

# Evidence estimation in finite and infinite mixture models

Christian P. Robert

U. Paris Dauphine & Warwick U.



Joint work with A. Hairault and J. Rousseau

BNP Monash, December 3, 2023

Ongoing 2023-2030 ERC funding for PhD positions and postdoctoral collaborations with

- ▶ Michael Jordan (Paris  $\geq$  Berkeley)
- ▶ Eric Moulines (Paris)
- ▶ Gareth Roberts (Warwick)
- ▶ myself (Paris  $\geq$  Warwick)

on OCEAN (On Intelligence And Networks) project



# Outline

- 1 Mixtures of distributions
- 2 Approximations to evidence
- 3 Dirichlet process mixtures
- 4 Distributed evidence evaluation



# Mixtures of distributions

Convex combination of densities

$x \sim f_j$  with probability  $p_j$ ,

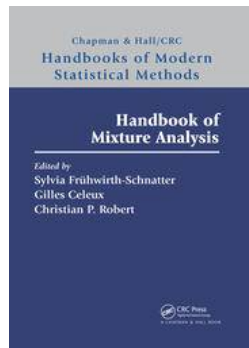
for  $j = 1, 2, \dots, k$ , with overall density

$$f^k(x; \mathbf{p}, \boldsymbol{\vartheta}) \equiv p_1 f_1(x) + \dots + p_k f_k(x)$$

Usual case: parameterised components

$$\sum_{i=1}^k p_i f(x|\vartheta_i) \quad \text{with} \quad \sum_{i=1}^n p_i = 1$$

where *weights*  $p_i$ 's are distinguished from other parameters



# Jeffreys priors for mixtures

True Jeffreys (1939) prior for mixtures of distributions defined from information matrix as

$$\left| \mathbb{E}_{\vartheta} [\nabla^{\top} \nabla \log f(X|\vartheta)] \right|^{1/2}$$

- ▶  $O(k)$  matrix
- ▶ unavailable in closed form except for special cases
- ▶ unidimensional integrals approximated by Monte Carlo tools

[Grazian & X, 2015]

# Difficulties

- ▶ complexity grows in  $O(k^3)$
- ▶ significant computing requirement (reduced by delayed acceptance)

[Banterle et al., 2014]

- ▶ differ from component-wise Jeffreys

[Diebolt & X, 1990; Stoneking, 2014]

- ▶ when is the posterior proper?
- ▶ how to check properness via MCMC outputs?

## Further reference priors

Reparameterisation of a location-scale mixture in terms of its global **mean**  $\mu$  and global **variance**  $\sigma^2$  as

$$\mu_i = \mu + \sigma\alpha_i \quad \text{and} \quad \sigma_i = \sigma\tau_i \quad 1 \leq i \leq k$$

where  $\tau_i > 0$  and  $\alpha_i \in \mathbb{R}$

Motivation: induced **compact** space on other parameters:

$$\sum_{i=1}^k p_i \alpha_i = 0 \quad \text{and} \quad \sum_{i=1}^k p_i \tau_i^2 + \sum_{i=1}^k p_i \alpha_i^2 = 1$$

© Posterior associated with prior  $\pi(\mu, \sigma) = 1/\sigma$  proper for Gaussian components for (at least) two observations in sample

[Kamary, Lee & X, 2018]

# Label switching paradox

$$p_1 f(x|\vartheta_1) + p_2 f(x|\vartheta_2) \equiv p_2 f(x|\vartheta_2) + p_1 f(x|\vartheta_1) \quad (!!!)$$

- ▶ Under exchangeability, **should** observe exchangeability of the components **[label switching]** to conclude about MCMC convergence
- ▶ If observed, how should we estimate parameters?
- ▶ If unobserved, uncertainty about MCMC convergence

[Celeux, Hurn & X, 2000; Frühwirth-Schnatter, 2001, 2004; Jasra & al., 2005]

[Unless adopting a point process perspective]

[Green, 2019]



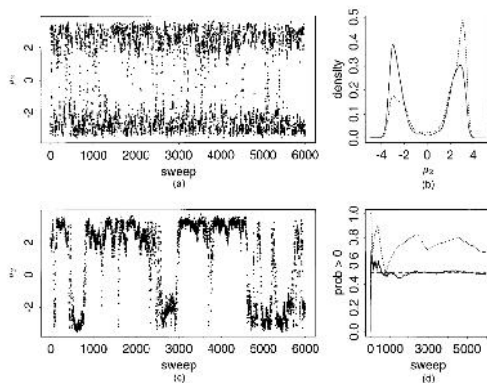


Fig. 9. Comparison of mixing of variable  $k$  and fixed  $k$  samplers: (a), (c) traces of  $\mu_j$  against sweep number; (b) posterior density estimates at the end of the runs; (d) sequences of estimates of  $p(\mu_k < 0|y)$ ,  $k = 3$  obtained as the runs proceed (—, variable  $k$  sampler)

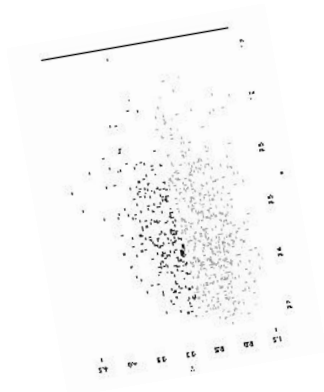
# Loss functions for mixture estimation

Global loss function that considers distance between predictives

$$L(\xi, \hat{\xi}) = \int_{\mathcal{X}} f_{\xi}(x) \log \{ f_{\xi}(x) / f_{\hat{\xi}}(x) \} dx$$

eliminates the labelling effect  
Similar solution for estimating clusters through allocation variables

$$L(z, \hat{z}) = \sum_{i < j} \left[ \mathbb{I}_{[z_i = z_j]} (1 - \mathbb{I}_{[\hat{z}_i = \hat{z}_j]}) + \mathbb{I}_{[\hat{z}_i = \hat{z}_j]} (1 - \mathbb{I}_{[z_i = z_j]}) \right] .$$



[Celeux, Hurn & X, 2000]

# Bayesian model choice

Comparison of models  $\mathfrak{M}_i$  by Bayesian methods:

probabilise the entire model/parameter space

- ▶ allocate probabilities  $p_i$  to all models  $\mathfrak{M}_i$
- ▶ define priors  $\pi_i(\vartheta_i)$  for each parameter space  $\Theta_i$
- ▶ compute

$$\pi(\mathfrak{M}_i|x) = \frac{p_i \int_{\Theta_i} f_i(x|\vartheta_i) \pi_i(\vartheta_i) d\vartheta_i}{\sum_j p_j \int_{\Theta_j} f_j(x|\vartheta_j) \pi_j(\vartheta_j) d\vartheta_j}$$

Computational difficulty on its own

[Chen, Shao & Ibrahim, 2000; Marin & X, 2007]

# Bayesian model choice

Comparison of models  $\mathfrak{M}_i$  by Bayesian methods:

Relies on marginals

$$m_k(\mathbf{x}) = \int_{\Theta_k} \pi_k(\vartheta_k) L_k(\vartheta_k | \mathbf{x}) d\vartheta_k,$$

aka the marginal likelihood.

