

THE LONDON SCHOOL OF ECONOMICS AND POLITICAL SCIENCE

# IMS 2022

Sequential Bayesian Learning for Hidden Semi-Markov Models (HSMM)

Patrick Aschermayr Dr. Kostas Kalogeropoulos

June 20, 2022

Department of Statistics

# **Motivation**

## **Economies move in cycles**



Credit: Economics fun.

## Inference

## State Space Models

► A State Space Model (SSM) with parameter  $\theta \in \mathbb{R}^{D}$  is a bivariate stochastic process  $\{E_t, S_t\}_{t=1,2,...}$ , with the following distributional form:

$$\begin{array}{l} \boldsymbol{\theta} \sim \boldsymbol{p}(\theta), \\ \boldsymbol{S}_0 \sim \boldsymbol{p}(\boldsymbol{s}_0 \mid \theta), \\ \boldsymbol{S}_t \sim \boldsymbol{p}(\boldsymbol{s}_t \mid \boldsymbol{s}_{0:t-1}, \theta), \\ \boldsymbol{E}_t \sim \boldsymbol{p}(\boldsymbol{e}_t \mid \boldsymbol{e}_{1:t-1}, \boldsymbol{s}_{0:t}, \theta) \end{array}$$

► The goal is to infer the full posterior distribution :

$$p(s_{0:T}, \theta \mid e_{1:T}) = \frac{p(e_{1:T} \mid s_{0:T}, \theta) p(s_{0:T} \mid \theta) p(\theta)}{p(e_{1:T})}.$$
 (1)

SSMs can handle structural breaks, shifts, or time-varying parameters of a model and still have an interpretable structure. They are generative, and allow for multi step forecasting, imputing missing data, and account for non-equal time steps.

#### Challenges

- ► Marginal likelihood p(e<sub>1:T</sub>) intractable, but computation can be avoided.
- ► Naive batch estimation of high dimensional full posterior distribution p(s<sub>0:T</sub>, θ | e<sub>1:T</sub>) computationally unfeasible.
- ► Need to efficiently evaluate full posterior distribution iteratively as  $p(e_{1:T}, s_{0:T} | \theta) = p(s_0 | \theta) \prod_{t=1}^{T} p(s_t | s_{0:t-1}, \theta) p(e_t | e_{1:t-1}, s_{0:t}, \theta).$
- Marginal posterior distribution p(θ | e<sub>1:T</sub>) difficult to compute, as p(e<sub>1:t</sub> | θ) = ∫ p(e<sub>1:T</sub>, s<sub>0:T</sub> | θ) ds<sub>0:T</sub> is costly to evaluate or intractable.

#### Focus - HSMMs

- ► In HMM case,  $P(S_{t+k} = j, S_{t+1:t+k-1} = i | s_t = i)$  is implicitly geometric.
- ► HSMMs \* explicitly describe state durations:

$$\begin{split} \bullet \ \ e_t &\sim \mathcal{N}(\mu_{s_t}, \sigma_{s_t}) \\ \bullet \ \ s_t &\sim \begin{cases} \delta(S_t = s_{t-1}) & d_{t-1} > 0, \,^{\star\star} \\ P(S_t \mid s_{t-1}, d_{t-1}) & d_{t-1} = 0. \end{cases} \\ \bullet \ \ d_t &\sim \begin{cases} \delta(S_t = s_{t-1}) & d_{t-1} > 0, \\ P(S_t \mid s_{t-1}, d_{t-1}) & d_{t-1} = 0. \end{cases} \end{aligned}$$

► Can compute  $\sum_{s_{0:T}, d_{0:T}} p(e_{1:T}, s_{0:T}, d_{0:T} \mid \theta)$  in  $\mathcal{O}(K^2(d_{max} - d_{min})^2 T)^{***}$  for HSMM instead of  $\sum_{s_{0:T}} p(e_{1:T}, s_{0:T} \mid \theta)$  in  $\mathcal{O}(K^2 T)^{****}$  for HMM.

<sup>(\*)</sup> see (Yu, 2010) and (Yu, 2016)

<sup>(\*\*)</sup>  $\delta(a, b)$  is the Kronecker product and equals 1 if a = b and 0 otherwise.

<sup>(\*\*\*)</sup>  $d_{min}$  = minimal state duration,  $d_{max}$  = maximal state duration, typically ( $d_{max} - d_{min}$ ) >> K

<sup>(\*\*\*\*)</sup> see (Baum and Petrie, 1966) . K = number of latent states, T = number of data points.

