

# On the use of ABC-MCMC with inflated tolerance and post-correction

---

Matti Vihola

Department of Mathematics and Statistics, University of Jyväskylä

One world ABC seminar, 10 December 2020

Joint work with Jordan Franks (Newcastle)



Model:

- $\text{pr}(\theta)$  — the prior
- $g(y \mid \theta)$  — the observation model
- $y^*$  — the data (observed values of  $y$ )
- posterior  $\pi(\theta) \propto \text{pr}(\theta)g(y^* \mid \theta)$



Model:

- $\text{pr}(\theta)$  — the prior
- $g(y \mid \theta)$  — the observation model
- $y^*$  — the data (observed values of  $y$ )
- posterior  $\pi(\theta) \propto \text{pr}(\theta)g(y^* \mid \theta)$

ABC is inference methods for models that

- × have intractable (or extremely expensive) likelihood  $g(y^* \mid \theta)$ , but



Model:

- $\text{pr}(\theta)$  — the prior
- $g(y \mid \theta)$  — the observation model
- $y^*$  — the data (observed values of  $y$ )
- posterior  $\pi(\theta) \propto \text{pr}(\theta)g(y^* \mid \theta)$

ABC is inference methods for models that

- ✗ have intractable (or extremely expensive) likelihood  $g(y^* \mid \theta)$ , but
- ✓ (cheap) simulations of pseudo-data  $Y \sim g(\cdot \mid \theta)$ .



ABC rejection sampler:

(R1) Draw  $\Theta \sim \text{pr}(\cdot)$  and  $Y \sim g(\cdot | \Theta)$ .

(R2) If  $\|s(Y) - s(y^*)\| \leq \epsilon$ , output  $(\Theta, Y)$ ; otherwise go to (R1).

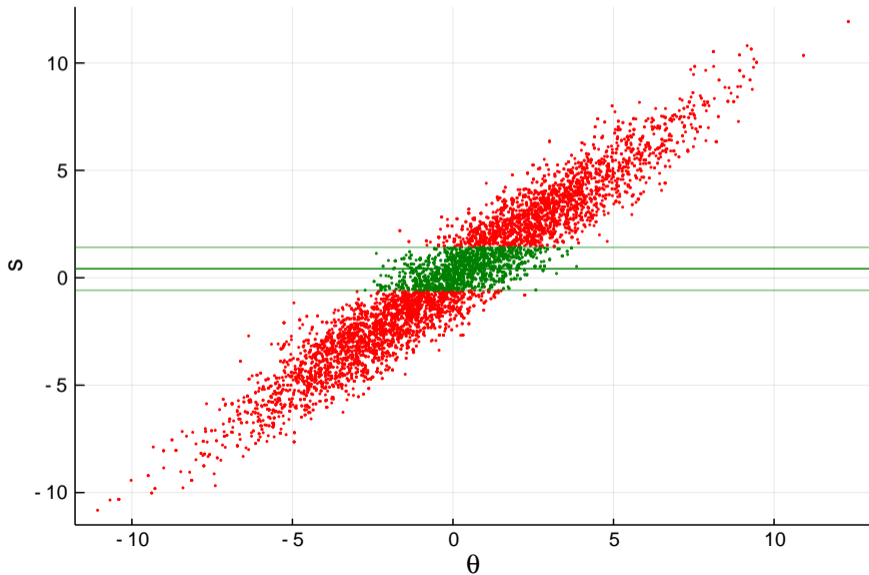
where

- $\|\cdot\|$  is Euclidean distance on  $\mathbb{R}^d$
- $s$  extracts  $d$  **summary statistics** from the (pseudo-)observations.

It is not hard to see that the distribution of  $(\Theta, Y)$  is

$$\tilde{\pi}_\epsilon(\theta, y) \propto \text{pr}(\theta)g(y | \theta)\mathbb{I}\{\|s(y) - s(y^*)\| \leq \epsilon\}.$$

**Example:**  $\text{pr}(\theta) = N(\theta; 0, 3^2)$ ,  $g(y | \theta) = N(y | \theta, 1)$ ,  $s(y) = y$ ,  $s(y^*) = 0.42$ ,  $\epsilon = 1.0$





