiv \inf_{\omega: f_\omega \in \mathcal{F}} \left\| \log \frac{\pi(\theta | X)}{\pi_{\widehat{\beta}_T}(\theta | X)} - f_\omega(X, \theta) \right\|_\infty \quad (7)$$

(2) The ability of the *generator* to approximate the average true posterior

$$\mathcal{A}_2(\mathcal{G}) \equiv \inf_{\beta: g_\beta \in \mathcal{G}} \left[ E_X \left\| \log \frac{\pi_\beta(\theta | X)}{\pi(\theta | X)} \right\|_\infty \right]^{1/2}, \quad (8)$$

(3) The *complexity* of the (generating and) critic function classes measured by the pseudo-dimension  $Pdim(\cdot)$ .

We denote with  $\mathcal{H} = \{h_{\omega, \beta} : h_{\omega, \beta}(Z, X) = f_\omega(g_\beta(Z, X), X)\}$  a structured composition of networks  $f_\omega \in \mathcal{F}$  and  $g_\beta \in \mathcal{G}$ .



# TV Bounds

Let  $\widehat{\beta}_T$  be as in (5) where  $\mathcal{F} = \{f : \|f\|_\infty \leq B\}$  for some  $B > 0$ .

Denote with  $E$  the expectation with respect to  $\{(X_j, \theta_j)\}_{j=1}^T \stackrel{\text{iid}}{\sim} \pi(X, \theta)$  and  $\{Z_j\}_{j=1}^T \stackrel{\text{iid}}{\sim} \pi_Z$  in the reference table.

**Prior Concentration:** Assume

$$\Pi[B_n(\theta_0; \epsilon)] \geq e^{-C_2 n \epsilon^2} \quad \text{for some } C_2 > 0 \text{ and } \epsilon > 0. \quad (9)$$

Then for  $T \geq Pdim(\mathcal{F}) \vee Pdim(\mathcal{H})$  we have for any  $C > 0$

$$P_{\theta_0}^{(n)} E d_{TV}^2(\pi(\theta | X_0), \pi_{\widehat{\beta}_T}(\theta | X_0)) \leq C_n^T(\widehat{\beta}_T, \epsilon, C),$$

where, for some  $\widetilde{C} > 0$  and  $Pmax \equiv Pdim(\mathcal{F}) \vee Pdim(\mathcal{H})$ ,

$$C_n^T(\widehat{\beta}_T, \epsilon, C) = \frac{1}{C^2 n \epsilon^2} + \frac{e^{(1+C_2+C)n\epsilon^2}}{4} \left[ 2\mathcal{A}_1(\mathcal{F}, \widehat{\beta}_T) + \frac{B\mathcal{A}_2(\mathcal{G})}{\sqrt{2}} + 4\widetilde{C}B\sqrt{\frac{\log T \times Pmax}{T}} \right].$$

# Implicit Variational Bayes

The goal of VB is to find a set of parameters  $\beta^*$  that maximize ELBO

$$\log \pi(X_0) \geq \mathcal{L}(\beta) \equiv \int \log \left( \frac{\pi(X_0, \theta)}{q_\beta(\theta | X_0)} \right) q_\beta(\theta | X_0) d\theta. \quad (10)$$

The tightness increases with expressiveness of  $q_\beta(\cdot)$ , where the equality occurs when  $q_\beta(\theta | X_0) = \pi(\theta | X_0)$ .

Implicit VB defines  $q_\beta(\theta | X_0)$  through a **push-forward mapping**  $g_\beta$ .

We can re-write the ELBO in terms of Kullback-Leibler discrepancy

$$\mathcal{L}(\beta) = -\text{KL}(q_\beta(\theta | X_0) | \pi(\theta | X_0)) + C$$

- ☹ We cannot evaluate the *conditional* density ratio in the ELBO
- ☺ We can estimate the ratio of *joint* distributions with a different conditional, given  $X$ , but the same marginal  $\pi(X)$ .

# Adversarial Variational Bayes

☺ Joint LRT trick: define

$$\frac{d_{g_\beta}^*(X, \theta)}{1 - d_{g_\beta}^*(X, \theta)} = \frac{\pi(X, \theta)}{q_\beta(\theta | X)\pi(X)}, \quad (11)$$

where  $d_{g_\beta}^* : (\mathcal{X} \times \Theta) \rightarrow (0, 1)$

The variational lower bound (10) can be re-written as

$$\mathcal{L}(\beta) \equiv E_{\theta \sim q_\beta(\theta | X_0)} \left[ \text{logit}(d_{g_\beta}^*(X_0, \theta)) \right] + C. \quad (12)$$

Note that  $d_{g_\beta}^*(\theta, X)$  is a solution to

$$d_{g_\beta}^*(\theta, X) = \arg \max_{d \in \mathcal{D}} D(g_\beta, d). \quad (13)$$

Adversarial VB is a max-max game!

