



**HAL**  
open science

# The Gauss-Galerkin approximation method in nonlinear filtering

Fabien Campillo

► **To cite this version:**

Fabien Campillo. The Gauss-Galerkin approximation method in nonlinear filtering. 2023. hal-03985941

**HAL Id: hal-03985941**

**<https://inria.hal.science/hal-03985941v1>**

Submitted on 13 Feb 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# The Gauss-Galerkin approximation method in nonlinear filtering

Fabien Campillo\*

January 22, 2023

## Abstract

We study an approximation method for the one-dimensional nonlinear filtering problem, with discrete time and continuous time observation. We first present the method applied to the Fokker-Planck equation. The convergence of the approximation is established. We finally present a numerical example.

**Keywords:** nonlinear filtering, moment method, particle approximation.

*This is the English translation of the paper: “Fabien Campillo. La méthode d’approximation de Gauss-Galerkin en filtrage non linéaire. *RAIRO M2AN*, 20(2):203–223, 1986” with some supplementary material, see Addendum page 23.*

## 1 Introduction

Usual methods for numerical solutions of partial differential equations typically involve a large number of space discretization points. Moreover, in their classical form, these methods use time-fixed discretization grids.

The method proposed by Donald A. Dawson [4] for the numerical solution of the Fokker-Planck equation, called the Gauss-Galerkin method, combines the Gauss quadrature and Galerkin approximation methods. This method, which can be considered as a particle method [10], has the double advantage of giving acceptable results even with a small number of unknown variables to calculate and a discretization grid able to adapt to the evolution of the solution of the partial differential equation considered. However, in its current form, the method is limited to the case of a single dimension of space.

We will study the behavior of this method, applied to the nonlinear filtering problem. In Section 2, we present the Gauss-Galerkin approximation method applied to the Fokker-Planck equation, and we establish a convergence result. The results in this section are a reworking and development of the work of Donald A. Dawson [4].

In Section 3, we first consider the nonlinear filtering problem with discrete time observation: we present the Gauss-Galerkin approximation and prove its convergence. We then

---

\*Fabien.Campillo@inria.fr – [website](#)– MathNeuro, Inria Montpellier, France

consider the nonlinear filtering problem with continuous time observation: In this case, before introducing the approximation, we go back to the previous case by discretizing the observation equation.

We define the following spaces:

$\mathcal{M}_1(\mathbb{R})$	probability measures on $\mathbb{R}$ ,
$\mathcal{M}_+(\mathbb{R})$	non-negative measures on $\mathbb{R}$ ,
$\mathcal{M}(\mathbb{R})$	signed measures on $\mathbb{R}$ ,
$\mathcal{C}_u(\mathbb{R})$	bounded and uniformly continuous functions $\mathbb{R} \rightarrow \mathbb{R}$ ,
$\mathcal{C}_c^\infty(\mathbb{R})$	continuous functions $\mathbb{R} \rightarrow \mathbb{R}$ of class $\mathcal{C}^\infty$ with compact support,
$\mathcal{C}_b^2(\mathbb{R})$	continuous and bounded functions $\mathbb{R} \rightarrow \mathbb{R}$ of class $\mathcal{C}^2$ ,
$\mathcal{C}[0, T]$	continuous functions $[0, T] \rightarrow \mathbb{R}$ ,
$\mathcal{P}_{2N-1}$	polynomial functions of degree at most $2N - 1$ .

## 2 Numerical solution of the Fokker-Planck equation

### 2.1 The Gauss-Galerkin approximation method

To introduce the Fokker-Planck equation, we consider the stochastic differential equation:

$$(1) \quad dX_t = b(X_t) dt + \sigma(X_t) dW_t, \quad 0 \leq t \leq T, X_0 \sim \mu_0,$$

where  $(X_t)_{t \leq T}$  takes values in  $\mathbb{R}$ ;  $(W_t)_{t \leq T}$  is a real standard Wiener process independent of  $X_0$ . Let  $a(x) = \sigma^2(x)$ ,  $a'(x) = da(x)/dx$ ,  $b'(x) = db(x)/dx$ , we make the following assumptions:

**(H1)**  $b, \sigma : \mathbb{R} \rightarrow \mathbb{R}$ , are measurable and bounded applications;

**(H2)**  $a' \in L^\infty(\mathbb{R})$  and there exists  $\underline{a} > 0$  such that  $a(x) \geq \underline{a}$ , for all  $x \in \mathbb{R}$ ;

**(H3)**  $b'$  is measurable bounded, and  $a'$  is continuous.

Under Assumptions **(H1)**-**(H2)**, Equation (1) admits a unique solution in the weak sense [12]. Hypothesis **(H3)** will be used in the following to demonstrate the convergence of the approximation.

Let  $\mu_t \in \mathcal{M}_1(\mathbb{R})$  be the distribution law of  $X_t$  on  $\mathbb{R}$ , for all  $0 \leq t \leq T$ :

$$\langle \mu_t, \varphi \rangle = \mathbb{E}(\varphi(X_t)), \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}),$$

where:

$$\langle \mu_t, \varphi \rangle := \int_{\mathbb{R}} \varphi(x) \mu_t(dx).$$

It results from the Itô's formula that  $(\mu_t)_{t \leq T}$  is a solution of the Fokker-Planck equation (written in weak form) :

$$(2) \quad \langle \mu_t, \varphi \rangle = \langle \mu_0, \varphi \rangle + \int_0^t \langle \mu_s, \mathcal{L}\varphi \rangle ds, \quad 0 \leq t \leq T, \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}),$$

where  $\mathcal{L}$  denotes the the infinitesimal generator of the Markov process  $X_t$ :

$$\mathcal{L}\varphi(x) := b(x) \varphi'(x) + \frac{1}{2} a(x) \varphi''(x).$$

For  $N \in \mathbb{N}$  given, the Gauss-Galerkin approximation method consists in approximating  $(\mu_t)_{t \leq T}$  by a family of probability measures  $(\mu_t^N)_{t \leq T}$  of the form :

$$\mu_t^N(\mathrm{d}x) = \sum_{i=1}^N w_t^{(i)} \delta_{x_t^{(i)}}(\mathrm{d}x) \in \mathcal{M}_1(\mathbb{R}).$$

The functions  $t \rightarrow w_t^{(i)}, x_t^{(i)}$  are determined by posing :

$$(3) \quad \langle \mu_t^N, \pi \rangle = \langle \mu_0, \pi \rangle + \int_0^t \langle \mu_s^N, \mathcal{L}\pi \rangle \mathrm{d}s, \quad 0 \leq t \leq T, \quad \forall \pi \in \mathcal{P}_{2N-1};$$

Note that:

$$\langle \mu_0^N, \pi \rangle = \langle \mu_0, \pi \rangle, \quad \text{for all } \pi \in \mathcal{P}_{2N-1},$$

i.e.  $\mu_0^N$  is the  $N$ -points Gauss-Christoffel approximation of the  $\mu_0$  (see Section 4.1.1).

## 2.2 Convergence of the approximation

Under an additional assumption, we will establish a convergence result.

**Lemma 2.1 (Billingsley [1])** *Let  $\mu \in \mathcal{M}_+(\mathbb{R})$  with finite moments of all orders  $m_n = \langle \mu, x^n \rangle$ . Suppose that the power series :*

$$\sum_{n \in \mathbb{N}} \frac{\theta^n}{n!} m_n$$

*admits a strictly positive radius of convergence, then if  $\nu \in \mathcal{M}_+(\mathbb{R})$  is s.t.  $\langle \nu, x^n \rangle = m_n$  for all  $n$  then  $\mu = \nu$ . In this case, we say that the moment problem for  $\mu$  is well posed. As of now, we make the abuse of notation  $x^n$  to designate the polynomial function  $x \rightarrow x^n$ .*

Let:

$$\begin{aligned} m_n(t) &= \langle \mu_t, x^n \rangle, & \dot{m}_n(t) &= \mathrm{d}m_n(t)/\mathrm{d}t, \\ m_n^N(t) &= \langle \mu_t^N, x^n \rangle, & \dot{m}_n^N(t) &= \mathrm{d}m_n^N(t)/\mathrm{d}t. \end{aligned}$$

We make the additional hypothesis:

$$(H4) \quad \limsup_{n \rightarrow \infty} \left( \frac{m_{2n}(0)}{(2n)!} \right)^{\frac{1}{2n}} < \infty.$$

This assumption ensures in particular the existence of moments of all orders for  $X_0$ , and thus for  $X_t$ , for all  $0 \leq t \leq T$ . Moreover, by using the Cauchy criterion on the convergence of series, (H4) implies that the power series  $\sum_{n \in \mathbb{N}} (\theta^n/n!) m_n(0)$  has a strictly positive radius of convergence, so that according to Lemma 2.1,  $\mu_0$  is the only nonnegative measure on  $\mathbb{R}$  admitting  $(m_n(0))_{n \in \mathbb{N}}$  as moments.

**Theorem 2.2** Under assumptions (H1)-(H4) the Gauss-Galerkin approximation is convergent:

$$\mu_t^N \xrightarrow[N \rightarrow \infty]{} \mu_t, \quad t \geq 0.$$

To prove this theorem, we use several lemmas. In the following we will reason for  $t \in [0, T]$ ; all results will be true for any  $T > 0$

**Lemma 2.3** There exist real numbers  $K_n, K'_n$  which do not depend on  $N$  such that :

- (i)  $|m_n^N(t)| \leq K_n$ , for all  $(n, N)$  s.t.  $n \leq 2N - 1, 0 \leq t \leq T$ ,
- (ii)  $|\dot{m}_n^N(t)| \leq K'_n$ , for all  $(n, N)$  s.t.  $n \leq 2N - 1, 0 \leq t \leq T$ ,
- (iii) the power series  $\sum_{n \in \mathbb{N}} (\theta^n / n!) K_n$  has a strictly positive radius of convergence.

**Proof** We show that there exist  $K_n$  and  $K'_n$  such that :

- (4)  $|m_n(t)| \leq K_n$ , for all  $n, 0 \leq t \leq T$ ,
- (5)  $|\dot{m}_n(t)| \leq K'_n$ , for all  $n, 0 \leq t \leq T$ ,
- (6) the power series  $\sum_{n \in \mathbb{N}} (\theta^n / n!) K_n$  has a strictly positive radius of convergence.

