# Bayesian model comparison with the Hyvärinen score

Stephane Shao*

November 21, 2017

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

[*] Department of Statistics, Harvard University
[†] School of Engineering and Applied Sciences, Harvard University

# Table of contents

# Table of contents

# Table of contents

# Table of contents

# Table of contents

# Table of contents

# Setting and examples

- Given $T$ observations $Y_{1:T} = (Y_1, ..., Y_T) \in (\mathbb{R}^{d_y})^T$.

- We want to compare candidate models from a finite set $\{M_1, ..., M_q\}$.

- Each model $M_j$ is a collection of distributions $p_j(dy_{1:T}|\theta_j)$ parametrized by $\theta_j \in \mathbb{T}_j$ with a prior distribution $p_j(d\theta_j)$ on the parameter.

- We are interested in settings where the prior may be vague or improper.

- Given $T$ observations $Y_{1:T} = (Y_1, ..., Y_T) \in (\mathbb{R}^{d_y})^T$.

- **We want to compare candidate models from a finite set $\{M_1, ..., M_q\}$.**

- Each model $M_j$ is a collection of distributions $p_j(dy_{1:T}|\theta_j)$ parametrized by $\theta_j \in \mathbb{T}_j$ with a prior distribution $p_j(d\theta_j)$ on the parameter.

- We are interested in settings where the prior may be vague or improper.

- Given $T$ observations $Y_{1:T} = (Y_1, ..., Y_T) \in (\mathbb{R}^{d_y})^T$.

- We want to compare candidate models from a finite set $\{M_1, ..., M_q\}$.

- Each model $M_j$ is a collection of distributions $p_j(dy_{1:T}|\theta_j)$ parametrized by $\theta_j \in \mathbb{T}_j$ with a prior distribution $p_j(d\theta_j)$ on the parameter.

- We are interested in settings where the prior may be vague or improper.

# Goal: compare models in a Bayesian framework

- Given $T$ observations $Y_{1:T} = (Y_1, ..., Y_T) \in (\mathbb{R}^{d_y})^T$.

- We want to compare candidate models from a finite set $\{M_1, ..., M_q\}$.

- Each model $M_j$ is a collection of distributions $p_j(dy_{1:T}|\theta_j)$ parametrized by $\theta_j \in \mathbb{T}_j$ with a prior distribution $p_j(d\theta_j)$ on the parameter.

- We are interested in settings where the prior may be vague or improper.

- **Example 1**. Normal i.i.d. models (*O'Hagan, 1995*)

| Model 1 | Model 2 |
|---|---|
| $Y_1, ..., Y_T \mid \theta_1 \overset{i.i.d.}{\sim} \mathcal{N}(\theta_1, 1)$ | $Y_1, ..., Y_T \mid \theta_2 \overset{i.i.d.}{\sim} \mathcal{N}(0, \theta_2)$ |
| $\theta_1 \sim \mathcal{N}(0, \sigma_0^2)$ | $\theta_2 \sim \text{Inv-}\chi^2(\nu_0, s_0^2)$ |

with known hyperparameters $\sigma_0 > 0$, $\nu_0 > 0$, and $s_0 > 0$.

- **Example 2**. Population dynamics of red kangaroos (*Knape & de Valpine, 2012*)



| Model 1 | Model 2 | Model 3 |
|---|---|---|
| Given parameters $(b, r, \sigma, \tau)$: | Given parameters $(r, \sigma, \tau)$: | Given parameters $(\sigma, \tau)$: |
| $X_1 \sim \text{LN}(0,5)$ | $X_1 \sim \text{LN}(0,5)$ | $X_1 \sim \text{LN}(0,5)$ |
| $dX_t/X_t = \left(r - bX_t + \frac{\sigma^2}{2}\right) dt + \sigma \, dW_t$ | $dX_t/X_t = \left(r + \frac{\sigma^2}{2}\right) dt + \sigma \, dW_t$ | $dX_t/X_t = \left(\frac{\sigma^2}{2}\right) dt + \sigma \, dW_t$ |
| $Y_{1,t}, Y_{2,t} \mid X_t \overset{\text{i.i.d.}}{\sim} \text{NB}(X_t, X_t + \tau X_t^2)$ | $Y_{1,t}, Y_{2,t} \mid X_t \overset{\text{i.i.d.}}{\sim} \text{NB}(X_t, X_t + \tau X_t^2)$ | $Y_{1,t}, Y_{2,t} \mid X_t \overset{\text{i.i.d.}}{\sim} \text{NB}(X_t, X_t + \tau X_t^2)$ |
| with independent priors: | with independent priors: | with independent priors: |
| $b, \sigma, \tau \sim \text{Unif}(0,10)$ | $\sigma, \tau \sim \text{Unif}(0,10)$ | $\sigma, \tau \sim \text{Unif}(0,10)$ |
| $r \sim \text{Unif}(-10,10)$ | $r \sim \text{Unif}(-10,10)$ | |

# Why not use Bayes factors ?

- Bayes factors favor the model $M_j$ with the largest *evidence*

$$p_j(y_{1:T}) = \int p_j(y_{1:T}|\theta_j) \, p_j(\theta_j) \, d\theta_j$$

(a.k.a. the *marginal likelihood* of $M_j$)

## Sensitivity to the choice of prior

- The evidence of any model can be made arbitrarily small by making the prior arbitrarily vague, regardless of the sample size.

