

Approximate Bayesian Computation with Statistical Distances for Model Selection

Clara Grazian

School of Mathematics and Statistics, University of Sydney

clara.grazian@sydney.edu.au

OWABI Seminar Series

Online

27 November 2025

The problem of model choice

The selection of a model that best explains a dataset or process is central to statistics and relevant across disciplines.

While no model can fully capture all nuances in complex phenomena, approximating the true data-generating process (DGP) can yield valuable insights [Molnar et al., 2022, Gelman and Shalizi, 2013].

However, many models involve complex interdependencies or latent variables that render their likelihood functions intractable [Martin et al., 2024] → [ABC](#).

- 1 To mitigate the curse of dimensionality, ABC often relies on low-dimensional summary statistics but this may be particularly problematic for model selection [Robert et al., 2011].
- 2 Forbes et al. [2022] introduced a framework that constructs surrogate posteriors using finite Gaussian mixtures
- 3 other ABC methods compare entire empirical distributions using statistical distances [Park et al., 2016, Bernton et al., 2019, Nguyen et al., 2020, Frazier, 2020, Drovandi and Frazier, 2022]
- 4 machine learning can be used to learn informative low-dimensional summaries automatically [Blum and François, 2010, Sheehan and Song, 2016, Jiang et al., 2017]

Data: Let $\mathbf{y} = (y_1, \dots, y_n)^T \in \mathcal{Y}^n \subseteq \mathbb{R}^n$ denote the observed data, generated from $P_{\boldsymbol{\theta}_0}^{(n)}$ within a model family $\mathcal{P} = \{P_{\boldsymbol{\theta}}^{(n)} : \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^d\}$.

Prior: Model M is parameterized by $\boldsymbol{\theta}$ with prior distribution $\pi(\boldsymbol{\theta})$

Likelihood: The likelihood function is $p(\mathbf{y}|\boldsymbol{\theta})$

Posterior: The posterior distribution is $\pi(\boldsymbol{\theta}|\mathbf{y})$

Marginal: The marginal likelihood is as $p(\mathbf{y}) = \int_{\Theta} p(\mathbf{y}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})d\boldsymbol{\theta}$.

Model choice in a Bayesian setting

- a set of K candidate models M_1, \dots, M_K
- each model M_k has parameters $\boldsymbol{\theta}_k \in \Theta_k \subseteq \mathbb{R}^{d_k}$
- prior distribution for parameters $\pi_k(\boldsymbol{\theta}_k)$ and prior probabilities for models $\pi(M_k) = \pi(M = M_k)$ such that $\sum_{k=1}^K \pi(M_k) = 1$
- likelihood function $p_k(\mathbf{y}|\boldsymbol{\theta}_k) = p(\mathbf{y}|\boldsymbol{\theta}_k, M_k)$
- the posterior probability for model M_k is:

$$\pi(M_k|\mathbf{y}) = \frac{p_k(\mathbf{y})\pi(M_k)}{\sum_{j=1}^K p_j(\mathbf{y})\pi(M_j)} \quad \forall k = 1, \dots, K,$$

The true DGP $P_{M_0, \boldsymbol{\theta}_0}^{(n)}$ is assumed to belong to the broader set

$$\mathcal{P} = \{P_{M_k, \boldsymbol{\theta}_k}^{(n)} \mid \boldsymbol{\theta}_k \in \Theta_k \subseteq \mathbb{R}^{d_k}, k \in \{1, \dots, K\}\}.$$

ABC for Model Choice

Model choice via ABC introduces challenges beyond standard inference [Robert et al., 2011]. Consider the Bayes factor between Model M_1 and M_2 ; the ABC approximation to B_{12} is

$$\hat{B}_{12}(\mathbf{y}) = \frac{\pi(M = M_2) \sum_{t=1}^T \mathbb{I}_{m^{(t)}=M_1} \mathbb{I}_{\rho\{\eta(\mathbf{z}^{(t)}), \mathbf{y}\} \leq \varepsilon}}}{\pi(M = M_1) \sum_{t=1}^T \mathbb{I}_{m^{(t)}=M_1} \mathbb{I}_{\rho\{\eta(\mathbf{z}^{(t)}), \mathbf{y}\} \leq \varepsilon}}$$

By letting $T \rightarrow \infty$ and $\varepsilon \rightarrow 0$, we have

$$B_{12}^{\eta}(\mathbf{y}) = \frac{\int \pi_1(\theta_1) f_1^{\eta}(\eta(\mathbf{y}|\theta_1)) d\theta_1}{\int \pi_2(\theta_2) f_2^{\eta}(\eta(\mathbf{y}|\theta_2)) d\theta_2}$$

But this means that there is a mismatch between the true Bayes factor and the ABC Bayes factor

$$B_{12}(\mathbf{y}) = \frac{g_1(\mathbf{y}) \int \pi_1(\theta_1) f_1^{\eta}(\eta(\mathbf{y}|\theta_1)) d\theta_1}{g_2(\mathbf{y}) \int \pi_2(\theta_2) f_2^{\eta}(\eta(\mathbf{y}|\theta_2)) d\theta_2} = \frac{g_1(\mathbf{y})}{g_2(\mathbf{y})} B_{12}^{\eta}(\mathbf{y}).$$

Alternatives: Marin et al. [2013], Barnes et al. [2012], Prangle et al. [2014].

To incorporate statistical distances into model selection, the ABC distance metric $\rho(\cdot, \cdot)$ is replaced with a discrepancy measure $\mathcal{D}(\cdot, \cdot)$ that operates on the probability measures of the observed data \mathbf{y} and simulated data \mathbf{z} .