Suppose that  $|m_{2n-2}(t)| \leq K_{2n-2}$ , for all  $0 \leq t \leq T$ , taking  $\varphi(x) = x^{2n}$  in (2) leads to:

$$m_{2n}(t) \leq m_{2n}(0) + c \int_0^t (m_{2n}(s) + n^2 K_{2n-2}) ds.$$

Then using Gronwall's lemma:

$$m_{2n}(t) \leq c (m_{2n}(0) + n^2 K_{2n-2}),$$

where  $c$  denotes a constant that depends on  $T$ , but not on  $n$ . Let  $K_0$  be such that  $|m_0(t)| \leq K_0$ , for all  $0 \leq t \leq T$ , we define by recurrence :

$$(7) \quad K_{2n} = c (m_{2n}(0) + n^2 K_{2n-2}),$$

then  $|m_{2n}(t)| \leq K_{2n}$ , for all  $n, 0 \leq t \leq T$ . Moreover:

$$|x|^{2n-1} \leq \frac{1}{2} \left( \frac{x^{2n}}{2n} + 2n x^{2n-2} \right),$$

then we can choose:

$$K_{2n-1} = \frac{1}{2} \left( \frac{K_{2n}}{2n} + 2n K_{2n-2} \right),$$

and (4) is thus proved. By explicitly writing  $K_{2n}$  from (7) we can show (6). (5) is verified without difficulty. To establish the lemma it suffices to note that the above argument remains valid for moments  $m_n^N(t)$  with the same constants  $K_n$  and  $K'_n$ .  $\square$

**Lemma 2.4** There exists a family of distribution laws  $(\nu_t; 0 \leq t \leq T)$ , and a subsequence  $(\nu_t^N; 0 \leq t \leq T)_{n \in \mathbb{N}}$  extracted from  $(\nu_t^N; 0 \leq t \leq T)_{n \in \mathbb{N}}$ , such that :

$$\nu_t^N \xrightarrow[N \rightarrow \infty]{} \nu_t, \quad 0 \leq t \leq T.$$

**Proof** According to Lemma 2.3 (i)-(ii), for all  $n$  fixed, the family  $(m_n^N(\cdot))$ ;  $N > (n + 1)/2$  is bounded and equicontinuous in  $\mathcal{C}[0, T]$ , and therefore relatively compact. By a Cantor diagonalization procedure we show that there exists an increasing sequence of integers  $(N_{n'})_{n' \in \mathbb{N}}$  and functions  $m_n^* \in \mathcal{C}[0, T]$ , such that:

$$(8) \quad m_n^{N_{n'}}(\cdot) \xrightarrow{n' \rightarrow \infty} m_n^*(\cdot) \quad \text{in } \mathcal{C}[0, T], \quad \forall n \in \mathbb{N}.$$

Moreover, we consider the following result [11]: Given a sequence of real numbers  $(m_p)_{p \in \mathbb{N}}$ , a necessary and sufficient condition for there to exist a non-negative measure which admits  $(m_p)_{p \in \mathbb{N}}$  for moments, is that

$$\forall P \in \mathbb{N}, C_0, C_1, \dots, C_P \in \mathbb{R} : \quad \left( \sum_{p=0}^P C_p x^p \geq 0, \quad \forall x \in \mathbb{R} \right) \Rightarrow \left( \sum_{p=0}^P C_p m_p \geq 0 \right).$$

This last property is satisfied by  $(m_n^{N_p}(t))_{n \in \mathbb{N}}$ , so is preserved at the limit  $p \rightarrow \infty$ . According to (8), there exists a nonnegative measure  $\nu_t$  which admits  $(m_n^*(t))_{n \in \mathbb{N}}$  for moments, for all  $n$  and  $0 \leq t \leq T$ . Hence for  $t \leq T$ :

$$(9) \quad \langle \mu_t^{N_p}, x^n \rangle \xrightarrow{p \rightarrow \infty} \langle \nu_t, x^n \rangle$$

According to Lemma 2.3, the power series  $\sum_{n \in \mathbb{N}} (\theta^n/n!) m_n^*(t)$  has a strictly positive radius of convergence  $\nu_t$  is the only law on  $\mathbb{R}$  which verifies (9) (see [1]), which makes it possible to assert that  $\nu_t^{N_p} \Rightarrow \nu_t$  as  $p \rightarrow \infty$  [2, p. 181].  $\square$

**Lemma 2.5** *Under the assumptions (H1)-(H3), the Fokker-Planck equation (2) has a unique solution  $t \rightarrow \mu_t$ , a function with values in  $\mathcal{M}_+(\mathbb{R})$ .*

**Proof** Using Itô's formula, we can easily verify that the law of  $X_t$  solves (2), hence the existence of a solution is proved. Let  $\psi(\cdot, \cdot) \in \mathcal{C}_b^{1,2}(\mathbb{R}_+ \times \mathbb{R})$  and  $\tilde{\mu}$  a solution of (2) with values in  $\mathcal{M}(\mathbb{R})$ . Then,

$$(10) \quad \langle \tilde{\mu}_t, \psi(t, \cdot) \rangle = \langle \tilde{\mu}_0, \psi(0, \cdot) \rangle + \int_0^t \langle \tilde{\mu}_s, \partial_s \psi(s, \cdot) + \mathcal{L} \psi(s, \cdot) \rangle ds.$$

Furthermore, we consider the backward partial differential equation:

$$(11) \quad \frac{\partial v(s, x)}{\partial s} + \mathcal{L} v(s, x) = 0, \quad s < t, \quad v(t, x) = \bar{v}(x), \quad \forall x \in \mathbb{R}$$

( $v'(s, x) := \partial v(s, x)/\partial s$ ). According to the assumptions made, and using regularity theorems for solutions of parabolic PDEs [8] we have : for all  $\bar{v} \in \mathcal{C}_c^\infty(\mathbb{R})$ , (11) admits a solution  $v \in \mathcal{C}_b^{1,2}([0, t] \times \mathbb{R})$ . After taking the difference between two solutions, to prove uniqueness it suffices to check that if  $\mu_0 = 0$  then  $\tilde{\mu}_t = 0$  for  $t \geq 0$ .

Let  $t \geq 0$  and  $\bar{v} \in \mathcal{C}_c^\infty(\mathbb{R})$ , by (11) we associate to  $\bar{v}$  an application  $v \in \mathcal{C}_b^{1,2}([0, t] \times \mathbb{R})$ . From (10), with  $\mu_0 = 0$ , and (11):

$$\langle \tilde{\mu}_t, v(t, \cdot) \rangle = \int_0^t \langle \tilde{\mu}_s, \partial_s v(s, \cdot) + \mathcal{L} v(s, \cdot) \rangle ds = 0,$$

so that  $\langle \tilde{\mu}_t, v(t) \rangle = \langle \tilde{\mu}_t, \bar{v} \rangle = 0$  for all  $\bar{v} \in \mathcal{C}_c^\infty(\mathbb{R})$ , hence  $\tilde{\mu}_t = 0$ .  $\square$

**Proof of Theorem 2.2** If we establish that

$$(12) \quad \langle \nu_t, \varphi \rangle = \langle \nu_0, \varphi \rangle + \int_0^t \langle \nu_s, \mathcal{L}\varphi \rangle ds, \quad t \leq T, \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}),$$

where  $(\nu_t)_{t \leq T}$  is the limit of a subsequence whose existence is guaranteed by Lemma 2.4, then by Lemma 2.5 we have  $\nu_t = \mu_t$ , for all  $0 \leq t \leq T$ . We deduce that a subsequence of  $\mu_t^N$  converges to  $\mu_t$ . But by redoing the demonstration, by uniqueness of the limit we show that the whole sequence converges. So we have to show that (12) is verified for all  $\varphi \in \mathcal{C}_b^2(\mathbb{R})$ . We will consider several steps.

*Step 1:* Suppose that  $\varphi$  is a polynomial function of degree  $d$ . For all  $N \geq (d+1)/2$ :

$$\langle \nu_t^N, \varphi \rangle = \langle \mu_0^N, \varphi \rangle + \int_0^t \langle \nu_s^N, \mathcal{L}\varphi \rangle ds,$$

so (12) is obtained by dominated convergence when  $N \rightarrow \infty$ .

*Step 2:* Suppose that  $\varphi(x) = e^{i\theta x} \pi(x)$ , where  $\theta \in \mathbb{R}$ ,  $\pi$  polynomial function, and  $i^2 = -1$ . Let us first take  $\varphi(x) = e^{i\theta x}$ ,  $|\theta| \leq \theta_1$ , where  $\theta_1$  is the radius of convergence given by Lemma 2.3-(iii). Let  $\varphi_n(x) = \sum_{k=0}^n (i\theta x)^k / k!$ ,  $\varphi_n$  verifies (12), so when  $n \rightarrow \infty$  we get  $\Phi(\theta) = 0$  for all  $|\theta| \leq \theta_1$  where:

$$\Phi(\theta) := \langle \nu_t, \varphi \rangle - \langle \mu_0, \varphi \rangle - \int_0^t \langle \nu_s, \mathcal{L}\varphi \rangle ds.$$

Thus for all  $j \geq 1$ ,  $\Phi^{(j)}(\theta) = 0$ ,  $|\theta| \leq \theta_1$ , where  $\Phi^{(j)}$  is the  $j$ th derivative of  $\Phi$  w.r.t.  $\theta$ , we deduce that (12) is true for any  $\varphi$  of the form  $e^{i\theta x} \pi(x)$ ,  $|\theta| \leq \theta_1$ ,  $\pi$  polynomial function. Using the inequality :

$$\left| \exp(i(\theta + \theta_1)x) + \exp(i\theta_1 x) \sum_{k=0}^n \frac{(i\theta x)^k}{k!} \right| \leq c \frac{|\theta|^{n+1}}{(n+1)!} |x|^{n+1},$$

and by the same argument, we show that (12) is verified for any  $\varphi$  of the form  $e^{i\theta x} \pi(x)$ ,  $|\theta| \leq 2\theta_1$ ,  $\pi$  polynomial function, and recursively for all  $\theta \in \mathbb{R}$ . Thus, Step 2 is proved.

