



**HAL**  
open science

# Comparative survey on non linear filtering methods: the quantization and the particle filtering approaches

Afef Sellami

► **To cite this version:**

Afef Sellami. Comparative survey on non linear filtering methods: the quantization and the particle filtering approaches. 2005. hal-00012274v1

**HAL Id: hal-00012274**

**<https://hal.science/hal-00012274v1>**

Preprint submitted on 18 Oct 2005 (v1), last revised 20 Oct 2005 (v2)

**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.

# Quantization based filtering method using first order approximation

Afef SELLAMI

Laboratoire de Probabilités et

Modèles Aléatoires

CNRS, UMR 7599

and Université Paris Dauphine

sellami@math.jussieu.fr

October 17, 2005

## Abstract

The quantization based filtering method (see [15], [16]) is a grid based approximation method to solve nonlinear filtering problems with discrete time observations. It relies on off-line preprocessing of some signal grids in order to construct fast recursive schemes for filter approximation. We give here an improvement of this method by taking advantage of the stationary quantizer property. The key ingredient is the use of vanishing correction terms to describe schemes based on piecewise linear approximations. Convergence results are given and comparison with sequential Monte Carlo methods is made. Numerical results are presented for the particular cases of linear Gaussian model and stochastic volatility models.

**Key words:** Quantization, nonlinear filtering, off-line preprocessing, stationary quantizer, particle filtering, stochastic volatility models.

## 1 Introduction

In several scientific fields, it is often required to estimate the changing state of a system using noisy observations of its evolution over time. A common manner to do this is the Bayesian approach which constructs the probability density function (pdf) of the state at a given date conditionally to all the available observations till this date.

In the Gaussian linear case, called also the Kalman case (KF) [8, 1], the required pdf is Gaussian and by computing sequentially its two first moments, we can determine it exactly. So in this case an explicit solution is provided. Unfortunately, except in this case, or in a few other cases like the discrete finite state space [1] and some other mixing Gaussian models [7], there is usually no closed expression to the problem solution. So, many numerical estimations have been suggested to represent and recursively produce approximations of the state pdf.

In this context, two different approaches can be mentioned: first, the required pdf is represented as a sample which would provide an approximation of the distribution when its size becomes very large [6], this includes for example bootstrap Bayesian method [9] or the interacting particle filter [13, 14]. Second, a quantization of the state space is used in order to come back to the discrete finite case. As the size of the quantizations grows to infinity, it is shown that we can asymptotically approach the continuous infinite state space case. Here, the deal will be in estimating some weights associated to some given *grid* points, which define a finite state discrete distribution. This distribution will approach the continuous space case as the *grid* size gets larger. The weighted Monte Carlo filter [1, 14, 6] using random samples to compute grids and the Kitagawa method [12] for linear non Gaussian models using predefined grids and optimal quantization filtering [15] using off line computations to produce an optimal quantization of the state process are examples of this approach.

The technique of optimal quantization of random vectors is especially useful in problems where many expectations or conditional expectations need to be computed. It appears as an efficient method to transform an integral into a finite weighted sum with a controlled approximation error. We can find some applications of this technique in [16, 2]. In [17], some numerical methods to construct optimal quantization grids for multidimensional Gaussian distributions are given.

Now for the pdf estimation problem we treat here, we use Kallianpur-Striebel formula [11] to derive a dynamic programming formula allowing to estimate the pdf recursively. Like in [15], this approach makes possible the use of quantization at each time step in order to compute conditional expectations. We will call the algorithm introduced in [15] the *zero order scheme*. In this paper, we are interested by first order approximation using optimal or at least stationary quantizers to estimate the required pdf. This approach was first introduced in [3] for solving optimal stopping time problems, namely multi-asset American option pricing. It improves the convergence rate of the method. In [16], a first sketch of this idea is presented for pdf estimation but with a pseudo-numerical scheme, which cannot be implemented in practice. Our aim here is to propose operating first order schemes which improve the convergence rate of the zero order schemes from both theoretical and practical viewpoints. We first present them in a backward way; this is the natural manner to devise them and the appropriate formulation to establish error estimates. Then, we show how to derive the forward formulation to be implemented in practice.

The paper is organized as follows: in the second section we give some brief preliminaries on quantization and filtering. The third and fourth sections will deal with the algorithms using first order schemes. Each one presents the approximation procedure, the schemes in their backward and forward formulation and finally convergence theorems. Then, the fifth section is dedicated to summarize the previous results, and enlarge them to the case of normalized filters. Finally, numerical results are presented in the sixth section, including comparison with particle methods for several models.

#### Notations:

$p \in (1, +\infty)$  is a fixed real number,  $|\cdot|$  and  $\|\cdot\|_p$  denote respectively Euclidean norm on  $\mathbb{R}^d$  and  $\mathbf{L}^p$ -norm.  $\mathcal{C}_{b,Lip}^1$  is the set of continuous differentiable functions  $\mathbb{R}^d \rightarrow \mathbb{R}$ , bounded with bounded Lipschitz continuous derivative and  $\mathcal{C}_b^k$  the set of continuous  $k$ -times dif-

ferentiable functions  $\mathbb{R}^d \rightarrow \mathbb{R}$ , bounded with bounded derivatives. We will also define  $\|f\|_\infty = \sup_{x \in \mathbb{R}^d} |f(x)|$  and  $[f]_{Lip} = \sup_{x \neq x'} \frac{|f(x) - f(x')|}{|x - x'|}$ .  $\alpha > 0$  denotes a generic constant,  $\langle \cdot, \cdot \rangle$  the Euclidean inner product on  $\mathbb{R}^d$ ,  $A'$  the transpose of the real matrix  $A$ . Finally,  $(e_i)_{1 \leq i \leq d}$  is the canonical orthonormal basis of  $\mathbb{R}^d$ .

## 2 Preliminaries

### 2.1 Quantization filtering schemes

We consider a fixed discrete horizon  $n \in \mathbb{N}^*$  and some probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . A signal process is an  $\mathbb{R}^d$ -valued discrete time hidden Markov chain  $(X_k)_{0 \leq k \leq n}$  evolving according to the following signal equation:

$$X_{k+1} = F_{k+1}(X_k, \varepsilon_{k+1}), \quad 0 \leq k \leq n-1, \quad (2.1)$$

where  $F_k : \mathbb{R}^d \times \mathbb{R}^q \rightarrow \mathbb{R}^d$ , is a Borel function and  $(\varepsilon_k)_{1 \leq k \leq n}$  is a sequence of iid  $\mathbb{R}^q$ -valued random variables, independent of  $X_0$ . The distribution  $\mu_0$  of  $X_0$  is supposed to be known. Furthermore,  $\mathbf{P}_k(x, dx')$  will denote the probability transition of  $X_k$ , and:

$$\mu_0 f = \int f(x) \mu_0(dx) \quad \text{and} \quad \mathbf{P}_k f(x) = \int f(x') \mathbf{P}_k(x, dx').$$

At each time step  $k$ ,  $Y_k$  an  $\mathbb{R}^d$ -valued noisy observation of  $X_k$  is made. The dynamics of the observation process  $(Y_k)_{0 \leq k \leq n}$  are driven by Borel functions  $G_k : \mathbb{R}^d \times \mathbb{R}^d \times \mathbb{R}^q \rightarrow \mathbb{R}^d$  so that:

$$Y_k = G_k(Y_{k-1}, X_k, \eta_k), \quad 1 \leq k \leq n, \quad (2.2)$$

where  $(\eta_k)$  is a sequence of iid  $\mathbb{R}^d$ -valued random variables, independent of  $\sigma(X_0, \varepsilon_k, k \geq 1)$ . We assume for convenience, that  $Y_0 = 0$  and that, for every  $1 \leq k \leq n$ , the distribution of  $Y_k$  given  $X_k$  and  $Y_{k-1}$  admits a continuous conditional pdf  $y \mapsto g_k(Y_{k-1}, X_k, y)$ . We suppose in addition that  $g_k$  satisfies the following Lipschitz assumption:

$$\forall x, x' \in \mathbb{R}^d, \quad \forall y, y' \in \mathbb{R}^d,$$

$$|g_k(y, x, y') - g_k(y, x', y')| \leq [g_k]_{Lip}^{y, y'} |x - x'| \quad \text{and} \quad \max_{0 \leq k \leq n} \sup_{x \in \mathbb{R}^d} |g_k(y, x, y')| \leq L^{y, y'} < +\infty.$$

**Remark 2.1** As the observation process is fixed, we will drop the dependency of  $[g_k]_{Lip}^{y, y'}$  and  $L^{y, y'}$  in  $(y, y')$  for notational convenience.

The problem we aim to solve is to compute

$$\Pi_n f = \mathbb{E}[f(X_n) | Y_1 = y_1, \dots, Y_n = y_n],$$

for any reasonable Borel function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and a given observation sequence  $y = (y_1, \dots, y_n)$ .

Using Kallianpur-Striebel formula [11], the problem can be reduced to the computation of the unnormalized filter  $\pi_n$  defined by:

$$\pi_n f = \mathbb{E}[f(X_n) \prod_{k=1}^n g_k(y_{k-1}, X_k, y_k)].$$

Then,  $\Pi_n f = \frac{\pi_n f}{\pi_n \mathbf{1}}$ .

**Remark 2.2** For convenience, the dependency of  $\Pi_n$  and  $\pi_n$  in the observation process has been omitted, as  $y$  is fixed. For the same reason, we will denote  $g_k(x) := g_k(y_{k-1}, x, y_k)$  for  $1 \leq k \leq n$ , and  $g_0 := \mathbf{1}$ .

By introducing the operators  $(H_k)_{0 \leq k \leq n}$  defined below, a sequential definition of the unnormalized filter  $\pi_n$  can be given.

Namely, if one defines, for every  $x \in \mathbb{R}^d$ :

$$\begin{cases} H_k f(x) = g_k(x) \mathbb{E}[f(X_{k+1}) | X_k = x], & 0 \leq k \leq n-1, \\ H_n^n f(x) = g_n(x) f(x), \end{cases} \quad (2.3)$$

then we have

$$\pi_n f = \mu_0 \circ H_0 \cdots \circ H_n^n f. \quad (2.4)$$

Consequently, we can write sequentially, either in the forward way:

$$U_0 = \mu_0 \circ H_0, \quad U_k = U_{k-1} \circ H_k, \quad 1 \leq k \leq n-1, \quad (2.5)$$

or in the backward way:

$$R_n = H_n^n, \quad R_k = H_k \circ R_{k+1}, \quad 0 \leq k \leq n-1, \quad (2.6)$$

so that  $\pi_n f = \mu_0 R_0 f = U_{n-1} \circ H_n^n f$ .

**Remark 2.3** Note that if  $G_k$  depends on  $X_{k-1}$  instead of  $X_k$  for  $1 \leq k \leq n$ , we are led to consider the conditional pdf of  $Y_k$ , given  $X_{k-1}$  and  $Y_{k-1}$ . We can then define differently the operators  $H_k$  so that  $\pi_n f$  still satisfy formally equation (2.4).

Namely,

$$\begin{cases} H_k f(x) = g_{k+1}(x) \mathbb{E}[f(X_{k+1}) | X_k = x], & 0 \leq k \leq n-1, \\ H_n^n f(x) = f(x). \end{cases} \quad (2.7)$$

Then, schemes (2.5) and (2.6), with this new definition of the  $(H_k)$  operators, are still valid.

**Remark 2.4** When  $G_k$  depends on both  $X_{k-1}$  and  $X_k$ , we can also adapt the scheme to the modified  $\mathbb{R}^{2d}$ -valued signal Markov chain  $Z_k = (X_{k-1}, X_k)$  and the same observation process  $Y_k$ . In this case we define the new observation dynamics:

$$\bar{G}_k(Y_{k-1}, Z_k, \eta_k) \stackrel{Def}{=} G_k(X_{k-1}, Y_{k-1}, X_k, \eta_k).$$

We succeed then to restore state equations of type (2.1) and (2.2). The point is that in this case, the signal dimension is twice the original one. This can be numerically constraining, particularly when using grid based approximation methods.

From the recursive definition of either  $U_k$  or  $R_k$ , it becomes clear that it will be useful to approximate  $X_k$  by a random variable  $\hat{X}_k$  taking a finite number of values, in order to transform conditional expectations in finite weighted sums. This operation is commonly called *quantization*, and is extensively used in signal processing fields (see [10, 2, 17]).

Temporarily, we suppose that we are able to construct such an approximation  $\hat{X}_k$ . We define the induced error  $\Delta_k := X_k - \hat{X}_k$ . Further details about the error modulus  $\|\Delta_k\|_p$ ,  $p \geq 1$  will be given in the next paragraph. In [15], these quantizations  $\hat{X}_k$  are used to produce a piecewise constant approximation of  $R_k$ . So, the *natural* approximation procedure by quantization, as defined in (2.8) below appears as a zero order scheme.

It is defined as follows:

$$\begin{cases} \hat{H}_k f(\hat{X}_k) = g_k(\hat{X}_k) \mathbb{E}[f(\hat{X}_{k+1}) | \hat{X}_k], & 0 \leq k \leq n-1, \\ \hat{H}_n^n f(\hat{X}_n) = g_n(\hat{X}_n) f(\hat{X}_n). \end{cases} \quad (2.8)$$

Defining  $\hat{\mu}_0$  the discrete distribution of  $\hat{X}_0$ , we have respectively the following forward and backward iterative zero order approximation schemes:

$$\hat{U}_0 = \hat{\mu}_0 \hat{H}_0, \quad \hat{U}_k = \hat{U}_{k-1} \circ \hat{H}_k, \quad 1 \leq k \leq n-1, \quad (2.9)$$

and

$$\hat{R}_n = H_n^n, \quad \hat{R}_k = \hat{H}_k \circ \hat{R}_{k+1}, \quad 0 \leq k \leq n-1, \quad (2.10)$$

so that  $\hat{\pi}_n f = \hat{\mu}_0 \hat{R}_0 f = \hat{U}_{n-1} \circ H_n^n f$ .

Formally, this scheme is slightly different from that presented in [15] (the definition of  $H_k$  operators is different inducing a shifted scheme structure). Nevertheless, the zero order quantization filter estimator itself remains the same. This form of the scheme allows to produce costlessly some error bounds for a wider class of test functions  $f$  than in the original theorem established in [15].

**Theorem 2.1** *Assume that the transition kernels  $\mathbf{P}_k$  of the signal Markov chain are  $K$ -Lipschitz operators i.e  $\forall f : \mathbb{R}^d \rightarrow \mathbb{R}$  Lipschitz,  $[\mathbf{P}_k f]_{Lip} \leq K[f]_{Lip}$ .*

*Then, for any  $f$  such that  $H_n^n f$  is bounded Lipschitz continuous, and  $0 \leq k \leq n$ , there exists a sequence of positive constants  $(C_j^{k,n})_{k \leq j \leq n}$  such that:*

$$\|R_k f(X_k) - \hat{R}_k f(\hat{X}_k)\|_p \leq \sum_{j=k}^n C_j^{k,n} \|\Delta_j\|_p$$

and  $C_j^{k,n} \leq \alpha(p, f) L^{n-k} \frac{K^{n-j+1} - 1}{K-1}$ .

**Proof.**

The proof of this result is easily adapted from [15] by considering the *shifted* scheme (2.10), based on the definition (2.3) of the  $H_k$  operators. We simply take in consideration that at the last date, we will have  $H_n^n f$  instead of  $f$ . For that reason, the Lipschitz bounded assumption is made on  $H_n^n f$  rather than on  $f$ . For a detailed proof, see [18].  $\square$

**Remark 2.5** This shifted structure (2.10) of the zero order scheme can be useful since regularity and boundedness assumptions have to be satisfied by  $H_n^n f$  instead of  $f$  (see [15]). This is an advantage, particularly when the conditional pdf  $g_k$  goes to zero very fast as  $|x| \rightarrow +\infty$ . For example, if  $g_n(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{|y_n - x|^2}{2})$ ,  $H_n^n f$  is Lipschitz continuous and bounded for  $f$  bounded Lipschitz continuous as well as for any Lipschitz function  $f$  such that  $|f(x)| = O(\exp(\frac{\alpha|x|^2}{2}))$  for some  $0 < \alpha < 1$ .

