

Final draft

Introduction

Chapter 1

Markov Processes and State-Space Models

1.1 Notation

Let us fix the basic notation that will be employed henceforth. Upper case letters, e.g. X_t, Y_t , denote random variables, while lower case letters denote the realizations. Finite sequences are made compact using the semi-colon notation, e.g. $x_{0:t} = (x_0, x_1, \dots, x_t)$ for some $t \geq 0$. The sets of outcomes are denoted by upper case calligraphic letters, e.g. \mathcal{X}, \mathcal{Y} . The σ -algebra of a set \mathcal{X} reads $\mathcal{B}(\mathcal{X})$. The pair $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ is a measurable space. The densities are represented with lower case cursive letters. For example, the density of X_t in x_t is $p_t^\theta(x_t)$, where $\theta \in \Theta$ is some (eventually known) vector, while the density of X_t in x_t , given that X_{t-1} has previously taken value x_{t-1} , is $p_t^\theta(x_t | x_{t-1})$. Note that $p_t^\theta(x_t | x_{t-1})$ is a member of the parametric family with parameters (θ, x_{t-1}) . The probability distribution of a random variable X_t is represented by $\mathbb{P}_t(dx_t)$, while the probability distribution of a sequence $X_{0:t}$ is represented by $\mathbb{P}_t(dx_{0:t})$. The Dirac measure, i.e., the measure that assigns probability 1 to the singleton $\{x\}$, is denoted by $\delta_x(dy)$. The expectation operator is denoted by \mathbb{E} , where the section of \mathbb{E} at \mathbb{Q} , i.e., $\mathbb{E}_{\mathbb{Q}}$ denotes the operator computed with respect to the probability distribution \mathbb{Q} . Let φ be a function defined on the random variable $X \propto \mathbb{Q}$. Then, we denote the expected value of $\varphi(X)$ with

$$\mathbb{E}_{\mathbb{Q}}[\varphi(X)] = \int_{\mathbb{X}} \varphi(x) \mathbb{Q}(dx).$$

1.2 Markov Processes

Let $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ and $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$ be two (eventually equal) measurable spaces.

Definition A function $P(x, dy) : (\mathcal{X}, \mathcal{B}(\mathcal{X})) \rightarrow [0, 1]$ such that

- for all $x \in \mathcal{X}$, $P(x, \cdot)$ is a probability measure on $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$, and
- for all subsets $A \in \mathcal{B}(\mathcal{Y})$, the map $x \mapsto P(x, A)$ is measurable in $\mathcal{B}(\mathcal{X})$,

is defined a probability kernel from $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ to $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$.

Notice that for any random variables X_1 and $X_0 = x_0$,

$$\mathbb{P}_1(dx_{0:1}) = \mathbb{P}_0(dx_0)P_1(x_0, dx_1),$$

where $P_1(x_0, dx_1)$ is a probability kernel from $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ to $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. This remark, provides some intuition about the link between probability kernels and conditional probabilities. As a matter of fact, it is possible to show that for given $X_{0:1}$ with probability distribution $\mathbb{P}_1(dx_{0:1})$,

$$\mathbb{P}_1(X_1 \in dx_1 \mid X_0 = x_0) = P_1(x_0, dx_1). \quad (1.1)$$

Probability kernels can be used to define Markov processes. Before moving to the definition, we will introduce some important concepts. Recall that given two measures \mathbb{M} and \mathbb{Q} defined on a measurable space $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$, \mathbb{Q} is absolutely continuous with respect to \mathbb{M} , that is $\mathbb{Q} \ll \mathbb{M}$, if for all sets $A \in \mathcal{B}(\mathcal{X})$ such that $\mathbb{M}(A) = 0$, then it must also hold that $\mathbb{Q}(A) = 0$. If $\mathbb{Q}(A)/\mathbb{M}(A)$ is well-defined for all A , then the Radon-Nikodym theorem guarantees that there exists a measurable function $w(x) \geq 0$ such that

$$w(x) = \frac{\mathbb{Q}(dx)}{\mathbb{M}(dx)},$$

meaning that,

$$\mathbb{Q}(A) = \int_{\mathcal{X}} w(x) \mathbb{M}(dx).$$

Let P_1, P_2, \dots, P_T be a finite measure of probability kernels from $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ to $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. Fix some $\mathbb{P}_0(dx)$ on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$.

Definition A discrete-time Markov process is a sequence $X_{0:T}$ of random variables whose joint distribution can be written as

$$\mathbb{P}_T(X_{0:T} \in dx_{0:T}) = \mathbb{P}_0(dx_0) \prod_{s=1}^T P_s(x_{s-1}, x_s). \quad (1.2)$$

The set \mathcal{X} is the state-space, \mathbb{P}_0 is the initial distribution and the probability kernel P_t is the transition kernel.

It is possible to show that for any $t \in T$,

$$\mathbb{P}_T(X_t \in dx_t \mid X_{0:t-1} \in dx_{0:t-1}) = \mathbb{P}_T(X_t \in dx_t \mid X_{t-1} = x_{t-1}) = \mathbb{P}_t(x_{t-1}, dx_t).$$

The first equality implies conditional independence. The second equality implies that it is possible to identify conditional distributions with transition kernels. Assume $\mathbb{P}_t, t \leq T$, is a sequence of probability measures. Then, the following proposition holds.

Proposition 1

For all $t \leq T$,

$$\mathbb{P}_T(dx_{0:t}) = \mathbb{P}_0(dx_0) \prod_{s=1}^t P_s(x_{s-1}, dx_s) = \mathbb{P}_t(dx_{0:t}).$$

This means that the marginal distribution of $X_{0:t}$ with respect to \mathbb{P}_T is \mathbb{P}_t .

1.3 State-Space Models

Chapter 2

Linear-Gaussian case

Chapter 3

Particle Filter

In the previous chapter we explained how the Kalman Filter is applied to the case of a linear model with Gaussian errors to find closed-form solution for the filtering and the predictive distributions. Nonetheless, the majority of time-series models have a non-linear structure and errors might follow non-Gaussian distributions.

In this chapter we describe particle filter as an algorithm that allows to update sequentially the filtering and the predictive distributions in the non-linear and non-Gaussian case. Additionally, the Kalman Filter can be seen as a special case of the more general particle filter.

In section 3.1 we introduce the fundamental notions of importance sampling, sequential importance sampling and resampling. Then, in section 3.2 we describe three particle filtering algorithms: the bootstrap, the guided and the {auxiliary particle filters}. Instead, in section 3.3 we describe the {Liu and West filter}, when particle filtering is used to estimate model's parameters. Finally, in section 3.4 we establish some general convergence results for Particle Filter - Monte Carlo estimators.

3.1 Introductory notions: sampling and resampling

Let $\varphi : \mathcal{X} \rightarrow \mathbb{R}$ be a measurable, integrable and bounded real-valued function, where $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbb{Q})$ is a probability space.¹ Suppose that the probability measure \mathbb{Q} admits a density $q(x)$ with respect to some measure $\nu(dx)$ and it is feasible to sample from that density.

¹ $\mathcal{B}(\mathcal{X})$ is the Borel sigma-algebra of \mathcal{X} .