- Bayes factors do not allow for improper priors.

Yet, vague or improper priors often stem from reasonable approaches (genuine non-informativeness, Jeffreys prior, … )

- Bayes factors favor the model $M_j$ with the largest *evidence*

$$p_j(y_{1:T}) = \int p_j(y_{1:T}|\theta_j)\, p_j(\theta_j)\, d\theta_j$$

(a.k.a. the *marginal likelihood* of $M_j$)

### Sensitivity to the choice of prior

- The evidence of any model can be made arbitrarily small by making the prior arbitrarily vague, regardless of the sample size.

- Bayes factors do not allow for improper priors.

Yet, vague or improper priors often stem from reasonable approaches
(genuine non-informativeness, Jeffreys prior, … )

- Bayes factors favor the model $M_j$ with the largest *evidence*

$$p_j(y_{1:T}) = \int p_j(y_{1:T}|\theta_j)\, p_j(\theta_j)\, d\theta_j$$

(a.k.a. the *marginal likelihood* of $M_j$)

### Sensitivity to the choice of prior

- The evidence of any model can be made arbitrarily small by making the prior arbitrarily vague, regardless of the sample size.

- Bayes factors do not allow for improper priors.

**Yet, vague or improper priors often stem from reasonable approaches**
(genuine non-informativeness, Jeffreys prior, … )

# A new criterion: the Hyvärinen score

- Bayes factors favor the model minimizing $-\log p(y_{1:T})$ or equivalently

$$\sum_{t=1}^{T} -\log p(y_t|y_{1:t-1})$$

- This is a particular case of a more general decision rule that favors the model with the smallest prequential score

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

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

Key idea: use the Hyvärinen score (H-score) instead *(Dawid & Musio, 2015)*

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

with $\mathcal{H}(y, p) = 2 \Delta \log p(y) + \|\nabla \log p(y)\|^2$, where the Laplacian $\Delta$ and the gradient $\nabla$ are taken with respect to the observation $y$

# Bayes factors relate to a particular choice of scoring rule

- Bayes factors favor the model minimizing $-\log p(y_{1:T})$ or equivalently

$$\sum_{t=1}^{T} -\log p(y_t|y_{1:t-1})$$

- This is a particular case of a more general decision rule that favors the model with the smallest prequential score

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

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

Key idea: use the Hyvärinen score (H-score) instead *(Dawid & Musio, 2015)*

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

with $\mathcal{H}(y, p) = 2\,\Delta \log p(y) + \|\nabla \log p(y)\|^2$, where the Laplacian $\Delta$ and the gradient $\nabla$ are taken with respect to the observation $y$

# Bayes factors relate to a particular choice of scoring rule

- Bayes factors favor the model minimizing $-\log p(y_{1:T})$ or equivalently

$$\sum_{t=1}^{T} -\log p(y_t|y_{1:t-1})$$

- This is a particular case of a more general decision rule that favors the model with the smallest prequential score

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

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

Key idea: use the Hyvärinen score (H-score) instead *(Dawid & Musio, 2015)*

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

with $\mathcal{H}(y, p) = 2\,\Delta \log p(y) + \|\nabla \log p(y)\|^2$, where the Laplacian $\Delta$ and the gradient $\nabla$ are taken with respect to the observation $y$

- **Example**. Given i.i.d. realizations $y_1, ..., y_{100}$ from a $\mathcal{N}(0,1)$, consider the model describing $Y_1, ..., Y_{100} \mid \mu \overset{i.i.d.}{\sim} \mathcal{N}(\mu, 1)$ with prior $\mu \sim \mathcal{N}(0, \sigma_0^2)$.

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - Proper: the expected loss $\mathbb{E}_\star[\mathcal{H}(Y, p)]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-local: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \ldots, p^{(m)}(y)$   (with $m = 2$)

- The H-score has the additional advantage of being 0-homogeneous
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \ldots, p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \ldots, T\}$. This calls for sequential estimation.

# Using H-scores is principled and justified non-asymptotically

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper**: the expected loss $\mathbb{E}_\star \left[ \mathcal{H}(Y, p) \right]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-local: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \ldots, p^{(m)}(y)$  (with $m = 2$)
- The H-score has the additional advantage of being 0-homogeneous
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \ldots, p^{(m)}(y)$ by any $\lambda > 0$
- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)
- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)
- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \ldots, T\}$. This calls for sequential estimation.

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t{}^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper**: the expected loss $\mathbb{E}_\star \left[ \mathcal{H}(Y, p) \right]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-**local**: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \dots, p^{(m)}(y)$    (with $m = 2$)

- The H-score has the additional advantage of being 0-homogeneous
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \dots, p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \dots, T\}$. This calls for sequential estimation.

