

Consistency of H-factors for model selection

Stephane Shao^{*}

September 12, 2018

in collaboration with Pierre E. Jacob^{}, Jie Ding[†] and Vahid Tarokh[‡]*

^{*}Department of Statistics, Harvard University

[†]School of Statistics, University of Minnesota

[‡]Department of Electrical and Computer Engineering, Duke University

Table of contents

1. An alternative to Bayes factors when priors are vague: H-factors
2. Consistent model selection with H-factors
3. Numerical illustrations

Table of contents

1. An alternative to Bayes factors when priors are vague: H-factors
2. Consistent model selection with H-factors
3. Numerical illustrations

Table of contents

1. An alternative to Bayes factors when priors are vague: H-factors
2. Consistent model selection with H-factors
3. Numerical illustrations

Why use H-factors ?

Goal: compare and select Bayesian models

- We want to select a model from a set $\{M_1, \dots, M_q\}$, given observations $Y_{1:T} = (Y_1, \dots, Y_T) \in (\mathbb{R}^{d_y})^T$ from a data generating process p_\star .
- Each model M_j is a collection of distributions $p_{\theta_j}(dy_{1:T})$ parametrized by $\theta_j \in \mathbb{T}_j \subseteq \mathbb{R}^{d_j}$ with a prior distribution $p(d\theta_j)$.
- We are interested in settings where the priors may be vague, and the models may be misspecified ($p_\star \notin M_j$).

Goal: compare and select Bayesian models

- We want to select a model from a set $\{M_1, \dots, M_q\}$, given observations $Y_{1:T} = (Y_1, \dots, Y_T) \in (\mathbb{R}^{d_y})^T$ from a data generating process p_\star .
- Each model M_j is a collection of distributions $p_{\theta_j}(dy_{1:T})$ parametrized by $\theta_j \in \mathbb{T}_j \subseteq \mathbb{R}^{d_j}$ with a prior distribution $p(d\theta_j)$.
- We are interested in settings where the priors may be vague, and the models may be misspecified ($p_\star \notin M_j$).

Goal: compare and select Bayesian models

- We want to select a model from a set $\{M_1, \dots, M_q\}$, given observations $Y_{1:T} = (Y_1, \dots, Y_T) \in (\mathbb{R}^{d_y})^T$ from a data generating process p_\star .
- Each model M_j is a collection of distributions $p_{\theta_j}(dy_{1:T})$ parametrized by $\theta_j \in \mathbb{T}_j \subseteq \mathbb{R}^{d_j}$ with a prior distribution $p(d\theta_j)$.
- We are interested in settings where the priors may be vague, and the models may be misspecified ($p_\star \notin M_j$).

What kind of models ? i.i.d., state-space models, ...

- **Example 1. Normal i.i.d. models** [O'Hagan, 1995]

Model 1

$$Y_{1:T} | \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1)$$

$$\theta_1 \sim \mathcal{N}(0, \sigma_0^2)$$

Model 2

$$Y_{1:T} | \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2)$$

$$\theta_2 \sim \text{Inv-}\chi^2(\nu_0, s_0^2)$$

What kind of models ? i.i.d., state-space models, ...

- **Example 2.** Lévy-driven stochastic volatility models for log-returns of financial assets [Barndorff-Nielsen & Shephard, 2001]

Given parameters (λ, ξ, ω) , generate random variables $(V_t, Z_t)_{t \geq 1}$ recursively as

$$\left. \begin{aligned} k &\sim \text{Poisson}(\lambda \xi^2 / \omega^2); & C_{1:k} &\stackrel{\text{i.i.d.}}{\sim} \text{Unif}(t-1, t); & E_{1:k} &\stackrel{\text{i.i.d.}}{\sim} \text{Exp}(\xi / \omega^2); & Z_0 &\sim \Gamma(\xi^2 / \omega^2, \xi / \omega^2) \\ Z_t &= e^{-\lambda} Z_{t-1} + \sum_{j=1}^k e^{-\lambda(t-C_j)} E_j; & V_t &= \lambda^{-1} (Z_{t-1} - Z_t + \sum_{j=1}^k E_j) \end{aligned} \right\} \quad (\text{gears})$$