*Step 3:* Suppose that  $\varphi \in \mathcal{C}^2(\mathbb{R})$  with compact support. There exists  $\varphi_n$  of the form:

$$\varphi_n(x) = \sum_{k=-n}^n a_k^n \exp(ib_k^n x) \quad \text{such that} \quad \|\varphi^{(j)} - \varphi_n^{(j)}\|_\infty \xrightarrow{n \rightarrow \infty} 0 \quad (j = 0, 1, 2),$$

where  $\varphi^{(j)}$  is the  $j$ th derivative of  $\varphi$ , (12) is verified for  $\varphi_n$  for all  $n$ , and therefore, by taking the limit  $n \rightarrow \infty$ , also for  $\varphi$ .

Step 4: Suppose  $\varphi \in \mathcal{C}_b^2(\mathbb{R})$ . Let  $\tau_n \in \mathcal{C}_b^2(\mathbb{R})$  s.t.

$$\begin{aligned} 0 &\leq \tau_n^{(j)} \leq 1, \quad j = 0, 1, 2, \\ \tau_n &= 1 \text{ on } [-n, n], \\ \tau_n &= 0 \text{ on } (-\infty, -n-1] \cup (n+1, \infty). \end{aligned}$$

Then we can apply the previous step to  $\varphi_n := \tau_n \varphi$  and by dominated convergence ( $n \rightarrow \infty$ ) we prove that  $\varphi$  satisfies (12), which ends the proof.  $\square$

**Remark 2.6** We proved the conservation of the Cauchy criterion :

$$\text{if } \lim_{n \rightarrow \infty} \left( \frac{m_{2n}(0)}{(2n)!} \right)^{\frac{1}{2n}} < \infty, \text{ then } \overline{\lim}_{n \rightarrow \infty} \left( \frac{m_{2n}(t)}{(2n)!} \right)^{\frac{1}{2n}} < \infty, \quad \forall t \in [0, T].$$

This result will be used to prove convergence in the case of nonlinear filtering.

### 3 Numerical solution of the Zakai equation

#### 3.1 Filtering with discrete time observation

We consider the system :

$$\begin{aligned} dX_t &= b(X_t) dt + \sigma(X_t) dW_t, \quad X_0 \sim \mu_0, \\ y_k &= h(X_{t_k}) + v_k, \end{aligned}$$

where  $0 \leq t \leq T$ ,  $0 < t_1 < \dots < t_K = T$  is a sequence of given instants, to simplify we take:

$$t_k = k \Delta, \text{ with } \Delta = \frac{T}{K} \text{ for some } K \in \mathbb{N}.$$

$(X_t)_{t \leq T}$ ,  $(W_t)_{t \leq T}$ ,  $(y_k)_{k \leq K}$  et  $(v_k)_{k \leq K}$  are processes with values in  $\mathbb{R}$ ;  $(v_k)_{k \leq K}$  is a sequence of independent Gaussian variables,  $v_k \sim N(0, R)$ ;  $(W_t)_{t \leq T}$  is a standard standard Wiener process independent of  $(v_k)_{k \leq K}$ ;  $X_0$  is independent of  $(W_t)_{t \leq T}$  and  $(v_k)_{k \leq K}$ . Note that the case where the observation  $y_k$  takes values in  $\mathbb{R}^d$  is treated in exactly the same way.

Let us assume Hypotheses (H1)-(H4) satisfied, as well as the hypothesis :

**(H5)**  $h : \mathbb{R} \rightarrow \mathbb{R}$  is measurable and bounded.

$X_t$  describes the evolution of a physical system, and  $y_k$  its discrete time observation. The filtering problem consists in determining  $\eta_t$ , the conditional law of  $X_t$  given  $(y_k)_{k; t_k \leq t} = (y_1, \dots, y_{\lfloor t/\Delta \rfloor})$ , that is:

$$\langle \nu_t, \varphi \rangle = \mathbb{E}(\varphi(X_t) | y_1, \dots, y_{\lfloor t/\Delta \rfloor}), \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}),$$

where  $\lfloor t/\Delta \rfloor$  is the integer part of  $t/\Delta$ . Between two moments of observation, i.e.  $t_{k-1} < t < t_k$ , the evolution of  $\eta_t$  is described by the (weak form of the) Fokker-Planck equation :

$$\frac{d}{dt} \langle \eta_t, \varphi \rangle = \langle \eta_t, \mathcal{L}\varphi \rangle, \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}).$$



At the time of observation  $t = t_k$ , by using the Bayes formula:

$$\langle \eta_{t_k}, \varphi \rangle = \frac{\langle \eta_{t_k^-}, f(\cdot, y_k) \varphi \rangle}{\langle \eta_{t_k^-}, f(\cdot, y_k) \rangle},$$

where

$$\langle \eta_{t_k^-}, \varphi \rangle := \lim_{\substack{t \rightarrow t_k \\ t < t_k}} \langle \eta_t, \varphi \rangle,$$

and  $f(x, y)$  is the local likelihood function:

$$f(x, y) := \exp\left(\frac{1}{R} h(x) y - \frac{1}{2R} h(x)^2\right).$$

Thus  $(\eta_t)_{t \leq T}$  is a solution of the equation :

$$(13) \quad \langle \eta_t, \varphi \rangle = \langle \mu_0, \varphi \rangle + \int_0^t \langle \eta_s, \mathcal{L}\varphi \rangle ds + \sum_{k=1}^{\lfloor t/\Delta \rfloor} \left\{ \frac{\langle \eta_{t_k^-}, f(\cdot, y_k) \varphi \rangle}{\langle \eta_{t_k^-}, f(\cdot, y_k) \rangle} - \langle \eta_{t_k^-}, \varphi \rangle \right\}, \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}).$$

We propose to approximate  $\eta_t$  by a probability measure of the form:

$$\eta_t^N(dx) = \sum_{i=1}^N w_t^{(i)} \delta_{x_t^{(i)}}(dx),$$

where the stochastic processes  $w_t^{(i)}$  and  $x_t^{(i)}$  are determined by posing:

$$(14) \quad \langle \eta_t^N, \pi \rangle = \langle \mu_0, \pi \rangle + \int_0^t \langle \eta_s^N, \mathcal{L}\pi \rangle ds + \sum_{k=1}^{\lfloor t/\Delta \rfloor} \left\{ \frac{\langle \eta_{t_k^-}^N, f(\cdot, y_k) \pi \rangle}{\langle \eta_{t_k^-}^N, f(\cdot, y_k) \rangle} - \langle \eta_{t_k^-}^N, \pi \rangle \right\}, \quad \forall \pi \in \mathcal{P}_{2N-1}.$$

**Theorem 3.1** *Under assumptions (H1)-(H5), for any given trajectory  $(y_1, \dots, y_K)$ , the Gauss-Galerkin approximation is convergent :  $\eta_t^N \Rightarrow \eta_t$  as  $N \rightarrow \infty$ , for all  $0 \leq t \leq T$ .*

**Proof** Let  $m_n(t) = \langle \eta_t, x^n \rangle$ , let's assume that the hypotheses:

$$(15) \quad \eta_t^N \Rightarrow \eta_t, \text{ as } N \rightarrow \infty,$$

$$(16) \quad \overline{\lim}_{n \rightarrow \infty} \left( \frac{m_{2n}(t)}{(2n)!} \right)^{\frac{1}{2n}} < \infty$$

are verified for  $t = t_{k-1}$ ; we will show that (15)-(16) are verified for  $t \in [t_{k-1}, t_k]$ . To prove the theorem it will be enough for us to establish (15)-(16) for  $t = 0$ .

For  $t \in (t_{k-1}, t_k)$ , the evolution of  $\eta_t$  is described by the Fokker-Planck equation, we deduce from Theorem 2.2, and from (16) in  $t = t_{k-1}$ , that (15) is satisfied for all  $t \in (t_{k-1}, t_k)$ . Since:

$$\langle \eta_{t_k}^N, \varphi \rangle = \frac{\langle \eta_{t_k}^N, f(\cdot, y_k) \varphi \rangle}{\langle \eta_{t_k}^N, f(\cdot, y_k) \rangle},$$

we deduce that (15) is also true for  $t = t_k$ . Moreover,

$$\begin{aligned} \overline{\lim}_{n \rightarrow \infty} \left( \frac{m_{2n}(t_k)}{(2n)!} \right)^{\frac{1}{2n}} &= \overline{\lim}_{n \rightarrow \infty} \left( \frac{1}{(2n)!} \frac{\langle \eta_{t_k}^-, f(\cdot, y_k) x^{2n} \rangle}{\langle \eta_{t_k}^-, f(\cdot, y_k) \rangle} \right)^{\frac{1}{2n}} \\ &\leq \overline{\lim}_{n \rightarrow \infty} \left( \frac{m_{2n}(t_k^-)}{(2n)!} \right)^{\frac{1}{2n}}. \end{aligned}$$

Using (15), for  $t = t_{k-1}$ , and Remark 2.6, we show that the latter expression is finite. We deduce that (16) is true for  $t = t_k$ . To end the demonstration, we just need to check (15)-(16) for  $t = 0$ . From Equation (14),  $\eta_0^N$  is the Gauss-Christoffel approximation of  $\eta_0 = \mu_0$ , and the convergence  $\eta_0^N \Rightarrow \eta_0$  can be deduced from Theorem 2.2. By Moreover (16) in  $t = 0$  is exactly Hypothesis (H4).  $\square$

### 3.2 Filtering with continuous time observation

We consider the nonlinear system :

$$(17) \quad \begin{cases} dX_t = b(X_t) dt + \sigma(X_t) dW_t, & X_0 \sim \mu_0, \\ dY_t = h(X_t) dt + dV_t, & Y_0 = 0, \end{cases}$$

for  $0 \leq t \leq T$ . The assumptions of the previous sections are assumed to be satisfied. The observation  $(Y_t)_{t \leq T}$  with values in  $\mathbb{R}$ , is here in continuous time,  $(V_t)_{t \leq T}$  is a standard Wiener process independent of  $X_0$  and  $(W_t)_{t \leq T}$ .

