

# Approximate Bayesian computation via the energy statistic

**Hien D. Nguyen<sup>1</sup>**

<sup>1</sup>Senior Lecturer, La Trobe University, Melbourne Australia.  
(Email: [h.nguyen5@latrobe.edu.au](mailto:h.nguyen5@latrobe.edu.au); Website: [hiendn.github.io](https://hiendn.github.io))

One World ABC Seminar, 18 June, 2020



# Acknowledgements



From left to right: Julyan Arbel, Florence Forbes, Hongliang Lü  
(STATIFY Team, Inria Grenoble Rhône-Alpes)

# Outline

- A *very brief* introduction to approximate Bayesian computation
- The energy distance and its estimator: **the energy statistic**
- The **energy statistic ABC** and its properties
- Numerical illustrations and comparisons to other methods

## Preliminary setup

- Let  $\mathbf{X}_n = \{\mathbf{X}_i\}_{i=1}^n$  be an IID sample of  $n$  replicates of the random variable  $\mathbf{X} \in \mathbb{X} \subseteq \mathbb{R}^q$ , where:
  - The DGP of  $\mathbf{X}$  has parametric PDF  $f(\mathbf{x}|\boldsymbol{\theta})$ .
  - The parameter  $\boldsymbol{\theta} \in \mathbb{T}$  has prior PDF  $\pi(\boldsymbol{\theta})$ .
- Conditional on  $\boldsymbol{\theta}$ , the likelihood of  $\mathbf{X}_n$  is

$$f(\mathbf{x}_n|\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{x}_i|\boldsymbol{\theta}).$$

- We aim to compute the PDF of the posterior distribution:

$$\pi(\boldsymbol{\theta}|\mathbf{x}_n) = c^{-1}(\mathbf{x}_n) f(\mathbf{x}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}),$$

where  $\mathbf{x}_n$  is a realization of  $\mathbf{X}_n$ , and

$$c(\mathbf{x}_n) = \int_{\mathbb{T}} f(\mathbf{x}_n|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

## Simulation-based inference

- Suppose that  $f(\mathbf{x}|\boldsymbol{\theta})$  is infeasible to compute, but it is feasible to simulate size  $m$  IID samples  $\mathbf{Y}_m = \{\mathbf{Y}_i\}_{i=1}^m$  with PDF

$$f(\mathbf{y}_n|\boldsymbol{\theta}) = \prod_{i=1}^m f(\mathbf{y}_i|\boldsymbol{\theta}).$$

- Let  $\boldsymbol{\theta}_k$  be simulated from PDF  $\pi(\boldsymbol{\theta})$ , and let  $\mathbf{Y}_{m,k}$  be the corresponding sample simulated from PDF  $f(\mathbf{y}_n|\boldsymbol{\theta}_k)$  and generate  $N$  pairs  $\mathbf{Z}_{m,k} = (\mathbf{Y}_{m,k}, \boldsymbol{\theta}_k)$ :

$$\mathbf{Z}_N = \{\mathbf{Z}_{m,k}\}_{k=1}^N.$$

- We wish to approximate the posterior  $\pi(\boldsymbol{\theta}|\mathbf{x}_n)$  using  $\mathbf{Z}_N$ .

## Approximating the posterior

- Let  $\mathcal{D}(\mathbf{x}_n, \mathbf{y}_m) \geq 0$  be a *discrepancy function* that measures the difference between the distributions of  $\mathbf{x}_n$  and  $\mathbf{y}_m$ .
- Let  $w(d, \varepsilon) \geq 0$  be a *weight function*, decreasing in  $d \geq 0$  and calibrated by  $\varepsilon > 0$ .
- Approximate the likelihood  $f(\mathbf{x}_n | \boldsymbol{\theta})$  by

$$L_{m,\varepsilon}(\mathbf{x}_n | \boldsymbol{\theta}) = \int_{\mathbb{X}^m} w(\mathcal{D}(\mathbf{x}_n, \mathbf{y}_m), \varepsilon) f(\mathbf{y}_m | \boldsymbol{\theta}) d\mathbf{y}_m.$$