#### Focus - HSMMs continued





(a) K-state Bayesian HMM, parameter  $\theta$  and hyper-parameter  $\{\beta,\gamma\}$ 

(b) K-state Bayesian HSMM, parameter  $\theta$  and hyper-parameter  $\{\alpha,\beta,\gamma\}$ 

- ► In HMM case, P(S<sub>t+k</sub> = j, S<sub>t+1:t+k-1</sub> = i | s<sub>t</sub> = i) is implicitly geometric, HSMMs can explicitly model state durations.
- ► Can compute  $\sum_{s_{0:T}, d_{0:T}} p(e_{1:T}, s_{0:T}, d_{0:T} | \theta)$  in  $\mathcal{O}(K^2(d_{max} d_{min})^2 T)^*$  for HSMM instead of  $\sum_{s_{0:T}} p(e_{1:T}, s_{0:T} | \theta)$  in  $\mathcal{O}(K^2 T)^{**}$  for HMM.

Inference

Applications

Impa

Appendix

## **Motivation and contribution**

#### Particle MCMC

- Independent of continuity of s<sub>t</sub>, can decompose problem into targetting p(s<sub>0:τ</sub> | e<sub>1:τ</sub>, θ) and p(θ | e<sub>1:τ</sub>, s<sub>0:τ</sub>):
  - ► Approximate  $p(s_{0:T} | e_{1:T}, \theta)$  via a particle filter (PF \*).
  - ► Target  $p(\theta | e_{1:T}, s_{0:T})$  via MCMC.
  - ► Formally known as Particle Gibbs \*\*.
- Can compute PF estimate for both p(s<sub>0:T</sub> | e<sub>1:T</sub>, θ) and p(e<sub>0:T</sub> | θ) in O(NT) \*\*\*.

<sup>(\*)</sup> see, e.g., Doucet and Johansen (2011)

<sup>(\*\*)</sup> see (Andrieu and Roberts, 2009), (Andrieu et al., 2010), (Lindsten et al., 2014) and (Lindsten et al., 2015)

<sup>(\*\*\*)</sup> N = number of particles, typically K << N < T.

## **Sequential Estimation**

#### ► Obtain **posterior predictive distribution** by integrating $s_{0:T} \& \theta$ :

$$p(e_{T+1} \mid e_{1:T}) = \int p(e_{T+1}, s_{T+1}, s_{0:T}, \theta \mid e_{1:T}) \, ds_{T+1}, s_{0:T}, \theta$$
  
= 
$$\int p(e_{T+1} \mid s_{T+1}, s_{0:T}, \theta, e_{1:T}) \, p(s_{T+1} \mid s_{0:T}, \theta, e_{1:T}) \, p(s_{0:T}, \theta \mid e_{1:T}) \, ds_{T+1}, s_{0:T}, \theta$$
  
(2)

- ▶ Sampling  $S_{T+1}$  and  $E_{T+1}$  trivial after  $p(s_{0:T}, \theta | e_{1:T})$  is obtained.
- ▶ Goal: sequentially explore  $p(s_{0:t}, \theta | e_{1:t})$  for t = 1, ..., T.

Inference

Applications

## Sequential Monte Carlo Squared \*

- ► Explore n sequences of distributions  $p(s_{0:t}^n, \theta^n | e_{1:t})$  for t = 1, ..., T.
  - ► Calculate  $p(e_t | e_{1:t-1}, \theta^n)$ ,  $p(e_{1:t} | \theta^n)$  and propagate  $s_{0:t}^n$  online via PF.
  - ▶ If  $p(e_t | e_{1:t-1}, \theta)$  estimates too noisy, jitter  $s_{0:t}^n, \theta^n$  via Particle Gibbs.
  - Almost real time.
- ► Obtain predictive distributions for e<sub>t+1</sub> and s<sub>t+1</sub> and an estimate for marginal likelihood p(e<sub>1:t</sub>) for each t = 1,...,T.

► Use CRPS \*\* to compare predictive distribution of models. For forecasts X<sub>i</sub>, i = 1,..., m and observation y, CRPS can be calculated as

$$CRPS(\hat{F}_m, y) = \frac{1}{m} \sum_{i=1}^m |X_i - y| - \frac{1}{2m^2} \sum_{i=1}^m \sum_{j=1}^m |X_i - X_j|.$$
(3)

(\*) see Chopin (2002) and Chopin et al. (2013)

<sup>(\*\*)</sup> see, e.g., (Jordan et al., 2019)

### State Space Models and Financial Data

#### ► Stylized financial facts \*:

- ▶ (u.1) Returns not iid, but show little serial correlation.
- ▶ (u.2) Extreme returns appear in clusters.
- ▶ (u.3) Returns have heavy tails.
- ▶ (u.4) Volatility clusters and varies over time.
- ▶ U.S. economic cycles widely vary in duration \*\*.