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper**: the expected loss $\mathbb{E}_\star [\mathcal{H}(Y, p)]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-**local**: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \ldots, p^{(m)}(y)$  (with $m = 2$)

- The H-score has the additional advantage of being $0$-homogeneous
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \ldots, p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \ldots, T\}$. This calls for sequential estimation.

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper**: the expected loss $\mathbb{E}_{\star}[\mathcal{H}(Y, p)]$ under $Y \sim p_{\star}$ is minimized at $p = p_{\star}$
  - $m$-**local**: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \dots, p^{(m)}(y)$  (with $m = 2$)

- The H-score has the additional advantage of being 0-**homogeneous**
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \dots, p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \dots, T\}$. This calls for sequential estimation.

# Using H-scores is principled and justified non-asymptotically

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper**: the expected loss $\mathbb{E}_\star [\mathcal{H}(Y, p)]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-**local**: $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), ..., p^{(m)}(y)$ (with $m = 2$)

- The H-score has the additional advantage of being 0-**homogeneous**
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), ..., p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, ..., T\}$. This calls for sequential estimation.

- For univariate observations, the H-score can be written explicitly as

$$\sum_{t=1}^{T} \left( 2 \frac{\partial^2 \log p(y_t|y_{1:t-1})}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1})}{\partial y_t} \right)^2 \right)$$

- Similarly to the log-score, the H-score is
  - **Proper:** the expected loss $\mathbb{E}_\star [\mathcal{H}(Y, p)]$ under $Y \sim p_\star$ is minimized at $p = p_\star$
  - $m$-**local:** $\mathcal{H}(y, p)$ is a function of $y$ and $p(y), p'(y), \ldots, p^{(m)}(y)$   (with $m = 2$)

- The H-score has the additional advantage of being 0-**homogeneous**
  - $\mathcal{H}(y, p)$ is unchanged when multiplying $p(y), p'(y), \ldots, p^{(m)}(y)$ by any $\lambda > 0$

- The H-score is the simplest scoring rule satisfying propriety, $m$-locality, and 0-homogeneity (*Parry, Dawid & Lauritzen, 2012*)

- It can be extended to discrete observations via finite differences, while preserving all the above properties (*Dawid, Parry & Lauritzen, 2012*)

- But it involves derivatives of typically intractable predictive densities $p(y_t|y_{1:t-1})$ at every $t \in \{1, \ldots, T\}$. This calls for sequential estimation.

# How to estimate H-scores ?

# H-scores via SMC if likelihoods can be evaluated

· **By differentiating under the integral sign, the H-score turns out to be**

$$\sum_{t=1}^{T}\left(2\,\mathbb{E}_t\left[\frac{\partial^2 \log p(y_t|y_{1:t-1},\Theta)}{\partial y_t^2} + \left(\frac{\partial \log p(y_t|y_{1:t-1},\Theta)}{\partial y_t}\right)^2\right] - \left(\mathbb{E}_t\left[\frac{\partial \log p(y_t|y_{1:t-1},\Theta)}{\partial y_t}\right]\right)^2\right)$$

**where $\mathbb{E}_t$ denotes posterior expectations with respect to $\Theta \sim p(d\theta|y_{1:t})$**

**Estimation using SMC** *(Chopin, 2002; Del Moral, Doucet & Jasra, 2006)*
H-scores can be consistently estimated using standard SMC samplers,
as long as one can evaluate the incremental likelihoods $p(y_t|y_{1:t-1},\Theta)$

· However, incremental likelihoods $p(y_t|y_{1:t-1},\Theta)$ are typically intractable
for general state-space models.

- By differentiating under the integral sign, the H-score turns out to be

$$\sum_{t=1}^{T} \left( 2\,\mathbb{E}_t \left[ \frac{\partial^2 \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t} \right)^2 \right] - \left( \mathbb{E}_t \left[ \frac{\partial \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t} \right] \right)^2 \right)$$

where $\mathbb{E}_t$ denotes posterior expectations with respect to $\Theta \sim p(d\theta|y_{1:t})$

**Estimation using SMC** (*Chopin, 2002; Del Moral, Doucet & Jasra, 2006*)

H-scores can be consistently estimated using standard SMC samplers, as long as one can evaluate the incremental likelihoods $p(y_t|y_{1:t-1}, \Theta)$

- However, incremental likelihoods $p(y_t|y_{1:t-1}, \Theta)$ are typically intractable for general state-space models.