If we only consider  $\Theta$ , its marginal is the 'pseudo-posterior'

$$\pi_\epsilon(\theta) = \frac{1}{Z_\epsilon} \text{pr}(\theta) \int g(y | \theta) \mathbb{I}\{\|s(y) - s(y^*)\| \leq \epsilon\} dy = \text{pr}(\theta) \frac{\ell_\epsilon(y^* | \theta)}{Z_\epsilon},$$

where  $\ell_\epsilon(y^* | \theta)$  is a 'pseudo-likelihood'

$$\ell_\epsilon(y^* | \theta) = \mathbb{P}(\|s(Y) - s(y^*)\| \leq \epsilon) \quad \text{where} \quad Y \sim g(\cdot | \theta)$$

How tolerance parameter  $\epsilon$  affects the pseudo-posterior?

- If  $\epsilon \rightarrow \infty$ , then  $\ell_\epsilon(y^* | \theta) \nearrow 1 \implies \pi_\epsilon(\theta) \rightarrow \text{pr}(\theta)$ .



If we only consider  $\Theta$ , its marginal is the 'pseudo-posterior'

$$\pi_\epsilon(\theta) = \frac{1}{Z_\epsilon} \text{pr}(\theta) \int g(y | \theta) \mathbb{I}\{\|s(y) - s(y^*)\| \leq \epsilon\} dy = \text{pr}(\theta) \frac{\ell_\epsilon(y^* | \theta)}{Z_\epsilon},$$

where  $\ell_\epsilon(y^* | \theta)$  is a 'pseudo-likelihood'

$$\ell_\epsilon(y^* | \theta) = \underbrace{\mathbb{P}(\|s(Y) - s(y^*)\| \leq \epsilon)}_{\approx 2\epsilon g(y^* | \theta) \text{ in the example}} \quad \text{where} \quad Y \sim g(\cdot | \theta)$$

How tolerance parameter  $\epsilon$  affects the pseudo-posterior?

- If  $\epsilon \rightarrow \infty$ , then  $\ell_\epsilon(y^* | \theta) \nearrow 1 \implies \pi_\epsilon(\theta) \rightarrow \text{pr}(\theta)$ .

If  $s(\cdot)$  are sufficient (& some regularity conditions):

- if  $\epsilon \rightarrow 0$ , then  $Z_\epsilon^{-1} \ell_\epsilon(y^* | \theta) \rightarrow Z^{-1} \ell(y^* | \theta) \implies \pi_\epsilon(\theta) \rightarrow \pi(\theta)$ .



- Markov process  $(X_t, Y_t)_{t \geq 0}$  with jump rates:

$$\theta_1 : X_t \rightarrow X_t + 1$$

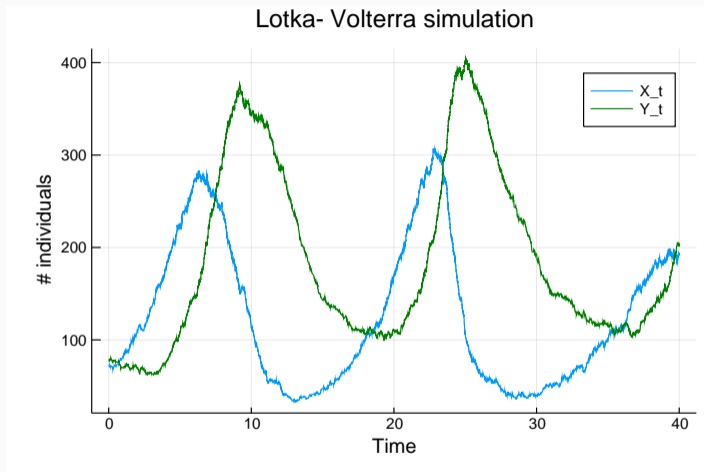
$$\theta_2 : X_t \rightarrow X_t - 1 \ \& \ Y_t \rightarrow Y_t + 1$$

$$\theta_3 : Y_t \rightarrow Y_t - 1$$

- Simulation with:

- $(X_0, Y_0) = (71, 79)$

- $\theta = (0.5, 0.0025, 0.3)$ .