**Corollary 2.1** *If  $\mathbf{P}_k$  is Lipschitz and  $H_n^n f$  is bounded Lipschitz continuous, then there exists a sequence of positive constants  $(C_j^n)_{0 \leq j \leq n}$  such that:*

$$|\pi_n f - \hat{\pi}_n f| \leq \sum_{j=0}^n C_j^n \|\Delta_j\|_p.$$

Let us now examine the *quantization error*  $\Delta_k$  and try to establish some convergence rate toward 0, in which case Corollary 2.1 will give a convergence rate of the zero order quantization filter estimation.

## 2.2 Background on quantization and optimal quantization

The aim of quantization is the definition of a random variable taking finite number of values in  $\mathbb{R}^d$  as an approximation of an  $\mathbb{R}^d$ -valued one. In this paragraph, we will present results useful to our work, further details can be found in [10, 17].

Let  $X : (\Omega, \mathcal{F}, \mathbb{P}) \rightarrow \mathbb{R}^d$  be a random vector and let  $\mathbb{P}_X$  denote its probability distribution. A positive integer  $N$  being fixed, let  $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$  be a Borel map such that  $|h(\mathbb{R}^d)| \leq N$ .

We say that  $h(X)$  is a  $N$ -quantization of  $X$  and that  $h(\mathbb{R}^d)$  is a  $N$ -quantizer. For convenience, the function  $h$  itself will be called  $N$ -quantizer.

Now, when  $X \in L^p(\Omega)$ , we aim to construct an  $L^p$ -optimal  $N$ -quantization of  $X$ . That is to determine the function  $h$ , if any, which minimizes the  $L^p$ -quantization error.

This amounts to solving the optimization problem:

$$\inf\{\|X - h(X)\|_p^p, h : \mathbb{R}^d \rightarrow \mathbb{R}^d, \text{ Borel map s.t. } |h(\mathbb{R}^d)| \leq N\}. \quad (2.11)$$

This optimization problem has (at least) one solution (see e.g [10]). Any such a solution  $h^*$  is called an  $L^p$ -optimal  $N$ -quantizer (or  $L^p$ -optimal  $N$ -codebook). Furthermore, one shows that  $L^p$ -optimal  $N$ -quantizers have full size i.e  $|h^*(\mathbb{R}^d)| = N$  and we denote  $\Gamma^* := h^*(\mathbb{R}^d) = \{x^1, \dots, x^N\}$ . It is clear that in this case,  $h^*$  will necessarily be a projection following the nearest neighbor rule on  $\Gamma^*$ . Namely:

$$h^*(\xi) = \sum_{i=1}^N x^i \mathbf{1}_{\mathbf{C}_i(\Gamma^*)}(\xi) \quad (2.12)$$

where  $(\mathbf{C}_i(\Gamma^*))_{1 \leq i \leq N}$ , called the Voronoi diagram of  $\Gamma^*$ , makes up a Borel partition of  $\mathbb{R}^d$  satisfying :

$$\mathbf{C}_i(\Gamma^*) \subset \{\xi \in \mathbb{R}^d \text{ s.t. } |\xi - x^i| = \min_{1 \leq k \leq N} |\xi - x^k|\}.$$

As a consequence, the induced  $L^p$ -mean quantization error (or  $L^p$ -distortion) reads:

$$\underline{\mathcal{D}}_N^{X,p} := \|X - h^*(X)\|_p^p = \left\| \min_{1 \leq i \leq N} |X - x^i| \right\|_p^p.$$

According to [10, 2],  $\underline{\mathcal{D}}_N^{X,p}$  is a (strictly) decreasing sequence converging to 0 when  $N \rightarrow +\infty$ . Furthermore, the rate of convergence of  $\underline{\mathcal{D}}_N^{X,p}$  toward 0 is ruled by Zador's Theorem:

**Theorem 2.2** (see [10, 2]) *Assume that  $\int_{\mathbb{R}^d} |\xi|^{p+\eta} \mathbb{P}_X(d\xi) < +\infty$  for some  $\eta > 0$ . Then*

$$\lim_N (N^{\frac{p}{d}} \underline{\mathcal{D}}_N^{X,p}) = J_{p,d} \|\varphi\|_{\frac{d}{d+p}}$$

where  $\mathbb{P}_X(d\xi) = \phi(\xi) \lambda_d(d\xi) + \bar{\mu}(d\xi)$ ,  $\bar{\mu} \perp \lambda_d$  ( $\lambda_d$  Lebesgue measure on  $\mathbb{R}^d$ ) and for every  $q \in \mathbb{R}_+^*$ ,  $\|g\|_q := (\int |g|^q(u) du)^{\frac{1}{q}}$ .

This theorem, combined with Corollary 2.1 establishes a convergence rate result for the quantization based zero order scheme (2.9).

Now let us introduce an important property of quadratic optimal quantizers:

**Proposition 2.1 (Stationary quantizer property)**

*If  $\hat{X}$  is a  $L^2$ -optimal  $N$ -quantization of  $X$ , then the stationary quantizer property is verified. Namely,*

$$\mathbb{E}[X|\hat{X}] = \hat{X}. \quad (2.13)$$

This property is of great help to appreciate the quality of some estimations. This is shown in further details in [17] for numerical integration and in [3] for optimal stopping problems. To illustrate this point by a short example, take the problem of approximating  $f(X)$  by  $f(\hat{X})$ , when  $f \in \mathcal{C}_b^2$ . We have for some  $\xi \in (X, \hat{X})$ :

$$f(X) - f(\hat{X}) = \langle Df(\hat{X}), \Delta \rangle + \frac{1}{2} \Delta' D^2 f(\xi) \Delta.$$

So, if  $\hat{X}$  is a stationary  $N$ -quantization of  $X$ , we have:

$$\begin{aligned} \mathbb{E}[f(X)|\hat{X}] - f(\hat{X}) &= \langle Df(\hat{X}), \mathbb{E}[\Delta|\hat{X}] \rangle + \frac{1}{2} \mathbb{E}[\Delta' D^2 f(\xi) \Delta|\hat{X}] \\ \|\mathbb{E}[f(X)|\hat{X}] - f(\hat{X})\|_p &\leq \frac{1}{2} \|D^2 f\|_\infty \|\langle \Delta, \Delta \rangle\|_p \leq \frac{1}{2} \|D^2 f\|_\infty \|\Delta\|_{2p}^2 \end{aligned}$$

We see that, owing to the stationary quantizer property (2.13) we succeed to gain one order in estimation costlessly.

Back to our filtering problem, we are interested in quantizing the Markov chain  $(X_k)_{0 \leq k \leq n}$ . We must settle at each step  $0 \leq k \leq n$ , a quantizer size  $N_k$  and an  $L^p$ -optimal  $N_k$ -quantizer of  $X_k$  denoted  $\Gamma_k = \{x_k^1, \dots, x_k^{N_k}\}$ . Consequently, we define  $(\hat{X}_k)$  an  $L^p$ -optimal  $(N_k)$ -quantization of the process  $(X_k)$  by:

$$\hat{X}_k = \sum_{i=1}^{N_k} x_k^i \mathbf{1}_{\mathbf{C}_i(\Gamma_k)}(X_k), \quad \text{for } 0 \leq k \leq n. \quad (2.14)$$

As the resulting process  $(\hat{X}_k)_{0 \leq k \leq n}$  is no longer a Markov chain, this procedure is called *marginal quantization*<sup>1</sup> of the process  $(X_k)$ .

Nevertheless, an approximation of the transition kernels  $\mathbf{P}_k$  of the chain is provided by the following *transition probability* terms:

$$p_k^{ij} = \mathbb{P}[\hat{X}_{k+1} = x_{k+1}^j | \hat{X}_k = x_k^i], \quad i \in \{1, \dots, N_k\} \text{ and } j \in \{1, \dots, N_{k+1}\}.$$

For  $0 \leq k < n$  and  $i \in \{1, \dots, N_k\}$ , we will denote

$$\hat{\mathbf{P}}_k f(x_k^i) = \mathbb{E}[f(\hat{X}_{k+1}) | \hat{X}_k = x_k^i] = \sum_{j=1}^{N_{k+1}} f(x_{k+1}^j) p_k^{ij}.$$

<sup>1</sup>More details on process quantization are given in [15].



### 2.3 Generic first order scheme

As Theorem 2.2 gives a convergence rate of  $\underline{D}_N^{X_k, p}$  toward zero, results such as Corollary 2.1 suggest that the quantization filter scheme would lead to better results if we succeed to upper bound the error by higher powers of  $\|\Delta_j\|_p$ . This leads us to the idea of mimicking first order Taylor expansions in the  $R_k$  approximation.

From now on,  $(\hat{X}_k)$  denotes a marginal stationary  $(N_k)$ -quantization of  $(X_k)$ , and we denote  $\hat{X}_k(\Omega) = \Gamma_k = \{x_k^1, \dots, x_k^{N_k}\}$ . So,  $\hat{X}_k = \sum_{i=1}^{N_k} x_k^i \mathbf{1}_{\mathbf{C}_i(\Gamma_k)}$ . Since  $\hat{X}_k$  is  $\sigma(X_k)$ -measurable, using the chaining rule for conditional expectation  $\mathbb{E}[\cdot | \hat{X}_k] = \mathbb{E}[\mathbb{E}[\cdot | X_k] | \hat{X}_k]$  yields:

$$\mathbb{E}[f(X_{k+1}) | \hat{X}_k] = \mathbb{E}[\mathbf{P}_k f(X_k) | \hat{X}_k]. \quad (2.15)$$

In view of Proposition 2.1, if  $\mathbf{D}(\mathbf{P}_k f)$  exists (and is Lipschitz) we can write:

$$\mathbb{E}[f(X_{k+1}) | \hat{X}_k] = \mathbf{P}_k f(\hat{X}_k) + \langle \mathbf{D}(\mathbf{P}_k f)(\hat{X}_k), \overbrace{\mathbb{E}[\Delta_k | \hat{X}_k]}^0 \rangle + O(|\Delta_k|^2). \quad (2.16)$$

We can then approach  $\mathbb{E}[f(X_{k+1}) | \hat{X}_k]$  by  $\mathbf{P}_k f(\hat{X}_k)$  with an  $L^1$ -estimation error of order  $O(\|\Delta_k\|_2^2)$ . This is the key idea for constructing first order quantization schemes. For such a purpose, we assume that:

**H 1** For any observation process  $y$ , all functions  $g_k$  lie in  $\mathcal{C}_{b, Lip}^1$  and there exists  $L > 0$  such that

$$\max_{0 \leq k \leq n} \{\|g_k\|_\infty, \|\mathbf{D}g_k\|_\infty, [\mathbf{D}g_k]_{Lip}\} \leq L.$$

and that:

**H 2**  $\mathbf{P}_k$  is  $K$ -Lipschitz and  $\forall f \in \mathcal{C}_{b, Lip}^1$ :

$$\mathbf{P}_k f \in \mathcal{C}_{b, Lip}^1 \quad \text{and} \quad [\mathbf{D}\mathbf{P}_k f]_{Lip} \leq K(\|\mathbf{D}f\|_\infty \vee [f]_{Lip}).$$

**Remark 2.6** Notice that under assumption **H2**, for  $f \in \mathcal{C}_{b, Lip}^1$  we have:

$$\|\mathbf{D}\mathbf{P}_k f\|_\infty = [\mathbf{P}_k f]_{Lip} \leq K[f]_{Lip} = K\|\mathbf{D}f\|_\infty$$

Under these assumptions, we can see that  $\forall f \in \mathcal{C}_{b, Lip}^1$ ,  $\forall 0 \leq k \leq n-1$ ,  $R_k f$  defined recursively by (2.6), is differentiable and:

$$\mathbf{D}R_k f = \mathbf{D}g_k \mathbf{P}_k R_{k+1} f + g_k \mathbf{D}\mathbf{P}_k R_{k+1} f \quad (2.17)$$

So, we can establish the following proposition, using a backward induction:

**Proposition 2.2** Assuming **H1** and **H2** involves:

$\forall f \in \mathcal{C}^1$  such that  $H_n^n f \in \mathcal{C}_{b, Lip}^1$ , we have  $\forall 0 \leq k \leq n-1$ ,  $R_k f \in \mathcal{C}_{b, Lip}^1$ .

Furthermore:

$$\begin{aligned} \|R_k f\|_\infty &\leq L^{n-k} \|H_n^n f\|_\infty \\ \|\mathbf{D}R_k f\|_\infty &\leq (LK)^{n-k} \|\mathbf{D}H_n^n f\|_\infty + L^{n-k} \|H_n^n f\|_\infty \frac{K^{n-k} - 1}{K - 1} \\ u_k &:= \|\mathbf{D}R_k f\|_\infty \vee [\mathbf{D}R_k f]_\infty \\ &\leq (3LK)^{n-k} u_n + L^{n-k} \|H_n^n f\|_\infty \frac{(3K)^{n-k} - 1}{3K - 1} \end{aligned}$$

with the convention  $\frac{K^m-1}{K-1} = m$  when  $K = 1$ .

**Proof.** The proof is based on an induction on  $k$ . Suppose for a given  $0 \leq k \leq n-1$ ,  $R_{k+1}f \in \mathcal{C}_{b,Lip}^1$ .

(Notice that  $H_n^n f \in \mathcal{C}_{b,Lip}^1$  by assumption).

By definition, we have  $R_k f = g_k \mathbf{P}_k R_{k+1} f$ .

According to **H1** and **H2**, we can establish easily that  $R_k f \in \mathcal{C}_{b,Lip}^1$ , through a backward induction.

Furthermore,

$$\begin{aligned} \|R_k f\|_\infty &\leq L \|\mathbf{P}_k R_{k+1} f\|_\infty \\ &\leq L \|R_{k+1} f\|_\infty \end{aligned} \quad (2.18)$$

From (2.17) and Remark 2.6, we have also:

$$\begin{aligned} \|DR_k f\|_\infty &\leq L \|\mathbf{P}_k R_{k+1} f\|_\infty + L \|\mathbf{D}\mathbf{P}_k R_{k+1} f\|_\infty \\ &\leq L \|R_{k+1} f\|_\infty + LK \|DR_{k+1} f\|_\infty \end{aligned} \quad (2.19)$$

In addition,

$$\begin{aligned} [DR_k f]_{Lip} &\leq L (\|R_{k+1} f\|_\infty + K \|DR_{k+1} f\|_\infty \\ &\quad + K u_{k+1} + K \|DR_{k+1} f\|_\infty) \end{aligned} \quad (2.20)$$

where  $u_{k+1} := \|DR_{k+1} f\|_\infty \vee [DR_{k+1} f]_{Lip}$ .

Noticing from (2.19) that also:

$$\|DR_k f\|_\infty \leq L (\|R_{k+1} f\|_\infty + K \|DR_{k+1} f\|_\infty + K u_{k+1} + K \|DR_{k+1} f\|_\infty),$$

we have:

$$\begin{aligned} u_k &\leq L (\|R_{k+1} f\|_\infty + K \|DR_{k+1} f\|_\infty + K u_{k+1} + K \|DR_{k+1} f\|_\infty) \\ &\leq 3LK u_{k+1} + L \|R_{k+1} f\|_\infty. \end{aligned} \quad (2.21)$$

Recursively we conclude the announced result.  $\square$

Now, applying the previous idea (from equations (2.15) and (2.16)) to the sequential filter estimation via quantization, when  $H_n^n f \in \mathcal{C}_{b,Lip}^1$ , a *generic first order scheme* can be designed as follows:

$$\left\{ \begin{array}{l} \widehat{R}_n f(\widehat{X}_n) = H_n^n f(\widehat{X}_n), \\ \widehat{DR}_n f(\widehat{X}_n) = \mathbf{D}H_n^n f(\widehat{X}_n), \\ \widehat{R}_k f(\widehat{X}_k) = g_k(\widehat{X}_k) \mathbb{E}[\widehat{R}_{k+1} f(\widehat{X}_{k+1}) + \langle \widehat{DR}_{k+1} f(\widehat{X}_{k+1}), \Delta_{k+1} \rangle | \widehat{X}_k], \\ 0 \leq k \leq n-1. \end{array} \right. \quad (2.22)$$

and then,  $\widehat{\pi}_n f = \widehat{\mu}_0 \widehat{R}_0 f$ .