Computational difficulty on its own

[Chen, Shao & Ibrahim, 2000; Marin & X, 2007]

# Bayesian model comparison

Bayes Factor consistent for selecting number of components

[Ishwaran et al., 2001; Casella & Moreno, 2009; Chib and Kuffner, 2016]

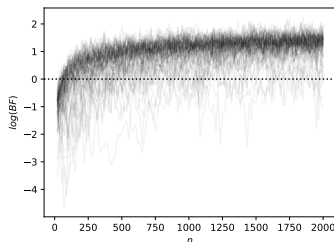
Bayes Factor consistent for testing parametric versus nonparametric alternatives

[Verdinelli & Wasserman, 1997; Dass & Lee, 2004; McVinish et al., 2009]

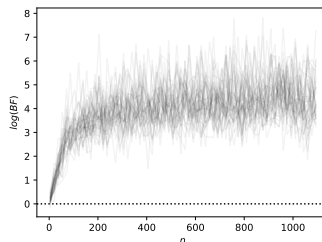
# Consistent evidence for location DPM

Consistency of Bayes factor comparing finite mixtures against (location) Dirichlet Process Mixture

$$H_0 : f_0 \in \mathfrak{M}_K \text{ vs. } H_1 : f_0 \notin \mathfrak{M}_K$$



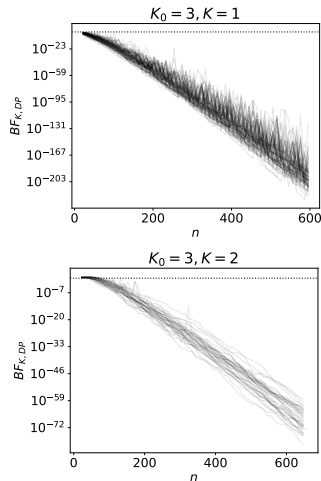
(left)  $K_0 = K = 1$



(right)  $K_0 = K = 3$

# Consistent evidence for location DPM

Consistency of Bayes factor comparing finite mixtures against (location) Dirichlet Process Mixture



# Consistent evidence for location DPM

Under assumptions [next], when  $x_1, \dots, x_n$  iid  $f_{p_0}$  with

$$p_0 = \sum_{j=1}^{k_0} p_j^0 \delta_{\vartheta_j^0}$$

and Dirichlet  $DP(M, G_0)$  prior on  $P$ , there exists  $t > 0$  such that for all  $\varepsilon > 0$

$$\mathbb{P}_{f_0} \left( m_{DP}(\mathbf{x}) > n^{-(k_0-1+dk_0+t)/2} \right) = o(1)$$

Moreover there exists  $q \geq 0$  such that

$$\Pi_{DP} \left( \|f_0 - f_p\|_1 \leq \frac{(\log n)^q}{\sqrt{n}} \middle| \mathbf{x} \right) = 1 + o_{P_{f_0}}(1)$$

[Hairault, X & Rousseau, 2022]



DATA IN FSL



European  
Research  
Council



# Consistent evidence for location DPM

**Assumption A1** [Regularity]

**Assumption A2** [Strong identifiability]

**Assumption A3** [Compactness]

**Assumption A4** [Existence of DP random mean]

**Assumption A5** [Truncated support of  $M$ , e.g. trunc'd  $\mathcal{G}_\alpha$ ]

(i) If  $f_{P_0} \in \mathfrak{M}_{k_0}$  satisfies Assumptions **A1–A5**, then

$$m_{k_0}(\mathbf{x})/m_{DP}(\mathbf{x}) \rightarrow \infty \text{ under } f_{P_0}$$

(ii) Moreover for all  $k \geq k_0$ , if Dirichlet parameter  $\alpha = \eta/k$  and  $\eta < kd/2$ , then

$$m_k(\mathbf{x})/m_{DP}(\mathbf{x}) \rightarrow \infty \text{ under } f_{P_0}$$

# Consistent evidence for location DPM

(i) If  $f_{P_0} \in \mathfrak{M}_{k_0}$  satisfies Assumptions **A1–A5**, then

$$m_{k_0}(\mathbf{x})/m_{DP}(\mathbf{x}) \rightarrow \infty \text{ under } f_{P_0}$$

(ii) Moreover for all  $k \geq k_0$ , if Dirichlet parameter  $\alpha = \eta/k$  and  $\eta < kd/2$ , then

$$m_k(\mathbf{x})/m_{DP}(\mathbf{x}) \rightarrow \infty \text{ under } f_{P_0}$$

(iii) If  $\inf_{f_P \in \mathfrak{M}_{k_0}} \text{KL}(f_{P_0}, f_P) > 0$  and the DP prior verifies  $\Pi_{DP}(\text{KL}(f_{P_0}, f_P) \leq \varepsilon) > 0$  for all  $\varepsilon > 0$ , then

$$m_{k_0}(\mathbf{y})/m_{DP}(\mathbf{y}) \rightarrow 0 \text{ under } f_{P_0}$$

# Outline

- 1 Mixtures of distributions
- 2 Approximations to evidence
- 3 Dirichlet process mixtures
- 4 Distributed evidence evaluation



# Chib's or candidate's representation

Direct application of Bayes' theorem: given  $\mathbf{x} \sim f_k(\mathbf{x}|\vartheta_k)$  and  $\vartheta_k \sim \pi_k(\vartheta_k)$ ,

$$\mathfrak{z}_k = \mathfrak{m}_k(\mathbf{x}) = \frac{f_k(\mathbf{x}|\vartheta_k) \pi_k(\vartheta_k)}{\pi_k(\vartheta_k|\mathbf{x})}$$

Replace with an approximation to the posterior

$$\hat{\mathfrak{z}}_k = \widehat{\mathfrak{m}}_k(\mathbf{x}) = \frac{f_k(\mathbf{x}|\vartheta_k^*) \pi_k(\vartheta_k^*)}{\hat{\pi}_k(\vartheta_k^*|\mathbf{x})}.$$

[Besag, 1989; Chib, 1995]

# Natural Rao-Blackwellisation

For missing variable  $\mathbf{z}$  as in mixture models, natural Rao-Blackwell (unbiased) estimate

$$\widehat{\pi}_k(\vartheta_k^*|\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \pi_k(\vartheta_k^*|\mathbf{x}, \mathbf{z}_k^{(t)}),$$

where the  $\mathbf{z}_k^{(t)}$ 's are Gibbs sampled latent variables

[Diebolt & X, 1990; Chib, 1995]

# Compensation for label switching

For mixture models,  $\mathbf{z}_k^{(t)}$  usually fails to visit all configurations, despite symmetry predicted by theory

Significant consequences on numerical approximation, biased by an order  $k!$

# Compensation for label switching

For mixture models,  $\mathbf{z}_k^{(t)}$  usually fails to visit all configurations, despite symmetry predicted by theory

