

Univerzita Karlova v Praze
Matematicko-fyzikální fakulta

DIPLOMOVÁ PRÁCE



David Coufal

Sekvenční metody Monte Carlo

Katedra pravděpodobnosti a matematické statistiky

Vedoucí diplomové práce: prof. RNDr. Viktor Beneš, DrSc.

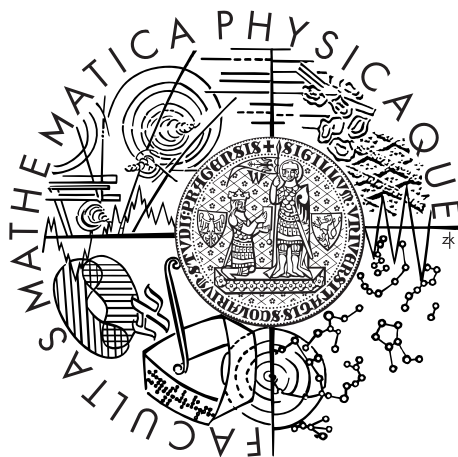
Studijní program: Matematika

Studijní obor: PMSE

Prague 2013

Charles University in Prague
Faculty of Mathematics and Physics

MASTER THESIS



David Coufal

Sequential Monte Carlo Methods

Department of Probability and Mathematical Statistics

Supervisor of the master thesis: prof. RNDr. Viktor Beneš, DrSc.

Study programme: Mathematics

Specialization: PMSE

Prague 2013

I would like to thank to my supervisor, Dr. Viktor Beneš, for providing me with an interesting topic of study and his helpful comments. I would like also to thank to the Institute of Computer Science of Academy of Sciences of the Czech Republic, as my home institute, for technical support during writing the thesis.

I declare that I carried out this master thesis independently, and only with the cited sources, literature and other professional sources.

I understand that my work relates to the rights and obligations under the Act No. 121/2000 Coll., the Copyright Act, as amended, in particular the fact that the Charles University in Prague has the right to conclude a license agreement on the use of this work as a school work pursuant to Section 60 paragraph 1 of the Copyright Act.

In Prague, 7.4.2013

David Coufal

Název práce: Sekvenční metody Monte Carlo

Autor: David Coufal

Katedra: Katedra pravděpodobnosti a matematické statistiky

Vedoucí diplomové práce: prof. RNDr. Viktor Beneš, DrSc.

Abstrakt: Práce shrnuje teoretické základy sekvenčních metod Monte Carlo se zaměřením na použití v oblasti částicových filtrů a základní výsledky z oblasti neparametrických jádrových odhadů hustot pravděpodobnostních rozdělání. Přehled výsledků tvoří základ ke zkoumání použití jádrových metod pro aproximaci hustot rozdělání částicových filtrů. Hlavními výsledky práce jsou důkaz konvergence jádrových odhadů k příslušným teoretickým hustotám a popis vývoje chyby aproximace v souvislosti s časovou evolucí filtru. Práce je doplněna experimentální částí demonstrující použití popsaných algoritmů formou simulací ve výpočetním prostředí MATLAB[®].

Klíčová slova: sekvenční metody Monte Carlo, částicové filtry, neparametrické jádrové odhady

Title: Sequential Monte Carlo Methods

Author: David Coufal

Department: Department of Probability and Mathematical Statistics

Supervisor: prof. RNDr. Viktor Beneš, DrSc.

Abstract: The thesis summarizes theoretical foundations of sequential Monte Carlo methods with a focus on the application in the area of particle filters; and basic results from the theory of nonparametric kernel density estimation. The summary creates the basis for investigation of application of kernel methods for approximation of densities of distributions generated by particle filters. The main results of the work are the proof of convergence of kernel estimates to related theoretical densities and the specification of the development of approximation error with respect to time evolution of a filter. The work is completed by an experimental part demonstrating the work of presented algorithms by simulations in the MATLAB[®] computational environment.

Keywords: sequential Monte Carlo methods, particle filters, nonparametric kernel estimates

To the Memory of Jan Wilda

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 1 |
| 1.1 | Motivation | 2 |
| 1.2 | Structure of the work and main results | 3 |
| 1.3 | Notation and typography | 5 |
| 2 | SMC methods | 7 |
| 2.1 | Monte Carlo methods | 7 |
| 2.2 | Importance Sampling | 9 |
| 2.3 | Sequential Importance Sampling | 12 |
| 2.4 | Resampling | 14 |
| 2.5 | SMC algorithm | 15 |
| 2.6 | Particle filters | 18 |
| 2.6.1 | Signal process | 18 |
| 2.6.2 | Observation process | 19 |
| 2.6.3 | Filtering distribution | 20 |
| 2.6.4 | SMC algorithm for particle filters | 23 |
| 2.6.5 | Recursive evolution of filtering distributions | 25 |
| 2.7 | Convergence results for particle filters | 29 |
| 3 | Kernel methods | 37 |
| 3.1 | Histograms | 38 |
| 3.2 | Kernel methods in one dimension | 41 |
| 3.2.1 | MISE analysis | 43 |
| 3.2.2 | AMISE analysis | 48 |
| 3.3 | Kernel methods in multiple dimensions | 49 |
| 3.3.1 | MISE analysis | 51 |
| 3.3.2 | AMISE analysis | 51 |
| 3.4 | Fourier analysis of kernel estimators | 53 |
| 3.4.1 | The first term of Fourier MISE decomposition | 56 |
| 3.4.2 | Other terms of Fourier MISE decomposition | 60 |
| 3.4.3 | Upper bound on the Fourier MISE formula | 60 |
| 4 | SMC and kernel methods | 63 |
| 4.1 | Convergence of SMC kernel estimates | 63 |
| 4.2 | Sobolev character of SMC particle filters | 68 |
| 4.3 | Lower bound on $\bar{\pi}_t g_t$ | 72 |
| 4.3.1 | Exact bound | 73 |

| | | |
|----------|---|------------|
| 4.3.2 | Sequential computation of $(\bar{\pi}_t g_t)^*$ | 79 |
| 4.3.3 | Approximate bounds | 81 |
| 5 | Experiments | 85 |
| 5.1 | Univariate Gaussian process | 85 |
| 5.1.1 | Univariate Kalman filter | 86 |
| 5.1.2 | Univariate Gaussian SMC filter | 86 |
| 5.1.3 | Sobolev character of univariate filter | 88 |
| 5.1.4 | Properties of univariate Gaussian kernel | 88 |
| 5.1.5 | MATLAB implementation | 89 |
| 5.1.6 | Experiments with univariate Gaussian SMC filter | 90 |
| 5.2 | Multivariate Gaussian process | 93 |
| 5.2.1 | Multivariate Kalman filter | 93 |
| 5.2.2 | Multivariate Gaussian SMC filter | 94 |
| 5.2.3 | Sobolev character of multivariate filter | 96 |
| 5.2.4 | Properties of multivariate Gaussian kernel | 97 |
| 5.2.5 | MATLAB implementation and experiments | 98 |
| 6 | Summary | 103 |
| A | MATLAB source codes | 105 |
| A.1 | uvsmc.m | 105 |
| A.2 | mvsmc.m | 108 |

List of symbols

- $\mathbb{N}, \mathbb{N}_0, \mathbb{Z}, \mathbb{R}, \mathbb{C}$ - sets of natural, $\mathbb{N} \cup \{0\}$, integer, real and complex numbers
- \mathbb{R}^d - d -dimensional Euclidean space, $\mathbb{R}^d = \mathbb{R} \times \cdots \times \mathbb{R}$
- $\mathcal{B}(\mathbb{R}^d)$ - σ -algebra of Borel subsets of \mathbb{R}^d
- $|x|$ - absolute value of a real or complex number x
- $\|f\|_\infty = \sup_{\mathbf{x}} |f(\mathbf{x})|$ - supremum norm of function f
- $\|X\| = (\mathbb{E}|X|^p)^{1/p}$ - L_p norm of random variable X
- L_1, L_2 - the class of L_1 or L_2 integrable random variables
- $B(\mathbb{R}^d), B(\mathbb{C}^d)$ - set of bounded real functions over $\mathbb{R}^d, \mathbb{C}^d$
- $C_b(\mathbb{R}^d), C_b(\mathbb{C}^d)$ - set of continuous and bounded functions over $\mathbb{R}^d, \mathbb{C}^d$
- $\delta(X_i) = \delta_{x_i}(dx)$ - the Dirac measure determined by $X_i(\omega) = x_i$
- $K_{t-1}(A, x_{t-1})$ - transition kernels of a Markov process for $t \geq 1$
- $g_t(y_t|x_t)$ - conditional densities of an observation process
- $\pi_{0:t|t}$ - conditional distribution of states $X_{0:t}$ given observations $Y_{1:t}$
- $\pi_{t|t}$ - marginal conditional distribution of X_t given observations $Y_{1:t}$
- $(\cdot)_+ = \max\{0, \cdot\}$ - positive part
- $\|\mathbf{x}\| = (\sum_{i=1}^d x_i^2)^{1/2}$ - the Euclidean norm of vector $\mathbf{x} \in \mathbb{R}^d$
- $\|f\|$ - L_2 norm of real function f , i.e., $\|f\| = (\int f(\mathbf{x}) d\mathbf{x})^{1/2}$
- $\langle \boldsymbol{\omega}, \mathbf{x} \rangle = \sum_{j=1}^d \omega_j x_j$ dot product of vectors $\boldsymbol{\omega}$ and \mathbf{x} from \mathbb{R}^d
- $L_1(\mathbb{R}^d), L_2(\mathbb{R}^d)$ - the class of $L_1(\mathbb{R}^d)$ or $L_2(\mathbb{R}^d)$ integrable real functions
- $\mathcal{N}(\mu, \sigma^2)$ - univariate normal distribution
- $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ - multivariate normal distribution
- \mathbf{I}_d - unit matrix of size $d \in \mathbb{N}$
- $\|\mathbf{A}\|_{spc}$ - spectral norm of matrix \mathbf{A}

1. Introduction

The research on description and control of dynamical systems constitutes a broad and live area. In the deterministic setup it is tightly connected with the development of calculus and solution of systems of generally partial differential equations. The level of rigor reached in the probability theory in 30-ties and 40-ties of the last century brought a strong theoretical background into the area of description of stochastic dynamical systems. An eruption of new ideas and powerful algorithms (e.g., Shannon's information theory, Bellman's dynamic programming, Kalman filtering, etc.) was tightly related to ongoing advances in semi-conductor industry in 50-ties and massive space research in 60-ties.

70-ties and generally the last third of 20-th century is marked by an enormous increase of computational power in terms of efficiency and costs of computation. This progress enables to maintain and process huge portion of data monitored and collected from processes of interest. Computers allow the application of mathematical methods which yield reasonable results only when a brute computational power is employed. Moreover, new paradigma have emerged based on the possibility of massive computer simulations.

Monte Carlo methods (MC methods) are empirical, but theoretically well-founded methods of convenient representation and manipulation of integral characteristics of complex measures which are hard to be processed analytically. MC methods rest on the strong law of large numbers which guarantees the soundness of their use. Basically, MC methods enables an effective computation (with a help of a computer) of integral characteristics of random variables of complicated distributions by generating samples from these distributions and averaging them to approximate related theoretical values.

The *nonparametric density estimation* is today well-established area of mathematical statistic. Basic works date back to 50-ties and 60-ties. A typical setup of the problem is that we are provided by a sample of data without apriori knowledge on the distribution driving the sample and we want to represent this distribution in terms of a density. Originally, these approaches were based on studying histograms of empirical data, but these as nonsmooth objects were replaced by kernels, which results into what is now known as *kernel methods* of nonparametric density estimation.

In practice, stochastic processes are of great interest to researchers as they are able to accommodate stochastic behavior of real dynamic systems. Of course, a certain level of simplification has to be applied when stating assumptions

on modeled systems. In the basic setup, a dynamic system evolves in discrete time steps and therefore the relevant stochastic process is considered as a chain. Markov chains are one of the most theoretically elaborated processes so it is natural that the application of Monte Carlo methods is studied in this context. The related research falls into the area of so-called *particle filters*.

In particle filters, empirical measures represented by groups of samples are generated sequentially to simulate time evolution of theoretical measures structurally assumed to form a Markov chain. *Sequential Monte Carlo methods* (SMC methods) enable to generate samples in an effective way when previous samples are parts of actual samples and a computational effort is reduced. However, the sequential setup brings new questions on convergence of resulting empirical measures and corresponding integral characteristics as a direct application of the strong law of large number is not possible here due to time evolution of a filter. These questions were solved successfully in terms of theorems proving the weak convergence of empirical measures to their theoretical counterparts, but what about the related densities?

To be precise, is there a constructive way how to create convergent approximations of densities of theoretical distributions we are interested in? This is the question which motivated my work presented in the thesis.

1.1 Motivation

In Monte Carlo methods one knows the probabilistic structure of a problem and wants to generate data according to this structure in order to compute integral characteristics of its interest. In a density estimation problem, we encounter the reverse task as we have empirical data at our disposal, but we do not know their probabilistic structure in terms of a density of their distribution. Knowing the density of a distribution is worth because it enables a convenient manipulation with the distribution, especially if the density representation has a suitable form as it is the case for kernel estimates. The canonical example when densities are employed is the computation of conditional expectations serving as regression functions.

The reason for employment of SMC methods is that they produce empirical distributions representing in limit case some theoretical distributions of interest. In the context of particle filters these theoretical distributions are called *filtering distributions*. The representation in the form of an empirical measure is suitable for computation of integral characteristics, but for other purposes it

would be worth to have also analytical representation in the form of a density. We are not able to state the density of a filtering distribution apriori - actually, this the reason why we employ SMC methods - but we can do it posteriori on the basis of created empirical distributions by means of nonparametric kernel density estimation.

It may seem that forming a kernel density estimate of a filtering distribution constitutes two consequent but independent tasks, however, this is not the case when the sequential setup of SMC methods is taken into account. The reason is that samples lose i.i.d. character during time evolution of a filter and therefore a research has to be taken on the behavior of kernel density estimates in the SMC context. The following questions are relevant here:

- Do kernel density estimates for SMC methods converge to the theoretical densities of related filtering distributions even though that the estimates are not based on i.i.d. samples?
- If so, what is the rate of convergence with respect to the number of samples employed?
- What is the effect of time evolution on the error of estimates.
- Finally, what is the relationship between the properties of Markov chain underlying a filter and the properties of the kernel underlying the estimate in order to the error of estimate can be established in a reasonable way?

The answers to the above questions represent the results of my study of combination of sequential Monte Carlo methods and kernel methods of nonparametric density estimation in the context of particle filters.

1.2 Structure of the work and main results

The thesis comprises of five chapters. After the first introductory chapter, Chapter 2 overviews the basic theory of Monte Carlo methods designed to work in the sequential setup. Mathematical backgrounds and algorithmic description of SMC methods when working in the particle filters context is presented. The special emphasis is put on the review of convergence results for this class of algorithms.

The chapter, as the overview one, is primarily written on the basis of comprehensive book by [Doucet et al. 2001], especially its theoretical part. Explanations and descriptions presented in journal papers [Crisan and Doucet 2002; Andrieu et al. 2010] were employed as well. The theory of particle filters is presented at the basic level in [Fristedt et al. 2007].

Chapter 3 summarizes basic results obtained in the area of kernel density estimation. The chapter starts by inspection of histograms and follows historical path of generalization to kernel methods. Univariate and multivariate methods are treated separately. The discussion is provided on the behavior of asymptotic error of approximation in both cases. The important part of the chapter is formed by the section related to the Fourier analysis of kernel estimates as the referred results are employed in our own research.

Classical works in the field of nonparametric estimations were used to compile Chapter 3. These are textbooks of [Silverman 1986; Tarter and Lock 1993] and [Wand and Jones 1995]. Multivariate kernel estimation is widely treated in [Scott 1992]. The Fourier analysis section is based on the book by [Tsybakov 2009].

The fourth chapter consists of main results of the thesis. First, there is proved the theorem showing that convergence of kernel density estimates to the theoretical densities of filtering distributions is retained. The rate of convergence with respect to the number of employed samples is quantified altogether with the description of time evolution of related error. Further, there is investigated the relationship between properties of Markov chain underlying the filter and the order of the kernel underlying the density estimation. Finally, exact and approximate bounds on the value of the certain normalizing integral are discussed as the integral affects the value of the error of density estimate.

Chapter 5 constitutes the experimental part of the thesis. In this part results of simulations performed with SMC methods are presented and compared with results obtained by Kalman filtering. The MATLAB[®] computational environment was employed for this purpose.

The thesis is concluded by Chapter 6 summarizing the performed work and obtained results in the context of the assignment of the thesis.

1.3 Notation and typography

The mathematical symbols employed through the thesis are summarized in the related list. Special symbols and notation are defined and mentioned before they are used for the first time.

Through the thesis, *italic* is used to emphasize new or important concepts.

Definitions, theorems and lemmas are also typeset in italic. The ends of proofs are denoted by the \square symbol.

The program codes written in the MATLAB[®] scripting language are printed in the `monospace font`.

The thesis is typeset in L^AT_EX_{2 ϵ} , MiKTeX 2.9 distribution, using standard Computer Modern Fonts with several enhancements, especially $\mathcal{A}\mathcal{M}\mathcal{S}$ -L^AT_EX mathematical symbols.

2. SMC methods

In this chapter we review the basics of *sequential Monte Carlo methods* (SMC methods), their application in the area of *particle filters* and we focus on relevant convergence results. We start by recalling the general idea of *Monte Carlo methods* (MC methods). Then we introduce the technique of *importance sampling* and its sequential version to enhance and make comfortable the application of MC methods in sequential setup. Incorporating the *resampling* step finishes the specification of the work of SMC methods. The algorithmic form of the specification is provided as well. Finally, we discuss the concept of *particle filters* which represents one of the most prominent areas of application of SMC methods. We show that the application of SMC algorithm is correct in the sense that increasing the number of samples (particles) leads to convergence of empirical measures and integrals to their theoretical counterparts.

In the review we follow the standard literature in the field. We have mainly employed [Doucet et al. 2001] and [Crisan and Doucet 2002; Andrieu et al. 2010; Fristedt et al. 2007].

2.1 Monte Carlo methods

Very roots of Monte Carlo methods stem from the strong law of large numbers (SLLN). As it is well-known, the law guarantees under certain assumptions that integral characteristics of a random variable can be approximated by averaging over empirical samples from its distribution. In the limit, the approximations coincide with the theoretical characteristics if the number of samples goes to infinity.

The reference to the city of Monte Carlo points out the process of random sampling, which constitutes the heart of the methods, in a reminiscence to random sampling a roulette wheel and to the location of the world's most famous casino established in Monaco in the second half of the 19-th century.

The statement of SLLN for a sequence of independent identically distributed (i.i.d.) random variables reads as follows.

Theorem 2.1. (Kolmogorov's SLLN) Let $\{X_n, n \geq 1\}$ be an i.i.d. sequence of random variables and set $S_n = \sum_{i=1}^n X_i$. There exists $c \in \mathbb{R}$ such that

$$\bar{X}_n = S_n/n \xrightarrow[n \rightarrow \infty]{} c \quad \text{a.s.}$$

iff $\mathbb{E}(|X_1|) < \infty$ in which case $c = \mathbb{E}(X_1)$.

Proof. See [Resnick 1999], p. 220.

As a corollary, for a Borel function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, $d \in \mathbb{N}$, and an i.i.d. sequence of d -dimensional random vectors $\{X_i\}_{i=1}^\infty$, we have that if $\mathbb{E}[|f(X_1)|] < \infty$, i.e., if $f(X_1) \in L_1$, then

$$\hat{I}_n(f) = \frac{1}{n} \sum_{i=1}^n f(X_i) \xrightarrow[n \rightarrow \infty]{} \int f(X_1) dP_{X_1} = I(f) \quad \text{a.s.} \quad (2.1)$$

and if further $\mathbb{E}[f(X_1)^2] < \infty$, i.e. if $f(X_1) \in L_2$, then

$$\frac{1}{n} \sum_{i=1}^n (f(X_i) - \hat{I}_n(f))^2 \xrightarrow[n \rightarrow \infty]{} \text{var}(f(X_1)) \quad \text{a.s.} \quad (2.2)$$

Formula (2.1) says that integral characteristics of X_1 for an integrable function f can be computed by means of random sampling from distribution of X_1 . Similarly the variance of $f(X_1)$ can be approximated by the sample variance with the desired convergence property as the number of samples goes to infinity.

Approximation $\hat{I}_n(f)$ is a random variable with the expected value $\mathbb{E}[f(X_1)]$. Thus, $\hat{I}_n(f)$ represents an unbiased consistent estimate of the corresponding expected value. Under the assumptions of the above corollary, the distribution of the estimate is asymptotically normal due to the central limit theorem, [Resnick 1999], p. 312,

$$\sqrt{n} \cdot (\hat{I}_n(f) - I(f)) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \mathcal{N}(0, \text{var}(f(X_1))). \quad (2.3)$$

Formula (2.1) can be reformulated in the framework of application of empirical random measures. Having a realization of an i.i.d. sequence $\{X_i(\omega) = x_i\}_{i=1}^\infty$ for $x_i \in \mathbb{R}^d$, we can associate each group of samples x_1, \dots, x_n with the empirical measure given as the uniformly weighted sum of corresponding Dirac measures

$$\delta_n(dx) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}(dx) = \frac{1}{n} \sum_{i=1}^n \delta(X_i). \quad (2.4)$$

The second expression points out the random character of $\delta_n(dx)$. Integration with respect to this empirical measure brings the entity

$$\int f(x) \delta_n(dx) = \frac{1}{n} \sum_{i=1}^n f(X_i) = \widehat{I}_n(f), \quad (2.5)$$

i.e., the same approximate integral $\widehat{I}_n(f)$ as in (2.1).

Concerning in (2.1) all continuous and bounded functions on \mathbb{R}^d , i.e., $f \in C_b(\mathbb{R}^d)$, we get from the assertion of SLLN the weak convergence of empirical measures $\delta_n(dx)$ to the distribution of X_1 a.s. [Billingsley 1995]. That is why we will treat $\delta_n(dx)$ empirical measures as approximations of the distribution the samples are taken from.

Due to the above reformulation, Monte Carlo methods can be understood primarily as the tool for approximation of distributions of random variables and only subsequently as a tool for approximation of integral characteristics of these distributions. Integral characteristics are then obtained via integration with respect to the presented empirical measures. Integrals, which are in fact sums, then converge with the increasing number of samples to theoretical entities.

An important observation is, that when using MC methods the dimensionality does not play negative role. This is very important because the curse of dimensionality makes serious problems in different branches of applied mathematics. Of course, the employment of SLLN as a limit law needs computational effort, but this is linear in number of samples, not exponential with increasing dimension.

The crucial assumption for using MC methods is that we are able to generate i.i.d. samples from the distribution of a certain random variable. An elegant approach which transforms the problem of sampling from a given “possibly weird” distribution to sampling from a “comfort distribution” is addressed in the next section.

2.2 Importance Sampling

The idea of importance sampling deals with the problem of how to sample from a given distribution $\pi(dx)$ possessing a density with respect to a certain basic measure. Typically, the basic measure is the corresponding (in sense of

dimension) Lebesgue measure. We denote the density of $\pi(dx)$ by $p(x)$, i.e., $\pi(dx) = p(x) dx$.

In what follows, we will not associate distributions with concrete random variables anymore. As a consequence we further denote distributions by small Greek letters, e.g., $\pi(dx)$, instead of P_X . Relevant densities will be then denoted by small Latin letters, e.g., by $p(x)$.

In real applications, it is a common case that the density and consequently the distribution is specified only up to a positive normalizing constant. This fact is referred to by the standard notation $p(x) \propto p^*(x)$ or $\pi(dx) \propto \pi^*(dx)$, where $p^*(x)$ and $\pi^*(dx)$ are known unnormalized density and distribution, respectively. Denoting the normalizing constant Z , $Z > 0$, we have $p(x) = Z^{-1}p^*(x)$ and $\pi(dx) = Z^{-1}\pi^*(dx)$.

The main idea of importance sampling is instead to sample directly from the distribution of interest $\pi(dx)$, which is from some reason uncomfortable or even impossible, to sample from the other so-called *proposal distribution* (or *importance sampling distribution*) where we are able to do this. A proposal distribution is again considered to be characterized by its density with respect to the basic measure. We denote this density by $q(x)$. A technical condition is that if $p(x) > 0$ then also $q(x) > 0$.

If $p(x)$ is specified, i.e., the normalizing constant is known, we have for any function f which is integrable with respect to $\pi(dx)$,

$$\mathbb{E}_\pi[f] = \int f(x)p(x) dx = \int f(x)\frac{p(x)}{q(x)}q(x) dx. \quad (2.6)$$

However, if $p(x)$ is known only up to a normalizing constant, then the formula is extended to

$$\mathbb{E}_\pi[f] = \int f(x)p(x) dx = \int f(x)\frac{p(x)}{q(x)}q(x) dx = \frac{\int f(x)\frac{p^*(x)}{q(x)}q(x) dx}{\int \frac{p^*(x)}{q(x)}q(x) dx}. \quad (2.7)$$

Clearly, the unknown normalizing constant is replaced by the term in the denominator.

Further, defining $w(x) = \frac{p^*(x)}{q(x)}$, we obtain (2.7) in more compact form of

$$\mathbb{E}_\pi[f] = \frac{\int f(x)w(x)q(x)dx}{\int w(x)q(x)dx} \quad (2.8)$$

where $w(x)$ are known as the *importance weights*.

If we are able to sample from $q(x)$ then the problem of specification of $\mathbb{E}_\pi[f]$ can be considered as an instance of MC methods for function $f(x)w(x)$ and density $q(x)$. The MC estimate of $\mathbb{E}_\pi[f]$ writes as

$$\widehat{I}_n(f) = \frac{\frac{1}{n} \sum_{i=1}^n f(x_i)w(x_i)}{\frac{1}{n} \sum_{i=1}^n w(x_i)} = \sum_{i=1}^n f(x_i)\tilde{w}(x_i) \quad (2.9)$$

where $\tilde{w}(x_i)$ are *normalised importance weights* computed as

$$\tilde{w}(x_i) = \frac{w(x_i)}{\sum_{j=1}^n w(x_j)}. \quad (2.10)$$

The price we pay for the unspecification of a normalizing constant is that we compute the ratio of two approximations which results in a no more unbiased estimate. However, as pointed out in [Doucet et al. 2001], asymptotically SLLN still applies, i.e., $\widehat{I}_n(f) \rightarrow \mathbb{E}_\pi[f]$ a.s., as $n \rightarrow \infty$, and also the convergence rate is still independent of the dimension.

Formula (2.9) can be again reformulated in terms of application of an empirical measure. Having an i.i.d. sequence of samples from a proposal distribution with density q , then defining

$$\widehat{\pi}_n(dx) = \sum_{i=1}^n \tilde{w}(x_i) \delta_{x_i}(dx) = \sum_{i=1}^n \tilde{w}(X_i) \delta(X_i) \quad (2.11)$$

we get $\widehat{\pi}_n(dx)$ as a non-uniformly weighted sum of Dirac measures, which is again a random measure. It is a probability measure because $\sum_i \tilde{w}(x_i) = 1$. Formula (2.9) can be then rewritten as the integral of f with respect to measure $\widehat{\pi}_n(dx)$ because

$$\int f(x) \widehat{\pi}_n(dx) = \sum_{i=1}^n \tilde{w}(x_i) f(x_i) = \widehat{I}_n(f). \quad (2.12)$$

The presented importance sampling technique shows the practical way how to compute approximations of integrals with respect to distributions which we are not able to sample from directly. At the basement, there is an approximation of the distribution of interest by a non-uniformly weighted empirical measure. In the following section we show how this technique can be enhanced when a sequence of integrals and therefore measures is needed to be computed and how the basic importance sampling method can be redesigned to fit the sequential setup.

2.3 Sequential Importance Sampling

The sequential setup of the importance sampling method is related to the problem of computing MC approximations with respect to a sequence of probability spaces of incrementing dimension. An effort here is to somehow beat the curse of dimensionality which is manifested by exponentially increasing number of required computations. The problem can be tackled by assuming certain structural relations on evolving spaces.

Let us consider a sequence of measurable spaces $\{(E_t, \mathcal{E}_t), t \in \mathbb{N}_0\}$, such that $(E_0, \mathcal{E}_0) = (E, \mathcal{E})$ and $(E_t, \mathcal{E}_t) = (E_{t-1}, \mathcal{E}_{t-1}) \otimes (E, \mathcal{E})$ for $t \geq 1$. Further we assume the existence of a related sequence of probability measures π_t specified on (E_t, \mathcal{E}_t) . Thus, we finally consider a sequence $\{(E_t, \mathcal{E}_t, \pi_t), t \geq 0\}$ of probabilistic spaces of increasing dimension as t increments.

Our aim is to build up sequentially approximations of π_t , denoted by $\hat{\pi}_t$, implemented as empirical measures and determined by random samples from appropriate proposal distributions.

The standard choice of (E, \mathcal{E}) is $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, i.e., E corresponds to the $d \in \mathbb{N}$ dimensional Euclidean space and \mathcal{E} to its Borel sigma-field. Incrementing t we have $E_t = (\mathbb{R}^d)^{t+1} = \mathbb{R}^{d(t+1)}$ and $\mathcal{E}_t = \otimes_{k=0}^t \mathcal{B}(\mathbb{R}^d) = \mathcal{B}(\mathbb{R}^{d(t+1)})$. Elements of $E_t = \mathbb{R}^{d(t+1)}$ will be denoted $x_{0:t}$.

We assume that measures $\pi_t(dx_{0:t})$ admit densities $p(x_{0:t})$ with respect to the Lebesgue measures $(\lambda^d)^{t+1}$, i.e., $\pi_t(dx_{0:t}) = p_t(x_{0:t}) dx_{0:t}$. The densities $p(x_{0:t})$ are assumed to be known up to normalizing constants Z_t , $p_t(x_{0:t}) = Z_t^{-1} p_t^*(x_{0:t})$, $p_t(x_{0:t}) \propto p_t^*(x_{0:t})$.

For fixed t , the approximation $\hat{\pi}_t^n(dx_{0:t})$ of $\pi_t(dx_{0:t})$ follows the importance sampling formula (2.11). For a given sample $\{X_{0:t}^i\}_{i=1}^n$ (from now on we use the lower index to reflect the order in the sequence and the upper index for indexing individual samples i or to referring to the number of samples n) we have

$$\hat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(x_{0:t}^i) \delta_{x_{0:t}^i}(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(X_{0:t}^i) \delta(X_{0:t}^i). \quad (2.13)$$

The idea of sequential setup is to somehow reuse samples $X_{0:t-1}^i$ when moving one step ahead to index t without increasing the computational effort due to the move to the higher dimension as t increases. More specifically, the sequencing is based on the decomposition of proposal distribution to a stream of conditional distributions. That is, we assume that density $q_t(x_{0:t})$ can be

decomposed as

$$q_t(x_{0:t}) = q_0(x_0) \prod_{k=1}^t q_{k-1}(x_k | x_{0:k-1}). \quad (2.14)$$

The conditional structure enables us to obtain a new sample from (E_t, \mathcal{E}_t) by first sampling x_t^i from conditional distribution $q_{t-1}(dx_t | x_{0:t-1})$, $t \geq 1$, and then constituting new sample $x_{0:t}^i$ by composition of x_t^i with $x_{0:t-1}^i$, i.e., $x_{0:t}^i = (x_{0:t-1}^i, x_t^i)$. Hence, we reuse $x_{0:t-1}^i$ twice in order to setup $x_{0:t}^i$ on the basis of $x_{0:t-1}^i$.

Let us incorporate the assumption on conditional decomposition of the proposal density into the general importance sampling formula. We have

$$\begin{aligned} w(x_{0:t}^i) &= \frac{p_t^*(x_{0:t}^i)}{q_t(x_{0:t}^i)} = \frac{p_t^*(x_{0:t}^i)}{q_{t-1}(x_t^i | x_{0:t-1}^i) q_{t-1}(x_{0:t-1}^i)} \\ &= \frac{p_{t-1}^*(x_{0:t-1}^i)}{q_{t-1}(x_{0:t-1}^i)} \cdot \frac{p_t^*(x_{0:t}^i)}{p_{t-1}^*(x_{0:t-1}^i) q_{t-1}(x_t^i | x_{0:t-1}^i)} \\ &= w(x_{0:t-1}^i) \cdot \frac{p_t^*(x_{0:t}^i)}{p_{t-1}^*(x_{0:t-1}^i) q_{t-1}(x_t^i | x_{0:t-1}^i)}. \end{aligned} \quad (2.15)$$

In the last product, the first term is the unnormalized weight for the i -th sample at time $t - 1$ and the second is an update factor based on both the present $x_t^i \sim q_{t-1}(dx_t | x_{0:t-1}^i)$ and the previous sample $x_{0:t-1}^i$, $x_{0:t}^i = (x_{0:t-1}^i, x_t^i)$. From importance weights we easily compute the normalized versions $\tilde{w}(x_{0:t}^i)$. Hence we update all needed ingredients in the importance sampling formula sequentially.

The sequential setup has its advantages. With an increasing dimension of probabilistic spaces we do not need to sample from the distributions of increasing complexity. We instead start from an initial distribution and go on sequentially with samples from conditional distributions of given complexity which are then used to extend previously obtained samples.

Unfortunately, as t increases the phenomenon of degeneracy arises. It is manifested by the fact that the majority of normalized weights goes to zero and only few have negligible value [Doucet et al. 2001]. The standard technique to overcome this problem is to shift from weighted empirical measures to uniformly weighted ones. This is performed practically by introducing the resampling step into the sampling algorithm.

2.4 Resampling

The incorporation of the *resampling step* is enforced by the phenomenon of degeneracy of weighted empirical distribution as t (dimension) increases. The origin of the degeneracy lies in the increasing variance of weights. To overcome the problem, the proposal is to reset the approximation by excluding the samples with small weights and multiply samples with large weights (relative to $1/n$).

Importance sampling in sequential setup provides the approximation of distribution of interest $\pi_t(dx_{0:t})$ in the form

$$\widehat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(x_{0:t}^i) \delta_{x_{0:t}^i}(dx_{0:t}). \quad (2.16)$$

To proceed, let us make the following important notation agreement. From now on we denote the samples employed in the importance sampling formula with the *bar notation*. That is, instead of x_t^i we use \bar{x}_t^i and similarly for $x_{0:t}^i$ and $\bar{x}_{0:t}^i$. The *plain notation* will be reserved for samples and measures resulting from the resampling step.

The reason for this change of notation is that by incorporating resampling step we start to have another empirical measure established and it must be somehow distinguished. Because the final measure of interest is the empirical measure resulting from the resampling step we reserve the plain notation for its samples. The carriers (samples) of intermediate measures (in the next section we see that they are in fact two) will be denoted by the bar notation. So saying it simply, *samples before resampling are denoted using the bar notation and after resampling by the plain notation*.

In the view of this agreement, the formula (2.16) rewrites as

$$\widehat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(\bar{x}_{0:t}^i) \delta_{\bar{x}_{0:t}^i}(dx_{0:t}). \quad (2.17)$$

The idea of resampling is to sample $\{X_{0:t}^i\}_{i=1}^n \sim \widehat{\pi}_t(dx_{0:t})$ to obtain an unweighted (more specifically uniformly weighted) empirical measure

$$\pi_t^n(dx_{0:t}) = \frac{1}{n} \sum_{i=1}^n \delta_{x_{0:t}^i}(dx_{0:t}) = \frac{1}{n} \sum_{i=1}^n \delta(X_{0:t}^i). \quad (2.18)$$

Formula (2.17) can be rewritten as

$$\widehat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \widetilde{w}(\bar{x}_{0:t}^i) \delta_{\bar{x}_{0:t}^i}(dx_{0:t}) = \sum_{i=1}^n \frac{\mathbb{E}[N_t^i]}{n} \delta_{\bar{x}_{0:t}^i}(dx_{0:t}) \quad (2.19)$$

where $(N_t^1, \dots, N_t^n) \sim \mathcal{M}(n, \widetilde{w}(\bar{x}_{0:t}^1), \dots, \widetilde{w}(\bar{x}_{0:t}^n))$. That is, vector (N_t^1, \dots, N_t^n) is sampled from the multinomial distribution of presented parameters. The reformulation of (2.19) is based on the property of the multinomial distribution stating (in the case of selected parameters) that $\mathbb{E}[N_t^i] = n \cdot \widetilde{w}(\bar{x}_{0:t}^i)$.

Now, the idea of resampling is straightforward. In order to obtain uniformly weighted samples from $\widehat{\pi}_t(dx_{0:t})$ we sample from the specified multinomial distribution and according to the values of N_t^i we replicate N_t^i -times sample $x_{0:t}^i$. If $N_t^i = 0$, then the sample is excluded from the population and it is not used anymore. Clearly, the number of samples in the population is retained. Denoting points of resampled population by $x_{0:t}^i$, or more generally as the realizations of $X_{0:t}^i$ random variables, we obtain the uniformly weighted measure π_t^n as required in formula (2.18).

Resampling leads to an unweighted empirical measure but **samples are no more independent**. However, integral characteristics are still unbiased, that is

$$\mathbb{E}_{\pi_t^n}[f] = \mathbb{E}_{\widehat{\pi}_t^n}[f]. \quad (2.20)$$

The multinomial resampling increases the variance with respect to integral characteristics, i.e., $\text{var}_{\pi_t^n}(f) \geq \text{var}_{\widehat{\pi}_t^n}(f)$. This fact led to the introduction of other resampling techniques as uniform resampling or tree-based branching mechanism [Doucet et al. 2001]. Some of them have equality in the variance formula. However, in the thesis we further consider only the multinomial resampling as it constitutes the basic resampling technique employed in SMC methods.

The incorporation of the resampling step finishes the theoretical specification of individual parts of SMC algorithm. In the next section we present its algorithmic description.

