

Distance-based ABC procedures

Hien Nguyen¹

¹La Trobe University

(Contact–Email: h.nguyen5@latrobe.edu.au; Website: hiendn.github.io)

ABC in Svalbard (in Melbourne)



Acknowledgements



(a) Florence Forbes.



(b) Julyan Arbel.



(c) Trung Tin Nguyen.

Figure: Highly esteemed collaborators.

Setup

- Let $(\Omega, \mathcal{F}, \Pr)$ be a sufficiently rich probability space and let λ denote the Lebesgue measure.
- Let $\Theta \in \mathbb{T}$ be a random parameter with **prior probability measure** defined by

$$d\Pi(\theta) = \pi(\theta) d\lambda(\theta).$$

- Let $\mathbf{X}_n = (X_i)_{i \in [n]}$ be a random sequence ($[n] = \{1, \dots, n\}$, $n \in \mathbb{N}$) , where each $X_i \in \mathbb{X}$, and \mathbf{X}_i , conditioned on $\Theta = \theta$ has probability measure defined by

$$dF(x_i | \theta) = f(x_i | \theta) d\lambda(x_i).$$

- The joint probability of the pair $W = (\mathbf{X}_n, \Theta)$ is thus given by the joint probability measure defined by

$$dF(x_i, \theta) = \pi(\theta) f(x_i | \theta) d\lambda(x_i, \theta).$$

Setup

- Suppose that we observe $\mathbf{X}_n = \mathbf{x}_n$ from a pair W , but not Θ , which we wish to estimate.
- Using the previous characterizations, we can draw inference regarding Θ , conditionally on $\mathbf{X}_n = \mathbf{x}_n$ by constructing the **posterior PDF**

$$\pi(\theta | \mathbf{x}_n) = \frac{\pi(\theta) f(\mathbf{x}_n | \theta)}{c(\mathbf{x}_n)},$$

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

- Often $f(\mathbf{x}_n | \theta)$ is not available in a tractable form, but it is still possible to simulate from $f(\mathbf{x}_n | \theta)$.

Distance-based ABC

- For two sequences $\mathbf{x}_n, \mathbf{y}_n \in \mathbb{X}^n$, define the **discrepancy**
 $D(\mathbf{x}_n, \mathbf{y}_n) \geq 0$.
- For $d \geq 0$ and $\varepsilon > 0$, define the **weight function** $w(d, \varepsilon)$,
where $w(\cdot, \varepsilon)$ is decreasing in d .
- From a simulation $\mathbf{W}_N = (W_j)_{j \in [N]}$, where $W_j = (\mathbf{Y}_{n,j}, \Theta_j)$
and $N \in \mathbb{N}$, from $f(\mathbf{y}_i, \theta)$, we can estimate the characteristics
of the so called **ABC pseudo-posterior PDF**

$$\pi_\varepsilon(\theta | \mathbf{x}_n) = \frac{\pi(\theta) L_\varepsilon(\mathbf{x}_n | \theta)}{c_\varepsilon(\mathbf{x}_n)},$$

where

$$L_\varepsilon(\mathbf{x}_n | \theta) = \int_{\mathbb{X}^n} w(D(\mathbf{x}_n, \mathbf{y}_n), \varepsilon) f(\mathbf{y}_n | \theta) d\lambda(\mathbf{y}_n)$$

and $c_\varepsilon(\mathbf{x}_n) = \int_{\mathbb{T}} \pi(\theta) L_\varepsilon(\mathbf{x}_n | \theta) d\lambda(\theta)$.