# H-scores via SMC if likelihoods can be evaluated

· By differentiating under the integral sign, the H-score turns out to be

$$\sum_{t=1}^{T} \left( 2\,\mathbb{E}_t \left[ \frac{\partial^2 \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t^2} + \left( \frac{\partial \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t} \right)^2 \right] - \left( \mathbb{E}_t \left[ \frac{\partial \log p(y_t|y_{1:t-1}, \Theta)}{\partial y_t} \right] \right)^2 \right)$$
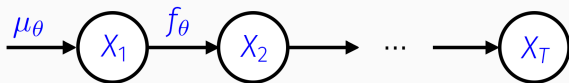
where $\mathbb{E}_t$ denotes posterior expectations with respect to $\Theta \sim p(d\theta|y_{1:t})$

**Estimation using SMC** (*Chopin, 2002; Del Moral, Doucet & Jasra, 2006*)

H-scores can be consistently estimated using standard SMC samplers, as long as one can evaluate the incremental likelihoods $p(y_t|y_{1:t-1}, \Theta)$

· However, incremental likelihoods $p(y_t|y_{1:t-1}, \Theta)$ are typically intractable for general state-space models.

$$\xrightarrow{\mu_\theta} \ (X_1) \xrightarrow{f_\theta} \ (X_2) \longrightarrow \ \cdots \ \longrightarrow \ (X_T)$$
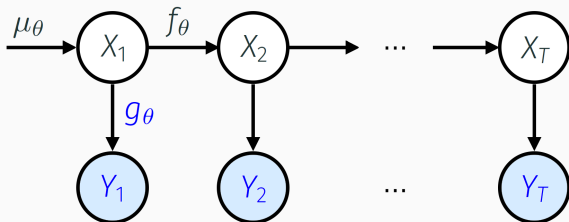
- Unobserved Markov chain of latent states $X_1, ..., X_T$ with
$$X_1 \sim \mu_\theta \quad \text{and} \quad X_t \,|\, X_{t-1} \sim f_\theta(\cdot | X_{t-1}) \quad \text{for } t \geq 2$$

- Observations $Y_1, ..., Y_T$ conditionally independent given $X_1, ..., X_T$ with
$$Y_t \,|\, X_t \sim g_\theta(\cdot | X_t) \quad \text{for } t \geq 1$$

The likelihood is generally intractable

$$p(y_{1:T}|\theta) = \int \mu_\theta(x_1) \prod_{t=2}^{T} f_\theta(x_t|x_{t-1}) \prod_{t=1}^{T} g_\theta(y_t|x_t) \, dx_{1:T}$$
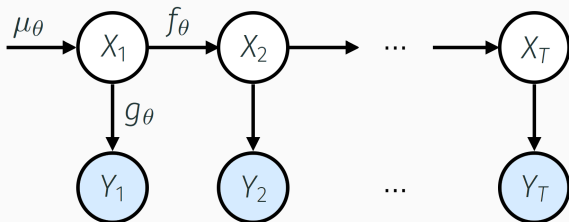
10

- Unobserved Markov chain of latent states $X_1, ..., X_T$ with
$$X_1 \sim \mu_\theta \quad \text{and} \quad X_t \,|\, X_{t-1} \sim f_\theta(\cdot \,|\, X_{t-1}) \quad \text{for } t \geq 2$$

- Observations $Y_1, ..., Y_T$ conditionally independent given $X_1, ..., X_T$ with
$$Y_t \,|\, X_t \sim g_\theta(\cdot \,|\, X_t) \quad \text{for } t \geq 1$$

The likelihood is generally intractable

$$p(y_{1:T} | \theta) = \int \mu_\theta(x_1) \prod_{t=2}^{T} f_\theta(x_t | x_{t-1}) \prod_{t=1}^{T} g_\theta(y_t | x_t) \, dx_{1:T}$$

- Unobserved Markov chain of latent states $X_1, ..., X_T$ with

$$X_1 \sim \mu_\theta \quad \text{and} \quad X_t \,|\, X_{t-1} \sim f_\theta(\cdot | X_{t-1}) \quad \text{for } t \geq 2$$

- Observations $Y_1, ..., Y_T$ conditionally independent given $X_1, ..., X_T$ with

$$Y_t \,|\, X_t \sim g_\theta(\cdot | X_t) \quad \text{for } t \geq 1$$

**The likelihood is generally intractable**

$$p(y_{1:T} | \theta) = \int \mu_\theta(x_1) \prod_{t=2}^{T} f_\theta(x_t | x_{t-1}) \prod_{t=1}^{T} g_\theta(y_t | x_t) \, dx_{1:T}$$

· **Under further integrability conditions, the H-score proves to equal**

$$\sum_{t=1}^{T} \left( 2\, \mathbb{E}_t \left[ \frac{\partial^2 \log g_\Theta(y_t|X_t)}{\partial y_t^2} + \left( \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right)^2 \right] - \left( \mathbb{E}_t \left[ \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right] \right)^2 \right)$$

**where the expectations $\mathbb{E}_t$ are now with respect to the joint posterior distributions $(\Theta, X_t) \sim p(d\theta|y_{1:t})p(dx_t|y_{1:t}, \theta)$**

**Estimation using SMC$^2$** *(Chopin, Jacob & Papaspiliopoulos, 2013)*
H-scores can be consistently estimated for state-space models, as long as one can simulate transitions from $f_\theta$ and evaluate $g_\theta$

· SMC and SMC$^2$ can be used similarly in the case of discrete observations.