2.5 SMC algorithm

The algorithm of sequential Monte Carlo methods (SMC algorithm) consists of three blocks - *initialization*, *the importance sampling step* and *the resampling step*. Initialization is performed once at the start of the algorithm. Importance

sampling and resampling steps are performed in a loop reflecting the evolution of approximated distributions from $t = 1$ to $t = T$ where $T \in \mathbb{N}$ is the selected computational horizon. The pseudocode of the algorithm writes as follows.

- **0. declarations**

n - number of samples,
 T - computational horizon,
 $p_t^*(x_{0:t})$ - unnormalized densities of interest,
 $q_0(x_0)$ - the density of an initial proposal distribution,
 $q_{t-1}(x_t|x_{0:t-1})$, $t = 1, \dots, T$, - conditional proposal densities.

- **1. initialization**

$t = 0$,
sample $\{\bar{x}_0^i \sim q_0(dx_0)\}_{i=1}^n$,
compute initial importance weights $\left\{w(\bar{x}_0^i) = \frac{p_0^*(\bar{x}_0^i)}{q_0(\bar{x}_0^i)}\right\}_{i=1}^n$,
normalize the weights to obtain $\{\tilde{w}(\bar{x}_0^i)\}_{i=1}^n$,
constitute $\hat{\pi}_0^n(dx_0) = \sum_{i=1}^n \tilde{w}(\bar{x}_0^i) \delta_{\bar{x}_0^i}(dx_0)$,
resample $\{x_0^i\}_{i=1}^n$ using $\mathcal{M}(n, \tilde{w}(\bar{x}_0^1), \dots, \tilde{w}(\bar{x}_0^n))$ to obtain $\pi_0^n(dx_0)$.

- **2. importance sampling**

$t = t + 1$,
sample $\{\bar{x}_t^i \sim q_{t-1}(dx_t^i|x_{0:t-1}^i)\}_{i=1}^n$,
compose $\{\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i)\}_{i=1}^n$,
for $i = 1:n$ update weights according to (2.15) to obtain $w(\bar{x}_{0:t}^i)$, i.e.,

$$w(\bar{x}_{0:t}^i) = w(x_{0:t-1}^i) \cdot \frac{p_t^*(\bar{x}_{0:t}^i)}{p_{t-1}^*(x_{0:t-1}^i)q_{t-1}(\bar{x}_t^i|x_{0:t-1}^i)},$$

normalize $w(\bar{x}_{0:t}^i)$ to obtain $\tilde{w}(\bar{x}_{0:t}^i)$,
constitute empirical measure $\hat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(\bar{x}_{0:t}^i) \delta_{\bar{x}_{0:t}^i}(dx_{0:t})$.

- **3. resampling**

resample $\{x_{0:t}^i\}_{i=1}^n$ using $\mathcal{M}(n, \tilde{w}(\bar{x}_{0:t}^1), \dots, \tilde{w}(\bar{x}_{0:t}^n))$ to obtain $\pi_t^n(dx_{0:t})$.

- **4.** if $t = T$ end, else go to step 2.

Algorithm 2.1: SMC algorithm.

The sequential character of the algorithm and the incorporation of the resampling step causes the work of the algorithm to be seen as an alternating

generation of three empirical measures denoted by $\bar{\pi}_t^n$, $\hat{\pi}_t^n$ and π_t^n . Let us inspect the presented pseudocode to see this fact.

We start at the beginning of the importance sampling step. At time t we have at our disposal uniformly weighted samples $x_{0:t-1}^i \in \mathbb{R}^{dt}$ from the resampling step performed at time $t-1$. Generating samples $\bar{x}_t^i \in \mathbb{R}^d$ from $q_{t-1}(dx_t|x_{0:t-1})$, where $q_{t-1}(dx_t|x_{0:t-1})$ denotes the conditional measure corresponding to density $q_{t-1}(x_t|x_{0:t-1})$, we constitute samples from $\mathbb{R}^{d(t+1)}$ space by composition of $x_{0:t-1}^i$ and \bar{x}_t^i . The composition gives samples $\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i)$. Remark that it is **not true** that $\bar{x}_{0:t}^i = (\bar{x}_{0:t-1}^i, \bar{x}_t^i)$ because $\bar{x}_{0:t-1}^i$ are resampled into $x_{0:t-1}^i$ and those resampled samples - not $\bar{x}_{0:t-1}^i$ - enter the importance sampling formula at time t .

The composed samples $\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i)$ when considered in the uniformly weighted arrangement, carry the first empirical measure

$$\bar{\pi}_t^n(dx_{0:t}) = \frac{1}{n} \sum_{i=1}^n \delta_{\bar{x}_{0:t}^i}(dx_{0:t}). \quad (2.21)$$

We call $\bar{\pi}_t^n$ as the *empirical prediction measure*, which will be explained in Section 2.6.5.

Reweighting the prediction measure using formula (2.15) we obtain the *empirical update measure* $\hat{\pi}_t^n$. That is, using the normalized weights $\tilde{w}(\bar{x}_{0:t}^i)$ which are obtained on the basis of updating presented in (2.15), we specify

$$\hat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(\bar{x}_{0:t}^i) \delta_{\bar{x}_{0:t}^i}(dx_{0:t}). \quad (2.22)$$

Finally, the resampling step transforms $\hat{\pi}_t^n$ into the uniformly weighted *empirical resampled measure* π_t^n which is carried by resampled samples $x_{0:t}^i \in \mathbb{R}^{d(t+1)}$. Clearly, this measure then enters the importance sampling step at the next iteration of the SMC algorithm. Schematically, the work of the algorithm can be expressed by the following stream:

$$\pi_0 \rightarrow \bar{\pi}_1^n \rightarrow \hat{\pi}_1^n \rightarrow \pi_1^n \rightarrow \dots \rightarrow \bar{\pi}_t^n \rightarrow \hat{\pi}_t^n \rightarrow \pi_t^n \rightarrow \dots$$

We comment on the schema again in Section 2.6.5.

The other important observation is that resampling has a simplification impact on the weight update step in formula (2.15) and consequently on the normalization of updated weights. In fact, when entering the formula at time t , $w(x_{0:t-1}^i) = 1/n$ due to resampling. So the updated weights $w(\bar{x}_{0:t}^i)$ read as $1/n$

times the update factor. The factor $1/n$ is erased by normalization, therefore normalized weights $\tilde{w}(\bar{x}_{0:t}^i)$ are determined only by normalization of update factors occurring in (2.15). To retain the clarity of the explanation the discussed simplification is not incorporated directly in the algorithm above. But we use it in the modification of the algorithm for particle filters.

This finalizes introduction of basics of SMC methods. In the next section we discuss their application in particle filters where more structural assumptions are set on the probability measures which are to be approximated.

2.6 Particle filters

Particle filters represent the most prominent area of application of SMC methods. The task of filtering is to compute conditional distributions of primarily unobservable variables on the basis of variables which are observable. Unobservable variables are traditionally called *signal*, observable ones *output* and the conditional distributions of interest which evolves over time *filtering distributions*.

Certain structural relations on the probability model driving behavior of both groups of variables are assumed. The signal is treated as a generally inhomogeneous Markov chain. Based on the actual value of the signal, the observation is determined according to a known formula typically including a *noise term*. The task is then as follows - on the basis what we have seen until now to present our best possible estimate on the actual value of the signal either for description or prediction purposes. In what follows we state the model of particle filters mathematically. Our presentation is mainly based on [Doucet et al. 2001; Crisan and Doucet 2002] and [Fristedt et al. 2007].

2.6.1 Signal process

Let (Ω, \mathcal{F}, P) be a probabilistic space. Let $\{X_t, t = 0, 1, \dots\}$ be an \mathbb{R}^{n_x} valued Markov chain defined on this space, $n_x \in \mathbb{N}$. Let \mathcal{F}_s be the natural filtration of the chain, i.e., $\mathcal{F}_s = \sigma(X_k, 0 \leq k \leq s)$. The Markov property writes as

$$P(X_t \in A | \mathcal{F}_s) = P(X_t \in A | X_s), \quad P\text{-a.s.}$$

for all $t \in \mathbb{N}_0$, $A \in \mathcal{B}(\mathbb{R}^{n_x})$ and $s \leq t$. Hence our best estimate on future given the whole past can be done only by the actual information.

The complete description of the probabilistic behavior of the chain in the form of finite-dimensional joint distributions is determined by an initial distribution of X_0 and by the set of transition kernels $\{K_t(A, x), t \in \mathbb{N}_0\}$. We denote the initial distribution of X_0 by π_0 , i.e., $X_0 \sim \pi_0$.

It is well known from the theory of Markov chains that a transition kernel $K_t(A, x)$ represents for any $t \in \mathbb{N}_0$ a function $K_t(A, x) : \mathcal{B}(\mathbb{R}^{n_x}) \times \mathbb{R}^{n_x} \rightarrow [0, 1]$ such that

- $K_t(\cdot, x)$ is a probability measure on $\mathcal{B}(\mathbb{R}^{n_x})$ for each $x \in \mathbb{R}^{n_x}$.
- $K_t(A, \cdot)$ is a Borel function for each $A \in \mathcal{B}(\mathbb{R}^{n_x})$.

It means that by fixing the second variable we get a measure representing the probability of shifting to the set of states A under the condition of being in state x at time t . Hence the measure $K_t(\cdot, x)$ represents a conditional distribution. In the following text we denote this distribution either as $K_t(dy|x)$ or more specifically as $K_t(dx_{t+1}|x_t)$.

In the context of our application, we will work with kernels which possess the Feller property, that is, for each f continuous and bounded on \mathbb{R}^{n_x} , i.e., $f \in C_b(\mathbb{R}^{n_x})$, the formula

$$g(x) = \int f(y) K_t(dy|x) \tag{2.23}$$

yields $g : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ also in $C_b(\mathbb{R}^{n_x})$ for each $t \in \mathbb{N}_0$.

Fixing a set in the specification of a transition kernel we get an integrable function. Its integral with respect to the distribution of X_t gives the marginal of being in set A at time $t + 1$; or, if the indicator of $X_t \in B$ is taken into the account the integral gives the joint probability $P(X_{t+1} \in A, X_t \in B)$ as follows from the definition of a conditional distribution.

2.6.2 Observation process

The observation process $\{Y_t, t = 1, 2, \dots\}$ is considered as n_y -dimensional stochastic process, i.e., $Y_t \in \mathbb{R}^{n_y}$, $n_y \in \mathbb{N}$. Observations are determined on the basis of actual values of the signal by using h_t modification functions. They are further considered to be corrupted by a noise. Formally the observation process is specified as

$$Y_t = h_t(X_t) + V_t, t \in \mathbb{N} \tag{2.24}$$

where $h_t : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_y}$ are continuous Borel functions and V_t n_y -dimensional noise terms, independent of X_t , with the probability laws absolutely continuous with respect to the Lebesgue measure on \mathbb{R}^{n_y} . For fixed t , the density of V_t is denoted by g_t and it is assumed to be bounded and continuous. The most prominent choice for V_t is when they are i.i.d. multivariate normal.

Y_t values are considered mutually and conditionally independent of other variables given the state X_t . This means that for all $t \in \mathbb{N}$ and $B \in \mathcal{B}(\mathbb{R}^{n_y})$,

$$P(Y_t \in B_t | X_{0:t}, Y_{1:t-1}) = P(Y_t \in B_t | X_t). \quad (2.25)$$

Further, for $n_y = 1$ we have $P(Y_t \leq b | X_t = x_t) = P(V_t \leq b - h(x_t))$, which implies

$$P(Y_t \leq b | X_t = x_t) = \int_{-\infty}^{b-h(x_t)} g_t(u) du = \int_{-\infty}^b g_t(y_t - h(x_t)) dy_t.$$

The similar formula can be derived in the multivariate case of $n_y > 1$. Denoting

$$g_t(y_t | x_t) = g_t(y_t - h_t(x_t)) \quad (2.26)$$

we get $g_t(y_t | x_t)$ as the conditional density of Y_t with respect to $X_t = x_t$.

2.6.3 Filtering distribution

Now we discuss the computation of the filtering distributions. The probabilistic structure and evolution of a particle filter is schematically presented in the following diagram:

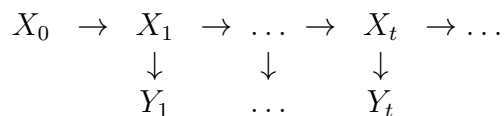


Figure 2.1: Evolution of a particle filter.

In the most general setting, filtering means the computation of conditional distribution of the vector of states $X_{0:t} = (X_0, \dots, X_t)$ given an observable history Y_1, \dots, Y_t . That is, we would like to compute for any $A \in \mathcal{B}(\mathbb{R}^{n_x(t+1)})$,

$$P(X_{0:t} \in A | \sigma(Y_1, \dots, Y_t)) = P(X_{0:t} \in A | Y_{1:t}).$$

Conditional distributions are typically computed from corresponding joint distributions by suitable integrations. We follow this standard way, so we start with the joint distribution of $X_{0:t}$ and $Y_{1:t}$.

The joint distribution of the signal is determined by the initial distribution of $X_0 \sim \pi_0$ and by transition kernels K_t , $t \in \mathbb{N}_0$. We further adopt $P(dx_{0:t})$ to denote the joint distribution of X_0, \dots, X_t . Similarly $P(dx_{0:t}y_{1:t})$ will denote the joint distribution of $X_0, \dots, X_t, Y_1, \dots, Y_t$. Let us compute the probability of $X_{0:t}$ being in $A_{0:t}$ for $A_0, \dots, A_t \in \mathcal{B}(\mathbb{R}^{n_x})$, i.e.,

$$P(X_{0:t} \in A_{0:t}) = P(X_0 \in A_0, X_1 \in A_1, \dots, X_t \in A_t) = \int_{A_{0:t}=A_0 \times \dots \times A_t} P(dx_{0:t}).$$

We have the following well known iterative expansion

$$\begin{aligned} P(X_0 \in A_0) &= \int_{A_0} \pi_0(dx_0), \\ P(X_{0:1} \in A_{0:1}) &= \int_{A_0} K_0(A_1, x_0) \pi_0(dx_0), \\ P(X_{0:2} \in A_{0:2}) &= \int_{A_0} \left(\int_{A_1} K_1(A_2, x_1) K_0(dx_1|x_0) \right) \pi_0(dx_0), \\ &\vdots \\ P(X_{0:t} \in A_{0:t}) &= \int_{A_0} \left(\int_{A_1} \dots \left(\int_{A_{t-1}} K_{t-1}(A_t, x_{t-1}) K_{t-2}(dx_{t-1}|x_{t-2}) \right) \dots \right. \\ &\quad \left. \dots K_0(dx_1|x_0) \right) \pi_0(dx_0). \end{aligned}$$

The above expansion presents the evolution of the joint distribution of a Markov chain over time.

Now we include also the observation process. Because of the independence (2.25) of Y_t we have

$$\begin{aligned} P(X_{0:t} \in A_{0:t}, Y_t \in B_t) &= \int_{A_{0:t}} P(Y_t \in B_t | X_{0:t} = x_{0:t}) P(dx_{0:t}) \\ &= \int_{A_{0:t}} P(Y_t \in B_t | X_t = x_t) P(dx_{0:t}), \\ P(X_{0:t} \in A_{0:t}, Y_t \in B_t) &= \int_{A_{0:t}} \int_{B_t} g_t(dy_t|x_t) P(dx_{0:t}) \end{aligned}$$

where $g_t(dy_t|x_t) = P(dy_t|x_t) = P(Y_t \in dy_t | X_t = x_t)$ is the conditional distribution of Y_t on the condition of $X_t = x_t$.

The combination with the expansion of $P(dx_{0:t})$ gives the following evolution of the joint distribution over time. First of all,

$$P(X_{0:1} \in A_{0:1}) = \int_{A_0} K_0(A_1, x_0) \pi(dx_0) = \int_{A_0} \int_{A_1} K_0(dx_1|x_0) \pi(dx_0).$$

From this we have

$$P(X_{0:1} \in A_{0:1}, Y_1 \in B_1) = \int_{A_0} \int_{A_1} \int_{B_1} g_1(dy_1|x_1) K_0(dx_1|x_0) \pi(dx_0)$$

and by a chain of iterations we finally get

$$P(X_{0:t} \in A_{0:t}, Y_{1:t} \in B_{1:t}) = \int_{A_0} \left(\int_{A_1} \int_{B_1} \dots \int_{A_t} \int_{B_t} g_t(dy_t|x_t) K_{t-1}(dx_t|x_{t-1}) \dots g_1(dy_1|x_1) K_0(dx_1|x_0) \right) \pi_0(dx_0).$$

Let us assume that the initial distribution π_0 , $K_{t-1}(dx_t|x_{t-1})$ and $g_t(dy_t|x_t)$ admit for $t \geq 1$ densities with respect to the respective Lebesgue measures, i.e.,

- $\pi_0(dx_0) = p_0(x_0) dx_0$,
- $K_{t-1}(dx_t|x_{t-1}) = K_{t-1}(x_t|x_{t-1}) dx_t$,
- $g_t(dy_t|x_t) = g_t(y_t|x_t) dy_t$.

This assumption leads to the joint density $p(x_{0:t}, y_{1:t})$ of $P(dx_{0:t}y_{1:t})$,

$$p(x_{0:t}, y_{1:t}) = p_0(x_0) \prod_{k=1}^t g_k(y_k|x_k) K_{k-1}(x_k|x_{k-1}). \quad (2.27)$$

Having the joint distribution at our disposal we would like to compute the conditional densities $p(x_{0:t}|y_{1:t})$ of the filtering distributions for $t \geq 1$. The standard formula based on the ratio of the joint and the respective marginal distribution is directly inapplicable as it requires the computation of a complicated normalization integral. That is,

$$p(x_{0:t}|y_{1:t}) = \frac{p(x_{0:t}, y_{1:t})}{\int p(x_{0:t}, y_{1:t}) dx_{0:t}} = \frac{p(x_{0:t}, y_{1:t})}{Z_t(y_{1:t})} \quad (2.28)$$

where $Z_t(y_{1:t}) = \int p(x_{0:t}, y_{1:t}) dx_{0:t}$ is generally analytically intractable. However, note that $Z_t = Z_t(y_1, \dots, y_t)$ is constant for given y_1, \dots, y_t . So the filtering density $p(x_{0:t}|y_{1:t})$ falls into the framework of importance sampling methods because we have its specification up to the normalizing constant $Z_t(y_{1:t})$ for

a given observation history. Therefore, this filtering problem seems to be a good candidate for an application of SMC methods.

In what follows, we denote the conditional distributions of interest, i.e., the filtering distributions, by the notation taken from [Crisan and Doucet 2002],

$$\pi_{k:l|m}(dx_{k:l}) = P(X_{k:l} \in dx_{k:l} | Y_{1:m} = y_{1:m}). \quad (2.29)$$

2.6.4 SMC algorithm for particle filters

In order to incorporate the particle filtering problem into the SMC methods' framework we have to recognize the sequence of increasing probability spaces, corresponding densities and a suitable proposal density as it was explained in Section 2.3.

Concerning the probabilistic spaces, we have the standard setup of $E_0 = \mathbb{R}^{n_x}$, $\mathcal{E}_0 = \mathcal{B}(\mathbb{R}^{n_x})$ and $E_t = (\mathbb{R}^{n_x})^{t+1}$, $\mathcal{E}_t = \otimes_{k=0}^t \mathcal{B}(\mathbb{R}^{n_x})$ for $t \geq 1$. We set π_0 as the initial distribution of the state process, i.e., $X_0 \sim \pi_0$, with density $p_0(x_0)$. For $t \geq 1$, π_t are the filtering distributions, i.e., the conditional distributions of state vector $X_{0:t}$ with respect to vector of observations $Y_{1:t}$. Using notation (2.29) we have $\pi_t = \pi_{0:t|t}$ with the corresponding densities $p(x_{0:t}|y_{1:t})$.

The densities of interest $p(x_{0:t}|y_{1:t})$ are known only up to the normalizing constants Z_t . According to (2.28) we have $p(x_{0:t}|y_{1:t}) \propto p(x_{0:t}, y_{1:t})$, i.e., $p^*(x_{0:t}|y_{1:t}) = p(x_{0:t}, y_{1:t})$, which writes in more details as

$$p_t^*(x_{0:t}|y_{1:t}) = p_0(x_0) \prod_{k=1}^t g_k(y_k|x_k) K_{k-1}(x_k|x_{k-1}).$$

Concerning a proposal density, we assume that we are able to sample from conditional measures $K_{t-1}(dx_t|x_{t-1})$. The proposal density at time t has then form

$$q_t(x_{0:t}) = p(x_0) \prod_{k=1}^t K_{k-1}(x_k|x_{k-1}).$$

and therefore $q_{k-1}(x_k|x_{0:k-1}) = K_{k-1}(x_k|x_{k-1})$ to match formula (2.14).

By the specification of the structure of proposal density we can determine the counterpart of the weight update formula (2.15) for particle filters. We have

for $t \geq 1$, ($w(x_0^i) = 1$ as $q_0(x_0) = p_0(x_0)$),

$$\begin{aligned}
w(x_{0:t}^i) &= \frac{p_t^*(x_{0:t}^i)}{q_t(x_{0:t}^i)} = w(x_{0:t-1}^i) \cdot \frac{p_t^*(x_{0:t}^i)}{p_{t-1}^*(x_{0:t-1}^i) q_{t-1}(x_t^i | x_{0:t-1}^i)} \\
&= w(x_{0:t-1}^i) \cdot \frac{p_0(x_0^i) \prod_{k=1}^t g_k(y_k | x_k^i) K_{k-1}(x_k^i | x_{k-1}^i)}{[p_0(x_0^i) \prod_{k=1}^{t-1} g_k(y_k | x_k^i) K_{k-1}(x_k^i | x_{k-1}^i)] K_{t-1}(x_t^i | x_{t-1}^i)} \\
&= w(x_{0:t-1}^i) \cdot g_t(y_t | x_t^i).
\end{aligned}$$

Thus weights are updated sequentially by factor $g_t(y_t | x_t^i)$. Remark that all elements of $g_t(y_t | x_t^i) = g_t(y_t - h_t(x_t^i))$ are known from the preceding step. Let us state explicitly the version of SMC algorithm for particle filters.

- **0. declarations**

n - number of samples,

T - computational horizon,

$p_0(x_0)$ - the initial density of $X_0 \sim \pi_0$,

$K_{t-1}(x_t | x_{0:t-1})$, $t = 1, \dots, T$, - conditional transition densities.

- **1. initialization**

$t = 0$,

sample $\{\bar{x}_0^i \sim p_0(dx_0)\}_{i=1}^n$,

constitute $\hat{\pi}_0^n(dx_0) = \frac{1}{n} \sum_{i=1}^n \delta_{\bar{x}_0^i}(dx_0)$, (*no reweighting is needed*)

set $\pi_0^n(dx_0) = \hat{\pi}_0^n(dx_0)$, i.e., $\{x_0^i = \bar{x}_0^i\}_{i=1}^n$. (*no resampling is needed*)

- **2. importance sampling**

$t = t + 1$,

sample $\{\bar{x}_t^i \sim K_{t-1}(dx_t | x_{0:t-1}^i)\}_{i=1}^n$,

compose $\{\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i)\}_{i=1}^n$,

for $i = 1:n$ compute (*simplification due to resampling*)

$$\tilde{w}(\bar{x}_{0:t}^i) = \frac{g_t(y_t - h_t(\bar{x}_t^i))}{\sum_i g_t(y_t - h_t(\bar{x}_t^i))},$$

constitute $\hat{\pi}_t^n(dx_{0:t}) = \sum_{i=1}^n \tilde{w}(\bar{x}_{0:t}^i) \delta_{\bar{x}_{0:t}^i}(dx_{0:t})$.

- **3. resampling**

resample $\{x_{0:t}^i\}_{i=1}^n$ using $\mathcal{M}(n, \tilde{w}(\bar{x}_{0:t}^1), \dots, \tilde{w}(\bar{x}_{0:t}^n))$ to obtain $\pi_t^n(dx_{0:t})$,

- **4.** if $t = T$ end, else go to step 2.

Algorithm 2.2: SMC algorithm in the particle filter design.

The presented pseudocode of the SMC algorithm for filtering problem can be straightforwardly implemented on a suitable software platform. This is done in Chapter 5 which is devoted to computer simulations.

2.6.5 Recursive evolution of filtering distributions

The application of SMC methods enables to overcome the problem with the analytic expression of conditional densities of interest in the filtering problem. More specifically, by applying the SMC algorithm we are able to compute integral characteristics of filtering distributions on the basis of computer simulations.

Inspecting the work of the SMC algorithm, we can identify again three sub-steps in a single cycle of operation. The first sub-step is generation of new data samples employing the transition distribution $K_{t-1}(dx_t|x_{t-1})$ yielding (after composition) data samples $\{\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i)\}_{i=1}^n$. The related empirical measure is the prediction measure denoted $\bar{\pi}_t^n$. The second sub-step is the generation of empirical update measure $\hat{\pi}_t^n$ from $\bar{\pi}_t^n$ through updating factors $g_t(y_t|\bar{x}_t^i)$. Finally, the third sub-step is the resampling step yielding uniformly weighted empirical measure π_t^n , which is the empirical counterpart of the filtering distribution $\pi_{0:t|t}$. But, only the first two sub-steps generates empirical measures related directly to the filtering distribution. Resampling is only a rearrangement of $\hat{\pi}_t^n$ measure to the uniformly weighted design.

Interestingly enough, the time evolution of filtering distributions $\pi_{0:t|t}$ and related densities $p(x_{0:t}|y_{1:t})$ can be also described in the analytic form by recursive alternations of true counterparts of $\bar{\pi}_t^n$ and $\hat{\pi}_t^n$ empirical distributions. Using the notation of (2.29), these counterparts read as $\pi_{0:t|t-1}$ for $\bar{\pi}_t^n$ and $\pi_{0:t|t}$ for $\hat{\pi}_t^n$. More specifically, the densities of these counterparts $p(x_{0:t}|y_{1:t-1})$ and $p(x_{0:t}|y_{1:t})$ can be expressed analytically, but only in a recursive design. That is, only one step ahead from a given time instant. (If it would be possible to express the densities non-recursively, then we would not need SMC methods).

The respective formulas are known in the literature as *the prediction* and *the update* formula [Doucet et al. 2001; Crisan and Doucet 2002].

Lemma 2.1. *Let the density of the joint distribution of a particle filter is specified according to formula (2.27), then*

$$p(x_{0:t}|y_{1:t-1}) = \int K_{t-1}(x_t|x_{t-1})p(x_{0:t-1}|y_{1:t-1}) dx_{0:t-1}. \quad (2.30)$$

Proof. The formula is actually an instance of the Chapman-Kolmogorov equation for Markov chains. Nevertheless, let us present the full proof here. To start we rewrite formula (2.30) in an equivalent form using related measures.

$$\begin{aligned} p(x_{0:t}|y_{1:t-1}) &= \int K_{t-1}(x_t|x_{t-1})p(x_{0:t-1}|y_{1:t-1}) dx_{0:t-1}, \\ p(x_{0:t}|y_{1:t-1}) &= \pi_{0:t-1|t-1}K_{t-1}(x_t|x_{t-1}), \\ P(dx_{0:t}|y_{1:t-1}) &= \pi_{0:t-1|t-1}K_{t-1}(dx_t|x_{t-1}), \\ \pi_{0:t|t-1} &= \pi_{0:t-1|t-1}K_{t-1}. \end{aligned}$$

Remark that $K_{t-1}(x_t|x_{t-1})$ denotes the conditional density of the conditional distribution $K_{t-1}(dx_t|x_{t-1})$ carried by the transition kernel K_{t-1} . Term $\pi_{0:t-1|t-1}K_{t-1}$ represents the measure given by the composition of the unconditional distribution $\pi_{0:t-1|t-1}$ with the conditional distribution $K_{t-1}(dx_t|x_{t-1})$. For any set $A \in \mathcal{B}(\mathbb{R}^{n_x(t+1)})$ the composition reads as $\pi_{0:t-1|t-1}K_{t-1}(A) = \int K_{t-1}(A, x) d\pi_{0:t-1|t-1}$.

For any bounded $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ the Markov property of $\{X_t, t \geq 0\}$ chain gives

$$\mathbb{E}[f(X_{0:t})|X_{0:t-1}](\omega) = \mathbb{E}[f(X_{0:t-1}(\omega), X_t)|X_{t-1}](\omega) = (K_{t-1}f)(X_{0:t-1}(\omega))$$

where $(K_{t-1}f)(x_{0:t-1}) = \int f(x_{0:t})K_{t-1}(dx_t|x_{0:t-1}) dx_t$.

Due to the independence of $V_{1:t}$ on $X_{0:t}$ we have

$$\mathbb{E}[f(X_{0:t})|X_{0:t-1}, V_{1:t-1}](\omega) = \mathbb{E}[f(X_{0:t})|X_{0:t-1}](\omega) = (K_{t-1}f)(X_{0:t-1}(\omega))$$

and hence the following chain of equalities holds

$$\begin{aligned} \pi_{0:t|t-1}f &= \mathbb{E}[f(X_{0:t})|Y_{1:t-1}] \\ &= \mathbb{E}[\mathbb{E}[f(X_{0:t})|X_{0:t-1}, V_{1:t-1}]|Y_{1:t-1}] \\ &= \mathbb{E}[(K_{t-1}f)(X_{0:t-1})|Y_{1:t-1}] \\ &= \pi_{0:t-1|t-1}K_{t-1}f. \end{aligned}$$

As this holds for any bounded f , then the measures $\pi_{0:t|t-1}$ and $\pi_{0:t-1|t-1}K_{t-1}$ coincide [Billingsley 1995], which proves (2.30). \square

Now we are going to prove what is known as the update formula.

Lemma 2.2. *Let the density of the joint distribution of a particle filter be specified according to formula (2.27), then*

$$p(x_{0:t}|y_{1:t}) = \frac{g_t(y_t|x_t)p(x_{0:t}|y_{1:t-1})}{\int g_t(y_t|x_t)p(x_{0:t}|y_{1:t-1}) dx_{0:t}}. \quad (2.31)$$

Proof. We start with the Bayes' rule and rearrange

$$\begin{aligned} p(x_{0:t}|y_{1:t}) &= \frac{p(y_{1:t}|x_{0:t})p(x_{0:t})}{p(y_{1:t})}, \\ p(x_{0:t}|y_{1:t}) &= \frac{p(y_t, y_{1:t-1}|x_{0:t})p(x_{0:t})}{p(y_t, y_{1:t-1})}, \\ p(x_{0:t}|y_{1:t}) &= \frac{p(y_t|y_{1:t-1}, x_{0:t})p(y_{1:t-1}|x_{0:t})p(x_{0:t})}{p(y_t|y_{1:t-1})p(y_{1:t-1})}. \end{aligned}$$

We again use Bayes' rule on $p(y_{1:t-1}|x_{0:t})$

$$p(x_{0:t}|y_{1:t}) = \frac{p(y_t|y_{1:t-1}, x_{0:t})p(x_{0:t}|y_{1:t-1})p(y_{1:t-1})p(x_{0:t})}{p(y_t|y_{1:t-1})p(y_{1:t-1})p(x_{0:t})}.$$

Taking into account conditional independence $p(y_t|y_{1:t-1}, x_{0:t}) = p(y_t|x_{0:t})$ and cancelling $p(y_{1:t-1})p(x_{0:t})$ terms gives the final formula

$$p(x_{0:t}|y_{1:t}) = \frac{p(y_t|x_{0:t})p(x_{0:t}|y_{1:t-1})}{p(y_t|y_{1:t-1})}.$$

The normalizing constant in denominator is computed by integration

$$p(y_t|y_{1:t-1}) = \int p(y_t|x_{0:t})p(x_{0:t}|y_{1:t-1}) dx_{0:t}.$$

As we have $p(y_t|x_{0:t}) = g_t(y_t|x_t)$ this finishes the proof. \square

The obtained formulas can be also rewritten in terms of corresponding measures. We have

$$\begin{aligned} \text{prediction :} \quad \pi_{0:t|t-1} &= \pi_{0:t-1|t-1}K_{t-1}, \\ \text{update :} \quad \frac{d\pi_{0:t|t}}{d\pi_{0:t|t-1}} &= \frac{g_t(y_t|x_t)}{\int g_t(y_t|x_t) d\pi_{0:t|t-1}}. \end{aligned}$$

The measure $\pi_{0:t|t-1}(dx_{0:t})$ is obtained by composition of $\pi_{0:t-1|t-1}(dx_{0:t-1})$ with the kernel's conditional measure $K_{t-1}(dx_t|x_{t-1})$. As the composed measure is

based on the preceding measure $\pi_{0:t-1|t-1}$ and not on the current observation y_t , it is called the *prediction measure*. The measure $\pi_{0:t|t}(dx_{0:t})$ has the Radon-Nikodým derivative with respect to $\pi_{0:t|t-1}(dx_{0:t})$ given by the presented fraction. The fraction is based on the current observation y_t and the measure is updated version of the prediction measure. That is why, it is called the *update measure*.

In the view of the presented formulas we can schematically express the workflow of related empirical and theoretical distributions in a SMC particle filter as follows:

$$\begin{array}{ccccccccccc} \pi_0 & \rightarrow & \bar{\pi}_1^n & \rightarrow & \hat{\pi}_1^n & \rightarrow & \pi_1^n & \rightarrow & \dots & \rightarrow & \bar{\pi}_t^n & \rightarrow & \hat{\pi}_t^n & \rightarrow & \pi_t^n & \rightarrow & \dots \\ \pi_0 & \rightarrow & \pi_{0:1|0} & \rightarrow & \pi_{0:1|1} & \rightarrow & \dots & \rightarrow & \pi_{0:t|t-1} & \rightarrow & \pi_{0:t|t} & \rightarrow & \dots \end{array}$$

Figure 2.2: The evolution of empirical and theoretical distributions.

Corresponding pairs of empirical and theoretical measures are $(\bar{\pi}_t^n, \pi_{0:t|t-1})$ and $(\pi_t^n, \pi_{0:t|t})$. In fact also $(\hat{\pi}_t^n, \pi_{0:t|t})$ correspond, but we are interested in the resampled version of $\hat{\pi}_t^n$, i.e., in π_t^n , therefore we primarily consider as matching pair $(\pi_t^n, \pi_{0:t|t})$.

To conclude let us note that the prediction and update formulas are valid also if only marginal conditional distributions are considered [Crisan and Doucet 2002]. That is, when we are interested only in time evolution of densities $p(x_t|y_{1:t})$ of conditional marginal distributions $\pi_{t|t}$. The corresponding formulas write as

$$\text{marg. prediction : } p(x_t|y_{1:t-1}) = \int K_{t-1}(x_t|x_{t-1})p(x_{t-1}|y_{1:t-1}) dx_{t-1}, \quad (2.32)$$

$$\text{marg. update : } p(x_t|y_{1:t}) = \frac{g(y_t|x_t)p(x_t|y_{1:t-1})}{\int g(y_t|x_t)p(x_t|y_{1:t-1}) dx_t}. \quad (2.33)$$

In terms of the corresponding measures this writes as

$$\text{marg. prediction : } \pi_{t|t-1} = \pi_{t-1|t-1}K_{t-1}, \quad (2.34)$$

$$\text{marg. update : } \frac{d\pi_{t|t}}{d\pi_{t|t-1}} = \frac{g_t(y_t|x_t)}{\int g_t(y_t|x_t) d\pi_{t|t-1}} = \frac{g_t(y_t|x_t)}{\int g_t(y_t|x_t) d\pi_{0:t|t-1}}. \quad (2.35)$$

We employ these formulas extensively in Chapter 4 of the thesis.

2.7 Convergence results for particle filters

This section covers the main convergence results for particle filters. We review the necessary and sufficient condition ensuring the convergence of empirical random measures, which are generated by the SMC algorithm, to the measures of the interest. The results are taken from [Doucet et al. 2001] and [Crisan and Doucet 2002].

Measures. To start, we stress that we will work only with measures from the class $\mathcal{P}(\mathbb{R}^d)$ of probabilistic measures specified on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, $d \in \mathbb{N}$ measurable space. For an integrable function f and a probabilistic measure μ we denote the integral (image) of the function f over μ as μf , i.e., $\mu f = \int f d\mu$.

One of the equivalent definitions of the notion of *weak convergence of measures* reads as follows [Billingsley 1999]. Let $\{\mu^n\}_{n=1}^\infty$, $\mu^n \in \mathcal{P}(\mathbb{R}^d)$ be a sequence of probabilistic measures and $\mu \in \mathcal{P}(\mathbb{R}^d)$. We say that the sequence $\{\mu^n\}_{n=1}^\infty$ converges weakly to the measure μ if for all $f \in \mathcal{C}^b(\mathbb{R}^d)$, $\lim_{n \rightarrow \infty} \mu^n f = \mu f$. Remark that from reasons which will be apparent in Chapter 4, we will consider functions from the more general space $\mathcal{C}^b(\mathbb{C}^d)$ in the theorems below.

To switch to the random case, consider a probabilistic space (Ω, \mathcal{F}, P) and a sequence of random measures $\{\mu^n\}_{n=1}^\infty$, $\mu^n : \Omega \rightarrow \mathcal{P}(\mathbb{R}^d)$ and let $\mu \in \mathcal{P}(\mathbb{R}^d)$ be a fixed probabilistic measure, then the following two types of convergence can be recognized:

1. - *in expectation*: $\lim_{n \rightarrow \infty} \mathbb{E}[|\mu^n f - \mu f|] = 0$, for all $f \in \mathcal{C}^b(\mathbb{C}^d)$
2. - *almost surely*: $\lim_{n \rightarrow \infty} \mu^n = \mu$, $P - a.s.$

We will study only the convergence in expectation because this is the most relevant to the topic of the thesis. The results concerning the other type of convergence are presented in [Doucet et al. 2001]. Let us only remark, that if $\lim_{n \rightarrow \infty} \mathbb{E}[|\mu^n f - \mu f|] = 0$, then there exists a subsequence $n(m)$ such that $\lim_{m \rightarrow \infty} \mu^{n(m)} = \mu$, $P - a.s.$

Evolution of a SMC filter. The time evolution of empirical and filtering distributions in a SMC particle filter was presented in Fig. 2.2. The natural question is if the empirical prediction and resampled update measures $\bar{\pi}_t^n$ and π_t^n converge in expectation to their theoretical counterparts $\bar{\pi}_t = \pi_{0:t|t-1}$ and $\pi_t = \pi_{0:t|t}$ as the number of samples (particles) goes to infinity.