In (2.22),  $\widehat{DR}_k f$  is a quantization based estimate for  $DR_k f$ . It needs to be specified to transform the above scheme into an implementable algorithm. In [16], the scheme (2.22) is introduced with no computational considerations concerning  $DR_k f$ . It is shown that under assumptions **H2** and **H1**, the quantization based unnormalized filter converges toward  $\pi_n f$  at a rate  $\sum_{k=1}^n \|\Delta_k\|_2^2$  (instead of  $\sum_{k=1}^n \|\Delta_k\|_2$  in the original zero order scheme from [15]). Our aim is to propose some estimate  $\widehat{DR}_k f$  for  $DR_k f$ , in order to combine computability skills and convergence rate improvement. In this aim, two methods will be exhibited:

- the first one is based on an induction: at each time step  $k$  we evaluate  $\{\widehat{DR}_k, \widehat{R}_k\}$  using  $\{\widehat{DR}_{k+1}, \widehat{R}_{k+1}\}$ . This approach leads to a one step recursive scheme and is investigated in Section 3;
- the second one is based on an integration by parts following an approach developed in [3]: the operator  $\widehat{DR}_k$  is defined as a weighted expectation of  $\widehat{R}_k$ . The scheme constructed by plugging  $\widehat{DR}_k f$  expression in (2.22) leads to a two step recursive scheme, details are investigated in Section 4.

### 3 One step first order iterative scheme

We introduce for this section the following assumption, in the spirit of **H2**, but in fact a bit more restrictive:

**H 2'** For each  $1 \leq k \leq n$ ,  $F_k$  admits a bounded, uniformly Lipschitz derivative with respect to its first variable. Namely,  $\forall x, x' \in \mathbb{R}^d, \forall \varepsilon \in \mathbb{R}^d$ :

$$|\partial_x F_k(x, \varepsilon) - \partial_x F_k(x', \varepsilon)| \leq [\partial_x F_k]_{Lip}^1 |x - x'| \quad \text{and} \quad \|\partial_x F\|_\infty := \max_{1 \leq k \leq n} \|\partial_x F_k\|_\infty < +\infty.$$

**Example 3.1** This assumption is e.g. satisfied by dynamics with an *additive noise*, typically for functions  $F_k : (x, u) \mapsto b_k(x) + \sigma_k u$ , where  $b_k$  is differentiable with bounded Lipschitz continuous derivative and  $\sigma_k \in \mathcal{M}(d, q)$ , or by dynamics where  $F_k$  satisfies:  $F_k(x, u) = b_k(x) + \sigma_k(x)u$ ,  $b_k, \sigma_k$  being differentiable with bounded Lipschitz continuous derivatives, applied to signal innovations  $\varepsilon_k$  with compactly supported pdf.

#### 3.1 Definition of the scheme

In this paragraph, we investigate the recursive approach to estimate  $DR_k$ . Under **H2'**, the probability transitions  $\mathbf{P}_k$  are  $K$ -Lipschitz with  $K = \|\partial_x F\|_\infty$ . Furthermore, the  $\mathbf{P}_k$  are differentiable in the following sense: for every  $f \in \mathcal{C}_{b, Lip}^1$ ,

$$D\mathbf{P}_k f = Q_k Df, \quad k = 0, \dots, n-1, \quad (3.1)$$

where, for every Borel map  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,

$$Q_k \varphi(x) = \mathbb{E}[\partial_x F_{k+1}(X_k, \varepsilon_{k+1})' \varphi(X_{k+1}) | X_k = x], \quad \text{for } x \in \mathbb{R}^d. \quad (3.2)$$

The quantization based estimate for  $D\mathbf{P}_k f$  is then naturally defined by:

$$\hat{Q}_k Df(x_k^i) = \mathbb{E}[\partial_x F_{k+1}(X_k, \varepsilon_{k+1})' Df(\hat{X}_{k+1}) | \hat{X}_k = x_k^i], \quad \text{for } i = 1, \dots, N_k. \quad (3.3)$$

Finally, following equation (2.17) one sets:

$$\widehat{DR}_k f(x_k^i) = Dg_k(x_k^i) \widehat{\mathbf{P}}_k \widehat{R}_{k+1} f(x_k^i) + g_k(x_k^i) \hat{Q}_k Df(x_k^i) \quad (3.4)$$

as a zero order approximation of  $DR_k f$  defined on  $\Gamma_k = \{x_k^1, \dots, x_k^{N_k}\}$ , for any  $k \in \{1, \dots, n-1\}$ .

**Remark 3.1** From a numerical point of view, it would be more natural to use  $\widehat{\mathbf{P}}_k \widehat{R}_{k+1}$  instead of  $\widehat{\mathbf{P}}_k \widehat{\widehat{R}}_{k+1}$ . In fact, the algorithm structure would be less complex. Our choice in (3.4) is motivated on one hand by theoretical need to take a zero order approximation for the differential term estimator. On the other hand, using  $\widehat{\mathbf{P}}_k \widehat{R}_{k+1}$  will introduce distortion terms in both  $\widehat{DR}_k$  and  $\widehat{R}_k$  which generates important numerical instability as emphasized by numerical tests in Figure 5.

Now, plugging (3.4) into the generic first order scheme (2.22) yields the following first order scheme:

**Scheme B:** BACKWARD EXPRESSION

$$\begin{cases} \widehat{R}_n f(\widehat{X}_n) = H_n^n f(\widehat{X}_n), \\ \widehat{DR}_n(\widehat{X}_n) f = DH_n^n f(\widehat{X}_n), \\ \widehat{R}_k f(\widehat{X}_k) = g_k(\widehat{X}_k) \mathbb{E}[\widehat{R}_{k+1} f(\widehat{X}_{k+1}) + \langle \widehat{DR}_{k+1} f(\widehat{X}_{k+1}), \Delta_{k+1} \rangle | \widehat{X}_k], \\ \widehat{DR}_k f(\widehat{X}_k) = Dg_k(\widehat{X}_k) \mathbb{E}[\widehat{R}_{k+1} f(\widehat{X}_{k+1}) | \widehat{X}_k] + g_k(\widehat{X}_k) \widehat{Q}_k \widehat{DR}_{k+1} f(\widehat{X}_k) \\ k = 0, \dots, n-1. \end{cases} \quad (3.5)$$

Note that this scheme is completely computable, as it can be rewritten easily using finite weighted sums. The quantizers  $\Gamma_k$  and the weights - which we call from now on *companion parameters* - can be computed off line and stored in an accessible codebook, so that the only on line computation cost will be the calculus of operators  $\widehat{R}_k$ ,  $\widehat{\widehat{R}}_k$  and  $\widehat{DR}_k$ .

The scheme can be reformulated *in distribution* as follows:

**Scheme B**

$$\begin{cases} \widehat{\widehat{R}}_n f(x_n^i) = H_n^n f(x_n^i), & i = 1, \dots, N_n, \\ \widehat{R}_n f(x_n^i) = H_n^n f(x_n^i), & i = 1, \dots, N_n, \\ \widehat{DR}_n f(x_n^i) = DH_n^n f(x_n^i), & i = 1, \dots, N_n, \\ \widehat{\widehat{R}}_k f(x_k^i) = g_k(x_k^i) \sum_{j=1}^{N_{k+1}} \widehat{\widehat{R}}_{k+1} f(x_{k+1}^j) p_k^{ij}, \\ \widehat{R}_k f(x_k^i) = g_k(x_k^i) \sum_{j=1}^{N_{k+1}} \left( \widehat{R}_{k+1} f(x_{k+1}^j) p_k^{ij} + \langle \widehat{DR}_{k+1} f(x_{k+1}^j), \delta_k^{ij} \rangle \right) \\ \widehat{DR}_k f(x_k^i) = Dg_k(x_k^i) \sum_{j=1}^{N_{k+1}} \widehat{\widehat{R}}_{k+1} f(x_{k+1}^j) p_k^{ij} + g_k(x_k^i) \sum_{j=1}^{N_{k+1}} \gamma_k^{ij} \widehat{DR}_{k+1} f(x_{k+1}^j), \\ i = 1, \dots, N_k, \quad 0 \leq k < n, \end{cases} \quad (3.6)$$

where the companion parameters,  $p_k^{ij}$ ,  $\gamma_k^{ij}$ , and  $\delta_k^{ij}$  are defined by:

$$\begin{aligned} p_k^{ij} &= \mathbb{E}[\mathbf{1}_{\{\widehat{X}_{k+1}=x_{k+1}^j\}} | \widehat{X}_k = x_k^i] \in \mathbb{R}, \\ \gamma_k^{ij} &= \mathbb{E}[\partial_x F_k(X_k, \varepsilon_{k+1})' \mathbf{1}_{\{\widehat{X}_{k+1}=x_{k+1}^j\}} | \widehat{X}_k = x_k^i] \in \mathcal{M}_d(\mathbb{R}) \\ \delta_k^{ij} &= \mathbb{E}[\Delta_{k+1} \mathbf{1}_{\{\widehat{X}_{k+1}=x_{k+1}^j\}} | \widehat{X}_k = x_k^i], \\ &= \mathbb{E}[(X_{k+1} - x_{k+1}^j) \mathbf{1}_{\{\widehat{X}_{k+1}=x_{k+1}^j\}} | \widehat{X}_k = x_k^i] \in \mathbb{R}^d. \end{aligned} \quad (3.7)$$

**FORWARD EXPRESSION OF SCHEME B**

For applications, it is crucial in terms of computational efficiency, to rewrite the scheme in a forward way. This allows us to compute costlessly intermediate estimations of  $\pi_k f$ ,

$1 \leq k \leq n-1$ , and to use different test functions  $f$  without recomputing the hole scheme. This forward form can be established as follows: one first checks that at each  $0 \leq k \leq n-1$ ,

the vector  $\begin{bmatrix} \hat{R}_k \\ \widehat{DR}_k \\ \hat{R}_k \end{bmatrix}$  satisfies the following one step induction:

$$\begin{bmatrix} \hat{R}_k \\ \widehat{DR}_k \\ \hat{R}_k \end{bmatrix} = \hat{\mathcal{H}}_k \begin{bmatrix} \hat{R}_{k+1} \\ \widehat{DR}_{k+1} \\ \hat{R}_{k+1} \end{bmatrix}, \quad (3.8)$$

where  $\hat{\mathcal{H}}_k$  is a lower triangular operator matrix defined by:  $\hat{\mathcal{H}}_k = \begin{pmatrix} \hat{H}_k^1 & 0 & 0 \\ \hat{H}_k^2 & \hat{H}_k^3 & 0 \\ 0 & \hat{H}_k^4 & \hat{H}_k^1 \end{pmatrix}$ , with

for  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,

$$\begin{aligned} \hat{H}_0^1 f(x) &= \mathbb{E}[f(\hat{X}_1) | \hat{X}_0 = x], \\ \hat{H}_0^2 f(x) &= 0 \in \mathbb{R}^d, \\ \hat{H}_0^3 \varphi(x) &= \mathbb{E}[\partial_x F_1(x, \varepsilon_1)' \varphi(\hat{X}_1) | \hat{X}_0 = x], \\ \hat{H}_0^4 \varphi(x) &= \mathbb{E}[\langle \varphi(\hat{X}_1), \Delta_1 \rangle | \hat{X}_0 = x], \end{aligned}$$

and for every  $1 \leq k < n$

$$\begin{aligned} \hat{H}_k^1 f(\hat{X}_k) &= g_k(\hat{X}_k) \mathbb{E}[f(\hat{X}_{k+1}) | \hat{X}_k], \\ \hat{H}_k^2 f(\hat{X}_k) &= Dg_k(\hat{X}_k) \mathbb{E}[f(\hat{X}_{k+1}) | \hat{X}_k], \\ \hat{H}_k^3 \varphi(\hat{X}_k) &= g_k(\hat{X}_k) \mathbb{E}[\partial_x F_{k+1}(X_k, \varepsilon_{k+1})' \varphi(\hat{X}_{k+1}) | \hat{X}_k], \\ \hat{H}_k^4 \varphi(\hat{X}_k) &= g_k(\hat{X}_k) \mathbb{E}[\langle \varphi(\hat{X}_{k+1}), \Delta_{k+1} \rangle | \hat{X}_k]. \end{aligned}$$

Notice that here  $\hat{H}_k^1 = \hat{H}_k$ .

Then, one can see from (3.8) that:

$$\begin{bmatrix} \hat{R}_0 \\ \widehat{DR}_0 \\ \hat{R}_0 \end{bmatrix} = \hat{\mathcal{H}}_0 \circ \hat{\mathcal{H}}_1 \circ \dots \circ \hat{\mathcal{H}}_{n-1} \begin{bmatrix} H_n^n \\ DH_n^n \\ H_n^n \end{bmatrix}.$$

Setting  $\hat{\mathcal{U}}_k = \hat{\mu}_0 \circ \hat{\mathcal{H}}_0 \circ \dots \circ \hat{\mathcal{H}}_k$ , the forward scheme satisfies the following recursive formula:

$$\hat{\mathcal{U}}_0 = \hat{\mu}_0 \hat{\mathcal{H}}_0 \quad \text{and} \quad \hat{\mathcal{U}}_k = \hat{\mathcal{U}}_{k-1} \hat{\mathcal{H}}_k \quad k = 1, \dots, n-1,$$

so that  $\hat{\pi}_n f = \langle \hat{\mathcal{U}}_{n-1} \begin{bmatrix} H_n^n f \\ DH_n^n f \\ H_n^n f \end{bmatrix}, e_3 \rangle$ .

### 3.2 Error bounds

The main result of this section is to establish a convergence result for scheme **B** better than the zero scheme rate. We recall that here,  $(\hat{X}_k)$  is a marginal, stationary  $(N_k)$ -quantization of  $(X_k)$ .

**Theorem 3.1** Assume **H1** and **H2'** and let  $f$  satisfying  $H_n f \in \mathcal{C}_{b,Lip}^1$ . Then, there exists a sequence of positive real constants  $(M_j^n)_{0 \leq j \leq n}$  such that:

$$|\pi_n f - \hat{\pi}_n f| \leq \sum_{j=0}^n M_j^n \|\Delta_j\|_{2p}^2$$

$$\text{with } M_j^n \leq \alpha(p, f) \frac{n+5}{2} L^n \left( \frac{(LK)^{j+1}-1}{LK-1} \right) \left( \frac{(3K)^{n-j+1}-1}{3K-1} \right) \left( \frac{(L)^{j+1}-1}{L-1} \right).$$

The key to prove the above error bound is to rely on the *backward* form of the scheme **B** (see 3.6). The main technical step is to produce some upper error bounds for the differential term approximation, namely  $\hat{A}_k = DR_k f(\hat{X}_k) - \widehat{DR}_k f(\hat{X}_k)$ .