- In the language of Jiang et al. (2018), we have the **pseudo-posterior** PDF approximation

$$\pi_{m,\varepsilon}(\boldsymbol{\theta} | \mathbf{x}_n) = c_{m,\varepsilon}^{-1}(\mathbf{x}_n) \pi(\boldsymbol{\theta}) L_{m,\varepsilon}(\mathbf{x}_n | \boldsymbol{\theta}),$$

where  $c_{m,\varepsilon}(\mathbf{x}_n) = \int_{\mathbb{T}} \pi(\boldsymbol{\theta}) L_{m,\varepsilon}(\mathbf{x}_n | \boldsymbol{\theta}) d\boldsymbol{\theta}$ .

# The ABC algorithm

**Input:** Data  $\mathbf{x}_n$ ; Discrepancy function  $\mathcal{D}$ ; Weight function  $w$ ;  
Calibration parameter  $\varepsilon$ ; Sample size  $m$ .

**For**  $k \in \{1, \dots, N\}$ ;

Simulate  $\boldsymbol{\theta}_k$  from PDF  $\pi(\boldsymbol{\theta})$ ;

Simulate  $\mathbf{Y}_{m,k}$  from PDF  $f(\mathbf{y}_m | \boldsymbol{\theta}_k)$ ;

Compute discrepancy  $D_k = \mathcal{D}(\mathbf{x}_n, \mathbf{Y}_{m,k})$ ;

Put  $\mathbf{Z}_{m,k} = (\mathbf{Y}_{m,k}, \boldsymbol{\theta}_k)$  into  $\mathbf{Z}_N$  and  $D_k$  into  $\mathbf{D}_N = \{D_k\}_{k=1}^N$ .

**Output:**  $\mathbf{Z}_N$  and  $\mathbf{D}_N$ ; Approximate the pseudo-posterior distribution by the discrete measure

$$\Pi_{m,\varepsilon}^N(\boldsymbol{\theta} | \mathbf{x}_n) = \left[ \sum_{k=1}^N w(D_k, \varepsilon) \right]^{-1} \sum_{k=1}^N w(D_k, \varepsilon) \mathbb{I}\{\boldsymbol{\theta} = \boldsymbol{\theta}_k\}.$$

## The energy distance

- For  $\mathbf{x}, \mathbf{y} \in \mathbb{X}$ , let  $\delta(\mathbf{x}, \mathbf{y})$  be a **semi-metric of negative type**, in the sense of Sejdinovic et al. (2013).
- For independent  $\mathbf{X}, \mathbf{X}'$  and  $\mathbf{Y}, \mathbf{Y}'$  with probability measures  $\Pi_X$  and  $\Pi_Y$ , respectively, the **energy distance** (ED) based on  $\delta$  is  $\mathcal{E}_\delta^{1/2}$ , where:

$$\mathcal{E}_\delta(\Pi_X, \Pi_Y) = 2\mathbb{E}[\delta(\mathbf{X}, \mathbf{Y})] - \mathbb{E}[\delta(\mathbf{X}, \mathbf{X}')] - \mathbb{E}[\delta(\mathbf{Y}, \mathbf{Y}')].$$

- Szekely and Rizzo (2017) states that the ED satisfies:
  1.  $\mathcal{E}_\delta(\Pi_X, \Pi_Y) \geq 0$ ,
  2.  $\mathcal{E}_\delta(\Pi_X, \Pi_Y) = 0 \iff \Pi_X = \Pi_Y$ ,

under the condition that

$$\mathbb{E}[\delta(\mathbf{X}, \mathbf{X}')] + \mathbb{E}[\delta(\mathbf{Y}, \mathbf{Y}')] < \infty.$$