Denote μ_{0, θ_0} as the measure associated with the true distribution $P_{M_0, \theta_0}^{(n)}$, and $\mu_{k^*, \theta_{k^*}}$ as the measure associated with the simulated data's distribution $P_{M_{k^*}, \theta_{k^*}}^{(n)} \in \mathcal{P}$.

The empirical distributions of \mathbf{y} and \mathbf{z} replace the true measures, and are similarly defined as $\hat{\mu}_{0, \theta_0} = n^{-1} \sum_{i=1}^n \delta_{y_i}$ and $\hat{\mu}_{k^*, \theta_{k^*}} = n^{-1} \sum_{i=1}^n \delta_{z_i}$ respectively.

Algorithm 1 Discrepancy-based ABC-MC

- 1: **Input:** Given a set of observations $\mathbf{y} = (y_1, \dots, y_n)^T$, a set of possible models $\{M_1, M_2, \dots, M_K\}$, a tolerance level ε , and a discrepancy metric $\mathcal{D}(\cdot, \cdot)$:
 - 2: **for** $i = 1, \dots, N$ **do**
 - 3: **repeat**
 - 4: Generate M_{k^*} from $\pi(M_k)$, $k = 1, \dots, K$.
 - 5: Generate $\boldsymbol{\theta}_{k^*}$ from $\pi_{k^*}(\boldsymbol{\theta}_{k^*})$.
 - 6: Generate $\mathbf{z} = (z_1, \dots, z_n)^T$ from $P_{M_{k^*}, \boldsymbol{\theta}_{k^*}}^{(n)}$.
 - 7: **until** $\mathcal{D}(\hat{\boldsymbol{\mu}}_0, \boldsymbol{\theta}_0, \hat{\boldsymbol{\mu}}_{k^*}, \boldsymbol{\theta}_{k^*}) \leq \varepsilon$,
 - 8: Set $M^{(i)} = M_{k^*}$ and $\boldsymbol{\theta}^{(i)} = \boldsymbol{\theta}_{k^*}$.
 - 9: **end for**
 - 10: **Output:** A set of values $(M^{(1)}, \dots, M^{(N)})$ and $(\boldsymbol{\theta}^{(1)}, \dots, \boldsymbol{\theta}^{(N)})$ from $\pi_\varepsilon(M_k | \mathbf{y})$ and $\pi_{\varepsilon, k}(\boldsymbol{\theta}_k | \mathbf{y}, M_k)$, for $k = 1, \dots, K$, respectively.
-

What distances?

The maximum mean discrepancy (MMD) [Park et al., 2016] compares embedded empirical distributions in a reproducing kernel Hilbert space, using a kernel $g : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$,

$$MMD^2(\mu_{0, \theta_0}, \mu_{k^*, \theta_{k^*}}) = \mathbb{E}[g(y_1, y_2)] + \mathbb{E}[g(z_1, z_2)] - 2\mathbb{E}[g(y_1, z_1)].$$

The unbiased empirical estimate is:

$$MMD^2(\hat{\mu}_{0, \theta_0}, \mu_{k^*, \theta_{k^*}}) = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n g(y_i, y_j) + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n g(z_i, z_j) + \frac{2}{n^2} \sum_{i=1}^n \sum_{j=1}^n g(y_i, z_j).$$

What distances?

Bernton et al. [2019] proposed using the Wasserstein distance in ABC.

For univariate data and assuming that $P_{M_0, \theta_0}^{(n)}$ and $P_{M_{k^*}, \theta_{k^*}}^{(n)}$ have finite p -th moment with $p \geq 1$, the p -Wasserstein distance is

$$\mathcal{W}_p(\mu_0, \theta_0, \mu_{k^*}, \theta_{k^*}) = \left(\int_0^1 |F_{\mu_0, \theta_0}^{-1}(\lambda) - F_{\mu_{k^*}, \theta_{k^*}}^{-1}(\lambda)|^p d\lambda \right)^{1/p},$$

where $F_{\mu_0, \theta_0}(\cdot)$ and $F_{\mu_{k^*}, \theta_{k^*}}(\cdot)$ are the cumulative distribution.

The empirical p -Wasserstein distance, with $p = 1$, is computed by comparing order statistics:

$$\mathcal{W}_1(\hat{\mu}_0, \theta_0, \hat{\mu}_{k^*}, \theta_{k^*}) = n^{-1} \sum_{i=1}^n |y_{(i)} - z_{(i)}|.$$

What distances?

The Cramér–von Mises (CvM) distance [Frazier, 2020] measures the L_2 difference between the empirical CDF $\hat{F}_{\mu_{k^*}, \theta_{k^*}}$ and the CDF of the theoretical distribution F_{μ_0, θ_0}

$$\mathcal{C}^2(\mu_0, \theta_0, \hat{\mu}_{k^*}, \theta_{k^*}) = \int_{\mathcal{Y}} \left[\hat{F}_{\mu_{k^*}, \theta_{k^*}}(y) - F_{\mu_0, \theta_0}(y) \right]^2 dF_{\mu_0, \theta_0}(y).$$

It is estimated by:

$$\mathcal{C}^2(\hat{\mu}_0, \theta_0, \hat{\mu}_{k^*}, \theta_{k^*}) = \frac{U}{2n^2} - \frac{4n^2 - 1}{12n},$$

where U involves the ranks of \mathbf{y} and \mathbf{z} in their pooled sample.