Given  $\beta^{(t)}$ : find  $\psi^{(t+1)}$  such that

$$\psi^{(t+1)} = \arg \max_{\psi} D(g_{\beta^{(t)}}, d_\psi).$$

Given  $\psi^{(t+1)}$ : find  $\beta^{(t+1)}$

$$\beta^{(t+1)} = \arg \max_{\beta} E_{\theta \sim q_\beta(\theta | X_0)} \left[ \text{logit}(d_{\psi^{(t+1)}}(\theta, X_0)) \right]$$

## ...and there are Wasserstein versions

Instead of KL, we can minimize Wasserstein distance between  $\pi(\theta | X_0)$  and  $q_\beta(\theta | X_0)$ :

$$\beta^* = \arg \min_{\beta: g_\beta \in \mathcal{G}} \sup_{f_\omega \in \mathcal{F}_W} \left| E_{\theta \sim q_\beta(\theta | X_0)} \left( \frac{\pi(\theta | X_0)}{q_\beta(\theta | X_0)} - 1 \right) f_\omega(\theta) \right|, \quad (14)$$

Using the ABC reference table  $\{(\theta_j, X_j)\}_{j=1}^T \stackrel{\text{iid}}{\sim} \pi(\theta, X)$ ,  $\{Z_j\}_{j=1}^T \stackrel{\text{iid}}{\sim} \pi_Z(\cdot)$ ,

- update  $\omega^{(t+1)}$ , given  $\beta^{(t)}$ ,

$$\omega^{(t+1)} = \arg \max_{\omega: f_\omega \in \mathcal{F}} \left[ \sum_{j=1}^T f_\omega(X_j, g_{\beta^{(t)}}(Z_j, X_j)) - \sum_{j=1}^T f_\omega(X_j, \theta_j) \right] \quad (15)$$

- update  $\beta^{(t+1)}$ , given  $\omega^{(t+1)}$ ,

$$\beta^{(t+1)} = \arg \min_{\beta: g_\beta \in \mathcal{G}} \left[ \sum_{j=1}^T f_{\omega^{(t+1)}}(X_0, g_\beta(Z_j, X_0)) + C \right], \quad (16)$$

where  $C$  does not depend on  $\beta$ , given the most recent update  $\omega^{(t+1)}$ .

# Lotka-Volterra Model

The Lotka-Volterra (LV) model describes population evolutions in ecosystems where **predators** interact with **prey**.

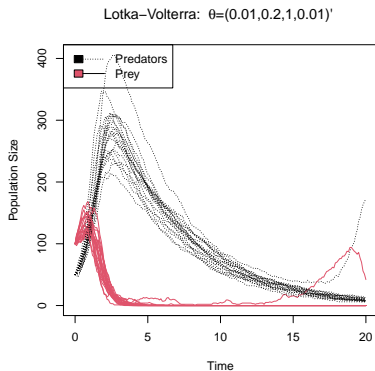
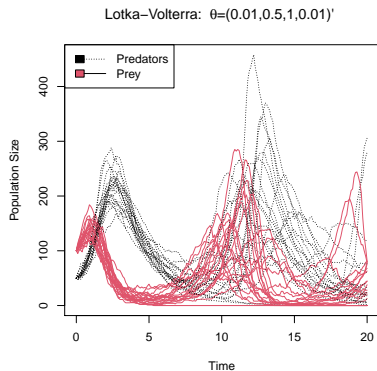
The model is deterministically prescribed via a system of first-order non-linear ODEs with four parameters  $\theta = (\theta_1, \theta_2, \theta_3, \theta_4)'$  controlling

- (1) the rate  $r_1^t = \theta_1 X_t Y_t$  of a predator being born,
- (2) the rate  $r_2^t = \theta_2 X_t$  of a predator dying,
- (3) the rate  $r_3^t = \theta_3 Y_t$  a prey being born and
- (4) the rate  $r_4^t = \theta_4 X_t Y_t$  a prey dying.