The proof of the first lemma below is left to the reader:

**Lemma 3.1** For any  $\varphi$  bounded Lipschitz continuous,  $Q_k \varphi : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is Lipschitz and

$$\|Q_k \varphi\|_{Lip} \leq \|\varphi\|_\infty [\partial_x F_{k+1}]_{Lip}^1 + \|\partial_x F\|_\infty^2 [\varphi]_{Lip}.$$

Then, the error bounds for  $\|\hat{A}_k\|_p$  are given in the lemma:

**Lemma 3.2** For  $f$  satisfying  $H_n f \in \mathcal{C}_{b,Lip}^1$ , there exists a non negative real sequence  $(D_j^{k,n})_{0 \leq k \leq j \leq n}$  such that:

$$\|\hat{A}_k\|_p \leq \sum_{j=k}^n D_j^{k,n} \|\Delta_j\|_p$$

$$\text{where } D_j^{k,n} \leq \alpha(p, f) L^{n-k} \left( \frac{(LK)^{j-k+1}-1}{LK-1} \right) \left( \frac{K^{n-j+1}-1}{K-1} \right).$$

**Proof.**

From equations (2.17), (3.4) and (3.1):

$$\begin{aligned} \hat{A}_k &= Dg_k(\hat{X}_k) \mathbf{P}_k R_{k+1} f(\hat{X}_k) + g_k(\hat{X}_k) Q_k DR_{k+1} f(\hat{X}_k) \\ &\quad - Dg_k(\hat{X}_k) \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k] - g_k(\hat{X}_k) \hat{Q}_k \widehat{DR}_{k+1} f(\hat{X}_k) \\ &= Dg_k(\hat{X}_k) \left( \mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k] \right) \\ &\quad + Dg_k(\hat{X}_k) \left( \mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k] \right) \\ &\quad + g_k(\hat{X}_k) \left( Q_k DR_{k+1} f(\hat{X}_k) - \mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k] \right) \\ &\quad + g_k(\hat{X}_k) \left( \mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k] - \hat{Q}_k \widehat{DR}_{k+1} f(\hat{X}_k) \right) \end{aligned}$$

Then, using **H1**, one gets:

$$\begin{aligned} \|\hat{A}_k\|_p &\leq L \|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k]\|_p \\ &\quad + L \|\mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p \\ &\quad + L \|Q_k DR_{k+1} f(\hat{X}_k) - \mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k]\|_p \\ &\quad + L \|\mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k] - \hat{Q}_k \widehat{DR}_{k+1} f(\hat{X}_k)\|_p. \end{aligned} \tag{3.9}$$

Now, the  $L^p$ -contraction property of conditional expectation implies that:

$$\begin{aligned} \|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k]\|_p &\leq \|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbf{P}_k R_{k+1} f(X_k)\|_p \\ &\leq [\mathbf{P}_k R_{k+1} f]_{Lip} \|\Delta_k\|_p \\ &\leq K \|DR_{k+1} f\|_\infty \|\Delta_k\|_p. \end{aligned}$$

For the second term in the right handside of inequality (3.9), we will use on one hand the chaining rule for conditional expectation (see equation (2.15)) and on the other hand its  $L^p$ -contraction property, to write:

$$\begin{aligned} \|\mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p &= \|\mathbb{E}[R_{k+1} f(X_{k+1}) | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p \\ &\leq \|R_{k+1} f(X_{k+1}) - \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p \end{aligned}$$

which implies, by Theorem 2.1:

$$\|\mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p \leq \sum_{j=k+1}^n C_j^{k+1,n} \|\Delta_j\|_p.$$

The same arguments on conditional expectations give:

$$\|Q_k DR_{k+1} f(\hat{X}_k) - \mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k]\|_p \leq [Q_k DR_{k+1} f]_{Lip} \|\Delta_k\|_p,$$

which by Lemma 3.1 writes:

$$\begin{aligned} \|Q_k DR_{k+1} f(\hat{X}_k) - \mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k]\|_p &\leq \\ &(\|DR_{k+1} f\|_\infty [\partial_x F_{k+1}]_{Lip}^1 + \|\partial_x F\|_\infty^2 [DR_{k+1} f]_{Lip}) \|\Delta_k\|_p \end{aligned}$$

since  $DR_{k+1}$  is bounded Lipschitz by Proposition 2.2.

Then, using the definition of  $\hat{Q}_k$  yields:

$$\begin{aligned} \|\mathbb{E}[Q_k DR_{k+1} f(X_k) | \hat{X}_k] - \hat{Q}_k \widehat{DR}_{k+1} f(\hat{X}_k)\|_p &\leq \|(\partial_x F_{k+1}(X_k, \varepsilon_{k+1}))' (DR_{k+1} f(X_{k+1}) - \widehat{DR}_{k+1} f(\hat{X}_{k+1}))\|_p \\ &\leq \|\partial_x F_{k+1}\|_\infty \left( \|\hat{A}_{k+1}\|_p + \|DR_{k+1} f(X_{k+1}) - DR_{k+1} f(\hat{X}_{k+1})\|_p \right) \\ &\leq \|\partial_x F_{k+1}\|_\infty \left( \|\hat{A}_{k+1}\|_p + [DR_{k+1} f]_{Lip} \|\Delta_{k+1}\|_p \right). \end{aligned}$$

Finally, using  $\|\partial_x F\|_\infty = K$  and Proposition 2.2, we derive:

$$\begin{aligned} \|\hat{A}_k\|_p &\leq L \left( [\partial_x F_{k+1}]_{Lip}^1 \|DR_{k+1} f\|_\infty + K^2 [DR_{k+1} f]_{Lip} + K \|DR_{k+1} f\|_\infty \right) \|\Delta_k\|_p \\ &\quad + L \left( C_{k+1}^{k+1,n} + K ([DR_{k+1} f]_{Lip}) \right) \|\Delta_{k+1}\|_p \\ &\quad + L \sum_{j=k+2}^n C_j^{k+1,n} \|\Delta_j\|_p + LK \|\hat{A}_{k+1}\|_p. \end{aligned} \tag{3.10}$$

The required result follows from a backward induction on  $k$ . See [18] for explicit upper bounds.  $\square$

**Proof of Theorem 3.1.**

Let  $V_k f$  denote the intermediate estimation error when considering the previous first order approximation scheme  $\mathbf{B}$  in its backward form :  $V_k f := \mathbb{E}[R_k f(X_k)|\hat{X}_k] - \hat{R}_k f(\hat{X}_k)$ .

Using triangular inequalities, we isolate three error sources in  $|V_k f|$ . If we set:

$$\begin{aligned}\bar{R}_k f(\hat{X}_k) &= g_k(\hat{X}_k)\mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k)|\hat{X}_k], \\ &= g_k(\hat{X}_k)\mathbb{E}[R_{k+1} f(X_{k+1})|\hat{X}_k],\end{aligned}$$

then we have:

$$\begin{aligned}|V_k f| &\leq |\mathbb{E}[R_k f(X_k)|\hat{X}_k] - R_k f(\hat{X}_k)| + |R_k f(\hat{X}_k) - \bar{R}_k f(\hat{X}_k)| \\ &\quad + |\bar{R}_k f(\hat{X}_k) - \hat{R}_k f(\hat{X}_k)|.\end{aligned}\tag{3.11}$$

Using a first order Taylor expansion, there exists  $\hat{\zeta}_k^1 \in (\hat{X}_k, X_k)$  such that

$$\begin{aligned}\mathbb{E}[R_k f(X_k)|\hat{X}_k] &= \mathbb{E}[R_k f(\hat{X}_k) + \langle DR_k f(\hat{\zeta}_k^1), \Delta_k \rangle | \hat{X}_k] \\ &= \mathbb{E}[R_k f(\hat{X}_k) + \langle DR_k f(\hat{X}_k), \Delta_k \rangle + \langle DR_k f(\hat{\zeta}_k^1) - DR_k f(\hat{X}_k), \Delta_k \rangle | \hat{X}_k]\end{aligned}$$

$\hat{X}_k$  being a stationary quantization of  $X_k$ , one derives from Proposition 2.1 that:

$$\mathbb{E}[\langle DR_k f(\hat{X}_k), \Delta_k \rangle | \hat{X}_k] = \langle DR_k f(\hat{X}_k), \mathbb{E}[\Delta_k | \hat{X}_k] \rangle = 0.$$

Then,

$$\begin{aligned}|\mathbb{E}[R_k f(X_k)|\hat{X}_k] - R_k f(\hat{X}_k)| &= |\mathbb{E}[\langle DR_k f(\hat{\zeta}_k^1) - DR_k f(\hat{X}_k), \Delta_k \rangle | \hat{X}_k]| \\ &\leq \mathbb{E}[|DR_k f(\hat{\zeta}_k^1) - DR_k f(\hat{X}_k)| |\Delta_k| | \hat{X}_k] \\ &\leq [DR_k f]_{Lip} \mathbb{E}[|\hat{X}_k - \hat{\zeta}_k^1| |\Delta_k| | \hat{X}_k] \\ &\leq [DR_k f]_{Lip} \mathbb{E}[|\Delta_k|^2 | \hat{X}_k].\end{aligned}\tag{3.12}$$

By Taylor expansion of  $\mathbf{P}_k R_{k+1} f$ , we analogically find  $\hat{\zeta}_k^2 \in (\hat{X}_k, X_k)$  such that:

$$\begin{aligned}\bar{R}_k f(\hat{X}_k) &= g_k(\hat{X}_k) \left( \mathbf{P}_k R_{k+1} f(\hat{X}_k) + \langle D\mathbf{P}_k R_{k+1} f(\hat{X}_k), \mathbb{E}[\Delta_k | \hat{X}_k] \rangle \right. \\ &\quad \left. + \mathbb{E}[\langle D\mathbf{P}_k R_{k+1} f(\hat{\zeta}_k^2) - D\mathbf{P}_k R_{k+1} f(\hat{X}_k), \Delta_k \rangle | \hat{X}_k] \right) \\ R_k f(\hat{X}_k) - \bar{R}_k f(\hat{X}_k) &= g_k(\hat{X}_k) \mathbb{E}[\langle D\mathbf{P}_k R_{k+1} f(\hat{\zeta}_k^2) - D\mathbf{P}_k R_{k+1} f(\hat{X}_k), \Delta_k \rangle | \hat{X}_k] \\ \text{Hence, } |R_k f(\hat{X}_k) - \bar{R}_k f(\hat{X}_k)| &\leq L [D\mathbf{P}_k R_{k+1} f]_{Lip} \mathbb{E}[|\Delta_k|^2 | \hat{X}_k] \\ &\leq LK ([DR_{k+1} f]_{Lip} \vee \|DR_{k+1} f\|_\infty) \mathbb{E}[|\Delta_k|^2 | \hat{X}_k]\end{aligned}\tag{3.13}$$

For the last term in the right handside of inequality (3.11), we have:

$$\begin{aligned}|\bar{R}_k f(\hat{X}_k) - \hat{R}_k f(\hat{X}_k)| &= \left| g_k(\hat{X}_k) \left( \mathbb{E}[R_{k+1} f(X_{k+1})|\hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1})|\hat{X}_k] \right. \right. \\ &\quad \left. \left. - \mathbb{E}[\langle \widehat{DR}_{k+1} f(\hat{X}_{k+1}), \Delta_{k+1} \rangle | \hat{X}_k] \right) \right| \\ &\leq L \left| \mathbb{E} \left[ R_{k+1} f(X_{k+1}) - \mathbb{E}[R_{k+1} f(X_{k+1})|\hat{X}_{k+1}] \mid \hat{X}_k \right] \right. \\ &\quad \left. - \mathbb{E}[\langle \widehat{DR}_{k+1} f(\hat{X}_{k+1}), \Delta_{k+1} \rangle | \hat{X}_k] \right| \\ &\quad + L \left| \mathbb{E} \left[ \mathbb{E}[R_{k+1} f(X_{k+1})|\hat{X}_{k+1}] - \hat{R}_{k+1} f(\hat{X}_{k+1}) \mid \hat{X}_k \right] \right|\end{aligned}$$



Furthermore, there exists  $\hat{\zeta}_{k+1}^3 \in (\hat{X}_{k+1}, X_{k+1})$  such that

$$\begin{aligned} R_{k+1}f(X_{k+1}) &= R_{k+1}f(\hat{X}_{k+1}) + \langle DR_{k+1}f(\hat{X}_{k+1}), \Delta_{k+1} \rangle \\ &\quad + \langle DR_{k+1}f(\hat{\zeta}_{k+1}^3) - DR_{k+1}f(\hat{X}_{k+1}), \Delta_{k+1} \rangle \\ \mathbb{E}[R_{k+1}f(X_{k+1})|\hat{X}_{k+1}] &= R_{k+1}f(\hat{X}_{k+1}) + \mathbb{E}[\langle DR_{k+1}f(\hat{\zeta}_{k+1}^3) - DR_{k+1}f(\hat{X}_{k+1}), \Delta_{k+1} \rangle|\hat{X}_{k+1}] \end{aligned}$$

Consequently:

$$\begin{aligned} |\bar{R}_k f(\hat{X}_k) - \hat{R}_k f(\hat{X}_k)| &\leq L|\mathbb{E}[V_{k+1}f|\hat{X}_k]| \tag{3.14} \\ &\quad + L|\mathbb{E}[\langle (DR_{k+1}f(\hat{X}_{k+1}) - \widehat{DR}_{k+1}f(\hat{X}_{k+1})), \Delta_{k+1} \rangle|\hat{X}_k]| \\ &\quad + L[DR_{k+1}f]_{Lip} \left( \mathbb{E}[|\Delta_{k+1}|^2|\hat{X}_k] + \mathbb{E} \left[ \mathbb{E}[|\Delta_{k+1}|^2|\hat{X}_{k+1}] \middle| \hat{X}_k \right] \right) \tag{3.15} \end{aligned}$$

Finally combining previous inequalities (3.12), (3.13), (3.14), we obtain by using  $L^p$ -contraction property of conditional expectation:

$$\begin{aligned} \|V_k f\|_p &\leq ([DR_k f]_{Lip} + LKu_{k+1}) \|\Delta_k\|_{2p}^2 \\ &\quad + 2L[DR_{k+1}f]_{Lip} \|\Delta_{k+1}\|_{2p}^2 + L\|\langle \hat{A}_{k+1}, \Delta_{k+1} \rangle\|_p \tag{3.16} \\ &\quad + L\|V_{k+1}f\|_p. \end{aligned}$$

Applying Holder inequality combined to Lemma 3.2 to the term  $\|\langle \hat{A}_{k+1}, \Delta_{k+1} \rangle\|_p$ , we have:

$$\begin{aligned} \|\langle \hat{A}_{k+1}, \Delta_{k+1} \rangle\|_p &\leq \|\hat{A}_{k+1}\| \|\Delta_{k+1}\|_p \\ &\leq \|\hat{A}_{k+1}\|_{2p} \|\Delta_{k+1}\|_{2p} \\ &\leq \sum_{j=k+1}^n D_j^{k+1,n} \|\Delta_j\|_{2p} \|\Delta_{k+1}\|_{2p} \\ &\leq \frac{1}{2} \sum_{j=k+1}^n D_j^{k+1,n} (\|\Delta_j\|_{2p}^2 + \|\Delta_{k+1}\|_{2p}^2) \tag{3.17} \end{aligned}$$

Plugging (3.17) into (3.16) yields:

$$\begin{aligned} \|V_k f\|_p &\leq ([DR_k f]_{Lip} + LKu_{k+1}) \|\Delta_k\|_{2p}^2 \\ &\quad + L \left( 2[DR_{k+1}f]_{Lip} + \frac{1}{2} \sum_{j=k+1}^n D_j^{k+1,n} \right) \|\Delta_{k+1}\|_{2p}^2 \tag{3.18} \\ &\quad + \frac{1}{2} L \sum_{j=k+1}^n D_j^{k+1,n} \|\Delta_j\|_{2p}^2 + L\|V_{k+1}f\|_p. \end{aligned}$$

Then, by induction taking  $k = 0$  and writing  $|\pi_n f - \hat{\pi}_n f| \leq \|V_0 f\|_p$  we derive the required result. See [18] for further details.  $\square$

Theorem 3.1, with  $\|\Delta_k\|_{2p}^2 = O(\|\Delta_k\|_{2p})$ , shows that the scheme **B** succeeds to embetter the zero-order convergence rate. The forthcoming section has been motivated by our wish to relax **H2'** and to preserve the convergence rate improvement.

## 4 Two step iterative first order scheme

To construct this second first order scheme, the idea is to represent  $\mathbf{DP}_k R_{k+1}f$  as a weighted conditional expectation of  $R_{k+1}f$  i.e.

$$\mathbf{DP}_k R_{k+1}f(x) = \mathbb{E}[R_{k+1}f(X_{k+1}) \times \text{Weight} | X_k = x],$$

and then to quantize this representation formula. This is achieved classically by the mean of an integration by parts formula.

Note that in all this section, we will assume  $q = d$ . Furthermore,  $F_k$  will be supposed to be differentiable.

### 4.1 Integration by parts formula

For notational convenience, we will temporarily drop the  $k$  indices in the notations of  $X_k$ ,  $F_k$  and  $\mathbf{P}_k$ . We will also temporarily assume  $f \in \mathcal{C}_b^1$ .