Theoretical justification: Consistency

Consider Algorithm 1 using

- the Wasserstein distance [Bernton et al., 2019],
- the candidate models are well-separated in the Wasserstein sense, and
- a sequence $\varepsilon_n \rightarrow 0$ (since with growing sample size the empirical distributions concentrate around their population laws),
- observations are i.i.d.

The ABC posterior probability of model $M_{k^*} \in \{M_1, \dots, M_K\} = \mathcal{M}$ is

$$\pi_{\varepsilon_n}(M_{k^*} | \mathbf{y}) \propto \pi(M_{k^*}) \int_{\Theta_{k^*}} \pi(\boldsymbol{\theta}_{k^*} | M_{k^*}) \mathbb{P} \{ \mathcal{W}_p(\hat{\boldsymbol{\mu}}_0, \boldsymbol{\theta}_0, \hat{\boldsymbol{\mu}}_{k^*}, \boldsymbol{\theta}_{k^*}) \leq \varepsilon_n \} d\boldsymbol{\theta}_{k^*}.$$

Under regularity conditions, ABC-Wass yields consistent parameter estimation [Bernton et al., 2019]. This property extends to model selection.

Theoretical justification: Consistency

Theorem (Model Selection Consistency of ABC-Wass)

Assume:

(A1) (Identifiability): For all $M_{k^*} \neq M_0$, there exists $\delta_{k^*} > 0$ such that $\inf_{\boldsymbol{\theta}_{k^*} \in \Theta_{k^*}} \mathcal{W}_p(P_{M_0, \boldsymbol{\theta}_0}^{(n)}, P_{M_{k^*}, \boldsymbol{\theta}_{k^*}}^{(n)}) > \delta_{k^*}$.

(A2) (Prior Positivity): $\pi(M_0) > 0$ and $\pi_{M_0}(\boldsymbol{\theta}_0) > 0$.

(A3) (Tail Behavior): For any sequence $\varepsilon_n \rightarrow 0$ with $\varepsilon_n < \delta_{k^*}/2$ for large n

$$\int_{\Theta_{k^*} \setminus \mathcal{B}_{\delta_{k^*}}} \pi_{M_{k^*}}(\boldsymbol{\theta}_{k^*}) \cdot \mathbb{P}(\mathcal{W}_p(\hat{\boldsymbol{\mu}}_{0, \boldsymbol{\theta}_0}, \hat{\boldsymbol{\mu}}_{k^*, \boldsymbol{\theta}_{k^*}}) \leq \varepsilon_n) d\boldsymbol{\theta}_{k^*} \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

(A4) (Empirical Convergence): $\mathcal{W}_p(\hat{\boldsymbol{\mu}}_{0, \boldsymbol{\theta}_0}, P_{M_0, \boldsymbol{\theta}_0}^{(n)}) \rightarrow 0$ and $\mathcal{W}_p(\hat{\boldsymbol{\mu}}_{k^*, \boldsymbol{\theta}_{k^*}}, P_{M_{k^*}, \boldsymbol{\theta}_{k^*}}^{(n)}) \rightarrow 0$ almost surely as $n \rightarrow \infty$, and $\Delta_n := \sup_{M_{k^*}, \boldsymbol{\theta}_{k^*}} \mathbb{E}[\mathcal{W}_p(\hat{\boldsymbol{\mu}}_{k^*, \boldsymbol{\theta}_{k^*}}, P_{M_{k^*}, \boldsymbol{\theta}_{k^*}}^{(n)})]$ satisfies $\Delta_n = o(\varepsilon_n)$.

Then, for any sequence $\varepsilon_n \rightarrow 0$ such that $\varepsilon_n < \min_{k^*} \delta_{k^*}/2$,

$$\pi_{\varepsilon_n}(M_0 \mid \mathbf{y}) \xrightarrow{P} 1 \quad \text{as } n \rightarrow \infty.$$

Theoretical justification: Consistency

- **Assumption (A1)** ensures that for any incorrect model $M_{k^*} \neq M_0$ and any parameter $\theta_{k^*} \in \Theta_{k^*}$, the Wasserstein distance between the simulated data distribution $P_{M_{k^*}, \theta_{k^*}}^{(n)}$ and the true distribution $P_{M_0, \theta_0}^{(n)}$ is bounded below by some fixed $\delta_{k^*} > 0$.
- **Assumption (A2)** requires that both the model prior and the parameter prior assign positive mass to neighborhoods of the true model and true parameter.
- **Assumption (A3)** ensures that parameter values far (in Wasserstein distance) to $P_{M_0, \theta_0}^{(n)}$ cannot produce empirical measures within the small threshold with non-negligible prior mass.
- **Assumption (A4)** holds for i.i.d. observations and in cases with a finite p -th moment.

Theorem (Robustness to Model Misspecification)

Let the true data distribution $P_{M_0, \boldsymbol{\theta}_0}^{(n)}$ be misspecified with respect to the candidate set

$$\{P_{M_k, \boldsymbol{\theta}_k}^{(n)} : M_k \in \mathcal{M}, \boldsymbol{\theta}_k \in \Theta_k\}.$$

Suppose Assumptions (A1)-(A4) hold, and let $\varepsilon_n \rightarrow 0$, at a rate slower than the stochastic convergence of the Wasserstein distance $\mathcal{W}_p(\hat{\mu}_0, P_{M_0, \boldsymbol{\theta}_0}^{(n)})$. Then the ABC-Wass posterior asymptotically concentrates on the model-parameter pair

$$(M^\dagger, \boldsymbol{\theta}^\dagger) = \arg \min_{M_k \in \mathcal{M}, \boldsymbol{\theta}_k \in \Theta_k} \mathcal{W}_p(P_{M_0, \boldsymbol{\theta}_0}^{(n)}, P_{M_k, \boldsymbol{\theta}_k}^{(n)}).$$