Discrepancies

- $D(\mathbf{x}_n, \mathbf{y}_n) = \|s(\mathbf{x}_n) - s(\mathbf{y}_n)\|_p$, $p \geq 1$, where $s(\cdot)$ is a **summary statistic**.
- Jiang et al. (2018) considered the **sample Kullback–Leibler divergence** of Perez-Cruz (2008), for $\mathbb{X} \subseteq \mathbb{R}^d$:

$$D(\mathbf{x}_n, \mathbf{y}_n) = \frac{d}{n} \sum_{i=1}^n \log \frac{\min_j \|\mathbf{x}_i - \mathbf{y}_j\|_2}{\min_{j \neq i} \|\mathbf{x}_i - \mathbf{y}_j\|_2} + \log \frac{n}{n-1}.$$

- Bernton et al. (2019) considered $D(\mathbf{x}_n, \mathbf{y}_n) = \mathfrak{W}_p(\mathbf{x}_n, \mathbf{y}_n)$, where \mathfrak{W}_p is the **p -Wasserstein distance** between multisets of n vectors.
- Nguyen et al. (2020) considered the **energy statistic** of Szekely and Rizzo (2004), for $\mathbb{X} \subseteq \mathbb{R}^d$:

$$D(\mathbf{x}_n, \mathbf{y}_n) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left\{ 2 \|\mathbf{x}_i - \mathbf{y}_j\|_2 - \|\mathbf{x}_i - \mathbf{x}_j\|_2 - \|\mathbf{y}_i - \mathbf{y}_j\|_2 \right\}.$$

Weight functions

- Typically we choose the weight function w to be

$$w(d, \varepsilon) = \mathbf{1}\{d \leq \varepsilon\},$$

where $\mathbf{1}\{\cdot\}$ is the indicator function.

- Here, the choice corresponds to the traditional **rejection-sampling** ABC algorithms.
- A more exotic choice is considered in [Park et al. \(2016\)](#) and [Miller and Dunson \(2019\)](#), where

$$w(d, \varepsilon) = \exp(-\varepsilon^{-1} d^p), p \geq 1.$$

- This corresponds to **importance sampling**-based ABC procedures, as considered in [Karabatsos and Leisen \(2018\)](#).

Behavior of the pseudo-posterior when $\varepsilon \rightarrow 0$

- A1 Use $w(d, \varepsilon) = \mathbf{1}\{d \leq \varepsilon\}$.
- A2 The PDF $f(\mathbf{x}_n | \theta)$ is continuous and uniformly bounded, for all $\theta \in \mathbb{T}$ and $\mathbf{x}_n \in \mathbb{X}^n$.
- A3 The discrepancy is continuous in the sense that $D(\mathbf{x}_n, \mathbf{y}_n) \rightarrow D(\mathbf{x}_n, \bar{\mathbf{y}}_n)$, when $\mathbf{y}_n \rightarrow \bar{\mathbf{y}}_n$.
- A4 $dF(\mathbf{x}_n | \theta)$ defines an exchangeable measure and $D(\mathbf{x}_n, \mathbf{y}_n) = 0$, if and only if $\{x_i\}_{i=1}^n = \{y_i\}_{i=1}^n$.

The following result is proved by Rubio and Johansen (2013) and Bernton et al. (2019).

Proposition 1

Under A1–A4, for fixed $\mathbf{x}_n \in \mathbb{X}_n$,

$$\int_{\mathbb{T}} |\pi_\varepsilon(\theta | \mathbf{x}_n) - \pi(\theta | \mathbf{x}_n)| d\lambda(\theta) \rightarrow 0, \text{ as } \varepsilon \rightarrow 0.$$

Behavior of the pseudo-posterior when $n \rightarrow \infty$

- B1 \mathbf{X}_n and \mathbf{Y}_n have probability measures defined by $dF(\mathbf{x}_n|\theta_0)$ and $dF(\mathbf{y}_n|\theta)$, respectively.
- B2 $D(\mathbf{X}_n, \mathbf{Y}_n) \rightarrow D(\theta_0, \theta)$, almost surely, as $n \rightarrow \infty$, where $D(\theta_0, \theta)$ is a function of $\theta_0, \theta \in \mathbb{T}$.
- B3 $w(d, \varepsilon)$ is piecewise continuous, and $w(d, \varepsilon) \leq a < \infty$.