To estimate the expected value of φ with respect to the density q , the basic Monte Carlo (MC) method would require to sample N times from the target density q and to compute the average of the function φ values. $\hat{\varphi}_{MC}$ is the standard MC estimator:

$$\hat{\varphi}_{MC} = \frac{1}{N} \sum_{n=1}^N \varphi(X^n) \approx \mathbb{E}_q(\varphi) \quad X^n \sim q(\cdot)$$

However, in many circumstances either it is impossible to sample from the target density or there exist another density - call it $m(\cdot)$ - different from q that produces more efficient estimates. In all the cases for which sampling takes place from a distribution different from the correct (i.e., the true) one, importance sampling and resampling methods can be applied (as augmented versions of the standard MC algorithm).

3.1.1 Importance sampling

Let us consider two probability spaces: $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbb{Q})$ and $(\mathcal{X}, \mathcal{B}(\mathcal{X}), \mathbb{M})$ and let the measurable function $w : \mathcal{X} \rightarrow \mathbb{R}$ be proportional to the Radom-Nikodym derivative between \mathbb{Q} and \mathbb{M} , i.e., $w(x) \propto \frac{\mathbb{Q}(dx)}{\mathbb{M}(dx)}$.²

The following result represents the starting point for the idea of importance sampling.

Lemma 1

For any measurable, integrable and bounded function $\varphi : \mathcal{X} \rightarrow \mathbb{R}$, the following holds:

$$\mathbb{E}_{\mathbb{M}}(\varphi \cdot w) = \mathbb{E}_{\mathbb{Q}}(\varphi) \cdot \mathbb{E}_{\mathbb{M}}(w) \tag{3.1}$$

Proof. Note that $\mathbb{Q}(dx) = \frac{w(x)}{\mathbb{E}_{\mathbb{M}}(w)} \mathbb{M}(dx)$. Hence:

$$\mathbb{E}_{\mathbb{Q}}(\varphi) = \int_{\mathcal{X}} \varphi(x) \mathbb{Q}(dx) = \int_{\mathcal{X}} \varphi(x) \frac{w(x)}{\mathbb{E}_{\mathbb{M}}(w)} \mathbb{M}(dx) = \frac{\mathbb{E}_{\mathbb{M}}(\varphi \cdot w)}{\mathbb{E}_{\mathbb{M}}(w)}$$

■

Lemma 1 represents the key theoretical foundation of importance sampling (IS), a numerical algorithm that allows to sample from a proposal density to estimate moments that are computed with respect to the target density, by means of re-weighting the values of the target function φ .

Adopting the previous notation, the target measure is \mathbb{Q} , while the proposal is \mathbb{M} . IS requires to

²From the Radom-Nikodym theorem a sufficient condition for w to exist is that \mathbb{Q} is absolutely continuous with respect to \mathbb{M} : $\mathbb{Q} \ll \mathbb{M}$.

sample N values from \mathbb{M} , then to find the associated normalized³ weights $W(X^n)$, finally to compute a weighted average of the values taken by the function φ . The IS is implemented according to the following algorithm.

1. Sample N values from the proposal measure: $X^n \sim \mathbb{M}$ for $n = 1, \dots, N$.
2. Compute the associated normalized weights: $W(X^n)$ for $n = 1, \dots, N$
3. Compute the estimator:

$$\hat{\varphi} = \frac{1}{N} \sum_{n=1}^N \varphi(X^n) W(X^n).$$

Note that $\hat{\varphi}$ is an unbiased estimator for the target of inference $\mathbb{E}_{\mathbb{Q}}(\varphi)$. In fact:

$$\mathbb{E}_{\mathbb{M}}(\hat{\varphi}) = \frac{1}{N} \sum_{n=1}^N \mathbb{E}_{\mathbb{M}}(\varphi \cdot W) = \mathbb{E}_{\mathbb{M}}(\varphi \cdot W) = \mathbb{E}_{\mathbb{Q}}(\varphi)$$

where the last equality follows from lemma 1 and the fact that weights are normalized to 1.

3.1.1.1 Auto-normalized IS

Usually, the Radom-Nikodym derivative is known up to a normalizing constant and the IS algorithm has to be described in more general terms. In the previous section the normalized weighting function $W(\cdot)$ was the Radom-Nikodym derivative of \mathbb{Q} with respect to \mathbb{M} . More generally, consider the function $w(x) \propto \frac{\mathbb{Q}(dx)}{\mathbb{M}(dx)}$, denoted as importance function, proportional to the derivative $W(\cdot)$.

In this case, weights have to be auto-normalized to obtain the IS - Monte Carlo estimator. The auto-normalized IS is implemented according to the following procedure.

1. Sample N values from the proposal measure: $X^n \sim \mathbb{M}$ for $n = 1, \dots, N$.
2. Compute the associated un-normalized weights: $w(X^n)$ for $n = 1, \dots, N$.
3. Normalize the weights, for $n = 1, \dots, N$:

$$W^n := \frac{w(X^n)}{\sum_{m=1}^N w(X^m)}$$

4. Compute the estimator:

$$\hat{\varphi}_{AN} = \sum_{n=1}^N \varphi(X^n) W^n$$

³Weights are normalized when $\mathbb{E}_{\mathbb{M}}(W) = 1$. In this case W is the Radom-Nikodym derivative between \mathbb{Q} and \mathbb{M} .

While the auto-normalized estimator is biased (since W^n involves a ratio of random variables), it is consistent (see 3.1.1.2 for details).

It is useful to interpret IS as an algorithm that allows to find an approximation for the probability measure \mathbb{Q} as a weighted average of Dirac measures for the sampled values (X^n) : $\delta_{X^n}(dx)$ where the weights are the auto-normalized ones (W^n) . In particular:

$$\mathbb{Q}^N(dx) := \sum_{n=1}^N W^n \delta_{X^n}(dx), \quad X^n \sim \mathbb{M} \quad (3.2)$$

Furthermore, the IS estimator is the expected value of the function φ with respect to the approximating probability measure \mathbb{Q}^N :

$$\begin{aligned} \mathbb{E}_{\mathbb{Q}^N}(\varphi) &= \int_{\mathcal{X}} \varphi(x) \mathbb{Q}^N(dx) = \int_{\mathcal{X}} \varphi(x) \sum_{n=1}^N W^n \delta_{X^n}(dx) = \\ &= \sum_{n=1}^N \int_{\mathcal{X}} \varphi(x) W^n \delta_{X^n}(dx) = \sum_{n=1}^N \varphi(X^n) W^n = \hat{\varphi}_{AN} \end{aligned}$$

3.1.1.2 Convergence of the IS estimator

Consider the auto-normalized IS estimator $\hat{\varphi} = \sum_{n=1}^N \varphi(X^n) W^n$. We can write:

$$\hat{\varphi}_{AN} - \mathbb{E}_{\mathbb{Q}}(\varphi) = \mathbb{E}_{\mathbb{Q}^N}(\varphi) - \mathbb{E}_{\mathbb{Q}}(\varphi) = \frac{N^{-1} \sum_{n=1}^N w(X^n) [\varphi(X^n) - \mathbb{E}_{\mathbb{Q}}(\varphi)]}{N^{-1} \sum_{n=1}^N w(X^n)}$$