## The Euclidean case

- The original paper of Szekely and Rizzo (2004) studied the  $\delta(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$  case, and we denote the squared ED, in this case, by  $\mathcal{E}$ .
- Let  $\varphi_X$  and  $\varphi_Y$  be the characteristic functions of  $\Pi_X$  and  $\Pi_Y$ , respectively. Then, under the condition that  $\mathbb{E}\|\mathbf{X}\| + \mathbb{E}\|\mathbf{Y}\| < \infty$ ,  $\mathcal{E}$  has the closed form:

$$\mathcal{E}(\Pi_X, \Pi_Y) = \frac{\Gamma\left(\frac{d+1}{2}\right)}{\pi^{(d+1)/2}} \int_{\mathbb{R}^d} \frac{|\varphi_X(\mathbf{t}) - \varphi_Y(\mathbf{t})|^2}{\|\mathbf{t}\|_2^{d+1}} d\mathbf{t}.$$

- Via the characterization,  $\mathcal{E}^{1/2}$  is a metric over set space of probability measures with finite first moment.

## The energy statistic

- Given samples  $\mathbf{X}_n$  and  $\mathbf{Y}_m$ , from measures  $\Pi_X$  and  $\Pi_Y$ , one can estimate  $\mathcal{E}_\delta(\Pi_X, \Pi_Y)$  by the energy V-statistic (ES):

$$\begin{aligned}\mathcal{V}_\delta(\mathbf{X}_n, \mathbf{Y}_m) &= \frac{2}{mn} \sum_{i=1}^n \sum_{j=1}^m \delta(\mathbf{X}_i, \mathbf{Y}_j) \\ &\quad - \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \delta(\mathbf{X}_i, \mathbf{X}_j) \\ &\quad - \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \delta(\mathbf{Y}_i, \mathbf{Y}_j).\end{aligned}$$

- For brevity, we write  $\mathcal{V}_\delta = \mathcal{V}$  when  $\delta$  is the Euclidean norm, and note that  $\mathcal{V}^{1/2}$  is a metric over the set of finite discrete measures (Szekely and Rizzo, 2017).
- We call the use of  $\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m) = \mathcal{V}(\mathbf{X}_n, \mathbf{Y}_m)$  the **ES ABC algorithm**.

## Relationship to the MMD ABC

- In Park et al. (2016), Jiang et al. (2018), and Bernton et al. (2019), the maximum mean discrepancy (MMD) ABC algorithms were considered, where one uses

$$\begin{aligned}\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m) &= -\frac{2}{mn} \sum_{i=1}^n \sum_{j=1}^m \chi(\mathbf{X}_i, \mathbf{Y}_j) \\ &\quad + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \chi(\mathbf{X}_i, \mathbf{X}_j) \\ &\quad + \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m \chi(\mathbf{Y}_i, \mathbf{Y}_j),\end{aligned}$$

for some Mercer kernel  $\chi$ .

- Sejdinovic et al. (2013) established that Mercer kernels generative relationship between Mercer kernels and semi-metrics of negative type.

## A general asymptotic result

Let  $\mathbf{X}_n$  and  $\mathbf{Y}_m$  be IID samples with PDFs  $f(\mathbf{x}_n|\boldsymbol{\theta}_0)$  and  $f(\mathbf{y}_m|\boldsymbol{\theta})$ . Assume that  $\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m)$  converges to  $\mathcal{D}_\infty(\boldsymbol{\theta}_0, \boldsymbol{\theta})$ , almost surely as  $n, m(n) \rightarrow \infty$ . If  $w(d, \varepsilon)$  is piecewise continuous and bounded for all  $d, \varepsilon > 0$ , and if  $w(\cdot, \varepsilon)$  is continuous at  $\mathcal{D}_\infty(\boldsymbol{\theta}_0, \boldsymbol{\theta})$ , then

$$\pi_{m(n), \varepsilon}(\boldsymbol{\theta}|\mathbf{X}_n) \rightarrow \frac{\pi(\boldsymbol{\theta}) w(\mathcal{D}_\infty(\boldsymbol{\theta}_0, \boldsymbol{\theta}), \varepsilon)}{\int_{\mathbb{T}} \pi(\boldsymbol{\theta}) w(\mathcal{D}_\infty(\boldsymbol{\theta}_0, \boldsymbol{\theta}), \varepsilon) d\boldsymbol{\theta}},$$

almost surely, as  $n \rightarrow \infty$ .