The following is proved by [Jiang et al. \(2018\)](#), [Miller and Dunson \(2019\)](#) and [Nguyen et al. \(2020\)](#).

Theorem 1

Under B1–B3, for any $\varepsilon > 0$, if $w(\cdot, \varepsilon)$ is continuous at $D(\theta_0, \theta)$, then

$$\pi_\varepsilon(\theta | \mathbf{X}_n) \xrightarrow[n \rightarrow \infty]{a.s.} \frac{\pi(\theta) w(D(\theta_0, \theta), \varepsilon)}{\int_{\mathbb{T}} \pi(\theta) w(D(\theta_0, \theta), \varepsilon) d\lambda(\theta)} = \pi_\varepsilon(\theta; \theta_0).$$

Behavior of the pseudo-posterior when $n \rightarrow \infty$

- Let $Q(u; \omega)$ and $Q_n(u; \omega)$ ($n \in \mathbb{N}$) be probability measures with respect to $u \in (\mathbb{U}, \mathcal{U})$, for each $\omega \in \Omega$.
- By Berti et al. (2006) and Grubel and Kabluchko (2016), say that the sequence Q_n **almost surely weakly** (a.s.w) converges to Q , if on a set $A \in \mathcal{F}$ with $\Pr(A) = 1$, for all $\omega \in A$, $Q_n(\cdot; \omega)$ weakly converges on \mathbb{U} to $Q(\cdot; \omega)$, as $n \rightarrow \infty$.
- Consider the pseudo-posterior and its limiting measure, defined by

$$d\Pi_\varepsilon(\theta | \mathbf{X}_n) = \pi_\varepsilon(\theta | \mathbf{X}_n) d\lambda(\theta), \quad d\Pi_\varepsilon(\theta; \theta_0) = \pi_\varepsilon(\theta; \theta_0) d\lambda(\theta).$$

Corollary 1

Under the assumptions of Theorem 1, if $\mathbb{T} \subseteq \mathbb{R}^q$, $q \in \mathbb{N}$, then $\Pi_\varepsilon(\cdot | \mathbf{X}_n)$ a.s.w. converges on \mathbb{T} to $\Pi_\varepsilon(\cdot; \theta_0)$, as $n \rightarrow \infty$.

Behavior of the pseudo-posterior when $n \rightarrow \infty$

- C1 Use $w(d, \varepsilon) = \mathbf{1}\{d \leq \varepsilon\}$.
- C2 $D(\theta_0, \theta)$ is a metric on the set $\mathbb{T} \subseteq \mathbb{R}^q$, $q \in \mathbb{N}$.
- C3 The identity function $\text{id} : (\mathbb{T}, D) \rightarrow \left(\mathbb{T}, \|\cdot - \cdot\|_p\right)$ is uniformly continuous, for some $p \geq 1$ (satisfied if $\|\theta_0 - \theta\|_p \leq bD(\theta_0, \theta)$, for some finite constant $b > 0$).

Corollary 2

Assume C1–C3 and the assumptions of Corollary 1. For any $\iota > 0$, there exists an $\varepsilon > 0$ (independent of θ_0), and a set $\mathbb{S}_\iota(\theta_0) \subset \mathbb{T}$ containing θ_0 , with $\lambda(\mathbb{S}_\iota(\theta_0)) = \iota$, such that

$$\liminf_{n \rightarrow \infty} \Pi_\varepsilon(\mathbb{S}_\iota(\theta_0) | \mathbf{X}_n) = 1, \text{ a.s.}$$

Illustration

- \mathbf{X}_n and \mathbf{Y}_n are IID sequences, such that $X_i \sim N(\theta_0, \sigma^2)$ and $Y_i \sim N(\theta, \sigma^2)$, with $\sigma^2 > 0$ fixed.
- $\theta \sim N(0, \tau^2)$, with $\tau^2 > 0$ fixed.
- Use $w(d, \varepsilon) = \mathbf{1}\{d \leq \varepsilon\}$.
- Use $D(\mathbf{X}_n, \mathbf{Y}_n) = \left| n^{-1} \sum_{i=1}^n X_i - n^{-1} \sum_{i=1}^n Y_i \right| \xrightarrow[n \rightarrow \infty]{\text{a.s.}} |\theta_0 - \theta|$.