The filtering problem consists in determining  $\nu_t$  the conditional distribution of  $X_t$  given  $\mathcal{F}_t := \sigma(Y_s; s \leq t)$ , that is:

$$\langle \nu_t, \varphi \rangle = \mathbb{E}(\varphi(X_t) | Y_s, 0 \leq s \leq t), \quad \forall \varphi \in \mathcal{C}_b^2(\mathbb{R}).$$

To characterize  $\nu_t$ , we can use the method of the reference probability. Let  $\overset{\circ}{\mathbb{P}}$  be the law determined by :

$$\frac{d\overset{\circ}{\mathbb{P}}}{d\mathbb{P}} = Z_T^{-1}, \quad \text{with} \quad Z_t = \exp \int_0^t \left( h(X_s) dY_s - \frac{1}{2} h(X_s)^2 ds \right).$$

The computation of the conditional distribution of  $X_t$  given  $\mathcal{F}_t$  under  $\mathbb{P}$ , is related to an expression computed under  $\overset{\circ}{\mathbb{P}}$  by the Kallianpur-Striebel formula :

$$\mathbb{E}(\varphi(X_t) | \mathcal{F}_t) = \frac{\overset{\circ}{\mathbb{E}}(\varphi(X_t) Z_t | \mathcal{F}_t)}{\overset{\circ}{\mathbb{E}}(Z_t | \mathcal{F}_t)}.$$

We define  $\tilde{\nu}_t$  the unnormalized conditional distribution of  $X_t$  given  $\mathcal{F}_t$ , by posing:

$$\langle \tilde{\nu}_t, \varphi \rangle := \mathring{\mathbb{E}}(\varphi(X_t) Z_t | \mathcal{F}_t).$$

$\tilde{\nu}_t$  is a solution of the (weak form) Zakai equation:

$$(18) \quad \langle \tilde{\nu}_t, \varphi \rangle = \langle \mu_0, \varphi \rangle + \int_0^t \langle \tilde{\nu}_s, \mathcal{L}\varphi \rangle ds + \int_0^t \langle \tilde{\nu}_s, h\varphi \rangle dY_s, \quad \forall \varphi \in C_b^2(\mathbb{R}).$$

We can determine  $\tilde{\nu}_t^N$  the Gauss-Galerkin approximation of  $\tilde{\nu}_t$ , but after discretization in time, the equation of  $\tilde{\nu}_t^N$  involves only discrete time observations. It is therefore preferable to discretize the observation equation in (17) directly:

$$y_k = h(X_{t_k}) + v_k,$$

with  $t_k = k\Delta$ , and where  $v_k := (V_{t_{k+1}} - V_{t_k})/\Delta$  and  $y_k$  is the approximation of  $(Y_{t_{k+1}} - Y_{t_k})/\Delta$ .

We define  $\mathcal{F}_t^\Delta := \sigma(y_1, \dots, y_{\lfloor t/\Delta \rfloor})$  et  ${}^\Delta\nu_t$  the conditional distribution of  $X_t$  given  $\mathcal{F}_t^\Delta$ . As we saw in Section 3.1, the evolution of  $({}^\Delta\nu_t)$  is described by the equation:

$$(19) \quad \langle {}^\Delta\nu_t, \varphi \rangle = \langle \mu_0, \varphi \rangle + \int_0^t \langle {}^\Delta\nu_s, \mathcal{L}\varphi \rangle ds \\ + \sum_{k=1}^{\lfloor t/\Delta \rfloor} \left\{ \frac{\langle {}^\Delta\nu_{t_k^-}, f_\Delta(\cdot, y_k) \varphi \rangle}{\langle {}^\Delta\nu_{t_k^-}, f_\Delta(\cdot, y_k) \rangle} - \langle {}^\Delta\nu_{t_k^-}, \varphi \rangle \right\}, \quad \forall \varphi \in C_b^2(\mathbb{R}),$$

with:

$$f_\Delta(x, y) := \exp\left(h(x)y\Delta - \frac{1}{2}h(x)^2\Delta\right).$$

We have the following result:

**Theorem 3.2** *In addition to assumptions (H1)-(H5), suppose that  $h \in C_b^2(\mathbb{R})$ , then for any observed trajectory  $(Y_s)_{s \leq t}$ :*

$${}^\Delta\nu_t \xrightarrow[\Delta \rightarrow 0]{} \nu_t, \quad 0 \leq t \leq T$$

(provided that  $\nu_t$  is defined in “robust form” cf. for example [9]).

**Proof** Consider a probability space  $(\Omega, \mathcal{F}, \mathring{\mathbb{P}})$  and the following SDE on this space:

$$\begin{cases} dX_t = b(X_t) dt + \sigma(X_t) dW_t, \\ dY_t = d\mathring{V}_t, \end{cases}$$

where  $(W_t, \mathring{V}_t)$  is a standard Wiener process with values in  $\mathbb{R} \times \mathbb{R}$  independent from  $X_0$  and  $Y_0 = 0$ .

For  $Y$  we adopt the canonical representation  $(\mathcal{C}[0, T], \mathcal{B}, \mathcal{W}, Y)$ , i.e.  $(\mathcal{C}[0, T], \mathcal{B})$  is the space of continuous functions  $[0, T] \rightarrow \mathbb{R}$  equipped with the Borel  $\sigma$ -algebra  $\mathcal{B}$ ,  $\mathcal{W}$  is

the Wiener measure on this space and  $Y$  is the canonical process: for all  $\omega \in \mathcal{C}[0, T]$ ,  $Y_t(\omega) := \omega(t)$ . Moreover, let  $\bar{\mathbb{P}}$  be the marginal distribution of  $X$  on a space  $(\bar{\Omega}, \bar{\mathbb{P}}, \bar{\mathcal{F}})$ .

Under  $\overset{\circ}{\mathbb{P}}$ ,  $X$  and  $Y$  are independent:

$$(20) \quad \overset{\circ}{\mathbb{P}}(dX, dY) = \bar{\mathbb{P}}(dX) \times \mathcal{W}(dY).$$

Let

$$\Delta_K := \frac{T}{K},$$

and  $t_k^K := k \Delta_K$  which we will denote  $t_k$ . Define also:

$$h^K(t, x) := h(t_k), \quad \text{for } t \in [t_k, t_{k+1}).$$

Consider the following  $\overset{\circ}{\mathbb{P}}$ -exponential martingales:

$$\begin{cases} Z_t := \exp \int_0^t \left( h(X_s) dY_s - \frac{1}{2} h(X_s)^2 ds \right), \\ Z_t^K := \exp \int_0^t \left( h^K(X_s) dY_s - \frac{1}{2} h^K(X_s)^2 ds \right). \end{cases}$$

Let:

$$dM_s := \mathcal{L}h(X_s) ds + (h' \sigma)(X_s) dW_s,$$

integration by part in the Itô integral leads to:

$$(21) \quad Z_t = \exp \left( h(X_t) Y_t - \int_0^t \left( h(X_s) dM_s - \frac{1}{2} h(X_s)^2 ds \right) \right).$$

In addition, as  $t \rightarrow h^K(t, x)$  is piecewise constant:

$$(22) \quad Z_t^K = \exp \left( h(X_{t_k}) (Y_t - Y_{t_k}) - \frac{1}{2} h(X_{t_k})^2 (t - t_k) \right. \\ \left. + \sum_{j=0}^{k-1} \left\{ h(X_{t_j}) (Y_{t_{j+1}} - Y_{t_j}) - \frac{1}{2} h(X_{t_j})^2 \Delta_k \right\} \right), \quad \text{for all } t \in [t_k, t_{k+1}).$$

Representations (21) and (22) allow to consider  $Z_t$  and  $Z_t^K$  for any fixed trajectory of  $Y$ .

We define the distribution:

$$(23) \quad \frac{d\mathbb{P}}{d\overset{\circ}{\mathbb{P}}} := Z_T, \quad \frac{d\mathbb{P}^K}{d\overset{\circ}{\mathbb{P}}} := Z_T^K,$$

and  $\mathbb{E}, \mathbb{E}^K$  the associated expectations. Under  $\mathbb{P}$  (resp.  $\mathbb{P}^K$ ),  $(X, Y)$  admits the representation:

$$\begin{cases} dX_t = b(X_t) dt + \sigma(X_t) dW_t, \\ dY_t = h(X_t) dt + dV_t, \end{cases} \quad (\text{resp. } dY_t = h^K(X_t) dt + dV_t^K)$$

where  $V$  (resp.  $V^K$ ) is a  $\mathbb{P}$  standard Wiener process (resp.  $\mathbb{P}^K$  standard Wiener process) defined by:

$$V_t := \overset{\circ}{V}_t + \int_0^t h(X_s) ds \quad (\text{resp. } V_t^K := \overset{\circ}{V}_t + \int_0^t h^K(s, X_s) ds).$$

Consider now the system with discrete time observation:

$$\begin{cases} dX_t = b(X_t) dt + \sigma(X_t) dW_t, \\ y_k^K = h(X_{t_k}) + v_k^K, \end{cases}$$

with

$$v_k^K := \frac{1}{\Delta_k} (V_{t_{k+1}} - V_{t_k}).$$

Clearly, under  $\mathbb{P}^K$ , the conditional distribution of  $X_t$  given  $\sigma(y_k^K; k \text{ s.t. } t_k \leq t)$  is equal to the conditional distribution of  $X_t$  given  $\mathcal{F}_t := \sigma(Y_s; s \leq t)$ . Our goal is therefore to demonstrate the convergence of expressions  $\mathbb{E}^K(\varphi(X_t)|\mathcal{F}_t)$  for any continuous and bounded function  $\varphi$ .