Model 1

Given parameters $(\lambda, \xi, \omega, \mu, \beta)$:

$$(V_t, Z_t) \sim (\text{gears})$$

$$X_t = (V_t, Z_t)$$

$$Y_t | X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$$

with independent priors:

$$\lambda \sim \text{Exp}(1); \quad \xi, \omega^2 \sim \text{Exp}(1/5); \quad \mu, \beta \sim \mathcal{N}(0, 10)$$

Model 2

Given parameters $(\lambda_1, \lambda_2, w_1, w_2, \xi, \omega, \mu, \beta)$:

$$(V_{1,t}, Z_{1,t}) \sim (\text{gears}) \text{ with } (\lambda_1, \xi w_1, \omega w_1)$$

$$(V_{2,t}, Z_{2,t}) \sim (\text{gears}) \text{ with } (\lambda_2, \xi w_2, \omega w_2)$$

$$X_t = (V_{1,t}, V_{2,t}, Z_{1,t}, Z_{2,t})$$

$$V_t = V_{1,t} + V_{2,t}$$

$$Y_t | X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$$

with independent priors:

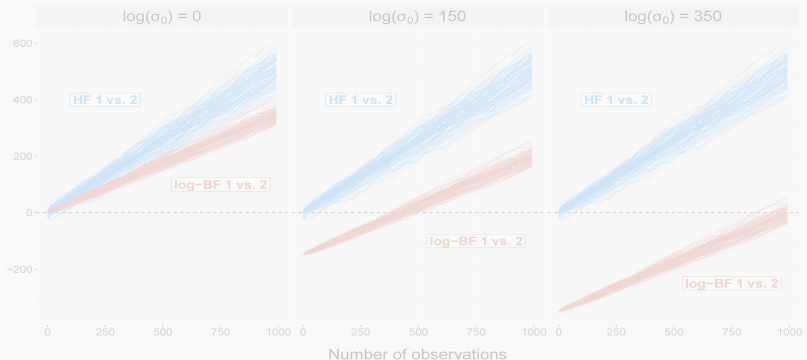
$$\begin{aligned} \lambda_1 &\sim \text{Exp}(1); & \lambda_2 - \lambda_1 &\sim \text{Exp}(1/2); & \xi, \omega^2 &\sim \text{Exp}(1/5) \\ 1 - w_2 = w_1 &\sim \text{Unif}(0, 1); & \mu, \beta &\sim \mathcal{N}(0, 10) \end{aligned}$$

Limitations of Bayes factors: sensitivity to vague priors

- Making the prior more vague effectively multiplies the evidence $p_j(y_{1:T})$ of a model M_j by an arbitrarily small constant, for any fixed sample size.
- Example 1. Observations generated as $Y_{1:1000} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(1, 1)$.

$$M_1 = \{Y_{1:T} \mid \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1) ; \theta_1 \sim \mathcal{N}(0, \sigma_0^2)\}$$

$$M_2 = \{Y_{1:T} \mid \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2) ; \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)\}$$

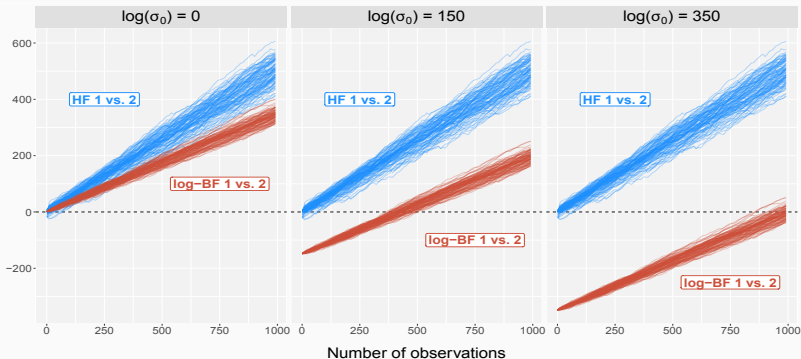


Limitations of Bayes factors: sensitivity to vague priors

- Making the prior more vague effectively multiplies the evidence $p_j(y_{1:T})$ of a model M_j by an arbitrarily small constant, for any fixed sample size.
- **Example 1.** Observations generated as $Y_{1:1000} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(1, 1)$.