Force predicted theoretical symmetry by using

$$\widetilde{\pi}_k(\vartheta_k^*|\mathbf{x}) = \frac{1}{T k!} \sum_{\sigma \in \mathfrak{S}_k} \sum_{t=1}^T \pi_k(\sigma(\vartheta_k^*)|\mathbf{x}, \mathbf{z}_k^{(t)}) .$$

for all  $\sigma$ 's in  $\mathfrak{S}_k$ , set of all permutations of  $\{1, \dots, k\}$

[Neal, 1999; Berkhof, Mechelen, & Gelman, 2003; Lee & X, 2018]

Benchmark galaxies for radial velocities of 82 galaxies

[Postman et al., 1986; Roeder, 1992; Raftery, 1996]

Conjugate priors for Gaussian components

$$\sigma_k^2 \sim \Gamma^{-1}(a_0, b_0)$$

$$\mu_k | \sigma_k^2 \sim \mathcal{N}(\mu_0, \sigma_k^2 / \lambda_0)$$





## Galaxy dataset (k)

Using Chib's estimate, with  $\vartheta_k^*$  as MAP estimator,

$$\log(\hat{\mathfrak{Z}}_k(\mathbf{x})) = -105.1396$$

for  $k = 3$ , while introducing permutations leads to

$$\log(\hat{\mathfrak{Z}}_k(\mathbf{x})) = -103.3479$$

Perfect difference:

$$-105.1396 + \log(3!) = -103.3479$$

k	2	3	4	5	6	7	8
$\mathfrak{Z}_k(\mathbf{x})$	-115.68	-103.35	-102.66	-101.93	-102.88	-105.48	-108.44

Estimations of the marginal likelihoods by the symmetrised Chib's approximation (based on  $10^5$  Gibbs iterations and, for  $k > 5$ , 100 permutations selected at random in  $\mathfrak{S}_k$ ).

# Rethinking Chib's solution

Alternate Rao–Blackwellisation by marginalising into partitions  
Apply **candidate's/Chib's formula** to a chosen partition:

$$m_k(\mathbf{x}) = \frac{f_k(\mathbf{x}|\mathfrak{C}^0)\pi_k(\mathfrak{C}^0)}{\pi_k(\mathfrak{C}^0|\mathbf{x})}$$

with

$$\pi_k(\mathfrak{C}(\mathbf{z})) = \frac{k!}{(k - k_+)!} \frac{\Gamma\left(\sum_{j=1}^k \alpha_j\right)}{\Gamma\left(\sum_{j=1}^k \alpha_j + n\right)} \prod_{j=1}^k \frac{\Gamma(n_j + \alpha_j)}{\Gamma(\alpha_j)}$$

$\mathfrak{C}(\mathbf{z})$  partition of  $\{1, \dots, n\}$  induced by cluster membership  $\mathbf{z}$

$n_j = \sum_{i=1}^n \mathbb{I}_{\{z_i=j\}}$  # observations assigned to cluster  $j$

$k_+ = \sum_{j=1}^k \mathbb{I}_{\{n_j>0\}}$  # non-empty clusters

# Rethinking Chib's solution

Alternate Rao–Blackwellisation by marginalising into partitions  
Apply **candidate's/Chib's formula** to a chosen partition:

$$m_k(\mathbf{x}) = \frac{f_k(\mathbf{x}|\mathfrak{C}^0)\pi_k(\mathfrak{C}^0)}{\pi_k(\mathfrak{C}^0|\mathbf{x})}$$

with

$$\pi_k(\mathfrak{C}(\mathbf{z})) = \frac{k!}{(k - k_+)!} \frac{\Gamma\left(\sum_{j=1}^k \alpha_j\right)}{\Gamma\left(\sum_{j=1}^k \alpha_j + n\right)} \prod_{j=1}^k \frac{\Gamma(n_j + \alpha_j)}{\Gamma(\alpha_j)}$$

$\mathfrak{C}(\mathbf{z})$  partition of  $\{1, \dots, n\}$  induced by cluster membership  $\mathbf{z}$

$n_j = \sum_{i=1}^n \mathbb{I}_{\{z_i=j\}}$  # observations assigned to cluster  $j$

$k_+ = \sum_{j=1}^k \mathbb{I}_{\{n_j > 0\}}$  # non-empty clusters

# Rethinking Chib's solution

Under conjugate prior  $G_0$  on  $\vartheta$ ,

$$f_k(\mathbf{x}|\mathfrak{C}(z)) = \prod_{j=1}^k \underbrace{\int_{\Theta} \prod_{i: z_i=k} f(\mathbf{x}_i|\vartheta) G_0(d\vartheta)}_{m(\mathfrak{C}_k(z))}$$

and

$$\hat{\pi}_k(\mathfrak{C}^0|\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \mathbb{I}_{\mathfrak{C}^0 \equiv \mathfrak{C}(z^{(t)})}$$

- ▶ considerably lower computational demand
- ▶ no label switching issue
- ▶ further Rao-Blackwellisation?

# Bridge sampling

Iterative bridge sampling:

$$\hat{\mathbf{e}}^{(t)}(\mathbf{k}) = \hat{\mathbf{e}}^{(t-1)}(\mathbf{k}) M_1^{-1} \sum_{l=1}^{M_1} \frac{\hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})}{M_1 q(\tilde{\vartheta}^l) + M_2 \hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})} / \\ M_2^{-1} \sum_{m=1}^{M_2} \frac{q(\hat{\vartheta}^m)}{M_1 q(\hat{\vartheta}^m) + M_2 \hat{\pi}(\hat{\vartheta}^m | \mathbf{x})}$$

[Gelman& Meng, 1998;Frühwirth-Schnatter, 2004]

where [for mixtures]

$$\tilde{\vartheta}^{1:M_1} \sim q(\vartheta) \quad \text{and} \quad \hat{\vartheta}^{1:M_2} \sim \pi(\vartheta)$$

# Bridge sampling

Iterative bridge sampling:

$$\hat{\mathbf{e}}^{(t)}(k) = \hat{\mathbf{e}}^{(t-1)}(k) M_1^{-1} \sum_{l=1}^{M_1} \frac{\hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})}{M_1 q(\tilde{\vartheta}^l) + M_2 \hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})} /$$
$$M_2^{-1} \sum_{m=1}^{M_2} \frac{q(\hat{\vartheta}^m)}{M_1 q(\hat{\vartheta}^m) + M_2 \hat{\pi}(\hat{\vartheta}^m | \mathbf{x})}$$

[Gelman & Meng, 1998; Frühwirth-Schnatter, 2004]

where

$$q(\vartheta) = \frac{1}{J_1} \sum_{j=1}^{J_1} p(\lambda | \mathbf{z}^{(j)}) \prod_{i=1}^k p(\xi_i | \xi_{i < j}^{(j)}, \xi_{i > i}^{(j-1)}, \mathbf{z}^{(j)}, \mathbf{x})$$