Let $\bar{\varphi}(X) := \varphi(X) - \mathbb{E}_{\mathbb{Q}}(\varphi)$. Applying the CLT to the numerator, assuming that $\mathbb{E}_{\mathbb{M}}(w^2 \bar{\varphi}^2) < +\infty$, we get:

$$\sqrt{N} \left\{ \frac{1}{N} \sum_{n=1}^N w(X^n) \bar{\varphi}(X^n) \right\} \Rightarrow \mathcal{N}(0, \mathbb{E}_{\mathbb{M}}(w^2 \bar{\varphi}^2)) \quad (3.3)$$

where \Rightarrow denotes convergence in distribution. Instead, for the denominator the strong LLN implies that:

$$\frac{1}{N} \sum_{n=1}^N w(X^n) \rightarrow \mathbb{E}_{\mathbb{M}}(w) \quad \text{a.s.} \quad (3.4)$$

Applying Slutsky's theorem, together with (3.3) and (3.4), we obtain:

$$\sqrt{N} [\hat{\varphi}_{AN} - \mathbb{E}_{\mathbb{Q}}(\varphi)] \Rightarrow \mathcal{N}\left(0, \frac{\mathbb{E}_{\mathbb{M}}(w^2 \bar{\varphi}^2)}{[\mathbb{E}_{\mathbb{M}}(w)]^2}\right) \quad (3.5)$$

Therefore, the IS (auto-normalized) estimator for φ is both consistent and asymptotically normal.

3.1.1.3 Effective Sample Size

A common measure of efficiency for IS is the effective sample size (ESS), defined as:

$$ESS(W^{1:N}) := \frac{1}{\sum_{n=1}^N (W^n)^2} = \frac{[\sum_{n=1}^N w(X^n)]^2}{\sum_{n=1}^N (w(X^n))^2}$$

Observe that the ESS lies in the interval $[1, N]$. For instance, if $W^n = 1/N$, as in the case of the standard MC method, $ESS = N$.

It is possible to relate the ESS to the variance of the importance function. In fact:

$$\frac{N}{ESS} = 1 + \frac{N^{-1} \sum_{n=1}^N (w(X^n))^2 - [N^{-1} \sum_{n=1}^N w(X^n)]^2}{[N^{-1} \sum_{n=1}^N w(X^n)]^2}$$

The second term in the numerator of the previous formula is the square of the coefficient of variation (i.e., the ratio between the variance and the squared mean) of the un-normalized weights. It is clear that a lower ESS is associated with a higher variance for the weights.

3.1.2 Sequential Importance Sampling

When dealing with state-space models, the importance sampling algorithm has to be applied dynamically and the importance weights are updated period-by-period. To illustrate the relevance of the sequential importance sampling (SIS) algorithm, we consider a two-period stochastic process.

Let us consider the approximation of the measure \mathbb{Q}_0 obtained at time 0 through IS from the measure \mathbb{M}_0 :

$$\mathbb{Q}_0^N(dx_0) = \sum_{n=1}^N W_0^n \delta_{X_0^n}(dx_0), \quad \text{with } X_0^n \sim \mathbb{M}_0, \quad \text{and } W_0^n := \frac{w_0(X_0^n)}{\sum_{m=1}^N w_0(X_0^m)}.$$

Moreover, let $M_1(x_0, dx_1)$ be the transition kernel. We want to update, sequentially, the approximating measure \mathbb{Q}_0 to obtain an approximation of the measure:

$$\mathbb{Q}_1(dx_{0:1}) = \mathbb{Q}_0(dx_0)M_1(x_0, dx_1)$$

SIS works according to the following algorithm.

1. Sample N values for X_0 from the proposal measure \mathbb{M}_0 : $X_0^n \sim \mathbb{M}_0$ for $n = 1, \dots, N$.
2. Sample N values for X_1 from the transition kernel M_1 : $X_1^n \sim M_1(X_0^n, dx_1)$ for $n = 1, \dots, N$.

3. Compute the auto-normalized weights. In this case, they are computed taking into account only time 0 sampled observations since X_1^n is sampled from the correct distribution: $W_1^n = W_0^n$.

Therefore, the particle approximation of $\mathbb{Q}_1(dx_{0:1})$ is:

$$\mathbb{Q}_1^N(dx_{0:1}) = \sum_{n=1}^N W_0^n \delta_{X_0^n}(dx_0) M_1(X_0^n, dx_1)$$

3.1.3 Sequential Importance Resampling

A second approach can be used when sampling has to be done sequentially: sequential importance resampling (SIR). Differently from SIS, an intermediate resampling stage takes place, in which time 0 values are sampled with replacement according to a given resampling strategy (usually, sampling indexes from a multinomial distribution with probabilities equal to the weights computed through IS at time 0).

While in the SIS all the particles (i.e., the sampled values) survive”, with therisk” that particles associated with small weights are carried-over (in the sense that sampling in period 1 occurs from the kernel $M_1(X_0^n, dx_1)$), the same does not occur in SIR. In this second algorithm, particles with larger weights are more likely to be resampled and be used to “generate” the particles in the next period.

We describe the SIR more formally, focusing on the previous two-period case.

1. Sample N values X_0 from the proposal measure \mathbb{M}_0 : $X_0^n \sim \mathbb{M}_0$ for $n = 1, \dots, N$.
2. Compute the weights W_0^n for $n = 1, \dots, N$.
3. Sample N indexes (A_1^1, \dots, A_1^N) from a multinomial distribution such that for all $j = 1, \dots, N$, $Pr(\{A_1^j = n\}) = W_0^n$, for all $n = 1, \dots, N$.
4. Build a new time 0 sample, using the values for X_0 associated to the previously extracted indexes $(A_1^n)_n$: $(X_0^{A_1^1}, \dots, X_0^{A_1^N})$.
5. Sample N values for X_1 from the transition kernel M_1 : $X_1^n \sim M_1(X_0^{A_1^n}, dx_1)$ for $n = 1, \dots, N$.

Steps 3-4 constitute the actual resampling ones. At this point, let us define the new two-period sequence $\tilde{X}_{0:1}^n := (X_0^{A_1^n}, X_1^n)$. Then, the particle approximation for the joint distribution of (X_0, X_1)

becomes:

$$\mathbb{Q}_1^N(dx_{0:1}) = \frac{1}{N} \sum_{n=1}^N \delta_{\tilde{X}_{0:1}^n}(dx_{0:1})$$

Note that, after resampling occurred, each particle $\tilde{X}_{0:1}^n := (X_0^{A_1^n}, X_1^n)$ has the same probability weight.

To understand the reason of the previous statement, let us assume that both \mathbb{Q}_0 and \mathbb{M}_0 admit densities q_0, m_0 with respect to the same dominating measure $\nu(dx_0)$ and that the transition kernel M_1 also admits the density m_1 . Consider the sequential updating for the weights:

$$w_1^n = \frac{\mathbb{Q}_1(\tilde{d}x_{0:1})}{\mathbb{M}_1(dx_{0:1})} \propto \frac{q_0(x_0^{A_1^n})}{Pr(A_1^n) \cdot m_0(x_0^{A_1^n})} \cdot \frac{m_1(x_1^n|x_0^{A_1^n})}{m_1(x_1^n|x_0^{A_1^n})} = \frac{q_0(x_0^{A_1^n})}{\underbrace{m_0(x_0^{A_1^n})}_{W_0^n}} \cdot \frac{1}{\underbrace{Pr(A_1^n)}_{1/W_0^n}} = 1 \quad (3.6)$$

where \mathbb{M}_1 denotes the probability measure for the extraction of the particle $\tilde{x}_{0:1}^n$ according to SIR. \mathbb{M}_1 (second equality sign in the previous formula) can be rewritten as the product between the kernel for x_1^n given $x_0^{A_1^n}$ and the probability of $x_0^{A_1^n}$ itself. This last probability can also be rewritten as the product between the proposal density m_0 from which X_0 is extracted and the probability that it is sampled again when resampling occurs.