Consider ABC with prior

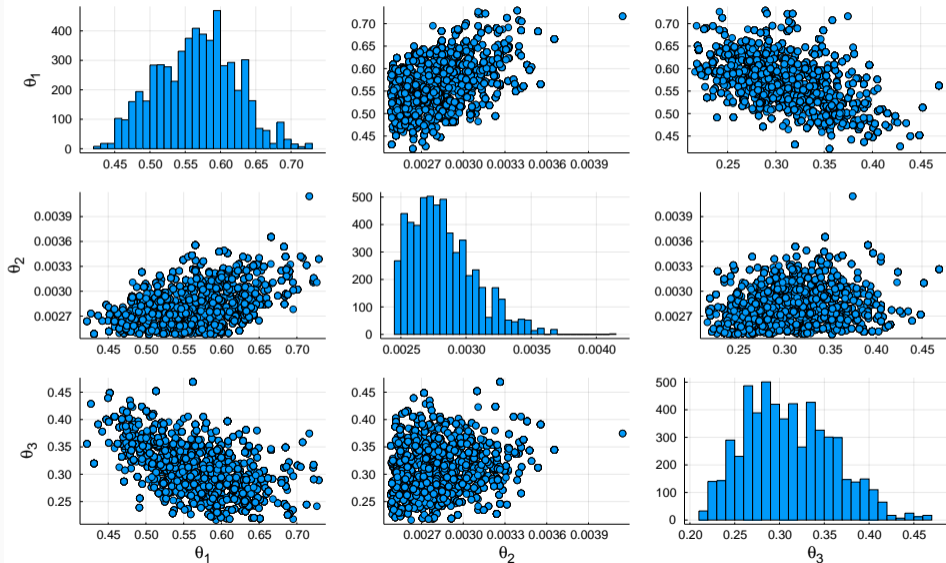
$$(\log \theta_1, \log \theta_2, \log \theta_3) \sim U([-6, 0]^3)$$

and (rather arbitrary?) summaries based on discrete-time observations  $\tilde{X}_k = X_{5k}$  and  $\tilde{Y}_k = Y_{5k}$ :

- sample autocorrelation of  $\tilde{X}_k$  at lag 2
- mean and standard deviations of  $\tilde{X}_k$  and  $\tilde{Y}_k$

Observed values  $(-46.32, 129.0, 88.64, 202.6, 140.8)$ .

# Lotka-Volterra: Inference with $\epsilon = 40$ (data from $\theta = (0.5, 0.0025, 0.3)$ )





Let  $q(y | x)$  be a proposal, start with some  $\Theta_0$  and  $L_0 = 1$ , and for  $k = 1, \dots, n$ :

- Draw  $\Theta'_k \sim q(\cdot | \Theta_{k-1})$  and  $Y'_k \sim g(\cdot | \Theta'_k)$ .
- Calculate  $L'_k := \mathbb{I}\{\|s(Y'_k) - s(y^*)\| \leq \epsilon\}$ .
- With probability

$$\min \left\{ 1, \frac{\text{pr}(\Theta'_k) L'_k q(\Theta_{k-1} | \Theta'_k)}{\text{pr}(\Theta_{k-1}) L_{k-1} q(\Theta'_k | \Theta_{k-1})} \right\},$$

accept and set  $(\Theta_k, L_k) \leftarrow (\Theta'_k, L'_k)$ ;

otherwise reject and set  $(\Theta_k, L_k) \leftarrow (\Theta_{k-1}, L_{k-1})$ .

ABC can be viewed as a pseudo-marginal MCMC (Andrieu & Roberts 2009)...



This is the (purely theoretical!) ‘marginal’ variant:

- Draw  $\Theta'_k \sim q(\cdot \mid \Theta_{k-1})$
- Calculate  $\ell_\epsilon(y^* \mid \Theta'_k) = \mathbb{P}(\|s(Y'_k) - s(y^*)\| \leq \epsilon)$ .
- With probability

$$\min \left\{ 1, \frac{\text{pr}(\Theta'_k) \ell_\epsilon(y^* \mid \Theta'_k) q(\Theta_{k-1} \mid \Theta'_k)}{\text{pr}(\Theta_{k-1}) \ell_\epsilon(y^* \mid \Theta_{k-1}) q(\Theta'_k \mid \Theta_{k-1})} \right\},$$

accept and set  $\Theta_k \leftarrow \Theta'_k$ ; otherwise reject and set  $\Theta_k \leftarrow \Theta_{k-1}$ .



ABC-MCMC is good because:

- Can use diffuse (uninformative, even improper) prior unlike rejection sampling.
- Admits convergence guarantees (to  $\pi_\epsilon$ ).
- Simple to use (hopefully!)