# Bridge sampling

Iterative bridge sampling:

$$\hat{\mathbf{e}}^{(t)}(\mathbf{x}) = \hat{\mathbf{e}}^{(t-1)}(\mathbf{x}) M_1^{-1} \sum_{l=1}^{M_1} \frac{\hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})}{M_1 q(\tilde{\vartheta}^l) + M_2 \hat{\pi}(\tilde{\vartheta}^l | \mathbf{x})} / \\ M_2^{-1} \sum_{m=1}^{M_2} \frac{q(\hat{\vartheta}^m)}{M_1 q(\hat{\vartheta}^m) + M_2 \hat{\pi}(\hat{\vartheta}^m | \mathbf{x})}$$

[Gelman & Meng, 1998; Frühwirth-Schnatter, 2004]

where

$$q(\vartheta) = \frac{1}{k!} \sum_{\sigma \in \mathfrak{S}(k)} p(\lambda | \sigma(\mathbf{z}^0)) \prod_{i=1}^k p(\xi_i | \sigma(\xi_{j \neq i}^0), \sigma(\mathbf{z}^0), \mathbf{x})$$

# Sparsity for permutations

Contribution of each term relative to  $q(\vartheta)$

$$\eta_{\sigma}(\vartheta) = \frac{h_{\sigma}(\vartheta)}{k!q(\vartheta)} = \frac{h_{\sigma_i}(\vartheta)}{\sum_{\sigma \in \mathfrak{S}_k} h_{\sigma}(\vartheta)}$$

and (unnormalised) importance of permutation  $\sigma$  evaluated by

$$\hat{\mathbb{E}}_{h_{\sigma_c}}[\eta_{\sigma_i}(\vartheta)] = \frac{1}{M} \sum_{l=1}^M \eta_{\sigma_i}(\vartheta^{(l)}) , \quad \vartheta^{(l)} \sim h_{\sigma_c}(\vartheta)$$

Approximate set  $\mathfrak{A}(k) \subseteq \mathfrak{S}(k)$  consist of  $[\sigma_1, \dots, \sigma_n]$  for the smallest  $n$  that satisfies the condition

$$\hat{\varphi}_n = \frac{1}{M} \sum_{l=1}^M \left| \tilde{q}_n(\vartheta^{(l)}) - q(\vartheta^{(l)}) \right| < \tau$$



# dual importance sampling with approximation

## DIS2A

- 1 Randomly select  $\{\mathbf{z}^{(j)}, \vartheta^{(j)}\}_{j=1}^J$  from Gibbs sample and un-switch  
Construct  $q(\vartheta)$
- 2 Choose  $h_{\sigma_c}(\vartheta)$  and generate particles  $\{\vartheta^{(t)}\}_{t=1}^T \sim h_{\sigma_c}(\vartheta)$
- 3 Construction of approximation  $\tilde{q}(\vartheta)$  using first  $M$ -sample
  - 3.1 Compute  $\widehat{\mathbb{E}}_{h_{\sigma_c}}[\eta_{\sigma_1}(\vartheta)], \dots, \widehat{\mathbb{E}}_{h_{\sigma_c}}[\eta_{\sigma_k}(\vartheta)]$
  - 3.2 Reorder the  $\sigma$ 's such that  
 $\widehat{\mathbb{E}}_{h_{\sigma_c}}[\eta_{\sigma_1}(\vartheta)] \geq \dots \geq \widehat{\mathbb{E}}_{h_{\sigma_c}}[\eta_{\sigma_k}(\vartheta)]$ .
  - 3.3 Initially set  $n = 1$  and compute  $\tilde{q}_n(\vartheta^{(t)})$ 's and  $\widehat{\varphi}_n$ . If  $\widehat{\varphi}_{n=1} < \tau$ , go to Step 4. Otherwise increase  $n = n + 1$
- 4 Replace  $q(\vartheta^{(1)}), \dots, q(\vartheta^{(T)})$  with  $\tilde{q}(\vartheta^{(1)}), \dots, \tilde{q}(\vartheta^{(T)})$  to estimate  $\widehat{\mathfrak{E}}$

[Lee & X, 2014]

k	k!	$ \overline{\mathfrak{A}(k)} $	$\overline{\Delta}(\mathfrak{A})$
3	6	1.0000	0.1675
4	24	2.7333	0.1148

Fishery data

k	k!	$ \overline{\mathfrak{A}(k)} $	$\overline{\Delta}(\mathfrak{A})$
3	6	1.000	0.1675
4	24	15.7000	0.6545
6	720	298.1200	0.4146

Galaxy data

Table: Mean estimates of approximate set sizes,  $|\mathfrak{A}(k)|$ , and the reduction rate of a number of evaluated h-terms  $\Delta(\mathfrak{A})$  for (a) fishery and (b) galaxy datasets

# Sequential Monte Carlo

Tempered sequence of targets ( $t = 1, \dots, T$ )

$$\pi_{kt}(\vartheta_k) \propto p_{kt}(\vartheta_k) = \pi_k(\vartheta_k) f_k(\mathbf{x}|\vartheta_k)^{\lambda_t} \quad \lambda_1 = 0 < \dots < \lambda_T = 1$$

particles (simulations) ( $i = 1, \dots, N_t$ )

$$\vartheta_t^i \stackrel{\text{i.i.d.}}{\sim} \pi_{kt}(\vartheta_k)$$

usually obtained by MCMC step

$$\vartheta_t^i \sim K_t(\vartheta_{t-1}^i, \vartheta)$$

with importance weights ( $i = 1, \dots, N_t$ )

$$\omega_i^t = f_k(\mathbf{x}|\vartheta_k)^{\lambda_t - \lambda_{t-1}}$$

[Del Moral et al., 2006; Buchholz et al., 2021]

# Sequential Monte Carlo

Tempered sequence of targets ( $t = 1, \dots, T$ )

$$\pi_{kt}(\vartheta_k) \propto p_{kt}(\vartheta_k) = \pi_k(\vartheta_k) f_k(\mathbf{x}|\vartheta_k)^{\lambda_t} \quad \lambda_1 = 0 < \dots < \lambda_T = 1$$

Produces approximation of evidence

$$\hat{\mathfrak{J}}_k = \prod_t \frac{1}{N_t} \sum_{i=1}^{N_t} \omega_i^t$$

[Del Moral et al., 2006; Buchholz et al., 2021]