Note explicitly that we use further only the shortcuts $\bar{\pi}_t$ and π_t for the prediction and filtering distributions, respectively.

In the theorem below, we prove that the convergence of the whole algorithm is assured, if there is assured convergence in each particular step. To present the theorem, let us remind from Section 2.6.3 that the prediction distribution $\bar{\pi}_t$ is determined by the composition of π_{t-1} distribution with the transition kernel K_{t-1} ; and the update distribution has the Radon-Nikodým derivative with respect to the prediction distribution. The same relations hold for empirical counterparts $\bar{\pi}_t^n$ and $\hat{\pi}_t^n$ as can be verified from the work of SMC algorithm. Formally this writes as

$$\bar{\pi}_t = \pi_{t-1}K_{t-1}, \quad \frac{d\pi_t}{d\bar{\pi}_t} = \frac{g_t}{\bar{\pi}_t g_t} \quad \text{and} \quad \bar{\pi}_t^n = \pi_{t-1}^n K_{t-1}, \quad \frac{d\hat{\pi}_t^n}{d\bar{\pi}_t^n} = \frac{g_t}{\bar{\pi}_t^n g_t}. \quad (2.36)$$

Having specified the relations between individual measures the convergence theorem writes as follows.

Theorem 2.2. *For $t = 1, 2, \dots$, the sequences $\bar{\pi}_t^n, \pi_t^n$ converge to $\bar{\pi}_t$ or π_t , respectively, with convergence in expectation if and only if the following three conditions are satisfied for all $t \geq 1$,*

- (i) for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$, $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_0^n f - \pi_0 f|] = 0$,
- (ii) for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$, $\lim_{n \rightarrow \infty} \mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|] = 0$,
- (iii) for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$, $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_t^n f - \hat{\pi}_t^n f|] = 0$.

Proof. The sufficiency is proved by induction over $t \geq 1$. We are going to prove that if π_{t-1}^n converge to π_{t-1} as $n \rightarrow \infty$ and (i)-(iii) hold, then also $\bar{\pi}_t^n, \pi_t^n$ converge to $\bar{\pi}_t, \pi_t$ for $t > 0$. The statement is true for π_0 due to (i).

We have $\bar{\pi}_t = \pi_{t-1}K_{t-1}$, thus by the triangle inequality for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$,

$$|\bar{\pi}_t^n f - \bar{\pi}_t f| \leq |\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f| + |\pi_{t-1}^n K_{t-1} f - \pi_{t-1} K_{t-1} f|. \quad (2.37)$$

Applying expectation on both sides, we get $\lim_{n \rightarrow \infty} \mathbb{E}[|\bar{\pi}_t^n f - \bar{\pi}_t f|] = 0$ for all $t > 0$, since the expectation of the first term converges to 0 due to (ii) and the expected value of the second term converges to 0 from the induction hypothesis as $K_{t-1} f \in \mathcal{C}_b(\mathbb{C}^{n_x t})$ for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$ due to the Feller property of the kernel.

Now, due to (2.36) we decompose

$$\begin{aligned}
|\widehat{\pi}_t^n f - \pi_t f| &= \left| \frac{\overline{\pi}_t^n f g_t}{\overline{\pi}_t^n g_t} - \frac{\overline{\pi}_t f g_t}{\overline{\pi}_t g_t} \right| \\
&\leq \left| \frac{\overline{\pi}_t^n f g_t}{\overline{\pi}_t^n g_t} - \frac{\overline{\pi}_t^n f g_t}{\overline{\pi}_t g_t} \right| + \left| \frac{\overline{\pi}_t^n f g_t}{\overline{\pi}_t g_t} - \frac{\overline{\pi}_t f g_t}{\overline{\pi}_t g_t} \right| \\
&\leq \|f\|_\infty \left| 1 - \frac{\overline{\pi}_t^n g_t}{\overline{\pi}_t g_t} \right| + \frac{1}{\overline{\pi}_t g_t} |\overline{\pi}_t^n f g_t - \overline{\pi}_t f g_t| \\
&= \frac{\|f\|_\infty}{\overline{\pi}_t g_t} |\overline{\pi}_t g_t - \overline{\pi}_t^n g_t| + \frac{1}{\overline{\pi}_t g_t} |\overline{\pi}_t^n f g_t - \overline{\pi}_t f g_t|,
\end{aligned}$$

therefore

$$\mathbb{E}[|\widehat{\pi}_t^n f - \pi_t f|] \leq \frac{\|f\|_\infty}{\overline{\pi}_t g_t} \mathbb{E}[|\overline{\pi}_t^n g_t - \overline{\pi}_t g_t|] + \frac{1}{\overline{\pi}_t g_t} \mathbb{E}[|\overline{\pi}_t^n f g_t - \overline{\pi}_t f g_t|] \quad (2.38)$$

and both terms on the right side converge to 0 due to (2.37).

Finally ($\widehat{\pi}_t^n$ - measure before resampling, π_t^n - measure after resampling),

$$|\pi_t^n f - \pi_t f| \leq |\pi_t^n f - \widehat{\pi}_t^n f| + |\widehat{\pi}_t^n f - \pi_t f|.$$

As the expected value of the first term on the right-hand-side converges to zero due to (iii), and the convergence of expectation for the second term is presented in (2.38), the expected value of the left-hand-side converges to 0 as well.

The necessity part. Assume that for $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$, $\lim_{n \rightarrow \infty} \mathbb{E}[|\overline{\pi}_t^n f - \overline{\pi}_t f|] = 0$ for $t > 0$ and $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_t^n f - \pi_t f|] = 0$ for $t \geq 0$. This implies (i) for $t = 0$. From (2.38) we have $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_t f - \widehat{\pi}_t^n f|] = 0$, and because

$$\mathbb{E}[|\pi_t^n f - \widehat{\pi}_t^n f|] \leq \mathbb{E}[|\pi_t^n f - \pi_t f|] + \mathbb{E}[|\pi_t f - \widehat{\pi}_t^n f|]$$

we obtain (iii). To conclude, as $\overline{\pi}_t = \pi_{t-1} K_{t-1}$, we have for all $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$,

$$\mathbb{E}[|\overline{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|] \leq \mathbb{E}[|\overline{\pi}_t^n f - \overline{\pi}_t f|] + \mathbb{E}[|\pi_{t-1} K_{t-1} f - \pi_{t-1}^n K_{t-1} f|] = 0$$

which implies (ii). \square

In the next theorem, we show that SMC algorithm introduced in Section 2.6.4 converges in expectation. It will be performed by validating the conditions (i)-(iii) of Theorem 2.2.

Theorem 2.3. Let $\bar{\pi}_t^n$ and π_t^n , $t \geq 1$ be the measure valued sequences produced by the SMC algorithm for particle filters presented in Section 2.6.4. Then, for all $t \geq 1$, $\lim_{n \rightarrow \infty} \mathbb{E}[|\bar{\pi}_t^n - \bar{\pi}_t|] = 0$ and $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_t^n - \pi_t|] = 0$.

Proof. In order to employ Theorem 2.2 we introduce two σ -algebras \mathcal{F}_t and $\bar{\mathcal{F}}_t$, which reflect the evolution of information based on conditionally sampled and resampled data:

$$\begin{aligned}\mathcal{F}_t &= \sigma(\bar{x}_s^i, x_s^i, s \leq t, i = 1, \dots, n), \\ \bar{\mathcal{F}}_t &= \sigma(\bar{x}_s^i, x_s^i, s < t, \bar{x}_t^i, i = 1, \dots, n).\end{aligned}$$

Clearly, \mathcal{F}_t stores information carried by conditionally sampled and resampled data until time t . In $\bar{\mathcal{F}}_t$ the information from resampled data at time t is excluded.

Since (i) is clearly satisfied we only need to show that (ii) and (iii) are satisfied as well. If $f \in \mathcal{C}_b(\mathbb{C}^{n \times (t+1)})$, then $(\bar{x}_{0:t}^i = (x_{0:t-1}^i, \bar{x}_t^i))$,

$$\mathbb{E}[f(\bar{x}_{0:t}^i) | \mathcal{F}_{t-1}] = K_{t-1} f(x_{0:t-1}^i), \quad \text{for } i = 1, \dots, n. \quad (2.39)$$

Note, that $K_{t-1} f$ is in $\mathcal{C}_b(\mathbb{R}^{n \times t})$ due to the Feller property of the kernel and $K_{t-1} f(x_{0:t-1}^i)$ denotes the value of this function at point $x_{0:t-1}^i$.

Further, $\bar{\pi}_t^n f = \frac{1}{n} \sum_i f(\bar{x}_{0:t}^i)$ so $\mathbb{E}[\bar{\pi}_t^n f | \mathcal{F}_{t-1}] = \frac{1}{n} \sum_i K_{t-1} f(x_{0:t-1}^i)$. We have also $\pi_{t-1}^n K_{t-1} f = \frac{1}{n} \sum_i K_{t-1} f(x_{0:t-1}^i)$ and therefore $\mathbb{E}[\bar{\pi}_t^n f | \mathcal{F}_{t-1}] = \pi_{t-1}^n K_{t-1} f$. Using the independence of conditionally generated samples, we get

$$\mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|^2 | \mathcal{F}_{t-1}] = \mathbb{E}[|\bar{\pi}_t^n f - \mathbb{E}[\bar{\pi}_t^n f | \mathcal{F}_{t-1}]|^2 | \mathcal{F}_{t-1}] = \text{var}[\bar{\pi}_t^n f | \mathcal{F}_{t-1}].$$

$$\begin{aligned}\text{var}[\bar{\pi}_t^n f | \mathcal{F}_{t-1}] &= \text{var} \left[\frac{1}{n} \sum_i f(\bar{x}_{0:t}^i) | \mathcal{F}_{t-1} \right] = \frac{1}{n^2} \sum_i \text{var}[f(\bar{x}_{0:t}^i) | \mathcal{F}_{t-1}] \\ &= \frac{1}{n^2} \sum_i [\mathbb{E}[|f(\bar{x}_{0:t}^i)|^2 | \mathcal{F}_{t-1}] - |\mathbb{E}[f(\bar{x}_{0:t}^i) | \mathcal{F}_{t-1}]|^2] \\ &= \frac{1}{n^2} \sum_i [K_{t-1} |f|^2(x_{0:t-1}^i) - |K_{t-1} f(x_{0:t-1}^i)|^2] \\ &= \frac{1}{n} \pi_{t-1}^n [K_{t-1} |f|^2 - |K_{t-1} f|^2] \\ &\leq \frac{1}{n} \pi_{t-1}^n K_{t-1} |f|^2 \leq \frac{\|f\|_\infty^2}{n}.\end{aligned}$$

The last inequality is by $K_{t-1}|f|^2(x) = \int |f|^2 K(dy|x) \leq \int \|f\|_\infty^2 K(dy|x) = \|f\|_\infty^2$ and $\pi_{t-1}^n \|f\|_\infty^2 = \frac{\sum_i^n \|f\|_\infty^2}{n} = \|f\|_\infty^2$. Thus we have

$$\begin{aligned} \mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|^2 | \mathcal{F}_{t-1}] &\leq \frac{\|f\|_\infty^2}{n} \\ \mathbb{E}[\mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|^2 | \mathcal{F}_{t-1}]] &\leq \mathbb{E}\left[\frac{\|f\|_\infty^2}{n}\right] \\ \mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|^2] &\leq \frac{\|f\|_\infty^2}{n} \end{aligned} \quad (2.40)$$

and $\lim_{n \rightarrow \infty} \mathbb{E}[|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f|^2] \rightarrow 0$, i.e., (ii) holds due to the Jensen's inequality. Since $\pi_t^n = \frac{1}{n} \sum_{i=1}^n N_t^i \delta_{\{\bar{x}_{0:t}^i\}}$, we have

$$\begin{aligned} \mathbb{E}[\pi_t^n f | \bar{\mathcal{F}}_t] &= \hat{\pi}_t^n f, \\ \mathbb{E}[|\pi_t^n f - \hat{\pi}_t^n f|^2 | \bar{\mathcal{F}}_t] &\leq \frac{1}{n^2} [(q_t^n)^T A_t^n q_t^n] \end{aligned}$$

where A_t^n is the covariance matrix of $(N_t^1, \dots, N_t^n) \sim \mathcal{M}(n, \tilde{w}(\bar{x}_{0:t}^1), \dots, \tilde{w}(\bar{x}_{0:t}^n))$, i.e., it is the matrix with entries given by $\text{var}(N_t^j) = n \cdot \tilde{w}(\bar{x}_{0:t}^j)(1 - \tilde{w}(\bar{x}_{0:t}^j))$, $j = 1, \dots, n$, and $\text{cov}(N_t^j, N_t^k) = -n \cdot \tilde{w}(\bar{x}_{0:t}^j) \tilde{w}(\bar{x}_{0:t}^k)$, $1 \leq j \neq k \leq n$; and q_t^n is the vector with entries $|f(\bar{x}_{0:t}^i)|$, i.e., $(q_t^n) = (|f(\bar{x}_{0:t}^1)|, \dots, |f(\bar{x}_{0:t}^n)|)^T$.

Let us show that $(q_t^n)^T A_t^n q_t^n \leq n$. Indeed (to simplify notation we use only \tilde{w}_t^j instead of full $\tilde{w}(\bar{x}_{0:t}^j)$),

$$\begin{aligned} [(q_t^n)^T A_t^n q_t^n] &= n \left[\sum_{j=1}^n \tilde{w}_t^j (1 - \tilde{w}_t^j) (q_t^j)^2 - 2 \sum_{j,k=1; j \neq k}^n \tilde{w}_t^j \tilde{w}_t^k q_t^j q_t^k \right] \\ &= n \left[\sum_{j=1}^n \tilde{w}_t^j (q_t^j)^2 - \frac{1}{n} \left(\sum_{i=1}^n \tilde{w}_t^i q_t^i \right)^2 \right] \\ &\leq n \left[\sum_{j=1}^n \tilde{w}_t^j (q_t^j)^2 \right] \leq n \|f\|_\infty^2 \left[\sum_{j=1}^n \tilde{w}_t^j \right] = n \|f\|_\infty^2. \end{aligned}$$

Therefore we have

$$\mathbb{E}[|\pi_t^n f - \hat{\pi}_t^n f|^2] = \mathbb{E}[\mathbb{E}[|\pi_t^n f - \hat{\pi}_t^n f|^2 | \bar{\mathcal{F}}_t]] \leq \frac{\|f\|_\infty^2}{n} \quad (2.41)$$

and $\lim_{n \rightarrow \infty} \mathbb{E}[|\pi_t^n f - \hat{\pi}_t^n f|^2] = 0$, i.e., (iii) is again satisfied due to the Jensen's inequality. \square

To end the section we state the theorem which is a direct corollary of the above two theorems.

Theorem 2.4. Let $\bar{\pi}_t^n$ and π_t^n , $t \geq 1$ be the measure valued sequences produced by the SMC algorithm for particle filters presented in Section 2.6.4. Then, for all $t \geq 1$ and $f \in \mathcal{C}_b(\mathbb{C}^{n_x(t+1)})$ we have

$$\mathbb{E}[|\pi_t^n f - \pi_t f|^2] \leq \frac{c_t^2 \|f\|_\infty^2}{n} \quad (2.42)$$

with $c_0 = 1$ and the recursion

$$c_t = c_{t-1} \left(1 + \frac{4 \|g_t\|_\infty}{\bar{\pi}_t g_t} \right). \quad (2.43)$$

Proof. Again, we prove the theorem by induction. Let $h \in \mathcal{C}_b(\mathbb{C}^d)$ and μ^n be an empirical measure given by n i.i.d. samples from μ , ($X_i \sim \mu$), then

$$\mathbb{E}[|\mu^n h - \mu h|^2] \leq \frac{\|h\|_\infty^2}{n}. \quad (2.44)$$

Indeed, consider $\mathbb{E}[|\frac{1}{n} \sum_i h(X_i) - \mathbb{E}[h(X_i)]|^2] = \frac{1}{n^2} \mathbb{E}[|\sum_i h(X_i) - n \mathbb{E}[h(X_i)]|^2] = \frac{1}{n^2} \text{var}(\sum_i h(X_i)) = \frac{n}{n^2} \text{var}(h(X_i))$ due to the i.i.d. character of the samples; and we have $\frac{1}{n} \text{var}(h(X_i)) = \frac{1}{n} (\mathbb{E}|h(X_i)|^2 - |\mathbb{E}[h(X_i)]|^2) \leq \frac{1}{n} \mathbb{E}|h(X_i)|^2 \leq \frac{\|h\|_\infty^2}{n}$.

Setting $h = f$ in (2.44) and $\mu = \pi_0$ gives $\mathbb{E}[|\pi_0^n f - \pi_0 f|^2] \leq \frac{\|f\|_\infty^2}{n}$, i.e., (2.42) holds for $c_0 = 1$.

Now, let $\mathbb{E}[|\pi_{t-1}^n h - \pi_{t-1} h|^2] \leq \frac{c_{t-1}^2 \|h\|_\infty^2}{n}$ and $c_{t-1} \geq 1$ hold for some $t-1 \geq 0$. For $h = K_{t-1} f$ this writes as

$$\mathbb{E}[|\pi_{t-1}^n K_{t-1} f - \pi_{t-1} K_{t-1} f|^2] \leq \frac{c_{t-1}^2 \|K_{t-1} f\|_\infty^2}{n} \leq \frac{c_{t-1}^2 \|f\|_\infty^2}{n} \quad (2.45)$$

The triangle inequality in L_2 norm $\|\cdot\| = (\mathbb{E}[|\cdot|^2])^{1/2}$ reads as

$$\|\bar{\pi}_t^n f - \bar{\pi}_t f\| \leq \|\bar{\pi}_t^n f - \pi_{t-1}^n K_{t-1} f\| + \|\pi_{t-1}^n K_{t-1} f - \pi_{t-1} K_{t-1} f\|.$$

From (2.40) and (2.45) we have ($c_{t-1} \geq 1$)

$$\|\bar{\pi}_t^n f - \bar{\pi}_t f\| \leq \frac{\|f\|_\infty}{\sqrt{n}} + \frac{c_{t-1} \|f\|_\infty}{\sqrt{n}} = \frac{(1 + c_{t-1}) \|f\|_\infty}{\sqrt{n}} \leq \frac{2c_{t-1} \|f\|_\infty}{\sqrt{n}}. \quad (2.46)$$

Formula (2.38) holds also in L_2 norm,

$$\|\widehat{\pi}_t^n f - \pi_t f\| \leq \frac{\|f\|_\infty}{\bar{\pi}_t g_t} \|\bar{\pi}_t g_t - \bar{\pi}_t^n g_t\| + \frac{1}{\bar{\pi}_t g_t} \|\bar{\pi}_t^n f g_t - \bar{\pi}_t f g_t\|,$$

which combined with (2.46) (for functions f and fg_t) gives

$$\|\widehat{\pi}_t^n f - \pi_t f\| \leq \frac{\|f\|_\infty}{\bar{\pi}_t g_t} \frac{2\|g_t\|_\infty c_{t-1}}{\sqrt{n}} + \frac{1}{\bar{\pi}_t g_t} \frac{2\|fg_t\|_\infty c_{t-1}}{\sqrt{n}} \leq \frac{4\|f\|_\infty \|g_t\|_\infty c_{t-1}}{\bar{\pi}_t g_t \sqrt{n}}.$$

Finally, using the triangle inequality and (2.41) ($c_{t-1} \geq 1$),

$$\begin{aligned} \|\pi_t^n f - \pi_t f\| &\leq \|\pi_t^n f - \widehat{\pi}_t^n f\| + \|\widehat{\pi}_t^n f - \pi_t f\| \\ &\leq \frac{\|f\|_\infty}{\sqrt{n}} + \frac{4\|f\|_\infty \|g_t\|_\infty c_{t-1}}{\bar{\pi}_t g_t \sqrt{n}} \\ &\leq \frac{c_{t-1}\|f\|_\infty}{\sqrt{n}} \cdot \left(1 + \frac{4\|g_t\|_\infty}{\bar{\pi}_t g_t}\right). \end{aligned} \quad (2.47)$$

Setting

$$c_t = c_{t-1} \left(1 + \frac{4\|g_t\|_\infty}{\bar{\pi}_t g_t}\right) \quad (2.48)$$

we see that if $c_{t-1} \geq 1$, then also $c_t \geq 1$. Incorporating (2.48) into (2.47) and squaring brings the result. \square

We conclude the chapter by two remarks concerning the proved theorem.

1) The presented inequality and therefore convergence holds for marginal filtering distributions $\pi_{t|t}$ as well. It follows immediately from the theorem considering function f in marginal form $f(x_{0:t}) = f(x_t)$, i.e., if we omit arguments x_s for $s < t$.

2) The sequence $\{c_t\}_{t=0}^\infty$ could be bounded by a geometric series if there would be some upper bound on $4\|g_t\|_\infty/\bar{\pi}_t g_t$ ratio. To establish such a bound we would need a lower bound on $\bar{\pi}_t g_t$ integral. The discussion of this issue is presented in Section 4.3 in Chapter 4.

3. Kernel methods

Kernel methods represent one of the most elaborated areas in the theory of nonparametric estimation. Roughly speaking, in the parametric case there is assumed some strong hypothesis on the distribution of analyzed data. Typically, data are assumed to be generated from a distribution being a member of some parametrized family. Employing appropriate statistical procedures we estimate the values of unknown parameters. The relevant notions in the parametric case are the notions of unbiasedness, consistency and efficiency of estimation. The whole theory of parametric estimation is employed here backed by all tremendous literature available on the topic.

In the nonparametric case, we work only with observed data without any strong hypothesis about their distribution. The only exception is that we adopt the assumption on the absolute continuity of the distribution with respect to the relevant Lebesgue measure, and usually on the i.i.d. character of data. The cornerstone of the nonparametric approach lies in the fact that the distribution of a random variable can be approximated by the sum of Dirac measures which are located at data points sampled from the distribution. The extension of this approach is not only to locate the probability mass point-wise, but also spread it around the observed data points. The techniques of effective spreading of the probability mass is what the kernel methods are concentrated on.

In the chapter we provide a basic overview of kernel methods. We present the result on the convergence of kernel density estimates to the true density of the data distribution. The quantification of the estimate discrepancy is provided by calculating the exact and asymptotic values of the *mean integrated squared error* - the MISE and AMISE analysis. Results are provided for both univariate and multivariate cases. Finally, we present the Fourier analysis approach to MISE calculation which is useful for our own research.

As the presented results are rather classical in the field, there is a bunch of related literature. From this reason not all proofs are supplied but only reference to the source. Our presentation is based on [Parzen 1962; Silverman 1986; Scott 1992; Tarter and Lock 1993; Wand and Jones 1995] and [Tsybakov 2009].

3.1 Histograms

Histograms represent the oldest and most widely used approach to a non-parametric representation of a distribution of empirical data. Histograms are simple and useful tools, but they exhibit several drawbacks which initiated the shift to kernel methods. Nevertheless, it is definitely worth to start with histograms when we are going to study the kernel methods.

An univariate histogram is characterized by the width and placement of its *bins*. Let $b > 0$ be a common width of bins, so-called *binwidth*. If $x_0 \in \mathbb{R}$ is a selected point on the real line, then bins correspond to a system of intervals $[x_0 + mb, x_0 + (m + 1)b)$, $m \in \mathbb{Z}$. That is, each bin is determined by the placement of an interval of length b with endpoints shifting along the real line. It is clear that the union of all bins exhausts the whole real line.

Provided by the data samples X_1, \dots, X_n , where X_i are univariate random variables with the common density f , each bin is associated with its *height* which is the averaged number of data falling into the bin and scaled to the binwidth. Relating the heights of individual bins to the corresponding intervals and representing them graphically in the form of a bar graph yields the typical picture of an univariate histogram. The union of top lines of bars (heights of bins) represents an empirical density function which is in the closed form specified by formula

$$\hat{f}_{\text{HIST}}(x) = \frac{\text{no. of } X_i \text{ in the bin containing } x}{nb}. \quad (3.1)$$

Examples of histograms are presented in Fig. 3.1.

The shape of a histogram, in fact of empirical density (3.1), is primarily affected by the value of binwidth b and consequently by the placement of bins which is controlled by the selection of an origin x_0 . The binwidth is actually the *smoothing parameter* - the larger is the value of b the smoother (less bumpy) is $\hat{f}_{\text{HIST}}(x)$. On the other hand, smoothing leads to passing over local characteristics of the distribution. Apparently, when b increases to infinity then only one bin accommodates all samples and \hat{f}_{HIST} is constant. If b goes to zero, then the number of bins increases to infinity and in the limit each bin contains either no or only one sample and \hat{f}_{HIST} is a highly jagged function. The “right” value of b creates a compromise between two limit cases. Some theory on the proper choice of b is provided in [Scott 1992], but the value of b is usually specified by trials and errors.

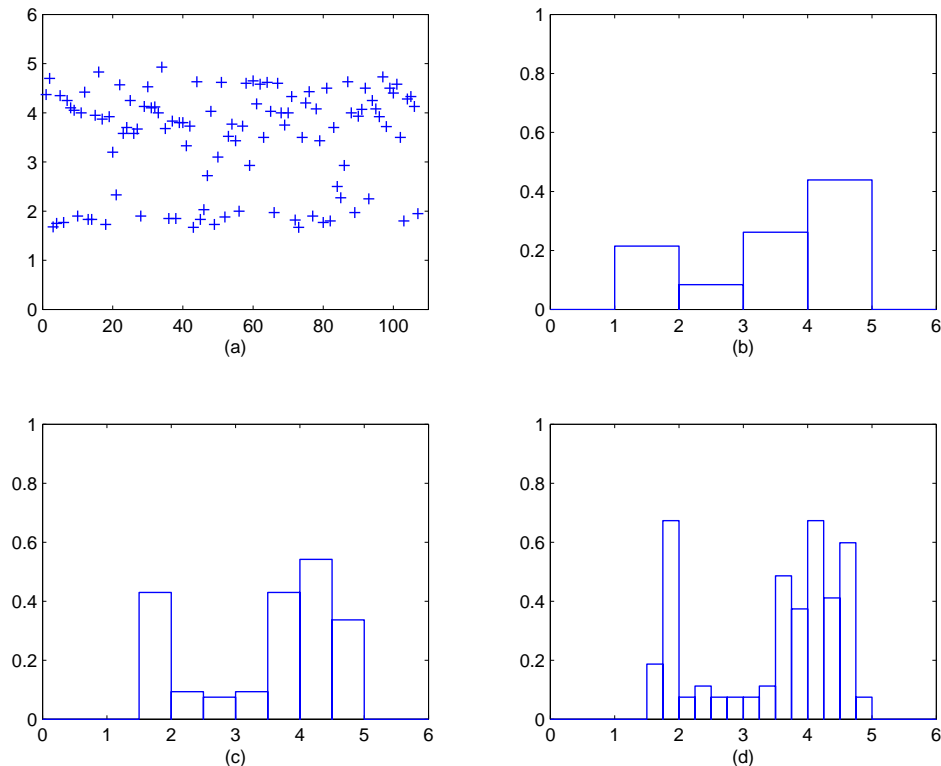


Figure 3.1: Old Faithful geyser data taken from [Silverman 1986] and examples of related histograms. (a) Lengths of $n = 107$ eruptions in minutes; (b) $x_0 = 0$, $b = 1$; (c) $x_0 = 0$, $b = 0.5$; (d) $x_0 = 0$, $b = 0.25$.

The criticism of the usage of histograms stems from the following three facts. 1. \hat{f}_{HIST} has stepwise character - it is only right continuous; 2. the need of a bin-width and placement specification; 3. a problematic extension to multivariate data. Let us comment on these issues.

The discontinuity is incorporated into the very heart of the concept of histogram due to the process of bins creation. In fact, bins are created by the application of the indicator function of semi-open interval, which causes the discontinuity when shifting over the real line. The way how the discontinuity can be removed is to replace the indicator function by some other continuous function. This leads to the introduction of kernel methods.

In order to a histogram be specified, the values of two parameters have to be selected. The binwidth is more influential on the global shape of \hat{f}_{HIST}

than the specification of the point of origin x_0 which affects the placement of bins. However, the placement still plays certain role and there is a natural requirement for methods which have as few degrees of freedom as possible in the given context. The introduction of the naive estimator shows how to circumvent the problem of x_0 selection.

Problems with an extension to the multivariate case raise from the variety of possibilities how a set of univariate intervals filling out the real line can be extended to a set of multidimensional compact sets filling out the respective space. For example, in two dimensions we can opt not only for rectangular bin's bases, but also for triangular or hexagonal ones. This issue is interconnected with the placement problem as well because the placement of bins' bases in multiple dimensions is controlled by a larger set of parameters. Hence again, it would be worth if the placement problem could be somehow eliminated.

The construction of the *naive estimator* is based on the following fact. If f is a density of a random variable X with an absolutely continuous distribution, then

$$f(x) = \lim_{h \rightarrow 0} \frac{1}{2h} P(x - h < X < x + h) \quad (3.2)$$

for each point $x \in \mathbb{R}$. Thus we can approximate the density f by choosing a small number $h > 0$ and look at each point x for the proportion of data falling into the vicinity of x , i.e., into the interval $(x - h, x + h)$. Formally,

$$\hat{f}_{\text{NE}}(x) = \frac{\text{no. of } X_i \in (x - h, x + h)}{2hn}. \quad (3.3)$$

The formula can be rewritten as

$$\hat{f}_{\text{NE}}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} w\left(\frac{x - X_i}{h}\right), \quad \text{where } w(u) = \frac{1}{2} I_{|u| < 1}(u), \quad (3.4)$$

i.e., w is the indicator function of the open interval $(-1, 1)$.

The naive estimator is constructed as a histogram with each point of the real line representing the central point of a bin. This construction solves the placement problem. In spite of the number of bins is now uncountable, the averaging by n remains the finite operation due to the finite number of samples. For binwidth we have clearly $b = 2h$. That is, the binwidth remains an optional parameter.

In Fig. 3.2 there are presented two naive estimates for Old Faithful geyser data introduced in Fig. 3.1. We see that obtained estimates are still jagged. Furthermore, the naive estimator still produces a discontinuous function (not

seen in the figure). The way leading to the removal of this obstacle is an adaptation of the w indicator function in formula (3.4).

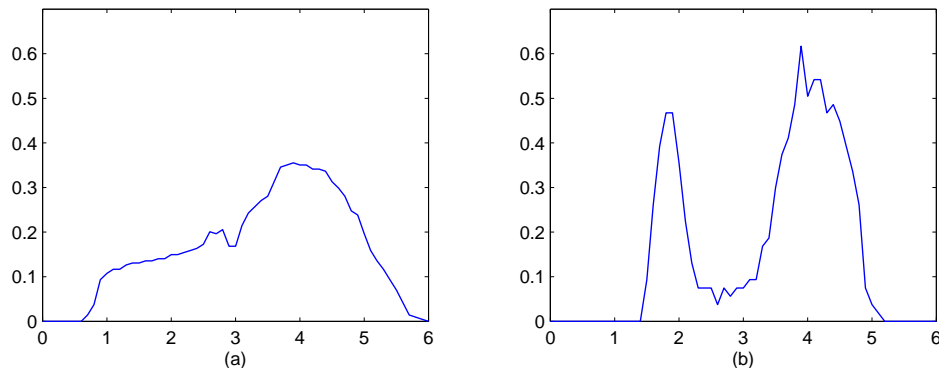


Figure 3.2: Naive estimates for Old Faithful geyser data (a) $h = 1$; (b) $h = 0.25$.

3.2 Kernel methods in one dimension

The development of kernel methods is driven by the ask for obtaining a continuous, sufficiently smooth empirical density which approximates the density of the distribution of a given data sample. The key idea is to approximate each data point by a piece of the probability mass spread around this point. This idea is consistent with the creation process of a histogram, especially under the perspective of the naive estimator introduced above. A bin in a histogram represents a piece of the probability mass related to the indicator function of the corresponding interval. The exact amount of the mass is estimated by the proportion of samples falling into the interval. The consistency of the estimate is guaranteed by SLLN.

We follow the classical approach to kernel methods introduced in [Parzen 1962]. The approach inherits histograms in a natural way. Let X_1, \dots, X_n be a set of independent random variables identically distributed as a random variable X . Let the distribution function of X , $F(x) = P(X \leq x)$, be absolutely continuous with respect to the Lebesgue measure. That is, $F(x) = \int_{-\infty}^x f(u) du$ where $f(u)$ is the corresponding density.

A natural estimate $\hat{F}_n(x)$ of the distribution function $F(x)$ at point x is

$$\hat{F}_n(x) = \frac{\text{no. of } X_i \text{ less or equal to } x}{n}. \quad (3.5)$$

For n and x fixed, $n\widehat{F}_n(x)$ is binomially distributed with $\mathbb{E}[\widehat{F}_n(x)] = F(x)$ and $\text{var}[\widehat{F}_n(x)] = F(x)/n \cdot (1 - F(x))$.

Following the idea of formula (3.2), the straightforward estimation of density f on the basis of $\widehat{F}_n(x)$ is

$$\hat{f}_n(x) = \frac{\widehat{F}_n(x+h) - \widehat{F}_n(x-h)}{2h} = \frac{\text{no. of } X_i \in (x-h, x+h]}{2nh}. \quad (3.6)$$

We see that the above estimate almost corresponds to the naive estimator (3.4). In fact, the naive estimator can be transformed to (3.6) by considering the function $K(u) = \bar{w}(u) = \frac{1}{2}I_{(-1,1]}(u)$, $u \in \mathbb{R}$, instead of function w in definition formula (3.4). The added value from introducing the empirical distribution function $\widehat{F}_n(x)$ is that \hat{f}_n admits the integral representation with respect to \widehat{F}_n . That is,

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right) = \int \frac{1}{h} K\left(\frac{x - s}{h}\right) d\widehat{F}_n(s). \quad (3.7)$$

The above considerations drive the idea of generalization of the naive estimator. Instead of using the discontinuous option $K = w$, we incorporate into formula (3.7) other suitable functions K . These functions are commonly called *kernels* and $h > 0$ is called *bandwidth*. We see that for a given data sample, the choice of K and h completely determines the estimator \hat{f}_n .

The basic requirement on \hat{f}_n bearing a general kernel K is, that \hat{f}_n constitutes a consistent unbiased estimator of $f(x)$. Because of the i.i.d. character of X_i we have

$$\mathbb{E}[\hat{f}_n(x)] = \frac{n}{n} \mathbb{E} \left[\frac{1}{h} K\left(\frac{x - X}{h}\right) \right] = \int \frac{1}{h} K\left(\frac{x - s}{h}\right) f(s) ds \quad (3.8)$$

and $\mathbb{E}[\hat{f}_n(x)]$ does not directly depend on n . If one wants the consistency with respect to n , the only choice is to made bandwidth h varying with n . Since (3.8) is the convolution integral we would expect consistency when $h^{-1}K((x-s)/h)$ goes to the Dirac delta function located at x as $n \rightarrow \infty$.

The conditions which are imposed on K and the evolution of $h(n)$ in order to the class of estimators (3.7) will be consistent are due to [Parzen 1962]. As referred in [Silverman 1986] p. 71, his assumptions on the kernel K are that K is a bounded Borel function, satisfying

$$1) \int |K(u)| du < \infty, \quad 2) \lim_{|u| \rightarrow \infty} |uK(u)| = 0, \quad 3) \int K(u) du = 1$$

and $h(n)$ is assumed to satisfy

$$1) \lim_{n \rightarrow \infty} h(n) = 0, \quad 2) \lim_{n \rightarrow \infty} nh(n) = \infty.$$

Under these conditions it was shown that, provided f is continuous at x , $\hat{f}_n(x) \rightarrow f(x)$ in probability. Thus $\hat{f}_n(x)$ is a consistent estimate of $f(x)$ at the points of continuity of f .

There is a plenty of kernels to choose from to setup a kernel estimator (3.7). Its worth to note that imposed conditions are commonly satisfied if K represents a density of an absolutely continuous distribution. The other common condition is that the kernel is symmetric, i.e., $K(u) = K(-u), u \in \mathbb{R}$. Some examples of kernels are provided in Table 3.1. In the table, $\mu_2(K) = \int u^2 K(u) du$ and $R(K) = \int K^2(u) du$.

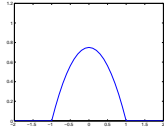
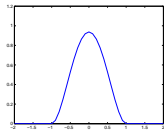
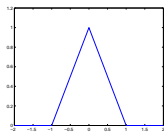
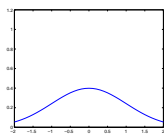
| <i>Kernel</i> | $K(u)$ | <i>Graph</i> | $\mu_2(K)$ | $R(K)$ |
|---------------|---|---|---------------|-------------------------|
| Epanechnikov | $\frac{3}{4}(1 - u^2)_+$ |  | $\frac{1}{5}$ | $\frac{3}{5}$ |
| Biweight | $\frac{15}{16}(1 - u^2)_+^2$ |  | $\frac{1}{7}$ | $\frac{5}{7}$ |
| Triangular | $(1 - u)_+$ |  | $\frac{1}{6}$ | $\frac{2}{3}$ |
| Gaussian | $\frac{1}{\sqrt{2\pi}}e^{-\frac{u^2}{2}}$ |  | 1 | $\frac{1}{2\sqrt{\pi}}$ |

Table 3.1: Examples of kernels.