- Under further integrability conditions, the H-score proves to equal

$$\sum_{t=1}^{T} \left( 2\, \mathbb{E}_t \left[ \frac{\partial^2 \log g_\Theta(y_t|X_t)}{\partial y_t^2} + \left( \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right)^2 \right] - \left( \mathbb{E}_t \left[ \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right] \right)^2 \right)$$

where the expectations $\mathbb{E}_t$ are now with respect to the joint posterior distributions $(\Theta, X_t) \sim p(d\theta|y_{1:t})p(dx_t|y_{1:t}, \theta)$

### Estimation using SMC$^2$ (*Chopin, Jacob & Papaspiliopoulos, 2013*)

H-scores can be consistently estimated for state-space models, as long as one can simulate transitions from $f_\theta$ and evaluate $g_\theta$

- SMC and SMC$^2$ can be used similarly in the case of discrete observations.

# H-scores via SMC$^2$ for general state-space models

- Under further integrability conditions, the H-score proves to equal

$$\sum_{t=1}^{T} \left( 2\, \mathbb{E}_t \left[ \frac{\partial^2 \log g_\Theta(y_t|X_t)}{\partial y_t^2} + \left( \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right)^2 \right] - \left( \mathbb{E}_t \left[ \frac{\partial \log g_\Theta(y_t|X_t)}{\partial y_t} \right] \right)^2 \right)$$

where the expectations $\mathbb{E}_t$ are now with respect to the joint posterior distributions $(\Theta, X_t) \sim p(d\theta|y_{1:t})p(dx_t|y_{1:t}, \theta)$

### Estimation using SMC$^2$ (*Chopin, Jacob & Papaspiliopoulos, 2013*)

H-scores can be consistently estimated for state-space models, as long as one can simulate transitions from $f_\theta$ and evaluate $g_\theta$

- SMC and SMC$^2$ can be used similarly in the case of discrete observations.

# Asymptotic guarantees

# Consistency of the H-score for model selection

- Assume the true data generating process is $Y_1, ..., Y_T \overset{i.i.d.}{\sim} p_\star$

- Consider two non-nested i.i.d. models $M_1$ and $M_2$ with respective posterior distributions concentrating around $\theta_1^\star$ and $\theta_2^\star$. Let $p_{\theta_j^\star} = p_j(dy|\theta_j^\star)$

- Under regularity conditions, the H-scores $\mathcal{H}_T$ of $M_1$ and $M_2$ satisfy

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

where $\Delta(p_\star, M_j) = \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$
with $\Delta(p_\star, M_j) \geq 0$ and $\Delta(p_\star, M_j) = 0$ if and only if $p_{\theta_j^\star} = p_\star$

- By analogy, the log-Bayes factor of $M_1$ against $M_2$ satisfies

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

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

- Similar results hold for dependent data and state-space models.

# Consistency of the H-score for model selection

- Assume the true data generating process is $Y_1, \ldots, Y_T \overset{\text{i.i.d.}}{\sim} p_\star$

- Consider two non-nested i.i.d. models $M_1$ and $M_2$ with respective posterior distributions concentrating around $\theta_1^\star$ and $\theta_2^\star$. Let $p_{\theta_j^\star} = p_j(dy|\theta_j^\star)$

- Under regularity conditions, the H-scores $\mathcal{H}_T$ of $M_1$ and $M_2$ satisfy

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

where $\Delta(p_\star, M_j) = \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$
with $\Delta(p_\star, M_j) \geq 0$ and $\Delta(p_\star, M_j) = 0$ if and only if $p_{\theta_j^\star} = p_\star$

- By analogy, the log-Bayes factor of $M_1$ against $M_2$ satisfies

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

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

- Similar results hold for dependent data and state-space models.

# Consistency of the H-score for model selection

- Assume the true data generating process is $Y_1, ..., Y_T \overset{\text{i.i.d.}}{\sim} p_\star$

- Consider two non-nested i.i.d. models $M_1$ and $M_2$ with respective posterior distributions concentrating around $\theta_1^\star$ and $\theta_2^\star$. Let $p_{\theta_j^\star} = p_j(dy|\theta_j^\star)$

- Under regularity conditions, the H-scores $\mathcal{H}_T$ of $M_1$ and $M_2$ satisfy

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

where $\Delta(p_\star, M_j) = \mathbb{E}_\star[\mathcal{H}(Y, p_{\theta_j^\star})] - \mathbb{E}_\star[\mathcal{H}(Y, p_\star)]$
with $\Delta(p_\star, M_j) \geq 0$ and $\Delta(p_\star, M_j) = 0$ if and only if $p_{\theta_j^\star} = p_\star$

- By analogy, the log-Bayes factor of $M_1$ against $M_2$ satisfies

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

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

- Similar results hold for dependent data and state-space models.