3.1.4 SIS vs SIR

The choice between the two sequential sampling strategies (SIS and SIR) has significant implications in terms of computational cost and variance of the importance weights. In fact, it is clear that SIR is more costly computationally since it requires an intermediate resampling step. Nonetheless, SIS suffers from the problem known as curse of dimensionality, i.e., the fact that the *ESS* collapses over time as the variance of the weights diverges.

Consider an extension of SIS to a T -period model and let us assume that the target and the proposal kernels admit densities, respectively, q_t and m_t with respect to the same dominating measures (this implies that the joint distributions⁴ admit a density too). Then, weights, for every $t = 1, \dots, T$, are updated according to:

$$w_t^n \propto \frac{q_t(x_{0:t}^n)}{m_t(x_{0:t}^n)} \propto \underbrace{\frac{q_{t-1}(x_{0:t-1}^n)}{m_{t-1}(x_{0:t-1}^n)} \cdot \frac{q_t(x_t|x_{t-1}^n)}{m_t(x_t|x_{t-1}^n)}}_{\propto w_{t-1}^n} \propto w_{t-1}^n \cdot \frac{q_t(x_t|x_{t-1}^n)}{m_t(x_t|x_{t-1}^n)} \quad (3.7)$$

Equation (3.7) suggests, intuitively, that over time new variance is added from the incremental weights (i.e., the updating factor $q_t(x_t|x_{t-1}^n)/m_t(x_t|x_{t-1}^n)$) determining the curse of dimensionality

⁴Let us denote the joint distributions with q_t and m_t too, with some abuse of notation

issue. Instead, when resampling takes place, since the past weights w_{t-1} are reset to be equal across them all the variance of w_t inherited from the past is shut-down and the only source of variability at each stage comes from the incremental weights.

3.2 Particle Filter

Let us consider the standard state-space model introduced in the previous chapters. Particle filter (PF) is a numerical algorithm that is used to obtain approximations for the filtering - $\mathbb{P}_t(dx_t|y_{1:t})$ - and the prediction - $\mathbb{P}_t(dx_t|y_{1:t-1})$ - distributions for inference in state-space models.

In a state-space model sampling sequentially from the transition kernels for the states is not enough to obtain good approximations since the observation sample $(y_{1:T})$ contains information about the underlying states through the likelihood function $f_t(y_t|x_t)$. Observations constitute signals for the model's states.

Therefore, the conditional distribution of the states, given the observed data - $\mathbb{P}_t(dx_{1:t}|y_{1:t})$ - changes over time as new data points become available. The PF algorithm allows to incorporate sequentially the new data to obtain new approximations for this distribution (and for the filtering and prediction ones).

In general, at each step of the algorithm, states are sampled sequentially from the proposal measure, call it \mathbb{M}_t , and importance auto-normalized weights are computed as the Radom-Nikodym derivative of the target \mathbb{P}_t with respect to the proposal \mathbb{M}_t . Moreover, the Markov properties of the sequence of the states allows to update sequentially the importance weights.

We will analyse three different types of particle filtering procedures: bootstrap, guided and auxiliary particle filters, that differ between them with respect to their scope and efficiency.

3.2.1 Degeneracy of the SIS in Particle Filter

Usually, PF requires a resampling step at each stage of the algorithm. Recall that, in general, resampling determines a trade-off between the increase in the variance at the resampling stage and a decrease in the variance when future states are sampled from the transition kernel.

Dead" particles (i.e., those with small weights) are likely to be removed when resampling takes place and past weights getreset" (usually, they are set to be equal to $1/N$). Otherwise, the curse of dimensionality problem will emerge and the variance of the weights will diverge over time as the

variance of past weights increases as new updates occur period-by-period. Figure ***** shows an example of a PF where SIS is used: note that the ESS drops dramatically after few periods, meaning that the variance of the importance weights diverges. This motivates the use of a SIR structure for the PF algorithm.

Lastly, it is now common practice to use a mixed approach with respect to resampling (recall that resampling has an additional computational cost). Generally, resampling is done adaptively, when some condition is met. A standard criterion that is used (and that we adopt in our implementations of the PF) is to resample when the ESS drops below some threshold $ESS_{min} \in [1, N]$:

$$ESS(W_{t-1}^{1:N}) < ESS_{min} \quad (3.8)$$

3.2.2 Bootstrap Particle Filter

The bootstrap particle filter (BPF) is the easiest form of an adaptive PF algorithm. In fact, future states are sampled from the correct transition kernel and the weights are updated sequentially using the likelihood $f_t(y_t|x_t)$.

We summarize the procedure with the following scheme.

1. Initial stage: For $n = 1, \dots, N$, at time 0 sample: $X_0^n \sim \mathbb{P}_0(dx_0)$ and set the un-normalized weights $w_0^n = 1$. Then, the normalized weights are: $W_0^n = \frac{1}{N}$.
2. For every $t = 1, \dots, T$:
 - (a) Resample if $ESS(W_{t-1}^{1:N}) < ESS_{min}$.⁵
 - i. Multinomial sampling for the indexes $(A_t^n)_n$ with probabilities given by $(W_{t-1}^n)_n$.
 - ii. Update weights as: $\hat{w}_{t-1}^n = 1$ for $n = 1, \dots, N$.

Otherwise,

⁵It is indifferent to implement this resampling step at the beginning or at the end of each iteration (as we do in our codes) as long as the initial stage $t = 0$ generates weights that are constant and equal to $1/N$. In fact, in this case, resampling will never take place at the beginning of the first iteration since $ESS(W_0^n = 1/N) = N$ that is larger than or equal to ESS_{min} . However, when we describe the algorithm, we chose to leave this resampling step at the beginning of each iteration. In fact, when considering the guided or auxiliary PF, we can, more generally, sample the initial values (time $t = 0$) for the state from a proposal distribution different from the target one. In this case, the weights at $t = 0$ might differ from $1/N$, hence resampling might take place already at the beginning of the first iteration.

- i. Set the indexes: $A_t^n = n$ for $n = 1, \dots, N$.
 - ii. Set the weights: $\hat{w}_{t-1}^n = w_{t-1}^n$ for $n = 1, \dots, N$.
- (b) Sample the future states from the transition kernel: $X_t^n \sim P_t(X_{t-1}^{A_t^n}, dx_t)$ for $n = 1, \dots, N$.
- (c) Update the un-normalized weights through the likelihood function for $n = 1, \dots, N$:

$$w_t^n = \hat{w}_{t-1}^n \cdot f_t(y_t | X_t^n).$$

- (d) Normalize the weights for $n = 1, \dots, N$: $W_t^n := \frac{w_t^n}{\sum_{m=1}^N w_t^m}$.

Therefore, at the initial step the pre-sample period states x_0^n are extracted and every sampled state is assigned an equal weight since there is no corresponding observation y_0 that gives additional information on the likelihood of that extraction.