Simulation study

- 3 ABC with distance metrics (ABC-CvM, ABC-MMD, and ABC-Wass)
- ABC-Stat
- NN, performing multi-task learning: classification (model selection) with categorical cross-entropy loss and regression (parameter estimation) using MSE loss. Loss weights are 1.0 (classification) and 0.5 (regression), prioritizing model identification. The architecture comprises five fully connected layers.
- ABC-QDA Prangle et al. [2014]
- ABC-SA [Gutmann et al., 2018]

Each experiment is repeated 100 times for sample sizes $n = 100$. For each observed dataset, 10^6 simulations.

Model priors: $\pi(M = M_k) = 1/K$ for $k = 1, \dots, K$.

The threshold ε is chosen as the q -th percentile of the simulated distances [Biau et al., 2015].

Normal example

Let $\mathbf{y} = (y_1, \dots, y_n)^T$ be i.i.d. with $y_i \sim N(\theta, \sigma^2)$, assuming $\sigma^2 = 1$. The competing models correspond to $H_0 : \theta = \tilde{\theta}$ versus $H_1 : \theta \neq \tilde{\theta}$.

Synthetic datasets are generated from the true model $\mathcal{N}(\theta_0, 1)$ with $\theta_0 \in \{0.0, 0.1, 0.2, 0.3, 0.4, 0.5\}$.

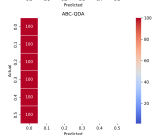
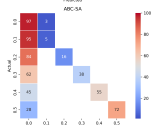
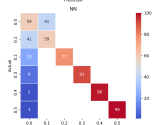
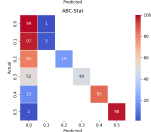
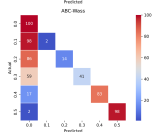
Under H_1 , the prior is $\theta \sim N(\tilde{\theta}, 100)$, providing a weakly informative prior; under H_0 , θ is fixed at $\tilde{\theta} = 0$.

For ABC-Stat, we use the sufficient statistic $\eta(\mathbf{y}) = \bar{y}$ and the Euclidean distance.

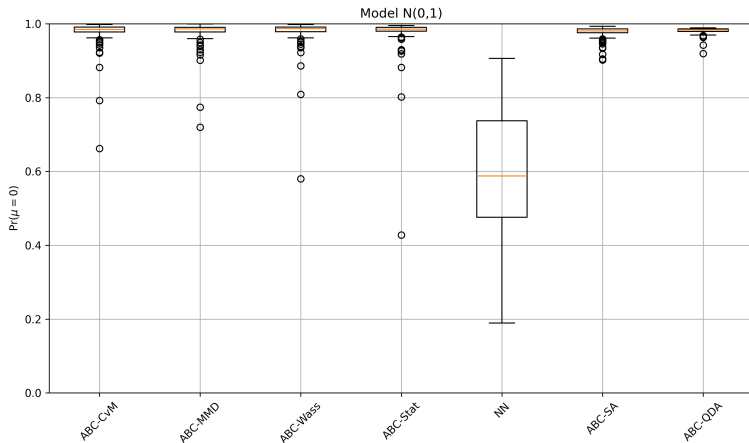
Normal example

		ABC				NN	ABC	
		CvM	MMD	Wass	Stat		SA	QDA
$\theta_0 = 0.0$	$\Pr(\theta = 0)$	0.975	0.976	0.975	0.973	0.599	0.976	0.981
	$\hat{\theta}$	0.000	0.000	0.000	0.002	0.042	0.000	0.000
	MSE	0.000	0.000	0.000	0.001	0.002	0.000	0.000
$\theta_0 = 0.1$	$\Pr(\theta = 0)$	0.931	0.933	0.929	0.412	0.925	0.948	0.974
	$\hat{\theta}$	0.007	0.006	0.007	0.010	0.052	0.000	0.000
	MSE	0.011	0.011	0.011	0.011	0.002	0.010	0.010
$\theta_0 = 0.2$	$\Pr(\theta = 0)$	0.788	0.796	0.789	0.760	0.261	0.847	0.955
	$\hat{\theta}$	0.046	0.042	0.042	0.048	0.155	0.010	0.000
	MSE	0.037	0.037	0.038	0.038	0.013	0.037	0.040
$\theta_0 = 0.3$	$\Pr(\theta = 0)$	0.522	0.543	0.519	0.459	0.084	0.671	0.934
	$\hat{\theta}$	0.147	0.144	0.149	0.168	0.273	0.141	0.000
	MSE	0.062	0.061	0.060	0.056	0.004	0.096	0.090
$\theta_0 = 0.4$	$\Pr(\theta = 0)$	0.210	0.236	0.203	0.157	0.005	0.450	0.902
	$\hat{\theta}$	0.344	0.335	0.345	0.364	0.404	0.285	0.000
	MSE	0.039	0.041	0.038	0.034	0.000	0.125	0.160
$\theta_0 = 0.5$	$\Pr(\theta = 0)$	0.056	0.067	0.050	0.031	0.014	0.263	0.870
	$\hat{\theta}$	0.494	0.489	0.497	0.506	0.484	0.359	0.000
	MSE	0.022	0.025	0.020	0.016	0.003	0.125	0.250