# Consistency of the H-score for model selection

- Assume the true data generating process is $Y_1, ..., Y_T \overset{\text{i.i.d.}}{\sim} p_\star$

- Consider two non-nested i.i.d. models $M_1$ and $M_2$ with respective posterior distributions concentrating around $\theta_1^\star$ and $\theta_2^\star$. Let $p_{\theta_j^\star} = p_j(dy|\theta_j^\star)$

- Under regularity conditions, the H-scores $\mathcal{H}_T$ of $M_1$ and $M_2$ satisfy

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

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

with $\Delta(p_\star, M_j) \geq 0$ and $\Delta(p_\star, M_j) = 0$ if and only if $p_{\theta_j^\star} = p_\star$

- **By analogy, the log-Bayes factor of $M_1$ against $M_2$ satisfies**

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

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

- Similar results hold for dependent data and state-space models.

# Consistency of the H-score for model selection

- Assume the true data generating process is $Y_1, ..., Y_T \overset{\text{i.i.d.}}{\sim} p_\star$

- Consider two non-nested i.i.d. models $M_1$ and $M_2$ with respective posterior distributions concentrating around $\theta_1^\star$ and $\theta_2^\star$. Let $p_{\theta_j^\star} = p_j(dy | \theta_j^\star)$

- Under regularity conditions, the H-scores $\mathcal{H}_T$ of $M_1$ and $M_2$ satisfy

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

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

with $\Delta(p_\star, M_j) \geq 0$ and $\Delta(p_\star, M_j) = 0$ if and only if $p_{\theta_j^\star} = p_\star$

- By analogy, the log-Bayes factor of $M_1$ against $M_2$ satisfies

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

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

- Similar results hold for dependent data and state-space models.

- **Example 1**. Given simulated $y_1, ..., y_{1000} \sim \mathcal{N}(\mu_\star, \sigma_\star^2)$, we compare

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

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

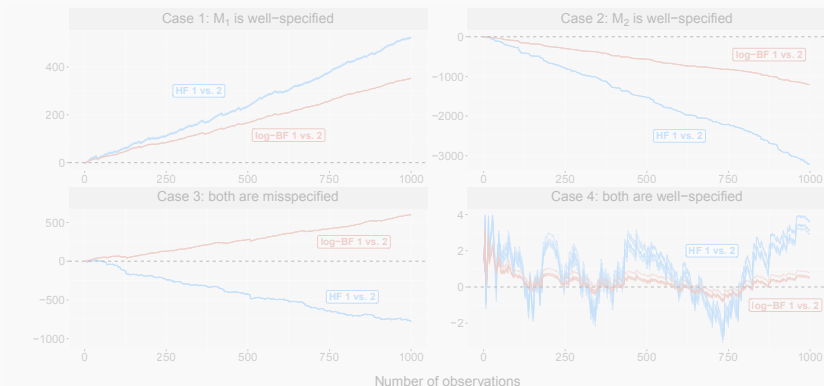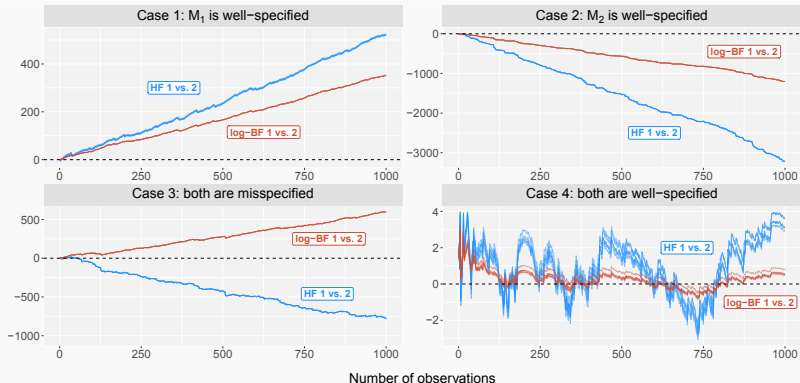in the following four cases $(\mu_\star, \sigma_\star^2) = (1, 1), (0, 5), (4, 3), (0, 1)$.



Number of observations

13

# Illustration of consistency for i.i.d. models

- **Example 1**. Given simulated $y_1, ..., y_{1000} \sim \mathcal{N}(\mu_\star, \sigma_\star^2)$, we compare