Resampling does never take place in $t = 1^6$ and the state in $t = 1$ is sampled from the correct transition kernel $P_1(x_0^n, dx_1)$. Now, the time 0 weights are updated incorporating the information from the data point y_1 in a very intuitive way: more likely states (given y_1) are given higher weights: $w_1^n = f_1(y_1 | x_1^n)$. These weights are used to approximate the time 1 filtering distribution.

At a subsequent stage, the *ESS* criterion should be checked and, in case it is satisfied, resampling takes place and time 1 weights are reset to be equal to $1/N$.

The key characteristic of the PF is the fact that the importance weights are updated sequentially, according to the “information” provided by the observed data (i.e., the sequence $y_{1:T}$). Recall that time t weight is the Radon-Nikodym derivative of the target with respect to the proposal distribution. In the BPF, the time t target is the measure $\mathbb{P}_t(dx_{0:t} | y_{1:t})$, while the proposal, denoted generically by \mathbb{M}_t , is $\mathbb{M}_t(dx_{0:t} | y_{1:t}) = \mathbb{M}_{t-1}(dx_{0:t-1} | y_{1:t-1}) \cdot P_t(x_{t-1}, dx_t)$. Hence:

$$\begin{aligned} w_t &\propto \frac{\mathbb{P}_t(dx_{0:t} | y_{1:t})}{\mathbb{M}_t(dx_{0:t} | y_{1:t})} \propto \frac{\mathbb{P}_{t-1}(dx_{0:t-1} | y_{1:t-1}) \cdot P_t(x_{t-1}, dx_t) \cdot f(y_t | x_t)}{\mathbb{M}_{t-1}(dx_{0:t-1} | y_{1:t-1}) \cdot P_t(x_{t-1}, dx_t)} = \\ &= \frac{\mathbb{P}_{t-1}(dx_{0:t-1} | y_{1:t-1})}{\underbrace{\mathbb{M}_{t-1}(dx_{0:t-1} | y_{1:t-1})}_{\hat{w}_{t-1}}} \cdot f(y_t | x_t) = \hat{w}_{t-1} \cdot f(y_t | x_t) \end{aligned} \quad (3.9)$$

Finally, the BPF is able to generate the two key objects for inference in state-space models:

⁶Since $ESS(W_0^n = 1/N) = N \geq ESS_{min}$.

- the predictive distribution:

$$\begin{aligned}\mathbb{P}_t^N(dx_t|y_{1:t-1}) &= \frac{1}{\sum_{n=1}^N \hat{w}_t^n} \sum_{n=1}^N \hat{w}_t^n \delta_{X_t^n}(dx_t) = \\ &= \begin{cases} \frac{1}{N} \sum_{n=1}^N \delta_{X_t^n}(dx_t) & \text{if resampling occurs} \\ \frac{1}{\sum_{n=1}^N w_{t-1}^n} \sum_{n=1}^N w_{t-1}^n \delta_{X_t^n}(dx_t) & \text{otherwise} \end{cases},\end{aligned}\quad (3.10)$$

- the filtering distribution:

$$\mathbb{P}_t^N(dx_t|y_{1:t}) = \sum_{n=1}^N W_t^n \delta_{X_t^n}(dx_t) \quad (3.11)$$

3.2.3 Guided Particle Filter

It is not always possible to sample sequentially from the correct transition kernel. Moreover, even if it were possible, there exist a proposal transition kernel that makes the algorithm more efficient. This motivates the guided particle filter (GPF), that can be interpreted as a more general version of the GPF. In fact, with the GPF it is possible, at each stage, to sample new states from a transition kernel that differs from the correct one.

Let $P_t(x_{t-1}, dx_t)$ be the correct kernel and consider a kernel $M_t(x_{t-1}, dx_t)$ defined on the same measurable space of P_t . Moreover, let us assume that that P_t is absolutely continuous with respect to M_t ⁷ (i.e., $P_t \ll M_t$). If sampling occurs sequentially from this kernel, the way in which weights are updated should differ from the BPF. Consider equation (3.9) and let us reformulate it for the GPF case. Here, the proposal distribution is not $\mathbb{M}_t(dx_{0:t}|y_{1:t}) = \mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1}) \cdot P_t(x_{t-1}, dx_t)$ anymore, but becomes $\mathbb{M}_t(dx_{0:t}|y_{1:t}) = \mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1}) \cdot M_t(x_{t-1}, dx_t)$. Hence:

$$\begin{aligned}w_t &\propto \frac{\mathbb{P}_t(dx_{0:t}|y_{1:t})}{\mathbb{M}_t(dx_{0:t}|y_{1:t})} \propto \frac{\mathbb{P}_{t-1}(dx_{0:t-1}|y_{1:t-1}) \cdot P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{\mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1}) \cdot M_t(x_{t-1}, dx_t)} = \\ &= \underbrace{\frac{\mathbb{P}_{t-1}(dx_{0:t-1}|y_{1:t-1})}{\mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1})}}_{\hat{w}_{t-1}} \cdot \underbrace{\frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t(x_{t-1}, dx_t)}}_{\text{updating factor}} = \hat{w}_{t-1} \cdot \underbrace{\frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t(x_{t-1}, dx_t)}}_{\text{updating factor}}\end{aligned}\quad (3.12)$$

Note that the updating factor is more complex than the one of the BPF: in addition to the likelihood function $f_t(y_t|x_t)$, the fact that sampling occurs from a kernel different from the correct ones requires the additional term $P_t(x_{t-1}, dx_t)/M_t(x_{t-1}, dx_t)$.

We are now ready to establish, informally, the optimality result about the transition kernel. Let us

⁷This condition is stronger than what would be needed for the existence of the importance weights in GPF, i.e., that $P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t) \ll M_t(x_{t-1}, dx_t)$.

denote by $G_t = [P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)dy_t]/M_t(x_{t-1}, dx_t)$ the incremental weights (i.e., the updating factor) and by $M_t^*(x_{t-1}, dx_t)$ the optimal kernel. It can be proved⁸ that the optimal kernel, i.e., the one that minimizes the variance of the incremental weights G_t (and of the time t weights themselves), is:

$$M_t^*(x_{t-1}, dx_t) = \frac{1}{\int_{\mathcal{X}} P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)} P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t) \quad (3.13)$$

Note that M_t^* is the conditional distribution of x_t given x_{t-1} and y_t .⁹ The intuition for this result is apparent when considering equation (3.12): if $M_t = M_t^*$, the updating factor collapses to $\int_{\mathcal{X}} P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)$ and information from the observed data are incorporated efficiently. As a matter of fact, it is usually impossible to sample from the optimal kernel and, in practice, it is replaced by its linear Gaussian approximation.

The GPF can be resumed by the following steps.

1. Initial stage: For $n = 1, \dots, N$, at time 0 sample: $X_0^n \sim \mathbb{M}_0(dx_0)$ and set the un-normalized weights:

$$w_0^n = \frac{\mathbb{P}_0(dx_0)}{\mathbb{M}_0(dx_0)}.$$

Then, the normalized weights are: $W_0^n := \frac{w_0^n}{\sum_{m=1}^N w_0^m}$.