$$M_1 = \{Y_{1:T} \mid \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1) ; \theta_1 \sim \mathcal{N}(0, \sigma_0^2)\}$$

$$M_2 = \{Y_{1:T} \mid \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2) ; \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)\}$$



Decision theoretic solution: change of scoring rule

- Bayes factors select models maximizing $\log p(y_{1:T}) = \sum_{t=1}^T \log p(y_t | y_{1:t-1})$.
- This corresponds to minimizing the prequential score [Dawid, 1984]

$$\sum_{t=1}^T \mathcal{S}(y_t, p(dy_t | y_{1:t-1}))$$

with the choice of scoring rule $\mathcal{S}(y, p) = -\log p(y)$ called the log-score.

- Each scoring rule has an associated divergence function

$$D_{\mathcal{S}}(p, q) = \mathbb{E}_{Y \sim p} [\mathcal{S}(Y, q) - \mathcal{S}(Y, p)]$$

\mathcal{S} is (strictly) proper if $q \mapsto D_{\mathcal{S}}(p, q)$ is (uniquely) minimized at $q = p$.

- The log-score is strictly proper and tied to the Kullback-Leibler divergence

$$KL(p, q) = \int [\log p(y) - \log q(y)] p(y) dy$$

Decision theoretic solution: change of scoring rule

- Bayes factors select models maximizing $\log p(y_{1:T}) = \sum_{t=1}^T \log p(y_t | y_{1:t-1})$.
- This corresponds to minimizing the **prequential score** [Dawid, 1984]

$$\sum_{t=1}^T \mathcal{S}(y_t, p(dy_t | y_{1:t-1}))$$

with the choice of **scoring rule** $\mathcal{S}(y, p) = -\log p(y)$ called the **log-score**.

- Each scoring rule has an associated divergence function

$$D_{\mathcal{S}}(p, q) = \mathbb{E}_{Y \sim p} [\mathcal{S}(Y, q) - \mathcal{S}(Y, p)]$$

\mathcal{S} is (strictly) proper if $q \mapsto D_{\mathcal{S}}(p, q)$ is (uniquely) minimized at $q = p$.

- The log-score is strictly proper and tied to the Kullback-Leibler divergence

$$\text{KL}(p, q) = \int [\log p(y) - \log q(y)] p(y) dy$$

Decision theoretic solution: change of scoring rule

- Bayes factors select models maximizing $\log p(y_{1:T}) = \sum_{t=1}^T \log p(y_t | y_{1:t-1})$.
- This corresponds to minimizing the prequential score [Dawid, 1984]

$$\sum_{t=1}^T \mathcal{S}(y_t, p(dy_t | y_{1:t-1}))$$

with the choice of scoring rule $\mathcal{S}(y, p) = -\log p(y)$ called the log-score.

- Each scoring rule has an associated divergence function

$$D_{\mathcal{S}}(p, q) = \mathbb{E}_{Y \sim p} [\mathcal{S}(Y, q) - \mathcal{S}(Y, p)]$$

\mathcal{S} is (strictly) proper if $q \mapsto D_{\mathcal{S}}(p, q)$ is (uniquely) minimized at $q = p$.

- The log-score is strictly proper and tied to the Kullback-Leibler divergence

$$KL(p, q) = \int [\log p(y) - \log q(y)] p(y) dy$$

Decision theoretic solution: change of scoring rule

- Bayes factors select models maximizing $\log p(y_{1:T}) = \sum_{t=1}^T \log p(y_t | y_{1:t-1})$.
- This corresponds to minimizing the prequential score [Dawid, 1984]

$$\sum_{t=1}^T \mathcal{S}(y_t, p(dy_t | y_{1:t-1}))$$

with the choice of scoring rule $\mathcal{S}(y, p) = -\log p(y)$ called the log-score.