# Sequential<sup>2</sup> imputation

For conjugate priors, (marginal) particle filter representation of a proposal:

$$\pi^*(\mathbf{z}|\mathbf{x}) = \pi(z_1|\mathbf{x}_1) \prod_{i=2}^n \pi(z_i|\mathbf{x}_{1:i}, \mathbf{z}_{1:i-1})$$

with importance weight

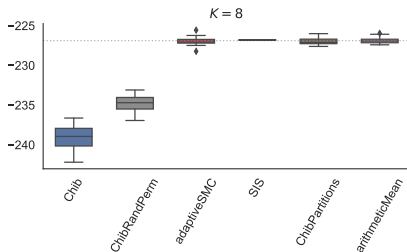
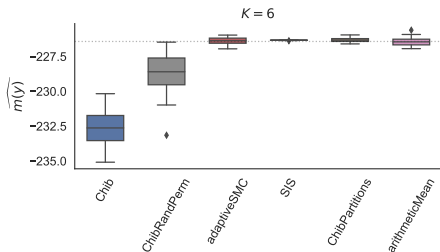
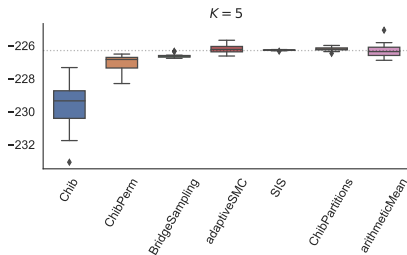
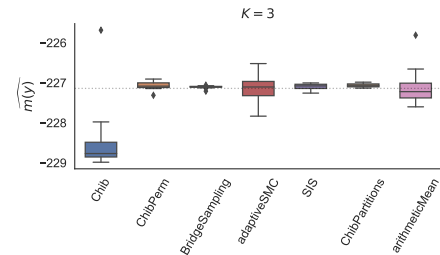
$$\frac{\pi(\mathbf{z}|\mathbf{x})}{\pi^*(\mathbf{z}|\mathbf{x})} = \frac{\pi(\mathbf{x}, \mathbf{z})}{m(\mathbf{x})} \frac{m(\mathbf{x}_1)}{\pi(z_1, \mathbf{x}_1)} \frac{m(z_1, \mathbf{x}_1, \mathbf{x}_2)}{\pi(z_1, \mathbf{x}_1, z_2, \mathbf{x}_2)} \dots \frac{\pi(\mathbf{z}_{1:n-1}, \mathbf{x})}{\pi(\mathbf{z}, \mathbf{x})} = \frac{w(\mathbf{z}, \mathbf{x})}{m(\mathbf{x})}$$

leading to unbiased estimator of evidence

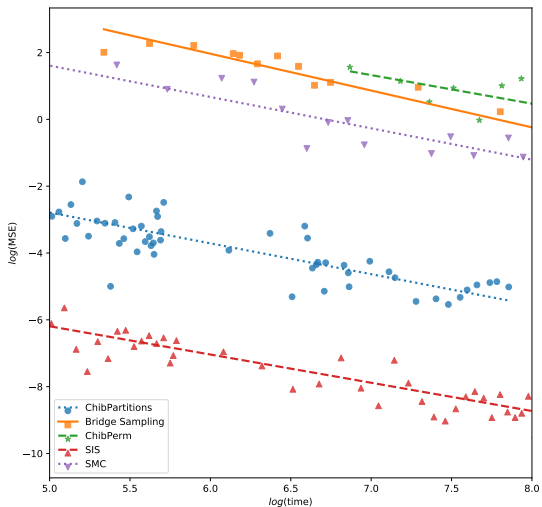
$$\hat{\mathfrak{Z}}_k(\mathbf{x}) = \frac{1}{T} \sum_{i=1}^T w(\mathbf{z}^{(i)}, \mathbf{x})$$

[Long, Liu & Wong, 1994; Carvalho et al., 2010]

# Galactic illustration



# Common illustration



# Empirical conclusions

- ▶ Bridge sampling, arithmetic mean and original Chib's method eventually fail to **scale with  $n$** , sample size
- ▶ Partition Chib's increasingly variable with  $k$ , number of components
- ▶ Adaptive SMC ultimately fails
- ▶ SIS remains most reliable method



- 1 Mixtures of distributions
- 2 Approximations to evidence
- 3 Dirichlet process mixtures
- 4 Distributed evidence evaluation



# Dirichlet process mixture (DPM)

Extension to the  $k = \infty$  (non-parametric) case

$$\begin{aligned}x_i|z_i, \boldsymbol{\vartheta} &\stackrel{\text{i.i.d}}{\sim} f(x_i|\vartheta_{x_i}), \quad i = 1, \dots, n \\ \mathbb{P}(Z_i = k) &= \pi_k, \quad k = 1, 2, \dots \\ \pi_1, \pi_2, \dots &\sim \text{GEM}(M) \quad M \sim \pi(M) \\ \vartheta_1, \vartheta_2, \dots &\stackrel{\text{i.i.d}}{\sim} G_0\end{aligned} \tag{1}$$

with GEM (**G**riffith-**E**ngen-**M**cCloskey) defined by the stick-breaking representation

$$\pi_k = v_k \prod_{i=1}^{k-1} (1 - v_i) \quad v_i \sim \text{Beta}(1, M)$$

[Sethuraman, 1994]

# Dirichlet process mixture (DPM)

Resulting in an infinite mixture

$$\mathbf{x} \sim \prod_{i=1}^n \sum_{i=1}^{\infty} \pi_i f(\mathbf{x}_i | \vartheta_i)$$

with (prior) cluster allocation

$$\pi(\mathbf{z} | \mathbf{M}) = \frac{\Gamma(\mathbf{M})}{\Gamma(\mathbf{M} + \mathbf{n})} \mathbf{M}^{K_+} \prod_{j=1}^{K_+} \Gamma(n_j)$$

and conditional likelihood

$$p(\mathbf{x} | \mathbf{z}, \mathbf{M}) = \prod_{j=1}^{K_+} \int \prod_{i: z_i=j} f(\mathbf{x}_i | \vartheta_j) dG_0(\vartheta_j)$$

available in closed form when  $G_0$  conjugate

# Approximating the evidence

Extension of Chib's formula by marginalising over  $\mathbf{z}$  and  $\vartheta$

$$m_{DP}(\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{M}^*, G_0)\pi(\mathbf{M}^*)}{\pi(\mathbf{M}^*|\mathbf{x})}$$

and using estimate

$$\hat{\pi}(\mathbf{M}^*|\mathbf{x}) = \frac{1}{T} \sum_{t=1}^T \pi(\mathbf{M}^*|\mathbf{x}, \eta^{(t)}, K_+^{(t)})$$

provided prior on  $M$  a  $\Gamma(a, b)$  distribution since

$$M|\mathbf{x}, \eta, K_+ \sim \omega \Gamma(a+K_+, b-\log(\eta)) + (1-\omega) \Gamma(a+K_+-1, b-\log(\eta))$$

with  $\omega = (a + K_+ - 1) / \{n(b - \log(\eta)) + a + K_+ - 1\}$  and  
 $\eta|\mathbf{x}, M \sim \text{Beta}(M + 1, n)$

# Approximating the likelihood

Intractable likelihood  $p(\mathbf{x}|M^*, G_0)$  approximated by sequential  
imputation importance sampling  
Generating  $\mathbf{z}$  from the proposal

$$\pi^*(\mathbf{z}|\mathbf{x}, M) = \prod_{i=1}^n \pi(z_i|\mathbf{x}_{1:i}, \mathbf{z}_{1:i-1}, M)$$

and using the approximation

$$\hat{L}(\mathbf{x}|M^*, G_0) = \frac{1}{T} \sum_{t=1}^T \hat{p}(\mathbf{x}_1|\mathbf{z}_1^{(t)}, G_0) \prod_{i=2}^n p(y_i|\mathbf{x}_{1:i-1} \mathbf{z}_{1:i-1}^{(t)}, G_0)$$