2. For every $t = 1, \dots, T$:

- (a) Resample if $ESS(W_{t-1}^{1:N}) < ESS_{min}$:

- i. Multinomial sampling for the indexes $(A_t^n)_n$ with probabilities given by $(W_{t-1}^n)_n$.
- ii. Update weights as: $\hat{w}_{t-1}^n = 1$ for $n = 1, \dots, N$.

Otherwise,

- i. Set the indexes: $A_t^n = n$ for $n = 1, \dots, N$
- ii. Set the weights: $\hat{w}_{t-1}^n = w_{t-1}^n$ for $n = 1, \dots, N$.

- (b) Sample the future states from the transition kernel: $X_t^n \sim M_t(X_{t-1}^{A_t^n}, dx_t)$ for $n = 1, \dots, N$.

- (c) Update the un-normalized weights, for $n = 1, \dots, N$:

$$w_t^n = \hat{w}_{t-1}^n \cdot \frac{P_t(X_{t-1}^{A_t^n}, dx_t) \cdot f_t(y_t|X_t^n)}{M_t(X_{t-1}^{A_t^n}, dx_t)}.$$

⁸For a formal proof see [?], th. 10.1 p. 142.

⁹In fact, assuming that the P_t admits a density p_t , applying Bayes rule, we have that $m_t(x_t|x_{t-1}, y_t) \propto p_t(x_t|x_{t-1}) \cdot f(y_t|x_t)$.

- (d) Normalize the weights, for $n = 1, \dots, N$: $W_t^n := \frac{w_t^n}{\sum_{m=1}^N w_t^m}$.

Note that $\mathbb{M}_0(dx_0)$ is the initial sampling distribution, that could differ from the correct one $\mathbb{P}_0(dx_0)$. To conclude, one of the issues related to the BPF is the fact that new states in t are sampled ignoring their likelihood given the future observation y_t . Hence, extractions from the “correct” transition kernel may be associated to low weights (small value for the likelihood). Instead, the GPF allows to solve for this problem, defining a proposal kernel that depends on the likelihood itself. Therefore, the extraction of future states already incorporates the information given by the data resulting in a more efficient procedure.

3.2.4 Auxiliary Particle Filter

The auxiliary particle filter (APF) constitutes a further extension of the BPF. The GPF allowed to sample from a transition kernel different from the correct ones. The APF, instead, allows to extract the ancestor variables from an arbitrary distribution at the resampling step.

Let $\eta_t : \mathcal{X} \rightarrow \mathbb{R}_+$ be a non-negative, real-valued positive function, called auxiliary function. At time t , when the resampling stage occurs, the multinomial distribution for the indexes to be extracted can use probabilities that differ from the weights in $t-1$, W_{t-1}^n . The new weights, called auxiliary weights, are computed as a transformation of the original ones, denoted as inferential weights, through the auxiliary function:

$$\tilde{W}_t^n = \frac{W_t^n \cdot \eta_t(X_t^n)}{\sum_{m=1}^N W_t^m \cdot \eta_t(X_t^m)} \quad (3.14)$$

Note that the computation of the auxiliary weights is done from the auto-normalized inferential ones. Additionally, note that imposing $\eta_t(x_t) = 1$ we recover the GPF and the BPF as special cases of the APF.

Clearly, inferential weights are updated by incorporating the fact that resampling occurs according to different weights (the auxiliary ones):

$$w_t^n = \frac{W_{t-1}^{A_t^n}}{\tilde{W}_{t-1}^{A_t^n}} G_t \quad (3.15)$$

where G_t is the incremental weight defined in the section on GPF.

To understand the logic of this sequential weight updating, consider (3.12). In particular, suppose

that multinomial resampling takes place at every stage and let $\tilde{w}_{t-1} := w_{t-1} \cdot \eta_{t-1}(x_{t-1})$. Then:

$$\begin{aligned}
 w_t &\propto \frac{\mathbb{P}_t(dx_{0:t}|y_{1:t})}{\mathbb{M}_t(dx_{0:t}|y_{1:t})} \propto \frac{\mathbb{P}_{t-1}(dx_{0:t-1}|y_{0:t-1})}{\mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1}) \cdot \tilde{w}_{t-1}} \cdot \frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t(x_{t-1}, dx_t)} \\
 &= \frac{\mathbb{P}_{t-1}(dx_{0:t-1}|y_{0:t-1})}{\underbrace{\mathbb{M}_{t-1}(dx_{0:t-1}|y_{1:t-1})}_{w_{t-1}}} \cdot \frac{1}{\tilde{w}_{t-1}} \cdot \frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t(x_{t-1}, dx_t)} \\
 &= \frac{w_{t-1}}{\tilde{w}_{t-1}} \cdot \frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t(x_{t-1}, dx_t)}
 \end{aligned} \tag{3.16}$$

It is possible to derive an optimality result for the auxiliary functions¹⁰. In the special case in which the transition kernel used for sampling corresponds with the optimal one M_t^* , the optimal auxiliary function (i.e., the auxiliary function that minimizes the variance of the inferential weights) becomes:

$$\eta_{t-1}^*(x_{t-1}) = \int_{\mathcal{X}} P_t(x_{t-1}, dx_t) f_t(y_t|x_t) \tag{3.17}$$

and is called perfectly adapted auxiliary function. In fact, if $M_t = M_t^*$, then, expression (3.16) becomes:

$$w_t \propto \frac{w_{t-1}}{\tilde{w}_{t-1}} \cdot \frac{P_t(x_{t-1}, dx_t) \cdot f_t(y_t|x_t)}{M_t^*(x_{t-1}, dx_t)} = \frac{1}{\eta_{t-1}(x_{t-1})} \cdot \int_{\mathcal{X}} P_t(x_{t-1}, dx_t) f_t(y_t|x_t)$$

From the previous equation, it is clear that if $\eta_t = \eta_t^*$, the inferential weights are constant and their variance is zero.

We can summarize the APF algorithm.

1. Initial stage:

- (a) For $n = 1, \dots, N$, at time 0 sample: $X_0^n \sim \mathbb{M}_0(dx_0)$ and set the un-normalized inferential weights:

$$w_0^n = \frac{\mathbb{P}_0(dx_0)}{\mathbb{M}_0(dx_0)}.$$

Then, the normalized inferential weights are: $W_0^n := \frac{w_0^n}{\sum_{m=1}^N w_0^m}$.

- (b) Compute the un-normalized auxiliary weights:

$$\tilde{w}_0^n = w_0^n \cdot \eta_0(X_0^n)$$

Then the normalized auxiliary weights are: $\tilde{W}_0^n := \frac{\tilde{w}_0^n}{\sum_{m=1}^N \tilde{w}_0^m}$.

2. For every $t = 1, \dots, T$:

¹⁰Details in ***** th 10.2, p. 148

(a) Resample if $ESS(W_{t-1}^{1:N}) < ESS_{min}$:

- i. Multinomial sampling for the indexes $(A_t^n)_n$ with probabilities given by $(\tilde{W}_{t-1}^n)_n$.
- ii. Update weights as: $\hat{w}_{t-1}^n = \frac{W_{t-1}^{A_t^n}}{\tilde{W}_{t-1}^{A_t^n}}$ for $n = 1, \dots, N$.