3.2.1 MISE analysis

The MISE analysis focuses on the evaluation of the quality of density function estimation at the global scale. To achieve this goal, a suitable measure of

discrepancy between empirical density \hat{f}_n and the actual density f must be specified. In the area of kernel methods the *mean integrated squared error* (MISE) is widely used for this purpose. We review the classical results on the analysis of this error.

We start with the *mean squared error* (MSE), which is the suitable measure for the evaluation of the quality of an estimation at a single point. MSE of an empirical density \hat{f}_n at a point x is specified as

$$\text{MSE}_x(\hat{f}_n) = \mathbb{E}[\hat{f}_n(x) - f(x)]^2. \quad (3.9)$$

By the standard properties of mean and variance ($\mathbb{E}X^2 = (\mathbb{E}X)^2 + \text{var}[X]$),

$$\begin{aligned} \text{MSE}_x(\hat{f}_n) &= \mathbb{E}[\hat{f}_n(x) - f(x)]^2 \\ &= (\mathbb{E}[\hat{f}_n(x) - f(x)])^2 + \text{var}[\hat{f}_n(x) - f(x)] \\ &= (\mathbb{E}[\hat{f}_n(x)] - f(x))^2 + \text{var}[\hat{f}_n(x)] \\ &= (b[\hat{f}_n(x)])^2 + \sigma^2[\hat{f}_n(x)]. \end{aligned} \quad (3.10)$$

The first term $b[\hat{f}_n(x)] = \mathbb{E}[\hat{f}_n(x)] - f(x)$ in the MSE decomposition is called the *bias* and the second $\sigma^2[\hat{f}_n(x)] = \text{var}[\hat{f}_n(x)]$ the *variance* of kernel estimate.

Since $\hat{f}_n(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x-X_i}{h}\right)$ and X_i are i.i.d. as X , the respective mean and variance at point x writes as

$$\mathbb{E}[\hat{f}_n(x)] = \frac{n}{n} \mathbb{E} \left[\frac{1}{h} K\left(\frac{x-X}{h}\right) \right] = \int \frac{1}{h} K\left(\frac{x-s}{h}\right) f(s) ds, \quad (3.11)$$

$$\text{var}[\hat{f}_n(x)] = \frac{1}{nh^2} \left(\mathbb{E} \left[K\left(\frac{x-X}{h}\right) \right]^2 - \left[\mathbb{E} K\left(\frac{x-X}{h}\right) \right]^2 \right). \quad (3.12)$$

Inspecting the first formula, we see that $\mathbb{E}[\hat{f}_n(x)]$ does not depend on the number of samples n . Therefore also the bias is independent on n . Contrary, the variance $\text{var}[\hat{f}_n(x)]$ depends on the number of samples via term n^{-1} . Thus, if $\text{var}[\hat{f}_n(x)] < \infty$, then increasing the number of data allows to decrease the variance to an arbitrarily low level.

We will investigate upper bounds on the bias and the variance. The following two lemmas are relevant.

Lemma 3.1. *Assume that the non-negative kernel K satisfies the following conditions*

$$\int K(u) du = 1, \quad \int uK(u) = 0, \quad \int u^2K(u) = \mu_2(K) < \infty$$

and the first derivative of density f is Lipschitz continuous, i.e.,

$$|f'(x) - f'(y)| \leq L_\alpha |x - y|$$

for $L_\alpha > 0$ and all $x, y \in \mathbb{R}$. Then for all $x \in \mathbb{R}$, $h > 0$ and $n \geq 1$ we have

$$|b[\hat{f}_n(x)]| \leq C_b h^2$$

where $C_b = L_\alpha \mu_2(K)$.

Proof. From (3.11),

$$b[\hat{f}_n(x)] = \frac{1}{h} \int K\left(\frac{x-s}{h}\right) f(s) ds - f(x) = \int K(u)[f(x-uh) - f(x)] du.$$

By the mean value theorem

$$f(x) - f(x-uh) = f'(x + \tau uh) \cdot uh$$

for some $-1 \leq \tau \leq 0$. Hence

$$\begin{aligned} b[\hat{f}_n(x)] &= \int K(u) f'(x + \tau uh) (-uh) du \\ &= \int K(u) (-uh) [f'(x + \tau uh) - f'(x)] du, \\ |b[\hat{f}_n(x)]| &\leq \int |K(u)| \cdot |-uh| \cdot |f'(x + \tau uh) - f'(x)| du \\ &\leq \int K(u) \cdot |uh| \cdot L_\alpha |\tau uh| du \\ &\leq L_\alpha \int K(u) \cdot (uh)^2 du = C_b h^2. \quad \square \end{aligned}$$

Lemma 3.2. Let the density f be bounded, i.e., $f(x) \leq \|f\|_\infty < \infty$ for all $x \in \mathbb{R}$ and the kernel K satisfies

$$\int K^2(u) du = R(K) < \infty.$$

Then for any $x \in \mathbb{R}$, $h > 0$ and $n \geq 1$ we have

$$\text{var}[\hat{f}_n(x)] \leq \frac{C_{\sigma^2}}{nh}$$

where $C_{\sigma^2} = \|f\|_\infty \cdot R(K)$.

Proof. We have

$$\mathbb{E} \left[K \left(\frac{x - X}{h} \right) \right]^2 = \int K^2 \left(\frac{x - s}{h} \right) f(s) ds.$$

Therefore from (3.12),

$$\text{var}[\hat{f}_n(x)] \leq \frac{1}{nh^2} \int K^2 \left(\frac{x - s}{h} \right) f(s) ds \leq \frac{\|f\|_\infty h}{nh^2} \int K^2(u) du = \frac{C_{\sigma^2}}{nh}. \quad \square$$

We see that if f and K satisfies the assumptions of the above lemmas, then MSE_x is upper bounded by constants C_b^2 and C_{σ^2} in the following form

$$\text{MSE}_x(\hat{f}_n) \leq C_b^2 h^4 + \frac{C_{\sigma^2}}{nh}. \quad (3.13)$$

The upper bound on the bias and the variance behaves differently with respect to the bandwidth h . For n fixed, if $h \rightarrow 0$ we have unbiasedness, i.e., $\lim_{h \rightarrow 0} \mathbb{E}[\hat{f}_n(x)] = f(x)$. If $h \rightarrow \infty$, i.e., if h increases, then the bound on the bias increases as well. For the variance, we have the opposite behavior. If $h \rightarrow 0$, then the bound on the variance increases and if $h \rightarrow \infty$, then the variance term diminishes.

This behavior tells us that there is the *bias-variance trade-off* we have to take into account when we want to minimize MSE_x by adjusting bandwidth h . In Fig. 3.3 there is presented the typical behavior of MSE_x .

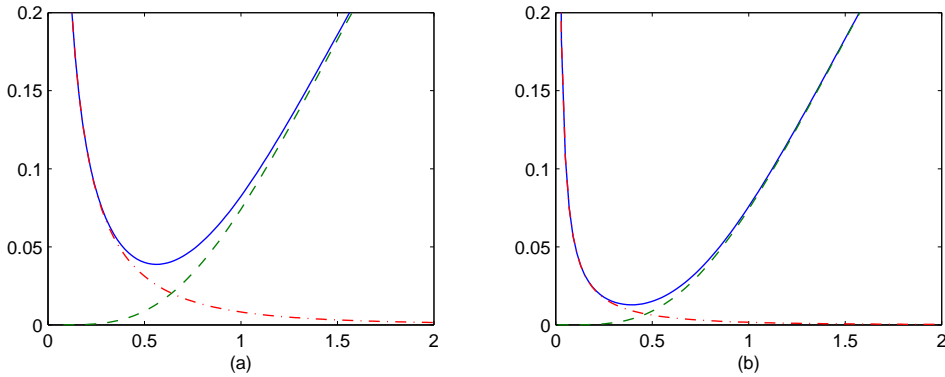


Figure 3.3: Demonstration of the bias (squared, dashed line) variance (dash-dot line) trade-off in MSE_x computation (solid line) with h varying along the x -axis for different choices of n ; (a) $n = 10$; (b) $n = 100$.

MSE_x reflects the quality of approximation at a fixed point. Concerning a global measure, it is natural to take for this purpose the integrated version of MSE_x which is called the *mean integrated squared error* (MISE) and defined as

$$\text{MISE}(\hat{f}_n) = \mathbb{E} \int (\hat{f}_n(x) - f(x))^2 dx. \quad (3.14)$$

Due to the non-negativity of the integrand we can employ Fubini theorem and (3.14) can be rewritten as

$$\begin{aligned} \text{MISE}(\hat{f}_n) &= \int \mathbb{E}[(\hat{f}_n(x) - f(x))^2] dx \\ &= \int \text{MSE}_x(\hat{f}_n) dx \\ &= \int (\mathbb{E}[\hat{f}_n(x)] - f(x))^2 dx + \int \text{var}[\hat{f}_n(x)] dx. \end{aligned} \quad (3.15)$$

The goal of the MISE analysis is to obtain the dependence of MISE on bandwidth h , kernel K and the number of data samples n in some tractable form. As n and K are usually given, the typical task is to search for h^* which minimizes MISE for n and K fixed. h^* is then considered as the optimal bandwidth which should be used for the construction of the kernel estimator based on K .

Unfortunately, the expression of above integrals in a closed form is generally impossible. Closed forms can be stated only in some special cases, for example, when the true density f and K correspond to the normal density functions. The standard approach for a general case is to investigate the asymptotic behavior of presented integrals when n goes to infinity and h changes appropriately. This will be discussed in the next section. However, before we do this let us review the exact MISE calculations for the Gaussian density.

Gaussian density and its mixtures. Let $\phi_\sigma(x)$ be the density of normal distribution $\mathcal{N}(0, \sigma^2)$. $\phi_\sigma(x - \mu)$ is then the density of general normal distribution $\mathcal{N}(\mu, \sigma^2)$. Let $K(u) = \phi_1(u)$, i.e., kernel K corresponds to the density of $\mathcal{N}(0, 1)$. The result of Fryer and Deheuvels referred in [Wand and Jones 1995] p. 25 or in [Silverman 1986] p. 37, states that

$$\text{MISE}(\hat{\phi}_\sigma(x - \mu)_n) = \frac{1}{2\sqrt{\pi}} \left[\frac{1}{nh} + \left(1 - \frac{1}{n}\right) \frac{1}{\sqrt{\sigma^2 + h^2}} + \frac{1}{\sigma} - \frac{2\sqrt{2}}{\sqrt{2\sigma^2 + h^2}} \right]. \quad (3.16)$$

The calculation can be extended to mixtures of normal densities which are the densities of form

$$f(x) = \sum_{j=1}^m w_j \phi_{\sigma_j}(x - \mu_j)$$

where weights w_j add to unity, i.e., $\sum_j w_j = 1$. In this case the exact MISE formula reads as

$$\text{MISE}(\hat{f}_n) = (2\pi^{1/2}nh)^{-1} \cdot \mathbf{w}^T [(1 - n^{-1})\mathbf{\Omega}_2 - 2\mathbf{\Omega}_1 + \mathbf{\Omega}_0]\mathbf{w} \quad (3.17)$$

where $\mathbf{w} = (w_1, \dots, w_m)^T$ is the vector of weights and $\mathbf{\Omega}_a$ is the $m \times m$ matrix having (j, j') entry equal to

$$\phi_{(ah^2 + \sigma_j^2 + \sigma_{j'}^2)^{-1/2}}(\mu_j - \mu_{j'}).$$

for $a \in \{0, 1, 2\}$.

A minimizer h^* of (3.16) and (3.17) can be found by standard optimization techniques.

3.2.2 AMISE analysis

As the computation of MISE in a closed form is generally intractable, the common approach is to examine the asymptotic behavior of the variables of interest. That is, to investigate the asymptotic behavior of $\text{MISE}(\hat{f}_n)$ formula (3.15) for $n \rightarrow \infty$ and $h = h(n)$ varying accordingly. This leads to the specification of AMISE (asymptotic MISE) characteristic of the quality of approximation by the empirical density \hat{f}_n .

The AMISE analysis is in details treated in [Scott 1992] and [Wand and Jones 1995]. The main result obtained here can be stated as follows. Let f , h and K satisfies the following conditions. 1. The density f is such that its second derivative f'' is continuous, square integrable and ultimately monotone, i.e., it is motone over both $(-\infty, -M)$ and (M, ∞) for some $M > 0$; 2. bandwidth h varies with n in such a way that

$$\lim_{n \rightarrow \infty} h(n) = 0, \quad \lim_{n \rightarrow \infty} nh(n) = \infty.$$

and 3. the kernel K corresponds to a non-negative function satisfying

$$\int K(u) du = 1, \quad \int uK(u) du = 0, \quad \int u^2K(u) du = \mu_2(K) < \infty.$$

Under these conditions $\text{MISE}(\hat{f}_n) = \text{AMISE}(\hat{f}_n) + o((nh)^{-1} + h^4)$ and

$$\begin{aligned} \text{AMISE}(\hat{f}_n) &= \frac{R(K)}{nh} + \frac{1}{4}\mu_2(K)^2 h^4 R(f''), \\ h_{\text{AMISE}}^* &= \left[\frac{R(K)}{\mu_2(K)^2 R(f'')} \right]^{1/5} \frac{1}{n^5}, \\ \text{AMISE}(h_{\text{AMISE}}^*) &= \frac{5}{4} [\mu_2(K)^2 R(K)^4 R(f'')]^{1/5} \frac{1}{n^{4/5}} \end{aligned} \quad (3.18)$$

where $R(K) = \int K^2(u) du < \infty$ and $R(f'') = \int (f''(u))^2 du < \infty$.

AMISE provides an useful large sample approximation to the MISE and enables to find directly its bandwidth minimizer h_{AMISE}^* . On the other hand, the main disadvantage of the asymptotic approach is that AMISE tightly depends on the properties of f via $R(f'')$ term which is usually not known apriori.

3.3 Kernel methods in multiple dimensions

The ideas underlying an application of kernel methods in multiple dimensions are similar to the univariate case. Basically, a multivariate kernel refers to a piece of the probability mass spread around a sample point. Individual kernels are additively combined and normalized to obtain an approximation of the actual density of the data distribution.

However, with the increasing dimension we get more degrees of freedom and several adjustment of univariate considerations has to be taken into account. It results into a generally broader discussion. Because of the lack of space we review here only the very basic results without details. The interested reader is referred to [Scott 1992; Wand and Jones 1995] which based our review. Nevertheless, the Fourier analysis approach presented in Section 3.4 relates well to the multivariate case with the presence of results relevant to our own research. Thus, the restricted extent does not cause any substantial limitations.

To start let us make a note concerning the notation. In the previous chapter we worked with elements from $(\mathbb{R}^d)^{(t+1)}$ space. These elements were referred by $x_{0:t}$ notation emphasizing the “time-dimension”, but not explicitly referring to the “space-dimension” d . In multivariate kernel methods we work with static samples, i.e., only with elements from \mathbb{R}^d space. In order to distinguish obtained results from the univariate ones we will write multivariate elements in the bold notation.

In what follows we will work with a random sample $\mathbf{X}_1, \dots, \mathbf{X}_n$ of n independent variables identically distributed as a d -dimensional random vector \mathbf{X} . The distribution of \mathbf{X} is assumed to be absolutely continuous with respect to the d -dimensional Lebesgue measure with the density $f(\mathbf{x}) = f(x_1, \dots, x_n)$. We approximate this density by a multivariate kernel estimator.

Following [Wand and Jones 1995], the most general form of a d -dimensional multivariate kernel estimator is

$$\hat{f}_n(\mathbf{x}; \mathbf{H}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(\mathbf{x} - \mathbf{X}_i), \quad (3.19)$$

where \mathbf{H} is a symmetric positive definite $d \times d$ matrix, called the *bandwidth matrix*, $K_{\mathbf{H}}(\mathbf{u}) = |\mathbf{H}|^{-1/2} K(\mathbf{H}^{-1/2} \mathbf{u})$ (in the formula $|\cdot|$ denotes the determinant and for the inverse matrix we have $\mathbf{H}^{-1} = \mathbf{H}^{-1/2} \mathbf{H}^{-1/2}$) and K is a d -variate kernel satisfying $\int K(\mathbf{u}) d\mathbf{u} = 1$.

There are two approaches to the generation of multivariate kernels from a symmetric univariate kernel κ . The first leads to so-called *product kernels* and the second to *radially symmetric kernels*. The definition formulas are

$$K^P(\mathbf{u}) = \prod_{i=1}^d \kappa(u_i), \quad K^R(\mathbf{u}) = c_{\kappa,d}^{-1} \cdot \kappa((\mathbf{u}^T \mathbf{u})^{1/2}), \quad (3.20)$$

where $c_{\kappa,d}^{-1}$ is the normalization constant, i.e., $c_{\kappa,d}^{-1} = \int \kappa((\mathbf{u}^T \mathbf{u})^{1/2}) d\mathbf{u}$.

A popular choice for a kernel K is the standard d -variate normal density

$$K(\mathbf{u}) = (2\pi)^{-d/2} \exp \left[-\frac{1}{2} \mathbf{u}^T \mathbf{u} \right] \quad (3.21)$$

in which case $K_{\mathbf{H}}(\mathbf{u} - \mathbf{X}_i)$ corresponds to the density of $\mathcal{N}(\mathbf{X}_i, \mathbf{H})$ in vector \mathbf{u} . This kernel is both product and radial and is constructed from univariate normal densities.

In practice, matrix \mathbf{H} is restricted to some simpler class in order to decrease the number of parameters determining a kernel. The common choices are $\mathbf{H} = \text{diag}(h_1^2, \dots, h_d^2)$, i.e., \mathbf{H} is a diagonal matrix or, which is the preferred case, $\mathbf{H} = h^2 \mathbf{I}$. In the last case, the explicit formula for the kernel estimator writes as

$$\hat{f}_n(\mathbf{x}; h^2 \mathbf{I}) = \frac{1}{nh^d} \sum_{i=1}^n K \left(\frac{\mathbf{x} - \mathbf{X}_i}{h} \right). \quad (3.22)$$

3.3.1 MISE analysis

The counterpart of univariate MISE formula (3.15) writes as

$$\text{MISE}(\hat{f}_n; \mathbf{H}) = \int (\mathbb{E}[\hat{f}_n(\mathbf{x}; \mathbf{H})] - f(\mathbf{x}))^2 d\mathbf{x} + \int \text{var}(\hat{f}_n(\mathbf{x}; \mathbf{H})) d\mathbf{x}. \quad (3.23)$$

As one would expect, the computation of the above multidimensional integrals in a closed form is generally impossible. The standard approach is to perform an asymptotic analysis or to consider some special cases. As in one dimension, such the special case is the computation of MISE for data generated from a mixture of Gaussian distributions. Let us review the relevant results.

Mixture of Gaussian densities. Consider f to be a mixture of multivariate normal densities

$$f(\mathbf{x}) = \sum_{j=1}^m w_j \phi_{\Sigma_j}(\mathbf{x} - \boldsymbol{\mu}_j),$$

where $\mathbf{w} = (w_1, \dots, w_m)$ is a weights vector, $\sum_{j=1}^m w_j = 1$, ϕ_{Σ_j} is the density of $\mathcal{N}(\mathbf{0}, \Sigma_j)$, i.e.,

$$\phi_{\Sigma_j}(\mathbf{u}) = (2\pi)^{-d/2} |\Sigma_j|^{-1/2} \exp \left[-\frac{1}{2} \mathbf{u}^T \Sigma_j^{-1} \mathbf{u} \right]$$

and $\boldsymbol{\mu}_j$ is the vector of means and Σ_j is a covariance matrix.

Let Ω_a , $a \in \{0, 1, 2\}$ be a $m \times m$ matrix with the (j, j') entry, $j, j' = 1, \dots, m$, equal to $\phi_{a\mathbf{H} + \Sigma_j + \Sigma_{j'}}(\boldsymbol{\mu}_j - \boldsymbol{\mu}_{j'})$. Then

$$\text{MISE}(\hat{f}_n; \mathbf{H}) = n^{-1} (4\pi)^{-d/2} |\mathbf{H}|^{-1/2} + \mathbf{w}^T [(1 - n^{-1})\Omega_2 - 2\Omega_1 + \Omega_0] \mathbf{w}. \quad (3.24)$$

This formula is the multivariate counterpart of univariate formula (3.17).

3.3.2 AMISE analysis

AMISE analysis is based on the application of the multidimensional Taylor's theorem on $\mathbb{E}[\hat{f}_n(\mathbf{x}; \mathbf{H})]$ and $\text{var}[\hat{f}_n(\mathbf{x}; \mathbf{H})]$ terms of (3.23). Further, certain integrability and limit conditions are imposed on f , K and \mathbf{H} . They include (for details see [Wand and Jones 1995] p. 95) the assumption that each entry of Hessian of f is piecewise continuous and square integrable or that K is a bounded, compactly supported d -variate kernel satisfying

$$\int K(\mathbf{u}) d\mathbf{u} = 1, \quad \int \mathbf{u} K(\mathbf{u}) d\mathbf{u} = 0, \quad \int \mathbf{u} \mathbf{u}^T K(\mathbf{u}) d\mathbf{u} = \mu_2(K) \mathbf{I}_d,$$

where $\mu_2(K) = \int u_j K(\mathbf{u}) d\mathbf{u} < \infty$ is independent of j , $j = 1, \dots, d$.

Under the mentioned conditions the following result for AMISE is achieved for the $\mathbf{H} = h^2 \mathbf{I}_d$ case

$$\text{AMISE}(\hat{f}_n, h^2 \mathbf{I}_d) = n^{-1} h^{-d} R(K) + \frac{1}{4} h^4 \mu_2(K)^2 \int (\nabla^2 f(\mathbf{x}))^2 d\mathbf{x}$$

where $R(K) = \int K^2(\mathbf{u}) d\mathbf{u} < \infty$ and $\nabla^2 f(\mathbf{x}) = \sum_{j=1}^d (\partial^2 / \partial x_j^2) f(\mathbf{x})$.

Dealing with (3.25), the optimal bandwidth h^* (and consequently the optimal matrix $\mathbf{H}^* = h^{*2} \mathbf{I}_d$) can be identified as

$$h_{\text{AMISE}}^* = \left[\frac{d \cdot R(K)}{n \cdot \mu_2(K)^2 \int (\nabla^2 f(\mathbf{x}))^2 d\mathbf{x}} \right]^{1/(d+4)}.$$

Plugging h_{AMISE}^* back into the error formula we get

$$\text{AMISE}(h_{\text{AMISE}}^*) = \frac{d+4}{4d} \left(\mu_2(K)^{2d} (dR(K))^4 \left[\int (\nabla^2 f(\mathbf{x}))^2 d\mathbf{x} \right]^d n^{-4} \right)^{1/(d+4)}. \quad (3.25)$$

By this result we conclude the review of classical results concerning the MISE and AMISE analysis of kernel estimators. Let us sum up what we referred to.

- In both cases of univariate and multivariate kernel estimations we were able to set exact MISE computations for i.i.d. data generated from Gaussian mixtures when the Gaussian kernel was employed.
- Asymptotic results were provided in both univariate and multivariate case. However, these results are based on the knowledge of certain analytic properties of the generating density. That is why [Tsybakov 2009] provides a criticism of the asymptotic approach. The criticism stems from the fact that the results of AMISE analysis are tightly related to the properties of only a single density f and not to common properties of a sufficiently broad class of densities the density f is assumed to belong to. He proposed a more general approach for the MISE analysis based on the Fourier analysis. Unfortunately, his results are presented only for univariate case. In the next section we review and extend these results to the multivariate case.
- In the univariate case, we have presented an upper bound on MSE_x for an univariate density, formula (3.13). Employing the Fourier analysis the similar upper bound can be stated directly on multivariate MISE (and therefore on the univariate one as a special case). Thus, there is another reason to inspect the Fourier analysis of kernel estimators.

3.4 Fourier analysis of kernel estimators

This section deals with the Fourier analysis of kernel estimators. Interestingly enough, the employment of characteristic functions of empirical distributions provides a tractable and accurate error analysis with no extra effort. The Fourier analysis approach thus complements and extends asymptotic results presented in the previous sections. Let us first review some basic facts.

Let \mathbf{X} be a d -dimensional random vector with a joint distribution $P_{\mathbf{X}}$. The characteristic function $\phi_{\mathbf{X}}(\mathbf{t}): \mathbb{R}^d \rightarrow \mathbb{C}$ of \mathbf{X} is defined as

$$\phi_{\mathbf{X}}(\mathbf{t}) = \mathbb{E}[e^{i\langle \mathbf{t}, \mathbf{X} \rangle}] = \int e^{i\langle \mathbf{t}, \mathbf{X} \rangle} dP_{\mathbf{X}}, \quad \mathbf{t} \in \mathbb{R}^d. \quad (3.26)$$

The definition formula shows that the characteristic function of a distribution is provided by an integral transform. It is well known [Lachout 2004] that this transform provides the full characterization of a given distribution and it can be seen as the transform from the space of probabilistic measures to the space of functions from \mathbb{R}^d to \mathbb{C} .

The other quite common view of the Fourier transform comes from the area of applied mathematics when for a given integrable function $f: \mathbb{R}^d \rightarrow \mathbb{R}$, i.e., $f \in L_1(\mathbb{R}^d)$, (a signal in electrical engineering) its Fourier transform is specified as

$$\mathcal{F}[f](\boldsymbol{\omega}) = \int e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) d\mathbf{x}, \quad \boldsymbol{\omega} \in \mathbb{R}^d. \quad (3.27)$$

Formula (3.27) can be treated as a special case of formula (3.26) when the distribution of \mathbf{X} is absolutely continuous with respect to the respective d -dimensional Lebesgue measure and has density $f(\mathbf{x})$, i.e., $dP_{\mathbf{X}} = f(\mathbf{x}) d\mathbf{x}$. On the other hand, in (3.27) f need not be necessarily a density. Only integrability is assumed.

Let $f, g \in L_1(\mathbb{R}^d)$, then the following properties of the Fourier transform are relevant to our research:

- boundedness: $|\mathcal{F}[f](\boldsymbol{\omega})| \leq 1$, for f being a density
- linearity: $\mathcal{F}[af + bg](\boldsymbol{\omega}) = a\mathcal{F}[f](\boldsymbol{\omega}) + b\mathcal{F}[g](\boldsymbol{\omega})$, $a, b \in \mathbb{R}$
- shifting: $\mathcal{F}[f(\mathbf{x} - \mathbf{s})](\boldsymbol{\omega}) = e^{i\langle \boldsymbol{\omega}, \mathbf{s} \rangle} \mathcal{F}[f](\boldsymbol{\omega})$, $\mathbf{s} \in \mathbb{R}^d$
- scaling: $\mathcal{F}[f(\mathbf{x}/h)/h^d](\boldsymbol{\omega}) = \mathcal{F}[f](h\boldsymbol{\omega})$, $h > 0$
- shifting & scaling: $\mathcal{F}[f((\mathbf{x} - \mathbf{s})/h)/h^d] = e^{i\langle \boldsymbol{\omega}, \mathbf{s} \rangle} \mathcal{F}[f](h\boldsymbol{\omega})$, $\mathbf{s} \in \mathbb{R}^d$

- complex conjugate: $\overline{\mathcal{F}[f](\boldsymbol{\omega})} = \mathcal{F}[f](-\boldsymbol{\omega})$
- convolution: $\mathcal{F}[f * g](\boldsymbol{\omega}) = \mathcal{F}[f](\boldsymbol{\omega})\mathcal{F}[g](\boldsymbol{\omega})$
- symmetry: if $f(-\mathbf{x}) = f(\mathbf{x})$, then $\mathcal{F}[f](-\boldsymbol{\omega}) = \mathcal{F}[f](\boldsymbol{\omega})$
- isometry, provided by the Plancherel's formula for $f \in L_2(\mathbb{R}^d)$:

$$\int_{\mathbb{R}^d} f^2(\mathbf{x}) d\mathbf{x} = \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} |\mathcal{F}[f](\boldsymbol{\omega})|^2 d\boldsymbol{\omega}.$$

In the text bellow, we extend to the multivariate case Tsybakov's results provided for the univariate case [Tsybakov 2009]. Tsybakov works with the engineering view of the transform, however, the Fourier transform of empirical measures he employs has to be defined via the characteristic function's point of view. The reason is that there are no densities of empirical measures (in fact of the Dirac measure) with respect to the corresponding Lebesgue measure.

Let $\widehat{F}_n(\mathbf{x}) = \delta_n(d\mathbf{x}) = \sum_{j=1}^n \delta_{\mathbf{x}_j}(d\mathbf{x})$ be the empirical distribution associated with an i.i.d. sample of particles $\{\mathbf{X}_j\}_{j=1}^n$. The characteristic function of this distribution writes as¹

$$\phi_n(\boldsymbol{\omega}) = \int e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} d\widehat{F}_n(\mathbf{x}) = \frac{1}{n} \sum_{j=1}^n e^{i\langle \boldsymbol{\omega}, \mathbf{X}_j \rangle}, \quad \boldsymbol{\omega} \in \mathbb{R}^d. \quad (3.28)$$

Note that $\phi_n(\boldsymbol{\omega})$ constitutes a random variable for $\boldsymbol{\omega}$ fixed.

Consider the standard multivariate kernel estimator

$$\widehat{f}_n(\mathbf{x}) = \frac{1}{nh^d} \sum_{j=1}^n K\left(\frac{\mathbf{x} - \mathbf{X}_j}{h}\right) \quad (3.29)$$

providing an empirical density \widehat{f}_n . By the properties of the Fourier transform

$$\mathcal{F}[\widehat{f}_n](\boldsymbol{\omega}) = \frac{1}{n} \sum_{j=1}^n \mathcal{F}\left[\frac{1}{h^d} K\left(\frac{\mathbf{x} - \mathbf{X}_j}{h}\right)\right] = \frac{1}{n} \sum_{j=1}^n e^{i\langle \boldsymbol{\omega}, \mathbf{X}_j \rangle} \mathcal{F}[K](h\boldsymbol{\omega}).$$

Writing $K_{\mathcal{F}}(\boldsymbol{\omega})$ for $\mathcal{F}[K](\boldsymbol{\omega})$ we obtain for the characteristic function of \widehat{f}_n the expression

$$\mathcal{F}[\widehat{f}_n](\boldsymbol{\omega}) = \phi_n(\boldsymbol{\omega}) K_{\mathcal{F}}(h\boldsymbol{\omega}). \quad (3.30)$$

In order to study the properties of MISE error the following lemma is relevant.

¹In what follows we use $\boldsymbol{\omega}$ for argument of characteristic function instead of \mathbf{t} .

Lemma 3.3. Let $\{\mathbf{X}_j\}_{j=1}^n$ be an i.i.d. sample from a distribution with the density f . Let the characteristic function of \mathbf{X}_j be $\phi(\boldsymbol{\omega})$. Then for ϕ_n of (3.28) we have

- (i) $\mathbb{E}[\phi_n(\boldsymbol{\omega})] = \phi(\boldsymbol{\omega})$
- (ii) $\mathbb{E}[|\phi_n(\boldsymbol{\omega})|^2] = \left(1 - \frac{1}{n}\right) |\phi(\boldsymbol{\omega})|^2 + \frac{1}{n}$
- (iii) $\mathbb{E}[|\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2] = \frac{1}{n}(1 - |\phi(\boldsymbol{\omega})|^2)$

Proof. To show (i) consider the i.i.d. character of \mathbf{X}_j .

$$\mathbb{E}[\phi_n(\boldsymbol{\omega})] = \frac{1}{n} \sum_{j=1}^n \int e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) d\mathbf{x} = \frac{1}{n} \sum_{j=1}^n \phi(\boldsymbol{\omega}) = \phi(\boldsymbol{\omega}). \quad (3.31)$$

To show (ii) note that

$$\begin{aligned} \mathbb{E}[|\phi_n(\boldsymbol{\omega})|^2] &= \mathbb{E}[\phi_n(\boldsymbol{\omega}) \overline{\phi_n(\boldsymbol{\omega})}] = \mathbb{E}[\phi_n(\boldsymbol{\omega}) \phi_n(-\boldsymbol{\omega})] \\ &= \mathbb{E} \left[\frac{1}{n^2} \sum_{j,k:j \neq k} e^{i\langle \boldsymbol{\omega}, \mathbf{X}_j \rangle} e^{-i\langle \boldsymbol{\omega}, \mathbf{X}_k \rangle} \right] + \frac{n}{n^2} \\ &= \frac{1}{n^2} \sum_{j,k:j \neq k} \mathbb{E} [e^{i\langle \boldsymbol{\omega}, \mathbf{X}_j \rangle}] \mathbb{E} [e^{-i\langle \boldsymbol{\omega}, \mathbf{X}_k \rangle}] + \frac{1}{n} \\ &= \frac{n^2 - n}{n^2} \phi(\boldsymbol{\omega}) \phi(-\boldsymbol{\omega}) + \frac{1}{n} \\ &= \left(1 - \frac{1}{n}\right) |\phi(\boldsymbol{\omega})|^2 + \frac{1}{n}. \end{aligned} \quad (3.32)$$

Case (iii) follows from (ii) a (i). Indeed,

$$\begin{aligned} \mathbb{E}[|\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2] &= \mathbb{E}[(\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})) \overline{(\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega}))}] \\ &= \mathbb{E}[(\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega}))(\phi_n(-\boldsymbol{\omega}) - \phi(-\boldsymbol{\omega}))] \\ &= \mathbb{E}[|\phi_n(\boldsymbol{\omega})|^2 - \phi_n(\boldsymbol{\omega})\phi(-\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})\phi_n(-\boldsymbol{\omega}) + |\phi(\boldsymbol{\omega})|^2] \\ &= \mathbb{E}[|\phi_n(\boldsymbol{\omega})|^2] - \phi(\boldsymbol{\omega})\phi(-\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})\phi(-\boldsymbol{\omega}) + |\phi(\boldsymbol{\omega})|^2 \\ &= \mathbb{E}[|\phi_n(\boldsymbol{\omega})|^2] - |\phi(\boldsymbol{\omega})|^2 \\ &= \left(1 - \frac{1}{n}\right) |\phi(\boldsymbol{\omega})|^2 + \frac{1}{n} - |\phi(\boldsymbol{\omega})|^2 \\ &= \frac{1}{n}(1 - |\phi(\boldsymbol{\omega})|^2). \quad \square \end{aligned} \quad (3.33)$$

Now we can proceed with the MISE computation for kernel estimate (3.29). We assume that both density f and kernel K belong to $L_2(\mathbb{R}^d)$. Employing

the Plancherel's theorem we have

$$\begin{aligned}
\text{MISE}(\hat{f}_n) &= \mathbb{E} \int (\hat{f}_n(\mathbf{x}) - f(\mathbf{x}))^2 d\mathbf{x} \\
&= \frac{1}{(2\pi)^d} \mathbb{E} \int |\mathcal{F}[\hat{f}_n](\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
&= \frac{1}{(2\pi)^d} \mathbb{E} \int |\phi_n(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}. \tag{3.34}
\end{aligned}$$

The following theorem provides the exact MISE computation of estimate \hat{f}_n for any fixed n .

Theorem 3.1. *Let the density f and kernel K be in $L^2(\mathbb{R}^d)$. Then for all $n \geq 1$ and $h > 0$ the mean integrated squared error of the kernel estimator \hat{f}_n has the form*

$$\begin{aligned}
\text{MISE}(\hat{f}_n) &= \frac{1}{(2\pi)^d} \left[\int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} + \frac{1}{n} \int |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \right] \\
&\quad - \frac{1}{(2\pi)^d} \frac{1}{n} \int |\phi(\boldsymbol{\omega})|^2 |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega}. \tag{3.35}
\end{aligned}$$

Proof. As $\phi, K \in L_2(\mathbb{R}^d)$ and $|\phi(\boldsymbol{\omega})| \leq 1$ for all $\boldsymbol{\omega} \in \mathbb{R}^d$, all the integrals are finite. To obtain the Fourier MISE formula it suffices to develop (3.34),

$$\begin{aligned}
&\mathbb{E} \int |\phi_n(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
&= \mathbb{E} \int |(\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega}))K_{\mathcal{F}}(h\boldsymbol{\omega}) - (1 - K_{\mathcal{F}}(h\boldsymbol{\omega}))\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
&= \int \mathbb{E}[|\phi_n(\boldsymbol{\omega}) - \phi(\boldsymbol{\omega})|^2] |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 + |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
&= \frac{1}{n} \int (1 - |\phi(\boldsymbol{\omega})|^2) |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} + \int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}
\end{aligned}$$

After rearranging we obtain the assertion of the theorem. \square

We see that the formula (3.35) is composed of three terms. Let us discuss these terms separately.

3.4.1 The first term of Fourier MISE decomposition

Concerning the properties of the first term of Fourier MISE formula (3.35) we are able to say something more concrete for the so-called *Sobolev class of*

densities. To make a preparation for our own research we extend the original Tsybakov's definition [Tsybakov 2009] to multiple dimensions for integer β s.

Definition 3.1. Let $\beta \geq 1$ be an integer and $L > 0$. The Sobolev class of densities $\mathcal{P}_{S(\beta,L)}$ consists of all probability density functions $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying

$$\int \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \leq (2\pi)^d L^2 \quad (3.36)$$

where $\phi = \mathcal{F}[f]$ and $\|\cdot\|$ is the Euclidean norm.

The condition (3.36) may look strange, but it actually concerns the boundedness of partial derivatives of f . Let us show this for the basic case of $\beta = 1$.

Lemma 3.4. Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$, $d \in \mathbb{N}$ be a multivariate density and $f_j = \partial f / \partial x_j$. Let $\int (f_j(\mathbf{x}))^2 d\mathbf{x} \leq L_j^2$ for some $L_j > 0$, $j = 1, \dots, d$. Then (3.36) holds for $\beta = 1$, i.e., $\int \|\boldsymbol{\omega}\|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \leq (2\pi)^d \|\mathbf{L}\|^2$.