- Each scoring rule has an associated divergence function

$$D_{\mathcal{S}}(p, q) = \mathbb{E}_{Y \sim p} [\mathcal{S}(Y, q) - \mathcal{S}(Y, p)]$$

\mathcal{S} is (strictly) proper if $q \mapsto D_{\mathcal{S}}(p, q)$ is (uniquely) minimized at $q = p$.

- The log-score is strictly proper and tied to the Kullback-Leibler divergence

$$KL(p, q) = \int [\log p(y) - \log q(y)] p(y) dy$$

Change of scoring rule: the Hyvärinen score

- Instead of the KL-divergence, Dawid & Musio [2015] propose to use

$$D_{\mathcal{H}}(p, q) = \int \left\| \nabla \log p(y) - \nabla \log q(y) \right\|^2 p(y) dy$$

sometimes called the **relative Fisher information divergence**.

- It induces a scoring rule known as the Hyvärinen score [Hyvärinen, 2005]

$$\mathcal{H}(y, p) = 2 \Delta \log p(y) + \|\nabla \log p(y)\|^2$$

where the Laplacian Δ and the gradient ∇ are with respect to y . This score is strictly proper, local, homogeneous [Parry, Dawid & Lauritzen, 2012].

Select M_j minimizing the prequential Hyvärinen score (H-score)

$$\mathcal{H}_T(M_j) = \sum_{t=1}^T \mathcal{H} \left(y_t, p_j(dy_t | y_{1:t-1}) \right)$$

which can be consistently estimated using SMC [Chopin, 2002; Del Moral, Doucet & Jasra, 2006] or SMC² [Chopin, Jacob & Papaspiliopoulos, 2013].

Change of scoring rule: the Hyvärinen score

- Instead of the KL-divergence, Dawid & Musio [2015] propose to use

$$D_{\mathcal{H}}(p, q) = \int \left\| \nabla \log p(y) - \nabla \log q(y) \right\|^2 p(y) dy$$

sometimes called the relative Fisher information divergence.

- It induces a scoring rule known as the **Hyvärinen score** [Hyvärinen, 2005]

$$\mathcal{H}(y, p) = 2 \Delta \log p(y) + \|\nabla \log p(y)\|^2$$

where the Laplacian Δ and the gradient ∇ are with respect to y . This score is **strictly proper, local, homogeneous** [Parry, Dawid & Lauritzen, 2012].

Select M_j minimizing the prequential Hyvärinen score (H-score)

$$\mathcal{H}_T(M_j) = \sum_{t=1}^T \mathcal{H} \left(y_t, p_j(dy_t | y_{1:t-1}) \right)$$

which can be consistently estimated using SMC [Chopin, 2002; Del Moral, Doucet & Jasra, 2006] or SMC² [Chopin, Jacob & Papaspiliopoulos, 2013].

Change of scoring rule: the Hyvärinen score

- Instead of the KL-divergence, Dawid & Musio [2015] propose to use

$$D_{\mathcal{H}}(p, q) = \int \left\| \nabla \log p(y) - \nabla \log q(y) \right\|^2 p(y) dy$$

sometimes called the relative Fisher information divergence.

- It induces a scoring rule known as the **Hyvärinen score** [Hyvärinen, 2005]

$$\mathcal{H}(y, p) = 2 \Delta \log p(y) + \|\nabla \log p(y)\|^2$$

where the Laplacian Δ and the gradient ∇ are with respect to y . This score is strictly proper, local, homogeneous [Parry, Dawid & Lauritzen, 2012].

Select M_j minimizing the prequential Hyvärinen score (**H-score**)

$$\mathcal{H}_T(M_j) = \sum_{t=1}^T \mathcal{H}(y_t, p_j(dy_t | y_{1:t-1}))$$

which can be consistently estimated using **SMC** [Chopin, 2002; Del Moral, Doucet & Jasra, 2006] or **SMC²** [Chopin, Jacob & Papaspiliopoulos, 2013].