Otherwise,

- i. Set the indexes: $A_t^n = n$ for $n = 1, \dots, N$.
- ii. Set the weights: $\hat{w}_{t-1}^n = w_{t-1}^n$ for $n = 1, \dots, N$.

(b) Sample the future states from the transition kernel: $X_t^n \sim M_t(X_{t-1}^{A_t^n}, dx_t)$ for $n = 1, \dots, N$.

(c) Update the un-normalized inferential weights, for $n = 1, \dots, N$:

$$w_t^n = \hat{w}_{t-1}^n \cdot \frac{P_t(X_{t-1}^{A_t^n}, dx_t) \cdot f_t(y_t | X_t^n)}{M_t(X_{t-1}^{A_t^n}, dx_t)}.$$

(d) Normalize the inferential weights, for $n = 1, \dots, N$: $W_t^n := \frac{w_t^n}{\sum_{m=1}^N w_t^m}$.

(e) Compute the un-normalized auxiliary weights, for $n = 1, \dots, N$:

$$\tilde{w}_t^n = w_t^n \cdot \eta_t(X_t^n)$$

(f) Normalize the auxiliary weights, for $n = 1, \dots, N$: $\tilde{W}_t^n := \frac{\tilde{w}_t^n}{\sum_{m=1}^N \tilde{w}_t^m}$.¹¹

However, it is worth noting that the optimality results both for the kernel in the GPF and of the auxiliary in the APF, has two main weaknesses:

- it is not clear that what is optimal at time t remains such at future stages. In other words, it is not truly clear how the minimization of the variance of the inferential weights at each stage affect the efficiency of the PF estimators;
- both the choice of M_t and of η_t affect the way in which weights are computed and, indirectly, the decision to resample in the future. It is difficult to take also this indirect effect into account when evaluating the overall efficiency of the algorithm.

¹¹It is worth mentioning that the original APF included a final resampling step at the end of each iteration. However, our formulation is more efficient since having two resampling steps might introduce additional Monte Carlo variance.
Se ***Johansen p. 120

3.3 Liu and West Filter

Let us consider a more general state-space model, where the state vector includes a time-constant parameter vector $\psi \in \Psi$, where Ψ is the parameter space. Since we can interpret the parameter as a state with the law of motion $\psi_t = \psi_{t-1} = \psi$, then it is possible to adopt a PF algorithm to estimate ψ .

However, being ψ constant over time, it is meaningless to sample its values sequentially: the first extraction of ψ at the initial stage 0 (i.e., ψ_0^n) generates a constant path of particles for the parameter ($\psi_t^n = \psi_0^n$).¹²

Therefore, Liu and West (2001) proposed a modified particle filter, called Liu and West filter (LWF) that allows to resample the parameter over time from a continuous distribution. In this way, at every time, the support of ψ is not limited to the initially sampled N values.

Here, we describe the LWF that makes use of the normal distribution to sample sequentially the parameter. Moreover, consider the simpler version of the bootstrap filter.¹³

Let the transition kernel and the likelihood depend on the parameter ψ (i.e., $P_t(x_{t-1}, dx_t; \psi)$ and $f_t(y_t|x_t; \psi)$) and let $\pi(\cdot)$ be the prior distribution for the parameter ψ .

The LWF can be described by the following algorithm.

1. Initial stage: For $n = 1, \dots, N$, at time 0 sample: $X_0^n \sim \mathbb{P}_0(dx_0)$ and $\psi_0^n \sim \pi(\psi_0)$.

Then, set the un-normalized weights $w_0^n = 1$.

Hence, the normalized weights are: $W_0^n = \frac{1}{N}$.

2. For every $t = 1, \dots, T$:

(a) Resample if $ESS(W_{t-1}^{1:N}) < ESS_{min}$:

- i. Multinomial sampling for the indexes $(A_t^n)_n$ with probabilities given by $(W_{t-1}^n)_n$.
- ii. Update weights as: $\hat{w}_{t-1}^n = 1$ for $n = 1, \dots, N$.

Otherwise,

- i. Set the indexes: $A_t^n = n$ for $n = 1, \dots, N$
- ii. Set the weights: $\hat{w}_{t-1}^n = w_{t-1}^n$ for $n = 1, \dots, N$.

¹²To be more precise, at the resampling step, also ψ^n is resampled, but the possible values that it may take are the N ones obtained at the initial time 0 extraction.

¹³Clearly, the LWF can be also used within the more general frameworks of the GPF and the APF.

- (b) Extracting new parameter values:
 - i. Compute the weighted average and the variance for the parameters, respectively:

$$\bar{\psi} = \sum_{n=1}^N W_{t-1}^n \psi^n, \Omega = \sum_{n=1}^N W_{t-1}^n (\psi^n - \bar{\psi})^2.$$
 - ii. Compute: $m^n = a \cdot \psi^n + (1 - a) \cdot \bar{\psi}$ for $n = 1, \dots, N$.
 - iii. Draw new parameters: $\psi^n \sim \mathcal{N}(m^{A_t^n}, h^2 \cdot \Omega)$ for $n = 1, \dots, N$.
- (c) Sample future states from the transition kernel: $X_t^n \sim P_t(X_{t-1}^{A_t^n}, dx_t; \psi^n)$ for $n = 1, \dots, N$.
- (d) Update the un-normalized weights through the likelihood function, for $n = 1, \dots, N$:

$$w_t^n = \hat{w}_{t-1}^n \cdot f_t(y_t | X_t^n; \psi^n).$$

- (e) Normalize the weights, for $n = 1, \dots, N$: $W_t^n := \frac{w_t^n}{\sum_{m=1}^N w_t^m}$.

At the initial step, N states and N parameter values are sampled independently, respectively, from \mathbb{P}_0 and π and the particle weights are set to be constant, as usual.

Then, at each iteration (time t), the decision to resample is done according to the *ESS* criterion and indexes are extracted from a multinomial distribution. However, in the LWF both previous periods states and parameters are resampled.

Meanwhile, to extract the next period states and parameters it is necessary, firstly, to sample the new N parameter values. Let $\bar{\psi}$ be the weighted average of the parameter values and Ω be their variance (the weights used to compute this mean and variance are the previous period inferential ones, i.e. W_{t-1}^n).

The algorithm suggests to sample new values for the parameters from a continuous distribution, i.e., a Normal centered at the (resampled) values $m^{A_t^n}$ with variance $h^2 \cdot \Omega$ where m^n is a convex combination of the sample mean $\bar{\psi}$ and the actual extraction of ψ^n with coefficient $a \in (0, 1)$ and h is a scalar such that $a^2 + h^2 = 1$. While it would be more intuitively immediate to sample new parameters from a Normal density centered at the previously sampled values (i.e., ψ^n), this would increase the variance of ψ over time.

In fact, let $\pi_{t-1}^N(\psi) = \sum_{n=1}^N \hat{W}_{t-1}^{A_t^n} \cdot \mathcal{N}(\psi^{A_t^n}, \Omega)^{14}$ be the empirical distribution of the parameters.

¹⁴Big hat-weights $\hat{W}_{t-1}^{A_t^n}$ are the normalized weights after the resampling stage, equal to $1/N$ if resampling takes place, W_{t-1}^n otherwise.