Proof. First of all note that $\|\boldsymbol{\omega}\|^2 = \sum_j |\omega_j|^2$. Further, let us show that $\mathcal{F}[f_j](\boldsymbol{\omega}) = (-i\omega_j)\phi(\boldsymbol{\omega})$. We have

$$\mathcal{F}[f_j](\boldsymbol{\omega}) = \int e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f_j(\mathbf{x}) d\mathbf{x} = \int_{k \neq j} \int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f_j(\mathbf{x}) dx_j dx_{k \neq j}.$$

By per partes: $u = e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle}$, $u' = \partial u / \partial x_j = i\omega_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle}$, $v' = f_j(\mathbf{x})$, $v = f(\mathbf{x})$, the inner integral writes as

$$\int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f_j(\mathbf{x}) dx_j = e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) \Big|_{x_j=-\infty}^{x_j=\infty} - \int_j (i\omega_j) e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) dx_j.$$

The first term is zero because $f(\mathbf{x})$ is a density, i.e.,

$$\begin{aligned} \int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f_j(\mathbf{x}) dx_j &= (-i\omega_j) \int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) dx_j, \\ \int_{k \neq j} \int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f_j(\mathbf{x}) dx_j dx_{k \neq j} &= (-i\omega_j) \int_{k \neq j} \int_j e^{i\langle \boldsymbol{\omega}, \mathbf{x} \rangle} f(\mathbf{x}) dx_j dx_{k \neq j}, \\ \mathcal{F}[f_j](\boldsymbol{\omega}) &= (-i\omega_j)\phi(\boldsymbol{\omega}). \end{aligned}$$

By the Plancherel's theorem, $\int |\mathcal{F}[f_j](\boldsymbol{\omega})|^2 d\boldsymbol{\omega} = (2\pi)^d \int (f_j(\mathbf{x}))^2 d\mathbf{x}$, we obtain

$$\begin{aligned} \int |\omega_j|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &= (2\pi)^d \int (f_j(\mathbf{x}))^2 d\mathbf{x} \leq (2\pi)^d L_j^2, \\ \sum_j \int |\omega_j|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq (2\pi)^d \sum_j L_j^2, \\ \int \left(\sum_j |\omega_j|^2 \right) |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq (2\pi)^d \sum_j L_j^2, \\ \int \|\boldsymbol{\omega}\|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq (2\pi)^d \|\mathbf{L}\|^2. \quad \square \end{aligned}$$

The other ingredients for setting up the properties of the Fourier MISE formula are the properties of the employed kernel with respect to the Sobolev class of densities. Let us start with the notion of the *order of kernel*.

Definition 3.2. Let $\ell \geq 1$ be an integer. We say that the kernel $K: \mathbb{R}^d \rightarrow \mathbb{R}$ is of order ℓ , if K is symmetric, $L_2(\mathbb{R}^d)$ -integrable and its Fourier transform $K_{\mathcal{F}}(\boldsymbol{\omega})$ satisfies $K_{\mathcal{F}}(\mathbf{0}) = 1$ and has continuous all partial derivatives $K_{\mathcal{F},i_1,\dots,i_m}^{(m)} = \partial^m K_{\mathcal{F}} / \partial_{i_1} \dots \partial_{i_m}$ up to the ℓ -th order such that $K_{\mathcal{F},i_1,\dots,i_m}^{(m)}(\mathbf{0}) = 0$ for all $m = 1 \dots \ell$.

Remark that the above definition generalizes the notion of the order of kernel as it is defined in [Tsybakov 2009]. The definition imposes the following conditions on an univariate kernel to be of order $\ell \geq 1$, $\ell \in \mathbb{N}$,

- $\int K(u) du = 1$
- $\int u^m K(u) du = 0$ for $m = 1, \dots, \ell$

It can be easily seen that these properties are implied by our definition. Indeed, from the definition of the Fourier transform and assumptions of Definition 3.2 we have $K_{\mathcal{F}}(0) = \int e^{i0u} K(u) du = \int K(u) du = 1$. For the m -th derivatives, $K_{\mathcal{F}}^{(m)}(\boldsymbol{\omega}) = \int (iu)^m e^{i\boldsymbol{\omega}u} K(u) du$ and therefore $K_{\mathcal{F}}^{(m)}(0) = i^m \int u^m K(u) du = 0$.

Theorem 3.2. Let $K: \mathbb{R}^d \rightarrow \mathbb{R}$ be a multivariate kernel of order $\ell \geq 1$, $\ell \in \mathbb{N}$. Then there exists a constant $A \geq 0$ such that

$$\sup_{\boldsymbol{\omega} \in \mathbb{R}^d \setminus \{0\}} \frac{|1 - K_{\mathcal{F}}(\boldsymbol{\omega})|}{\|\boldsymbol{\omega}\|^\ell} \leq A \quad (3.37)$$

and for any density f with the Fourier transform $\phi(\boldsymbol{\omega})$ and $h > 0$,

$$\int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})| |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \leq A^2 h^{2\ell} \int \|\boldsymbol{\omega}\|^{2\ell} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}. \quad (3.38)$$

Proof. We start the proof by the analysis of the univariate case. Let us aim on the limit of $(1 - K_{\mathcal{F}}(\omega))/\omega^\ell$ at $\omega = 0$, i.e., we will investigate $\lim_{\omega \rightarrow 0} (1 - K_{\mathcal{F}}(\omega))/\omega^\ell$. As the kernel is of order ℓ we have zero derivatives at the origin, i.e., $K_{\mathcal{F}}^{(j)}(0) = 0$ for $j = 1, \dots, \ell$. Applying ℓ -times the L'Hospital rule we get for the limit of interest

$$\lim_{\omega \rightarrow 0} \frac{1 - K_{\mathcal{F}}(\omega)}{\omega^\ell} = \frac{-K_{\mathcal{F}}^{(\ell)}(0)}{\ell!} = 0.$$

As for any function $g : \mathbb{R} \rightarrow \mathbb{R}$, $\lim_{x \rightarrow 0} g(x) = 0$ if and only if $\lim_{x \rightarrow 0} |g(x)| = 0$, we have also $(|\omega|^\ell = |\omega^\ell|$ for $\ell \in \mathbb{N}$)

$$\lim_{\omega \rightarrow 0} \frac{|1 - K_{\mathcal{F}}(\omega)|}{|\omega|^\ell} = 0.$$

Having the limit equal to zero, then for any $\delta > 0$ there exists $\omega_0 > 0$ such that

$$\frac{|1 - K_{\mathcal{F}}(\omega)|}{|\omega|^\ell} \leq \max\{\delta, 1/|\omega_0|^\ell\} = A, \quad \omega > 0. \quad (3.39)$$

Indeed, from the definition of the notion of the limit we have for any $\delta > 0$ such $\omega_0 > 0$ that (3.39) holds for $\omega \in (0, \omega_0)$. The numerator in (3.39) is bounded, i.e., $|1 - K_{\mathcal{F}}(\omega)| \leq 1$, and $|\omega|^\ell$ is increasing hence the left-hand-side of (3.39) is less or equal to $1/|\omega_0|^\ell$ on the interval $[\omega_0, \infty)$.

In order to prove the multivariate case, we rely on the multidimensional Taylor's theorem [Brabec and Hruza 1986] p. 118. Under the assumption that the kernel is of order ℓ , we have its Fourier transform to be real because the kernel is symmetric and the Taylor's theorem specifies the values of $K_{\mathcal{F}}(\boldsymbol{\omega})$ as

$$K_{\mathcal{F}}(\boldsymbol{\omega}) = K_{\mathcal{F}}(\mathbf{0}) + \frac{1}{1!} \sum_{i=1}^d K'_{\mathcal{F}}(\mathbf{0})\omega_i + \dots + \frac{1}{\ell!} \sum_{i_1, \dots, i_\ell=1}^d K_{\mathcal{F}, i_1, \dots, i_\ell}^{(\ell)}(\mathbf{0})\omega_{i_1} \dots \omega_{i_\ell} + R_\ell(\boldsymbol{\omega})$$

with the remainder satisfying the limit property $\lim_{\boldsymbol{\omega} \rightarrow \mathbf{0}} R_\ell(\boldsymbol{\omega})/||\boldsymbol{\omega}||^\ell = 0$.

As all partial derivatives equal zero, the remainder writes as $R_\ell(\boldsymbol{\omega}) = K_{\mathcal{F}}(\boldsymbol{\omega}) - K_{\mathcal{F}}(\mathbf{0}) = K_{\mathcal{F}}(\boldsymbol{\omega}) - 1$. Thus the above Taylor's formula gives

$$\lim_{\boldsymbol{\omega} \rightarrow \mathbf{0}} \frac{|1 - K_{\mathcal{F}}(\boldsymbol{\omega})|}{||\boldsymbol{\omega}||^\ell} = \lim_{\boldsymbol{\omega} \rightarrow \mathbf{0}} \frac{K_{\mathcal{F}}(\boldsymbol{\omega}) - 1}{||\boldsymbol{\omega}||^\ell} = 0.$$

Now, the same arguments as in the univariate case, see the discussion concerning formula (3.39), lead to the first assertion of the theorem (with the absolute $|\omega|$ replaced by the norm $||\boldsymbol{\omega}||$).

It is clear that the univariate case is a special case of the multivariate one, but we have used the split into the cases and a slightly longer proof from methodological reasons.

The inequality (3.38) follows from (3.37) by the following chain:

$$\begin{aligned}
\sup_{\boldsymbol{\omega} \in \mathbb{R}^d \setminus \{0\}} \frac{|1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|}{\|h\boldsymbol{\omega}\|^\ell} &\leq A, \\
|1 - K_{\mathcal{F}}(h\boldsymbol{\omega})| &\leq A\|h\boldsymbol{\omega}\|^\ell, \\
|1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 &\leq A^2\|h\boldsymbol{\omega}\|^{2\ell}, \\
|1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 &\leq A^2 h^{2\ell} \|\boldsymbol{\omega}\|^{2\ell} |\phi(\boldsymbol{\omega})|^2, \\
\int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq A^2 h^{2\ell} \int \|\boldsymbol{\omega}\|^{2\ell} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}. \quad \square
\end{aligned}$$

3.4.2 Other terms of Fourier MISE decomposition

The other terms in formula (3.35) refers to individual properties of the kernel and the density under considerations. We mention only two straightforward observations.

The second term can be directly translated from the frequency to the “time” domain by the Plancherel’s theorem and the scaling property of the Fourier transform

$$\frac{1}{n} \int |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} = \frac{(2\pi)^d}{nh^{2d}} \int K^2(\boldsymbol{x}/h) d\boldsymbol{x} = \frac{(2\pi)^d}{nh^d} \int K^2(\boldsymbol{u}) d\boldsymbol{u}. \quad (3.40)$$

For the correction term we have the following inequality

$$\begin{aligned}
\frac{1}{(2\pi)^d} \frac{1}{n} \int |\phi(\boldsymbol{\omega})|^2 |K(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq \frac{\|K_{\mathcal{F}}\|_\infty^2}{(2\pi)^d n} \int |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
&\leq \frac{\|K_{\mathcal{F}}\|_\infty^2}{n} \int f^2(\boldsymbol{x}) d\boldsymbol{x}. \quad (3.41)
\end{aligned}$$

3.4.3 Upper bound on the Fourier MISE formula

Concerning an upper bound on the Fourier MISE decomposition formula (3.35) we actually sum up the results obtained in the preceding sections. First of all, to obtain the upper bound we can omit the correction (the third) term in (3.35) formula. The second term is solely determined by the properties of the

kernel, which is expressed by formula (3.40). Finally, to obtain the limit on the first term, the properties of the density the data are sampled from and the properties of the kernel have to be matched somehow. The matching is provided by fitting the order of the employed kernel with the Sobolev character of the estimated density. The next theorem provides the final result.

Theorem 3.3. *Let n be the number of i.i.d. samples from a distribution with a density $f : \mathbb{R}^d \rightarrow [0, \infty)$ which is β -Sobolev for some $\beta \in \mathbb{N}$ and $L > 0$, i.e., $f \in \mathcal{P}_S(\beta, L)$. Let K be a symmetric, $L_2(\mathbb{R}^d)$ -integrable kernel of order β . Assume that inequality (3.37) holds for some constant $A > 0$. Fix $\alpha > 0$ and take $h = \alpha n^{-\frac{1}{2\beta+d}}$ where $n \in \mathbb{N}$ is the number of samples. Then for any $n \geq 1$ the kernel estimator \hat{f}_n satisfies*

$$\sup_{f \in \mathcal{P}_S(\beta, L)} \mathbb{E} \int (\hat{f}_n(\mathbf{x}) - f(\mathbf{x}))^2 d\mathbf{x} \leq C \cdot n^{-\frac{2\beta}{2\beta+d}} \quad (3.42)$$

where $C > 0$ is a constant depending only on α, d, A, L and the kernel K .

Proof. By the preceding lemma and from the definition of the Sobolev class of densities we have

$$\int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \leq A^2 h^{2\beta} \int \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \leq (2\pi)^d A^2 h^{2\beta} L^2.$$

Plugging this into the Fourier MISE decomposition formula (3.35) and employing $\frac{1}{(2\pi)^d n} \int |K(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} = \frac{1}{nh^d} \int K^2(\mathbf{u}) d\mathbf{u}$ we get for $h = \alpha n^{-\frac{1}{2\beta+d}}$,

$$h^{2\beta} = \alpha^{2\beta} n^{-\frac{2\beta}{2\beta+d}}, \quad (nh^d)^{-1} = n^{-1} \alpha^{-d} n^{\frac{d}{2\beta+d}} = \alpha^{-d} n^{-\frac{2\beta}{2\beta+d}}$$

and

$$\begin{aligned} \text{MISE} &\leq \frac{1}{(2\pi)^d} \left[\int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} + \frac{1}{n} \int |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \right] \\ &\leq A^2 h^{2\beta} L^2 + \frac{1}{nh^d} \int K^2(\mathbf{u}) d\mathbf{u}, \\ &\leq (AL)^2 \alpha^{2\beta} n^{-\frac{2\beta}{2\beta+d}} + \alpha^{-d} n^{-\frac{2\beta}{2\beta+d}} \int K^2(\mathbf{u}) d\mathbf{u}, \\ &\leq \left[(AL)^2 \alpha^{2\beta} + \alpha^{-d} \int K^2(\mathbf{u}) d\mathbf{u} \right] \cdot n^{-\frac{2\beta}{2\beta+d}}, \\ &\leq C(\alpha, d, A, L, K) \cdot n^{-\frac{2\beta}{2\beta+d}}. \quad \square \end{aligned}$$

The theorem says that we are able to state an upper bound on the MISE of a kernel density estimator if we match the order of the employed kernel with the Sobolev character of the density the analyzed data are sampled from.

4. SMC and kernel methods

This chapter presents our own research in the area of the combination of SMC and kernel methods. The basic idea of our approach is the following. A SMC particle filter generates an approximation of a filtering distribution in the form of an empirical measure. This representation enables an effective computation of integral characteristics of the filtering distribution. However, we can also be interested in other characteristics which are commonly computed from the analytical representation of a density of a distribution. The canonical example is the computation of conditional expectations. Hence, our idea is to compute sequentially *kernel density estimates of filtering distributions*. The estimates are based on the empirical distributions generated by a SMC particle filter.

Our main results are threefold. First, we show that the kernel density estimates formed on the basis of empirical distributions converge to the true densities of filtering distributions. We identify the rate of the convergence and show how the number of particles and the bandwidth have to be selected in order to the MISE of the approximation can be controlled. The proof of the result is based on the Fourier analysis of the convergence result for SMC particle filters.

The other results concern a deeper analysis of the obtained convergence formula. First, the convergence result is based on the assumption on the Sobolev character of involved distributions. We present a result addressing this issue. Second, the convergence formula contains factors which are constant with respect to the number of sampled particles, but evolve with time on the basis of the values of observation process. We discuss the character of this evolution and present exact and approximate lower bounds related to these factors.

4.1 Convergence of SMC kernel estimates

In the theorem below, we prove that at each time instant, the kernel density estimates created on the basis of the operation of a SMC particle filter, converge to the density of the related filtering distribution. The convergence is provided in terms of decreasing MISE when the number of employed particles goes to infinity. The rate of the convergence is further linked to the selected bandwidth, time-evolving character of densities of underlying distributions and to the properties of the employed kernel.

To present the theorem, let us stress that the result is primarily related to kernel density estimates of **marginal filtering distributions**.

The SMC particle filter introduced in Chapter 2, generates empirical distributions on $(\mathbb{R}^{n_x(t+1)}, \mathcal{B}(\mathbb{R}^{n_x(t+1)}))$ spaces which approximate filtering distributions over time. The filtering distribution is the joint conditional distribution of states $X_{0:t}$ conditioned by an observed history $Y_{1:t}$. The most widely used integral characteristic of this joint filtering distribution is the conditional expected value of the actual state X_t conditioned by observations $Y_{1:t}$, i.e., $\mathbb{E}[X_t|Y_{1:t}]$. This integral characteristic is in fact the characteristic of the conditional distribution of X_t conditioned by $Y_{1:t}$. That is, it is the integral characteristic of $\pi_{t|t}$ distribution if we employ the (2.29) notation of Chapter 2. The $\pi_{t|t}$ distribution is the marginal distribution of the joint filtering distribution $\pi_t = \pi_{0:t|t}$ at time t , and it is the probability measure specified on $(\mathbb{R}^{n_x}, \mathcal{B}(\mathbb{R}^{n_x}))$ space for each $t \in \mathbb{N}_0$. The analogous considerations hold for the empirical counterparts $\pi_{t|t}^n$ and π_t^n of filtering distributions, which are produced by a SMC particle filter.

Working with joint distributions π_t^n and π_t instead of marginal ones $\pi_{t|t}^n$ and $\pi_{t|t}$ forces the kernel density estimates to be based on kernels of increasing dimension. The related explanation and notation becomes more complicated without any substantial gain in the presented theory. Moreover, the marginal distributions are practically the most desired and employed. The restriction to the marginal distributions is in fact not a serious simplification because it does not harm the employed proof techniques. Thus, relating the presented results primarily to the marginal distributions is a reasonable and practical decision.

To present the theorem, remark that filtering distributions are specified for $t \in \mathbb{N}_0$ where $\pi_{0|0} = \pi_0$ is the initial distribution of the state process. The empirical distributions are generated for $t \in \mathbb{N}$ because the observation process starts from time $t = 1$. Further, in what follows we use d letter to denote the dimension of the state process instead of former n_x . The reason is to maintain consistency with the notation of the preceding chapter.

Theorem 4.1. *Let $\{\pi_{t|t}^n, t \in \mathbb{N}\}$ be the sequence of marginal empirical measures generated by a SMC particle filter on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, $d \in \mathbb{N}$ space. Let $p_t^n: \mathbb{R}^d \rightarrow \mathbb{R}$ be the sequence of related kernel density estimates with the bandwidth varying as $h(n) = \alpha n^{-\frac{1}{2\beta+d}}$ for some $\alpha > 0$. Let $\{\pi_{t|t}, t \in \mathbb{N}_0\}$ be the sequence of corresponding marginal filtering distributions with densities $p_t: \mathbb{R}^d \rightarrow \mathbb{R}$. Let $p_t \in \mathcal{P}_{S(\beta, L_t)}$ for some $\beta \in \mathbb{N}$ and $L_t > 0$, $t = 0, 1, \dots$, and the employed kernel be of order β . Then we have the following evolution of MISE of kernel density*

estimates over time $t \in \mathbb{N}$,

$$\mathbb{E} \left[\int (p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x}_t \right] \leq C_t^2 \cdot n^{-\frac{2\beta}{2\beta+d}} \quad (4.1)$$

where

$$C_t = (AL_t\alpha^\beta + c_t\alpha^{-d/2}\|K\|) \quad (4.2)$$

with the positive constant A of Theorem 3.2 and the sequence c_t following the formula (2.43) of Theorem 2.4.

Proof. The proof is based on the employment of the Fourier transform. We start by the assertion of Theorem 2.4,

$$\mathbb{E}[|\pi_t^n f - \pi_t f|^2] \leq \frac{c_t^2 \|f\|_\infty^2}{n} \quad (4.3)$$

where we replace a general function f by the complex exponential specified on \mathbb{R}^d .

Let $f(x_{0:t}) = f(\mathbf{x}_t) = e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle}$, then $\pi_t^n f = \pi_{t|t}^n f$, $\pi_t f = \pi_{t|t} f$ and $\|f\|_\infty = 1$. Denoting $\psi_t^n = \mathcal{F}[\pi_{t|t}^n]$ and $\psi_t = \mathcal{F}[\pi_{t|t}]$ we have from the above

$$\begin{aligned} \mathbb{E}[|\psi_t^n(\boldsymbol{\omega}) - \psi_t(\boldsymbol{\omega})|^2] &\leq \frac{c_t^2}{n}, \quad (4.4) \\ |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 \cdot \mathbb{E}[|\psi_t^n(\boldsymbol{\omega}) - \psi_t(\boldsymbol{\omega})|^2] &\leq |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 \cdot \frac{c_t^2}{n}, \\ \mathbb{E}[|\psi_t^n(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega}) - \psi_t(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega})|^2] &\leq |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 \cdot \frac{c_t^2}{n}, \\ \mathbb{E} \left[\int |\psi_t^n(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega}) - \psi_t(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \right] &\leq \frac{c_t^2}{n} \int |K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\ \mathbb{E} \left[\int (p_t^n(\mathbf{x}_t) - p_t^*(\mathbf{x}_t))^2 d\mathbf{x}_t \right] &\leq \frac{c_t^2}{nh^d} \int K^2(\mathbf{u}) d\mathbf{u}. \quad (4.5) \end{aligned}$$

For any density p_t and its convolved version $p_t^* = p_t * h^{-d}K(\cdot/h)$,

$$\begin{aligned} \int (p_t^*(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x} &= \frac{1}{(2\pi)^d} \int |\psi_t(\boldsymbol{\omega})K_{\mathcal{F}}(h\boldsymbol{\omega}) - \psi_t(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\ &= \frac{1}{(2\pi)^d} \int |1 - K_{\mathcal{F}}(h\boldsymbol{\omega})|^2 |\psi_t(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}. \quad (4.6) \end{aligned}$$

As we assume that $p_t \in \mathcal{P}_{S(\beta, L_t)}$ and the employed kernel has order β , then according to Theorem 3.2 the right-hand-side of (4.6) is bounded and because there is nothing random we can apply the expectation with no effect to obtain

$$\mathbb{E} \left[\int (p_t^*(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x}_t \right] \leq A^2 h^{2\beta} L_t^2. \quad (4.7)$$

To proceed, let us consider the product measure $\lambda^d \otimes P$ with the corresponding norm $\|\cdot\|_{2, \lambda^d \otimes P} = [\int \int |\cdot|^2 d(\lambda^d \otimes P)]^{1/2}$, then we have

$$\|p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t)\|_{2, \lambda^d \otimes P} \leq Ah^\beta L_t + \frac{c_t}{(nh^d)^{1/2}} \|K\| \quad (4.8)$$

on the basis of (4.5), (4.6) and the triangle inequality for $\|\cdot\|_{2, \lambda^d \otimes P}$.

Let the bandwidth h develop with n as $h(n) = \alpha n^{-\frac{1}{2\beta+d}}$ for some $\alpha > 0$. We have $h^\beta = \alpha^\beta n^{-\frac{\beta}{2\beta+d}}$. Further,

$$(nh^d)^{-1} = n^{-1} \alpha^{-d} n^{\frac{d}{2\beta+d}} = \alpha^{-d} n^{-\frac{2\beta}{2\beta+d}} \quad \text{and therefore} \quad (nh^d)^{-1/2} = \alpha^{-d/2} n^{-\frac{\beta}{2\beta+d}}.$$

Inequality (4.8) then reads as

$$\begin{aligned} \|p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t)\|_{2, \lambda^d \otimes P} &\leq AL_t \alpha^\beta n^{-\frac{\beta}{2\beta+d}} + c_t \alpha^{-d/2} n^{-\frac{\beta}{2\beta+d}} \|K\| \\ &\leq (AL_t \alpha^\beta + c_t \alpha^{-d/2} \|K\|) \cdot n^{-\frac{\beta}{2\beta+d}}. \end{aligned}$$

Squaring to obtain MISE we get

$$\mathbb{E} \int (p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x}_t \leq (AL_t \alpha^\beta + c_t \alpha^{-d/2} \|K\|)^2 \cdot n^{-\frac{2\beta}{2\beta+d}}$$

or in more compact form

$$\mathbb{E} \int (p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x}_t \leq C_t^2 \cdot n^{-\frac{2\beta}{2\beta+d}}$$

for $C_t = AL_t \alpha^\beta + c_t \alpha^{-d/2} \|K\|$ what we wanted to prove. \square

Let us discuss the theorem.

0) First of all, the theorem is proved **without assumption on i.i.d. character of samples** constituting the empirical measures $\pi_{t|t}^n$. This is the crucial observation, as we know that due to the resampling step the generated samples are not i.i.d. Further, we see that the Fourier analysis of kernel estimates is

superior to the standard AMISE analysis as this is based on the assumption on the i.i.d. character of sampled data.

1) **Convergence.** For $t \in \mathbb{N}_0$ fixed, we immediately see from (4.1) that the MISE of kernel estimates goes to zero if the number of samples (particles) increases and the bandwidth decreases accordingly, i.e., $\lim_{n \rightarrow \infty} \mathbb{E} \int (p_t^n(\mathbf{x}_t) - p_t(\mathbf{x}_t))^2 d\mathbf{x}_t = 0$. Therefore we can generate kernel density estimates of marginal filtering densities with the MISE convergence assured at each time instant $t \in \mathbb{N}$.

2) The dimension matters. For $d_1 < d_2$, $n^{-\frac{2\beta}{2\beta+d_1}} < n^{-\frac{2\beta}{2\beta+d_2}}$, and therefore we have to increase the number of particles in order to satisfy a required accuracy as the dimension increases.

3) For $\alpha = 1$, the specification of C_t simplifies to $C_t = AL_t + c_t||K||$ and C_t is composed of four terms. Two of them, A and $||K|| = [\int K^2(\mathbf{u}) d\mathbf{u}]^{1/2}$ are the constants determined by the employed kernel. The other two, L_t and c_t , develop with time. Let us first discuss the L_t term.

4) The theorem assumes that the true marginal filtering densities p_t are β -Sobolev for some $L_t > 0$, $t \in \mathbb{N}_0$ and β being constant over time. It is the question if this assumption is true. In the next section we show that generally the Sobolev character of a SMC particle filter is retained under certain conditions on the transition kernel of the filter.

5) The other entity developing with time is c_t . Its values are computed recursively according to Theorem 2.4 as $c_t = c_{t-1} \left(1 + \frac{4||g_t||_\infty}{\bar{\pi}_t g_t}\right)$, $c_0 = 1$. Integral $\bar{\pi}_t g_t$ depends on the values of the observation process and it is generally hard to state any reasonable lower bound on it. Section 4.3 discusses this issue.

6) The order β of a kernel is actually optional parameter as there are techniques how to construct kernels of arbitrary orders [Tsybakov 2009; Scott 1992].

7) In the presented proof, if the complex exponential would not be restricted, i.e., if we would use $f(x_{0:t}) = e^{i\langle \omega, x_{0:t} \rangle}$, then the proof remains valid with d replaced by $d(t+1)$, \mathbf{x}_t replaced by $x_{0:t}$ and the kernel employed for density estimates being $d(t+1)$ dimensional - the version of the theorem for the joint filtering distributions.

4.2 Sobolev character of SMC particle filters

In this section we show that if the density p_0 of π_0 is β -Sobolev, i.e., if $p_0 \in \mathcal{P}(\beta, L_0)$, $\beta \in \mathbb{N}$, $L_0 > 0$, then also densities p_t are β -Sobolev and we compute the explicit expressions for L_t . We show the result for time homogeneous filter, when the transition kernels K_t , observation densities g_t and also modification functions h_t do not change with time. Let us start with the lemma showing that if a density is β -Sobolev then is L_2 -integrable.

Lemma 4.1. *Let a multivariate density $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be β -Sobolev for $d, \beta \in \mathbb{N}$, i.e., $f \in \mathcal{P}_{S(\beta, L)}$. Then $\int (f(\mathbf{x}))^2 d\mathbf{x} \leq L^2 + (2\pi)^{-d} V_d$, where V_d is the volume of the unit ball B_d in \mathbb{R}^d , $B_d = \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| \leq 1\}$.*

Proof. To start remind that the Fourier transform $\phi(\boldsymbol{\omega})$ of a density f is a characteristic function, i.e., $|\phi(\boldsymbol{\omega})| \leq 1$ for any $\boldsymbol{\omega} \in \mathbb{R}^d$.

From the definition of the Sobolev class of densities (3.36) we get

$$\begin{aligned}
 (2\pi)^d L^2 &\geq \int \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} = \int_{\|\boldsymbol{\omega}\| \leq 1} \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2 + \int_{\|\boldsymbol{\omega}\| > 1} \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2, \\
 (2\pi)^d L^2 &\geq \int_{\|\boldsymbol{\omega}\| > 1} \|\boldsymbol{\omega}\|^{2\beta} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} \\
 &\geq \int_{\|\boldsymbol{\omega}\| > 1} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\
 (2\pi)^d L^2 + \int_{\|\boldsymbol{\omega}\| \leq 1} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\geq \int_{\|\boldsymbol{\omega}\| \leq 1} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} + \int_{\|\boldsymbol{\omega}\| > 1} |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\
 (2\pi)^d L^2 + \int_{\|\boldsymbol{\omega}\| \leq 1} 1 d\boldsymbol{\omega} &\geq \int |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\
 (2\pi)^d L^2 + V_d &\geq \int |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\
 L^2 + (2\pi)^{-d} V_d &\geq (2\pi)^{-d} \int |\phi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega}, \\
 L^2 + (2\pi)^{-d} V_d &\geq \int (f(\mathbf{x}))^2 d\mathbf{x}. \quad \square
 \end{aligned}$$

We proceed with the marginal prediction and update formulas (2.32) and (2.33) of Section 2.6.5. We use the bold notation to be consistent with the actual

notation and omit t subscript in K and g to reflect time homogeneity.

$$\begin{aligned} \text{marg. prediction : } p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) &= \int K(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1}, \\ \text{marg. update : } p(\mathbf{x}_t|\mathbf{y}_{1:t}) &= \frac{g(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t-1})}{\int g(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t}. \end{aligned}$$

We rewrite the formulas in more compact form

$$\text{marg. prediction : } \bar{p}_t(\mathbf{x}_t) = \int K(\mathbf{x}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \quad (4.9)$$

$$\text{marg. update : } p_t(\mathbf{x}_t) = \frac{g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)}{\bar{\pi}_{t|t}g_t} = \frac{g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)}{\bar{\pi}_t g_t} \quad (4.10)$$

using shortcuts $\bar{p}_t(\mathbf{x}_t) = p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$, $p_t(\mathbf{x}_t) = p(\mathbf{x}_t|\mathbf{y}_{1:t})$, $g_t(\mathbf{x}_t) = g(\mathbf{y}_t|\mathbf{x}_t)$ and $\bar{\pi}_{t|t}g_t = \int g(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t = \int g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t$.

Integral $\bar{\pi}_{t|t}g_t$ coincides with the integral $\bar{\pi}_t g_t = \int g(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_{0:t}|\mathbf{y}_{1:t-1}) d\mathbf{x}_{0:t}$ of the constant c_t of the last theorem of Chapter 2. The reason is that the function $g(\mathbf{y}_t|\mathbf{x}_t) = g(\mathbf{y}_t - h(\mathbf{x}_t))$ is the function only of actual state \mathbf{x}_t and not of whole history $\mathbf{x}_{0:t}$. Due to this coincidence and for the purposes of the next section we use instead of symbol $\bar{\pi}_{t|t}g_t$ the symbol $\bar{\pi}_t g_t$.

Definition 4.1. Let $K(\mathbf{x}_t|\mathbf{x}_{t-1})$ be the transition kernel of a SMC filter. As the conditional characteristic function $K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})$ of the transition kernel K we denote the characteristic function of the conditional distribution determined by the kernel, i.e.,

$$K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1}) = \int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} K(\mathbf{x}_t|\mathbf{x}_{t-1}) d\mathbf{x}_t.$$

Theorem 4.2. For a SMC particle filter, let $p_0 \in \mathcal{P}_{S(\beta, L_0)}$. Let the transition kernel of the filter be K and its conditional characteristic function $K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})$ be either (i) bounded by some β -Sobolev function $K_b: \mathbb{R}^d \rightarrow \mathbb{C}$, $K_b \in \mathcal{P}_{S(\beta, L_{K_b})}$ in such a way that for any $\mathbf{x}_{t-1} \in \mathbb{R}^d$ and $\boldsymbol{\omega} \in \mathbb{R}^d$,

$$|K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})| \leq |K_b(\boldsymbol{\omega})|, \quad (4.11)$$

or (ii) $L_2(\mathbb{R}^d)$ -integrable for some $\mu K > 0$ in the following sense

$$\iint \|\boldsymbol{\omega}\|^{2\beta} |K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})|^2 d\mathbf{x}_{t-1} d\boldsymbol{\omega} \leq (2\pi)^d \mu K^2. \quad (4.12)$$

Then the densities p_t of $\pi_{t|t}$ are β -Sobolev for all $t \in \mathbb{N}_0$, i.e., $p_t \in \mathcal{P}_{S(\beta, L_t)}$, with the recurrence for $L_t, t \in \mathbb{N}_0$ written for the case (i) or (ii) as

$$(i) \quad L_t = \frac{g_{\max} L_{K_b}}{\bar{\pi}_t g_t}, \quad \text{or} \quad (ii) \quad L_t = \frac{g_{\max} \mu K \sqrt{L_{t-1}^2 + (2\pi)^{-d} V_d}}{\bar{\pi}_t g_t}, \quad (4.13)$$

respectively.

Proof. We prove the theorem by induction over $t \in \mathbb{N}_0$. The proof has at certain spot two branches according to if there holds either the property (i) or (ii). We assume that the theorem holds for p_0 . Let $p_{t-1} \in \mathcal{P}_{(\beta, L_{t-1})}$, $t \geq 1$, then $\int (p_{t-1}(\mathbf{x}))^2 d\mathbf{x} \leq L_{t-1}^2 + (2\pi)^{-d} V_d$ by Lemma 4.1. Employing the prediction formula, we get $(\bar{\psi}_t(\boldsymbol{\omega}))$ is the Fourier transform of $\bar{p}_t(\mathbf{x}_t)$,

$$\begin{aligned} \bar{p}_t(\mathbf{x}_t) &= \int K(\mathbf{x}_t | \mathbf{x}_{t-1}) p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \\ e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} \bar{p}_t(\mathbf{x}_t) &= e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} \int K(\mathbf{x}_t | \mathbf{x}_{t-1}) p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \\ \int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} \bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t &= \int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} \int K(\mathbf{x}_t | \mathbf{x}_{t-1}) p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} d\mathbf{x}_t \\ &= \int p_{t-1}(\mathbf{x}_{t-1}) \left(\int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} K(\mathbf{x}_t | \mathbf{x}_{t-1}) d\mathbf{x}_t \right) d\mathbf{x}_{t-1}, \\ |\bar{\psi}_t(\boldsymbol{\omega})|^2 &= \left| \int p_{t-1}(\mathbf{x}_{t-1}) K_{\mathcal{F}}(\boldsymbol{\omega} | \mathbf{x}_{t-1}) d\mathbf{x}_{t-1} \right|^2. \end{aligned} \quad (4.14)$$

Now, let us assume that the property (i) of $K_{\mathcal{F}}(\boldsymbol{\omega} | \mathbf{x}_{t-1})$ holds. We have

$$\begin{aligned} |\bar{\psi}_t(\boldsymbol{\omega})|^2 &\leq \left(\int p_{t-1}(\mathbf{x}_{t-1}) |K_{\mathcal{F}}(\boldsymbol{\omega} | \mathbf{x}_{t-1})| d\mathbf{x}_{t-1} \right)^2 \\ &\leq \left(|K_b(\boldsymbol{\omega})| \int p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} \right)^2 = |K_b(\boldsymbol{\omega})|^2, \\ \int \|\boldsymbol{\omega}\|^{2\beta} |\bar{\psi}_t(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq \int \|\boldsymbol{\omega}\|^{2\beta} |K_b(\boldsymbol{\omega})|^2 \leq (2\pi)^d L_{K_b}^2. \end{aligned} \quad (4.15)$$

If the property (ii) of $K_{\mathcal{F}}(\boldsymbol{\omega} | \mathbf{x}_{t-1})$ holds, then we employ the Cauchy-Schwarz

inequality in (4.14) and Lemma 4.1. We have

$$\begin{aligned}
|\bar{\psi}_t(\boldsymbol{\omega})|^2 &\leq \int (p_{t-1}(\mathbf{x}_{t-1}))^2 d\mathbf{x}_{t-1} \cdot \int |K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})|^2 d\mathbf{x}_{t-1}, \\
\int \|\boldsymbol{\omega}\|^{2\beta} |\bar{\psi}_t(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq \int (p_{t-1}(\mathbf{x}_{t-1}))^2 d\mathbf{x}_{t-1} \cdot \\
&\quad \cdot \int \int \|\boldsymbol{\omega}\|^{2\beta} |K_{\mathcal{F}}(\boldsymbol{\omega}|\mathbf{x}_{t-1})|^2 d\mathbf{x}_{t-1} d\boldsymbol{\omega}, \\
\int \|\boldsymbol{\omega}\|^{2\beta} |\bar{\psi}_t(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq (2\pi)^d (L_{t-1}^2 + (2\pi)^{-d} V_d) \cdot \mu K^2. \tag{4.16}
\end{aligned}$$

The above formulas show that $\bar{p}_t \in \mathcal{P}_{(\beta, \bar{P}_t)}$ with \bar{P}_t specified on the basis of (4.15) or (4.16). In the first case $\bar{P}_t = L_{K_b}$ and \bar{P}_t is constant over time. In the second case $\bar{P}_t = \sqrt{(L_{t-1}^2 + (2\pi)^{-d} V_d)} \cdot \mu K$ and \bar{P}_t depends on the Sobolev constant L_{t-1} of p_{t-1} . Let us proceed to the specification of the Sobolev constant L_t of p_t on the basis of knowledge of \bar{P}_t .