Illustration

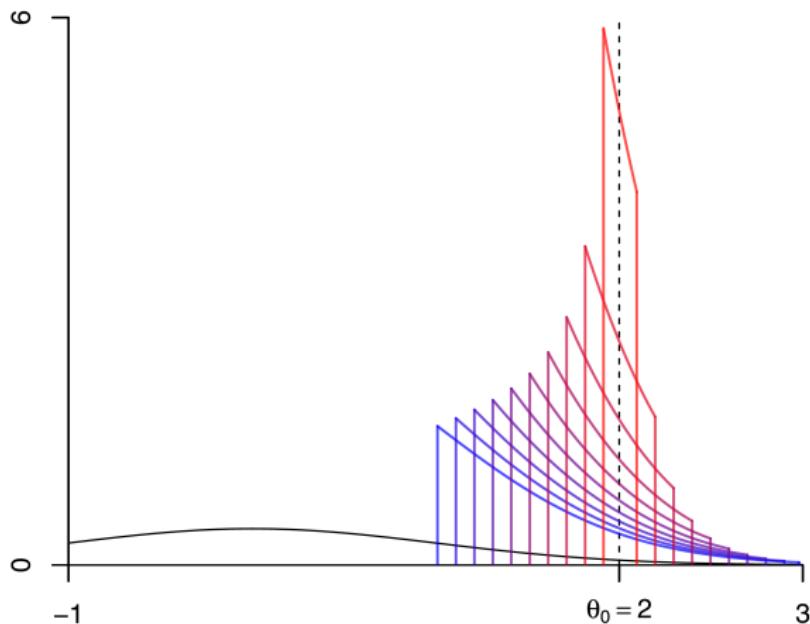


Figure: Limiting pseudo-posteriors $\pi_\varepsilon(\theta; \theta_0)$ with $\varepsilon \in [0.1, 1]$, when $\theta_0 = 2$. The black curve represents the prior density, and the red-blue curves represent the pseudo-posteriors for increasing ε .

The energy distance

- Recall that for $\mathbb{X} \subseteq \mathbb{R}^d$, we define the **energy statistic** between two samples \mathbf{x}_n and \mathbf{y}_n as:

$$D(\mathbf{x}_n, \mathbf{y}_n) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \left\{ 2 \|\mathbf{x}_i - \mathbf{y}_j\|_2 - \|\mathbf{x}_i - \mathbf{x}_j\|_2 - \|\mathbf{y}_i - \mathbf{y}_j\|_2 \right\}.$$

- From [Szekely and Rizzo \(2017\)](#), we have the fact $\sqrt{D(\mathbf{x}_n, \mathbf{y}_n)}$ is a metric on multisets of n vectors.
- From [Nguyen et al. \(2020\)](#), for IID \mathbf{X}_n and \mathbf{Y}_n with parameters θ_0 and θ , and $E\left[\|\mathbf{X}_1\|_2^2\right] + E\left[\|\mathbf{Y}_1\|_2^2\right] < \infty$, we have $D(\mathbf{X}_n, \mathbf{Y}_n)$ a.s. converges, as $n \rightarrow \infty$, to

$$D(\theta_0, \theta) = \frac{\Gamma\left(\frac{d+1}{2}\right)}{\pi^{\frac{d+1}{2}}} \int_{\mathbb{R}^d} \frac{|E \exp(it^\top \mathbf{X}_1) - E \exp(it^\top \mathbf{Y}_1)|^2}{\|t\|_2^{d+1}} d\lambda(t).$$

The energy distance

Write $\mathbf{Z}_n = (Z_i)_{i \in [n]}$, where $Z_i = (X_i, Y_i)$.