Consistency of the H-score

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

- Consider a generic model $M : Y_{1:T} | \theta \stackrel{\text{i.i.d.}}{\sim} p_\theta ; \theta \sim p(d\theta)$.
- By differentiating under the integral sign, the H-score $\mathcal{H}_T(M)$ equals

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_\Theta) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_\Theta(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

where the expectations and variances are with respect to $\Theta \sim p(d\theta|y_{1:t})$.

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

- Consider a generic model $M : Y_{1:T} | \theta \stackrel{\text{i.i.d.}}{\sim} p_\theta$; $\theta \sim p(d\theta)$.
- By differentiating under the integral sign, the H-score $\mathcal{H}_T(M)$ equals

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_\Theta) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_\Theta(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

where the expectations and variances are with respect to $\Theta \sim p(d\theta|y_{1:t})$.

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

Concentration of the posterior

+ uniform integrability

+ equicontinuity

$$\mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \xrightarrow[t \rightarrow +\infty]{} \mathcal{H}(Y_t, p_{\theta^*}) \quad \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \xrightarrow[t \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} 0$$

Césaro's theorem

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \xrightarrow[T \rightarrow +\infty]{} \frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \quad \frac{1}{T} \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} 0$$

Law of large numbers

$$\frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} \mathbb{E}_* \left[\mathcal{H}(Y, p_{\theta^*}) \right]$$

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

Concentration of the posterior

+ uniform integrability

+ equicontinuity

$$\mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \underset{t \rightarrow +\infty}{\approx} \mathcal{H}(Y_t, p_{\theta^*}) \qquad \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \underset{t \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} 0$$

Césaro's theorem

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \underset{T \rightarrow +\infty}{\approx} \frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \qquad \frac{1}{T} \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \underset{T \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} 0$$

Law of large numbers

$$\frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \underset{T \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} \mathbb{E}_* \left[\mathcal{H}(Y, p_{\theta^*}) \right]$$

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

Concentration of the posterior

+ uniform integrability

+ equicontinuity

$$\mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \underset{t \rightarrow +\infty}{\approx} \mathcal{H}(Y_t, p_{\theta^*}) \qquad \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \underset{t \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} 0$$

Césaro's theorem

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \underset{T \rightarrow +\infty}{\approx} \frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \qquad \frac{1}{T} \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \underset{T \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} 0$$

Law of large numbers

$$\frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \underset{T \rightarrow +\infty}{\xrightarrow{\mathbb{P}_* \text{ a.s.}}} \mathbb{E}_* \left[\mathcal{H}(Y, p_{\theta^*}) \right]$$

Elements of proof in the univariate i.i.d. case ($d_y = 1$)

$$\sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] + \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right)$$

Concentration of the posterior

+ uniform integrability

+ equicontinuity

$$\mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \xrightarrow[t \rightarrow +\infty]{} \mathcal{H}(Y_t, p_{\theta^*}) \quad \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \xrightarrow[t \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} 0$$

Césaro's theorem

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[\mathcal{H}(Y_t, p_{\Theta}) \middle| Y_{1:t} \right] \xrightarrow[T \rightarrow +\infty]{} \frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \quad \frac{1}{T} \sum_{t=1}^T \text{Var} \left(\frac{\partial \log p_{\Theta}(Y_t)}{\partial y} \middle| Y_{1:t} \right) \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} 0$$

Law of large numbers

$$\frac{1}{T} \sum_{t=1}^T \mathcal{H}(Y_t, p_{\theta^*}) \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_* \text{ a.s.}} \mathbb{E}_* \left[\mathcal{H}(Y, p_{\theta^*}) \right]$$

Consistency result in the univariate i.i.d. case

- Under regularity conditions, the H-factor of M_1 vs. M_2 satisfies


$$\frac{1}{T} \left[\mathcal{H}_T(M_2) - \mathcal{H}_T(M_1) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} D_{\mathcal{H}}(p_\star, M_2) - D_{\mathcal{H}}(p_\star, M_1)$$

where $D_{\mathcal{H}}(p_\star, M_j) := \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$.