We start by a transformation of the problem, via differentiation. For that, we need first to assume the following:

**H 3**  $\forall 0 \leq k \leq n, \exists c_k > 0$  such that for any  $x \in \mathbb{R}^d$  and  $\varepsilon \in \mathbb{R}^d$ :

$$(\partial_\varepsilon F_k(x, \varepsilon))(\partial_\varepsilon F_k(x, \varepsilon))' \geq c_k I_d.$$

We have then, for any  $x, \varepsilon \in \mathbb{R}^d$ :

$$\begin{aligned} \partial_x(f \circ F)(x, \varepsilon) &= \partial_x F(x, \varepsilon)' (Df) \circ F(x, \varepsilon), \\ \partial_\varepsilon(f \circ F)(x, \varepsilon) &= \partial_\varepsilon F(x, \varepsilon)' (Df) \circ F(x, \varepsilon). \end{aligned}$$

Assuming **H3** yields  $\partial_x(f \circ F)(x, \varepsilon) = \mathcal{G}_x(\varepsilon) \partial_\varepsilon(f \circ F)(x, \varepsilon)$ , where:

$$\begin{aligned} \mathcal{G}_x : \mathbb{R}^d &\rightarrow \mathcal{M}_d(\mathbb{R}) \\ \varepsilon &\mapsto (\partial_\varepsilon F(x, \varepsilon)^{-1} \partial_x F(x, \varepsilon))'. \end{aligned}$$

Now, in order to allow a differentiation under the integral sign and then apply integration by parts, we will assume the following technical hypothesis:

**H 4** Assume that signal innovations  $\varepsilon_k$  distribution is absolutely continuous toward Lebesgue measure, with a differentiable density  $\mathbf{p}$  satisfying for all  $x \in \mathbb{R}^d$ ,

$$\int_{\mathbb{R}^d} |\partial_x F(x, \varepsilon)| \mathbf{p}(\varepsilon) d\varepsilon < +\infty \quad \text{and} \quad \lim_{|\varepsilon| \rightarrow +\infty} \mathcal{G}_x(\varepsilon) \mathbf{p}(\varepsilon) = 0.$$

Then, the  $i$ -th component of  $\mathbf{DP}f(x)$  for a given index  $1 \leq i \leq d$  reads:

$$\frac{\partial \mathbf{P}f}{\partial x^i}(x) = \int_{\mathbb{R}^d} \langle \mathcal{G}_x^i(\varepsilon), \partial_\varepsilon(f \circ F)(x, \varepsilon) \rangle \mathbf{p}(\varepsilon) d\varepsilon \quad (4.1)$$

$$\begin{aligned} \text{where: } \mathcal{G}_x^i : \mathbb{R}^d &\rightarrow \mathbb{R}^d \\ \varepsilon &\mapsto (\mathcal{G}_x(\varepsilon))' e_i. \end{aligned}$$

Furthermore, performing an integration by parts formula on (4.1), and taking in account **H4** yields:

$$\frac{\partial \mathbf{P}f}{\partial x^i}(x) = - \int_{\mathbb{R}^q} (f \circ F(x, \varepsilon) + C(x)) \Psi^i(x, \varepsilon) \mathbf{p}(\varepsilon) d\varepsilon \quad (4.2)$$

$$\begin{aligned} \text{where: } \Psi^i : \mathbb{R}^d \times \mathbb{R}^d &\rightarrow \mathbb{R} \\ (x, \varepsilon) &\mapsto \operatorname{div} \mathcal{G}_x^i(\varepsilon) + \frac{1}{\mathbf{p}(\varepsilon)} \langle \mathcal{G}_x^i(\varepsilon), \mathbf{D}\mathbf{p}(\varepsilon) \rangle. \end{aligned}$$

Finally, defining the weight vector  $\Psi(x, \varepsilon) := (\Psi^i(x, \varepsilon))_{0 \leq i \leq d}$ , we obtain the generalization of equation (4.2):  $\mathbf{D}(\mathbf{P}_k f)(x) = -\mathbb{E}[(f(F_{k+1}(x, \varepsilon_{k+1})) + C^k(x)) \Psi_k(x, \varepsilon_{k+1})]$ .

In a Monte Carlo method context, the constant  $C^k$  is tuned in order to minimize the variance of a probabilistic estimator of  $\mathbf{D}(\mathbf{P}_k f)(x)$ . In our quantization context, as the variance problem does not occur, a natural value for  $C^k$  would be zero. It is at least the choice that minimizes computation cost and provides satisfactory numerical results (see [4] for a discussion about  $C^k$  for an American option pricing problem).

From now on, we will take  $C^k = 0$ .

## 4.2 Numerical scheme

Consider now a test function  $f$  satisfying  $H_n^n f \in \mathcal{C}_{b, Lip}^1$ . Then, according to Proposition 2.2,  $R_k f \in \mathcal{C}_{b, Lip}^1$ . Using results of the previous paragraph, we can write, for each  $0 \leq k \leq n-1$ :  $\mathbf{D}\mathbf{P}_k R_{k+1} f(x) = -\mathbb{E}[R_{k+1} f(X_{k+1}) \Psi_{k+1}(x, \varepsilon_{k+1}) | X_k = x]$ .

So,  $(\hat{X}_k)$  still being a stationary marginal  $(N_k)$ -quantization of  $(X_k)$ , an approximation of  $\mathbf{D}R_k f$  would be:

$$\widehat{\mathbf{D}R}_k f(\hat{X}_k) = \mathbf{D}g_k(\hat{X}_k) \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k] - g_k(\hat{X}_k) \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) \Psi_k(X_k, \varepsilon_{k+1}) | \hat{X}_k].$$

If one replaces this expression in (2.22), it results in the following two step recursive scheme formulated in a backward way:

**Scheme A** BACKWARD FORMULATION

$$\left\{ \begin{aligned} \hat{R}_n f(\hat{X}_n) &= H_n^n f(\hat{X}_n), \\ \hat{R}_{n-1} f(\hat{X}_{n-1}) &= g_{n-1}(\hat{X}_{n-1}) \mathbb{E}[H_n^n f(\hat{X}_n) + \langle \mathbf{D}H_n^n f(\hat{X}_n), \Delta_n \rangle | \hat{X}_{n-1}], \\ \hat{R}_k f(\hat{X}_k) &= g_k(\hat{X}_k) \hat{\mathbf{P}}_k \hat{R}_{k+1} f(\hat{X}_{k+1}) + g_k(\hat{X}_k) \times \\ &\quad \left( \mathbb{E}[\langle \mathbf{D}g_{k+1}(\hat{X}_{k+1}) \hat{\mathbf{P}}_{k+1} \hat{R}_{k+2} f(\hat{X}_{k+2}), \Delta_{k+1} \rangle | \hat{X}_k] - \mathbb{E}[\langle g_{k+1}(\hat{X}_{k+1}) \times \right. \\ &\quad \left. \mathbb{E}[\hat{R}_{k+2} f(\hat{X}_{k+2}) \Psi_{k+1}(X_{k+1}, \varepsilon_{k+2}) | \hat{X}_{k+1}], \Delta_{k+1} \rangle | \hat{X}_k] \right), \\ &0 \leq k \leq n-2. \end{aligned} \right. \quad (4.3)$$

This scheme **A** can be rewritten *in distribution* using finite weighted sums. As for the previous scheme, the weights are to be computed simultaneously with the optimal quantizers. Consequently, the implemented algorithm reads as follows:

### Scheme A

$$\left\{ \begin{array}{l} \hat{R}_n f(x_n^i) = H_n^n f(x_n^i), \quad i = 1, \dots, N_n, \\ \hat{R}_n f(x_n^i) = H_n^n f(x_n^i), \quad i = 1, \dots, N_n, \\ \hat{R}_{n-1} f(x_{n-1}^i) = g_{n-1}(x_{n-1}^i) \sum_{j=1}^{N_n} \left( H_n^n f(x_n^j) p_{n-1}^{ij} + \langle \text{D}H_n^n f(x_n^j), \delta_{n-1}^{ij} \rangle \right), \\ \quad i = 1, \dots, N_{n-1}, \\ \hat{R}_k f(x_k^i) = g_k(x_k^i) \sum_{j=1}^{N_{k+1}} \hat{R}_{k+1} f(x_{k+1}^j) p_k^{ij}, \quad i = 1, \dots, N_k, \quad 0 \leq k < n, \\ \hat{R}_k f(x_k^i) = g_k(x_k^i) \sum_{j=1}^{N_{k+1}} \hat{R}_{k+1} f(x_{k+1}^j) p_k^{ij} + g_k(x_k^i) \times \\ \quad \sum_{j=1}^{N_{k+1}} \sum_{l=1}^{N_{k+2}} \left( \hat{R}_{k+2} f(x_{k+2}^l) p_{k+1}^{jl} \langle \text{D}g_k(x_{k+1}^j), \delta_k^{ij} \rangle \right. \\ \quad \left. - g_{k+1}(x_{k+1}^j) \hat{R}_{k+2} f(x_{k+2}^l) \langle \lambda_{k+1}^{jl}, \delta_k^{ij} \rangle \right), \\ \quad i = 1, \dots, N_k, \quad 0 \leq k \leq n-2 \end{array} \right. \quad (4.4)$$

where the companion parameters,  $p_k^{ij}$ ,  $\lambda_k^{ij}$ , and  $\delta_k^{ij}$  are defined by:

$$p_k^{ij} = \mathbb{E}[\mathbf{1}_{\{\hat{X}_{k+1}=x_{k+1}^j\}} | \hat{X}_k = x_k^i] \in \mathbb{R}, \quad (4.5)$$

$$\lambda_k^{ij} = \mathbb{E}[\Psi_k(X_k, \varepsilon_{k+1}) \mathbf{1}_{\{\hat{X}_{k+1}=x_{k+1}^j\}} | \hat{X}_k = x_k^i] \in \mathbb{R}^d, \quad (4.6)$$

$$\begin{aligned} \delta_k^{ij} &= \mathbb{E}[\Delta_{k+1} \mathbf{1}_{\{\hat{X}_{k+1}=x_{k+1}^j\}} | \hat{X}_k = x_k^i], \\ &= \mathbb{E}[(X_{k+1} - x_{k+1}^j) \mathbf{1}_{\{\hat{X}_{k+1}=x_{k+1}^j\}} | \hat{X}_k = x_k^i] \in \mathbb{R}^d. \end{aligned} \quad (4.7)$$

It is important to recall, that the interest of such an approach lies in the possibility of carrying out the computation of all the above companion parameters off line. Once the state equations are fixed and the noise distribution is simulatable, the quantizers and companion parameters can be kept off line. On line computation cost will then be reduced to the sequential determination of  $\hat{R}_k$  and  $\hat{R}_k$  on each grid. Compared to the previous case, scheme **A** is more demanding in on line memory capacity, as it involves two step computations, but, it is worth noting that the new companion parameters  $\lambda_k^{ij}$  are of lower dimension, which compensates the two step recursion effect while considering the algorithm complexity or the storage capacity dedicated to codebooks.

Here also, as for scheme **B**, we can see that this backward definition can be rewritten in a forward form. For  $0 \leq k \leq n$ , let  $\hat{H}_k$  be the operator defined on any function  $f : \Gamma_{k+2} \rightarrow \mathbb{R}$ , such that:

$$\begin{aligned} \hat{H}_k f(x_k^i) &= g_k(x_k^i) \mathbb{E}[\langle \mathbb{E}[f(\hat{X}_{k+2}) | \hat{X}_{k+1}] \text{D}g_{k+1}(\hat{X}_{k+1}) \\ &\quad - g_{k+1}(\hat{X}_{k+1}) \mathbb{E}[f(\hat{X}_{k+2}) \Psi_{k+1}(X_{k+1}, \varepsilon_{k+2}) | \hat{X}_{k+1}], \Delta_{k+1} \rangle | \hat{X}_k = x_k^i]. \end{aligned}$$

For a time step  $0 \leq k \leq n-2$ , we have the following one step transition system:

$$\left\{ \begin{array}{l} \hat{R}_k = \hat{H}_k \hat{R}_{k+1}, \\ \hat{R}_k = \hat{H}_k \hat{R}_{k+1} + \hat{H}_k \hat{R}_{k+2}. \end{array} \right. \quad (4.8)$$

Introducing  $\hat{U}_k$  in addition to  $\hat{U}_k$  we can define the following forward scheme:

**Scheme A:** FORWARD EXPRESSION

$$\left\{ \begin{array}{l} \hat{U}_0 = \hat{\mu}_0 \circ \hat{H}_0, \\ \hat{U}_2 = \hat{\mu}_0 \circ \hat{H}_0, \\ \text{for any } 0 \leq k \leq n-3, \\ \hat{U}_{k+1} = \hat{U}_k \circ \hat{H}_{k+1}, \\ \hat{U}_{k+3} = \hat{U}_{k+2} \circ \hat{H}_{k+2} + \hat{U}_k \circ \hat{H}_{k+1}. \end{array} \right. \quad (4.9)$$

Finally, given the final conditions:

$$\left\{ \begin{array}{l} \hat{R}_n f(\hat{X}_n) = H_n^n f(\hat{X}_n), \\ \hat{R}_n(\hat{X}_n) = H_n^n(\hat{X}_n), \\ \hat{R}_{n-1}(\hat{X}_{n-1}) = g_{n-1}(\hat{X}_{n-1}) \mathbb{E}[H_n^n f(\hat{X}_n) + \langle DH_n^n f(\hat{X}_n), \Delta_n \rangle | \hat{X}_{n-1}], \end{array} \right.$$

we have for any  $n > 1$ ,

$$\left\{ \begin{array}{l} \hat{\mu}_0 \hat{R}_0 = \hat{U}_{n-1} \circ H_n^n = \hat{\pi}_n, \\ \hat{\mu}_0 \hat{R}_0 = \hat{U}_{n-2} \hat{R}_{n-1} + \hat{U}_n H_n^n = \hat{\pi}_n. \end{array} \right.$$

### 4.3 Error bounds

The main result of this paragraph is the following theorem, providing a convergence rate of the unnormalized filter approximation error for the two step recursive scheme **A**.

**Theorem 4.1** *Let  $(\hat{X}_k)$  be a marginal stationary  $(N_k)$ -quantization of  $(X_k)$ ,  $f$  satisfying  $H_n^n f \in C_{b,Lip}^1$ . Assume **H1**, **H2**, **H3**, **H4**,  $q = d$  and furthermore :*

**H 5** *There exists a constant  $\psi_p > 0$  and  $\bar{s} > 1$  such that:*

$$\max_{0 \leq k \leq n-1} \|\Psi_k(X_k, \varepsilon_{k+1})\|_{\bar{s}p} \leq \psi_p < +\infty.$$

Hence, there exists a non negative real sequence of constants  $(M_j^n)_{0 \leq j \leq n}$  such that:

$$|\pi_n f - \hat{\pi}_n f| \leq \sum_{j=0}^n M_j^n \|\Delta_j\|_{\max\{stp, \bar{t}p, 2p\}}^2$$

where  $s = \frac{\bar{s}}{\bar{s}-1}$ ,  $t > 0$ ,  $\frac{1}{t} + \frac{1}{\bar{s}} = 1$  and  $M_j^n \leq \alpha(p, f)(n+1)L^n \left(\frac{(L)^{j+1}-1}{3K-1}\right) \left(\frac{(3K)^{n-j+1}-1}{3K-1}\right)$ .

**Example 4.2** Assume that  $F_k : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  reads:

$$F_k(x, \varepsilon) = b_k(x) + \sigma_k(x)\varepsilon, \quad (4.10)$$

where  $\sigma_k$  and  $b_k$  are differentiable with bounded derivatives and  $\forall x \in \mathbb{R}^d$ ,  $\sigma_k(x) > c_k$ . Then:

$$\Psi_k(x, \varepsilon) = \frac{\sigma'_k(x) + \frac{p'}{p}(\varepsilon)(\varepsilon \sigma'_k(x) + b'_k(x))}{\sigma_k(x)}.$$

- When  $\varepsilon_k \sim \mathcal{N}(0, 1)$ , it is the natural framework to study the Euler scheme of a Brownian diffusion. In this case, the previous hypothesis **H5** is satisfied.
- When  $\varepsilon_k$  distribution is centered Laplace of parameter  $\lambda > 0$ , or  $\varepsilon_k + m \sim \text{Gamma}(m, 1)$  with  $m > 1$ , hypothesis **H5** is also satisfied.
- In a more general case, when  $\varepsilon_k \in L^{p+\eta}$  for some  $\eta > 0$  the following assumption:  
**H 5'** There exists a constant  $\psi_p > 0$  and  $\bar{s}' > 1$  such that  $\|\mathbf{P}'_{\mathbf{p}}(\varepsilon_1)\|_{\bar{s}'p} \leq \psi_p < +\infty$ .  
could replace **H5** and gives more explicit conditions on the signal innovation distribution.