In Section 2.6.2, there was shown that function $g_t(\mathbf{x}_t)$ of update formula (4.10) has form $g_t(\mathbf{x}_t) = g(\mathbf{y}_t - h(\mathbf{x}_t))$. Function g is the density of the noise term of observation process, so it is bounded. Thus, regardless of the form of function h , we have $\sup_{\mathbf{x}_t, \mathbf{y}_t} |g_t|^2 \leq \sup_{\mathbf{u}} |g(\mathbf{u})|^2 = \|g\|_{\infty}^2 = \max_{\mathbf{u}} |g(\mathbf{u})|^2 = g_{\max}^2$.

The update formula then gives

$$\begin{aligned}
(\bar{\pi}_t g_t) p_t(\mathbf{x}_t) &= g_t(\mathbf{x}_t) \bar{p}_t(\mathbf{x}_t), \\
(\bar{\pi}_t g_t) \int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} p_t(\mathbf{x}_t) d\mathbf{x}_t &\leq g_{\max} \int e^{i\langle \boldsymbol{\omega}, \mathbf{x}_t \rangle} \bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t, \\
(\bar{\pi}_t g_t)^2 |\psi(\boldsymbol{\omega})|^2 &\leq g_{\max}^2 |\bar{\psi}(\boldsymbol{\omega})|^2, \\
(\bar{\pi}_t g_t)^2 (2\pi)^{-d} \|\boldsymbol{\omega}\|^{2\beta} |\psi(\boldsymbol{\omega})|^2 &\leq g_{\max}^2 (2\pi)^{-d} \|\boldsymbol{\omega}\|^{2\beta} |\bar{\psi}(\boldsymbol{\omega})|^2, \\
(2\pi)^{-d} \int \|\boldsymbol{\omega}\|^{2\beta} |\psi(\boldsymbol{\omega})|^2 d\boldsymbol{\omega} &\leq \frac{g_{\max}^2 \bar{P}_t^2}{(\bar{\pi}_t g_t)^2} = L_t^2. \quad \square \tag{4.17}
\end{aligned}$$

The theorem tells us that during the evolution of a SMC particle filter with a bounded or L_2 -integrable transition kernel $K(\mathbf{x}_t|\mathbf{x}_{t-1})$, in the sense presented above, the β -Sobolev character of the density p_0 of the initial distribution π_0 is preserved over time.

The preservation of the Sobolev property is determined by the properties of the transition kernel. Since the prediction and update formulas have the identical structure for marginal and joint filtering distributions, the theorem can

be proved in the same fashion for the densities of joint prediction and update measures $\bar{\pi}_t$ and π_t . The only difference is that we finally consider term $(2\pi)^{-d(t+1)}V_{d(t+1)}$ instead of $(2\pi)^{-d}V_d$ in (4.13). Further, the extension to time inhomogeneous filters is straightforward with L_{K_b} , μK and g_{\max} possibly varying, given the properties of K_t and g_t , over time.

Concerning the evolution of L_t constant in case (ii), it exhibits asymptotically ($L_t \gg (2\pi)^{-d}V_d$ as $t \rightarrow \infty$) the geometric growth with the quotient $q = (g_{\max}\mu K)/(\bar{\pi}g)^*$, under the assumptions that $\bar{\pi}_t g_t$ is bounded from below by some $(\bar{\pi}g)^* > 0$ and $q > 1$. Nevertheless, such a reasonable lower bound is generally hard to state as the value of $\bar{\pi}_t g_t$ integral depends on the values of observation process. On the other hand, we see that $\bar{\pi}_t g_t$ entity arises not only in the specification of L_t but also in the specification of c_t of (2.43) and therefore it substantially affects the value of C_t constant of Theorem 4.1. This implies that the accuracy of kernel density estimates depend on $\bar{\pi}_t g_t$ integral and it is worth to study its properties.

4.3 Lower bound on $\bar{\pi}_t g_t$

In Theorem 4.1, we have stated an upper bound on the MISE of kernel density estimates of marginal filtering distributions in a SMC particle filter. The bound depends on the constant¹ C_t which further depends on constants c_t of (2.43) and L_t of (4.17). The explicit formulas write as

$$c_t = c_{t-1} \left(1 + \frac{4\|g_t\|_\infty}{\bar{\pi}_t g_t} \right), c_0 = 1 \quad \text{and} \quad L_t = \frac{g_{\max} \bar{P}_t}{\bar{\pi}_t g_t}. \quad (4.18)$$

The number-of-particles constants depend on the value of integral

$$\bar{\pi}_t g_t = \int g(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) d\mathbf{x}_t = \int g_t(\mathbf{x}_t) \bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t. \quad (4.19)$$

We use this specification of $\bar{\pi}_t g_t$ integral instead of the original one based on the $p(\mathbf{x}_{0:t} | \mathbf{y}_{1:t-1})$ density because the investigation of the properties of the integral is more convenient. Especially, the dimension of the involved prediction density $\bar{p}_t(\mathbf{x}_t)$ is fixed to d for all $t \geq 1$ in (4.19) irrespective of time.

As in both cases the integral occurs in the denominator, we have to work with its lower bound in order to state some upper bounds on the c_t and L_t

¹We are taking about C_t , c_t and L_t as about constants because they are constant with respect to a number of particles, but we still have in mind that they generally evolve with time, which is indicated by the subscript.

constants (and consequently on the related errors). In this section we state such a lower bound for each time instant $t \in \mathbb{N}$ in two variants: *the exact bound* which is unfortunately inapplicable for higher dimensions and times because it goes rapidly to zero, and *approximate bounds* based on the application of the strong law of large numbers (SLLN). The specification of lower bounds is further based on the assumption of the Lipschitz continuity of the theoretical marginal prediction and update densities related to the filter.

4.3.1 Exact bound

Before we start we stress that in this section lower bounds are generally denoted using the star notation, i.e., for some real entity $b \in \mathbb{R}$ we have $b^* \leq b$. Further, we use interchangeably “lower bound of b ” and “lower bound on b ”.

The main idea of our approach to lay down a lower bound $(\bar{\pi}_t g_t)^*$ on $\bar{\pi}_t g_t$, i.e., $(\bar{\pi}_t g_t)^* \leq \bar{\pi}_t g_t$, is based on the fact that in (4.19) the integrand is a non-negative function as both g_t and \bar{p}_t are densities. Due to the non-negativity, any integral on a subset of whole space forms a lower bound of the original integral. Moreover, having specified the value of the integrand at some point $\mathbf{a} \in \mathbb{R}^d$ and assuming its Lipschitz continuity we can explicitly compute a lower bound by integration over any ball $\|\mathbf{x} - \mathbf{a}\| \leq r$ with diameter $r > 0$. This estimate can be further optimized over r in order to be as maximal as possible. The next lemma presents the idea in a more accurate way.

Lemma 4.2. *Let a non-negative function $f: \mathbb{R}^d \rightarrow [0, \infty)$ be Lipschitz continuous on \mathbb{R}^d with constant $f_\alpha > 0$ in the following sense² $|f(\mathbf{x}) - f(\mathbf{y})| \leq f_\alpha \|\mathbf{x} - \mathbf{y}\|$ for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Then for any point $\mathbf{a} \in \mathbb{R}^d$,*

$$\int f(\mathbf{x}) d\mathbf{x} \geq \bar{f}(\mathbf{a}) = (f(\mathbf{a}))^{d+1} \frac{D_d V_d}{f_\alpha^d} \geq 0$$

where $D_d = \left(\frac{1}{d+1}\right) \left(\frac{d}{d+1}\right)^d$, $V_d = \frac{\pi^{d/2}}{\Gamma(d/2+1)}$ is the volume of the unit ball in \mathbb{R}^d and

$$\bar{f}(\mathbf{a}) = \max_{r>0} \int_{\|\mathbf{x}-\mathbf{a}\|\leq r} f(\mathbf{x}) d\mathbf{x}.$$

²Here we use for the definition of Lipschitz continuity the Euclidean norm instead of usual L_1 norm. Clearly, due to the equivalence of norms on \mathbb{R}^d this does not makes any problems w.r.t. to the usual definition.

Proof: First of all, as f is non-negative we have for any $\mathbf{a} \in \mathbb{R}^d$ and $r > 0$,

$$\int f(\mathbf{x}) d\mathbf{x} \geq \int_{\|\mathbf{x}-\mathbf{a}\| \leq r} f(\mathbf{x}) d\mathbf{x}.$$

Further, as f is Lipschitz continuous we have also

$$\begin{aligned} |f(\mathbf{x}) - f(\mathbf{a})| &\leq f_\alpha \|\mathbf{x} - \mathbf{a}\|, \\ f(\mathbf{a}) - f(\mathbf{x}) &\leq f_\alpha \|\mathbf{x} - \mathbf{a}\|, \\ f(\mathbf{x}) &\geq f(\mathbf{a}) - f_\alpha \|\mathbf{x} - \mathbf{a}\|, \\ \min_{\{\mathbf{x}: \|\mathbf{x}-\mathbf{a}\| \leq r\}} \{f(\mathbf{x})\} &\geq f(\mathbf{a}) - f_\alpha r. \end{aligned}$$

This implies

$$\int_{\|\mathbf{x}-\mathbf{a}\| \leq r} f(\mathbf{x}) d\mathbf{x} \geq (f(\mathbf{a}) - f_\alpha r) \int_{\|\mathbf{x}-\mathbf{a}\| \leq r} 1 d\mathbf{x} = (f(\mathbf{a}) - f_\alpha r) V_d(r).$$

We have $V_d(r) = V_d(1)r^d = V_d r^d$ for the volume of a d -dimensional ball with the radius $r > 0$. Let us search for the maximum of $b(r) = r^d(f(\mathbf{a}) - f_\alpha r)V_d$ with respect to $r > 0$. Setting $b'(r)$ equal to zero we get

$$\begin{aligned} dr^{d-1}(f(\mathbf{a}) - f_\alpha r)V_d + r^d(-f_\alpha)V_d &= 0, \\ d(f(\mathbf{a}) - f_\alpha r) + r(-f_\alpha) &= 0, \\ -(d+1)f_\alpha r &= -d f(\mathbf{a}), \\ r^* &= (d f(\mathbf{a})) / ((d+1)f_\alpha). \end{aligned}$$

The solution of $b'(r) = 0$ is denoted r^* , i.e., $b'(r^*) = 0$. As $b''(r^*) < 0$, r^* is the point of the local and also global maxima. The explicit formula for $b(r^*)$ reads as

$$b(r^*) = (f(\mathbf{a}))^{d+1} \left(\frac{1}{d+1} \right) \left(\frac{d}{d+1} \right)^d \frac{V_d}{f_\alpha^d} = (f(\mathbf{a}))^{d+1} \frac{D_d V_d}{f_\alpha^d}. \quad (4.20)$$

Clearly, the expression (4.20) is always non-negative. \square

Let $(g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t))^*$ be a lower bound on the value of integrand $g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)$ in (4.19) at some selected point $\mathbf{x}_t = \mathbf{a}_t \in \mathbb{R}^d$, i.e., $(g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t))^* \leq g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t)$. The reason why we generally use a lower bound is that we are not able to compute the exact value of the integrand because we do not know exact form of \bar{p}_t over time. Employing Lemma 4.2 we can immediately state the lower bound on $\bar{\pi}_t g_t$ at $\mathbf{a}_t \in \mathbb{R}^d$ as

$$(\bar{\pi}_t g_t)^* = [(g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t))^*]^{d+1} \frac{D_d V_d}{(g\bar{p}_\alpha^t)^d} \quad (4.21)$$

where $g\bar{p}_\alpha^t$ is the Lipschitz constant of the integrand $g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)$. In order to have (4.21) fully determined we have to

- select some suitable point $\mathbf{a}_t \in \mathbb{R}^d$,
- state the value of the integrand at the selected point, i.e., the value of $g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t)$ or at least its lower bound $(g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t))^*$,
- specify a Lipschitz constant $g\bar{p}_\alpha^t$ of the integrand $g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)$.

We discuss **the first two items** together. The evolution of the prediction density \bar{p}_t is driven by the prediction and update formulas of the employed SMC particle filter. In the compact notation of (4.9) and (4.10) these write as

$$\text{marg. prediction : } \quad \bar{p}_t(\mathbf{x}_t) = \int K(\mathbf{x}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \quad (4.22)$$

$$\text{marg. update : } \quad p_t(\mathbf{x}_t) = \frac{g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)}{\bar{\pi}_t g_t} \quad (4.23)$$

and since $g_t(\mathbf{x}_t) = g(\mathbf{y}_t|\mathbf{x}_t) = g(\mathbf{y}_t - h(\mathbf{x}_t))$,

$$\bar{\pi}_t g_t = \int g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t = \int g(\mathbf{y}_t - h(\mathbf{x}_t))\bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t.$$

Let us proceed sequentially. We assume that at time t we have at our disposal a lower bound $p_{t-1}^*(\mathbf{a}_{t-1})$ of the update density p_{t-1} computed at some point $\mathbf{a}_{t-1} \in \mathbb{R}^d$, i.e., $p_{t-1}^*(\mathbf{a}_{t-1}) \leq p_{t-1}(\mathbf{a}_{t-1})$ for some $\mathbf{a}_{t-1} \in \mathbb{R}^d$.

Now, for any $\mathbf{a} \in \mathbb{R}^d$ and $\bar{p}_t^*(\mathbf{a}) \leq \bar{p}_t(\mathbf{a})$ we have from (4.23) the lower bound

$$p_t(\mathbf{a}) = \frac{g(\mathbf{y}_t - h(\mathbf{a}))\bar{p}_t(\mathbf{a})}{\int g(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t-1}) d\mathbf{x}_t} \geq \frac{g(\mathbf{y}_t - h(\mathbf{a}))\bar{p}_t^*(\mathbf{a})}{g_{\max}} = p_t^*(\mathbf{a}). \quad (4.24)$$

We would like to have $p_t^*(\mathbf{a})$ as large as possible in order to the estimate would be meaningful. In (4.24), we assume that the density g is bounded from above by $g_{\max} = \|g\|_\infty$ and that there exists \mathbf{a}_t such that $g(\mathbf{y}_t - h(\mathbf{a}_t)) = g_{\max}$. That is, let \mathbf{a}_t be driven by \mathbf{y}_t in such a way that $\mathbf{y}_t - h(\mathbf{a}_t)$ localizes the point of the maxima of density g . Under this assumption (4.24) gives for $\mathbf{a} = \mathbf{a}_t$ the value of the lower bound $p_t^*(\mathbf{a}_t)$ as $p_t^*(\mathbf{a}_t) = \bar{p}_t^*(\mathbf{a}_t) \leq \bar{p}_t(\mathbf{a}_t)$ where $\bar{p}_t(\mathbf{a}_t) = \int K(\mathbf{a}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}$. Thus, the lower bounds $\bar{p}_t^*(\mathbf{a}_t)$ and $p_t^*(\mathbf{a}_t)$ coincide.

We compute the value of $\bar{p}_t^*(\mathbf{a}_t)$ under the assumption of the Lipschitz continuity of $K(\mathbf{a}_t|\mathbf{x}_{t-1})$ and $p_{t-1}(\mathbf{x}_{t-1})$ with respect to \mathbf{x}_{t-1} . Let us assume that $K(\mathbf{a}_t|\mathbf{x}_{t-1})$ is Lipschitz continuous with constant $K_{2\alpha}$ irrespective of

the value of \mathbf{a}_t (the Lipschitz continuity of the transition kernel in the second variable), and p_{t-1} with constant p_α^{t-1} . Let $K(\mathbf{a}_t|\mathbf{x}_{t-1})$ be bounded by $K_{2\max} = \|K(\mathbf{a}_t|\mathbf{x}_{t-1})\|_\infty$ irrespective of the value of \mathbf{a}_t (the boundedness of the transition kernel in the second variable), and $p_{t-1}(\mathbf{x}_{t-1})$ by $p_{\max}^{t-1} = \|p^{t-1}\|_\infty$. Then $K(\mathbf{a}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1})$ is Lipschitz continuous with respect to \mathbf{x}_{t-1} with constant

$$Kp_\alpha^t = p_{\max}^{t-1}K_{2\alpha} + K_{2\max}p_\alpha^{t-1}. \quad (4.25)$$

To see this, consider the Lipschitz continuity of the product of two Lipschitz continuous functions and the subtraction-addition trick used in the update formula part of the proof of Lemma 4.3.

The function $K(\mathbf{a}_t|\mathbf{x}_{t-1})p_t(\mathbf{x}_{t-1})$ is non-negative hence we can employ Lemma 4.2 to compute

$$\begin{aligned} \bar{p}_t(\mathbf{a}_t) &= \int K(\mathbf{a}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} \\ &\geq \int_{\|\mathbf{x}_{t-1}-\mathbf{a}_{t-1}\|\leq r} K(\mathbf{a}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} \\ &\geq (K(\mathbf{a}_t|\mathbf{a}_{t-1})p_{t-1}(\mathbf{a}_{t-1}))^{d+1} \frac{D_d V_d}{(Kp_\alpha^t)^d} \\ &\geq (p_{t-1}^*(\mathbf{a}_{t-1}))^{d+1} \frac{(K(\mathbf{a}_t|\mathbf{a}_{t-1}))^{d+1} D_d V_d}{(Kp_\alpha^t)^d} = \bar{p}_t^*(\mathbf{a}_t). \end{aligned} \quad (4.26)$$

Thus (4.26) determines $\bar{p}_t^*(\mathbf{a}_t)$ on the basis of $p_{t-1}^*(\mathbf{a}_{t-1})$ which is assumed to be known from the preceding step. However, as we know that $p_t^*(\mathbf{a}_t)$ and $\bar{p}_t^*(\mathbf{a}_t)$ coincide, we actually need only to focus on computation of $\bar{p}_t^*(\mathbf{a}_t)$ and $p_{t-1}^*(\mathbf{a}_{t-1})$ can be replaced by $\bar{p}_{t-1}^*(\mathbf{a}_{t-1})$ in (4.26) for $t \geq 2$.

The introduced procedure sequentially computes points $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_t, \dots$ from equation $g(\mathbf{y}_t - h(\mathbf{a}_t)) = g_{\max}$ on the basis of observed values $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_t, \dots$, $t \in \mathbb{N}$. For these points it further computes the sequence of lower bounds $\bar{p}_t^*(\mathbf{a}_t)$ by formula (4.26) where we use $p_{t-1}^*(\mathbf{a}_t) = \bar{p}_{t-1}^*(\mathbf{a}_t)$ on the basis of (4.24), i.e., we actually compute only the sequence of $\bar{p}_t^*(\mathbf{a}_t)$ bounds. Concerning the start of the procedure and computation of $\bar{p}_1^*(\mathbf{a}_1)$ we set in (4.26) \mathbf{a}_0 as the point of maxima of p_0 , which is assumed to be known, i.e., $p_0(\mathbf{a}_0) = p_{\max}^0 = \|p_0\|_\infty$ and $p_0^*(\mathbf{a}_0) = p_{\max}^0$; or, if it is possible, we can compute directly $\bar{p}_1^*(\mathbf{a}_1) = \bar{p}_1(\mathbf{a}_1) = \int K(\mathbf{a}_1|\mathbf{x}_0)p_0(\mathbf{x}_0) d\mathbf{x}_0$. In order to the procedure be applicable we assume that the constants of (4.25) are at our disposal. This finishes the discussion of the first item of our list.

Concerning **the second item**, as $g_t(\mathbf{a}_t) = g_{\max}$, it is straightforward to have

$$(g_t(\mathbf{a}_t)\bar{p}_t(\mathbf{a}_t))^* = g_{\max}\bar{p}_t^*(\mathbf{a}_t). \quad (4.27)$$

The third item. To complete the requests of our itemized list, we have to show the Lipschitz continuity of both densities $\bar{p}_t(\mathbf{x}_t) = p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$ and $p_t(\mathbf{x}_t) = p(\mathbf{x}_t|\mathbf{y}_{1:t})$ for any $t \in \mathbb{N}$. The next lemma does the job.

Lemma 4.3. *In the model described by the prediction and update formulas (4.22) and (4.23), let the transition kernel K be Lipschitz continuous and bounded in the first variable irrespective of the value of $\mathbf{x}_{t-1} \in \mathbb{R}^d$ with constants $K_{1\alpha} > 0$ and $K_{1\max} = \|K(\cdot|\mathbf{x}_{t-1})\|_\infty < \infty$. Let the density g be Lipschitz continuous and bounded with constants $g_\alpha > 0$ and $g_{\max} = \|g\|_\infty < \infty$. Let the function h used for specification of $g_t(\mathbf{x}_t) = g(\mathbf{y}_t|\mathbf{x}_t) = g(\mathbf{y}_t - h(\mathbf{x}_t))$ be Lipschitz continuous with constant $h_\alpha > 0$. Then for each $t \in \mathbb{N}$,*

- $\bar{p}_t(\mathbf{x}_t)$ is Lipschitz continuous with constant $\bar{p}_\alpha^t = K_{1\alpha}$,
- $\bar{p}_t(\mathbf{x}_t)$ is bounded from above by $\bar{p}_{\max}^t = K_{1\max}$,
- $p_t(\mathbf{x}_t)$ is Lipschitz with constant $p_\alpha^t = (g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha)/(\bar{\pi}_t g_t)^*$,
- $p_t(\mathbf{x}_t)$ is bounded from above by $p_{\max}^t = (g_{\max}K_{1\alpha})/(\bar{\pi}_t g_t)^*$

Proof: In the case of prediction formula we have for any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$,

$$\begin{aligned}\bar{p}_t(\mathbf{u}) - \bar{p}_t(\mathbf{v}) &= \int (K(\mathbf{u}|\mathbf{x}_{t-1}) - K(\mathbf{v}|\mathbf{x}_{t-1}))p_t(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \\ |\bar{p}_t(\mathbf{u}) - \bar{p}_t(\mathbf{v})| &\leq \int |K(\mathbf{u}|\mathbf{x}_{t-1}) - K(\mathbf{v}|\mathbf{x}_{t-1})|p_t(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1}, \\ |\bar{p}_t(\mathbf{u}) - \bar{p}_t(\mathbf{v})| &\leq K_{1\alpha}\|\mathbf{u} - \mathbf{v}\| \int p_t(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} = K_{1\alpha}\|\mathbf{u} - \mathbf{v}\|.\end{aligned}$$

Thus, the prediction density \bar{p}_t is Lipschitz continuous with the same constant as the transition kernel in the first variable.

The boundedness is straightforward as $\bar{p}_t(\mathbf{x}_t) \leq K_{1\max} \int p_t(\mathbf{x}_{t-1}) d\mathbf{x}_{t-1} = K_{1\max}$.

For the update formula the proof is slightly more difficult. For any $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$,

$$|p_t(\mathbf{u}) - p_t(\mathbf{v})| = \frac{|g_t(\mathbf{u})\bar{p}_t(\mathbf{u}) - g_t(\mathbf{v})\bar{p}_t(\mathbf{v})|}{\bar{\pi}_t g_t}.$$

Let us denote $g_u = g_t(\mathbf{u})$, $p_u = \bar{p}_t(\mathbf{u})$, $g_v = g_t(\mathbf{v})$, $p_v = \bar{p}_t(\mathbf{v})$. The above equation then writes as

$$\begin{aligned}|p_t(\mathbf{u}) - p_t(\mathbf{v})| &= \frac{|g_u p_u - g_v p_v|}{\bar{\pi}_t g_t} = \frac{|g_u p_u - g_u p_v + g_u p_v - g_v p_v|}{\bar{\pi}_t g_t}, \\ |p_t(\mathbf{u}) - p_t(\mathbf{v})| &\leq \frac{g_u |p_u - p_v| + p_v |g_u - g_v|}{\bar{\pi}_t g_t}.\end{aligned}$$

Now, $g_u = g_t(\mathbf{u}) = g(\mathbf{y}_t|\mathbf{u}) = g(\mathbf{y}_t - h(\mathbf{u}))$. By the assumptions of the lemma both functions g and h are Lipschitz continuous with constants g_α and h_α . Therefore $g(\mathbf{y}_t - h(\cdot))$ is also Lipschitz continuous with constant $g_\alpha h_\alpha$ (the Lipschitz continuity of a compound function). As both g and K are bounded, we have $g_u \leq g_{\max}$, $p_u \leq K_{1\max}$, the latter was just proved above, and therefore

$$\begin{aligned} |p_t(\mathbf{u}) - p_t(\mathbf{v})| &\leq \frac{g_{\max}|p_u - p_v| + K_{1\max}|g_u - g_v|}{\bar{\pi}_t g_t}, \\ |p_t(\mathbf{u}) - p_t(\mathbf{v})| &\leq \frac{g_{\max}K_{1\alpha}\|\mathbf{u} - \mathbf{v}\| + K_{1\max}g_\alpha h_\alpha\|\mathbf{u} - \mathbf{v}\|}{\bar{\pi}_t g_t}, \\ |p_t(\mathbf{u}) - p_t(\mathbf{v})| &\leq \frac{g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha}{\bar{\pi}_t g_t} \|\mathbf{u} - \mathbf{v}\|, \\ |p_t(\mathbf{u}) - p_t(\mathbf{v})| &\leq \frac{g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha}{(\bar{\pi}_t g_t)^*} \|\mathbf{u} - \mathbf{v}\|. \end{aligned}$$

Concerning the boundedness, from update formula (4.23) and the boundedness of the prediction density \bar{p}_t we have

$$p_t(\mathbf{x}_t) = \frac{g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)}{\bar{\pi}_t g_t} \leq \frac{g_{\max}K_{1\max}}{(\bar{\pi}_t g_t)^*}. \quad \square$$

We see that under the assumptions on the Lipschitz continuity of the transition kernel in the first variable with constant $K_{1\alpha}$ and its independence of the value of the previous state \mathbf{x}_{t-1} , the Lipschitz constant \bar{p}_α^t of the prediction density $\bar{p}_t = p(\mathbf{x}_t|\mathbf{y}_{1:t-1})$ coincides with $K_{1\alpha}$. The constant is time invariant and independent of observations $\mathbf{y}_1, \dots, \mathbf{y}_{t-1}$, that is why we further use only the symbol \bar{p}_α instead of \bar{p}_α^t .

The Lipschitz constant p_α^t of the update density $p_t(\mathbf{x}_t) = p(\mathbf{x}_t|\mathbf{y}_{1:t})$ is more complicated. The constant is given by a ratio. The numerator is under the assumptions of the lemma independent of time, however, the denominator evolves with time so the Lipschitz constant does it so. In fact, it is indirectly dependent on observations through $(\bar{\pi}_t g_t)^*$. From the proof, it is clear that we could define p_α^t using the exact value of $\bar{\pi}_t g_t$, but as we do not know it, we use a lower bound on $\bar{\pi}_t g_t$ in the definition of p_α^t .

Inspecting the proof of Lemma 4.3, we immediately obtain for the Lipschitz constant of $g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)$,

$$g\bar{p}_\alpha^t = g\bar{p}_\alpha = g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha, \quad (4.28)$$

i.e., the constant does not depend on time and we can remove the time index in the notation.

4.3.2 Sequential computation of $(\bar{\pi}_t g_t)^*$

In the previous section, we have presented the derivation of how to state a lower bound $(\bar{\pi}_t g_t)^*$ on $\bar{\pi}_t g_t$ integrals for $t \geq 1$. In order to use the introduced formulas in an algorithmic way let us summarize the involved entities. We start with constants.

- The Lipschitz constants $K_{1\alpha}$, $K_{2\alpha}$ of the transition kernel in both variables

$$\begin{aligned} |K(\mathbf{u}|\mathbf{x}_{t-1}) - K(\mathbf{v}|\mathbf{x}_{t-1})| &\leq K_{1\alpha} \|\mathbf{u} - \mathbf{v}\| \text{ for any } \mathbf{u}, \mathbf{v}, \mathbf{x}_{t-1} \in \mathbb{R}^d, \\ |K(\mathbf{x}_t|\mathbf{u}) - K(\mathbf{x}_t|\mathbf{v})| &\leq K_{2\alpha} \|\mathbf{u} - \mathbf{v}\| \text{ for any } \mathbf{u}, \mathbf{v}, \mathbf{x}_t \in \mathbb{R}^d. \end{aligned}$$

- The bounding constants $K_{1\max}$, $K_{2\max}$ on the transition kernel in both variables

$$\begin{aligned} K(\mathbf{x}_t|\mathbf{x}_{t-1}) &\leq K_{1\max} = \|K(\cdot|\mathbf{x}_{t-1})\|_\infty < \infty \text{ for any } \mathbf{x}_{t-1} \in \mathbb{R}^d \text{ fixed,} \\ K(\mathbf{x}_t|\mathbf{x}_{t-1}) &\leq K_{2\max} = \|K(\mathbf{x}_t|\cdot)\|_\infty < \infty \text{ for any } \mathbf{x}_t \in \mathbb{R}^d \text{ fixed.} \end{aligned}$$

- The Lipschitz constant g_α of g , i.e., $|g(\mathbf{u}) - g(\mathbf{v})| \leq g_\alpha \|\mathbf{u} - \mathbf{v}\|$ for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$.
- The bounding constant g_{\max} on observation density, $g_{\max} = \|g\|_\infty < \infty$.
- The Lipschitz constant h_α of h , i.e., $|h(\mathbf{u}) - h(\mathbf{v})| \leq h_\alpha \|\mathbf{u} - \mathbf{v}\|$ for $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$.

Further we have the following computed entities.

- The bounding constants on the prediction and update densities

$$\bar{p}_{\max} = K_{1\max}, \quad p_{\max}^t = \frac{g_{\max} K_{1\max}}{(\bar{\pi}_t g_t)^*}. \quad (4.29)$$

- The Lipschitz constants of the prediction and update densities

$$\bar{p}_\alpha = K_{1\alpha}, \quad p_\alpha^t = (g_{\max} K_{1\alpha} + K_{1\max} g_\alpha h_\alpha) / (\bar{\pi}_t g_t)^*. \quad (4.30)$$

- The Lipschitz constant $g\bar{p}_\alpha$ of $g_t(\mathbf{x}_t)\bar{p}_t(\mathbf{x}_t)$

$$g\bar{p}_\alpha = g_{\max} K_{1\alpha} + K_{1\max} g_\alpha h_\alpha. \quad (4.31)$$

- The Lipschitz constant Kp_α^t of $K(\mathbf{a}_t|\mathbf{x}_{t-1})p_{t-1}(\mathbf{x}_{t-1})$

$$Kp_\alpha^t = p_{\max}^{t-1}K_{2\alpha} + K_{2\max}p_\alpha^{t-1}. \quad (4.32)$$

- The auxiliary constant

$$M = (g_{\max}^{d+1}D_dV_d)/(g\bar{p}_\alpha)^d. \quad (4.33)$$

After the summary is presented, we can proceed to the algorithm of the sequential computation of lower bounds $(\bar{\pi}_t g_t)^*$ for $t \geq 1$.

- **0. declarations**

set up $K_{1\alpha}, K_{2\alpha}, K_{1\max}, K_{2\max}, g_\alpha, g_{\max}, h_\alpha$,
 set up - D_d, V_d ,
 set up p_0 - the initial density,
 set up $T \in \mathbb{N}$ - the computation horizon.

- **1. initialization**

$t = 0$,
 set $\bar{p}_{\max} = K_{1\max}, \bar{p}_\alpha = K_{1\alpha}$,
 set $g\bar{p}_\alpha = g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha$,
 set $M = (g_{\max}^{d+1}D_dV_d)/(g\bar{p}_\alpha)^d$,
 set \mathbf{a}_0 as the point of maxima of p_0 ,
 set p_{\max}^0 as the value of maxima p_0 ,
 set p_α^0 as the Lipschitz constant of p_0 ,
 set $\bar{p}_0^*(\mathbf{a}_0) = p_{\max}^0$.

- **2. sequential computation**

$t = t + 1$,
 compute \mathbf{a}_t as the solution of $g(\mathbf{y}_t - h(\mathbf{a}_t)) = g_{\max}$,
 $Kp_\alpha^t = p_{\max}^{t-1}K_{2\alpha} + K_{2\max}p_\alpha^{t-1}$,
 $\bar{p}_t^*(\mathbf{a}_t) = (\bar{p}_{t-1}^*(\mathbf{a}_{t-1}))^{d+1} \frac{(K(\mathbf{a}_t|\mathbf{a}_{t-1}))^{d+1}D_dV_d}{(Kp_\alpha^t)^d}$,
 $(\bar{\pi}_t g_t)^* = (g_{\max}\bar{p}_t^*(\mathbf{a}_t))^{d+1} \frac{D_dV_d}{(g\bar{p}_\alpha)^d} = (\bar{p}_t^*(\mathbf{a}_t))^{d+1}M$,
 set $p_{\max}^t = \frac{g_{\max}K_{1\max}}{(\bar{\pi}_t g_t)^*}$,
 set $p_\alpha^t = \frac{g_{\max}K_{1\alpha} + K_{1\max}g_\alpha h_\alpha}{(\bar{\pi}_t g_t)^*} = \frac{g\bar{p}_\alpha}{(\bar{\pi}_t g_t)^*}$.

- **3. end** if $t = T$ end, else go to step 2

Algorithm 4.1: Computation of exact lower bound on $\bar{\pi}_t g_t$.

The presented procedure of $(\bar{\pi}_t g_t)^*$ specification is mathematically correct in the sense that it is based on the proved assertions of the employed theorems and lemmas. On the other hand, its practical applicability is problematic for higher dimensions and times as the exact lower bound goes rapidly to zero with the dimension and time increasing. This can be seen by an inspection of the recursive formulas in part 2 of the algorithm, when typically $K(\mathbf{a}_t|\mathbf{a}_{t-1})$ is less than one and any increase in dimension d dramatically accelerates the rate of the convergence of $(\bar{\pi}_t g_t)^*$ to zero.

4.3.3 Approximate bounds

In this section we introduce three approximations of $\bar{\pi}_t g_t$ integral or its lower bound. The approximations are based on the empirical distributions generated during the operation of a SMC particle filter.

1. The most straightforward idea is to approximate

$$\bar{\pi}_t g_t \approx \bar{\pi}_t^n g_t = \int g_t(\mathbf{x}_t) d\bar{\pi}_t^n = \frac{1}{n} \sum_{i=1}^n g(\mathbf{y}_t - h(\mathbf{x}_t^i)). \quad (4.34)$$

That is, we employ for the approximation the empirical distribution $\bar{\pi}_t^n$ which is computed during the operation of a filter. The rationality of the approximation is ensured by the convergence of empirical measure $\bar{\pi}_t^n$ to the theoretical measure $\bar{\pi}_t$. Clearly, the quality of the approximation increases with the increasing number n of employed particles.

2. In the second approach we do not approximate $\bar{\pi}_t g_t$ directly, but we approximate the value of the prediction density at point \mathbf{a}_t given as the solution of $g(\mathbf{y}_t - h(\mathbf{a}_t)) = g_{\max}$ equation. The approximation writes as

$$\bar{p}_t(\mathbf{a}_t) \approx \bar{p}_t^n(\mathbf{a}_t) = \pi_{t-1}^n K(\mathbf{a}_t|\mathbf{x}_{t-1}) = \frac{1}{n} \sum_{i=1}^n K(\mathbf{a}_t|\mathbf{x}_{t-1}^i). \quad (4.35)$$

Further, we use Lemma 4.2 and set the lower bound on $\bar{\pi}_t g_t$ according to formula (4.21), i.e.,

$$(\bar{\pi}_t g_t)^{*n} = (g_t(\mathbf{a}_t) \bar{p}_t^n(\mathbf{a}_t))^{d+1} \frac{D_d V_d}{(g \bar{p}_\alpha)^d} = (\bar{p}_t^n(\mathbf{a}_t))^{d+1} \frac{g_{\max}^{d+1} D_d V_d}{(g \bar{p}_\alpha)^d} = (\bar{p}_t^n(\mathbf{a}_t))^{d+1} M. \quad (4.36)$$

The advantage of this approach is that we work only with samples from π_{t-1}^n empirical distribution. It could happen that we would need to increase the

number of employed particles in the filter when going from time $t - 1$ to time t . The natural point in the SMC algorithm where this may be performed is the resampling step; and in fact searching for a suitable number of particles is an iterative procedure including iterative computations of $\bar{\pi}_t g_t$ integral or its lower bound. Employing the first approach we would need to work with both empirical distributions π_{t-1}^n and $\bar{\pi}_t^n$. In the second approach, we work only with the empirical distribution π_{t-1}^n and we do not need iteratively recreate distribution $\bar{\pi}_t^n$ during the search for suitable n . However, the problem of changing the number of particles is out of scope of this thesis and will not be developed here in more details. We only wanted to mention a situation when the second approximative approach is advantageous.

3. The third approach is slightly more elaborated. In the preceding approximations we replaced the theoretical expected values by their empirical estimates. In the first approach we actually state that $\pi_t g_t = \pi_t^n g_t$ and in the second $\bar{p}_t(\mathbf{a}_t) = \bar{p}_t^n(\mathbf{a}_t)$, which is equivalent to the statement of $|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)| = 0$. Let us work further with the second approach and weaken the statement by assuming only that

$$|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)| \leq \epsilon \quad (4.37)$$

for some $\epsilon \geq 0$. We will search for a suitable ϵ_0 and check if the assumption is reasonable for the given number of samples n .