- In contrast, the log-Bayes factor of M_1 vs. M_2 satisfies

$$\frac{1}{T} \left[\left(-\log p_2(Y_{1:T}) \right) - \left(-\log p_1(Y_{1:T}) \right) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} \text{KL}(p_\star, M_2) - \text{KL}(p_\star, M_1)$$

where $\text{KL}(p_\star, M_j) := \mathbb{E}_\star[-\log p_{\theta_j^\star}(Y)] - \mathbb{E}_\star[-\log p_\star(Y)]$.

- This extends to state-space models and dependent data, with additional technicalities (e.g. forgetting properties and ergodic theorems).
- The limit is meaningless if p_\star belongs to both models (e.g. nested well-specified setting): we need higher order Bayesian asymptotics, i.e. Bernstein-von-Mises-type results ... which are non-trivial for state-space models 

Consistency result in the univariate i.i.d. case

- Under regularity conditions, the H-factor of M_1 vs. M_2 satisfies


$$\frac{1}{T} \left[\mathcal{H}_T(M_2) - \mathcal{H}_T(M_1) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} D_{\mathcal{H}}(p_\star, M_2) - D_{\mathcal{H}}(p_\star, M_1)$$

where $D_{\mathcal{H}}(p_\star, M_j) := \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$.

- In contrast, the log-Bayes factor of M_1 vs. M_2 satisfies

$$\frac{1}{T} \left[\left(-\log p_2(Y_{1:T}) \right) - \left(-\log p_1(Y_{1:T}) \right) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} \text{KL}(p_\star, M_2) - \text{KL}(p_\star, M_1)$$

where $\text{KL}(p_\star, M_j) := \mathbb{E}_\star[-\log p_{\theta_j^\star}(Y)] - \mathbb{E}_\star[-\log p_\star(Y)]$.

- This extends to state-space models and dependent data, with additional technicalities (e.g. forgetting properties and ergodic theorems).
- The limit is meaningless if p_\star belongs to both models (e.g. nested well-specified setting): we need higher order Bayesian asymptotics, i.e. Bernstein-von-Mises-type results ... which are non-trivial for state-space models 

Consistency result in the univariate i.i.d. case

- Under regularity conditions, the H-factor of M_1 vs. M_2 satisfies


$$\frac{1}{T} \left[\mathcal{H}_T(M_2) - \mathcal{H}_T(M_1) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} D_{\mathcal{H}}(p_\star, M_2) - D_{\mathcal{H}}(p_\star, M_1)$$

where $D_{\mathcal{H}}(p_\star, M_j) := \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$.

- In contrast, the log-Bayes factor of M_1 vs. M_2 satisfies

$$\frac{1}{T} \left[\left(-\log p_2(Y_{1:T}) \right) - \left(-\log p_1(Y_{1:T}) \right) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} \text{KL}(p_\star, M_2) - \text{KL}(p_\star, M_1)$$

where $\text{KL}(p_\star, M_j) := \mathbb{E}_\star[-\log p_{\theta_j^\star}(Y)] - \mathbb{E}_\star[-\log p_\star(Y)]$.

- This extends to state-space models and dependent data, with additional technicalities (e.g. forgetting properties and ergodic theorems).
- The limit is meaningless if p_\star belongs to both models (e.g. nested well-specified setting): we need higher order Bayesian asymptotics, i.e. Bernstein-von-Mises-type results ... which are non-trivial for state-space models 

Consistency result in the univariate i.i.d. case

- Under regularity conditions, the H-factor of M_1 vs. M_2 satisfies

$$\frac{1}{T} \left[\mathcal{H}_T(M_2) - \mathcal{H}_T(M_1) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} D_{\mathcal{H}}(p_\star, M_2) - D_{\mathcal{H}}(p_\star, M_1)$$

where $D_{\mathcal{H}}(p_\star, M_j) := \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$.

- In contrast, the log-Bayes factor of M_1 vs. M_2 satisfies

$$\frac{1}{T} \left[\left(-\log p_2(Y_{1:T}) \right) - \left(-\log p_1(Y_{1:T}) \right) \right] \xrightarrow[T \rightarrow +\infty]{\mathbb{P}_\star - a.s.} \text{KL}(p_\star, M_2) - \text{KL}(p_\star, M_1)$$

where $\text{KL}(p_\star, M_j) := \mathbb{E}_\star[-\log p_{\theta_j^\star}(Y)] - \mathbb{E}_\star[-\log p_\star(Y)]$.