[Kong, Lu & Wong, 1994; Basu & Chib, 2003]

# Approximating the evidence (bis)

**Reverse logistic regression** applies to DPM:

Importance function

$$\pi_1(\mathbf{z}, \mathbf{M}) := \pi^*(\mathbf{z}|\mathbf{x}, \mathbf{M})\pi(\mathbf{M}) \quad \text{and} \quad \pi_2(\mathbf{z}, \mathbf{M}) = \frac{\pi(\mathbf{z}, \mathbf{M}|\mathbf{x})}{m(\mathbf{y})}$$

$\{\mathbf{z}^{(1,j)}, \mathbf{M}^{(1,j)}\}_{j=1}^T$  and  $\{\mathbf{z}^{(2,j)}, \mathbf{M}^{(2,j)}\}_{j=1}^T$  samples from  $\pi_1$  and  $\pi_2$

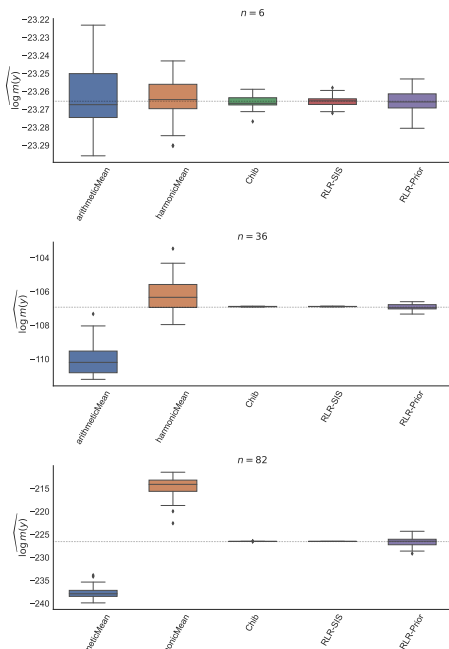
Marginal likelihood  $m(\mathbf{y})$  estimated as intercept of logistic regression with covariate

$$\log\{\pi_1(\mathbf{z}, \mathbf{M})/\tilde{\pi}_2(\mathbf{z}, \mathbf{M})\}$$

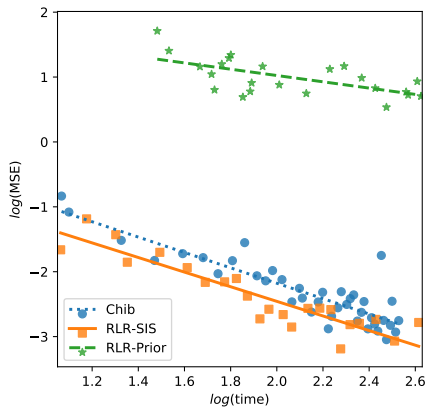
on merged sample

[Geyer, 1994; Chen & Shao, 1997]

# Galactic illustration

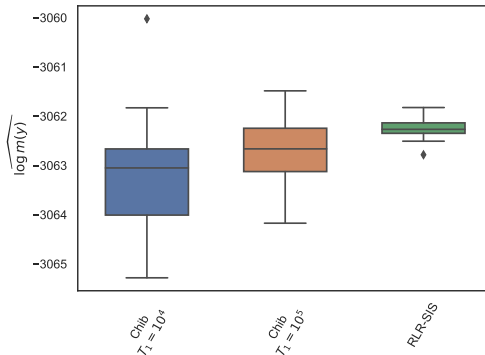


# Galactic illustration





# large data illustration



20 repetitions, synthetic data,  $n = 1000$  and  $K_0 = 6$ . Chib (left) :  $T_1 = 10^4$ , burnIn =  $10^3$ ,  $T_2 = 600$ , budget : 04:09:51. Chib (right)  $T_1 = 10^5$ , burnIn =  $10^4$ ,  $T_2 = 2000$ , budget : 34:01:18. RLR-SIS :  $T_1 = 10^4$ , burnIn =  $10^3$ ,  $T_2 = 600$ , budget : 06:23:12

- 1 Mixtures of distributions
- 2 Approximations to evidence
- 3 Dirichlet process mixtures
- 4 Distributed evidence evaluation



## Distributed Computation for Marginal Likelihood based Model Choice\*

Alexander Buchholz<sup>1,†</sup>, Daniel Abfack<sup>1,†</sup>, and Sylvia Richardson<sup>4</sup>

**Abstract.** We propose a general method for distributed Bayesian model choice, using the marginal likelihood, where a data set is split in non-overlapping subsets. These subsets are only accessed locally by individual workers and no data is shared between the workers. We approximate the model evidence for the full data set through Monte Carlo sampling from the posterior on every subset generating a model evidence per subset. The results are combined using a novel approach which corrects for the splitting using summary statistics of the generated samples. Our divide-and-conquer approach enables Bayesian model choice in the large data setting, exploiting all available information but limiting communication between workers. We derive theoretical error bounds that quantify the resulting trade-off between computational gain and loss in precision. The embarrassingly parallel nature yields important speed-ups when used on massive data sets as illustrated by our real world experiments. In addition, we show how the suggested approach can be extended to model choice within a reversible jump setting that explores multiple feature combinations within one run.

[Buchholz et al., 2022]

# Divide & Conquer

1. data  $\mathbf{y}$  divided into  $S$  batches  $\mathbf{y}_1, \dots, \mathbf{y}_S$  with

$$\begin{aligned}\pi(\vartheta|\mathbf{y}) &\propto p(\mathbf{y}|\vartheta)\pi(\vartheta) = \prod_{s=1}^S p(\mathbf{y}_s|\vartheta)\pi(\vartheta)^{1/S} \\ &= \prod_{s=1}^S p(\mathbf{y}_s|\vartheta)\tilde{\pi}(\vartheta) \propto \prod_{s=1}^S \tilde{\pi}(\vartheta|\mathbf{y}_s)\end{aligned}$$

2. infer with  $\tilde{\pi}(\vartheta|\mathbf{y}_s)$ , sub-posterior distributions, in parallel by MCMC
3. recombine all sub-posterior samples

[Buchholz et al., 2022]

# Connecting bits

While

$$m(\mathbf{y}) = \int \prod_{s=1}^S p(\mathbf{y}_s | \vartheta) \tilde{\pi}(\vartheta) d\vartheta \neq \prod_{s=1}^S \int p(\mathbf{y}_s | \vartheta) \tilde{\pi}(\vartheta) d\vartheta = \prod_{s=1}^S \tilde{m}(\mathbf{y}_s)$$

they can be connected as

$$m(\mathbf{y}) = \mathfrak{Z}^S \prod_{s=1}^S \tilde{m}(\mathbf{y}_s) \int \prod_{s=1}^S \tilde{\pi}(\vartheta | \mathbf{y}_s) d\vartheta$$