- Jiang et al. (2018) proved the rejection case:

$$w(d, \varepsilon) = \mathbb{I}\{d < \varepsilon\}.$$

## Verification of the ES case

- A sufficient condition to verify the almost sure convergence of  $\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m)$  to  $\mathcal{D}_\infty(\boldsymbol{\theta}_0, \boldsymbol{\theta})$  is for  $\Pi_X$  and  $\Pi_Y$  to have finite second moments:

$$\mathbb{E}\left(\|\mathbf{X}\|_2^2\right) + \mathbb{E}\left(\|\mathbf{Y}\|_2^2\right) < \infty.$$

- In this case, we have:

$$\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m) \rightarrow \frac{\Gamma\left(\frac{d+1}{2}\right)}{\pi^{(d+1)/2}} \int_{\mathbb{R}^d} \frac{|\varphi(\mathbf{t}; \boldsymbol{\theta}_0) - \varphi(\mathbf{t}; \boldsymbol{\theta})|^2}{\|\mathbf{t}\|_2^{d+1}} d\mathbf{t},$$

almost surely, as  $n \rightarrow \infty$ , where  $\Pi_X$  and  $\Pi_Y$  are characterized by  $\varphi(\mathbf{t}; \boldsymbol{\theta}_0)$  and  $\varphi(\mathbf{t}; \boldsymbol{\theta})$ .

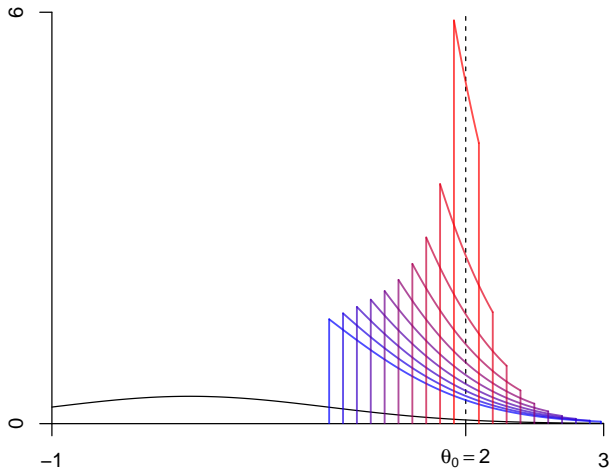
## A toy example

- Suppose that  $\theta \sim N(0, \tau^2)$
- We observe a sample  $\mathbf{X}_n$  of IID replications of  $X|\theta_0 \sim N(\theta_0, \sigma^2)$ , with  $\theta_0 = 2$ .
- The analytic solution for the posterior is that  $\theta|\mathbf{X}_n \sim N(\hat{\theta}, \hat{\sigma}^2)$ , where

$$\hat{\theta} = \frac{n\bar{X}_n}{n + \sigma^2/\tau^2}, \quad \hat{\sigma}^{-2} = \frac{1}{\tau^2} + \frac{n}{\sigma^2}.$$

- Note that  $\mathcal{D}_\infty(\theta_0, \theta) = (\theta_0 - \theta)^2$  in the ES case.
- We consider the ES ABC pseudo-posterior, when we use simulated samples  $\mathbf{Y}_m$  of IID replications of  $Y|\theta \sim N(\theta, \sigma^2)$ .

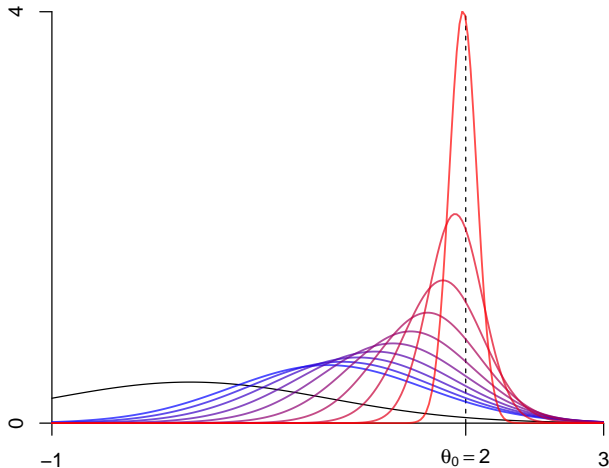
## A toy example



ES ABC pseudo-posterior limit using the rejection weights:

$$w(d, \varepsilon) = \mathbb{I}\{d < \varepsilon\}.$$

## A toy example



ES ABC pseudo-posterior limit using the Gaussian weights:

$$w(d, \epsilon) = \exp(-d^2/2\epsilon^2).$$

## A finite sample result

Let  $n = m$ , and  $w(d, \varepsilon) = \mathbb{I}\{d < \varepsilon\}$ . Assume that  $f(\mathbf{x}_n | \boldsymbol{\theta})$  is continuous and exchangeable, and that

$$\sup_{\boldsymbol{\theta} \in \mathbb{T} \setminus \{\Theta \subset \mathbb{T} : \int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0\}} f(\mathbf{x}_n | \boldsymbol{\theta}) < \infty,$$

and

$$\sup_{\boldsymbol{\theta} \in \mathbb{T} \setminus \{\Theta \subset \mathbb{T} : \int_{\Theta} \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} = 0\}} \sup_{\{\mathbf{y}_n : \mathcal{D}(\mathbf{x}_n, \mathbf{y}_n) \leq \bar{\varepsilon}\}} f(\mathbf{y}_n | \boldsymbol{\theta}) < \infty,$$

for some  $\bar{\varepsilon} > 0$ . Then, for fixed  $\mathbf{x}_n$ , the ES ABC pseudo-posterior PDF converges strongly to  $\pi(\boldsymbol{\theta} | \mathbf{x}_n)$ , as  $\varepsilon \rightarrow 0$ .

- This result along with others regarding convergence when  $\varepsilon \rightarrow 0$  and  $n \rightarrow \infty$ , simultaneously, can be taken directly from Bernton et al. (2019), due to the metric property of  $\gamma^{1/2}$ .

## Numerical setup

Set  $w(d, \varepsilon) = \mathbb{I}\{d < \varepsilon\}$  and  $\mathcal{D} = \mathcal{V}$ ;

Set  $N = 10^5$ ;

Set  $m = n$ ;

Set  $\varepsilon = \mathcal{Q}_N(0.05)$ , where  $\mathcal{Q}_N$  is the quantile empirical quantile function of  $\mathbf{D}_N$ .

For comparison, we also consider the following alternative for  $\mathcal{D}$ :

- the MMD discrepancy, with Gaussian kernel  $\chi$  (Park et al., 2016),
- the Kullback–Leibler (KL) discrepancy (Jiang et al., 2018),
- the “swapping distance”, approximation to the  $\mathcal{L}_2$  Wasserstein (WA) discrepancy (Bernton et al., 2019).

## Numerical study: Gaussian mixture

- We observe realization  $\mathbf{x}_n$  of  $\mathbf{X}_n$ ,  $n = 500$ , containing observations

$$\mathbf{X}_i \sim p^* \mathbf{N}(\boldsymbol{\mu}_0^*, \boldsymbol{\Sigma}_0) + (1 - p^*) \mathbf{N}(\boldsymbol{\mu}_1^*, \boldsymbol{\Sigma}_1),$$

where  $p^* = 0.3$ ,  $\boldsymbol{\mu}_0^* = (0.7, 0.7)$ ,  $\boldsymbol{\mu}_1^* = -\boldsymbol{\mu}_0^*$ ,