Despite **easy to sample from** (using the Gillespie algorithm), the likelihood for this model is **unavailable** which makes this model a natural candidate for ABC

*The pseudo-marginal approach far from straightforward, if at all possible 😊*

# Lotka-Volterra: A Closer Look

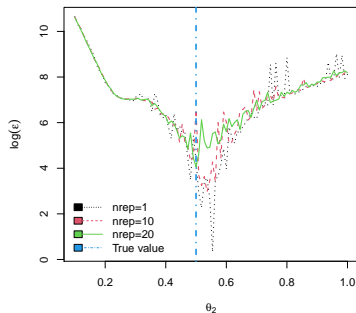


Simulation is started at  $X_0 = 50$  and  $Y_0 = 100$  simulated over 20 time units and recorded observations every 0.1 time units, resulting in a series of  $T = 201$  observations each.

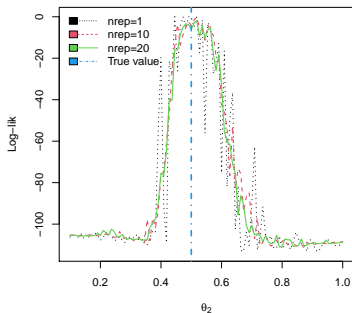
True values  $\theta^0 = (0.01, 0.5, 1, 0.01)$

# Prepping for ABC

ABC Discrepancy for  $\theta=(0.01, \dots, 1, 0.01)'$



Lotka-Volterra:  $\theta=(0.01, \dots, 1, 0.01)'$



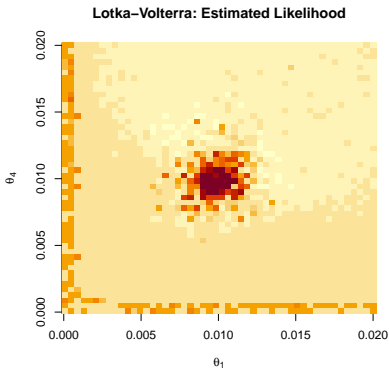
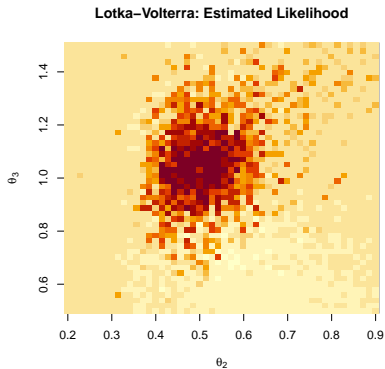
**LEFT:**

ABC tolerance (based on summary statistics)

**RIGHT:**

Classification-based log-lik estimator running LASSO (`glmnet`)

# Likelihood is Spiky!

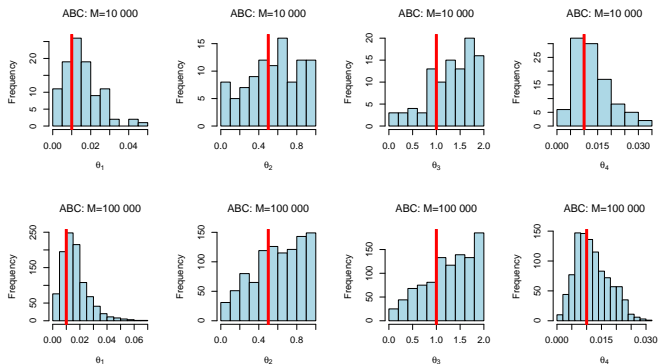


↷ ABC will need a *very* informative prior



# ABC Results

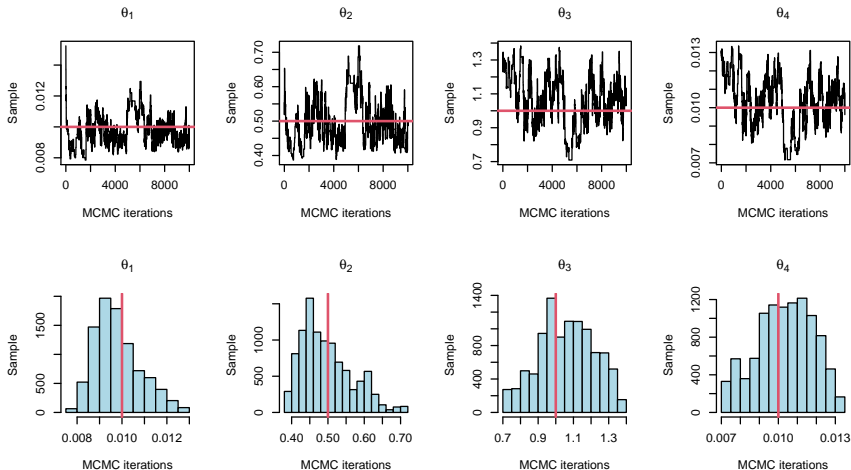
**Uniform Prior:** on  $[0, 0.1] \times [0, 1] \times [0, 2] \times [0, 0.1]$



*Upper panel:*  $M = 10\,000$  and  $r = 100$   
*Lower panel:*  $M = 100\,000$  and  $r = 1\,000$ .