[Buchholz et al., 2022]

$$m(\mathbf{y}) = \mathfrak{Z}^S \prod_{s=1}^S \tilde{m}(\mathbf{y}_s) \int \prod_{s=1}^S \tilde{\pi}(\vartheta|\mathbf{y}_s) d\vartheta$$

where

$$\tilde{\pi}(\vartheta|\mathbf{y}_s) \propto p(\mathbf{y}_s|\vartheta)\tilde{\pi}(\vartheta),$$

$$\tilde{m}(\mathbf{y}_s) = \int p(\mathbf{y}_s|\vartheta)\tilde{\pi}(\vartheta) d\vartheta,$$

$$\mathfrak{Z} = \int \pi(\vartheta)^{1/S} d\vartheta$$

[Buchholz et al., 2022]

# Label unswitching worries

While  $\mathfrak{J}$  usually closed-form,

$$\mathfrak{J} = \int \prod_{s=1}^S \tilde{\pi}(\vartheta | \mathbf{y}_s) d\vartheta$$

is not and need be evaluated as

$$\hat{\mathfrak{J}} = \frac{1}{T} \sum_{t=1}^T \int \prod_{s=1}^S \tilde{\pi}(\vartheta | \mathbf{z}_s^{(t)}, \mathbf{y}_s) d\vartheta$$

when

$$\tilde{\pi}(\vartheta | \mathbf{y}_s) = \int \tilde{\pi}(\vartheta | \mathbf{z}_s, \mathbf{y}_s) \tilde{\pi}(\mathbf{z}_s | \mathbf{y}_s) d\mathbf{z}_s$$

# Label unswitching worries

$$\tilde{\pi}(\vartheta|\mathbf{y}_s) = \int \tilde{\pi}(\vartheta|\mathbf{z}_s, \mathbf{y}_s) \tilde{\pi}(\mathbf{z}_s|\mathbf{y}_s) d\mathbf{z}_s$$

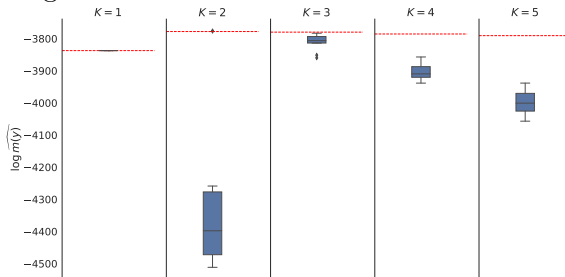
**Issue:** with distributed computing, shards  $\mathbf{z}_s$  are unrelated and corresponding clusters disconnected.



# Label unswitching worries

$$\tilde{\pi}(\vartheta|\mathbf{y}_s) = \int \tilde{\pi}(\vartheta|\mathbf{z}_s, \mathbf{y}_s) \tilde{\pi}(\mathbf{z}_s|\mathbf{y}_s) d\mathbf{z}_s$$

**Issue:** with distributed computing, shards  $\mathbf{z}_s$  are unrelated and corresponding clusters disconnected.



$10^3$  Gaussian observations, 10 repetitions, same number of Gibbs steps for all methods

# Label switching imposition

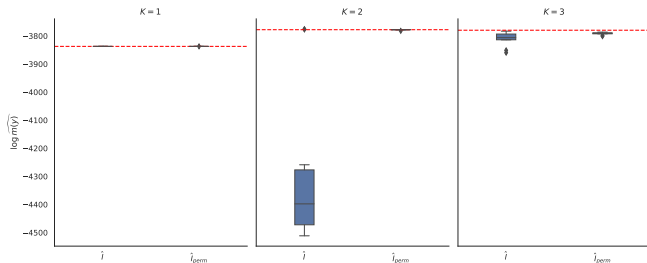
Returning to averaging across permutations, identity

$$\hat{\mathcal{J}}_{\text{perm}} = \frac{1}{TK!^{S-1}} \sum_{t=1}^T \sum_{\sigma_2, \dots, \sigma_S \in \mathfrak{S}_K} \int \tilde{\pi}(\vartheta | \mathbf{z}_1^{(t)}, \mathbf{y}_1) \prod_{s=2}^S \tilde{\pi}(\vartheta | \sigma_s(\mathbf{z}_s^{(t)}), \mathbf{y}_s) d\vartheta$$

# Label switching imposition

Returning to averaging across permutations, identity

$$\hat{\mathbf{j}}_{\text{perm}} = \frac{1}{TK!^{S-1}} \sum_{t=1}^T \sum_{\sigma_2, \dots, \sigma_S \in \mathfrak{S}_K} \int \tilde{\pi}(\vartheta | \mathbf{z}_1^{(t)}, \mathbf{y}_1) \prod_{s=2}^S \tilde{\pi}(\vartheta | \sigma_s(\mathbf{z}_s^{(t)}), \mathbf{y}_s) d\vartheta$$



# Label switching imposition

Returning to averaging across permutations, identity

$$\hat{\mathcal{J}}_{\text{perm}} = \frac{1}{TK!^{S-1}} \sum_{t=1}^T \sum_{\sigma_2, \dots, \sigma_S \in \mathfrak{S}_K} \int \tilde{\pi}(\vartheta | \mathbf{z}_1^{(t)}, \mathbf{y}_1) \prod_{s=2}^S \tilde{\pi}(\vartheta | \sigma_s(\mathbf{z}_s^{(t)}), \mathbf{y}_s) d\vartheta$$

Obtained at heavy computational cost:  $\mathcal{O}(T)$  for  $\hat{\mathcal{J}}$  versus  $\mathcal{O}(TK!^{S-1})$  for  $\hat{\mathcal{J}}_{\text{perm}}$

# Importance sampling version

Obtained at heavy computational cost:  $\mathcal{O}(T)$  for  $\hat{\mathcal{J}}$  versus  $\mathcal{O}(TK!^{S-1})$  for  $\hat{\mathcal{J}}_{\text{perm}}$

Avoid enumeration of permutations by using simulated values of parameter for the reference sub-posterior as anchors towards coherent labeling of clusters

[Celeux, 1998; Stephens, 2000]

# Importance sampling version

For each batch  $s = 2, \dots, S$ , define *matching* matrix

$$P_s = \begin{pmatrix} p_{s11} & \cdots & p_{s1K} \\ \vdots & \vdots & \vdots \\ p_{sK1} & \cdots & p_{sKK} \end{pmatrix}$$

where

$$p_{s1k} = \prod_{i: z_{si}=1} p(y_{si} | \vartheta_k)$$

used in creating proposals

$$q_s(\sigma) \propto \prod_{k=1}^K p_{sk\sigma(k)}$$

that reflect probabilities that each cluster  $k$  of batch  $s$  is well-matched with cluster  $\sigma(k)$  of batch 1