Normal example



Normal example



Normal example

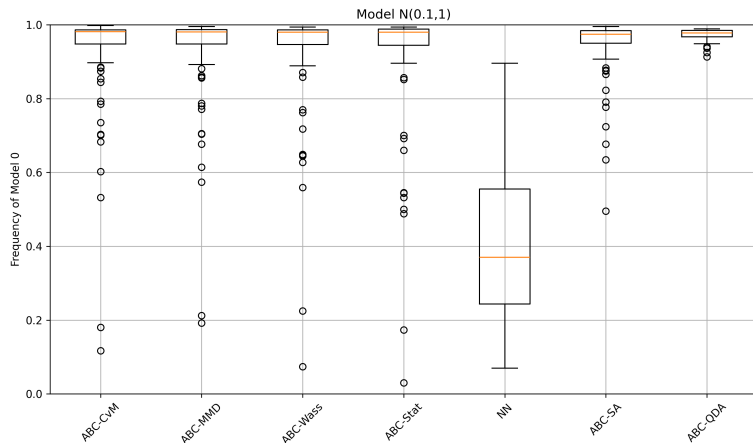


Figure: Boxplots of the estimated posterior probability of H_0 across 100 datasets generated from $N(0.1,1)$.

Normal example

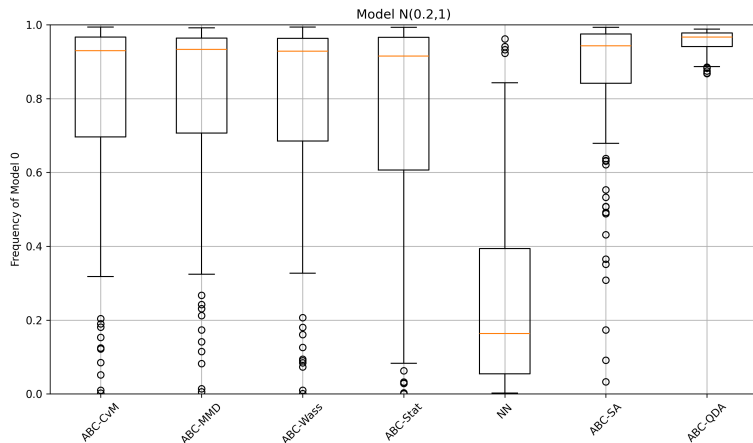


Figure: Boxplots of the estimated posterior probability of H_0 across 100 datasets generated from $N(0.2,1)$.

Normal example

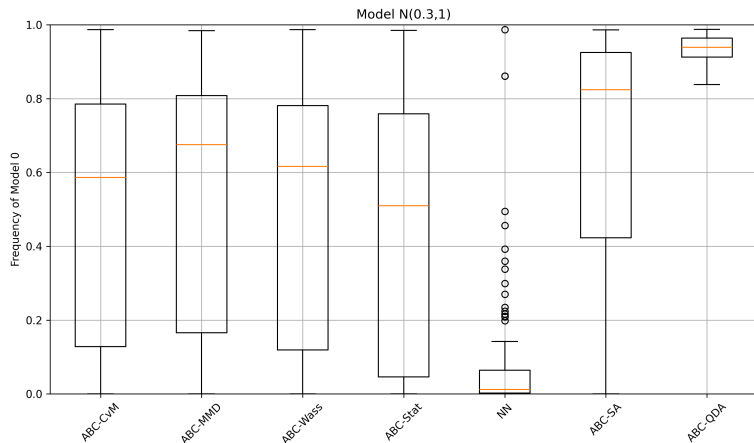


Figure: Boxplots of the estimated posterior probability of H_0 across 100 datasets generated from $N(0.3,1)$.

Normal example

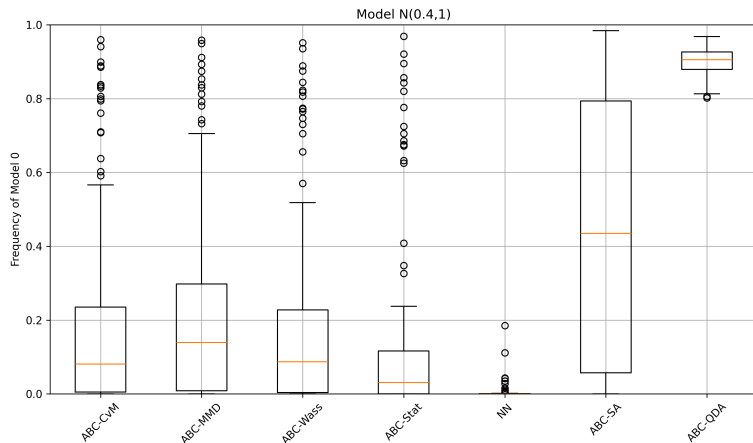


Figure: Boxplots of the estimated posterior probability of H_0 across 100 datasets generated from $N(0.4, 1)$.

Normal example

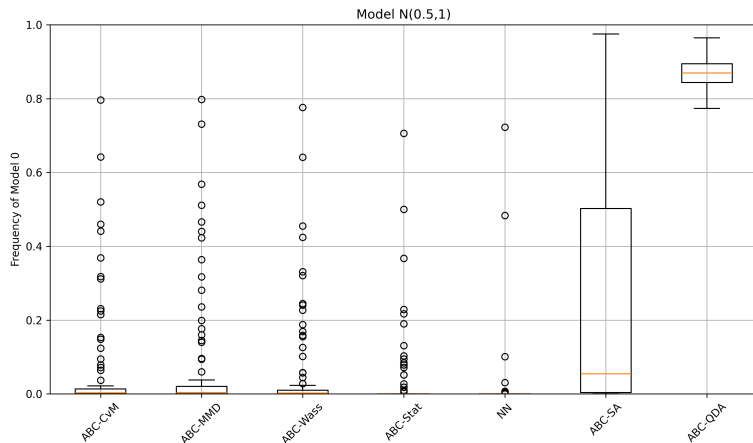


Figure: Boxplots of the estimated posterior probability of H_0 across 100 datasets generated from $N(0.5,1)$.