#### ► Apply SMC<sup>2</sup> for HSMM on financial data:

- are model parameter constant across time?
- how does HSMM fare against other SSM?

<sup>(\*)</sup> see, e.g., (McNeil et al., 2005)

<sup>(\*\*)</sup> see, e.g., https://www.nber.org/cycles.html

### **Results - HSMM**



### **Results - Prediction and Model Comparison**



CRPS score for various models:

•	Model	CRPS Score
	HSMM	228.13
	SV	229.43
	HMM	230.15

More analysis needed!

# Impact

## **Research contribution**

### ► Model and Applications contribution by

► providing alternative ways for parameter estimation on HSSMs.

#### ► Algorithmic contribution by

- providing a toolbox for estimation and further inference on SSMs with arbitrary state and observation dependency that will be open sourced over the next months.
- providing ideas for automatic adaption of SMC<sup>2</sup> tuning parameter, such as the number of jittering steps.

Motivation

Inferenc

Applications

Impact

Appendix

# Discussion

Appendix

## HMM as mixture distribution

HMM as sequential mixture:



$$\begin{aligned} \mathcal{P}(e_{t+1} \mid s_t = k) &= \sum_{s_{t+1}} \mathcal{P}(e_{t+1}, s_{t+1} \mid s_t = k) \\ &= \sum_{s_{t+1}} \mathcal{P}(s_{t+1} \mid s_t = k) \mathcal{P}(e_{t+1} \mid s_{t+1}) \end{aligned}$$

For a discrete 2-state, homogenous Markov chain, using the chain rule and the Markov assumption, it holds:

$$P(S_{t+3} = j, S_{t+2} = i, S_{t+1} = i \mid S_t = i) = P(S_{t+3} = j \mid S_{t+2} = i)P(S_{t+2} = i, \mid S_{t+1} = i)$$
$$= (1 - \mathcal{T}_{ii}) * \mathcal{T}_{ii}^2$$

In general, for t + k steps:

$$P(S_{t+k} = j, \dots, S_{t+1} = i \mid S_t = i) = (1 - \mathcal{T}_{ii}) * \mathcal{T}_{ii}^{k-1}$$
$$= Geometric_{\mathcal{T}_{ii}},$$

where the geometric distribution has to be interpreted as the length of state duration up to and including the transition to the other state.

#### **HSMM** distribution

In an EDHMM, transitions are allowed only at the end of each state, resulting in the following distributional forms:

$$S_{t} \mid s_{t-1}, d_{t-1} \sim P(S_{t} \mid s_{t-1}, d_{t-1}) = \begin{cases} \delta(S_{t} = s_{t-1}) & d_{t-1} > 0 \\ \mathcal{T}_{s_{t-1}, \dots} & d_{t-1} = 0 \end{cases}$$
(4)

$$D_t \mid s_t, d_{t-1} \sim P(D_t \mid s_t, d_{t-1}) = \begin{cases} \delta(D_t = d_{t-1} - 1) & d_{t-1} > 0\\ \mathcal{D}_{s_t} & d_{t-1} = 0 \end{cases}$$
(5)

$$E_t \mid s_t \sim \mathcal{O}_{s_t}$$
 (6)

where  $\delta(a, b)$  is the Kronecker product and equals 1 if a = b and 0 otherwise. Given equation 4 and 5, we can write the distribution for  $Z_t = \{S_t, D_t\}$ 

$$P(Z_t \mid z_{t-1}) = P(S_t \mid s_{t-1}, d_{t-1})P(D_t \mid s_t, d_{t-1})$$
  
= 1 \langle \mathcal{T}\_{s\_{t-1}, ..} \mathcal{D}\_{s\_t} (7)

The joint distribution of an EDHMM given the parameter corresponding to the graphical model can be stated as

$$P(S, D, E \mid \theta) = P(s_0 \mid \mathcal{T}_0)P(d_0 \mid \mathcal{D}_0) \prod_{t=1}^T P(s_t \mid s_{t-1}, d_{t-1}, \mathcal{T})P(d_t \mid s_t, d_{t-1}, \mathcal{D})P(e_t \mid s_t, \mathcal{O})$$
(8)