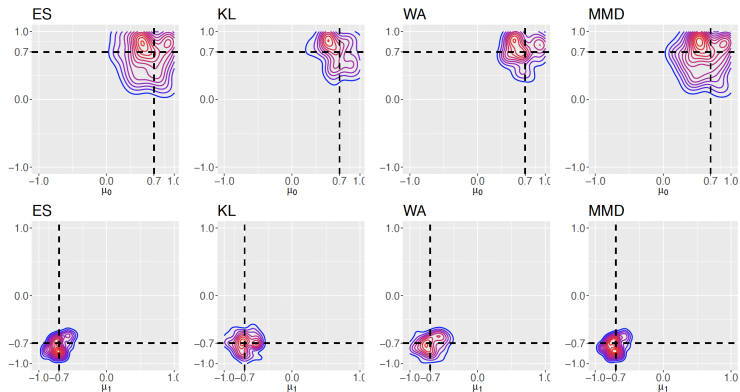
$$\boldsymbol{\Sigma}_0 = \begin{bmatrix} 0.5 & -0.3 \\ -0.3 & 0.5 \end{bmatrix}, \text{ and } \boldsymbol{\Sigma}_1 = \begin{bmatrix} 0.25 & 0 \\ 0 & 0.25 \end{bmatrix}.$$

- We wish to approximate  $\pi(\boldsymbol{\theta} | \mathbf{x}_n)$ , where  $\boldsymbol{\theta} = (p, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1)$ .

**Simulate**  $p_k \sim \text{U}(0, 1)$ ;  $\boldsymbol{\mu}_{0,k}, \boldsymbol{\mu}_{1,k} \sim \text{U}(-1, 1)^2$ ;

**Simulate**  $\mathbf{Y}_{n,k} \sim p_k \mathbf{N}(\boldsymbol{\mu}_{0,k}, \boldsymbol{\Sigma}_0) + (1 - p_k) \mathbf{N}(\boldsymbol{\mu}_{1,k}, \boldsymbol{\Sigma}_1)$ .

## Numerical study: Gaussian mixture



Kernel density estimates of the ABC-obtained posterior samples.

## Numerical study: moving average

- We observe realization  $\mathbf{x}_n$  of  $\mathbf{X}_n$ ,  $n = 200$ , containing observations  $\mathbf{X}_j \in \mathbb{R}^{10}$ , such that

$$X_{i,t} = Z_t + \sum_{j=1}^2 \theta_j^* Z_{t-j}, \quad (1)$$

where  $\{Z_t\}_{t \in \mathbb{Z}}$  is an IID sequence of Student- $t$  noise, with 5 degrees of freedom, and

$$\boldsymbol{\theta}^* = (\theta_1^*, \theta_2^*) = (0.6, 0.2).$$

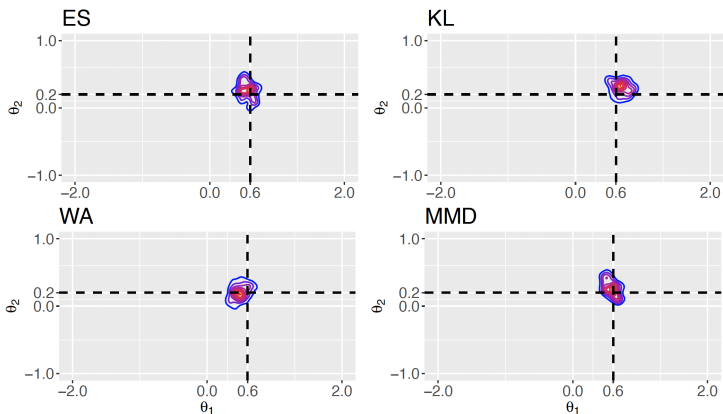
- We wish to approximate  $\pi(\boldsymbol{\theta} | \mathbf{x}_n)$ .

**Simulate**  $\theta_{1,k} \sim U(-2, 2)$ ;  $\theta_{2,k} \sim U(-1, 1)$ ;

**Simulate**  $\mathbf{Y}_{n,k} \in \mathbb{R}^{10}$  from (1), with  $\boldsymbol{\theta}^*$  replaced by

$$\boldsymbol{\theta}_k = (\theta_{1,k}, \theta_{2,k}).$$

## Numerical study: moving average



Kernel density estimates of the ABC-obtained posterior samples.