Exponential family example

Consider three models from the exponential family.

- Model M_1 : $y_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{Exp}(\theta)$, $i = 1, \dots, n$, with prior $\theta \sim \mathcal{Exp}(1)$.
- Model M_2 : $y_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{LN}(\theta, 1)$, with prior $\theta \sim N(0, 1)$.
- Model M_3 : $y_i \stackrel{\text{i.i.d.}}{\sim} \text{Ga}(2, \theta)$, with prior $\theta \sim \mathcal{Exp}(1)$.

Marin et al. [2016] shows that no information is lost by replacing the full dataset with these summaries:

$$\boldsymbol{\eta}(\mathbf{y}) = \left(\sum_{i=1}^n y_i, \sum_{i=1}^n \log y_i, \sum_{i=1}^n \log^2 y_i \right).$$

Exponential family example

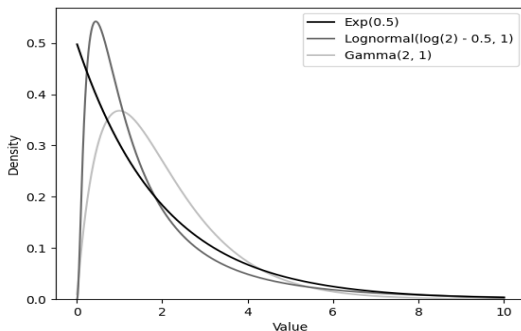


Figure: Comparison of the true data distributions for each model, with parameters set to $\theta = 1/2$ for M_1 , $\theta = \log 2 - 1/2$ for M_2 , and $\theta = 1$ for M_3 . These values give each model a mean of 2.

Exponential family example

		ABC						Stat	NN	ABC	
		CvM	MMD	MMD (log)	Wass	Wass (log)	SA			QDA	
$\mathcal{E}_{\text{xp}}(0.5)$	$\Pr(M = M_1)$	0.883	0.836	0.953	0.850	0.948	0.952	0.554	0.772	0.922	
	$\hat{\theta}$	0.512	0.506	0.511	0.515	0.514	0.513	0.454	0.510	0.522	
	MSE	0.003	0.003	0.003	0.003	0.003	0.003	0.021	0.003	0.0040	
$\mathcal{L}\mathcal{N}(0.193, 1)$	$\Pr(M = M_2)$	0.896	0.829	0.952	0.882	0.956	0.954	0.653	0.787	0.875	
	$\hat{\theta}$	0.188	0.196	0.187	0.189	0.187	0.191	0.337	0.204	0.187	
	MSE	0.009	0.008	0.008	0.008	0.009	0.010	0.034	0.013	0.140	
$\mathcal{G}\mathcal{a}(2, 1)$	$\Pr(M = M_3)$	0.952	0.966	0.971	0.984	0.987	0.987	0.967	0.955	0.819	
	$\hat{\theta}$	1.000	0.995	0.998	1.004	1.003	1.002	0.957	1.017	1.000	
	MSE	0.004	0.004	0.004	0.004	0.004	0.004	0.008	0.008	0.275	

Table: Results for the exponential family models across repetitions. For each model, we report the average posterior probability of selecting the correct model, the average estimate of the parameter θ , and the MSE across repetitions. Logarithm function is taken following Drovandi and Frazier [2022].

Exponential family example

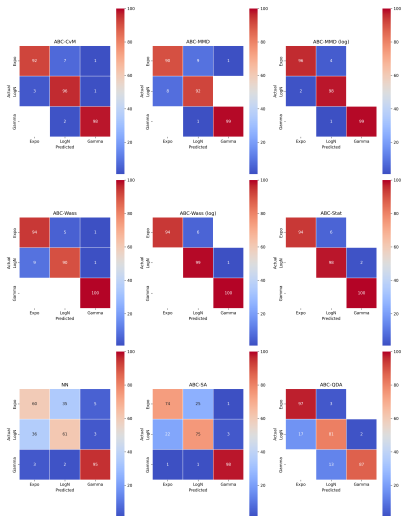


Figure: Confusion matrices for model selection in the exponential family example with $n = 100$.

Exponential family example

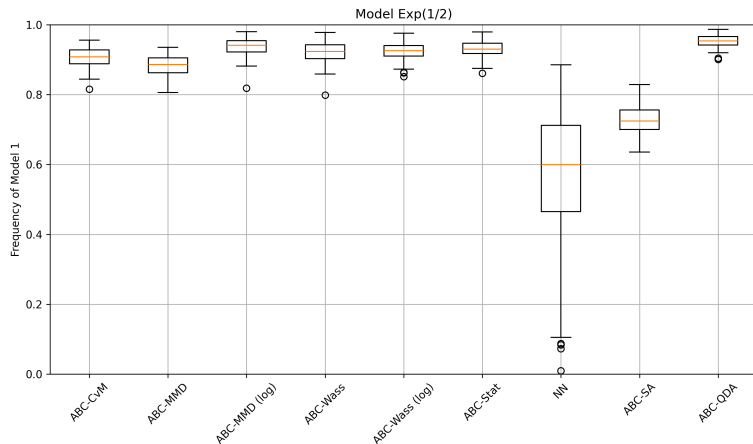


Figure: Boxplots of the estimated posterior probability of model M_1 across 100 datasets generated from $\mathcal{E}^{xp}(1/2)$.