D1 \mathbf{Z}_n is a contiguous subsequence of $\mathbf{Z}_{\mathbb{Z}} = (Z_i)_{i \in \mathbb{Z}}$.

D2 The sequence $\mathbf{Z}_{\mathbb{Z}}$ is stationary and ergodic, where $\mathbf{X}_{\mathbb{Z}}$ and $\mathbf{Y}_{\mathbb{Z}}$ are independent.

D3 X_1 and Y_1 have marginal measures defined by $dF_1(x|\theta_0)$ and $dF_1(y|\theta)$, and $E\|X_1\|_2 + E\|Y_1\|_2 < \infty$.

We have the following result, using the ergodic U -statistics theorem of [Aaronson et al. \(1996\)](#).

Proposition 2

Under assumptions D1–D3, $D(\mathbf{X}_n, \mathbf{Y}_n)$ a.s. converges, as $n \rightarrow \infty$, to $D(\theta_0, \theta)$.

Example application

- Let \mathbf{X}_n arise from an ARMA(1,1) model with $\theta_0 = (\text{AR}_0, \text{MA}_0)$ with standard normal noise, $n = 1000$.
- Let \mathbf{Y}_n arise from an ARMA(1,1) model with $\Theta \sim \text{Uniform}([-1, 1]^2)$ with standard normal noise.
- Construct the lagged vector sequences $\tilde{\mathbf{X}}_n = ((X_i, X_{i+1}))_{i \in [n-1]}$ and $\tilde{\mathbf{Y}}_n = ((Y_i, Y_{i+1}))_{i \in [n-1]}$.
- Generate $\mathcal{W}_N = ((\mathbf{Y}_{n,j}, \Theta_j))_{j \in [N]}$ consisting of $N = 10000$ realizations.
- Compute the energy statistics $\mathbf{D}_N = \left(D(\tilde{\mathbf{X}}_n, \tilde{\mathbf{Y}}_{n,j}) \right)_{j \in [N]}$.
- Determine the sets of \mathbf{D}_N that satisfy $D(\tilde{\mathbf{X}}_n, \tilde{\mathbf{Y}}_{n,j}) < \varepsilon$, for various levels of ε , determined by the quantiles of \mathbf{D}_N .

Example application

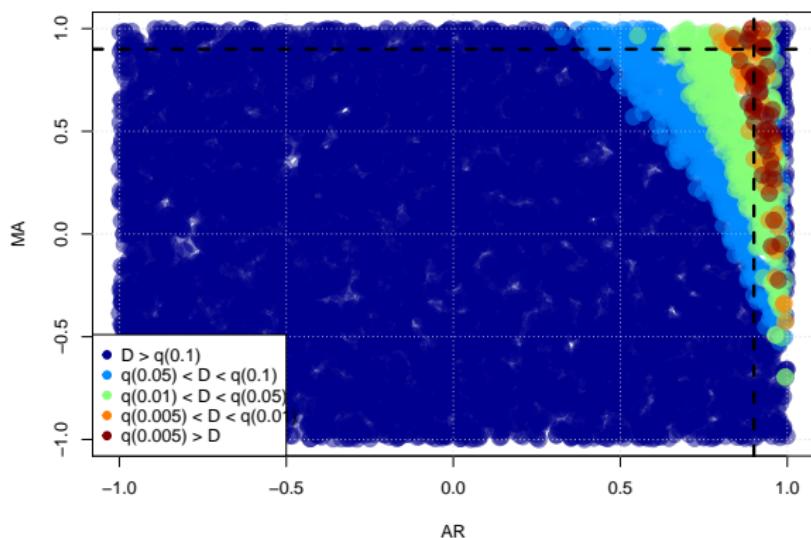


Figure: $\theta_0 = (\text{AR}_0, \text{MA}_0) = (0.9, 0.9)$.