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

$$M_2: \quad Y_1, ..., Y_T \mid \theta_2 \overset{\text{i.i.d.}}{\sim} \mathcal{N}(0, \theta_2), \quad \theta_2 \sim \text{Inv-}\chi^2(0.1, 1)$$
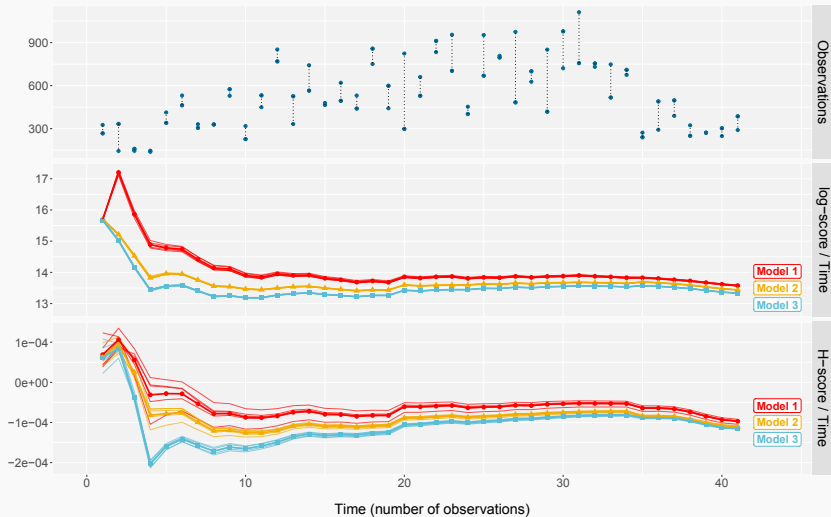
in the following four cases $(\mu_\star, \sigma_\star^2) = (1, 1), (0, 5), (4, 3), (0, 1)$.



Number of observations

13

# Applications and discussion

- **Example 2**. Population dynamics of red kangaroos (*Knape & de Valpine, 2012*)

## Advantages of using the H-score

- Robust to prior vagueness and allows for improper priors
- Justified non-asymptotically and generally consistent asymptotically
- Can be estimated sequentially for a wide class of parametric models by using SMC methods

## Limitations and avenues for future research

- Requires additional smoothness conditions on the densities
- Tends to require a larger number of particles than the log-evidence for accurate estimation via SMC methods
- Extension to nonparametric models ?

- More details in the manuscript: arxiv.org/pdf/1711.00136
- 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.
**SMC$^2$: an efficient algorithm for sequential analysis of state-space models.**
*Journal of the Royal Statistical Society*, 75 (3):397–426, 2013.

A. P. Dawid, S. Lauritzen, and M. Parry.
**Proper local scoring rules on discrete sample spaces.**
*The Annals of Statistics*, 40 (1):593–608, 2012.

A. P. Dawid and M. Musio.
**Bayesian model selection based on proper scoring rules.**
*Bayesian Analysis*, 10 (2):479–499, 2015.

# References (2/2)

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.

J. Knape and P. D. Valpine.
**Fitting complex population models by combining particle filters with markov chain monte carlo.**
*Ecology,* 93 (2):256–263, 2012.

F. B. Lempers.
*Posterior Probabilities of Alternative Linear Models.*
Rotterdam University Press, 1971.

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.

- **Example 3**. Lévy-driven stochastic volatility models for log-returns of financial assets (*Barndorff-Nielsen & Shephard, 2001*)

Given parameters $(\lambda, \xi, \omega)$, generate random variables $(V_t, Z_t)_{t \geq 1}$ recursively as

$$k \sim \text{Poisson}(\lambda \xi^2/\omega^2)\,;\quad C_{1:k} \overset{\text{i.i.d.}}{\sim} \text{Unif}(t-1,t)\,;\quad E_{1:k} \overset{\text{i.i.d.}}{\sim} \text{Exp}(\xi/\omega^2)\,;\quad Z_0 \sim \Gamma(\xi^2/\omega^2, \xi/\omega^2)$$

$$Z_t = e^{-\lambda} Z_{t-1} + \sum_{j=1}^{k} e^{-\lambda(t-C_j)} E_j\,;\quad V_t = \lambda^{-1}\left(Z_{t-1} - Z_t + \sum_{j=1}^{k} E_j\right) \quad \right\} \; (\text{⚙})$$

| Model 1 | Model 2 |
|---|---|
| Given parameters $(\lambda, \xi, \omega, \mu, \beta)$: | Given parameters $(\lambda_1, \lambda_2, w_1, w_2, \xi, \omega, \mu, \beta)$: |
| $(V_t, Z_t) \sim (\text{⚙})$ | $(V_{1,t}, Z_{1,t}) \sim (\text{⚙})$ with $(\lambda_1, \xi w_1, \omega w_1)$ |
| $X_t = (V_t, Z_t)$ | $(V_{2,t}, Z_{2,t}) \sim (\text{⚙})$ with $(\lambda_2, \xi w_2, \omega w_2)$ |
| $Y_t \mid X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$ | $X_t = (V_{1,t}, V_{2,t}, Z_{1,t}, Z_{2,t})$ |
| | $V_t = V_{1,t} + V_{2,t}$ |
| | $Y_t \mid X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$ |
| with independent priors: | with independent priors: |
| $\lambda \sim \text{Exp}(1)\,;\quad \xi, \omega^2 \sim \text{Exp}(1/5)\,;\quad \mu, \beta \sim \mathcal{N}(0,10)$ | $\lambda_1 \sim \text{Exp}(1)\,;\quad \lambda_2 - \lambda_1 \sim \text{Exp}(1/2)\,;\quad \xi, \omega^2 \sim \text{Exp}(1/5)$ |
| | $1 - w_2 = w_1 \sim \text{Unif}(0,1)\,;\quad \mu, \beta \sim \mathcal{N}(0,10)$ |