Thanks to the Kallianpur-Striebel formula, (23) gives:

$$\mathbb{E}(\varphi(X_t)|\mathcal{F}_t) = \frac{\mathbb{E}(\varphi(X_t) Z_t | \mathcal{F}_t)}{\mathbb{E}(Z_t | \mathcal{F}_t)}, \quad \mathbb{E}^K(\varphi(X_t)|\mathcal{F}_t) = \frac{\mathbb{E}(\varphi(X_t) Z_t^K | \mathcal{F}_t)}{\mathbb{E}(Z_t^K | \mathcal{F}_t)}, \quad \overset{\circ}{\mathbb{P}}\text{-a.s.}$$

But, according to (20):

$$\mathbb{E}(\varphi(X_t) Z_t | \mathcal{F}_t) = \bar{\mathbb{E}}(\varphi(X_t) Z_t), \quad \mathbb{E}(\varphi(X_t) Z_t^K | \mathcal{F}_t) = \bar{\mathbb{E}}(\varphi(X_t) Z_t^K), \quad \mathcal{W}\text{-a.s.}$$

For a given trajectory  $(Y_s; s \leq t)$  of the observation process, it is thus necessary to prove:

$$\mathbb{E}(\varphi(X_t) Z_t^K | \mathcal{F}_t) \xrightarrow{K \rightarrow \infty} \mathbb{E}(\varphi(X_t) Z_t | \mathcal{F}_t) \quad \text{a.s.}$$

Since  $\varphi$  is bounded, it is sufficient to show that:

$$(24) \quad Z_t^K \xrightarrow{K \rightarrow \infty} Z_t \quad \text{in } L^1(\bar{\Omega}, \bar{\mathcal{F}}, \bar{\mathbb{P}})$$

For any given  $t$ ,  $Z_t^K$  and  $Z_t$  are positive random variables with mean 1, for all  $K$ , so a sufficient condition for (24) is:

$$Z_t^K \xrightarrow{K \rightarrow \infty} Z_t \quad \text{in } \bar{\mathbb{P}}\text{-probability,}$$

this result can be deduced from definitions (21) and (22) of  $Z_t$  and  $Z_t^K$ , which completes the proof of Theorem 3.2.  $\square$

We use the Gauss-Galerkin method to approximate  $\Delta \nu_t$  by a probability measure  $\Delta \nu_t^N$  of the form:

$$\Delta \nu_t^N(dx) = \sum_{i=1}^N w_t^{(i)} \delta_{x_t^{(i)}}(dx),$$

where the stochastic processes  $(w_t^{(i)})_{t \leq T}$  and  $(x_t^{(i)})_{t \leq T}$ , which depend on  $\Delta$  and  $N$ , are determined by posing:

$$\langle \Delta \nu_t^N, \pi \rangle = \langle \mu_0, \pi \rangle + \int_0^t \langle \Delta \nu_s^N, \mathcal{L}\pi \rangle ds + \sum_{k=1}^{\lfloor t/\Delta \rfloor} \frac{\langle \Delta \nu_{t_k^-}^N, (f_\Delta(\cdot, y_k) - 1) \pi \rangle}{\langle \Delta \nu_{t_k^-}^N, f_\Delta(\cdot, y_k) \rangle}, \quad \forall \pi \in \mathcal{P}_{2N-1}.$$

According to Theorem 3.1, for any  $\Delta$ , we have the following convergence:

$$(25) \quad \Delta \nu_t^N(\omega) \xrightarrow[N \rightarrow \infty]{} \Delta \nu_t(\omega), \quad \text{for almost all } \omega, \text{ and } 0 \leq t \leq T.$$

Let  $(f_p)_{p \in \mathbb{N}}$  be a dense sequence in  $\mathcal{C}_b(\mathbb{R})$ , the set of bounded and uniformly continuous functions. We define :

$$d(\mu, \nu) := \sum_{p \in \mathbb{N}} \frac{1}{2^p} \frac{|\langle \mu, f_p \rangle - \langle \nu, f_p \rangle|}{\|f_p\|_\infty},$$

with  $\|f\|_\infty = \sup\{|f(x)|; x \in \mathbb{R}\}$ ;  $d(\cdot, \cdot)$  is a metric on  $\mathcal{M}_+(\mathbb{R})$ , which induces a topology equivalent to the one induced by the weak convergence of measures [12]. Thus (25) and Theorem 3.2 implies that for all  $\Delta > 0$  we can associate  $N(\Delta) \in \mathbb{N}$  such that:

$$\Delta \nu_t^{N(\Delta)}(\omega) \xrightarrow[\Delta \rightarrow 0]{} \nu_t(\omega), \quad \text{for all } \omega \text{ a.s. and } 0 \leq t \leq T.$$

This last convergence result is not entirely satisfactory, we do not know how to explicitly choose  $N(\Delta)$ , but as in practice the observation equation is always in discrete time, for a given discretization step  $\Delta$  the convergence (25) is satisfactory.

We could have obtained a root mean square convergence in the case of continuous time observations, however this is not of great interest.

To obtain a convergence for each observed trajectory, one could think of using the ‘‘robust form’’ of the Zakai equation, this was not feasible, because the multiplication by  $\exp(-h(x) Y_t)$  brings out of the space of polynomials of degree at most equal to  $2N - 1$ .

## 4 Numerical study

### 4.1 Presentation of the algorithm

We will use the following notations:

$$\begin{array}{lll} (w_k^{(1:N)}, x_k^{(1:N)}) & \text{will denote} & (w_k^{(i)}, x_k^{(i)})_{i=1, \dots, N}, \\ w_{0:K}^{(i)} & \dots & (w_k^{(i)})_{k=0, \dots, K}, \\ w_k^{(1:N)} = \tilde{w}_k^{(1:N)} & \dots & w_k^{(i)} = \tilde{w}_k^{(i)} \text{ for } i = 1, \dots, N, \\ i = 1 : N & \dots & i = 1, \dots, N, \quad \text{etc.} \end{array}$$

#### 4.1.1 A reminder on Gauss-Christoffel quadrature methods

All the results of this section come from Wheeler [13] and Gautschi [6]. Given  $N \in \mathbb{N}$  and a nonnegative measure  $\nu \in \mathcal{M}_+(\mathbb{R})$ , we want to find  $(w_{1:N}, x_{1:N})$  such that:

$$(26) \quad \sum_{i=1}^N w_i \pi(x_i) = \langle \nu, \pi \rangle, \quad \forall \pi \in \mathcal{P}_{2N-1}.$$

It is well known that, if  $x \rightarrow \nu((-\infty, x])$  admits at least  $N$  increasing points, then (26) admits a unique solution, where the particle  $x_i$  are two by two distinct and the weights are strictly positive,  $w_i > 0$ . The empirical measure:

$$\nu^N \stackrel{\text{def}}{=} \sum_{i=1}^N w_i \delta_{x_i}$$

is the Gauss-Christoffel approximation of  $\nu$ ;  $\nu^N$  and  $\nu$  have the same  $2N$  first moments:

$$\langle \nu^N, x^p \rangle = m_p := \langle \nu, x^p \rangle, \quad p = 0 : 2N - 1.$$

To compute  $(w_{1:N}, x_{1:N})$  from the moments  $m_{0:2N-1}$  we use a classical method. We introduce  $\pi_{0:2N-1}$ , the family of orthogonal polynomial functions relative to the measure  $\nu$ , i.e.  $\pi_p$  is of degree  $p$  and  $\langle \nu, \pi_p \pi_q \rangle = 0$  if  $p \neq q$ . These polynomial functions are defined up to a multiplicative constant, we can decide for example that the coefficient of the highest degree monomial in  $\pi_p$  is 1. In this case the family  $\pi_{0:2N-1}$  satisfies a recurrence relation of the form:

$$\begin{aligned} \pi_{-1}(x) &= 0, & & \text{(by convention),} \\ \pi_0(x) &= 1, \\ \pi_{p+1}(x) &= (x - \alpha_p) \pi_p(x) - \beta_p \pi_{p-1}(x), & & p = 0 : 2N - 2, \end{aligned}$$

for all  $x \in \mathbb{R}$ , for some  $(\alpha_{0:2N-2}, \beta_{0:2N-2}) \in \mathbb{R}^{2(2N-1)}$  with  $\beta_p > 0$  for  $p \geq 1$  and  $\beta_0 = 0$ .

The calculation of  $(w_{1:N}, x_{1:N})$  is reduced to the calculation of the coefficients  $(\alpha_{0:N-1}, \beta_{0:N-1})$ , with by convention  $\beta_0 = 0$ , in the following way, let:

$$(27) \quad J_N = \begin{pmatrix} \alpha_0 & \sqrt{\beta_1} & & & \\ \sqrt{\beta_1} & \alpha_1 & \sqrt{\beta_2} & & (0) \\ & \ddots & \ddots & \ddots & \\ & & \sqrt{\beta_{N-2}} & \alpha_{N-2} & \sqrt{\beta_{N-1}} \\ (0) & & & \sqrt{\beta_{N-1}} & \alpha_{N-1} \end{pmatrix}.$$

$J_N$  has  $N$  real eigenvalues  $\lambda_{1:N}$ , two by two distinct; let  $\mathbf{v}_{1:N}$  the respectively associated orthonormal eigenvectors. We have the following result :

$$(28) \quad (w_{1:N}, x_{1:N}) = (\mathbf{v}_{1,1:N}^2, \lambda_{1:N})$$

where  $\mathbf{v}_{1,i}$  denotes the first component of the vector  $\mathbf{v}_i$ .

Next, we have to notice that using the standard moments  $m_p = \langle \nu, x^p \rangle$  is numerically not a good idea, it leads to ill-conditioned algorithms. A classical method is instead to use *modified moments*, that is:

$$\tilde{m}_p \stackrel{\text{def}}{=} \langle \nu, \tilde{\pi}_p \rangle, \quad p = 0 : 2N - 1.$$

where  $\tilde{\pi}_{0:2N-1}$  is a *given* basis of  $\mathcal{P}_{2N-1}$  formed by orthogonal vectors, these kind of polynomial functions are defined by a recurrence:

$$\begin{aligned} \tilde{\pi}_{-1}(x) &= 0, & (\text{by convention}), \\ \tilde{\pi}_0(x) &= 1, \\ \tilde{\pi}_{p+1}(x) &= (x - \tilde{\alpha}_p) \tilde{\pi}_p(x) - \tilde{\beta}_p \tilde{\pi}_{p-1}(x), & p = 1 : 2N - 2, \end{aligned}$$

(for all  $x \in \mathbb{R}$ ) where the recurrence coefficients  $(\tilde{\alpha}_{0:2N-2}, \tilde{\beta}_{0:2N-2})$  are given with  $\tilde{\beta}_p > 0$  for  $p \geq 1$  and  $\tilde{\beta}_0 \equiv 0$ . In practice, we can use the Hermite polynomials.