Compared to Example 3.1 given for the one step iterative scheme, we see that hypothesis **H5** (or **H5'**) allows to relax the boundedness constraint on  $\partial_\varepsilon F_k$  in **H2'** to involve some other constraints on the signal innovations distribution.

The structure of the proof of Theorem 4.1 is the same as that of the previous section. We first study the error induced by the differential term estimation. Let us reconsider for  $0 \leq k \leq n-1$  and the test function  $f$ :  $\hat{A}_k := DR_k f(\hat{X}_k) - \widehat{DR}_k f(\hat{X}_k)$ . The error bounds for  $\|\hat{A}_k\|_p$  with the new definition of the differential term approximation  $\widehat{DR}_k f$  are given by the following lemma:

**Lemma 4.1** *With assumption **H5** on the weight function  $\Psi_k$  and  $f$  such that  $H_n^n f \in \mathcal{C}_{b,Lip}^1$ , there exists a non negative real sequence  $(D_j^{k,n})_{0 \leq k \leq j \leq n}$  such that:*

$$\|\hat{A}_k\|_p \leq \sum_{j=k}^n D_j^{k,n} \|\Delta_j\|_{sp}$$

where  $s = \frac{\bar{s}}{\bar{s}-1}$  and  $D_j^{k,n} \leq \alpha(p, f) L^{n-k} \frac{(3K)^{n-j+1} - 1}{3K-1}$ .

**Proof.**

We redefine the operators  $Q_k$  and  $\hat{Q}_k$  for  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  as follows:

$$\begin{aligned} Q_k f(X_k) &= -\mathbb{E}[f(X_{k+1})\Psi_k(X_k, \varepsilon_{k+1})|X_k], \\ \hat{Q}_k f(\hat{X}_k) &= -\mathbb{E}[f(\hat{X}_{k+1})\Psi_k(X_k, \varepsilon_{k+1})|\hat{X}_k]. \end{aligned}$$

Then  $Q_k f = \mathbf{D}\mathbf{P}_k f$ , so that:

$$\begin{aligned} DR_k f(\hat{X}_k) &= Dg_k(\hat{X}_k)\mathbf{P}_k R_{k+1} f(\hat{X}_k) + g_k(\hat{X}_k)Q_k R_{k+1} f(\hat{X}_k), \\ \widehat{DR}_k f(\hat{X}_k) &= Dg_k(\hat{X}_k)\hat{\mathbf{P}}_k \hat{R}_{k+1} f(\hat{X}_k) + g_k(\hat{X}_k)\hat{Q}_k \hat{R}_{k+1} f(\hat{X}_k). \end{aligned}$$

Consequently,  $\hat{A}_k$  can be written as:

$$\hat{A}_k = Dg_k(\hat{X}_k) \left[ \mathbf{P}_k R_{k+1} f(\hat{X}_k) - \hat{\mathbf{P}}_k \hat{R}_{k+1} f(\hat{X}_k) \right] + g_k(\hat{X}_k) \left[ Q_k R_{k+1} f(\hat{X}_k) - \hat{Q}_k \hat{R}_{k+1} f(\hat{X}_k) \right],$$

so that using **H1** and that  $\hat{\mathbf{P}}_k \hat{R}_{k+1} f(\hat{X}_k) = \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1})|\hat{X}_k]$ , we have:

$$\|\hat{A}_k\|_p \leq L \|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1})|\hat{X}_k]\|_p + L \|Q_k R_{k+1} f(\hat{X}_k) - \hat{Q}_k \hat{R}_{k+1} f(\hat{X}_k)\|_p. \quad (4.11)$$

Since conditional expectation is an  $L^p$ -contraction,, the first term on the right hand side of inequality (4.11) writes: we have:

$$\begin{aligned}
\|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p & \\
&\leq \|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) | \hat{X}_k]\|_p \\
&\quad + \|\mathbb{E}[\mathbf{P}_k R_{k+1} f(X_k) - \hat{R}_{k+1} f(\hat{X}_{k+1}) | \hat{X}_k]\|_p \\
&\leq [\mathbf{P}_k R_{k+1} f]_{Lip} \|\Delta_k\|_p + \|R_{k+1} f(X_{k+1}) - \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p.
\end{aligned} \tag{4.12}$$

It follows from Theorem 2.1 that:

$$\|\mathbf{P}_k R_{k+1} f(\hat{X}_k) - \hat{\mathbf{P}}_k \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p \leq \sum_{j=k+1}^n C_j^{k+1,n} \|\Delta_j\|_p + K \|DR_{k+1} f\|_\infty \|\Delta_k\|_p.$$

Moreover, the second term on the right hand side of inequality (4.11) gives:

$$\begin{aligned}
&\|Q_k R_{k+1} f(\hat{X}_k) - \hat{Q}_k \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p \\
&\leq \|Q_k R_{k+1} f(\hat{X}_k) - \mathbb{E}[Q_k R_{k+1} f(X_k) | \hat{X}_k]\|_p + \|\mathbb{E}[Q_k R_{k+1} f(X_k) | \hat{X}_k] - \hat{Q}_k \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p \\
&\leq \|Q_k R_{k+1} f(\hat{X}_k) - Q_k R_{k+1} f(X_k)\|_p \\
&\quad + \|\mathbb{E}[\mathbb{E}[R_{k+1} f(X_{k+1}) \Psi_k(X_k, \varepsilon_{k+1}) | X_k] | \hat{X}_k] - \mathbb{E}[\hat{R}_{k+1} f(\hat{X}_{k+1}) \Psi_k(X_k, \varepsilon_{k+1}) | \hat{X}_k]\|_p \\
&\leq \|Q_k R_{k+1} f(\hat{X}_k) - Q_k R_{k+1} f(X_k)\|_p + \|\Psi_k(X_k, \varepsilon_{k+1}) (R_{k+1} f(X_{k+1}) - \hat{R}_{k+1} f(\hat{X}_{k+1}))\|_p.
\end{aligned} \tag{4.13}$$

But,  $Q_k R_{k+1} f(X_k) = D\mathbf{P}_k R_{k+1} f(X_k)$ , so hypothesis **H2** on  $\mathbf{P}_k$  implies that:

$$[Q_k R_{k+1} f]_{Lip} = [D\mathbf{P}_k R_{k+1} f]_{Lip} \leq K ([DR_{k+1} f]_{Lip} \vee \|DR_{k+1} f\|_\infty) = K u_{k+1}.$$

Hence,

$$\|Q_k R_{k+1} f(\hat{X}_k) - Q_k R_{k+1} f(X_k)\|_p \leq K u_{k+1} \|\Delta_k\|_p. \tag{4.14}$$

Using Holder inequality, with  $s = \frac{\bar{s}}{\bar{s}-1} \geq 1$  we get:

$$\begin{aligned}
&\|\Psi_k(X_k, \varepsilon_{k+1}) (R_{k+1} f(X_{k+1}) - \hat{R}_{k+1} f(\hat{X}_{k+1}))\|_p \\
&\leq \|\Psi_k(X_k, \varepsilon_{k+1})\|_{\bar{s}p} \|R_{k+1} f(X_{k+1}) - \hat{R}_{k+1} f(\hat{X}_{k+1})\|_{sp}
\end{aligned} \tag{4.15}$$

It follows from Theorem 2.1 and hypothesis **H5** by combining terms (4.13), (4.14) and (4.15) that:

$$\begin{aligned}
\|Q_k R_{k+1} f(\hat{X}_k) - \hat{Q}_k \hat{R}_{k+1} f(\hat{X}_{k+1})\|_p & \\
&\leq \sum_{j=k+1}^n \psi_p C_j^{k+1,n} \|\Delta_j\|_{sp} + K u_{k+1} \|\Delta_k\|_p,
\end{aligned} \tag{4.16}$$

$$\begin{aligned}
\text{and } \|\hat{A}_k\|_p &\leq L(\psi_p + 1) \sum_{j=k+1}^n C_j^{k+1,n} \|\Delta_j\|_{sp} + LK (u_{k+1} + \|DR_{k+1} f\|_\infty) \|\Delta_k\|_{sp} \\
&\leq \sum_{j=k+1}^n D_j^{k,n} \|\Delta_j\|_{sp}.
\end{aligned} \tag{4.17}$$

Then, explicit upper bounds for  $D_j^{k,n}$  can easily be established. (see [18])  $\square$

**Proof of Theorem 4.1.** Reconsider  $V_k f = \mathbb{E}[R_k f(X_k)|\hat{X}_k] - \hat{R}_k f(\hat{X}_k)$  for  $0 \leq k \leq n$ . The proof can be carried out as in the previous case of Theorem 3.1. The unique difference lies in the term  $\hat{A}_k$ . Using Lemma 4.1 combined with Holder inequality for some  $t > 1$  and its conjugate  $\bar{t} = \frac{t}{t-1}$  we have:

$$\|\langle \hat{A}_{k+1}, \Delta_{k+1} \rangle\|_p \leq \|\hat{A}_{k+1}\|_{tp} \|\Delta_{k+1}\|_{\bar{t}p} \leq \frac{1}{2} \sum_{j=k+1}^n D_j^{k+1,n} (\|\Delta_j\|_{stp}^2 + \|\Delta_{k+1}\|_{\bar{t}p}^2).$$

Then inequality (3.16) writes:

$$\begin{aligned} \|V_k f\|_p &\leq ([DR_k f]_{Lip} + LK u_{k+1}) \|\Delta_k\|_{2p}^2 + 2L[DR_{k+1} f]_{Lip} \|\Delta_{k+1}\|_{2p}^2 \\ &\quad + L \frac{1}{2} \sum_{j=k+1}^n D_j^{k+1,n} (\|\Delta_j\|_{stp}^2 + \|\Delta_{k+1}\|_{\bar{t}p}^2) + L \|V_{k+1} f\|_p \\ &\leq ([DR_k f]_{Lip} + LK u_{k+1}) \|\Delta_k\|_{\max\{stp, \bar{t}p, 2p\}}^2 \\ &\quad + L \left( [DR_{k+1} f]_{Lip} + D_{k+1}^{k+1,n} + \frac{1}{2} \sum_{j=k+2}^n D_j^{k+1,n} \right) \|\Delta_{k+1}\|_{\max\{stp, \bar{t}p, 2p\}}^2 \\ &\quad + \frac{1}{2} L \sum_{j=k+2}^n D_j^{k+1,n} \|\Delta_j\|_{\max\{stp, \bar{t}p, 2p\}}^2 + L \|V_{k+1} f\|_p \end{aligned} \quad (4.18)$$

By induction, we derive:  $\|V_k f\|_p \leq \sum_{j=k}^n M_j^{k,n} \|\Delta_j\|_{\max\{stp, \bar{t}p, 2p\}}^2$ .

Taking  $k = 0$  and writing  $|\pi_n f - \hat{\pi}_n f| \leq \|V_0 f\|_p$  we establish the announced result. See [18] for a detailed proof of the explicit expressions of  $(M_j^{k,n})$ .  $\square$

#### 4.4 The case of regularizing kernels

In this paragraph we deal with an interesting skill of the two step recursive first order scheme, which allows to establish first order schemes for non differentiable test functions  $f$ , more precisely with no differentiability assumption on  $H_n^n f$ .

**Proposition 4.1 H2''** Assume  $\mathbf{P}_k$  is  $K$ -Lipschitz such that for all  $f$  bounded Lipschitz continuous,  $\mathbf{P}_k f \in \mathcal{C}_{b,Lip}^1$ .

If  $f$  is a test function such that  $H_n^n f$  is bounded Lipschitz continuous, then  $R_k f \in \mathcal{C}_{b,Lip}^1$  for all  $0 \leq k \leq n-1$ .

This proposition is easily proved, using equation (2.17) and an induction on  $k$ . Furthermore, it allows to define an alternative scheme to scheme **A**, taking into account the non differentiability of  $H_n^n f$ :



### Scheme A'

$$\left\{ \begin{array}{l} \widehat{R}_n f(\widehat{X}_n) = H_n^n f(\widehat{X}_n) = \widehat{R}_n f(\widehat{X}_n), \\ \widehat{R}_{n-1} f(\widehat{X}_{n-1}) = g_{n-1}(\widehat{X}_{n-1}) \mathbb{E}[H_n^n f(\widehat{X}_n) | \widehat{X}_{n-1}] = \widehat{R}_{n-1} f(\widehat{X}_{n-1}), \\ \widehat{R}_k f(\widehat{X}_k) = g_k(\widehat{X}_k) \widehat{\mathbf{P}}_k \widehat{R}_{k+1} f(\widehat{X}_{k+1}) + g_k(\widehat{X}_k) \times \\ \quad \left( \mathbb{E}[\langle \mathbf{D}g_{k+1}(\widehat{X}_{k+1}) \widehat{\mathbf{P}}_{k+1} \widehat{R}_{k+2} f(\widehat{X}_{k+2}), \Delta_{k+1} \rangle | \widehat{X}_k] - \mathbb{E}[\langle g_{k+1}(\widehat{X}_{k+1}) \times \right. \\ \quad \left. \mathbb{E}[\widehat{R}_{k+2} f(\widehat{X}_{k+2}) \Psi_{k+1}(X_{k+1}, \varepsilon_{k+2}) | \widehat{X}_{k+1}], \Delta_{k+1} \rangle | \widehat{X}_k] \right), \\ 0 \leq k \leq n-2. \end{array} \right. \quad (4.19)$$

We then define the first order unnormalized filter estimator by  $\widehat{\pi}_n f = \mathbb{E}[\widehat{R}_0 f(\widehat{X}_0)]$  generated from scheme A'. The error induced by such an estimator introduces additional zero order type terms as we need one single backward iteration to be able to use first order correctors. This can be seen clearly in the the following theorem which proof is detailed in [18].

**Theorem 4.2** *Let  $(\widehat{X}_k)$  be a stationary  $(N_k)$ -quantization of  $(X_k)$ ,  $f$  satisfying  $H_n^n f$  is bounded Lipschitz continuous. Assume **H1**, **H2''**, **H3**, **H4**, **H5** and  $q = d$ , then, there exists a non negative real sequence of constants  $(\bar{M}_j^n)_{0 \leq j \leq n}$  such that:*

$$|\pi_n f - \widehat{\pi}_n f| \leq \sum_{j=0}^n \bar{M}_j^n \|\Delta_j\|_{\max\{stp, \bar{t}p, 2p\}}^2 + C^1 \|\Delta_{n-1}\|_p + C^2 \|\Delta_n\|_p$$

where  $s = \frac{\bar{s}}{\bar{s}-1}$ ,  $t > 0$ ,  $\frac{1}{t} + \frac{1}{\bar{t}} = 1$ .

In practice, the regularizing effect can be viewed in the case of the Euler scheme of a diffusion implemented with a Gaussian noise (see Example 4.2 and Appendix B). This is the case studied in [3] for pricing American options with first order schemes. It is shown that  $\mathbf{P}_k$  satisfies **H2''**. Nevertheless, a special attention have to be given to the Lipschitz constants dependency in time discretization step, and consequently in our filtering problem, to the observtaion horizon. Namely, if  $f$  is Lipschitz continuous, then according to Proposition 2 in [3] we have  $[\mathbf{D}\mathbf{P}_k f]_{Lip} \leq C[f]_{Lip} \sqrt{n}$ . This result alters  $\bar{M}_j^n$  dependency in  $n$  and consequently the filter estimator convergence for high observation horizons.

**Remark 4.1** For numerical implementation, we can compensate the error bounds deterioration in Theorem 4.2 by bigger quantizers in the two last observation dates  $n-1$  and  $n$ .

## 5 Convergence result for the normalized filter

Let  $f$  be such that  $H_n^n f \in \mathcal{C}_{b, Lip}^1$ . Owing to Theorem 3.1 and Theorem 4.1, we have seen that the estimation error on the unnormalized filter  $\pi_n$ , using stationary  $(N_k)$ -quantizations  $(\widehat{X}_k)$ , can be written:

$$|\pi_n f - \widehat{\pi}_n f| \leq \sum_{j=0}^n M_j^n(f, \alpha) \|\Delta_j\|_{2\alpha p}^2,$$

where  $\alpha = \max\{\frac{st}{2}, \frac{\bar{t}}{2}, 1\} \geq 1$  or  $\alpha = 1$  depending on whether we are using scheme **A** or **B**, and  $\Delta_j = X_j - \hat{X}_j$ .