Example application

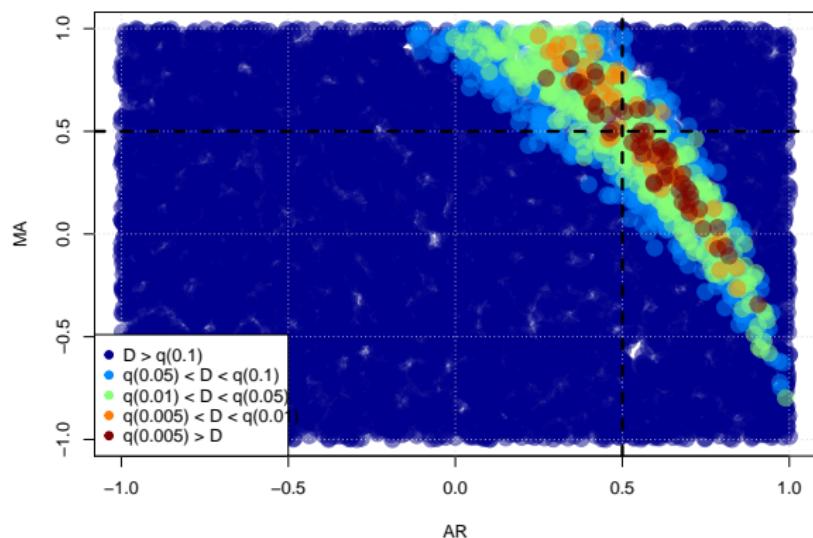


Figure: $\theta_0 = (\text{AR}_0, \text{MA}_0) = (0.5, 0.5)$.

Example application

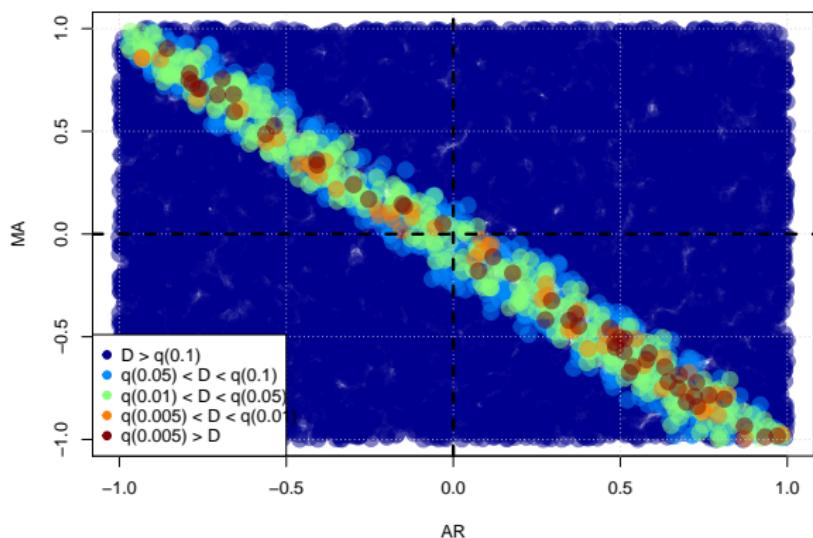


Figure: $\theta_0 = (AR_0, MA_0) = (0, 0)$.

A better example

- Let \mathbf{X}_n arise from an ARMA(0,2) model with $\theta_0 = (\text{MA1}_0, \text{MA2}_0)$ with standard normal noise, $n = 1000$.
- Let \mathbf{Y}_n arise from an ARMA(0,2) model with $\Theta \sim \text{Uniform}([-1, 1]^2)$ with standard normal noise.
- Construct the lagged vector sequences $\tilde{\mathbf{X}}_n = ((X_i, X_{i+1}))_{i \in [n-1]}$ and $\tilde{\mathbf{Y}}_n = ((Y_i, Y_{i+1}))_{i \in [n-1]}$.
- Generate $\mathcal{W}_N = ((\mathbf{Y}_{n,j}, \Theta_j))_{j \in [N]}$ consisting of $N = 10000$ realizations.
- Compute the energy statistics $\mathbf{D}_N = \left(D(\tilde{\mathbf{X}}_n, \tilde{\mathbf{Y}}_{n,j}) \right)_{j \in [N]}$.
- Determine the sets of \mathbf{D}_N that satisfy $D(\tilde{\mathbf{X}}_n, \tilde{\mathbf{Y}}_{n,j}) < \varepsilon$, for various levels of ε , determined by the quantiles of \mathbf{D}_N .

