

Wolf-Jürgen Beyn
Raphael Kruse

Numerical Methods for Stochastic Processes

– Vorlesungsskript –
Universität Bielefeld
Sommersemester 2011

August 26, 2011

Preface

These lecture notes grew out of a course *Numerical Methods for Stochastic Processes* that the authors taught at Bielefeld University during the summer term 2011. The text contains material for about 30 two-hour lectures and includes a series of exercises most of which were assigned during the course. We assume that readers/participants have a firm basis in measure theory and probability theory as usually provided during a one-semester course. To a lesser extent experience with basic numerical methods from a one-semester course is also very helpful.

The lecture notes address Bachelor students in their third year and Master students in their first year. For Bachelor students the topics may be taken as a basis for writing a Bachelor Thesis while for Master students they may serve as a starting point for a specialization in numerical methods for stochastic ordinary and partial differential equations.

Contents

1	Introduction	1
1.1	Financial Options	1
1.2	Asset Price Model	3
1.3	Black-Scholes Formula	8
1.4	Monte Carlo Methods for Financial Option Valuation	10
2	Preliminaries from Probability Theory	15
2.1	Probability Spaces and Random Variables	15
2.2	Independence and Distributions of Random Variables	17
2.3	Integrability and Moments of Random Variables	19
2.4	The Transformation Theorem for Integrals	23
2.5	Some Standard Distributions	24
2.6	Conditional Expectations	25
2.7	Limit Theorems	27
3	Generating Random Numbers	33
3.1	Motivation	33
3.2	Pseudo-Random Number Generators	34
3.3	Linear Congruential Generators	37
3.4	Empirical Tests	40
3.5	The Mersenne Twister	45
4	Generating Random Variables with Non-Uniform Distribution	49
4.1	Inversion Method	49
4.2	Rejection Method	50
4.3	The Box-Muller Method	53
4.4	Marsaglia's Ziggurat Method	55

5	Monte Carlo Methods	61
5.1	Statistical Analysis of Simulation Output	61
5.2	Monte Carlo Integration	65
5.3	Variance Reduction Techniques	67
5.4	Approximation of Multiple Integrals	78
6	Theory of Continuous Time Stochastic Processes and Itô-Integrals .	87
6.1	Continuous Time Stochastic Processes	87
6.2	Martingales	90
6.3	Brownian Motion	93
6.4	The Itô-Integral	103
6.5	Itô's Formula	114
7	Stochastic Ordinary Differential Equations	117
8	Numerical Solution of SODEs	119
9	Weak Approximation of SODEs	121
10	Monte Carlo Methods for SODEs	123
	References	125
	Index	129

Symbols and Acronyms

\emptyset	the empty set
Ω	set of elementary events
A^c	complement $\Omega \setminus A$ of a subset $A \subset \Omega$
$\mathbb{1}_A$	indicator function of the set A see (2.2)
\mathbb{R}	set of real numbers
\mathbb{R}^d	vector space of d -dimensional tuples (x_1, \dots, x_d) with $x_i \in \mathbb{R}, i = 1, \dots, d$
$\mathbb{R}^{m,d}$	vector space of real $m \times d$ -matrices
$\mathcal{B}(\mathbb{R}^d)$	Borel- σ -algebra on \mathbb{R}^d
$\overset{\circ}{M}$	interior of a subset $M \subset \mathbb{R}^d$
λ^d	Lebesgue measure on \mathbb{R}^d
$ x $	absolute value of x if $x \in \mathbb{R}$ or Euclidean norm of x if $x \in \mathbb{R}^d$
$\mathcal{M}^2([0, T])$	Banach space of continuous, square-integrable, \mathbb{R}^m -valued (\mathcal{F}_t) -martingales (also written as $\mathcal{M}^2([0, T]; \mathbb{R}^m)$)
$V_a^b(X)$	bounded variation of a stochastic process X on $[a, b]$, see (6.5)
$\langle X \rangle_{[a,b]}$	quadratic variation of a stochastic process X on $[a, b]$, see (6.6)
a.e.	almost everywhere, synonymous with a.s.
a.s.	almost surely, or with probability 1
i.i.d.	independent and identically distributed
c.d.f.	cumulative distribution function
CLT	central limit theorem
p.d.f.	probability density function
LCG	linear congruential generator
LLN	law of large numbers
ODE	ordinary differential equation
PDE	partial differential equation
PRNG	pseudo-random number generator
SODE	stochastic ordinary differential equation

Chapter 1

Introduction

In this chapter we present the basic ideas of the option valuation theory as a motivating example. Mathematical finance, and the Black-Scholes model in particular, is easily accessible and provides a range of typical and non-trivial applications of the different numerical methods which we will discuss later.

Here, we give an overview of the standard Black-Scholes model. At some places we already use terminology which will be introduced in later chapters. On first reading we recommend that the reader simply skip over unknown technical terms.

The content of this chapter is based on [12, 13, 14, 22, 32].

1.1 Financial Options

A financial derivative is a contract which value depends on the expected price movement of an underlying *asset*. An asset is anything tangible or intangible which can be owned and traded for cash. For example, this includes commodities such as oil or gold, as well as shares or bonds which are traded on the stock market.

In this chapter we are interested in a very specific financial derivative, the so called European call option, which we define in the same way as in [14].

Definition 1.1. A *European call option* gives its *holder* the right (but not the obligation) to purchase from the *writer* a prescribed asset S (the *underlying*) for a prescribed price $E > 0$ (the *exercise price*) at a prescribed time $T > 0$ (the *expiry date*) in the future.