ABC-MCMC is good because:

- Can use diffuse (uninformative, even improper) prior unlike rejection sampling.
- Admits convergence guarantees (to  $\pi_\epsilon$ ).
- Simple to use (hopefully!)

Limitations of ABC-MCMC:

- Cheap-to-simulate models: MCMC requires many simulations.
- Moderate dimensional parameter (and summary).

Other likelihood-free methods may be more useful in more challenging settings!



Need to specify:

- Initial value  $\theta_0$
- The (log-)prior  $\text{pr}(\theta)$
- The observed summary statistic  $s(y^*)$
- Simulator of summaries  $s(Y_\theta)$
- Number of MCMC iterations



Need to specify:

- Initial value  $\theta_0$  `-ones(3)`
- The (log-)prior  $\text{pr}(\theta)$  `function prior(th)`
- The observed summary statistic  $s(y^*)$  `-s_obs`
- Simulator of summaries  $s(Y_\theta)$  `function sim!(s, th, rng)`
- Number of MCMC iterations `30000`

```
julia> using AdaptiveToleranceABC_MCMC
julia> out = abc_mcmc(-ones(3), prior, s_obs, sim!, 30000)
```

[https://github.com/mvihola/AdaptiveToleranceABC\\_MCMC.jl](https://github.com/mvihola/AdaptiveToleranceABC_MCMC.jl)



Adaptive state  $\xi_k = (\mu_k, \Gamma_k, \log \delta_k)$ :

- Draw  $\Theta'_k \sim q_{\Gamma_{k-1}}(\cdot | \Theta_{k-1})$  and  $Y'_k \sim g(\cdot | \Theta'_k)$ .
- Calculate  $L'_k := \mathbb{I}\{\|s(Y'_k) - s(y^*)\| \leq \delta_{k-1}\}$ .
- With probability

$$\alpha_k := \min \left\{ 1, \frac{\text{pr}(\Theta'_k) L'_k q(\Theta_{k-1} | \Theta'_k)}{\text{pr}(\Theta_{k-1}) L_{k-1} q(\Theta'_k | \Theta_{k-1})} \right\},$$

accept and set  $(\Theta_k, L_k) \leftarrow (\Theta'_k, L'_k)$ ;

otherwise reject and set  $(\Theta_k, L_k) \leftarrow (\Theta_{k-1}, L_{k-1})$ .

- Adapt  $\xi_k = \xi_{k-1} + \eta_k H(\xi_{k-1}; \alpha_k, \Theta_k)$ .



- $q_{\Gamma}(y | x) = N(y; x, \frac{2.38^2}{n_{\text{par}}}\Gamma)$  — random-walk proposal with covariance  $\propto \Gamma$ .
- Adaptation:

$$\mu_k = (1 - \eta_k)\mu_{k-1} + \eta_k\Theta_k$$

$$\Gamma_k = (1 - \eta_k)\Gamma_{k-1} + \eta_k(\Theta_k - \mu_{k-1})(\Theta_k - \mu_{k-1})^T$$

$$\log \delta_k = \log \delta_{k-1} + \eta_k(\alpha^* - \alpha_k)$$

where  $\eta_k = k^{-2/3}$  are step sizes.



- $q_{\Gamma}(y | x) = N(y; x, \frac{2.38^2}{n_{\text{par}}}\Gamma)$  — random-walk proposal with covariance  $\propto \Gamma$ .
- Adaptation:

$$\mu_k = (1 - \eta_k)\mu_{k-1} + \eta_k\Theta_k$$

$$\Gamma_k = (1 - \eta_k)\Gamma_{k-1} + \eta_k(\Theta_k - \mu_{k-1})(\Theta_k - \mu_{k-1})^T$$

$$\log \delta_k = \log \delta_{k-1} + \eta_k(\alpha^* - \alpha_k)$$

where  $\eta_k = k^{-2/3}$  are step sizes. That is,

- **Adaptive Metropolis** (Haario, Saksman, Tamminen 2001; Andrieu & Moulines 2006) mean-covariance update.