# MHC (Kaji and Rockova (2021) Results

Initialized at posterior mean from a pilot ABC run.



MCMC trace plots (with  $M = 10\,000$ ) and histograms (without a burnin 1 000)

# MHC: Posterior Summary Statistics

Method	$\theta_1^0 = 0.01$			$\theta_2^0 = 0.5$			$\theta_3 = 1$			$\theta_4 = 0.01$		
	$\bar{\theta}$	$l$	$u$	$\bar{\theta}$	$l$	$u$	$\bar{\theta}$	$l$	$u$	$\bar{\theta}$	$l$	$u$
ABC1	0.015	0.003	0.038	0.554	0.037	0.985	1.315	0.189	1.955	0.012	0.004	0.029
ABC2	0.016	0.003	0.042	0.604	0.087	0.980	1.259	0.205	1.971	0.013	0.003	0.024
MHC	0.01	0.008	0.014	0.531	0.41	0.685	1.029	0.791	1.301	0.010	0.007	0.014

**ABC1:**  $M = 10\,000$  and  $r = 100$  (accepted samples)

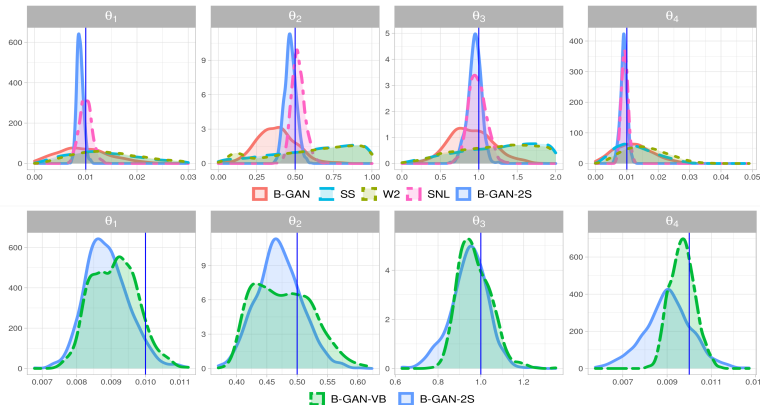
**ABC2:**  $M = 100\,000$  and  $r = 1\,000$  (accepted samples)

**MHC:**  $M = 10\,000$  with burnin 1 000

$\bar{\theta}$  denotes posterior mean,  $l$  and  $u$  denote the lower and upper boundaries of 95% credible intervals.

# What about $n = 1$ ?

We compare B-GAN with Sequential Neural Likelihood (SNL), W2 ABC, and Summary Statistics ABC



Sequential refinement and VB refinement work well.

# B-GAN Performance

Summary statistics of the approximated posteriors (averaged over 10 repetitions).

(scale)	$\theta_1 = 0.01$		$\theta_2 = 0.5$		$\theta_3 = 1.0$		$\theta_4 = 0.01$	
	bias ( $\times 10^{-3}$ )	CI width ( $\times 10^{-2}$ )	bias ( $\times 10^{-1}$ )	CI width	bias	CI width	bias ( $\times 10^{-2}$ )	CI width ( $\times 10^{-2}$ )
B-GAN	4.15	1.89	1.09	0.45	0.24	1.00	0.49	2.18
B-GAN-2S	<b>0.70</b>	<b>0.21 (0.9)</b>	0.42	<b>0.10 (0.7)</b>	<b>0.11</b>	0.33 (0.9)	0.13	0.34 (0.8)
B-GAN-VB	1.02	0.25 (0.7)	<b>0.38</b>	0.11 (0.9)	<b>0.11</b>	<b>0.29 (0.8)</b>	<b>0.12</b>	<b>0.29 (0.7)</b>
SNL	1.05	0.44	0.45	0.17	0.13	0.48	0.15	0.52
SS	9.58	3.80	2.49	0.91	0.49	1.76	0.68	2.72
W2	10.99	4.02 (0.9)	2.42	0.84	0.47	1.73	0.79	2.82

Bold fonts mark the best model of each column. The coverage of the 95% credible intervals are 1 unless otherwise noted in the parentheses.

Wang, Y. and Rockova, V. (2022) *Adversarial Bayesian Simulation*

**Thank you!**