Exponential family example

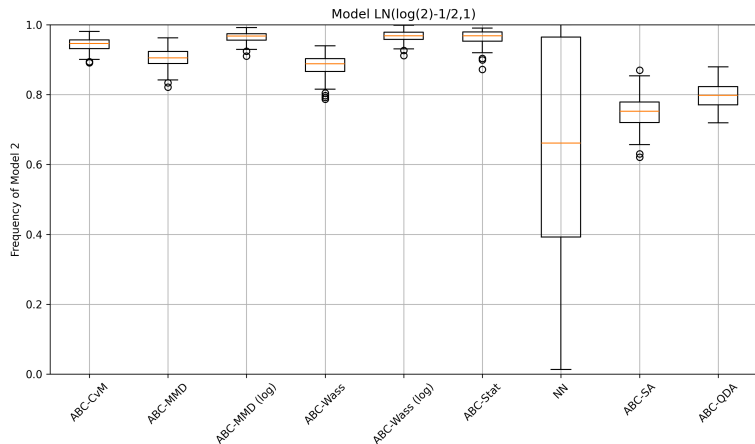


Figure: Boxplots of the estimated posterior probability of model M_2 across 100 datasets generated from $\mathcal{LN}(\log(2) - 1/2, 1)$.

Exponential family example

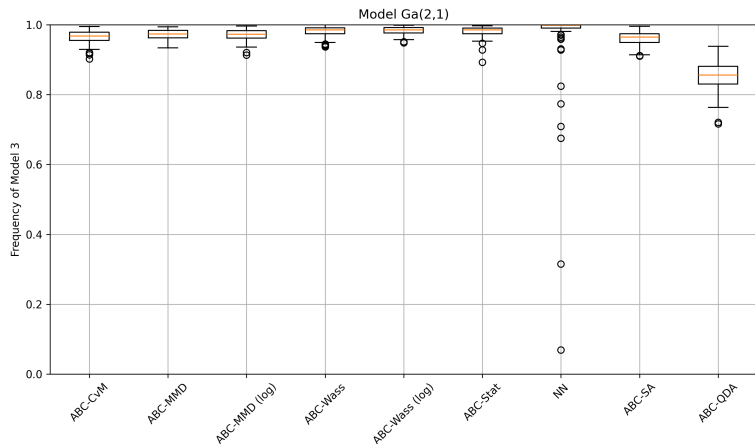
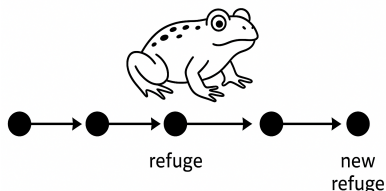


Figure: Boxplots of the estimated posterior probability of model M_3 across 100 datasets generated from $Ga(2,1)$.

Real data example: Toad movement [Drovandi and Frazier, 2022]

Toad Moving Between Refuges

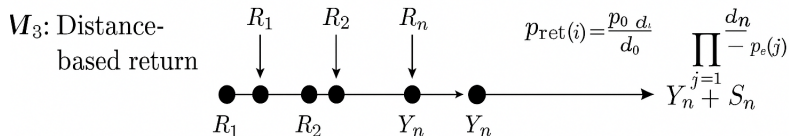
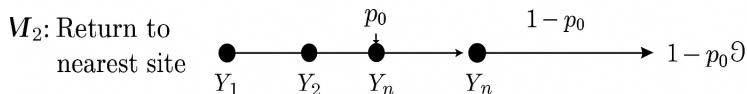
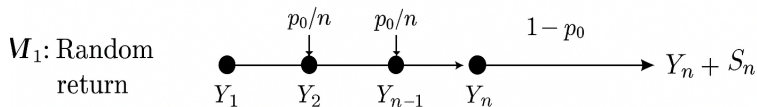


Following Marchand et al. [2017], the toad perform a displacement S_n from their last refuge location Y_n .

The toad either moves to a new location, $Y_{n+1} = Y_n + S_n$, or returns to a previous daytime refuge site (Y_i for some $i = 1, \dots, n-1$).

The displacement follows a symmetric, zero-centered stable distribution, $S_n \sim S(\alpha, \gamma)$

Toad Movement Models



M_3 : Distance-based return

Real data example: Toad movement

- Model M_1 is the random return model

$$Y_{n+1} = \begin{cases} Y_n + S_n & \text{with prob. } 1 - p_0, \\ Y_i & \text{with prob. } p_0/n \quad \forall i = 1, \dots, n. \end{cases}$$

- Model M_2 assumes that, when the toad returns, it always chooses the nearest refuge site.

$$Y_{n+1} = \begin{cases} Y_n + S_n & \text{with prob. } 1 - p_0, \\ \min_{Y \in \{Y_1, \dots, Y_n\}} |Y_{n+1} - Y| & \text{with prob. } p_0. \end{cases}$$

- Model M_3 is a distance-based return model

$$Y_{n+1} = \begin{cases} Y_n + S_n & \text{with prob. } \prod_{j=1}^{A_{n+1}} (1 - p_{ret(j)}), \\ R_i & \text{with prob. } p_i \quad \forall i = 1, \dots, A_{n+1}, \end{cases}$$

Real data example: Toad movement

	$\Pr(M = M_1)$	$\Pr(M = M_2)$	$\Pr(M = M_3)$
ABC-CvM	0.08	0.00	0.92
ABC-MMD	0.29	0.00	0.71
ABC-MMD (log)	0.07	0.00	0.93
ABC-Wass	0.14	0.00	0.86
ABC-Wass (log)	0.00	0.00	1.00
ABC-Stat	0.40	0.00	0.60
ABC-Stat [Marchand et al., 2017]	0.15	0.00	0.85
NN	0.32	0.02	0.66
ABC-SA	0.36	0.00	0.64
ABC-QDA	0.13	0.09	0.79

Table: Estimated posterior probabilities of the three toad movement models, comparing different methods with the original results of Marchand et al. [2017].