- $q_{\Gamma}(y | x) = N(y; x, \frac{2.38^2}{n_{\text{par}}}\Gamma)$  — random-walk proposal with covariance  $\propto \Gamma$ .
- Adaptation:

$$\mu_k = (1 - \eta_k)\mu_{k-1} + \eta_k\Theta_k$$

$$\Gamma_k = (1 - \eta_k)\Gamma_{k-1} + \eta_k(\Theta_k - \mu_{k-1})(\Theta_k - \mu_{k-1})^T$$

$$\log \delta_k = \log \delta_{k-1} + \eta_k(\alpha^* - \alpha_k)$$

where  $\eta_k = k^{-2/3}$  are step sizes. That is,

- Adaptive Metropolis (Haario, Saksman, Tamminen 2001; Andrieu & Moulines 2006) mean-covariance update.
- [Acceptance-rate adaptation](#) in the spirit of Andrieu & Robert (2001).



Must have a reasonable acceptance rate, otherwise ABC-MCMC will not mix.

- Optimise  $\delta_k \rightarrow \delta_*$  such that the ABC-MCMC has a desired acceptance rate  $\alpha^*$ .
- Observation:  $\delta \mapsto \alpha_\delta := \mathbb{E}_{\pi_\delta}[\alpha_k]$  (mean accept rate w/  $\delta$ ) is **increasing**.
  - If too small accept rate, increase  $\delta_k$  (and vice versa).

$\rightsquigarrow$  Convergent adaptation (for fixed  $\Gamma$ ). But how to choose  $\alpha^*$ ?



## Acceptance rate adaptation: rationale

Must have a reasonable acceptance rate, otherwise ABC-MCMC will not mix.

- Optimise  $\delta_k \rightarrow \delta_*$  such that the ABC-MCMC has a desired acceptance rate  $\alpha^*$ .
- Observation:  $\delta \mapsto \alpha_\delta := \mathbb{E}_{\pi_\delta}[\alpha_k]$  (mean accept rate w/  $\delta$ ) is **increasing**.
  - If too small accept rate, increase  $\delta_k$  (and vice versa).

$\rightsquigarrow$  Convergent adaptation (for fixed  $\Gamma$ ). But how to choose  $\alpha^*$ ?

If  $\{\pi_\delta\}_{\delta>0}$  were Gaussian, then AM adaptation **for fixed  $\delta$**  would lead to mean accept rate  $\leq \alpha_{\text{opt}}^{\text{RWM}}$  where  $\alpha_{\text{opt}}^{\text{RWM}} \in [0.234, 0.44]$  (depending on  $n_{\text{par}} \dots$ )

- $\delta \rightarrow \infty \implies$  random-walk Metropolis with accept rate  $\approx \alpha_{\text{opt}}^{\text{RWM}}$ .
- For any  $\delta$  the 'marginal' variant of ABC-MCMC has accept rate  $\approx \alpha_{\text{opt}}^{\text{RWM}}$ .  
The accept rate of pseudo-marginal (ABC-MCMC) is less.
- When  $\delta \rightarrow 0$ , accept rate  $\rightarrow 0$ .

This suggests  $\alpha^* \in (0, 0.234) \subset (0, \alpha_{\text{opt}}^{\text{RWM}})$ .

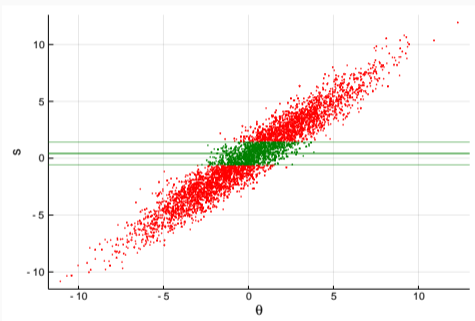


- We use  $\alpha^* = 0.10$  as 'balanced' acceptance rate.
- Initialise  $\delta_0 = \|s(Y_{\theta_0}) - s(y^*)\|$   
 $\implies \delta_0 \gg \delta_*$  — this ensures acceptances during the initial phase;  
 $\sim$  the explicit annealing scheme of Ratmann et al. (2007)?
- Changing tolerance changes target distribution  
 $\rightsquigarrow$  stop tolerance adaptation after burn-in.
- Continue covariance adaptation also after burn-in.

This leads to adaptive ABC-MCMC with automatically chosen tolerance  $\delta \approx \delta_*$ ...