The assumption (4.37) has three consequences:

- $\bar{p}_t^n(\mathbf{a}_t) - \epsilon$ forms the lower bound on $\bar{p}_t(\mathbf{a}_t)$, i.e., $\bar{p}_t(\mathbf{a}_t) \geq (\bar{p}_t^n(\mathbf{a}_t) - \epsilon)$,
- $(\bar{\pi}_t g_t)^{*n}(\epsilon) = (\bar{p}_t^n(\mathbf{a}_t) - \epsilon)^{d+1} M$, with M given by (4.33), gives the lower bound on $\bar{\pi}_t g_t$ by Lemma 4.2,
- finally, from (4.37) we have also

$$\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|] \leq \epsilon. \quad (4.38)$$

In what follows we investigate if (4.38) holds.

Lemma 4.4. *In a SMC particle filter, let $\mathbf{a}_t \in \mathbb{R}^d$, $t \in \mathbb{N}$ be such a sequence of points that $g(\mathbf{y}_t - h(\mathbf{a}_t)) = g_{\max}$ hold all t . Let the transition kernel of the filter K be bounded in the second variable by $K_{2\max} < \infty$. Let $\bar{p}_t(\mathbf{a}_t)$ be computed according to formula (4.35). Let*

$$(\bar{\pi}_t g_t)^{*n}(\epsilon_0) = (\bar{p}_t^n(\mathbf{a}_t) - \epsilon_0)^{d+1} M \quad \text{for} \quad \epsilon_0 = \frac{\bar{p}_t(\mathbf{a}_t)^n}{d+2}. \quad (4.39)$$

Then for each $t \in \mathbb{N}$, there exists such $n_t^0 \in \mathbb{N}$ that $\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|] \leq \epsilon_0$.

Proof. To simplify the notation we abbreviate $\bar{p}_t^n(\mathbf{a}_t)$ to \bar{p}_t^n in the following paragraphs. First of all note, that the lower bound $(\bar{\pi}_t g_t)^{*n}(\epsilon_0)$ specified according to formula (4.39) is reasonable because it is positive. This is given by the fact that $\epsilon_0 < \bar{p}_t^n$ and therefore $\bar{p}_t^n - \epsilon_0 > 0$.

In (4.39), the value of ϵ_0 is specified as the point of maxima of term $\epsilon \cdot (\bar{p}_t^n - \epsilon)^{d+1}$ on interval $[0, \bar{p}_t^n]$. Clearly, the term is non-negative continuous on $[0, \bar{p}_t^n]$ with zero values at endpoints, hence there is some maximum reached. Setting the derivative of the term to zero writes as

$$\begin{aligned} (\bar{p}_t^n - \epsilon)^{d+1} - \epsilon(d+1)(\bar{p}_t^n - \epsilon)^d &= 0, \\ \bar{p}_t^n - \epsilon - \epsilon(d+1) &= 0, \\ \bar{p}_t^n - \epsilon(d+2) &= 0, \\ \epsilon_0 &= \bar{p}_t^n / (d+2). \end{aligned}$$

At the point of maxima we have

$$\left(\bar{p}_t^n - \frac{\bar{p}_t^n}{d+2} \right) = \frac{d\bar{p}_t^n + 2\bar{p}_t^n - \bar{p}_t^n}{d+2} = \bar{p}_t^n \frac{d+1}{d+2}$$

and therefore $\epsilon_0 (\bar{p}_t^n - \epsilon_0)^{d+1}$ reads as

$$\frac{\bar{p}_t^n}{d+2} (\bar{p}_t^n)^{d+1} \left(\frac{d+1}{d+2} \right)^{d+1} = (\bar{p}_t^n)^{d+2} \left(\frac{1}{d+2} \right) \left(\frac{d+1}{d+2} \right)^{d+1} = (\bar{p}_t^n)^{d+2} D_{d+1}.$$

Similarly,

$$(\bar{\pi}_t g_t)^{*n}(\epsilon_0) = M \left(\bar{p}_t^n - \frac{\bar{p}_t^n}{d+2} \right)^{d+1} = M (\bar{p}_t^n)^{d+1} (d+2) D_{d+1}.$$

To proceed note that

$$\epsilon_0 = \frac{M \epsilon_0 (\bar{p}_t^n - \epsilon_0)^{d+1}}{M (\bar{p}_t^n - \epsilon_0)^{d+1}} = \frac{M (\bar{p}_t^n)^{d+2} D_{d+1}}{(\bar{\pi}_t g_t)^{*n}(\epsilon_0)}.$$

Now, let n_t^0 be such a minimal $n \in \mathbb{N}$ that

$$c_{t-1} \frac{5g_{\max}}{(\bar{\pi}_t g_t)^{*n}(\epsilon_0)} \frac{K_{2\max}}{\sqrt{n_t^0}} \leq \epsilon_0 = \frac{M (\bar{p}_t^n)^{d+2} D_{d+1}}{(\bar{\pi}_t g_t)^{*n}(\epsilon_0)}, \quad (4.40)$$

$$c_{t-1}^2 \left[\frac{5g_{\max} K_{2\max}}{M (\bar{p}_t^n)^{d+2} D_{d+1}} \right]^2 \leq n_t^0. \quad (4.41)$$

Such n_t^0 exists because on the left-hand-side of (4.41) we have constants and the random variable $\bar{p}_t^n = \bar{p}_t^n(\mathbf{a}_t)$ converges (stabilizes) a.s. as $n \rightarrow \infty$.

Let us assume that the number of particles n employed in the filter is greater or equal to n_t^0 for all $t \in \mathbb{N}$. Let us show that in this case the assumption (4.37) is consistent with its consequence (4.38) given the operation of the filter.

First of all, $\bar{\pi}_t g_t = \int g_t(\mathbf{x}_t) \bar{p}_t(\mathbf{x}_t) d\mathbf{x}_t \leq g_{\max}$. The SMC filter's convergence formula (2.43) applied to function $f(\mathbf{x}_{t-1}) = K(\mathbf{a}_t | \mathbf{x}_{t-1})$ writes as

$$\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|^2] \leq c_{t-1}^2 \left[1 + \frac{4g_{\max}}{\bar{\pi}_t g_t} \right]^2 \frac{K_{2\max}^2}{n} \leq c_{t-1}^2 \left[\frac{5g_{\max}}{\bar{\pi}_t g_t} \right]^2 \frac{K_{2\max}^2}{n}.$$

Using the Jensen's inequality this gives

$$\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|] \leq c_{t-1} \left[\frac{5g_{\max}}{\bar{\pi}_t g_t} \right] \frac{K_{2\max}}{\sqrt{n}} \leq c_{t-1} \left[\frac{5g_{\max}}{(\bar{\pi}_t g_t)^*} \right] \frac{K_{2\max}}{\sqrt{n}}$$

for any lower bound on $\bar{\pi}_t g_t$. Since $(\bar{\pi}_t g_t)^*(\epsilon_0)$ is the valid lower bound by assumption (4.37) and we assume that n is so large that $n \geq n_t^0$ for all t , the inequality (4.40) applies and we have

$$\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|] \leq c_{t-1} \left[\frac{5g_{\max}}{(\bar{\pi}_t g_t)^*(\epsilon_0)} \right] \frac{K_{2\max}}{\sqrt{n}} \leq \epsilon_0,$$

what we wanted to show. \square

Let us review the presented approach. We assume that $|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)| \leq \epsilon$ for some $\epsilon > 0$. We search for a reasonable value of ϵ which we denote ϵ_0 . The assumption implies that $\bar{p}_t^n(\mathbf{a}_t) - \epsilon$ is a lower bound on $\bar{p}_t(\mathbf{a}_t)$. In order to a lower bound be non-negative it must be $\epsilon_0 < \bar{p}_t^n(\mathbf{a}_t)$. The selection of ϵ_0 determines the minimal number of samples n_t^0 which assures that $\mathbb{E}[|\bar{p}_t^n(\mathbf{a}_t) - \bar{p}_t(\mathbf{a}_t)|] \leq \epsilon_0$. That is, we require the consistency of the assumption with its consequence under the operation of the filter. In fact, this is the criterion which drives the justification of the assumption.

One may ask why ϵ_0 is selected in the presented way. In fact the selection maximizes the denominator in (4.41), which minimizes n_t^0 . Indeed, the denominator writes as $M(\bar{p}_t^n)^{d+2} D_{d+1} = M\epsilon_0 (\bar{p}_t^n - \epsilon_0)^{d+1}$ which maximizes $\epsilon (\bar{p}_t^n - \epsilon)^{d+1} M$ on the interval $[0, \bar{p}_t^n(\mathbf{a}_t)]$.

5. Experiments

Here we present software implementations of the theory developed in the previous chapters together with outputs from the related computer experiments. As our research was not driven by any concrete application, we apply the SMC filter and kernel density estimation methodologies on the filtering problem of a Gaussian process. This problem has the analytical solution - the well know Kalman filter [Kalman 1960; Fristedt et al. 2007; Pollock 1999]. The purpose of this decision is to check if empirical results from computer simulations follow the analytic counterpart. By replacing the Gaussian transition kernel and Gaussian observation density by general entities we can build an appropriate SMC filter for a general Markov process, but without the possibility of checking against the analytical solution. For this reason we chose the Gaussian process as the first to test our results.

5.1 Univariate Gaussian process

Let us start with the univariate case. The signal and observation processes of Section 2.6 are specified as

$$X_t = aX_{t-1} + b + cW_t, \quad Y_t = hX_t + gV_t, \quad t \geq 1. \quad (5.1)$$

In (5.1), a, b, c, h, g are some fixed constants, $c, g > 0$, and $X_0, W_1, V_1, W_2, V_2, \dots$ are univariate real-valued Gaussian random variables with $X_0 \sim \mathcal{N}(\mu_0, \sigma_0^2)$ and $W_1, V_1, W_2, V_2, \dots$ are i.i.d. standard normal, i.e., $W_t, V_t \sim \mathcal{N}(0, 1)$. From these assumptions we see that $\{X_t, t \geq 0\}$ forms a Markov chain with the Gaussian transition kernel and the initial distribution being normal, i.e., $X_0 \sim \mathcal{N}(\mu_0, \sigma_0^2)$.

Watching a stream of observations $Y_s, 1 \leq s \leq t$, the filtering task is to estimate the current value X_t of the signal process in terms of the conditional expected value $\mathbb{E}[X_t|Y_1, \dots, Y_t] = \mathbb{E}[X_t|Y_{1:t}]$. The expected value is one of the integral characteristics of the conditional distribution of X_t conditioned by Y_1, \dots, Y_t . It can be computed analytically and the result is known as the univariate Kalman filter; or it can be stated empirically on the basis of samples from the related conditional distribution provided by the Gaussian SMC particle filter. Before we state the components of the filter let us remind the recursive formulas of the analytical solution.

5.1.1 Univariate Kalman filter

We follow [Fristedt et al. 2007] with the modification that we start the observation process at $t = 1$ in order to be compatible with our framework for particle filters. In [Fristedt et al. 2007] on p. 133, there is shown that the process $Z_t = (X_t, Y_t)$ is Gaussian and the vector (X_t, Y_1, \dots, Y_t) has a density. The density is multivariate normal and therefore the conditional density of X_t conditioned by Y_1, \dots, Y_t is univariate normal for fixed $Y_1 = y_1, \dots, Y_t = y_t$. The task is to find the parameters of this conditional density in terms of mean and variance, i.e., the values of $\hat{\mu}_t = \mathbb{E}[X_t | Y_{1:t}]$ and $\hat{\sigma}_t^2 = \mathbb{E}[(X_t - \hat{\mu}_t)^2 | Y_{1:t}]$.

The solution of the above task has the form of the following formulas for $t \geq 1$,

$$\hat{\mu}_t = \frac{g^2}{v^2}b + \frac{c^2}{v^2}hy_t + \frac{g^2}{v^2}a \frac{\hat{\sigma}_{t-1}^2 ah(y_t - bh) + v^2 \hat{\mu}_{t-1}}{h^2 \hat{\sigma}_{t-1}^2 a^2 + v^2}, \quad (5.2)$$

$$\hat{\sigma}_t^2 = \frac{g^4}{v^2}a^2 \frac{\hat{\sigma}_{t-1}^2}{h^2 \hat{\sigma}_{t-1}^2 a^2 + v^2} + \frac{g^2 c^2}{v^2} \quad (5.3)$$

where $v^2 = c^2 h^2 + g^2$, $\hat{\mu}_0 = \mu_0$ and $\hat{\sigma}_0^2 = \sigma_0^2$.

Using recursively the above so-called *Kalman's equations* we obtain sequentially the integral characteristics, i.e., mean and variance of the distribution of X_t conditioned on what was observed. Note that the formula for variance is deterministic, i.e., it does not depend on observations.

5.1.2 Univariate Gaussian SMC filter

The schema (5.1) can be easily transformed into the SMC particle filter design. First of all, from the specification (5.1) we see that the signal process forms a Markov chain with the time-homogeneous Gaussian transition kernel

$$K(x_t | x_{t-1}) = \frac{1}{c\sqrt{2\pi}} \exp \left[-\frac{(x_t - ax_{t-1} - b)^2}{2c^2} \right]. \quad (5.4)$$

The function h_t of (2.24) is also time-homogeneous, linear and without an absolute term, i.e., $h(u) = hu$. The density of gV_t is the density of $\mathcal{N}(0, g^2)$. Hence function g_t of (2.26) writes for all $t \geq 1$ as

$$g(y_t | x_t) = \frac{1}{g\sqrt{2\pi}} \exp \left[-\frac{(y_t - hx_t)^2}{2g^2} \right]. \quad (5.5)$$

Note that we are exposed here to a mild indiscipline in the notation as we use h and g letters to denote constants of (5.1) and functions of the general SMC

particle filters' framework. However, from the context it will be clear which role is applied.

In order to implement methodology presented in Chapter 4, we have to specify certain constants and check the Lipschitz and Sobolev character of relevant densities. As the transition kernel (5.4) and density (5.5) are Gaussian, the constants and the character of densities relate to properties of the density function f of the general univariate normal distribution $\mathcal{N}(\mu, \sigma^2)$,

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]. \quad (5.6)$$

Let us show that the density of $\mathcal{N}(\mu, \sigma^2)$ is bounded and Lipschitz on \mathbb{R} . The derivative of (5.6) reads as

$$f'(x) = -\frac{(x-\mu)}{\sigma^3\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$

with zero obtained at point $x^* = \mu$. This point is the point of maxima of f and its value $f(\mu) = \frac{1}{\sigma\sqrt{2\pi}}$ does not depend on the value of μ .

To show the Lipschitz continuity we search for an upper bound on $|f'|$. We have

$$f''(x) = \frac{1}{\sigma^3\sqrt{2\pi}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right] \cdot \left[\frac{(x-\mu)^2}{\sigma^2} - 1\right]$$

with zero reached at points $x^* = \mu \pm \sigma$. From the expression for f' we see that the absolute value of extrema is $|f'(x^*)| = \frac{1}{\sigma^2\sqrt{2\pi}} \exp[-\frac{1}{2}]$. As the first derivative is bounded by $|f'(x^*)|$ the density f is Lipschitz with constant $f_\alpha = |f'(x^*)|$ by the mean value theorem.

The specification of transition kernel (5.4), when considered as a function of x_t , coincides with the density function of $\mathcal{N}(\mu = ax_{t-1} + b, \sigma^2 = c^2)$. Hence the transition kernel is bounded and Lipschitz in x_t variable with the related constants specified as $K_{1\max} = \max_{x_t} \{K(x_t|x_{t-1})\} = \frac{1}{c\sqrt{2\pi}}$ and $K_{1\alpha} = \frac{1}{c^2\sqrt{2\pi}} \exp[-\frac{1}{2}]$. Concerning the transition kernel (5.4) as the function of x_{t-1} variable, we employ the fact that $K(x_t|x_{t-1}) = f(ax_{t-1})$ for f being the density of $\mathcal{N}(\mu = x_t - b, \sigma^2 = c^2)$. This immediately gives the boundedness with the same constant as above, i.e., $K_{2\max} = \max_{x_{t-1}} \{K(x_t|x_{t-1})\} = \frac{1}{c\sqrt{2\pi}}$. The Lipschitz constant in the second variable is then given by product $K_{2\alpha} = \frac{|a|}{c^2\sqrt{2\pi}} \exp[-\frac{1}{2}]$ because the Lipschitz constant of function ax_{t-1} is $|a|$.

The observation density $g(\cdot)$ corresponds here to the density of $\mathcal{N}(0, g^2)$ distribution. Therefore it is bounded and Lipschitz with the respective constants of $g_{\max} = \frac{1}{g\sqrt{2\pi}}$ and $g_\alpha = \frac{1}{g^2\sqrt{2\pi}} \exp[-\frac{1}{2}]$. Further, $h_\alpha = |h|$ because $h(x) = hx$.

5.1.3 Sobolev character of univariate filter

Here we specify the values of entities related to the Sobolev character of the filter. The distribution of X_0 is univariate normal, i.e., $X_0 \sim \mathcal{N}(\mu_0, \sigma_0^2)$, $\mu_0 \in \mathbb{R}$, $\sigma_0 > 0$.

For the corresponding density p_0 we have

$$\begin{aligned} \int (p'_0(x))^2 dx &= \int \frac{(x - \mu_0)^2}{\sigma_0^6 2\pi} \exp\left[-\frac{(x - \mu_0)^2}{\sigma_0^2}\right] dx = \frac{\sigma_0}{\sigma_0^4 2\pi} \int z^2 \exp(-z^2) dz \\ &= \frac{1}{\sigma_0^3 4\sqrt{\pi}}. \end{aligned}$$

Therefore p_0 is 1-Sobolev for $L_0 = 1/(\sigma_0^3 4\sqrt{\pi})$ by Lemma 3.4.

In order to present the Sobolev character of the filter we show that for the Gaussian transition kernel there exists 1-Sobolev bounding function. That is, we show that the case (i) of Theorem 4.2 applies.

The Fourier transform of the transition kernel (5.4) writes as

$$K_{\mathcal{F}}(\omega) = \mathcal{F}[\mathcal{N}(ax_{t-1} + b, c^2)] = e^{i\omega(ax_{t-1} + b) - \frac{1}{2}(\omega c)^2}.$$

Clearly, $|K_{\mathcal{F}}(\omega)| = e^{-\frac{1}{2}(\omega c)^2} |e^{i\omega(ax_{t-1} + b)}| \leq e^{-\frac{1}{2}(\omega c)^2}$.

For the upper bound $e^{-\frac{1}{2}(\omega c)^2}$ and $\beta = 1$ we have

$$(2\pi)^{-1} \int \omega^{2\beta} |e^{-\frac{1}{2}(\omega c)^2}|^2 d\omega = (2\pi)^{-1} \int \omega^2 e^{-(\omega c)^2} d\omega = \frac{1}{2\pi} \frac{\sqrt{\pi}}{2c^3} = \frac{1}{c^3 4\sqrt{\pi}}.$$

Thus, $|K_{\mathcal{F}}(\omega|x_{t-1})|^2 \leq |K_b(\omega)|$ with the bounding function $K_b(\omega) = e^{-\frac{1}{2}(\omega c)^2}$ which is $\beta = 1$ Sobolev with $L_{K_b} = 1/(c^3 4\sqrt{\pi})$.

5.1.4 Properties of univariate Gaussian kernel

Concerning the kernel density estimations, it is natural to choose the standard Gaussian kernel

$$K(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) \quad (5.7)$$

to create the estimator.

The order of the Gaussian kernel is $\ell = \beta = 1$. Indeed, $K_{\mathcal{F}}(\omega) = e^{-0.5\omega^2}$, i.e., $K_{\mathcal{F}}(0) = 1$. $K'_{\mathcal{F}}(\omega) = -\omega e^{-0.5\omega^2}$, i.e., $K'_{\mathcal{F}}(0) = 0$. $K''_{\mathcal{F}}(\omega) = (\omega^2 - 1)e^{-0.5\omega^2}$,

i.e., $K_{\mathcal{F}}''(0) = -1 \neq 0$. Searching for the constant A of Theorem 3.2 leads to the solution of equation $(\omega^2 + 1)e^{-0.5\omega^2} = 1$ which localizes the point of maxima of $(1 - K_{\mathcal{F}}(\omega))/\omega$ for $\omega \in (0, \infty)$. The solution has to be searched numerically¹ and writes as $\omega^* \approx 1.5852$. This gives $|1 - K_{\mathcal{F}}(\omega^*)|/|\omega^*| \approx 0.4513$ and therefore a safe choice for A is $A = 0.5$.

The last entity to specify is the L_2 norm of the kernel. We have

$$\|K\|^2 = \frac{1}{2\pi} \int \exp(-u^2) du = \frac{\sqrt{\pi}}{2\pi} = \frac{1}{2\sqrt{\pi}} \quad \text{hence} \quad \|K\| = \frac{1}{\sqrt{2\pi^{1/4}}}. \quad (5.8)$$

By this we have specified all ingredients needed to run the univariate Gaussian SMC filter and to compute the respective characteristics of interest.

5.1.5 MATLAB implementation

We implemented univariate SMC particle and Kalman filters in a form of the function in the MATLAB[®] computational environment R2012a (ver. 7.14). The implementation of the SMC filter follows the pseudocode of Algorithm 2.2 presented in Chapter 2. The Kalman filter is based on Kalman's equations (5.2) and (5.3).

The interface of the function writes as

```
[SMCm, SMCvar, KFm, KFvar] = uvsmc(HMM, T, n)
```

Inputs:

- **HMM** - is the row vector of seven parameters of an univariate state process. It reads as **HMM** = [**a**, **b**, **c**, **g**, **h**, **m0**, **s0**] where **a**, **b**, **c**, **g**, **h** are the parameters of the Markov chain introduced in formula (5.1), and **m0**, **s0** are the parameters of the initial normal distribution of the chain.
- **T** - is the computational horizon.
- **n** - is the number of particles.

Outputs:

- **SMCm** - is the empirical expected value computed on the basis of the distribution generated by the SMC filter at time **T**.

¹`solve (x^2 + 1) * exp(-0.5 * x^2) = 1, x > 0` at <http://www.wolframalpha.com>

- **SMCvar** - is the empirical variance computed on the basis of the distribution generated by the SMC filter at time T .
- **KFm** - is the theoretical expected value of marginal filtering distribution computed by the Kalman filter for time T .
- **KFvar** - is the theoretical variance of marginal filtering distribution computed by the Kalman filter for time T .

The script further produces a graphical representation of the kernel density estimate of filtering density altogether with its theoretical version.

The source code of the script is presented in Appendix A.

5.1.6 Experiments with univariate Gaussian SMC filter

We have performed several experiments with the filter. The parameters of the process (5.1) were set as $a = 1, b = 1, c = 1, h = 2, g = 1, \mu_0 = 0, \sigma_0 = 1$. The computational horizon was $T = 100$ and we ran the filter for $n = 10, 100$ and $n = 1000$ particles. The results of the experiments are presented in Table 5.1.

| $T = 100$ | $\hat{\mu}_T$ | μ_T | $ \hat{\mu}_T - \mu_T $ | $\hat{\sigma}_T - \text{SMC}$ | $\sigma_T - \text{KF}$ | $ \hat{\sigma}_T - \sigma_T $ |
|-----------|---------------|---------|-------------------------|-------------------------------|------------------------|-------------------------------|
| $n=10$ | 81.57 | 81.44 | 0.13 | 0.3485 | 0.2071 | 0.1414 |
| $n=100$ | 97.84 | 97.74 | 0.10 | 0.1976 | 0.2071 | 0.0095 |
| $n=1000$ | 95.10 | 95.08 | 0.02 | 0.2033 | 0.2071 | 0.0038 |

Table 5.1: Comparison of univariate Gaussian SMC and Kalman filter.

Inspecting the table, we see that the values of empirical integral characteristics are in a good agreement with their theoretical counterparts. In Fig. 5.1 there are graphically presented kernel estimates of true densities of filtering distributions at computational horizon T . In the figure, we also present the values of $\bar{\pi}_t^n g_t$ approximate integral computed according to formula (4.34) and corresponding lower bounds $(\bar{\pi}_t g_t)^{*n}(\epsilon_0)$ given by formula (4.39).

Let us discuss the values of $\bar{\pi}_t^n g_t$ integral. By inspection of Fig. 5.1 we see that the value of the integral varies around some average value. This value is around 0.1. As we have $\|g\|_\infty = g_{\max} = 1/\sqrt{2\pi} = 0.4$ we see that the quotient $q = (1 + 4g_{\max}/\bar{\pi}_t^n g_t)$ of (2.48) has average value $q = (1 + (4 \cdot 0.4)/0.1) = 17$. Hence the exponential growth of c_t applies substantially here. Therefore we cannot directly specify some **reasonable** number of particles in MISE formula (4.1), when we require certain prescribed precision, because C_t grows

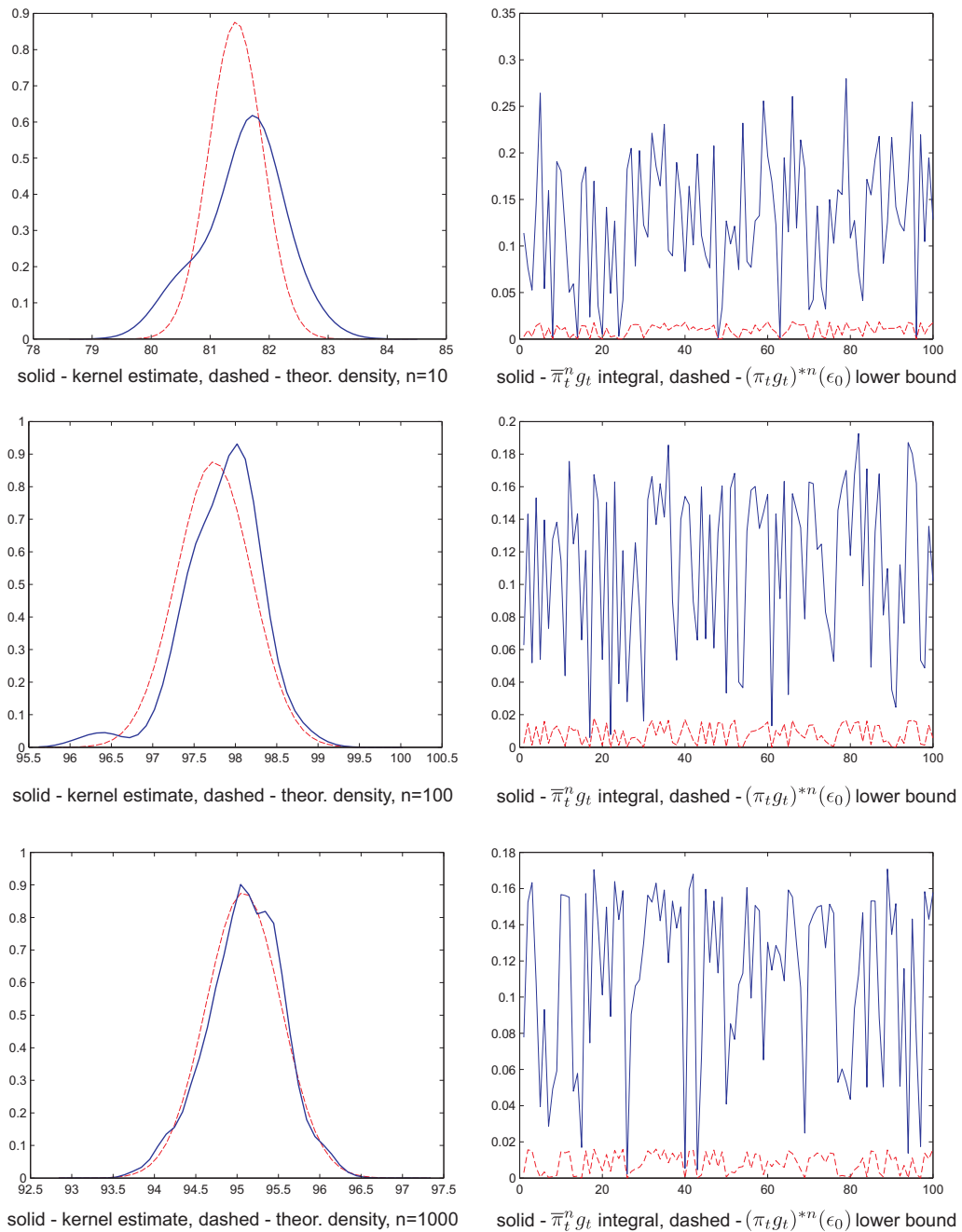


Figure 5.1: Outputs from univariate Gaussian SMC filter.

also exponentially. Using approximate lower bounds makes the situation even worse as the average value of $(\bar{\pi}_t g_t)^{*n}(\epsilon_0)$ is around 0.01. The exact bound of Section 4.3.1 is also inapplicable here due to the long computational horizon.

The problem with the exponential growth of c_t constant is reflected in the literature. Let us cite from [Doucet et al. 2001] p. 87, “The sequence c_t increases exponentially with time, however, so one must also increase the number of particles exponentially with time in order to ensure a given precision at time t . Nevertheless, it has been shown in [Morrall and Guionnet 2001] that, under additional assumptions on the theoretical optimal filter, one can obtain a **uniform convergence result on the marginal distribution** at time t , which is a bound independent of time.”

That is, under some additional assumptions, there exists $c > 0$ independent of number of particles n such that, for all $t \geq 0$ and for any $f \in B(\mathbb{R}^d)$, $\mathbb{E}[|\pi_{t|t}^n f - \pi_{t|t} f|^2] \leq c \cdot n^{-\alpha} \|f\|_\infty^2$ where $\alpha \leq 1$. The uniform convergence result is based on the fact that the theoretical optimal filter, that is $\pi_{t|t}$, exponentially forgets its initial distribution $\pi_0(dx_0)$. This fact is closely linked to the ergodicity of the underlying dynamic model.

So the solution to the presented problem of c_t exponential growth would be the uniform convergence of the filter. Unfortunately, we do not have at our disposal such result for the case of the Gaussian process.

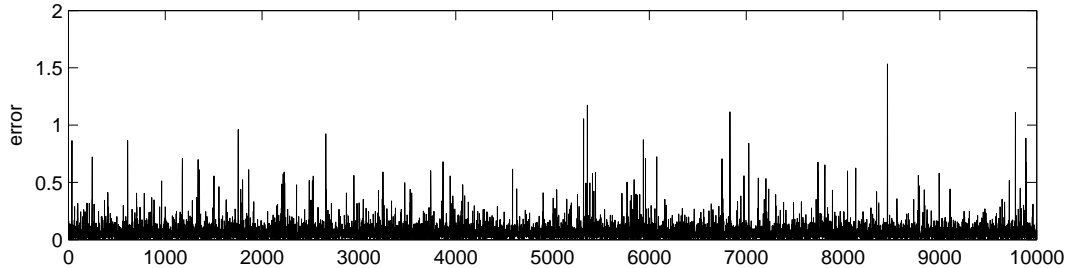


Figure 5.2: Difference of means from univariate SMC and Kalman filter.

On the other hand, the approach presented in [Heine and Crisan 2008] entitles us to consider the convergence of the Gaussian SMC filter to be uniform in spite of the absence of the explicit expression for the constant of the uniform convergence. The reason for this is presented in Fig. 5.2

In the figure, there is presented the graph of error $|\hat{\mu}_t - \mu_t|$ over long time period $t = 1, \dots, 10000$ for $n = 100$ particles employed. We see that the error is uniform in time, which suggest the uniform convergence of the filter. However, the explicit specification of the related constant is likely a matter for further research as we were not able to find it in the available literature.

5.2 Multivariate Gaussian process

After the univariate case of a Gaussian Markov process was discussed, we move to the general dimension $d \geq 1$. The counterpart of the univariate state/observation formulas writes as

$$\mathbf{X}_t = \mathbf{F}\mathbf{X}_{t-1} + \mathbf{W}_t, \quad \mathbf{Y}_t = \mathbf{H}\mathbf{X}_t + \mathbf{V}_t, \quad t \geq 1 \quad (5.9)$$

where \mathbf{F} , \mathbf{H} are $d \times d$ regular matrices and $W_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q})$, $V_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R})$ are multivariate Gaussian noise terms with $d \times d$ covariance matrices \mathbf{Q} and \mathbf{R} . The state process $\{\mathbf{X}_t, t \geq 0\}$ forms a multivariate Markov chain with the Gaussian transition kernel. The initial distribution is considered also multivariate normal, i.e., $\mathbf{X}_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$, $\boldsymbol{\mu}_0 \in \mathbb{R}^d$ and $\boldsymbol{\Sigma}_0$ is a $d \times d$ covariance matrix.

The multivariate filtering task has the same structure as the univariate one. We want to state the best estimate (in L_2 sense) of the current state on the basis of an observed history. Mathematically, this requires specification of conditional expected values $\mathbb{E}[\mathbf{X}_t | \mathbf{Y}_1, \dots, \mathbf{Y}_t]$ for $t \geq 1$. At a given time instant t , the conditional expected value is the integral characteristic of the respective conditional distribution which actually represents the filtering distribution we are searching for.

The vector $(\mathbf{X}_0, \mathbf{X}_1, \mathbf{Y}_1, \dots, \mathbf{X}_t, \mathbf{Y}_t)$ is multivariate Gaussian because it is given by a linear transformation of the Gaussian vector $(\mathbf{X}_0, \mathbf{W}_1, \mathbf{V}_1, \dots, \mathbf{W}_t, \mathbf{V}_t)$ [Fristedt et al. 2007]. Therefore the filtering distribution has also a multivariate normal distribution determined by a vector of means $\boldsymbol{\mu}_t$ and a covariance matrix $\boldsymbol{\Sigma}_t$. The preservation of the normal character of the filtering distribution over time allows the analytic expression of its parameters. The result is known as the *multivariate Kalman filter*.

5.2.1 Multivariate Kalman filter

The results of the theoretical analysis presented in [Pollock 1999] pp. 244 and 246 give the following recursive Kalman's equations which are computed in several steps using some auxiliary variables for $t \geq 1$:

$$\begin{aligned}
\boldsymbol{\mu}_{t|t-1} &= \mathbf{F}\boldsymbol{\mu}_{t-1} && (\text{state prediction}) \\
\boldsymbol{\Sigma}_{t|t-1} &= \mathbf{F}\boldsymbol{\Sigma}_{t-1}\mathbf{F}^T + \mathbf{Q} && (\text{covariance prediction}) \\
\mathbf{K}_t &= \boldsymbol{\Sigma}_{t|t-1}\mathbf{H}^T(\mathbf{H}\boldsymbol{\Sigma}_{t|t-1}\mathbf{H}^T + \mathbf{R})^{-1} && (\text{Kalman gain}) \\
\boldsymbol{\mu}_t &= \boldsymbol{\mu}_{t|t-1} + \mathbf{K}_t(\mathbf{Y}_t - \mathbf{H}\boldsymbol{\mu}_{t|t-1}) && (5.10) \\
\boldsymbol{\Sigma}_t &= (\mathbf{I}_d - \mathbf{K}_t\mathbf{H})\boldsymbol{\Sigma}_{t|t-1} && (5.11)
\end{aligned}$$

Using the above formulas, we can recursively compute the determining characteristics of the filtering distribution over time. Due to the normal character of the distribution we have apparently $\mathbb{E}[\mathbf{X}_t|\mathbf{Y}_1, \dots, \mathbf{Y}_t] = \boldsymbol{\mu}_t$. Further, similarly as in the univariate case, the formula for the evolution of the covariance matrix is deterministic. That is, it is not affected by observations.

5.2.2 Multivariate Gaussian SMC filter

The incorporation of schema (5.9) into the SMC filter's framework stems from the specification of the Gaussian transition kernel. This specification reads as

$$K(\mathbf{x}_t|\mathbf{x}_{t-1}) = (2\pi)^{-\frac{d}{2}}|\mathbf{Q}|^{-\frac{1}{2}} \exp \left[-(\mathbf{x} - \mathbf{F}\mathbf{x}_{t-1})^T \mathbf{Q}^{-1}(\mathbf{x} - \mathbf{F}\mathbf{x}_{t-1}) \right]. \quad (5.12)$$

The above formula reflects the multivariate normal character of noise term \mathbf{W}_t in (5.9) and in fact corresponds to the specification of the density of multivariate normal distribution $\mathcal{N}(\mathbf{F}\mathbf{X}_{t-1}, \mathbf{Q})$.

The analysis of properties of a general multivariate density helps us to specify bounding constants $K_{1\max}$, $K_{2\max}$ and Lipschitz constants $K_{1\alpha}$, $K_{2\alpha}$ required for computations related to the Gaussian SMC filter.

The density of a multivariate normal distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with mean $\boldsymbol{\mu} \in \mathbb{R}^d$ and a $d \times d$ symmetric positive semidefinite covariance matrix $\boldsymbol{\Sigma}$ has form

$$p(\mathbf{x}) = (2\pi)^{-\frac{d}{2}}|\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right] \quad (5.13)$$

with the gradient ∇p given by respective partial derivatives

$$\nabla p(\mathbf{x}) = \frac{\partial p}{\partial \mathbf{x}} = p(\mathbf{x}) \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \mathbf{x}). \quad (5.14)$$

It is well known, and can be easily shown from (5.14), that maximum of $p(\mathbf{x})$ is reached at $\mathbf{x}^* = \boldsymbol{\mu}$ with the value of $p(\mathbf{x}^*) = p_{\max} = (2\pi)^{-\frac{d}{2}}|\boldsymbol{\Sigma}|^{-\frac{1}{2}}$.