#### Initialization

$$\begin{aligned} \sigma_{-1,0} &\leftarrow 0 \\ \sigma_{0,p} &\leftarrow \tilde{m}_p, \quad p = 0 : 2N - 1 \\ \alpha_0 &\leftarrow \tilde{\alpha}_0 + \tilde{m}_1 / \tilde{m}_0 \\ \beta_0 &\leftarrow 0 \end{aligned}$$

#### Iterations

$$\begin{aligned} &\text{for } p = 1 : N - 1 \text{ do} \\ &\quad \text{for } q = p : 2N - p + 1 \text{ do} \\ &\quad\quad \sigma_{p,q} \leftarrow \sigma_{p-1,q+1} - (\alpha_{p-1} - \tilde{\alpha}_q) \sigma_{p-1,q} - \beta_{p-1} \sigma_{p-2,q} + \tilde{\beta}_q \sigma_{p-1,q-1} \\ &\quad \text{end for} \\ &\quad \alpha_p \leftarrow \tilde{\alpha}_p - \frac{\sigma_{p-1,p}}{\sigma_{p-1,p-1}} + \frac{\sigma_{p,p+1}}{\sigma_{p,p}} \\ &\quad \beta_p \leftarrow \frac{\sigma_{p,p}}{\sigma_{p-1,p-1}} \\ &\text{end for} \end{aligned}$$

Algorithm 1: *This modified Chebyshev algorithm allows us to compute  $(\alpha_{0:N-1}, \beta_{0:N-1})$  from  $(\tilde{m}_{0:2N-1}, \tilde{\alpha}_{0:2N-2}, \tilde{\beta}_{0:2N-2})$ , see [13].*

Finally, the computation of  $(\alpha_{0:N-1}, \beta_{0:N-1})$  from  $(\tilde{m}_{0:2N-1}, \tilde{\alpha}_{0:2N-2}, \tilde{\beta}_{0:2N-2})$  is performed using the modified Chebyshev Algorithm 1.

#### 4.1.2 Fokker-Planck equation

We first consider the approximation algorithm of the Fokker-Planck equation. The practical implementation of this algorithm requires a time discretization of the equation:

$$(29) \quad \langle \mu_t^N, \tilde{\pi}_p \rangle = \langle \mu_0, \tilde{\pi}_p \rangle + \int_0^t \langle \mu_s^N, \mathcal{L} \tilde{\pi}_p \rangle ds, \quad 0 \leq t \leq T, \quad p = 0 : 2N - 1,$$

where  $\tilde{\pi}_{0:2N-1}$  denotes a basis of  $\mathcal{P}_{2N-1}$ . In  $t = 0$ , (29) leads to the fact that  $\mu_0^N$  is the Gauss-Christoffel approximation of  $\mu_0$  so that  $\mu_0$  can be replaced by  $\mu_0^N$ .



All time discretization schemes could be considered, but in order to simplify the presentation we will use the Euler scheme with a time step  $\delta = T/L$ , with  $L \in \mathbb{N}$ . In order to simplify the notations, in the case of nonlinear filtering discussed later, we will assume that  $L$  is a multiple of  $K$ , so that the observation instants  $t_k = k \Delta$  are included in  $(\ell\delta)_{\ell=0:L}$ . Also to simplify the notation,  $\mu_{\ell\delta}^N, w_{\ell\delta}^{(i)}$  (etc.) will be noted  $\mu_\ell^N, w_\ell^{(i)}$  (etc.).

The time-discretized equation (29) is thus written:

$$(30) \quad \begin{aligned} & \langle \mu_0^N, \tilde{\pi}_p \rangle \leftarrow \langle \mu_0, \tilde{\pi}_p \rangle, \quad p = 0 : 2N - 1 \\ & \mathbf{for} \ell = 1 : L \mathbf{ do} \\ & \quad \langle \mu_\ell^N, \tilde{\pi}_p \rangle \leftarrow \langle \mu_{\ell-1}^N, \tilde{\pi}_p + \mathcal{L}\tilde{\pi}_p \delta \rangle, \quad p = 0 : 2N - 1 \\ & \mathbf{end for} \end{aligned}$$

where:

$$\mu_\ell^N(dx) = \sum_{i=1}^N w_\ell^{(i)} \delta_{x_\ell^{(i)}}(dx)$$

is the approximation of  $\mu_t^N$  at time  $t = \ell\delta$ .

From  $(w_{\ell-1}^{(1:N)}, x_{\ell-1}^{(1:N)})$ , the recurrence (30) allows us to approximate the modified moments of  $\mu_\ell^N$  with respect to the basis  $\tilde{\pi}_{0:2N-1}$ , that is :

$$\tilde{m}_p(\ell) \simeq \tilde{m}_p(\ell-1) + \delta \sum_{i=1}^N w_{\ell-1}^{(i)} \mathcal{L}\tilde{\pi}_p(x_{\ell-1}^{(i)}) \quad p = 0 : 2N - 1.$$

Given  $\tilde{m}_{0:2N-1}(\ell)$ , we now want to calculate  $(w_\ell^{(1:N)}, x_\ell^{(1:N)})$  such that :

$$(31) \quad \sum_{i=1}^N w_\ell^{(i)} \tilde{\pi}_p(x_\ell^{(i)}) = \tilde{m}_p(\ell), \quad p = 0 : 2N - 1.$$

To solve this problem, we use the Gauss-Christoffel quadrature method presented in Section 4.1.1. The Gauss-Galerkin approximation algorithm, for the Fokker-Planck equation, is given by see Algorithm 2.

### 4.1.3 Nonlinear filtering equation

For the nonlinear filtering problem, we have to solve numerically an equation of the form (cf. (14) et (18)) :

$$(32) \quad \begin{aligned} \langle \nu_t^N, \tilde{\pi}_p \rangle &= \langle \mu_0, \tilde{\pi}_p \rangle + \int_0^t \langle \nu_s^N, \mathcal{L}\tilde{\pi}_p \rangle ds \\ &+ \sum_{k=1}^{\lfloor t/\Delta \rfloor} \left\{ \frac{\langle \nu_{t_k^-}^N, f(\cdot, y_k) \tilde{\pi}_p \rangle}{\langle \nu_{t_k^-}^N, f(\cdot, y_k) \rangle} - \langle \nu_{t_k^-}^N, \tilde{\pi}_p \rangle \right\}, \quad p = 0 : 2N - 1. \end{aligned}$$

**Inputs**

$$\begin{aligned} \tilde{\alpha}_p, \tilde{\beta}_p, \quad p = 0 : 2N - 2 \\ \tilde{m}_p(0) := \langle \mu_0, \tilde{\pi}_p \rangle, \quad p = 0 : 2N - 1 \end{aligned}$$

**Iterations****for**  $\ell = 0 : L$  **do**Computation of  $(\alpha_{0:N-1}, \beta_{0:N-1})$  from  $(\tilde{m}_{0:2N-1}(\ell), \tilde{\alpha}_{0:2N-2}, \tilde{\beta}_{0:2N-2})$  (cf. Algo. 1)Computation of eigenvalues and orthonormal eigenvectors of  $J_N$  defined in (27)Computation of  $(w_\ell^{(1:N)}, x_\ell^{(1:N)})$  from (28)

$$\tilde{m}_p(\ell + 1) \leftarrow \tilde{m}_p(\ell) + \delta \sum_{i=1}^N w_\ell^{(i)} \mathcal{L} \tilde{\pi}_p(x_\ell^{(i)}), \quad p = 0 : 2N - 1$$

**end for**

Algorithm 2: *The Gauss-Galerkin approximation algorithm for the Fokker-Planck equation (2) presented with an Euler scheme (any other scheme can be used, see Section 4.1.4).*

Equation (32) after discretization using the Euler scheme, is written :

$$\begin{aligned} \langle \nu_0^N, \tilde{\pi}_p \rangle &\leftarrow \langle \mu_0, \tilde{\pi}_p \rangle \text{ for } p = 0 : 2N - 1 \\ \text{for } \ell = 1 : L \text{ do} \\ \langle \nu_\ell^N, \tilde{\pi}_p \rangle &\leftarrow \langle \nu_{\ell-1}^N, \tilde{\pi}_p + \Delta \mathcal{L} \tilde{\pi}_p \rangle \text{ for } p = 0 : 2N - 1 \quad \{\text{prediction}\} \\ \text{if } (\ell \bmod \frac{L}{K}) = 0 \text{ then} \\ k &\leftarrow \ell K / L \quad \{\text{observation index}\} \\ \langle \nu_\ell^N, \tilde{\pi}_p \rangle &\leftarrow \frac{\langle \nu_\ell^N, f(\cdot, y_k) \tilde{\pi}_p \rangle}{\langle \nu_\ell^N, f(\cdot, y_k) \rangle} \text{ for } p = 0 : 2N - 1 \quad \{\text{correction}\} \\ \text{end if} \\ \text{end for} \end{aligned}$$

The complete algorithm is then equivalent to the one presented for the of Fokker-Planck.

**4.1.4 Numerical tools**

For the approximation of the Fokker-Planck equation and of the prediction part of the nonlinear filter, we use a Runge-Kutta algorithm of order 2; one could of course use more efficient schemes if the nature of the considered problem requires it.

In practice, the basis  $\tilde{\pi}_{0:2N-1}$  of  $\mathcal{P}_{2N-1}$  used is that of the Hermite polynomial functions. For the computation of the eigenvalues of  $J_N$ , we used a variant of the *QL* algorithm for symmetric and tridiagonal matrices from the EISPACK software library [5].

**4.2 Example**

We present an example of application of the Gauss-Galerkin method in nonlinear filtering. The computations have been done on a VAX 730 computer in double precision FORTRAN 77. The approximation method applied to the Fokker-Planck equation on nonlinear examples gave good results up to  $N = 10$  ( $N$ : number of Gauss points). Beyond that, we run into problems of ill-conditioning. For the filtering problem, we first tested the method on linear examples. We compared the results obtained with those given by a Kalman-Bucy

filter. Here again, we obtained good results, even with very few Gauss points ( $N = 3$  or  $4$ ). We now present a numerical example; let us consider the nonlinear filtering problem:

$$(33) \quad \begin{cases} dX_t = -X_t dt + \sqrt{2} dW_t, & X_0 \sim N(0, 1), \\ dY_t = \exp(iX_t) dt + \rho dV_t, & Y_0 = 0, \end{cases}$$

$0 \leq t \leq T$ . The standard Wiener process  $(V_t)_{t \leq T}$  and the observation process  $(Y_t)_{t \leq T}$  take values in the complex plan ( $i^2 = -1$ );  $(W_t)_{t \leq T}$  is a real standard Wiener process independent of  $(V_t)_{t \leq T}$ ;  $X_0$  is independent of  $(W_t)_{t \leq T}$  and  $(V_t)_{t \leq T}$ . Let  $\nu_t$  be the conditional distribution of  $X_t$  given  $\mathcal{F}_t = \sigma(Y_s; s \leq t)$ . We implemented three methods of approximation of  $\nu_t$ :

**GGA Gauss-Galerkin approximation :**  $\nu_t$  is approximated by a distribution law of the form  $\nu_t^{\text{GGA}} = \sum_{i=1}^N w_t^{(i)} \delta_{x_t^{(i)}}$ , where  $N$  is the number of Gauss points. The calculation of  $\nu_t^{\text{GGA}}$  was presented in Section 4.1.3.