A better example

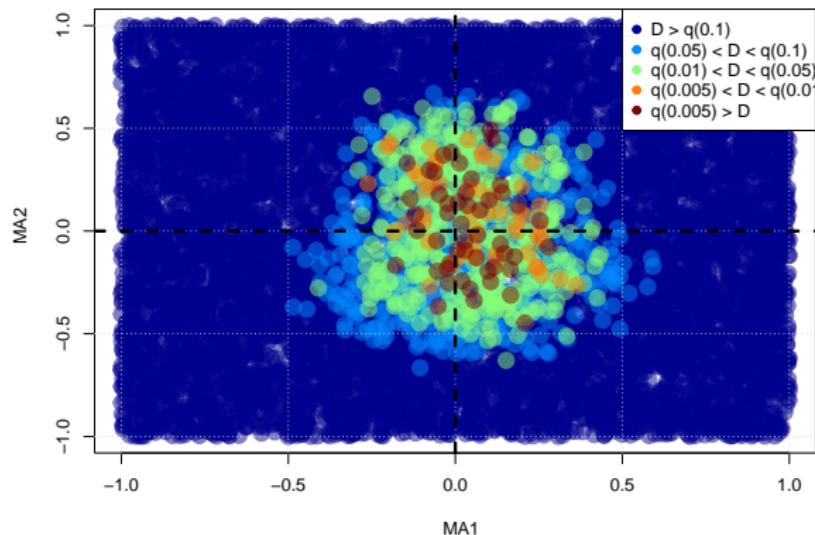


Figure: $\theta_0 = (\text{MA1}_0, \text{MA2}_0) = (0, 0)$.

A better example

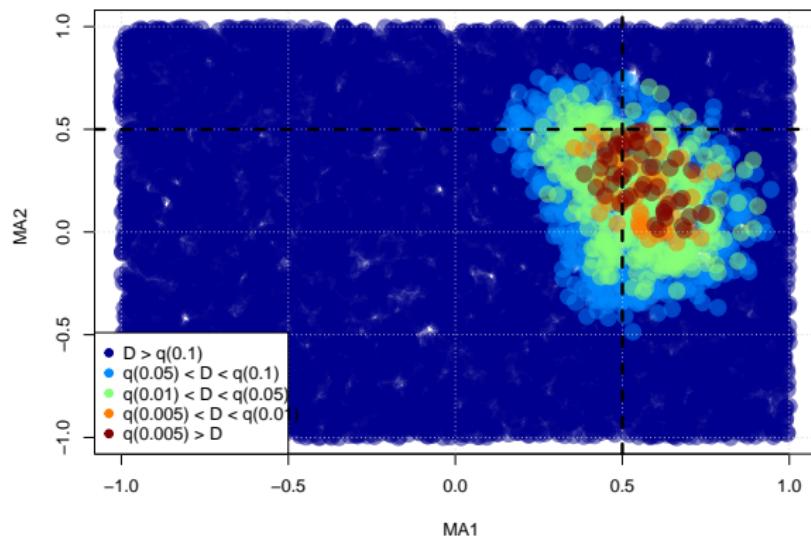


Figure: $\theta_0 = (\text{MA1}_0, \text{MA2}_0) = (0.5, 0.5)$.