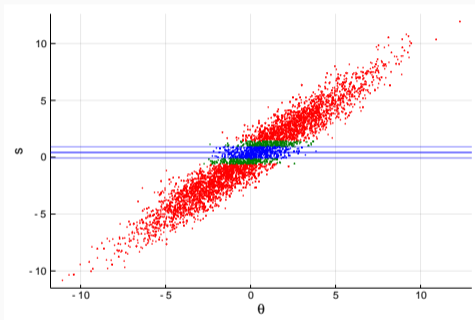
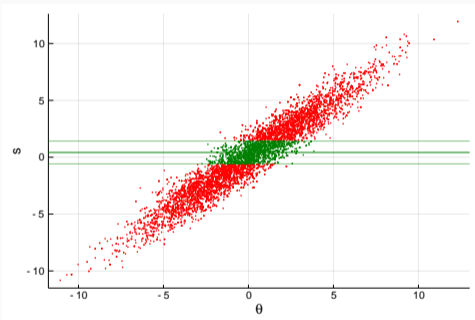


- Rejection sampling: Tolerance  $\delta = 1 \implies$  samples



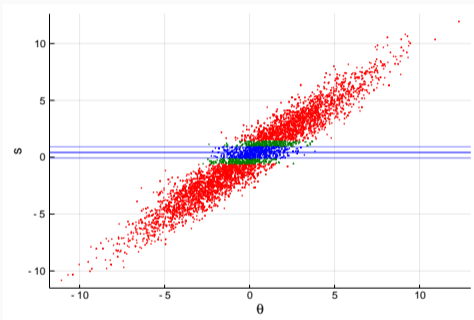
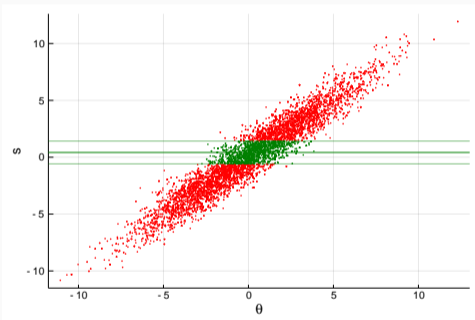


- Rejection sampling: Tolerance  $\delta = 1 \implies$  **samples**  $\implies$  **refined** for  $\epsilon = 0.5 < \delta$

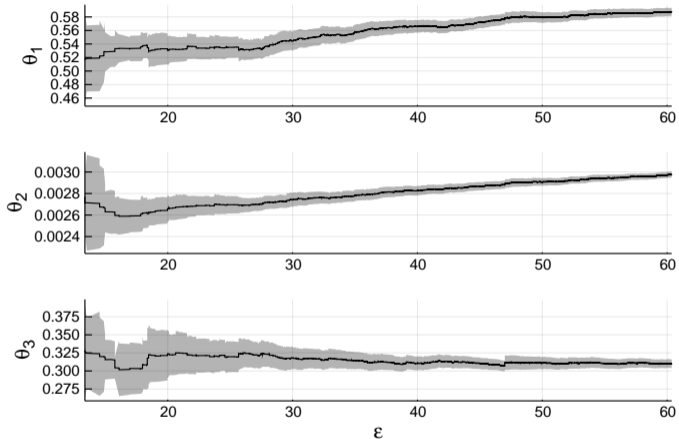




- Rejection sampling: Tolerance  $\delta = 1 \implies$  samples  $\implies$  refined for  $\epsilon = 0.5 < \delta$



- The same can be done with ABC-MCMC (Wegmann et al. *Genetics*, 2009)... but how does it affect posterior mean uncertainty?



```
julia> est = abc_postprocess(out, x -> exp.(x))
julia> plot_abc(est; labels=["\\theta_$(i)" for i in 1:3])
```



MCMC output  $(\Theta_k, D_k)_{k=1, \dots, n}$  with  $D_k = \|s(Y_k) - s(y^*)\|$ ;  
test function  $f : \mathbb{R}^{n_{\text{par}}} \rightarrow \mathbb{R}$ .