### **Particle Filter Algorithm**

input : Observation  $e_{1,T}$ , importance distribution  $\pi$ , parameter  $\theta = \{ \mathcal{Z} = \{ \mathcal{D}, T \}, \mathcal{O} \}$ , ParticleNumber N output: Particles  $X_{1,T}^{1:N}$ , Weights  $w_{1,T}^{1:N}$ , Weights.Normalized  $W_{1,T}^{1:N}$ // Initialize particles and weights for  $n \leftarrow 1$  to N do Sample particle  $X_1^n \sim Z_0$ ; Compute  $w_1^n(x_1^n) = \frac{\mathcal{O}_{\chi_1^n}(e_1)\mathcal{Z}_0(x_1^n)}{\frac{\pi(\chi_1^n)e_1}{\pi(\chi_1^n)e_1}}$ end Resample  $(x_1^n, w_1^n)_{n=1 \cdot N}$  with replacement to get equally weighted particles  $(\tilde{x}_1^n, \frac{1}{N})_{n=1 \cdot N}$ ; // Recursively calculate probabilities of interest for  $t \leftarrow 2$  to T do for  $n \leftarrow 1$  to N do 1. Sample  $X_t^n \sim \pi(X_t \mid \tilde{x}_{t-1}^n, e_t)$ ; 2. Set  $X_{1:t}^n = (\tilde{x}_{1:t-1}^n, x_t^n)$ ; 3. Calculate weights:  $w_t^n(X_{1:t}^n) \propto w_{t-1}^n(X_{1:t-1}^n) \frac{\mathcal{O}_{x_t^n(e_t)} \mathcal{Z}_{x_{t-1}^n, x_t^n}}{\pi(x_{t-1}^n)}$ end i. Normalize all N weights:  $W_t^n = \frac{w_t^n(X_{1:t}^n)}{\sum_i w_i^i(X_{1:t}^i)}$ ii. Resample  $(x_{1:t}^n, w_t^n)_{n=1:N}$  with replacement to get equally weighted particles  $(\tilde{x}_{1:t}^n, \frac{1}{N})_{n=1:N}$ ; end

Algorithm 1: General particle filter algorithm

### Particle Filter time complexity

**Forward backward algorithm in basic HMM:**  $\mathcal{O}(K^2T)$ , where K is the number of states and T is the time. At each time point t, one needs to evaluate both, the forward and the backward probabilities for all hidden states. If one would not use this iterative procedure and just try a brute force method to find all possible state sequences, one would have a time complexity of  $\mathcal{O}(K^TT)$ .

Forward backward algorithm in HSMM: In addition to the basic HMM complexity, one needs to truncate the sequence to a minimum and maximum duration,  $d_{min}$  and  $d_{max}$ . The computational complexity then becomes  $\mathcal{O}(K^2(d_{max} - d_{min})^2T)$ , where typically  $(d_{max} - d_{min}) >> K$  might go from 0 to *T*. Particle Filter algorithm in basic HMM and HSMM: Computational complexity both linear in time *T* and in number of particles *N*, so complexity is  $\mathcal{O}(NT)$ . However, if I would do forward filtering an backward something as well, the complexity would then also be  $\mathcal{O}(N^2T)$ . Usually, N >> K, but if K is growing (Infinite HMM/HSMM), particle filter might be faster than forward-backward algorithms.

**Conditional Particle Filter**: is a special case, where a reference trajectory guides the particle filter, making it possible to use very few particles and also use backward smoothing efficiently.

```
input : Proposal distribution Q, iterationNumber N,
               Particle filter proposal \pi, particleNumber M,
               observation e1.T
output: (\theta^i, Z_{1,\tau}^i)_{i=1:N}
Initialize \theta:
Run particle filter \rightarrow get \hat{P}(e_{1:T} \mid \theta).
for i \leftarrow 1 to N do
        1. Propose a new \theta^*, \theta^* \sim Q(\theta^* \mid \theta);
       2. Run particle filter \rightarrow get \hat{P}(e_{1:T} \mid \theta^*) and Z_{1:T}^*;
       3. Accept the pair (\theta^*, Z^*_{1,\tau}) with probability:
                                           \min(1, \frac{\hat{P}(e_{1:T} \mid \theta^{\star})}{\hat{P}(e_{1:T} \mid \theta)} \frac{P(\theta^{\star})}{P(\theta)} \frac{Q(\theta \mid \theta^{\star})}{Q(\theta^{\star} \mid \theta)})
         4. If accepted, set \hat{P}(e_{1:T} \mid \theta) = \hat{P}(e_{1:T} \mid \theta^*) and \theta = \theta^*.
end
```