From the law of iterated expectation and total variance, we obtain:¹⁵

$$\begin{aligned}\mathbb{E}_{\pi_{t-1}^N}(\psi) &= \mathbb{E}[\mathbb{E}(\psi|A_t^n)] = \mathbb{E}[\psi^{A_t^n}] = \bar{\psi} \\ \text{Var}_{\pi_{t-1}^N}(\psi) &= \text{Var}[\mathbb{E}(\psi|A_t^n)] + \mathbb{E}[\text{Var}(\psi|A_t^n)] = \text{Var}[\psi^{A_t^n}] + \mathbb{E}[\Omega] = 2 \cdot \Omega\end{aligned}$$

Hence, while the mean is preserved, the variance increases at each iteration. Instead, sampling from the distribution proposed by Liu and West $\mathcal{N}(m^{A_t^n}, h^2 \cdot \Omega)$ preserves both the mean and the variance of the parameters over time.

Let $\tilde{\pi}_{t-1}^N(\psi) = \sum_{n=1}^N \hat{W}_{t-1}^{A_t^n} \cdot \mathcal{N}(m^{A_t^n}, h^2 \cdot \Omega)$. Hence:

$$\begin{aligned}\mathbb{E}_{\tilde{\pi}_{t-1}^N}(\psi) &= \mathbb{E}[\mathbb{E}(\psi|A_t^n)] = \mathbb{E}[m^{A_t^n}] = \mathbb{E}[a \cdot \psi^{A_t^n} + (1-a) \cdot \bar{\psi}] = \bar{\psi} \\ \text{Var}_{\tilde{\pi}_{t-1}^N}(\psi) &= \text{Var}[\mathbb{E}(\psi|A_t^n)] + \mathbb{E}[\text{Var}(\psi|A_t^n)] = \text{Var}[m^{A_t^n}] + h^2 \cdot \mathbb{E}[\Omega] = \underbrace{(a^2 + h^2)}_{=1} \cdot \Omega = \Omega\end{aligned}$$

Finally, after sampling the new N parameter values, conditional on these values, new states are also sampled from the transition kernel: $X_t^n \sim P_t(X_{t-1}^{A_t^n}, dx_t; \psi^n)$. Lastly, weights are updated through the incremental weights.

To conclude, it is worth noting that the LWF produces, together with the usual filtering and prediction distributions, a posterior distribution for the parameter:

$$\pi_t^N(\psi) = \pi^N(\psi|y_{1:t}) = \sum_{n=1}^N W_t^n \delta_{\psi^n}(\psi) \quad (3.18)$$

3.4 Convergence

¹⁵For completeness, we specify that the first expected value in the law of iterated expectations and total variances is with respect to the distribution of the indexes and the second one with respect to each of the Normal distributions centered at $\psi^{A_t^n}$.

Chapter 4

Applications to Stochastic Volatility Models

4.1 Comparing Constant and Stochastic Volatility Models of Equity Returns

In this section, we apply the particle filter in a model-comparison exercise, ultimately underlining the importance of taking stochastic volatility into account in modeling equity returns.

We consider a stochastic volatility model

$$\mathcal{M}_{SV} : \begin{cases} r_t = \alpha_r + \sigma_{r,t} \cdot v_t \\ \log \sigma_{r,t}^2 = \alpha_\sigma + \beta_\sigma \cdot \log \sigma_{r,t-1}^2 + \sigma_\sigma \cdot w_t \end{cases} \quad (v_t, w_t) \sim \mathcal{N}_2(0, I_2)$$

and a constant volatility iid model

$$\mathcal{M}_{CV} : (r_t) \stackrel{iid}{\sim} \mathcal{N}(\alpha, \sigma^2).$$

We calibrate the model parameters using the whole sample, and proceed by setting the return mean of both models to equal the sample mean, σ^2 equal to the return variance, and the parameters of the volatility equation in \mathcal{M}_{SV} by running an $AR(1)$ regression of the log squared residuals.

In order to compare the predictive performance of the two, we compare the cumulative log likelihood of the observed returns up to time t given the considered model. In formulas, the cumulative likelihood can be computed as

$$p(y_{1:t}|\mathcal{M}_i) = \prod_{s=0}^{t-1} p(y_{s+1}|y_{1:s}, \mathcal{M}_i), i \in \{SV, CV\},$$

adopting the convention $y_{1:0} = \emptyset$. This is trivial to obtain in the iid case; considering the stochastic volatility model, the predictive likelihoods can be derived as

$$\begin{aligned} p(y_{t+1}|y_{1:t}, \mathcal{M}_{SV}) &= \int p(y_{t+1}|\sigma_{r,t+1}, y_{1:t}, \mathcal{M}_{SV})p(\sigma_{r,t+1}|y_{1:t}, \mathcal{M}_{SV})d(\sigma_{r,t+1}) \\ &= \int p(y_{t+1}|\sigma_{r,t+1}, \mathcal{M}_{SV})p(\sigma_{r,t+1}|y_{1:t}, \mathcal{M}_{SV})d(\sigma_{r,t+1}) \end{aligned}$$

Whereas the emission distribution $p(y_{t+1}|\sigma_{r,t+1}, \mathcal{M}_{SV})$ is known, we approximate the predictive state distribution $p(\sigma_{r,t+1}|y_{1:t}, \mathcal{M}_{SV})$ with weighted particles obtained with particle filtering.

Figure XX reports the time series of the cumulative log-likelihood given the constant volatility model \mathcal{M}_{CV} relative to the one given \mathcal{M}_{SV} , obtained as $\log p(y_{1:t}|\mathcal{M}_{CV}) - \log p(y_{1:t}|\mathcal{M}_{SV})$. Notice that the first difference of this series is the relative log-likelihood of the incoming observation, $p(y_t|y_{1:t-1}, \mathcal{M}_{CV}) - p(y_t|y_{1:t-1}, \mathcal{M}_{SV})$.

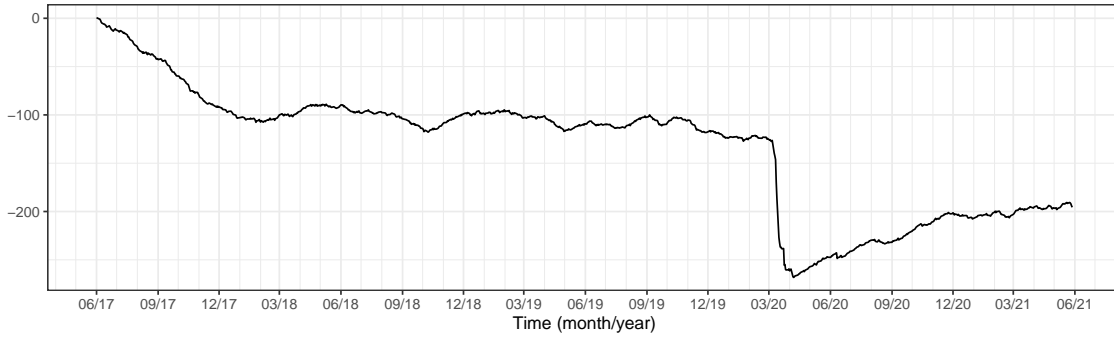


Figure 4.1: Relative log-likelihood

Overall, this analysis shows how a model with stochastic volatility is better at describing and predicting returns. As can be seen from the graph and perhaps unsurprisingly, the density forecast given the stochastic volatility model represents a particularly significant improvement in the period of extreme values of returns in March 2020.