This work has investigated ABC as a tool for model selection in settings with intractable likelihoods.

- full data ABC approaches yield posterior distributions closely matching ABC with sufficient statistics (when they are sufficient for model choice problems) or provides high support to the correct model
- the Wasserstein distance consistently produced the most accurate model selection results
- theoretical results establish both the consistency of ABC-Wass in selecting the correct model and its robustness to model misspecification
- future work will extend these results by leveraging the general framework of Legramanti et al. [2025]
- curse of dimensionality?

References I

- Chris Barnes, Sarah Filippi, Michael Stumpf, and Thomas Thorne. Considerate approaches to constructing summary statistics for ABC model selection. *Statistics and Computing*, 22:1181–1197, 11 2012.
- Espen Bernton, Pierre E. Jacob, Mathieu Gerber, and Christian P. Robert. Approximate Bayesian Computation with the Wasserstein distance. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 81(2):235–269, February 2019. ISSN 1467-9868.
- G erard Biau, Fr ed eric C erou, and Arnaud Guyader. New insights into approximate bayesian computation. *Annales de l'Institut Henri Poincar e*, 51(1):376–403, 2015.
- Michael GB Blum and Olivier Fran ois. Non-linear regression models for approximate bayesian computation. *Statistics and Computing*, 20:63–73, 2010.
- C. Drovandi and D.F. Frazier. A comparison of likelihood-free methods with and without summary statistics. *Statistics and Computing*, 32(42), 2022.
- Florence Forbes, Hien Duy Nguyen, TrungTin Nguyen, and Julyan Arbel. Summary statistics and discrepancy measures for approximate bayesian computation via surrogate posteriors. *Statistics and Computing*, 32(5):85, 2022.
- David T. Frazier. Robust and efficient approximate Bayesian computation: A minimum distance approach. arXiv:2006.14126, 2020.

References II

- Andrew Gelman and Cosma Rohilla Shalizi. Philosophy and the practice of Bayesian statistics. *British Journal of Mathematical and Statistical Psychology*, 66(1):8–38, 2013.
- Michael U Gutmann, Ritabrata Dutta, Samuel Kaski, and Jukka Corander. Likelihood-free inference via classification. *Statistics and Computing*, 28:411–425, 2018.
- Bai Jiang, Tung-Yu Wu, Charles Zheng, and Wing H Wong. Learning summary statistic for approximate bayesian computation via deep neural network. *Statistica Sinica*, pages 1595–1618, 2017.
- Sirio Legramanti, Daniele Durante, and Pierre Alquier. Concentration of discrepancy-based ABC via Rademacher complexity. 2025.
- Philippe Marchand, Morgan Boenke, and David M. Green. A stochastic movement model reproduces patterns of site fidelity and long-distance dispersal in a population of Fowler’s toads (*Anaxyrus fowleri*). *Ecological Modelling*, 360:63–69, 2017. ISSN 0304-3800.
- Jean-Michel Marin, Natesh S. Pillai, Christian P. Robert, and Judith Rousseau. Relevant Statistics for Bayesian Model Choice. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 76(5):833–859, 12 2013.

References III

- Jean-Michel Marin, Pierre Pudlo, Arnaud Estoup, and Christian P. Robert. *Likelihood-free Model Choice*, chapter 6. CRC Press Taylor & Francis Group, 2016.
- Gael M. Martin, David T. Frazier, and Christian P. Robert. Approximating Bayes in the 21st century. *Statistical Science*, 39(1):20–45, 2024.
- Christoph Molnar, Gunnar König, Julia Herbinger, Timo Freiesleben, Susanne Dandl, Christian A. Scholbeck, Giuseppe Casalicchio, Moritz Grosse-Wentrup, and Bernd Bischl. *General Pitfalls of Model-Agnostic Interpretation Methods for Machine Learning Models*, pages 39–68. Springer International Publishing, 2022.
- Hien D. Nguyen, Julyan Arbel, Hongliang Lü, and Florence Forbes. Approximate bayesian computation via the energy statistic. *IEEE Access*, 8:131683–131698, 2020.
- Mijung Park, Wittawat Jitkrittum, and Dino Sejdinovic. K2-ABC: Approximate Bayesian computation with kernel embeddings. *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics, PMLR*, 51:398–407, 2016.
- Dennis Prangle, Paul Fearnhead, Murray P. Cox, Patrick J. Biggs, and Nigel P. French. Semi-automatic selection of summary statistics for ABC model choice. *Statistical Applications in Genetics and Molecular Biology*, 13(1):67–82, 2014.

- C.P. Robert, J. Cornuet, J. Marin, and N.S. Pillai. Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences*, 108(37):15112–15117, 2011.
- Sara Sheehan and Yun S Song. Deep learning for population genetic inference. *PLoS Computational Biology*, 12(3):e1004845, 2016.