Algorithm 2: Particle Metropolis Hastings algorithm

#### Hamiltonian Monte Carlo Primer

(1) To draw from posterior of interest, introduce auxiliary momentum variable  $\rho$  and draw from the joint density  $p(\rho, \theta) = p(\rho \mid \theta)p(\theta)$ . Usually,  $\rho$  does not depend on  $\theta$ .  $p(\rho, \theta)$  define a Hamiltonian

$$H(\rho, \theta) = -\log p(\rho, \theta)$$
  
=  $-\log p(\rho \mid \theta) - \log p(\theta)$  (9)  
=  $T(\rho \mid \theta) + V(\theta),$ 

where  $T(\rho \mid \theta)$  is called "kinetic energy", and  $V(\theta)$  "potential energy" ( $\infty$  - log posterior). (2) Joint system { $\rho, \theta$ } evolves via Hamiltonian equations

$$\frac{d\theta}{dt} = +\frac{\partial H}{\partial \rho} = +\frac{\partial T}{\partial \rho}$$

$$\frac{d\rho}{dt} = -\frac{\partial H}{\partial \theta} = -\frac{\partial T}{\partial \theta} - \frac{\partial V}{\partial \theta}$$
(10)

(3) To solve two-state differential equations in (2), can use, e.g., leapfrog integrator.

(a) First, sample ρ ~ MvNormal(0, M)

(b) Alternate half-step updates of the momentum and full-step updates of the position L times (for some discretization size  $\epsilon$ ):

(i) 
$$\rho \leftarrow \rho - \frac{\epsilon}{2} \frac{\partial V}{\partial \theta}$$
  
(ii)  $\theta \leftarrow \theta + \epsilon M^{-1} \rho$   
(iii)  $\rho \leftarrow \rho - \frac{\epsilon}{2} \frac{\partial V}{\partial \theta}$ 

(c) half-step back for the momentum variable  $\rho$ 

(d) Apply a Metropolis acceptance step to account for numerical errors, which can be stated via the Hamiltonians,

$$\alpha = \min(1, \exp(H(\rho, \theta) - H(\rho^*, \theta^*))) \qquad (11)$$

(4) Once (3) is finished, discard momentum variable to have a draw from the posterior via HMC. Many different variations available, also ways to automate tuning of  $\epsilon$  or L (or both if integration time dynamic) and M.

## References

- Andrieu, C., Doucet, A., and Holenstein, R. (2010). Particle markov chain monte carlo methods. *Journal of the Royal Statistical Society Series B*, 72(3):269–342.
- Andrieu, C. and Roberts, G. O. (2009). The pseudo-marginal approach for efficient monte carlo computations. *Ann. Statist.*, 37(2):697–725.
- Baum, L. and Petrie, T. (1966). Statistical inference for probabilistic functions of finite state markov chains. *The Annals of Mathematical Statistics*, 37:1554–1563.

#### References ii

- Chopin, N. (2002). A sequential particle filter method for static models. *Biometrika*, 89(3):539–552.
- Chopin, N., Jacob, P. E., and Papaspiliopoulos, O. (2013). Smc2: an efficient algorithm for sequential analysis of state-space models.
- Doucet, A. and Johansen, A. (2011). A tutorial on particle filtering and smoothing: Fifteen years later.
- Jordan, A., Krüger, F., and Lerch, S. (2019). Evaluating probabilistic forecasts with scoringrules. *Journal of Statistical Software, Articles*, 90(12):1–37.
- Lindsten, F., Bunch, P., Singh, S. S., and Schön, T. B. (2015). Particle ancestor sampling for near-degenerate or intractable state transition models.

- Lindsten, F., Jordan, M. I., and Schön, T. B. (2014). Particle gibbs with ancestor sampling.
- McNeil, A., Frey, R., and Embrechts, P. (2005). *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton University Press, USA.
- Yu, S. (2016). *Hidden Semi-Markov Models: Theory, Algorithms and Applications*. Elsevier, Boston.
- Yu, S.-Z. (2010). Hidden semi-markov models. *Artificial Intelligence*, 174(2):215–243. Special Review Issue.