**FD Finite Differences :** we use a finite difference scheme in space, in order to solve numerically the Zakai equation for the unnormalized conditional density of  $\nu_t$ .  $\nu_t$  is thus approximated by a law  $\nu_t^{\text{FD}}$  of the form  $\nu_t^{\text{FD}}(dx) = p(t, x) dx$ ; for details of this method cf. Le Gland [7].

**EKF Extended Kalman filter :**  $\nu_t$  is approximated by the Gaussian distribution  $\nu_t^{\text{EKF}} = N(\hat{X}_t^{\text{EKF}}, Q_t^{\text{EKF}})$  where  $\hat{X}_t^{\text{EKF}}$  and  $Q_t^{\text{EKF}}$  are the outputs of the extended Kalman filter associated to (33).

**Remarks 4.1** (i) *These three methods are in fact implemented after discretization in time of the system (33).*

(ii) *The initial condition  $X_0$  as well as the Wiener processes (in discrete time: the Gaussian white noise)  $W_t$  and  $V_t$  have been simulated on a computer.*

(iii) *The FD method is used as a reference method: we will compare the conditional moments computed by GGA with those computed by FD. However, FD has the disadvantage that it cannot be applied in a simple way in the case where the support of  $\nu_t$  does not remain, when  $t$  varies, in a bounded and fixed domain of  $\mathbb{R}$ . Indeed in FD the conditional density  $p(t, x)$  is computed on a domain  $[-M, M]$  fixed in advance.*

For the simulation we take  $T = 10$  and  $\Delta = 0.01$ . In a first set of simulations we take  $N = 10$  and  $\rho = 0.5$ , see Figures 1-4. In a second set of simulations we take  $N = 2$  and  $\rho = 1$ , see Figure 5. In view of the numerical examples (two of which are presented at the end of the section) we can make several observations:

(i) The estimators  $\hat{X}_t^{\text{method}} := \langle \nu_t^{\text{method}}, x \rangle$  (method = GGA, FD, EKF) of  $X_t$ , given by the three methods, are equivalent (cf. Fig. 1). On the other hand, contrary to GGA, EKF gives a poor estimate of the conditional variance  $Q(t) := \langle \nu_t, x^2 \rangle - \langle \nu_t, x \rangle^2$  (cf. Fig. 2).

(ii) GGA correctly follows the evolution of the conditional moments for the first set of parameters ( $N = 10, \rho = 0.5$ ), the first 14 moments are estimated in a satisfactory way).

(iii) Even for a small number of Gauss points ( $N = 2$  in the second set of parameters), GGA gives significant results (cf. Fig. 5).

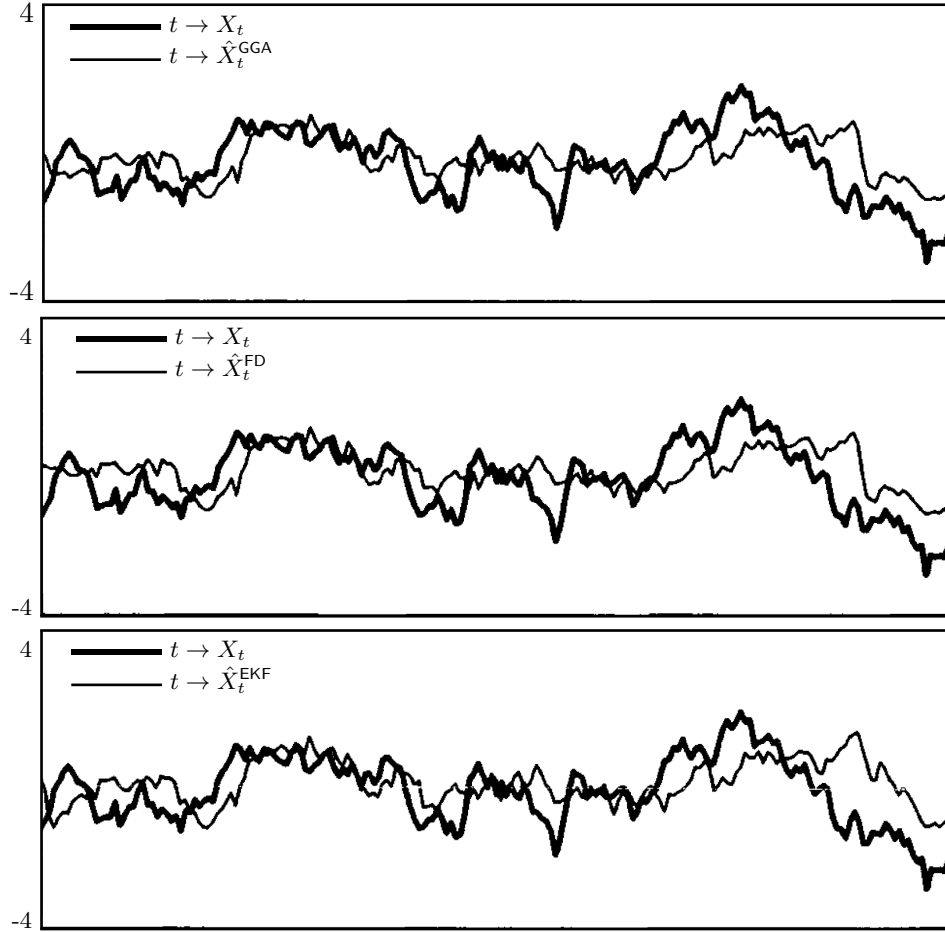


Figure 1: First set of parameters ( $N = 10$ ,  $\rho = 2$ ); comparison of the real state trajectory  $t \rightarrow X_t$  and of the estimators  $t \rightarrow \hat{X}_t^{\text{method}} := \langle \nu_t^{\text{method}}, x \rangle$  with method = GGA, FD, EKF.

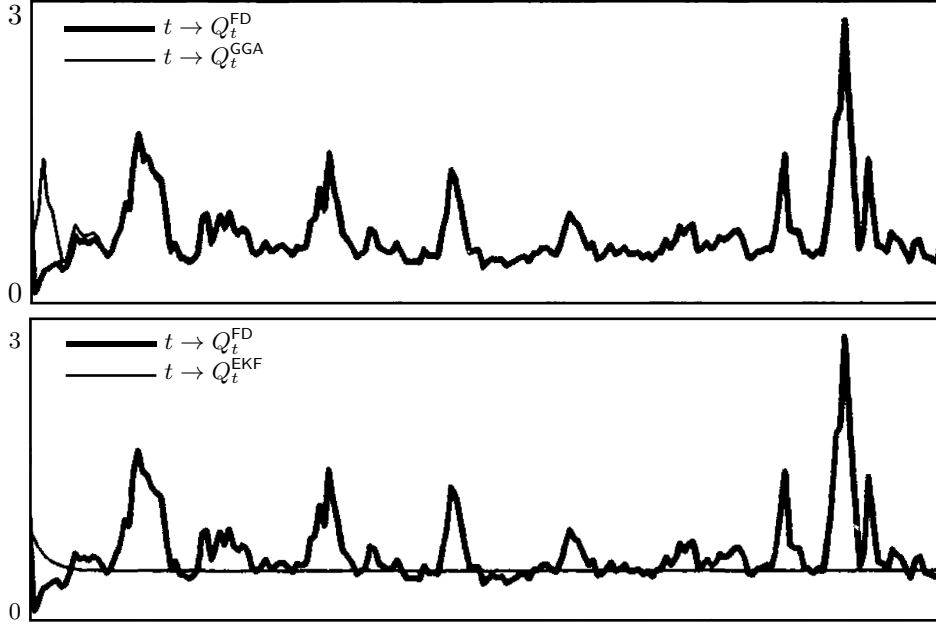


Figure 2: First set of parameters ( $N = 10, \rho = 2$ ); comparison of conditional variances,  $t \rightarrow Q^{\text{method}}(t) := \langle \nu_t^{\text{method}}, x^2 \rangle - \langle \nu_t^{\text{method}}, x \rangle^2$  with method = GGA, FD, EKF.

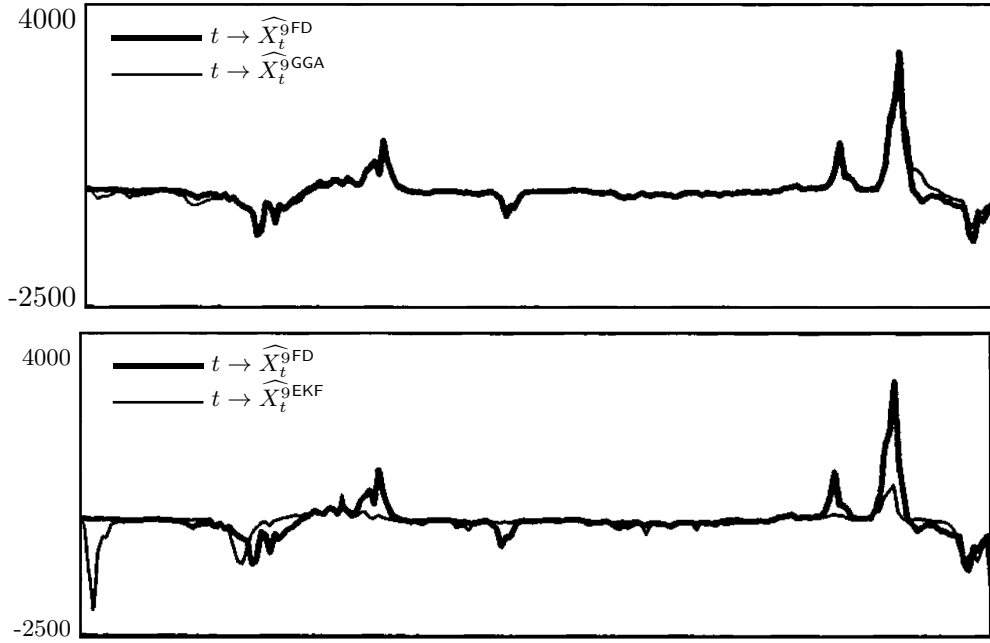


Figure 3: First set of parameters ( $N = 10, \rho = 2$ ); comparison of conditional moments of order 9,  $t \rightarrow \widehat{X}_t^{9\text{method}} = \langle \nu_t^{\text{method}}, x^9 \rangle$  with method = GGA, FD, EKF.