Now, we derive results on the normalized first order quantization filter estimator  $\hat{\Pi}_n$ , defined by Kallianpur-Striebel formula as  $\hat{\Pi}_n f = \frac{\hat{\pi}_n f}{\hat{\pi}_n \mathbf{1}}$ .

Thus, the estimation error will be:

$$\begin{aligned} |\Pi_n f - \hat{\Pi}_n f| &\leq \left| \frac{\pi_n f}{\pi_n \mathbf{1}} - \frac{\pi_n f}{\hat{\pi}_n \mathbf{1}} \right| + \left| \frac{\pi_n f - \hat{\pi}_n f}{\hat{\pi}_n \mathbf{1}} \right| \\ &\leq \frac{\|H_n^n f\|_\infty \pi_{n-1} \mathbf{1}}{\pi_n \mathbf{1} \hat{\pi}_n \mathbf{1}} |\pi_n \mathbf{1} - \hat{\pi}_n \mathbf{1}| + \frac{1}{\hat{\pi}_n \mathbf{1}} |\pi_n f - \hat{\pi}_n f| \\ &\leq \sum_{j=0}^n \frac{M_j^n(f, \alpha) + c^y M_j^n(\mathbf{1}, \alpha) \|H_n^n f\|_\infty}{\hat{\pi}_n \mathbf{1}} \|\Delta_j\|_{2\alpha p}^2 \end{aligned} \quad (5.1)$$

Since  $\alpha = 1$  in Theorem 3.1, the convergence rate improvement obtained for the unnormalized filter is preserved by the normalization.

When  $\alpha > 1$ , which is the case for Theorem 4.1, further results are needed to establish a convergence rate improvement. In fact, from inequality (5.1) it comes out that we need to describe the  $L^{2\alpha p}$ -behavior of sequences of  $L^{2p}$ -optimal quantizers. In this direction, a rather satisfactory result can be established using Zador Theorem 2.2 and Holder inequality. Namely, if  $X \in L^{r'}(\mathbb{R}^d)$  for every  $r' > 0$ , then  $\|X - h_N^*(X)\|_s = O(N^{-\frac{\rho}{d}})$  for any  $\rho \in (0, \frac{r}{s})$ .

This allows to establish the following theorem, for  $\hat{\Pi}_n$  obtained from the two step recursive scheme:

**Theorem 5.1** *Assume that  $\bar{s}$  in **H5** satisfies  $\bar{s} > \frac{3}{2}$  and that for  $0 \leq k \leq n$  and all  $r > 0$   $X_k \in L^r(\mathbb{R}^d)$ . Let  $(\hat{X}_k)$  be an  $L^2$ -optimal  $(N_k)$ -quantization of  $(X_k)$ .*

*Then, there exists  $\rho \in (\frac{1}{2}, 1]$  such that for all  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  satisfying  $H_n^n f \in \mathcal{C}_{b,Lip}^1$  we have:*

$$|\Pi_n f - \hat{\Pi}_n f| \leq \sum_{j=0}^n c_j(\rho, p, d) \frac{M_j^n(f, \alpha) + c^y M_j^n(\mathbf{1}, \alpha) \|H_n^n f\|_\infty}{\hat{\pi}_n \mathbf{1}} N_j^{-\frac{2\rho}{d}}.$$

**Proof.** If  $\bar{s} > \frac{3}{2}$ , then  $1 < s < 3$  and there exists  $\frac{4}{3} < t < \frac{4}{s}$ . For such a number  $t > 1$  we will have  $\bar{t} < 4$  and  $st < 4$  so that inequality (5.1) is satisfied for  $\alpha = \max\{\frac{\bar{t}}{2}, \frac{st}{2}, 1\} \in [1, 2]$ . Hence, for some  $\rho \in (\frac{1}{2}, \frac{1}{\alpha}) \subset (0, \frac{1}{\alpha})$ , we can write:  $\|\Delta_k\|_{2\alpha p} = O(N_k^{-\frac{2\rho}{d}})$ .

Consequently from (5.1),  $|\Pi_n f - \hat{\Pi}_n f| \leq \sum_{j=0}^n c_j(\rho, p, d) \frac{M_j^n(f, \alpha) + c^y M_j^n(\mathbf{1}, \alpha) \|H_n^n f\|_\infty}{\hat{\pi}_n \mathbf{1}} N_j^{-\frac{2\rho}{d}}$ .  $\square$

**Remark 5.1** A conjecture has been made recently by H. Luschgy and G. Pagès to describe the  $L^{s'}$ -behavior of sequences of  $L^r$ -optimal quantizers of an  $\mathbb{R}^d$ -valued random vector for some  $0 < r < s' < r + d$ :

If  $X \in L^r(\mathbb{R}^d)$  such that  $\mathbb{P}_X(d\xi) = \varphi(\xi)\lambda_d(d\xi)$  and  $\int \varphi^{1-\frac{s'}{r+d}} d\lambda_d < +\infty$ , then any sequence  $(h_N^*)$  of  $L^r$ -optimal  $N$ -quantizers most likely satisfies  $\|X - h_N^*(X)\|_{s'} = O(N^{-\frac{1}{d}})$ .

This allows to establish an equivalent of Theorem 5.1 where  $\rho = 1$  and  $\bar{s}$  is assumed to satisfy  $\bar{s} > 1 + \frac{1}{d}$  and then to rise convergence rate order of two step recursive schemes to the one step recursive one.

## 6 Numerical illustrations

Previous filter approximation methods will be applied to estimate  $\Pi_n f_1$  and  $\Pi_n f_2$ , where  $f_1(x) = x$  and  $f_2(x) = \exp(-|x|)$ . Elements of comparison with alternative filter estimation methods will be given, namely particle filtering methods:

**SIS** Sequential Importance Sampling [1, 6] which is based on a weighted Monte Carlo approach. This method can be considered as close to the quantization method in the sense that it uses weight transformations in the updating step. Unfortunately, it is known to suffer from weights degenerescense.

**SIR** Sequential Importance Re-sampling [9, 6] which adds a re-sampling step to the previous algorithm in order to avoid weights degenerescense.

We will test estimations for different fixed observation sets and so we denote by  $\hat{\Pi}_{y,n}$ , the estimation filter associated to the observation process  $y = (y_0, \dots, y_n)$ .

In all the following examples, we choose to study stationary signal processes in order to simplify the off line procedure of computing the quantizers. In fact, as we marginally quantize the signal process, we can just expand the grids of the centered reduced corresponding distribution. The obtained quantizers are no longer optimal, some further manipulations are necessary to save the quantizer stationarity property especially in the multidimensional cases. This paragraph is an overview of selected numerical experiments illustrating first order schemes behaviour. Further numerical tests will be detailed in a forthcoming paper [19], especially concerning the comparison with the particle filtering approach. Results obtained with infinite dimension filters [7, 5] will also be presented.

### 6.1 Kalman filter

Both signal and observation equations are linear with Gaussian independent noises. It is known, that the filter in this case has a Gaussian distribution which parameters (mean and variance) can be computed sequentially via a deterministic algorithm (KF), (see [8]).

$$\text{We set: } \begin{cases} X_k = \rho X_{k-1} + \theta \varepsilon_{k+1}, \\ Y_k = X_k + \alpha \eta_k, \\ \varepsilon_k \text{ and } \eta_k \text{ iid } \sim \mathcal{N}(0, I_d), \\ \rho, \theta, \alpha \in \mathcal{M}_d(\mathbb{R}). \end{cases} \quad (6.1)$$

#### 6.1.1 One dimensional case: d=1

We choose  $-1 < \rho < 1$  and  $X_0 \sim \mathcal{N}(0, \frac{\theta^2}{1-\rho^2})$ , so that for any  $0 \leq k \leq n$ , we have  $X_k \sim \mathcal{N}(0, \frac{\theta^2}{1-\rho^2})$ . In this particular case, we could first compute<sup>2</sup>  $\Gamma$  an  $L^2$ -optimal quantizer of the centered reduced Gaussian distribution and the companion parameters for a single transition step. The quantizers  $\Gamma_k$  are then deduced by an expansion  $\Gamma_k = \frac{\theta^2}{1-\rho^2} \times \Gamma$ .

The two first order schemes are compared to the zero order one with  $N_k = 200$ ,  $0 \leq k \leq n$ .

---

<sup>2</sup>Optimal quantizers for the Gaussian distribution are downloadable on <http://www.proba.jussieu.fr/pageperso/pages/>

Exact values are computed via the Kalman-Bucy recursive filter algorithm. Particles methods are also tested for the sake of comparison.

$(\rho, \theta, \alpha)$	(0.65,1.0,0.1)		(0.8,1.0,0.1)	
	$\hat{\Pi}_{y,25}f_1$	$\hat{\Pi}_{y,25}f_2$	$\hat{\Pi}_{y,25}f_1$	$\hat{\Pi}_{y,25}f_2$
KF(Ref. Value)	-3.239	0.039	1.754	0.17394
SIS (5000 pts)	-3.244	0.039	1.7487	0.17489
SIR (5000 pts)	-3.2398	0.039	1.7542	0.1739
QF Or0 (200 pts)	-3.2394	0.0393	1.7522	0.17425
QF Or1 1-step (200 pts)	-3.2381	0.039431	1.7524	0.17422
QF Or1 2-step (200 pts)	-3.2381	0.039431	1.7524	0.17422

Table 1: One dimensional Kalman filter case.

### 6.1.2 Multidimensional case: $d=2$

Although the quantization based filter schemes presented previously depend on the signal dimension  $d$ , for both complexity and convergence rate, it remains interesting to compute estimations for medium signal dimensions. We reconsider equation (6.1) with parameters:

$$\rho = \begin{pmatrix} 0.996 & 0 \\ 0 & 0.996 \end{pmatrix}, \theta = \begin{pmatrix} 0.05 & -0.01 \\ -0.01 & 0.02 \end{pmatrix} \quad \text{and} \quad \alpha = 0.5I_d.$$

The initial signal distribution is centered, Gaussian with covariance matrix:

$$\Sigma_0 = \begin{pmatrix} 0.325 & -0.087 \\ -0.087 & 0.0626 \end{pmatrix}.$$

The chosen prior distribution is the stationary one. For signal quantization, we take  $\Gamma = \{z^1, \dots, z^N\}$  the  $L^2$ -optimal  $N$ -quantizer of a centered reduced Gaussian distribution. At  $0 \leq k \leq N$ ,  $X_k \sim \mathcal{N}(0, \Sigma_0)$  and we define the marginal stationary  $(N_k)$ -quantizer of  $(X_k)$  as follows:

$$\hat{X}_k = \sum_{i=1}^N \Sigma_0^{\frac{1}{2}} z^i \mathbf{1}_{\{X_k \in \Sigma_0^{\frac{1}{2}} \mathbf{C}_i(\Gamma)\}}$$

Although quantizers are not optimal, we obtain satisfactory convergence results. Convergence errors are represented in Figure 4. From the log-log scale representation in Figure 4, we can evaluate the convergence rate improvement using a regression. Table 6.1.2 summarizes the computed slopes of the regressions.

Or0	Or1 1-step	Or1 2-step
-0.45	-1.1	-1.04

Table 2: Regression slopes on the log-log scale representation ( $d=2$ )

We observe nearly the expected theoretical results. The convergence rate for the zero order scheme is close of  $\frac{1}{d} = 0.5$ . For first order schemes, the slope is slightly better than the theoretical one  $\frac{2}{d} = 1$ .

## 6.2 Canonical stochastic volatility model (SVM)

We introduce now a non linearity in the observation equation. We consider the following state equations in  $\mathbb{R}$ :

$$\begin{cases} X_k = \beta X_{k-1} + \sigma \varepsilon_{k+1}, \\ Y_k = \exp(\frac{X_k}{2}) \eta_k, \\ \varepsilon_k \text{ and } \eta_k \text{ iid } \sim \mathcal{N}(0, 1), \\ -1 < \beta < 1 \text{ and } \sigma \in \mathbb{R}_+^*. \end{cases} \quad (6.2)$$

**Remark 6.1** This is the time discretization of a continuous diffusion model introduced in finance as a model of an asset dynamics with stochastic volatility. The stock price  $S_t$  and its volatility  $\sigma_t$  solve the following stochastic differential system:

$$\begin{cases} dS_t = \mu_t S_t dt + \sigma_t S_t dW_t, \\ d(\ln(\sigma_t^2)) = -\lambda \ln(\sigma_t^2) dt + \tau dW_t. \end{cases} \quad (6.3)$$

The stock price is supposed to be observable so that the filtering problem corresponds to a volatility estimation problem, given the set of observed past prices. Taking a time discretization step  $\Delta$ , the Euler scheme writes:

$$\begin{cases} \ln(\frac{S_{k+1}}{S_k}) = (\mu_k - \frac{1}{2}\sigma_k^2)\Delta + \sigma_k \sqrt{\Delta} \eta_k, \\ \ln(\sigma_{k+1}^2) = (1 - \lambda\Delta)\ln(\sigma_k^2) + \tau\sqrt{\Delta}\varepsilon_{k+1}. \end{cases} \quad (6.4)$$

Now, taking  $Y_k = \ln(\frac{S_{k+1}}{S_k})$ ,  $X_k = \ln(\sigma_k^2)$ ,  $\eta_k$  and  $\varepsilon_k$  iid  $\mathcal{N}(0, 1)$  conducts to the state equations adopted for the illustration.

Here also we choose  $X_0 \sim \mathcal{N}(0, \frac{\sigma^2}{1-\beta^2})$ , in order to use the same grid at each time step  $k$ . The choice of the triplet  $(\lambda, \tau, \Delta)$  will determine the discrete time model parameters  $(\beta, \sigma)$ . The exact filter value is not computable for such model, so Figure 1 shows the convergence behavior of the quantization filters. The first order schemes clearly converge faster. Comparison with particle methods is made possible by computing some confidence interval through the 5% and 95% centiles over 4000 realizations of the particle filter estimator. In Figure 2 are depicted this interval bounds and one realization of the random estimator as functions of the particle number. For a comparison between the two methods (particles and quantization), we represent in Figure 3 quantization based filters in the confidence interval of 10000 particles.

## 6.3 Numerical stability

Two stability aspects have been studied through numerical applications. The implemented state equations are those of the previous section (see equation (6.2)) when we model stochastic

volatility.

The first point we will be interested in, is degeneration of intuitive scheme devised in Remark 3.1. An illustration of such a problem is represented by Figure 5.

The second point is the stability of our estimations in time. This is a recurrent problem in filtering methods. Even if we considered a fixed observation horizon all over this work, it is important to study the estimation behaviour when  $n$  grows. As the constants  $M_j^n$  are exponentially depending of the observation horizon, we have been interested in verifying that this does not alter the numerical performances of our filter estimators. (see Figure 6). Note that the chosen state equations and the stationarity assumption give that  $K = \beta < 1$ .

**Acknowledgment:** I am grateful to Pr. Gilles Pagès and Pr. Huyên Pham who supported this work and contributed to its development by enriching suggestions.

## References

- [1] S. Arulampalam, T. Clapp, N. Gordon, and S. Maskall. A tutorial on particle filters for On-line Non-linear/Non Gaussian Bayesian tracking. QinetiQ Ltd, DSTO, IEEE, 2001.
- [2] V. Bally and G. Pagès. A quantization algorithm for solving discrete time multi-dimensional optimal stopping problems. *Bernoulli*, 9:1003–1049, 2003.
- [3] V. Bally, G. Pagès, and J. Printems. First order schemes in the numerical quantization method. *Mathematical finance*, 13(1):1–16, 2003.
- [4] V. Bally, G. Pagès, and J. Printems. A quantization tree method for pricing and hedging multidimensional american options. *Mathematical Finance*, 15(1):119–168, 2005.
- [5] M. Chaleyat-Maurel and V. Genon-Catalot. Computable infinite dimensional filters with applications to discretized diffusion processes. Preprint of Laboratoire de Probabilités et Modèles Aléatoires - PMA-989, 2005.
- [6] A. Doucet, N. de Freitas, and N. Gordon, editors. *Sequential Monte Carlo methods in Practice*. Springer, 1st edition, 2001.
- [7] V. Genon-Catalot. A non linear explicit filter. *Statist. and Prob. letters*, 61:145–154, 2003.
- [8] F. Le Gland. *Introduction au filtrage en temps discret. Filtre de Kalman, Modèles de Markov cachés*. IRISA/INRIA, 2002-2003.
- [9] N. Gordon, D.J. Salmond, and A.F.M. Smith. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. *IEEE Proceedings*, 140(2):107–113, April 1993.
- [10] S. Graf and H. Luschgy. *Foundations of quantization for probability distributions*. Lecture Notes in Mathematics. Springer, 2000.
- [11] G. Kallianpur and C. Striebel. Estimation of stochastic systems: Arbitrary system process with additive white noise observation errors. *Ann. Math. Statist.*, 39(2):785–801, 1968.