- **Example 3**. Lévy-driven stochastic volatility models, given $T = 1000$ observations simulated from a single-factor model with $\lambda = 0.01$, $\xi = 0.5$, $\omega^2 = 0.0625$, $\mu = 0$, and $\beta = 0$ (*Barndorff-Nielsen & Shephard, 2001*)

| Model 1 (single-factor) | Model 2 (multi-factor) |
|---|---|
| Given parameters $(\lambda, \xi, \omega, \mu, \beta)$: | Given parameters $(\lambda_1, \lambda_2, w_1, w_2, \xi, \omega, \mu, \beta)$: |
| $(V_t, Z_t) \sim (\text{⚙})$ | $(V_{1,t}, Z_{1,t}) \sim (\text{⚙})$ with $(\lambda_1, \xi w_1, \omega w_1)$ |
| $X_t = (V_t, Z_t)$ | $(V_{2,t}, Z_{2,t}) \sim (\text{⚙})$ with $(\lambda_2, \xi w_2, \omega w_2)$ |
| $Y_t \mid X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$ | $X_t = (V_{1,t}, V_{2,t}, Z_{1,t}, Z_{2,t})$ |
| | $V_t = V_{1,t} + V_{2,t}$ |
| | $Y_t \mid X_t \sim \mathcal{N}(\mu + \beta V_t, V_t)$ |
| with independent priors: | with independent priors: |
| $\lambda \sim \text{Exp}(1)$; $\quad \xi, \omega^2 \sim \text{Exp}(1/5)$; $\quad \mu, \beta \sim \mathcal{N}(0,10)$ | $\lambda_1 \sim \text{Exp}(1)$; $\quad \lambda_2 - \lambda_1 \sim \text{Exp}(1/2)$; $\quad \xi, \omega^2 \sim \text{Exp}(1/5)$ |
| | $1 - w_2 = w_1 \sim \text{Unif}(0,1)$; $\quad \mu, \beta \sim \mathcal{N}(0,10)$ |

# Illustration of consistency for state-space models

- **Example 3**. Lévy-driven stochastic volatility models, given $T = 1000$ observations simulated from a single-factor model with $\lambda = 0.01$, $\xi = 0.5$, $\omega^2 = 0.0625$, $\mu = 0$, and $\beta = 0$ (*Barndorff-Nielsen & Shephard, 2001*)

- Notice that, unlike for the log scoring rule, here we have:

$$\mathcal{H}\left(y_{1:T}, p(dy_{1:T})\right) \quad \neq \quad \sum_{t=1}^{T} \mathcal{H}\left(y_t, p(dy_t|y_{1:t-1})\right)$$

- Batch version:
    - Easier to compute, as it only requires to estimate final evidence $p(y_{1:T})$
    - But typically inconsistent for model selection

- Prequential version:
    - Generally consistent for model selection
    - Requires to estimate all the intermediary predictive densities $p(dy_t|y_{1:t-1})$, but this can be achieved using algorithms such as SMC or SMC[2]

# Partial Bayes factors (Lempers, 1971)

- Split the data $y_{1:T}$ into a training set $y_{1:m}$ and another set $y_{m+1:T}$ for some choice of $m$

- Idea: condition on the training set to make the prior proper (or less vague) then compute the Bayes factor on the remaining data

- Essentially we replace the prior $p(\theta|M)$ by the posterior given the training set $p(\theta|y_{1:m}, M)$, and compute the usual Bayes factor on the remaining data set $y_{m+1:T}$

- The partial Bayes factor between Models $M_1$ and $M_2$ is defined as:

$$\frac{p(y_{m+1:T}|y_{1:m}, M_1)}{p(y_{m+1:T}|y_{1:m}, M_2)}$$

- Drawback: choice of $m$ is a bit ad-hoc, undesirable to "waste" data for the training set especially in settings where the number of observations is small (e.g. Example 2 where $T = 41$)

# Fractional Bayes factors (O'Hagan, 1995)

- In the setting of partial Bayes factors, if $m$ and $T$ are both large, the likelihood $p(y_{1:m}|\theta, M)$ of the training set will approximate (at least in the i.i.d. case) the full likelihood raised to a power $b \equiv m/T$

- For a given model $M$ we define:

$$q_b(y_{1:T}|M) := \frac{\int p(\theta|M)p(y_{1:T}|\theta, M)d\theta}{\int p(\theta|M)p(y_{1:T}|\theta, M)^b d\theta}$$

which approximates $p(y_{m+1:T}|y_{1:m}, M)$ for large $m$ and $T$

- The fractional Bayes factor between Models $M_1$ and $M_2$ is defined as:

$$\frac{q_b(y_{1:T}|M_1)}{q_b(y_{1:T}|M_2)}$$

- Drawback: choice of $b$ is a bit ad-hoc, not very principled for small sample size since the main justification relies on asymptotics