- Calculate self-normalised IS estimator:

$$U_k^{(\delta, \epsilon)} := \frac{\mathbb{I}\{D_k \leq \epsilon\}}{\mathbb{I}\{D_k \leq \delta\}}, \quad W_k^{(\delta, \epsilon)} := \frac{U_k^{(\delta, \epsilon)}}{\sum_{j=1}^n U_j^{(\delta, \epsilon)}},$$

$$E_{\delta, \epsilon}(f) := \sum_{k=1}^n W_k^{(\delta, \epsilon)} f(\Theta_k), \quad S_{\delta, \epsilon}(f) := \sum_{k=1}^n (W_k^{(\delta, \epsilon)})^2 \{f(\Theta_k) - E_{\delta, \epsilon}(f)\}^2.$$

- Calculate  $\hat{\tau}_\delta(f) = \text{IACT}(f(X_1), \dots, f(X_n))$ .
- Report (say) 95% approximate CI

$$\left[ E_{\delta, \epsilon}(f) \pm 1.96 S_{\delta, \epsilon}(f) \hat{\tau}_\delta(f) \right] \quad \text{for any} \quad \epsilon \in \left[ \min_k(D_k), \delta \right]$$

(In the case of simple cut-off, the CI may be calculated for all  $\epsilon$  in  $O(n \log n)$  time.)



- Under general conditions, for  $\epsilon \leq \delta$

$$\underbrace{nS_{\delta,\epsilon}(f)}_{\rightarrow v_{\delta,\epsilon,\text{IS}}^2} \underbrace{\hat{\tau}_{\delta}(f)}_{\rightarrow \tau_{\delta}(f)} \xrightarrow{n \rightarrow \infty} \check{\sigma}^2(f) \gtrsim \sigma_{\delta,\epsilon}^2(f) = v_{\delta,\epsilon,\text{IS}}^2 \tau_{\delta,\epsilon}(f),$$

where  $\tau_{\delta,\epsilon}(f)$  is integrated autocorrelation of  $W_k^{(\delta,\epsilon)} f(\Theta_k)$ .

- The approximate upper bound is due to

$$\tau_{\delta}(f) \gtrsim \tau_{\delta,\epsilon}(f),$$

by continuity for  $\delta \approx \epsilon$ , and loosely analytically justified for  $\epsilon \ll \delta$ .

- Direct estimation of the asymptotic variance of  $W_k^{(\delta,\epsilon)} f(\Theta_k)$  is also possible, but
  - likely unstable for smaller  $\epsilon$  (lots of zeros...), and
  - computationally more demanding...



Above, we had “ $\|s(Y) - s(y^*)\|$ ” and “ $\mathbb{I}\{\cdot \leq \epsilon\}$ ”, but in general we might have:

- Any ‘dissimilarity’ function “ $D(Y, y^*)$ ”
- Any ‘kernel’ “ $k_\epsilon(\cdot)$ ” with values in  $[0, 1]$
- Or we might merge both and define directly  $K_\epsilon(Y, y^*) \in [0, 1] \dots$

The adaptive tolerance ABC-MCMC and the post-processing generalise directly.<sup>1</sup>

It is also possible to accommodate to ABC-MCMC with averaged pseudo-data:

$$L'_k = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{\|s(Y_k^{(i)}) - s(y^*)\| \leq \epsilon\}, \quad Y_k^{(1)}, \dots, Y_k^{(m)} \stackrel{\text{i.i.d.}}{\sim} g(\cdot \mid \Theta'_k)$$

Theoretical findings of Bornn et al. (2017) suggest that  $m = 1$  might be enough...

---

<sup>1</sup>The loose theoretical justification of CI is lost, but empirical results remain promising.



- *Regression correction/adjustment* is often useful with ABC.
- Instead of looking at  $(\Theta_k) \sim \pi_\delta$  only, consider  $(\Theta_k, S_k)$  where  $S_k = s(Y_k)$  are summaries of accepted pseudo-data.
- Fit a regression model such as:

$$\theta_k = \beta_0 + \beta^T (S_k - s(y^*)) + \eta_k, \quad \text{where } \eta_k \text{ are residuals}$$