- [12] G. Kitagawa. Non-Gaussian state-space modeling of nonstationary time series. *Journal of the American statistical association*, 82(400):1032–1063, 1987.
- [13] P. Del Moral, J. Jacod, and P. Protter. The Monte Carlo method for filtering with discrete-time observations. *Probab. Theory and Relat. Fields*, 120(2):346–368, June 2001.
- [14] N. Oudjane. *Stabilité et approximations particulières en filtrage non linéaire : Application au pistage*. PhD thesis, Université de Rennes, 2000.
- [15] G. Pagès and H. Pham. Optimal quantization methods for nonlinear filtering with discrete time observations. To appear in *Bernoulli*, 2003.
- [16] G. Pagès, H. Pham, and J. Printems. Optimal quantization methods and applications to numerical problems in finance. In *Handbook of Computational and Numerical Methods in Finance*. S.T. Rachev, Birkhauser, Boston, 2004.
- [17] G. Pagès and J. Printems. Optimal quadratic quantization for numerics : the Gaussian case. *Monte Carlo methods and applications*, 9(2):135–168, 2003.
- [18] A. Sellami. Optimal quantization methods for filtering and applications to finance. In progress PhD thesis, Université Paris Dauphine and Laboratoire de Probabilités et Modèles Aléatoires - Université Paris VI.
- [19] A. Sellami. Comparative survey on non linear filtering methods : the quantization and the particle filtering approaches. Preprint of Laboratoire de Probabilités et Modèles Aléatoires, 2005.

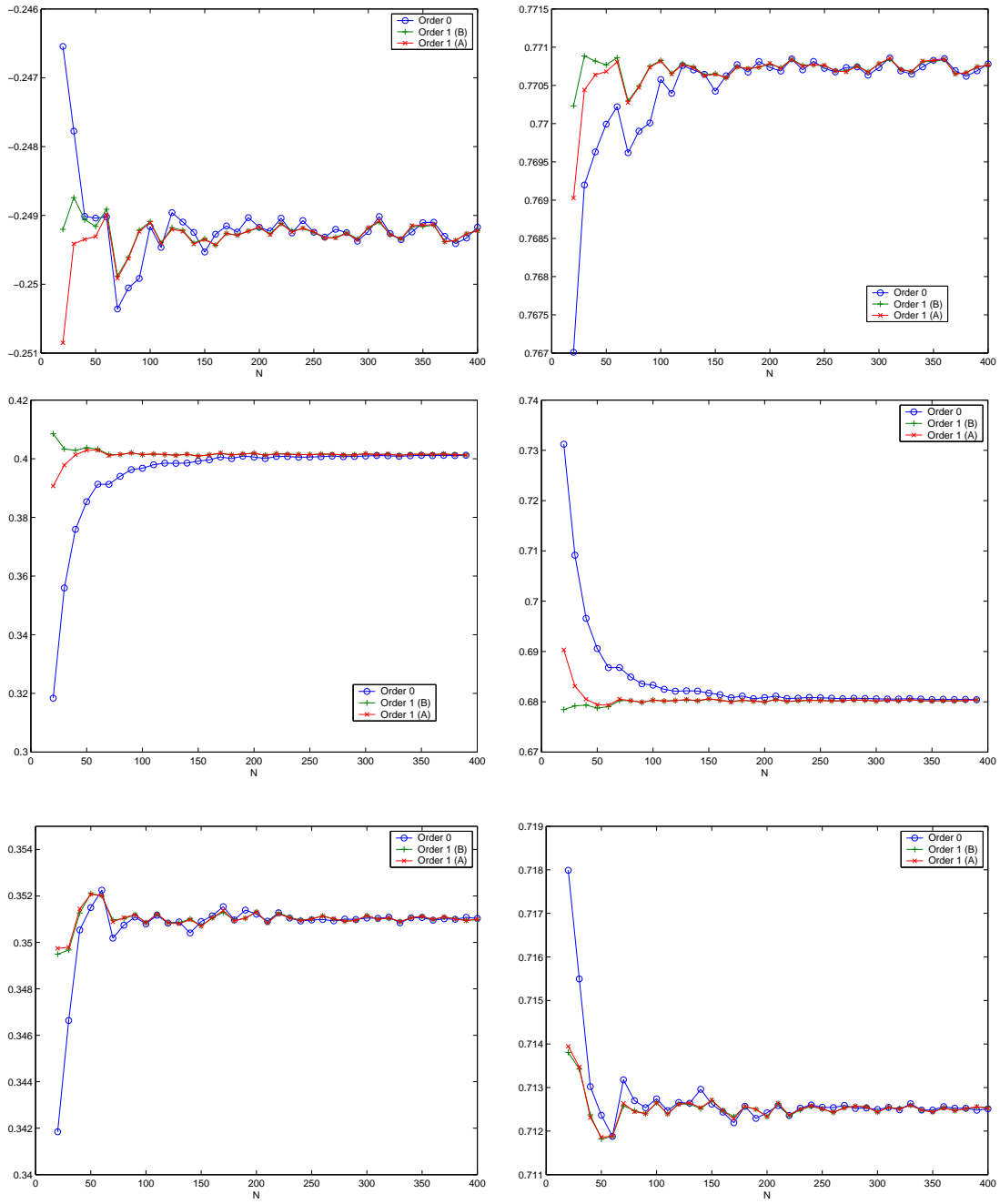


Figure 1: Quantization filter approximations for SVM as a function of the quantizer size  $N_k$  - three different observation 50-tuples (right:  $\Pi_{y,50}f_1$ , left:  $\Pi_{y,50}f_2$ ) -  $(\beta, \sigma) = (0.996, 0.0316)$ .



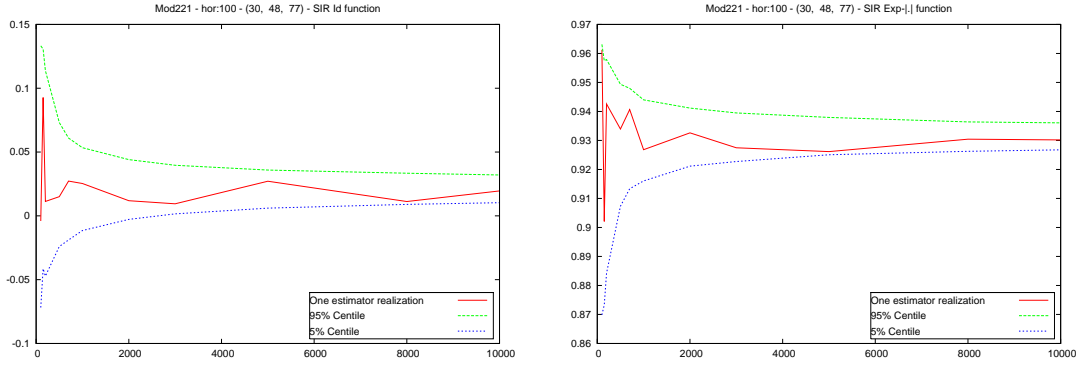


Figure 2: Particle filter approximations for SVM as a functions of particle number using SIR algorithm (left:  $\Pi_{100}f_1$ , right:  $\Pi_{100}f_2$ ) -  $(\beta, \sigma) = (0.995, 0.01)$  - Centiles over 4000 realizations.

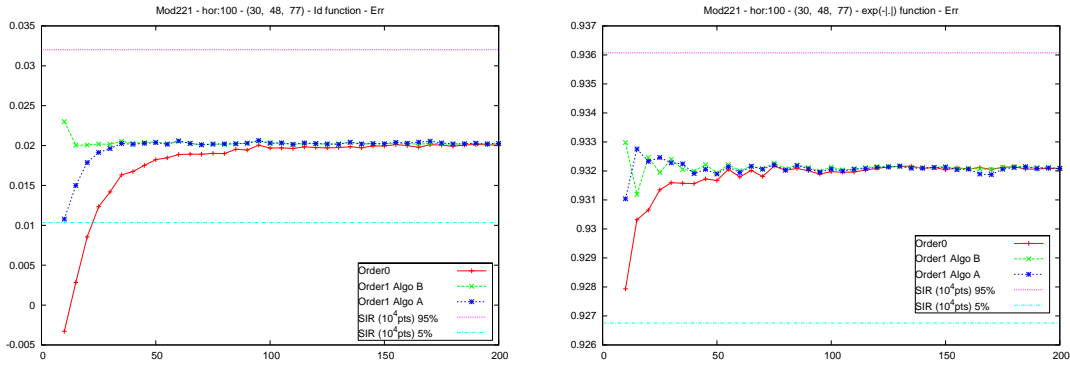


Figure 3: Quantization filter estimator as functions of quantizer size, in the SIR confidence interval with  $10^4$  particles (right:  $\hat{\Pi}_{100}f_1$ , left:  $\hat{\Pi}_{100}f_2$ ) -  $(\beta, \sigma) = (0.995, 0.01)$ .

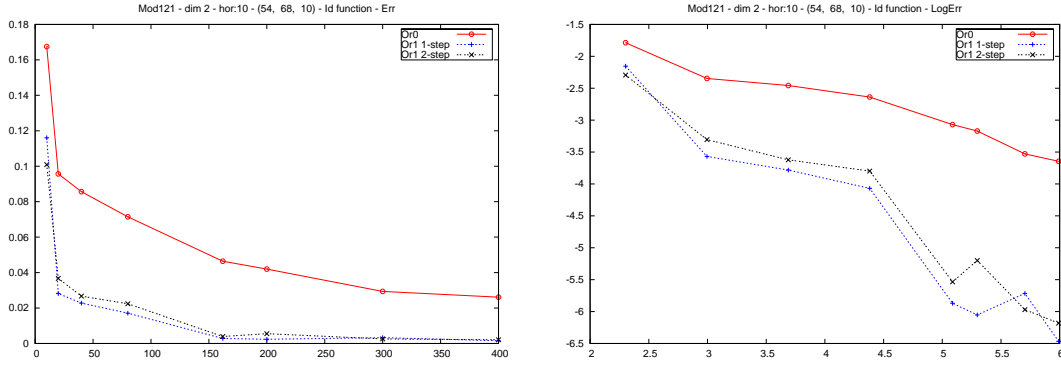


Figure 4: Quantization filter estimator errors for 2-dimensional Kalman case as a function of the quantizer size  $N_k$  (left:  $\|\Pi_{10}f_1 - \hat{\Pi}_{10}f_1\|_2$ , right: log-log scale representation).

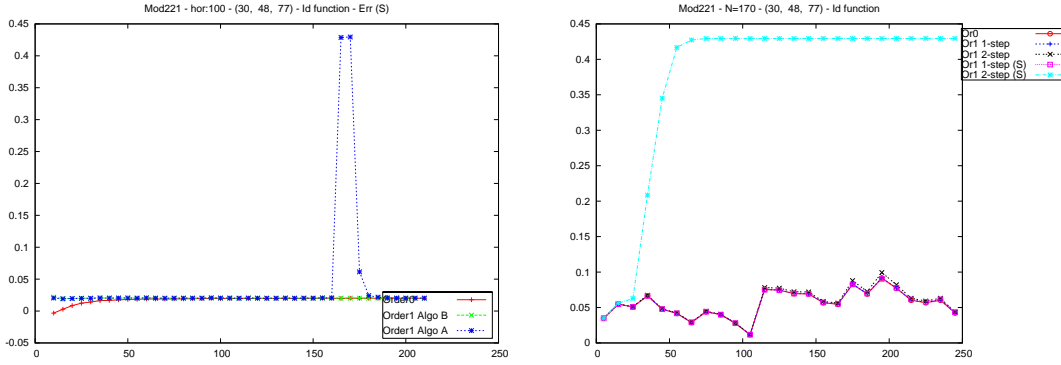


Figure 5: Quantization filter estimator for SVM using *intuitive* first order schemes as function of quantizer size (right:  $\|\hat{\Pi}_{100}f_1\|_2$  as a function of quantizer size  $N$  for  $n = 100$ , left:  $\|\hat{\Pi}_n f_1\|_2$  as a function of  $n$  for  $N = 170$ ) -  $(\beta, \sigma) = (0.995, 0.01)$ .

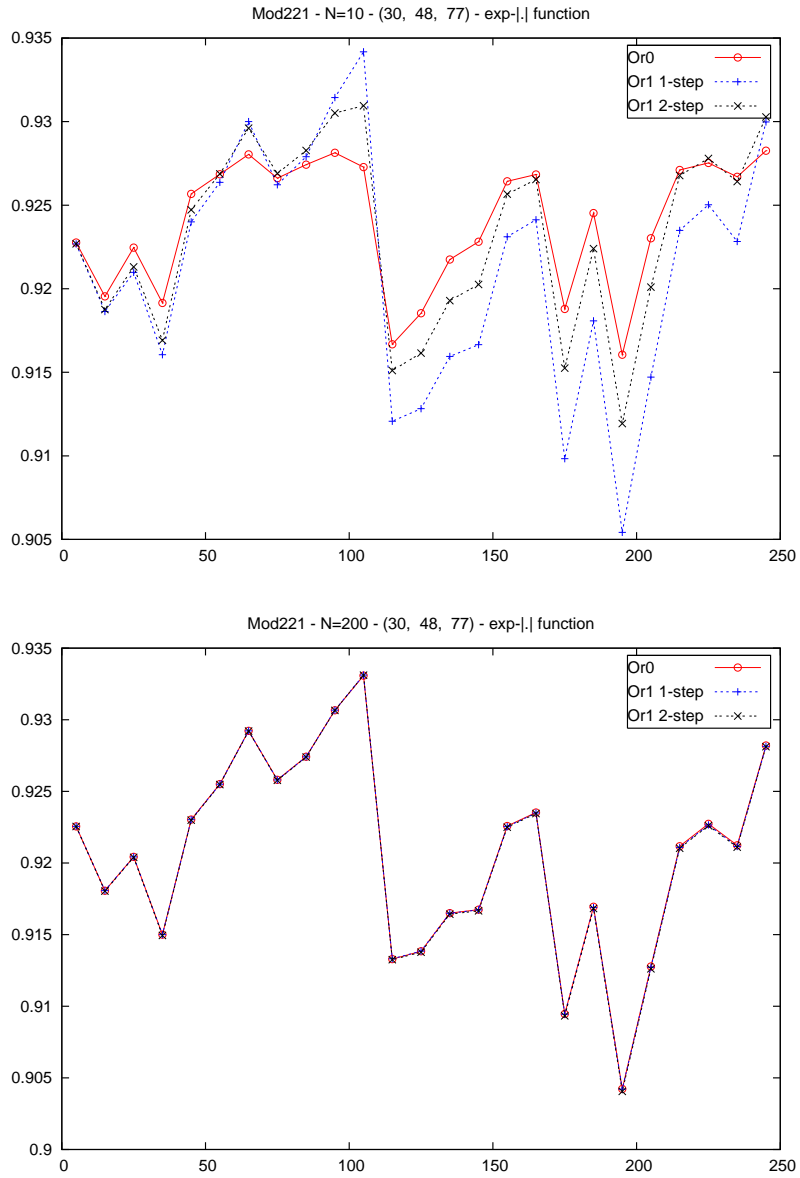


Figure 6: Horizon varying effect on quantization based filters for SVM (top:  $\hat{\Pi}.f_2$  for  $N_k = 10$  as a function of  $n$ , bottom:  $\hat{\Pi}.f_2$  for  $N_k = 200$  as a function of  $n$ ) -  $(\beta, \sigma) = (0.995, 0.01)$ .