# Importance sampling version

Considerably reduced computational cost compared to

$\hat{m}_{\hat{I}_{\text{perm}}}(\mathbf{y})$

- ▶ At each iteration  $t$ , total cost of  $O(Kn/S)$  for evaluating  $P_s$
- ▶ computing  $K!$  weights of discrete importance distribution  $q_{\sigma_s}$  requires  $K!$  operations
- ▶ sampling from the global discrete importance distribution requires  $M^{(t)}$  basic operations

Global cost of

$$O(T(Kn/S + K! + \bar{M}))$$

for  $\bar{M}$  maximum number of importance simulations

# Importance sampling version

Resulting estimator

$$\hat{\mathfrak{J}}_{\text{IS}} = \frac{1}{TK!^{S-1}} \sum_{t=1}^T \frac{1}{M^{(t)}} \sum_{m=1}^{M^{(t)}} \frac{\chi(\mathbf{z}^{(t)}; \sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}{\pi_{\sigma}(\sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}$$

where

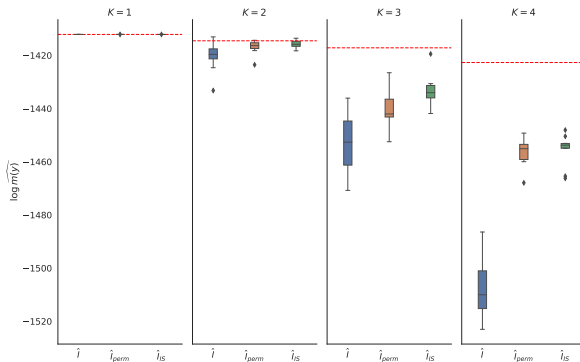
$$\chi(\mathbf{z}^{(t)}; \sigma_2, \dots, \sigma_S) := \int \tilde{\pi}(\vartheta | \mathbf{z}_1^{(t)}, \mathbf{y}_1) \prod_{s=2}^S \tilde{\pi}(\vartheta | \sigma_s(\mathbf{z}_s^{(t)}), \mathbf{y}_s) d\vartheta$$



# Importance sampling version

Resulting estimator

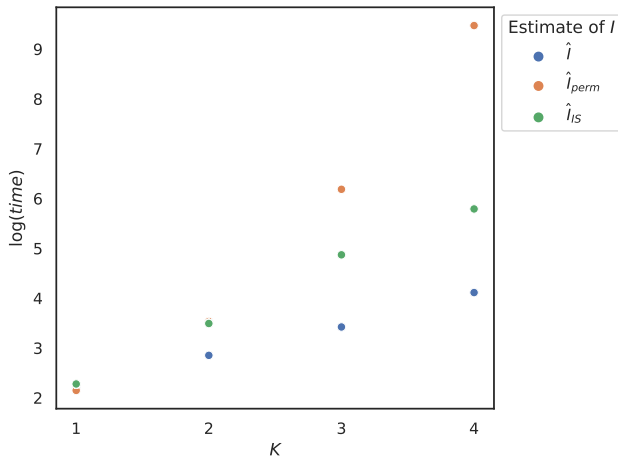
$$\hat{\mathfrak{J}}_{\text{IS}} = \frac{1}{TK!S-1} \sum_{t=1}^T \frac{1}{M^{(t)}} \sum_{m=1}^{M^{(t)}} \frac{\chi(z^{(t)}; \sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}{\pi_{\sigma}(\sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}$$



# Importance sampling version

Resulting estimator

$$\hat{\mathcal{J}}_{IS} = \frac{1}{TK!^{S-1}} \sum_{t=1}^T \frac{1}{M^{(t)}} \sum_{m=1}^{M^{(t)}} \frac{\chi(\mathbf{z}^{(t)}; \sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}{\pi_{\sigma}(\sigma_2^{(t,m)}, \dots, \sigma_S^{(t,m)})}$$



# Sequential importance sampling

Define

$$\tilde{\pi}_s(\vartheta) = \frac{\prod_{l=1}^s \tilde{\pi}(\vartheta | \mathbf{y}_l)}{z_s}$$

where

$$z_s = \int \prod_{l=1}^s \tilde{\pi}(\vartheta | \mathbf{y}_l)$$

then

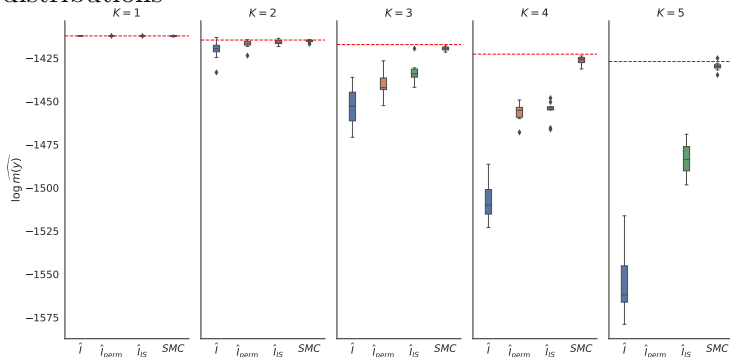
$$m(\mathbf{y}) = z^S \times m(\mathbf{y}_1) \times \prod_{s=2}^S \int \pi_{s-1}(\vartheta) p(\mathbf{y}_s | \vartheta) \tilde{\pi}(\vartheta) d\vartheta$$

# Sequential importance sampling

Calls for standard sequential importance sampling strategy making use of the successive distributions  $\pi_s(\vartheta)$  as importance distributions

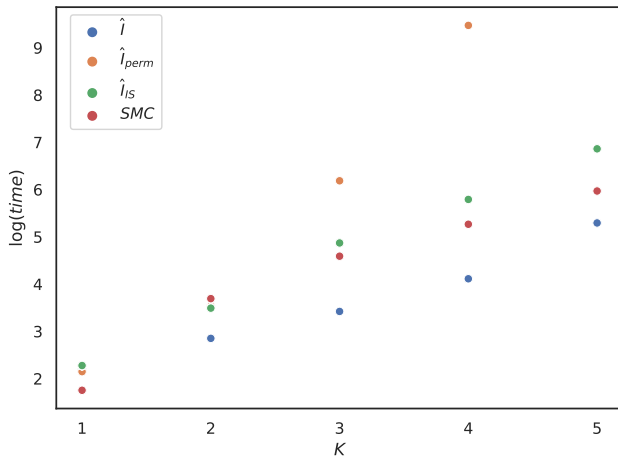
# Sequential importance sampling

Calls for standard sequential importance sampling strategy making use of the successive distributions  $\pi_s(\vartheta)$  as importance distributions



# Sequential importance sampling

Calls for standard sequential importance sampling strategy making use of the successive distributions  $\pi_s(\vartheta)$  as importance distributions



# Conclusion

- ▶ Buchholz et al. 2022 not applicable to finite mixture models
- ▶ adapted version with reasonable computational time.
- ▶ new identity bridging gap between full and batch marginal likelihoods straightforward to implement by SMC
- ▶ valid for all kind of parametric models while relaxes conjugacy