- This extends to state-space models and dependent data, with additional technicalities (e.g. forgetting properties and ergodic theorems).
- The limit is meaningless if p_\star belongs to both models (e.g. nested well-specified setting): we need higher order Bayesian asymptotics, i.e. Bernstein-von-Mises-type results ... which are non-trivial for state-space models 

Illustration of consistency for i.i.d. observations

- **Example 1.** Given simulated $Y_1, \dots, Y_{1000} \sim \mathcal{N}(\mu_\star, \sigma_\star^2)$, we compare

$$M_1: Y_1, \dots, Y_T | \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1), \quad \theta_1 \sim \mathcal{N}(0, 10)$$

$$M_2: Y_1, \dots, Y_T | \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2), \quad \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)$$

in the following four cases $(\mu_\star, \sigma_\star^2) = (1, 1), (0, 5), (4, 3), (0, 1)$.

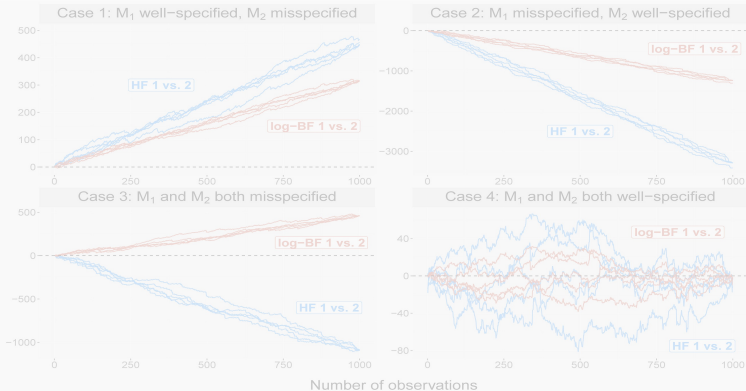


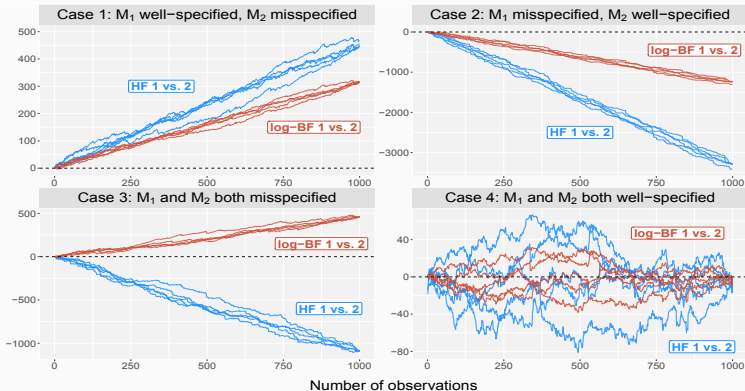
Illustration of consistency for i.i.d. observations

- **Example 1.** Given simulated $Y_1, \dots, Y_{1000} \sim \mathcal{N}(\mu_*, \sigma_*^2)$, we compare

$$M_1: Y_1, \dots, Y_T | \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1), \quad \theta_1 \sim \mathcal{N}(0, 10)$$

$$M_2: Y_1, \dots, Y_T | \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2), \quad \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)$$

in the following four cases $(\mu_*, \sigma_*^2) = (1, 1), (0, 5), (4, 3), (0, 1)$.



Summary

The H-score has the advantage of being ...

- Robust to vagueness of priors and allows for improper priors
- Justified non-asymptotically and also generally consistent
- Applicable to a wide range of parametric models via SMC methods






... albeit at the cost of more regularity on the candidate densities and more expensive computation in practice.

Avenues for future research

- Confidence intervals using unbiased MCMC [Jacob, O'Leary, Atchadé, 2018]
- Posterior consistency + asymptotic Normality for state-space models
- More details in Shao, Jacob, Ding & Tarokh (2018)
- R package available at: github.com/pierrejacob/bayeshscore

Questions ?