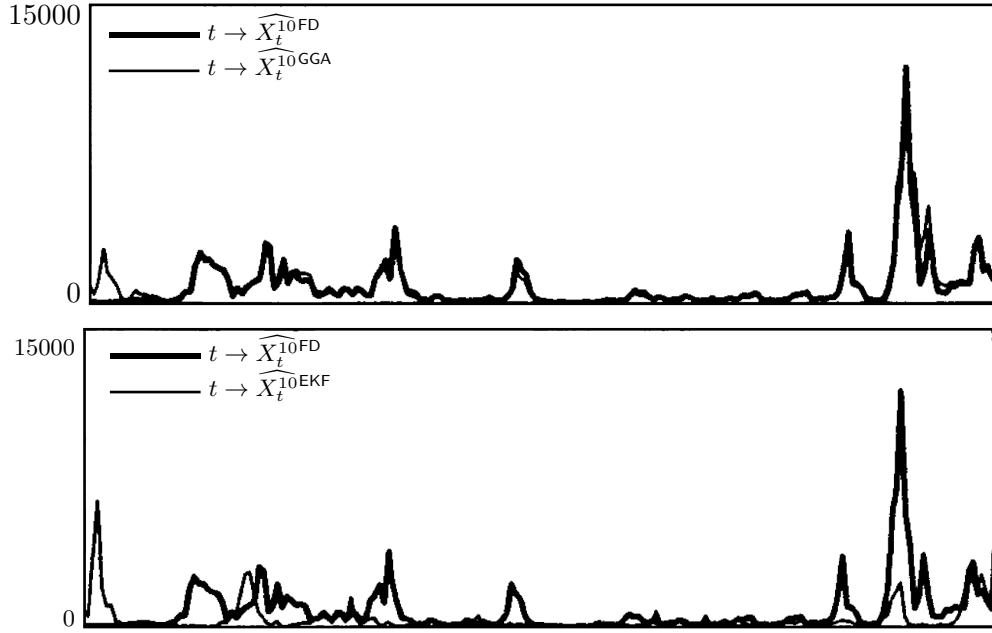


Figure 4: First set of parameters ( $N = 10, \rho = 2$ ); comparison of conditional moments of order 10,  $t \rightarrow \widehat{X}_t^{10\text{method}} = \langle \nu_t^{\text{method}}, x^{10} \rangle$  with method = GGA, FD, EKF.

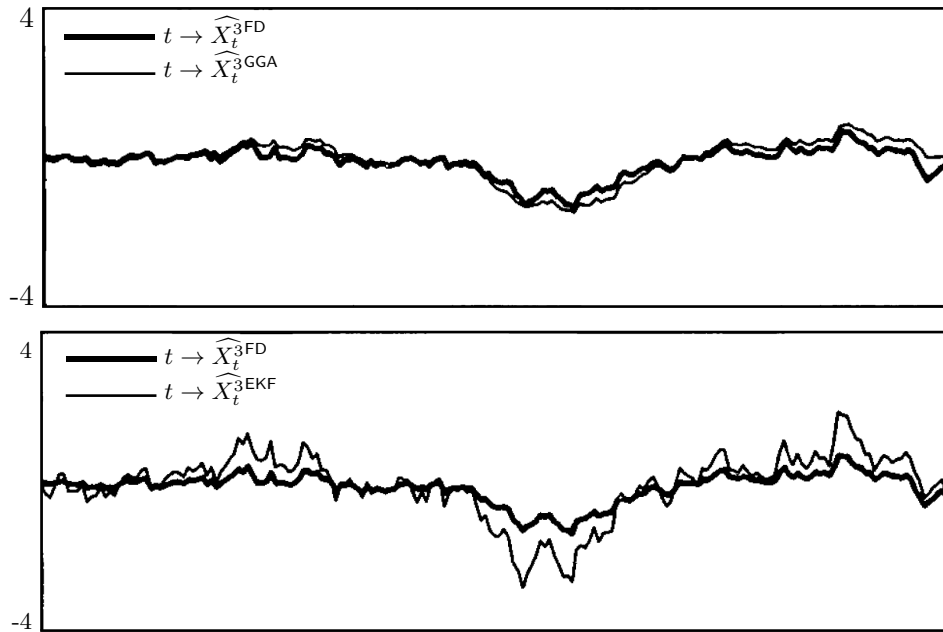


Figure 5: Second set of parameters ( $N = 2, \rho = 1$ ); comparison of conditional moments of order 3,  $t \rightarrow \widehat{X}_t^{3\text{method}} = \langle \nu_t^{\text{method}}, x^3 \rangle$  with method = GGA, FD, EKF.

## References

- [1] Patrick Billingsley. *Probability and Measure*. John Wiley & Sons, 1979.
- [2] Leo Breiman. *Probability*. Classics in Applied Mathematics. Addison Wesley, Philadelphia, 1968. First edition in 1968.
- [3] Fabien Campillo. *Filtrage et Détection de Ruptures de Processus Partiellement Observés*. PhD thesis, Thèse de Troisième Cycle, Université de Provence, Marseille, 1984.
- [4] Donald A. Dawson. Galerkin approximation of nonlinear Markov processes. In *Statistics and related topics (Ottawa, Ont., 1980)*, pages 317–339. North-Holland, Amsterdam, 1981.
- [5] Burton S. Garbow, James M. Boyle, Jack Dongarra, and Cleve B. Moler. *Matrix Eigensystem Routines - EISPACK Guide Extension*. Lecture Notes in Computer Science (LNCS, volume 51). Springer Verlag, 1977.
- [6] Walter Gautschi. On generating orthogonal polynomials. *SIAM Journal on Scientific and Statistical Computing*, 3(3):289–317, 1982.
- [7] François Le Gland. *Estimation de Paramètres dans les Processus Stochastiques, en Observation Incomplète — Applications à un Problème de Radio-Astronomie*. PhD thesis, Thèse de Docteur-Ingénieur, Université de Paris IX – Dauphine, 1981.
- [8] Olga A. Ladyzhenskaya and Nina N. Uraltseva. *Linear and quasilinear elliptic equations*. Academic Press, 1968.
- [9] Étienne Pardoux. Équations du filtrage non linéaire de la prédiction et du lissage. *Stochastics*, 6(3-4):193–231, 1982.
- [10] Pierre-Arnaud Raviart. An analysis of particle methods. In F. Brezzi, editor, *Numerical Methods in Fluid Dynamics – Lectures given at the 3rd 1983 Session of the Centro Internazionale Matematico Estivo (CIME) held at Como, Italy, July 7-15, 1983*, Lecture Notes in Mathematics vol. 1127. Springer Verlag, Berlin, 1985.
- [11] James A. Shohat and Jacob D. Tamarkin. *The Problem of Moments*. American Mathematical Society, 1950.
- [12] Daniel W. Stroock and S. R. Srinivasa Varadhan. *Multidimensional Diffusion Processes*. Springer-Verlag, 1979.
- [13] John C. Wheeler. Modified moments and Gaussian quadratures. *Rocky Mountain Journal of Mathematics*, 4(2):287–296, 1974.

## Addendum

This is the English translation of the paper [B]. A number of typos (many...) have been corrected, some notations and demonstrations have been clarified. Also some elements from the original works [3], that was not detailed or present in the 1986 version, such as the proof of Theorem 3.2, are developed here in order to obtain a self-contained version.

This article contains, to my knowledge, the first occurrence of the term “**particle approximation**” in the context of nonlinear filtering. Indeed, the conditional law  $\eta_t$  of the state given the observations is approximated by an empirical law of the form:

$$\eta_t(dx) \simeq \eta_t^N(dx) = \sum_{i=1}^N w_t^{(i)} \delta_{x_t^{(i)}}(dx),$$

where  $w_t^{(i)} \geq 0$ ,  $\sum_{i=1}^N w_t^{(i)} = 1$ , and  $\delta_{x_t^{(i)}}(dx)$  is the Dirac measure on the particle  $x_t^{(i)}$ .

I coined this term in reference to the recent work at the time of Pierre-Arnaud Raviart on the approximation of solutions of first-order PDEs: “the exact solution is approximated by a linear combination of Dirac measures in the space variables” [C] and [10].

In the following years we proposed another particle approximation method in nonlinear filtering limited to the noise-free state equation case. In this case, the infinitesimal generator  $\mathcal{L}$  is of first order making it possible to use the particle approximation methods proposed by P.A. Raviart [C]. Although proposed in a rather limited case, the proposed approximation method constitutes one of the premises of what will be called later “particle filtering” or “sequential Monte Carlo”. In our approach, a crucial step was however missing, the famous bootstrap step! This idea, in the context of nonlinear filtering, came to the table later, in the beginning of the 90s [E-F].

- [A] F. Campillo. *La méthode d'approximation de Gauss-Galerkin – Application à l'équation du filtrage non linéaire*, Master Thesis, Université de Provence, 1982 [PDF]
- [B] F. Campillo. *La méthode d'approximation de Gauss-Galerkin en filtrage non linéaire*. RAIRO M2AN, 20(2):203–223, 1986. [PDF]
- [C] P.A. Raviart, *Particle approximation of first order systems*, Journal of Computational Mathematics, 1(4):50-61, 1986.
- [D] F. Campillo, F. Legland, *Approximation particulière en filtrage non linéaire. Application à la trajectographie*, 22ème Congrès National d'Analyse Numérique, Loctudy, 1990. [PDF]
- [E] P. Del Moral, J.C. Noyer, G. Rigal, G. Salut, *Traitement non-linéaire du signal par réseau particulière: Application radar*, 14ème Colloque sur le Traitement du Signal et des Images (GRETSI), Juan les Pins 1993.
- [F] N.J Gordon, D.J. Salmond, A.F.M. Smith, *Novel approach to nonlinear/non-Gaussian Bayesian state estimation*, IEE Proceedings, Part F, 2(140):107–113, 1993.

*In [B], I regrettably forgot to thank Walter Gautschi. Summer of 1984, a few months before the defense of my thesis, I indeed needed some additional elements concerning the Gauss quadrature methods using orthogonal polynomial functions. As Walter Gautschi was visiting Europe, I had invited him to Marseille. He completely clarified the situation for me. To thank him I proposed him to visit Aix-en-Provence... but my car broke down on the highway, Walter Gautschi finally got to visit Aix-en-Provence in record time ! I warmly thank him.*