A better example

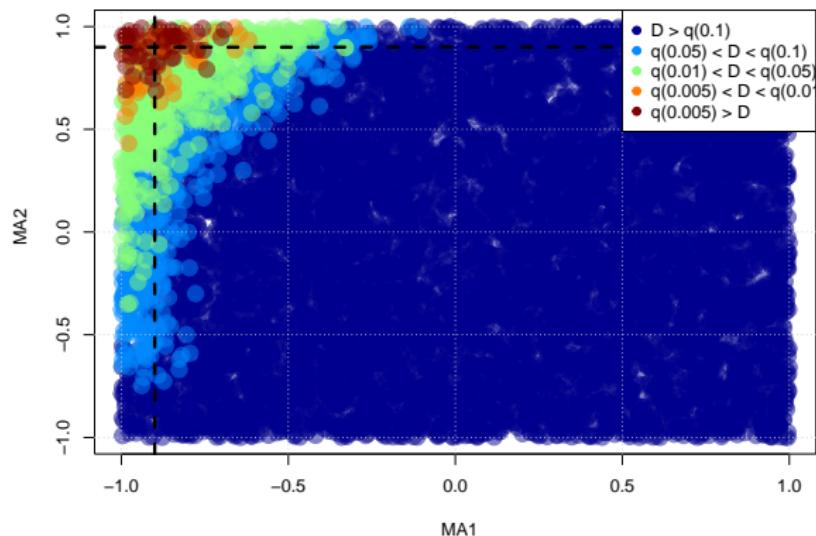


Figure: $\theta_0 = (\text{MA1}_0, \text{MA2}_0) = (-0.9, 0.9)$.

References I

- Aaronson, J., Burton, R., Dehling, H., Gilat, D., Hill, T., and Weiss, B. (1996). Strong laws for L- and U-statistics. *Transactions of the American Mathematical Society*, 348:2845–2866. [15](#)
- 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 B*, 8:235–269. [6](#), [8](#)
- Berti, P., Pratelli, L., and Rigo, P. (2006). Almost sure weak convergence of random probability measures. *Stochastics and Stochastics Reports*, 78:91–97. [10](#)
- Grubel, R. and Kabluchko, Z. (2016). A functional central limit theorem for branching random walks, almost sure weak convergence and applications to random trees. *Annals of Applied Probability*, 26:3659–3698. [10](#)

References II

- Jiang, B., Wu, T.-Y., and Wong, W. H. (2018). Approximate Bayesian computation with Kullback-Leibler divergence as data discrepancy. In *AISTATS*. [6](#), [9](#)
- Karabatsos, G. and Leisen, F. (2018). An approximate likelihood perspective on ABC methods. *Statistics Surveys*, 12:66–104. [7](#)
- Miller, J. W. and Dunson, D. B. (2019). Robust Bayesian inference via coarsening. *Journal of the American Statistical Association*, 114:1113–1125. [7](#), [9](#)
- Nguyen, H. D., Arbel, J., Lu, H., and Forbes, F. (2020). Approximate Bayesian computation via the energy statistic. *IEEE Access*, 8:131683–131698. [6](#), [9](#), [14](#)

References III

- Park, M., Jitkrittum, W., and Sejdinovic, D. (2016). K2-ABC: approximate Bayesian computation with kernel embeddings. In *AISTATS*. [7](#)
- Perez-Cruz, F. (2008). Kullback-Leibler divergence estimation of continuous distributions. In *ISIT 2008*. [6](#)
- Rubio, F. J. and Johansen, A. M. (2013). A simple approach to maximum intractable likelihood estimation. *Electronic Journal of Statistics*, 7:1632–1654. [8](#)
- Szekely, G. J. and Rizzo, M. L. (2004). Testing for equal distributions in high dimension. *InterStat*, 5:1249–1272. [6](#)
- Szekely, G. J. and Rizzo, M. L. (2017). The energy of data. *Annual Review of Statistics and Its Application*, 4:447–479. [14](#)

Thank you!

Most recent paper:

<https://doi.org/10.1109/ACCESS.2020.3009878>

Email: **h.nguyen5@latrobe.edu.au**

Website: **hiendn.github.io**