The value of maxima p_{\max} does not depend on the value of expectation $\boldsymbol{\mu}$. This immediately gives us the bounding constants for transition kernel (5.12) as

$$K_{1\max} = K_{2\max} = (2\pi)^{-\frac{d}{2}} |\mathbf{Q}|^{-\frac{1}{2}}. \quad (5.15)$$

As $\boldsymbol{\Sigma}^{-1}$ is also symmetric positive semidefinite matrix, it has non-negative real eigenvalues $\lambda_1, \dots, \lambda_d$. Denoting λ_{\min} or λ_{\max} minimal or maximal eigenvalue, respectively, the well known inequality for positive semidefinite matrices gives

$$\lambda_{\min} \|\mathbf{u}\|^2 \leq \mathbf{u}^T \boldsymbol{\Sigma}^{-1} \mathbf{u} \leq \lambda_{\max} \|\mathbf{u}\|^2. \quad (5.16)$$

Therefore we have

$$\exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \leq \exp \left[-\frac{1}{2} \lambda_{\min} \|\mathbf{x} - \boldsymbol{\mu}\|^2 \right].$$

Since the spectral norm of $\boldsymbol{\Sigma}^{-1}$ equals to λ_{\max} , i.e., $\|\boldsymbol{\Sigma}^{-1}\|_{\text{spc}} = \lambda_{\max}$ the expression (5.14) for the gradient of $p(\mathbf{x})$ implies

$$\begin{aligned} \|\nabla p\| &\leq | -p(\mathbf{x}) | \cdot \|\boldsymbol{\Sigma}^{-1}\|_{\text{spc}} \cdot \|\mathbf{x} - \boldsymbol{\mu}\|, \\ \|\nabla p\| &\leq \frac{\lambda_{\max}}{(2\pi)^d |\boldsymbol{\Sigma}|^{-\frac{1}{2}}} \exp \left[-\frac{1}{2} \lambda_{\min} \|\mathbf{x} - \boldsymbol{\mu}\|^2 \right] \cdot \|\mathbf{x} - \boldsymbol{\mu}\|. \end{aligned}$$

Searching for the maximum of function $f(z) = z \exp[-0.5\lambda_{\min}z^2]$, $z \geq 0$ gives the point of maxima $z^* = 1/\sqrt{\lambda_{\min}}$ and $f(z^*) = 1/\sqrt{\lambda_{\min}} \exp[-0.5]$. Thus

$$\|\nabla p\| \leq \frac{\lambda_{\max}}{\lambda_{\min} (2\pi)^d |\boldsymbol{\Sigma}|^{-\frac{1}{2}}} \exp \left[-\frac{1}{2} \right] = K_{\alpha}. \quad (5.17)$$

We employ the above formula for the specification of Lipschitz constants of transition kernel (5.12). The multivariate Taylor's theorem applied to a general density (5.13) states that for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ there exists some $\boldsymbol{\xi} = \mathbf{x} + \vartheta(\mathbf{y} - \mathbf{x})$, $\vartheta \in (0, 1)$ such that $p(\mathbf{x}) - p(\mathbf{y}) = (\nabla p)^T(\boldsymbol{\xi})(\mathbf{x} - \mathbf{y})$. Employing (5.17) this implies that

$$|p(\mathbf{x}) - p(\mathbf{y})| \leq \|\nabla p^T(\boldsymbol{\xi})\| \cdot \|\mathbf{x} - \mathbf{y}\| \leq K_{\alpha} \cdot \|\mathbf{x} - \mathbf{y}\|,$$

i.e., $p(\mathbf{x})$ is Lipschitz with constant K_{α} given by (5.17).

Since K_{α} does not depend on the value of expectation, the constant K_{α} also determines the Lipschitz constant of the transition kernel $K(\mathbf{x}_t | \mathbf{x}_{t-1})$ of (5.12) in the first variable. That is

$$K_{1\alpha} = \frac{\lambda_{\max}}{\lambda_{\min} (2\pi)^d |\mathbf{Q}|^{-\frac{1}{2}}} \exp \left[-\frac{1}{2} \right].$$

Concerning the Lipschitz constant of the transition kernel in the second variable we may consider $K(\mathbf{x}_t|\mathbf{x}_{t-1})$ of (5.12) as the compound function of the density of $\mathcal{N}(\cdot, \mathbf{Q})$ distribution and linear function $\mathbf{F}\mathbf{x}_{t-1}$. Employing the properties of spectral norm we have $\|\mathbf{F}\mathbf{y}_{t-1} - \mathbf{F}\mathbf{x}_{t-1}\| \leq \|\mathbf{F}\|_{spc} \|\mathbf{y}_{t-1} - \mathbf{x}_{t-1}\|$, i.e., linear function $\mathbf{F}\mathbf{x}_{t-1}$ is Lipschitz with constant $\|\mathbf{F}\|_{spc}$. The compound function determining the transition kernel as the function of \mathbf{x}_{t-1} variable is then Lipschitz with constant

$$K_{2\alpha} = \frac{\lambda_{\max}\|\mathbf{F}\|_{spc}}{\lambda_{\min}(2\pi)^d|\mathbf{Q}|^{-\frac{1}{2}}} \exp\left[-\frac{1}{2}\right].$$

The identical considerations then specify the Lipschitz constant for function $g(\mathbf{y}_t - \mathbf{H}\mathbf{x}_t)$ where g is the density of $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ distribution. The Lipschitz constant of this compound function with respect to \mathbf{x}_t variable reads as $g_\alpha h_\alpha$ where

$$g_\alpha h_\alpha = \frac{\|\mathbf{H}\|_{spc}}{(2\pi)^d} \exp\left[-\frac{1}{2}\right].$$

5.2.3 Sobolev character of multivariate filter

The Sobolev character of the multivariate filter is given by the Sobolev character of the multivariate Gaussian transition kernel. Let us show that the conditional characteristic function of the kernel is bounded, which implies the Sobolev character of the filter.

We have

$$\begin{aligned} \mathcal{F}_K(\boldsymbol{\omega}|\mathbf{x}_{t-1}) &= \int e^{i\langle\boldsymbol{\omega}, \mathbf{x}_t\rangle} K(\mathbf{x}_t|\mathbf{x}_{t-1}) d\mathbf{x}_t = \mathcal{F}[\mathcal{N}(\mathbf{F}\mathbf{x}_{t-1}, \mathbf{Q})], \\ \mathcal{F}_K(\boldsymbol{\omega}|\mathbf{x}_{t-1}) &= e^{i\langle\boldsymbol{\omega}, \mathbf{F}\mathbf{x}_{t-1}\rangle} \exp\left[-\frac{1}{2}\boldsymbol{\omega}^T \mathbf{Q} \boldsymbol{\omega}\right], \end{aligned}$$

and therefore

$$|\mathcal{F}_K(\boldsymbol{\omega}|\mathbf{x}_{t-1})| \leq \exp\left[-\frac{1}{2}\boldsymbol{\omega}^T \mathbf{Q} \boldsymbol{\omega}\right] \leq \exp\left[-\frac{1}{2}\lambda_{\min}\|\boldsymbol{\omega}\|^2\right] = K_b(\boldsymbol{\omega})$$

where λ_{\min} is the minimal eigenvalue of the covariance matrix \mathbf{Q} .

Let us compute the Sobolev constant of $K_b(\boldsymbol{\omega})$ for $\beta = 1$. We have

$$(2\pi)^{-d} \int \|\boldsymbol{\omega}\|^{2\beta} |K_b(\boldsymbol{\omega})|^2 = (2\pi)^{-d} \int \|\boldsymbol{\omega}\|^2 \exp[-\lambda_{\min}\|\boldsymbol{\omega}\|^2] = \frac{\pi^{-\frac{d}{2}}}{4^d(\sqrt{\lambda_{\min}})^{d+2}}.$$

From this result we immediately have also that any initial distribution with covariance matrix Σ_0 is 1-Sobolev with the constant $L_0 = 1/[4^d(\sqrt{\lambda_{\min}^0})^{d+2}\pi^{\frac{d}{2}}]$ where λ_{\min}^0 is the minimal eigenvalue of Σ_0 . The obtained result is consistent with the fact that all densities in the presented multivariate process (5.9) are normal, i.e., the character of the involved densities does not change.

5.2.4 Properties of multivariate Gaussian kernel

Kernel density estimates in the multivariate SMC Gaussian filter are performed using the multivariate standard normal kernel

$$K(\mathbf{u}) = (2\pi)^{-\frac{d}{2}} \exp\left[-\frac{1}{2}\mathbf{u}^T \mathbf{I}_d^{-1} \mathbf{u}\right] = (2\pi)^{-\frac{d}{2}} \exp\left[-\frac{1}{2}\|\mathbf{u}\|^2\right]. \quad (5.18)$$

The specification of the L_2 norm of the kernel is straightforward. We have

$$\|K\|^2 = (2\pi)^{-d} \int \exp(-\|\mathbf{u}\|^2) d\mathbf{u} = (4\pi)^{-\frac{d}{2}} \quad \text{hence} \quad \|K\| = (4\pi)^{-\frac{d}{4}}.$$

Concerning the A constant of the Theorem 3.2, we start with the Fourier transform of the standard multivariate Gaussian kernel which corresponds to the characteristic function of $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ distribution. That is $K_{\mathcal{F}}(\boldsymbol{\omega}) = e^{-\boldsymbol{\omega}^T \mathbf{I}_d \boldsymbol{\omega}} = e^{-\frac{1}{2}\|\boldsymbol{\omega}\|^2}$. In order to specify some constant A we need to determine a bound on the spectral norm of the Hessian of $K_{\mathcal{F}}$. The entries of the Hessian matrix $\mathcal{H}(K_{\mathcal{F}})$ reads as

$$\frac{\partial^2 K_{\mathcal{F}}}{\partial \omega_j^2} = (\omega_j^2 - 1)K_{\mathcal{F}}(\boldsymbol{\omega}), \quad \frac{\partial K_{\mathcal{F}}}{\partial \omega_j \partial \omega_k} = K_{\mathcal{F}}(\boldsymbol{\omega})\omega_j \omega_k, \quad j \neq k.$$

In the matrix notation the Hessian writes as $\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\omega}) = K_{\mathcal{F}}(\boldsymbol{\omega})(\boldsymbol{\omega}\boldsymbol{\omega}^T - \mathbf{I}_d)$. Employing the spectral norm we get

$$\begin{aligned} \|\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\omega})\|_{spc} &\leq K_{\mathcal{F}}(\boldsymbol{\omega})\|\boldsymbol{\omega}\boldsymbol{\omega}^T - \mathbf{I}_d\|_{spc} \\ &\leq K_{\mathcal{F}}(\boldsymbol{\omega})(\|\boldsymbol{\omega}\boldsymbol{\omega}^T\|_{spc} + \|\mathbf{I}_d\|_{spc}) \\ &\leq K_{\mathcal{F}}(\boldsymbol{\omega})(\|\boldsymbol{\omega}^T\| \|\boldsymbol{\omega}\| + 1) \\ &\leq K_{\mathcal{F}}(\boldsymbol{\omega})(\|\boldsymbol{\omega}\|^2 + 1). \end{aligned}$$

Note that for a vector $\boldsymbol{\omega}$ it is $\|\boldsymbol{\omega}\|_{spc} = \|\boldsymbol{\omega}\|$. Let $\boldsymbol{\omega} = \boldsymbol{\xi}$ such that $\|\boldsymbol{\xi}\| \leq 1$. Then we clearly have $\|\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\xi})\| \leq 2$ as $K_{\mathcal{F}}(\boldsymbol{\xi}) \leq 1$.

Now, the multidimensional Taylor's theorem for $K_{\mathcal{F}}$ writes as

$$K_{\mathcal{F}}(\boldsymbol{\omega}) = K_{\mathcal{F}}(\mathbf{0}) + \frac{1}{1!} \sum_{j=1}^d K'_{\mathcal{F}}(\mathbf{0})\omega_j + \cdots + \frac{1}{2} \boldsymbol{\omega}^T [\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\xi})] \boldsymbol{\omega}$$

As the standard Gaussian kernel is of order $\ell = 1$ we have the first partial derivatives equal to zero. Since $K_{\mathcal{F}}(\mathbf{0}) = 1$ then for any $\|\boldsymbol{\omega}\| \leq 1$ the Taylor's theorem gives

$$\begin{aligned} K_{\mathcal{F}}(\boldsymbol{\omega}) - 1 &= \frac{1}{2} \boldsymbol{\omega}^T [\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\xi})] \boldsymbol{\omega}, \\ |K_{\mathcal{F}}(\boldsymbol{\omega}) - 1| &\leq \frac{1}{2} \|\boldsymbol{\omega}^T\| \cdot \|\mathcal{H}(K_{\mathcal{F}})(\boldsymbol{\xi})\|_{\text{spc}} \cdot \|\boldsymbol{\omega}\|, \\ \frac{|K_{\mathcal{F}}(\boldsymbol{\omega}) - 1|}{\|\boldsymbol{\omega}\|} &\leq \|\boldsymbol{\omega}^T\| = \|\boldsymbol{\omega}\|. \end{aligned}$$

Further $|K_{\mathcal{F}}(\boldsymbol{\omega}) - 1| \leq 1$ for all $\boldsymbol{\omega}$ and therefore $|K_{\mathcal{F}}(\boldsymbol{\omega}) - 1|/\|\boldsymbol{\omega}\| \leq 1$ for $\|\boldsymbol{\omega}\| > 1$. Hence joining the two inequalities we finally get

$$\frac{|K_{\mathcal{F}}(\boldsymbol{\omega}) - 1|}{\|\boldsymbol{\omega}\|} \leq \max\{1, 1\} = 1.$$

So the A constant equals to 1, i.e., $A = 1$.

The above considerations helps us to state the order of the kernel. The Fourier transform of the kernel is $K_{\mathcal{F}}(\boldsymbol{\omega}) = e^{-\frac{1}{2}\|\boldsymbol{\omega}\|^2}$, so $K_{\mathcal{F}}(\mathbf{0}) = \mathbf{1}$. The related gradient writes as $\nabla K_{\mathcal{F}}(\boldsymbol{\omega}) = -e^{-\frac{1}{2}\|\boldsymbol{\omega}\|^2} \boldsymbol{\omega}$, thus $\nabla K_{\mathcal{F}}(\mathbf{0}) = \mathbf{0}$. For the Hessian of $K_{\mathcal{F}}(\boldsymbol{\omega})$ we have $\text{diag}(\mathcal{H}(K_{\mathcal{F}})(\mathbf{0})) = -\mathbf{1}$. Hence the order of the kernel is $\ell = \beta = 1$.

5.2.5 MATLAB implementation and experiments

The multivariate SMC filter is again implemented in the form of a function in the MATLAB computational environment. The implementation of the filter follows the pseudocode of Algorithm 2.2 presented in Chapter 2. The multivariate Kalman's equations are employed to obtain the parameters of filtering distributions.

The interface of the function writes as:

```
[SMCm, SMCcov, KFm, KFcov] = mvsmc(F, Q, H, R, T, n)
```

Inputs:

- $\mathbf{F}, \mathbf{Q}, \mathbf{H}, \mathbf{R}$ - are $d \times d$ matrices determining the model of the Gaussian process specified in (5.9).
- $\mathbf{M}_0, \mathbf{S}_0$ - are the parameters of the initial distribution of $\mathbf{X}_0 \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.
- T - is the computational horizon.
- n - is the number of particles.

Outputs:

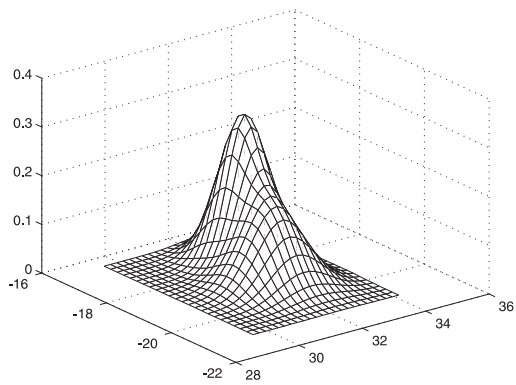
- \mathbf{SMCm} - is the vector of empirical expected values computed on the basis of the distribution generated by the SMC filter at time T .
- \mathbf{SMCcov} - is the empirical covariance matrix of the distribution generated by the SMC filter at time T .
- \mathbf{KFm} - is the vector of theoretical expected values obtained from the Kalman filter for time T .
- \mathbf{KFcov} - is the theoretical covariance matrix obtained from the Kalman filter for time T .

If the dimension of state process is $d = 2$, the script produces graphs of kernel density estimate and theoretical density of the marginal filtering distribution.

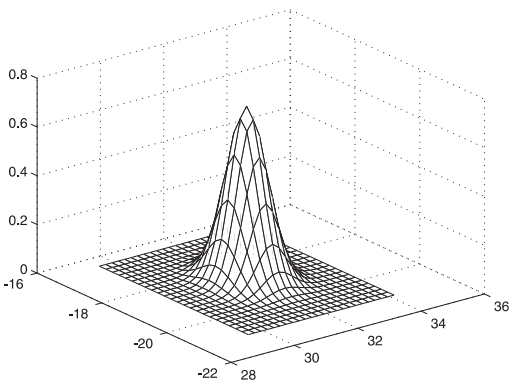
The source code of the script is presented in Appendix A.

We have performed several experiments in order to check if the computational behavior of the Gaussian SMC filter coincides with the analytical results. The experiments were performed for the following setting of parameters: $\mathbf{F} = \mathbf{I}_d$, $\mathbf{Q} = 2\mathbf{I}_d$, $\mathbf{H} = 2\mathbf{I}_d$, $\mathbf{R} = \mathbf{I}_d$. Computational horizon was set to $T = 100$. The results of three experiments with different number of particles $n = 10, 100$ and $n = 1000$ are presented in Table 5.2. Graphically, the obtained density estimates and theoretical filtering densities are presented in Fig. 5.3.

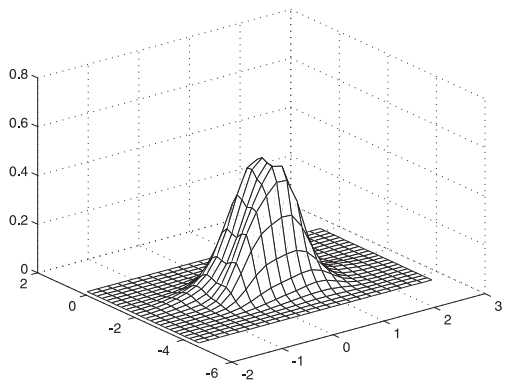
On the basis of the inspection of numerical results presented in Table 5.2, we can state a good agreement of numerical characteristics delivered by the Gaussian SMC filter with the theoretical characteristics of filtering distributions.



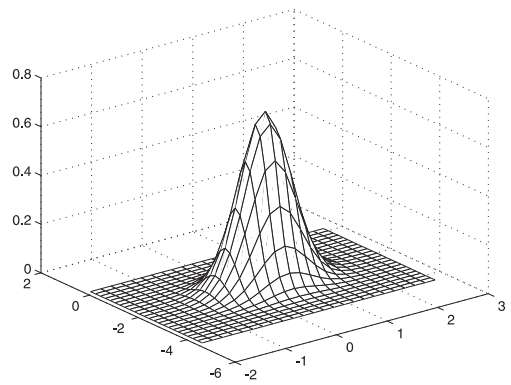
SMC filter, $T=100$, $n=10$



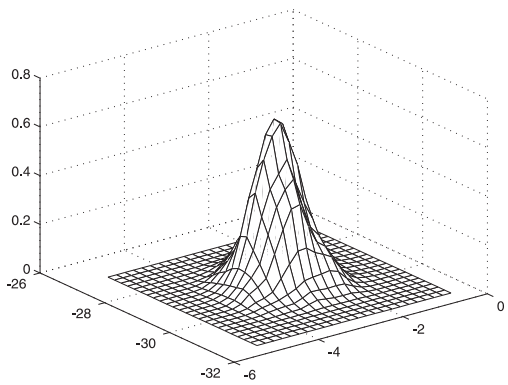
Kalman filter, $T=100$, $n=10$



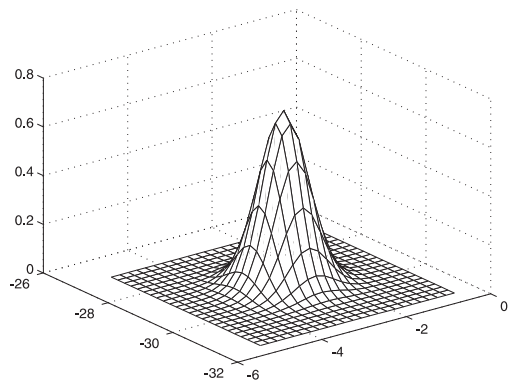
SMC filter, $T=100$, $n=100$



Kalman filter, $T=100$, $n=100$



SMC filter, $T=100$, $n=1000$



Kalman filter, $T=100$, $n=1000$

Figure 5.3: Kernel density estimates generated by a Gaussian SMC particle filter and theoretical filtering densities for a two-dimensional Gaussian process.

| $T=100$ | $\hat{\boldsymbol{\mu}}_T$ | $\boldsymbol{\mu}_T$ | $\ \hat{\boldsymbol{\mu}}_T - \boldsymbol{\mu}_T\ $ | $\hat{\boldsymbol{\Sigma}}_T$ - SMC | | $\boldsymbol{\Sigma}_T$ - KF | |
|----------|----------------------------|----------------------|---|-------------------------------------|---------|------------------------------|--------|
| $n=10$ | 32.25 | 31.92 | 0.41 | 0.1472 | 0.0992 | 0.2247 | 0 |
| | -18.43 | -18.65 | | 0.0092 | 0.4290 | 0 | 0.2247 |
| $n=100$ | 0.46 | 0.48 | 0.12 | 0.1557 | -0.0212 | 0.2247 | 0 |
| | -2.16 | -2.04 | | -0.0212 | 0.2144 | 0 | 0.2247 |
| $n=1000$ | -2.76 | -2.75 | 0.01 | 0.2207 | -0.0036 | 0.2247 | 0 |
| | -29.18 | -29.18 | | -0.0036 | 0.2206 | 0 | 0.2247 |

Table 5.2: Comparison of two-variate Gaussian SMC and Kalman filter.

The discussion concerning the uniform convergence of the Gaussian SMC particle filter remains valid here. Again, our suggestion for the uniform convergence is supported by the graph of time evolution of the norm of difference between empirical and theoretical expectations $\|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}_t\|$. In the experiment we employed $n = 100$ particles and the computational horizon was set to $T = 10000$. The graph is presented in Fig. 5.4.

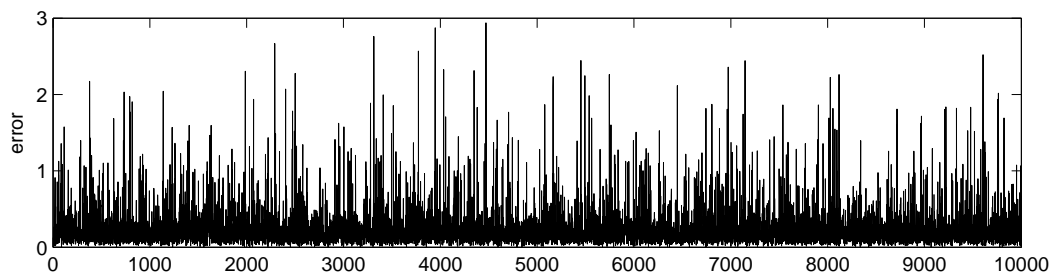


Figure 5.4: Time evolution of error $\|\hat{\boldsymbol{\mu}}_t - \boldsymbol{\mu}_t\|$ for a two-dimensional Gaussian process.

6. Summary

The thesis assignment reads as follows: “The student will study the theory of sequential Monte Carlo methods including their limit behavior. He will implement related algorithms and demonstrate their work on simulated data.” Let us comment on the assignment in the view of the thesis composition and the presented results.

The thesis opens with the short introductory chapter. The idea and the algorithm of sequential Monte Carlo methods (SMC methods) are reviewed in details in the second chapter. The application of SMC methods in the context of particle filters is also presented here. We review the convergence analysis of a SMC particle filter as it is presented in [Doucet et al. 2001]. The main result of the chapter shows that integral characteristics of empirical distributions generated by a SMC particle filter converge in expectation to integral characteristics of theoretical filtering distributions if the number of particles (samples) goes to infinity. The proofs of the related theorems are provided in full details.

In the third chapter, we review the basics of the theory of nonparametric kernel density estimation. We present results of the classical approach based on the asymptotic MISE analysis and also results based on the Fourier analysis. The review of the Fourier analysis approach is based on the book by [Tsybakov 2009]. We have extended his univariate results to the multivariate case. The Fourier analysis proves to be advantageous because kernel density estimates are related to the operation of convolution which is comfortably handled in the frequency domain. In order to the estimates on MISE (the mean integrated squared error) be bounded, certain assumptions on the character of estimated densities and employed kernel have to be made. This leads to the notions of the Sobolev class of densities and the order of kernel.

Our original results are presented in Chapter 4. It deals with the kernel density estimates of marginal empirical distributions generated by a SMC particle filter. The main result we have proved shows that kernel density estimates converge to theoretical filtering densities as the number of particles increases to infinity. The result is proved without assumption on the i.i.d. character of generated data. This is crucial in the context of SMC particle filters because generated samples are not i.i.d. due to resampling. We consider this result as to be the most important one achieved in the thesis.

The second result shows that the Sobolev character of a SMC particle filter is retained over time. This is an important assumption in the theorem on the convergence of kernel density estimates. The result is obtained by the Fourier analysis of the prediction and update formulas which recursively describe the evolution of theoretical filtering distributions.

The last group of presented results deals with computations of the integral which determines the values of number-of-particles constants occurring in the upper bound on the error of kernel density estimates. We were able to specify the exact lower bound on this integral (it occurs in denominators), however, for higher times and dimensions this bound decreases very quickly to zero. We further present approximate bounds based on empirical integrals which are computed during the operation of a SMC filter.

The computational algorithm of a SMC particle filter was programmed and tested in the MATLAB computational environment. Both univariate and multivariate cases were dealt with. We studied whether the results of the operation of the filter coincide with the analytical results which can be obtained for Gaussian processes. The analytical solution for this case is known as the Kalman filter. We found a good agreement of experimental and analytical results.

The results of performed experiments are presented in Chapter 5. In the chapter we also investigate the properties of the Gaussian transition kernel and normal observation density. The obtained results can be directly applied for the practical specification of kernel density estimates' error under the assumption of the uniform convergence of the filter and the knowledge of the constant driving the convergence. The specification of the explicit value of this constant constitutes a practically and theoretically interesting question. In fact, the issue of the uniform convergence of particle filters is nowadays a vivid research area.

In conclusion we can say that the thesis provides an interested reader with a basic introduction to the application of nonparametric kernel methods in the area of SMC particle filters. The introduction has a sufficient level of rigor for the presented combination of these two methodologies can be further mathematically investigated and developed.

A. MATLAB source codes

A.1 uvsmc.m

```
1 function [SMCm,SMCvar,KFm,KFvar]=uvsmc(HMM,T,n);
2
3 %---HMM---
4 a=HMM(1);b=HMM(2);c=HMM(3);
5 h=HMM(4);g=HMM(5);
6 m0=HMM(6);s0=HMM(7);
7 X0=m0+s0*randn(1,1);
8 X=zeros(1,T);Y=X;
9 X(1)=a*X0+b+c*randn(1,1);
10 Y(1)=h*X(1)+g*randn(1,1);
11 for t=2:T,
12     X(t)=a*X(t-1)+b+c*randn(1,1);
13     Y(t)=h*X(t)+g*randn(1,1);
14 end;
15
16 %---Kalman filter---
17 m=zeros(1,T);s2=zeros(1,T);
18 v2=c^2*h^2+g^2;
19 m(1)=g^2/v2*b+c^2/v2*h*Y(1);
20 m(1)=m(1)+g^2/v2*a*(s0^2*a*h*(Y(1)-b*h)+v2*m0)/(h^2*s0^2*a^2+v2);
21 s2(1)=g^4/v2*a^2*s0^2/(h^2*s0^2*a^2+v2)+g^2*c^2/v2;
22 for t=2:T,
23     m(t)=g^2/v2*b+c^2/v2*h*Y(t);
24     mupd=g^2/v2*a*(s2(t-1)*a*h*(Y(t)-b*h)+v2*m(t-1))/(h^2*s2(t-1)*a^2+v2);
25     m(t)=m(t)+mupd;
26     s2(t)=g^4/v2*a^2*s2(t-1)/(h^2*s2(t-1)*a^2+v2)+g^2*c^2/v2;
27 end;
28 KFm=m(T);
29 KFvar=s2(T);
30
31 %---SMC filter & approximate bounds---
32 d=1;
33 Dd=(1/(d+1))*((d)/(d+1))^(d);
34 Dd1=(1/(d+2))*((d+1)/(d+2))^(d+1);
```

```

35 Vd=pi^(d/2)/gamma(d/2+1);
36 K1max=1/(c*sqrt(2*pi));
37 K1a=1/(c^2*sqrt(2*pi))*exp(-0.5);
38 K2max=1/(c*sqrt(2*pi));
39 K2a=abs(a)/(c^2*sqrt(2*pi))*exp(-0.5);
40 gmax=1/(g*sqrt(2*pi));
41 ga=1/(g^2*sqrt(2*pi))*exp(-0.5);
42 ha=abs(h);
43 gpa=gmax*K1a+K1max*ga*ha;
44 M=(gmax^(d+1)*Dd*Vd)/gpa^d;
45
46 P0=m0+s0*randn(n,1);
47 err=zeros(1,T);
48 ippg=zeros(1,T);
49 appg=zeros(1,T);
50 aa=zeros(1,T); ppa=aa; nt0=aa;
51 for t=1:T,
52     disp(t);
53     aa(t)=Y(t)/h;
54     if (t==1)
55         ppa(t)=1/n*sum(normpdf(aa(t),a*P0+b,c));
56         nt0(t)=ceil(((5*gmax*K2max)/(M*ppa(t)^(d+2)*Dd1))^2);
57         pp=a*P0+b+c*randn(n,1);
58         P=[];
59     else
60         ppa(t)=1/n*sum(normpdf(aa(t),a*P(:,t-1)+b,c));
61         nt0(t)=ceil((cc(t-1)*(5*gmax*K2max)/(M*ppa(1)^(d+2)*Dd1))^2);
62         pp=a*P(:,t-1)+b+c*randn(n,1);
63     end;
64     Ppp=[P pp];
65     w=normpdf(Y(t)-h*pp,0,g);
66     ippg(t)=1/n*sum(w);
67     appg(t)=M*ppa(t)^(d+1)*(d+2)*Dd1;
68     if (t==1) cc(t)=(1+(4*gmax)/ippg(t));
69     else cc(t)=cc(t-1)*(1+(4*gmax)/ippg(t)); end;
70     wn=w/sum(w);
71     mn=randsample(n,n,true,wn);
72     P=Ppp(mn,:); P0=P0(mn,:);
73     err(t)=abs(mean(pp(mn,:))-m(t));
74 end;

```



```

75  pT=P(:,T);
76  SMCm=mean(pT);
77  SMCvar=var(pT);
78
79  close(gcf);
80  %---approximate integrals---
81  plot(1:T,ippg,1:T,appg,'--r');
82  pause;
83
84  %---mean error--
85  [abs(SMCm-KFm) sum(err)/T]
86  plot(1:T,err,'k');
87  pause;
88
89  %---kernel density estimate---
90  alpha=1;beta=1;
91  hn=alpha*n^(-1/(2*beta+d));
92  xx=[mean(pT)-5*std(pT):0.1:mean(pT)+5*std(pT)];
93  fx=zeros(1,length(xx));
94  for i=1:n,
95    fx=fx+1/(n*hn^d)*1/sqrt(2*pi)*exp(-(xx-pT(i)).^2/(2*hn^2));
96  end;
97  plot(xx,fx,'b',xx,normpdf(xx,m(T),sqrt(s2(T))), '--r');

```

A.2 mvsmc.m

```
1 function [SMCm,SMCcov,KFm,KFcov]=mvsmc(F,Q,H,R,T,n);
2
3 %---HMM---
4 d=size(Q,1);
5 m0=zeros(d,1);S0=eye(d);
6 X0=mvnrnd(m0',S0)';
7 X=zeros(d,T);Y=X;
8 X(:,1)=F*X0+mvnrnd(zeros(1,d),Q)';
9 Y(:,1)=H*X(:,1)+mvnrnd(zeros(1,d),R)';
10 for t=2:T,
11     W=mvnrnd(zeros(1,d),Q)';
12     V=mvnrnd(zeros(1,d),R)';
13     X(:,t)=F*X(:,t-1)+W;
14     Y(:,t)=H*X(:,t)+V;
15 end;
16
17 %---Kalman filter---
18 M=zeros(d,T);
19 m=m0;S=S0;
20 for t=1:T,
21     m1=F*m;
22     S1=F*S*F'+Q;
23     K=S1*H'*inv(H*S1*H'+R);
24     m=m1+K*(Y(:,t)-H*m1);
25     S=(eye(d)-K*H)*S1;
26     M(:,t)=m;
27 end;
28 KFm=M(:,T);
29 KFcov=S;
30
31 %---SMC filter---
32 P=mvnrnd(m0',S0,n)';
33 err=zeros(1,T);
34 for t=1:T,
35     disp(t);
36     pp=zeros(d,n);w=zeros(1,n);
37     for j=1:n;
38         pp(:,j)=F*P(:,j)+mvnrnd(zeros(1,d),Q)';
```

```

39     w(j)=mvnpdf((Y(:,t)-H*pp(:,j))',zeros(1,d),R);
40     end;
41     wn=w/sum(w);
42     mn=randsample(n,n,true,wn);
43     P=pp(:,mn);
44     err(t)=norm(mean(P')'-M(:,t));
45     end;
46     PT=P;
47     SMCm=mean(PT')';
48     SMCcov=cov(PT');
49
50     close(gcf);
51     %---mean error---
52     [norm(SMCm-KFm) sum(err)/T]
53     plot(1:T,err,'k');
54     pause;
55
56     %---kernel density estimate for general dimension at the origin---
57     alpha=1;beta=1;
58     hn=alpha*n^(-1/(2*beta+d));
59     x0=zeros(d,1);
60     fx0=0;
61     for j=1:n
62         fx0=fx0+1/(n*hn^d)*mvnpdf((x0'-PT(:,j))')/hn,zeros(1,d),eye(d));
63     end;
64
65     %---kernel density estimate for d=2 with the graphical output---
66     if d==2,
67         alpha=1;beta=1;
68         hn=alpha*n^(-1/(2*beta+d));
69         x1=[KFm(1)-5*sqrt(KFcov(1,1)):0.2:KFm(1)+5*sqrt(KFcov(1,1))];
70         x2=[KFm(2)-5*sqrt(KFcov(2,2)):0.2:KFm(2)+5*sqrt(KFcov(2,2))];
71         [X1,X2]=meshgrid(x1,x2);
72         Xr=[X1(:) X2(:)];
73         nXr=size(Xr,1);
74         fxr=zeros(nXr,1);
75         for j=1:n,
76             mvn=mvnpdf((Xr-ones(nXr,1)*PT(:,j))')/hn,zeros(1,d),eye(d));
77             fxr=fxr+1/(n*hn^d)*mvn;
78         end;

```

```
79 colormap([0 0 0]);
80 mesh(X1,X2,reshape(fxr,length(x2),length(x1)));
81 pause;
82 pr=mvnpdf(Xr,KFm',KFcov);
83 mesh(X1,X2,reshape(pr,length(x2),length(x1)));
84 end;
```

Bibliography

- Ch. Andrieu, A. Doucet, and R. Holenstein. Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society B*, 72(3):269–342, 2010.
- P. Billingsley. *Probability and Measure*. Wiley-Interscience; 3rd edition, 1995.
- P. Billingsley. *Convergence of Probability Measures*. Wiley-Interscience; 2nd edition, 1999.
- J. Brabec and B. Hruža. *Matematická analýza II (Mathematical Analysis II, in Czech)*. SNTL/ALFA, 1986.
- D. Crisan and A. Doucet. A Survey of Convergence Results on Particle Filtering Methods for Practitioners. *IEEE Transactions on Signal Processing*, 50(3):736–746, 2002.
- A. Doucet, N. de Freitas, and N. Gordon (Eds.). *Sequential Monte Carlo Methods in Practice*. Springer-Verlag, New York, LLC, 2001.
- B. Fristedt, N. Jain, and N. Krylov. *Filtering and Prediction: A Primer*. American Mathematical Society, 2007.
- K. Heine and D. Crisan. Uniform Approximations of Discrete-Time Filters. *Advances in Applied Probability*, 40(4):979–1001, 2008.
- R. E. Kalman. A New Approach to Linear Filtering and Prediction Problems. *Transactions of the ASME—Journal of Basic Engineering*, 82(Series D):35–45, 1960.
- P. Lachout. *Teorie pravděpodobnosti (Theory of Probability, in Czech)*. Karolinum, Univerzita Karlova v Praze, 2004.
- P. Del Moral and A. Guionnet. On the Stability of Interacting Processes with Applications to Filtering and Genetic Algorithms. *Annales de l’institut Henri Poincaré (B) Probabilités et Statistiques*, 37(2):155–194, 2001.
- E. Parzen. On Estimation of a Probability Density Function and Mode. *Annals of Mathematical Statistics*, 33(3):1065–1076, 1962.
- D.S.G. Pollock. *Handbook of Time Series Analysis, Signal Processing, and Dynamics*. Academic Press, 1999.

- S. Resnick. *A Probability Path*. Birkhäuser, 1999.
- D. W. Scott. *Multivariate Density Estimation: Theory, Practice, and Visualization*. John Wiley & Sons, Inc., 1992.
- B. W. Silverman. *Density Estimation for Statistics and Data Analysis*. Chapman and Hall/CRC, London, New York, 1986.
- M. E. Tarter and M. D. Lock. *Model-Free Curve Estimation*. Chapman and Hall/CRC, London, New York, 1993.
- A. B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer, 2009.
- M. P. Wand and M. C. Jones. *Kernel Smoothing*. Chapman and Hall/CRC, London, New York, 1995.