$\implies$  least squares estimates  $\hat{\beta}_0, \hat{\beta}$

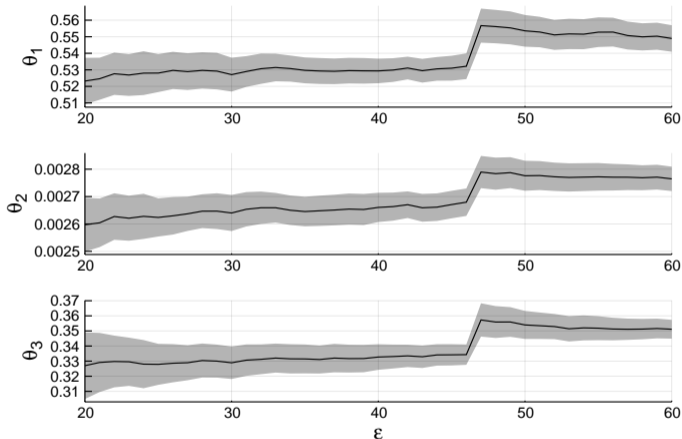
$\implies$  empirical residuals  $\hat{\eta}_k = \theta_k - \hat{\beta}_0 - \hat{\beta}^T (S_k - s(y^*))$

- The regression corrected samples are

$$\hat{\theta}_k := \hat{\beta}_0 + \hat{\eta}_k = \theta_k - \hat{\beta}^T (S_k - s(y^*));$$

(if model and estimates correct,  $\hat{\theta}_k$  follow  $\theta \mid S = s(y^*) \dots$ )

- Post-correction  $\delta \rightarrow \epsilon \implies$  weighted regression  $\implies$  correction.  
(Experimental CI analogous to plain ABC...)



```
julia> est_r = abc_postprocess(out, x->exp.(x), 20:out.cutoff.tol; regress=true)
julia> plot_abc(est_r; labels=["\\theta_$(i)" for i in 1:3])
```



- Easy-to-use ABC-MCMC
  - The only (mandatory) user-defined parameter is  $n$  — the number of MCMC steps.
  - Julia package — give it a try!
  - (Julia because it is good for implementing a *fast model simulator*...)
- Implicit assumptions:
  - Large  $n$ .
  - Low-dimensional parameter & summary.
  - Roughly unimodal (pseudo-)posterior — because of random-walk MCMC.
- No pre-defined ‘target’ tolerance  $\epsilon$ .
- Instead, automatically chosen ‘simulation tolerance’  $\delta$   
and possibility to inspect a range of  $\epsilon \leq \delta$ .
  - Can reveal bias due to still-too-large  $\delta$ ...

# References

- **Marjoram, P., Molitor, J., Plagnol, V., & Tavaré, S. (2003)**  
Markov chain Monte Carlo without likelihoods.  
*Proceedings of the National Academy of Sciences*, 100(26), 15324–15328.
- **Andrieu, C., & Roberts, G. O. (2009)**  
The pseudo-marginal approach for efficient Monte Carlo computations.  
*The Annals of Statistics*, 37(2), 697–725.
- **Wegmann, D., Leuenberger, C., & Excoffier, L. (2009)**  
Efficient approximate Bayesian computation coupled with Markov chain Monte Carlo without likelihood.  
*Genetics*, 182(4), 1207–1218.
- **Ratmann, O., Jørgensen, O., Hinkley, T., Stumpf, M., Richardson, S., & Wiuf, C. (2007).**  
Using likelihood-free inference to compare evolutionary dynamics of the protein networks of *H. pylori* and *P. falciparum*.  
*PLoS Comput Biol*, 3(11), e230.
- **Bornn, L., Pillai, N. S., Smith, A., & Woodard, D. (2017)**  
The use of a single pseudo-sample in approximate Bayesian computation.  
*Statistics and Computing*, 27(3), 583–590.
- **Haario, H., Saksman, E. & Tamminen, J. (2001)**  
An adaptive Metropolis algorithm.  
*Bernoulli*, 7(2), 223–242.
- **Andrieu, C., & Moulines, É. (2006).**  
On the ergodicity properties of some adaptive MCMC algorithms.  
*The Annals of Applied Probability*, 16(3), 1462–1505.
- **Andrieu, C., & Robert, C. P. (2001)**  
Controlled MCMC for optimal sampling.  
Tech. Rep. Ceremade 0125, Université Paris Dauphine.
- **Vihola, M., & Franks, J. (2020)**  
On the use of approximate Bayesian computation Markov chain Monte Carlo with inflated tolerance and post-correction.  
*Biometrika*, 107(2), 381–395.

Julia package: [https://github.com/mvihola/AdaptiveToleranceABC\\_MCMC.jl](https://github.com/mvihola/AdaptiveToleranceABC_MCMC.jl)