## Numerical study: bivariate beta distribution

- We observe realization  $\mathbf{x}_n$  of  $\mathbf{X}_n$ ,  $n = 500$ , containing observations  $\mathbf{X}_i \in \mathbb{R}^2$ , such that

$$X_{i,1} \sim \text{Beta}(\theta_1^* + \theta_3^*, \theta_4^* + \theta_5^*), \quad X_{i,2} \sim \text{Beta}(\theta_2^* + \theta_4^*, \theta_3^* + \theta_5^*), \quad (2)$$

where  $\boldsymbol{\theta}^* = (\theta_1^*, \dots, \theta_5^*) = \mathbf{1}$ .

- We wish to approximate  $\pi(\boldsymbol{\theta} | \mathbf{x}_n)$ .

**Simulate**  $\boldsymbol{\theta}_k = (\theta_{k,1}, \dots, \theta_{k,5}) \sim U(0,5)^5$ ;

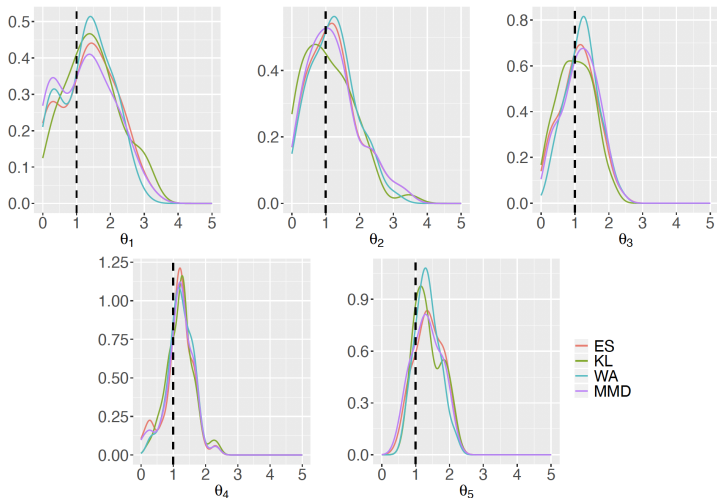
**Simulate**  $\mathbf{Y}_{n,k} \in \mathbb{R}^2$  from (2), with  $\boldsymbol{\theta}^*$  replaced by  $\boldsymbol{\theta}_k$ , by simulating  $U_{k,j} \sim \text{Gamma}(\theta_{k,j}, 1)$  and setting

$$Y_{n,k,1} = V_{k,1} / (1 + V_{k,1}), \quad \text{and} \quad Y_{n,k,2} = V_{k,2} / (1 + V_{k,2})$$

where

$$V_{k,1} = \frac{U_{k,1} + U_{k,3}}{U_{k,4} + U_{k,5}}, \quad \text{and} \quad V_{k,2} = \frac{U_{k,2} + U_{k,4}}{U_{k,3} + U_{k,5}}.$$

## Numerical study: bivariate beta distribution

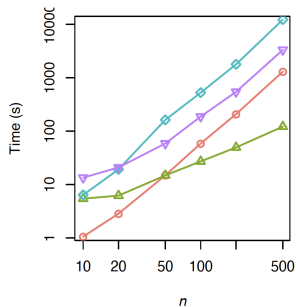
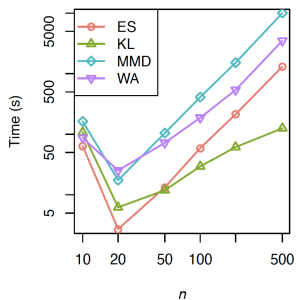


Kernel density estimates of the ABC-obtained posterior samples.

## Computational complexities

Discrepancy $\mathcal{D}$	Complexity
All methods (univariate)	$O((n+m)\log(n+m))$
KL (multi.)	$O((n+m)\log(n+m))$
ES/MMD, approx. WA (multi.)	$O((n+m)^2)$
WA (multi.)	$O((n+m)^{5/2}\log(n+m))$

# Timing simulations



**Left:** results for Gaussian mixture study. **Right:** results for moving average study.

## Why not KL?

Jiang et al. (2018) proposed to use the discrepancy

$$\mathcal{D}(\mathbf{x}_n, \mathbf{y}_m) = \frac{q}{n} \sum_{i=1}^n \log \left( \frac{\min_{j \in \{1, \dots, n\}} \|\mathbf{x}_i - \mathbf{y}_j\|_2}{\min_{j \neq i} \|\mathbf{x}_i - \mathbf{x}_j\|_2} \right) + \log \frac{m}{n-1},$$

where  $\mathcal{D}(\mathbf{X}_n, \mathbf{Y}_m)$  converges, almost surely, to  $\text{KL}(\Pi_X \parallel \Pi_Y)$ .