References (1/2)

-  O. E. Barndorff-Nielsen and N. Shephard.
Non-Gaussian Ornstein–Uhlenbeck-based models and some of their uses in financial economics.
Journal of the Royal Statistical Society: Series B, 63(2):167–241, 2001.
-  N. Chopin.
A sequential particle filter method for static models.
89:539–552, 2002.
-  N. Chopin, P. E. Jacob, and O. Papaspiliopoulos.
SMC²: an efficient algorithm for sequential analysis of state-space models.
Journal of the Royal Statistical Society, 75 (3):397–426, 2013.
-  A. P. Dawid and M. Musio.
Bayesian model selection based on proper scoring rules.
Bayesian Analysis, 10 (2):479–499, 2015.
-  P. Del Moral, A. Doucet, and A. Jasra.
Sequential Monte Carlo samplers.
Journal of the Royal Statistical Society: Series B (Statistical Methodology), 68(3):411–436, 2006.

References (2/2)



J. Knappe and P. D. Valpine.

Fitting complex population models by combining particle filters with markov chain monte carlo.

Ecology, 93 (2):256–263, 2012.



A. O'Hagan.

Fractional bayes factor for model comparison.

Journal of the Royal Statistical Society, 57 (1):99–138, 1995.



M. Parry, A. P. Dawid, and S. Lauritzen.

Proper local scoring rules.

The Annals of Statistics, 40 (1):561–592, 2012.



S. Shao, P. E. Jacob, J. Ding, and V. Tarokh.

Bayesian model comparison with the Hyvärinen score: computation and consistency.

Journal of the American Statistical Association. DOI: 10.1080/01621459.2018.1518237, 2018.

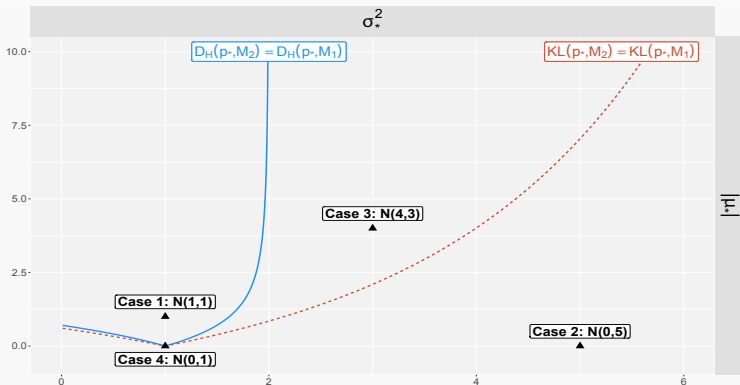
Illustration of consistency for i.i.d. observations (continued)

- **Example 1.** Given simulated $Y_1, \dots, Y_{1000} \sim \mathcal{N}(\mu_*, \sigma_*^2)$, we compare

$$M_1: Y_1, \dots, Y_T | \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\theta_1, 1), \quad \theta_1 \sim \mathcal{N}(0, 10)$$

$$M_2: Y_1, \dots, Y_T | \theta_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2), \quad \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)$$

in the following four cases $(\mu_*, \sigma_*^2) = (1, 1), (0, 5), (4, 3), (0, 1)$.



Nested models in the univariate i.i.d. case

- **Example 1.** Given simulated $Y_1, \dots, Y_{1000} \sim \mathcal{N}(\mu_*, \sigma_*^2)$, we compare

$$M_1: Y_1, \dots, Y_T | \theta_1 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_1, 1), \quad \mu_1 \sim \mathcal{N}(0, 10)$$

$$M_2: Y_1, \dots, Y_T | \mu_2, \sigma_2^2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mu_2, \sigma_2^2), \quad (\mu_2, \sigma_2^2) \sim \mathcal{N}\text{-Inv-}\chi^2(0, 1, 0.1, 1)$$

in the following two cases $(\mu_*, \sigma_*^2) = (0, 5), (0, 1)$.