- Cannot be used for discrete distributions, since  $\min_{j \neq i} \|\mathbf{x}_i - \mathbf{x}_j\|_2 = 0$  with non-zero probability.
- Cannot be used for quantized real data, for the same reason.

## A linear time estimator for ES/MMD

We can write

$$\mathcal{V}_\delta(\mathbf{x}_n, \mathbf{y}_m) = \sum_{i_1=1}^n \sum_{i_2=1}^n \sum_{j_1=1}^m \sum_{j_2=1}^m \frac{\kappa_\delta(\mathbf{x}_{i_1}, \mathbf{x}_{i_2}; \mathbf{y}_{j_1}, \mathbf{y}_{j_2})}{m^2 n^2}, \text{ where}$$

$$\begin{aligned} \kappa_\delta(\mathbf{x}_{i_1}, \mathbf{x}_{i_2}; \mathbf{y}_{j_1}, \mathbf{y}_{j_2}) &= \delta(\mathbf{x}_{i_1}, \mathbf{y}_{j_1}) + \delta(\mathbf{x}_{i_2}, \mathbf{y}_{j_2}) \\ &\quad - \delta(\mathbf{x}_{i_1}, \mathbf{y}_{j_2}) - \delta(\mathbf{x}_{i_2}, \mathbf{y}_{j_1}). \end{aligned}$$

- When  $m = n$ , Gretton et al. (2012) observed we can estimate  $\mathcal{E}_\delta$ , unbiasedly, by

$$\mathcal{U}_\delta(\mathbf{X}_n, \mathbf{Y}_m) = \lfloor n/2 \rfloor^{-1} \sum_{i=1}^{\lfloor n/2 \rfloor} \kappa_\delta(\mathbf{X}_{2i-1}, \mathbf{X}_{2i}; \mathbf{Y}_{2i-1}, \mathbf{Y}_{2i}).$$

- We can construct a biased estimator  $\mathcal{B}_\delta \geq 0$  that converges to  $\mathcal{E}_\delta$ , almost surely, by setting

$$\mathcal{B}_\delta(\mathbf{X}_n, \mathbf{Y}_m) = \max\{0, \mathcal{U}_\delta(\mathbf{X}_n, \mathbf{Y}_m)\}.$$

## References I

- Bernton, E., Jacob, P. E., Gerber, M., and Robert, C. P. (2019). Approximate Bayesian computation with the Wasserstein distance. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 81:235–269.
- Gretton, A., Bogwardt, K. M., Rasch, M. J., Scholkopf, B., and Smola, A. (2012). A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773.
- Jiang, B., Wu, T.-Y., and Wong, W. H. (2018). Approximate Bayesian computation with Kullback-Leibler divergence as data discrepancy. In *Conference on Artificial Intelligence and Statistics (AISTATS)*.

## References II

- Park, M., Jitkrittum, W., and Sejdinovic, D. (2016). K2-ABC: approximate Bayesian computation with kernel embeddings. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS)*.
- Sejdinovic, D., Sriperumbudur, B., Gretton, A., and Fukumizu, K. (2013). Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *Annals of Statistics*, 41:2263–2291.
- Szekely, G. J. and Rizzo, M. L. (2004). Testing for equal distributions in high dimension. *InterStat*, 5:1–16.
- Szekely, G. J. and Rizzo, M. L. (2017). The energy of data. *Annual Review of Statistics and Its Application*, 4(447-479).

Thank you!

Preprint: **<https://arxiv.org/abs/1905.05884>**

Email: **[h.nguyen5@latrobe.edu.au](mailto:h.nguyen5@latrobe.edu.au)**

Website: **<https://hiendn.github.io>**