For example, I (the writer) may offer you (the holder) the right to buy a share of the Apple Inc. (the underlying) for 300 USD (the exercise price) three months from now (expiry date). In three months there are two possible scenarios:

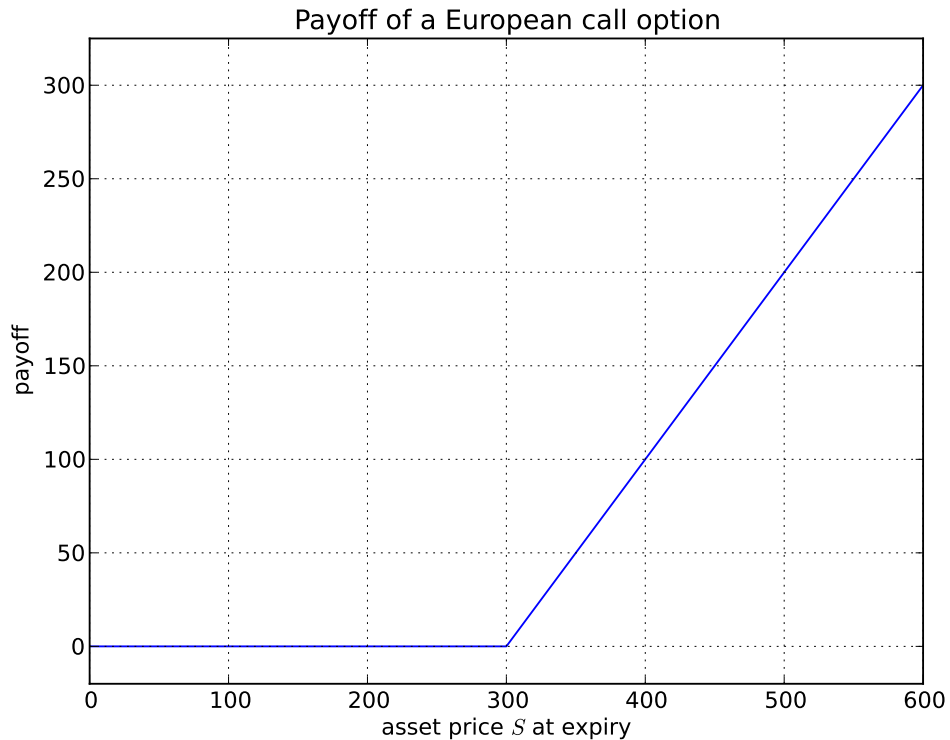


Fig. 1.1 The payoff diagram of the European call option with $E = 300$.

- a) The market price of the Apple Inc. share is higher than 300 USD. Then it makes sense for you to exercise the option and, if you immediately sell the share, you will gain the positive difference between the exercise price and the market price.
- b) If the market price is lower than 300 USD then you simply let the option expire and you may buy the share on the open stock market for the lower price.

In none of the scenarios, as it is indicated in Figure 1.1, you are loosing money but in case a) your gain is potentially unlimited. For me as the holder, however, I am facing a loss in scenario a) and no gain in b). So, in exchange for the option I may ask you for a compensation.

This chapter is devoted to answer the question how do we determine a *fair* value of this option. But before we follow in the footsteps of the Nobel prize winning paper by Black and Scholes [6] let us briefly note why this question is important.

On the financial markets options and other derivatives, such as swaps and futures, have become increasingly popular. In some cases more money is invested

in the derivatives than in the underlying asset. The two most common motivations for investors to buy options are *hedging* and *speculations*.

Investors, who buy options in order to hedge a risk, use them in the same way as an insurance policy. For example, consider an European investor who is planning to buy a factory in the USA which costs 10 million USD. The amount of money is payable in three months. If the investor is worried about a devaluation of the currency rate of the Euro he may consider to buy an options which gives him the right to buy 10 million USD for a fixed exchange rate in three month. Thus, the exercise price can be interpreted as a worst-case scenario.

On the other hand, investors may also buy options to speculate on price movements of the underlying assets. If the price of the underlying climbs by one percent the value of a corresponding European call usually climbs by a much larger amount. Therefore, the investor makes a larger profit by investing the same amount of money in the option instead of directly buying the underlying asset. Of course, the same holds true for possible losses if the price of the underlying asset moves in the wrong direction.

1.2 Asset Price Model

In this section we formulate our assumptions and derive a mathematical model for the price movements of the underlying asset.

Assumption 1.2 (Bank account). We postulate the existence of a risk-free bank account with continuously compounded *interest rate* $r \geq 0$. We are allowed to deposit or borrow any arbitrary amount of money at any time. The interest rate is assumed to be constant. By $B(t)$ we denote the balance of the bank account.

In practice, Assumption 1.2 is not satisfied for several reasons: The interest rate is usually not fixed and banks often charge a higher interest rate on credits than they pay on saving accounts. However, we interpret our assumption as an approximation of the reality on short time intervals.

If we put an amount of B_0 on the account at time t_0 then, at time $t_1 > t_0$, we have

$$B(t_1) = e^{r(t_1-t_0)} B_0. \quad (1.1)$$

Therefore, the balance process $B(t)$ follows the linear ordinary differential equation

$$\frac{d}{dt} B(t) = rB(t), \quad B(t_0) = B_0. \quad (1.2)$$

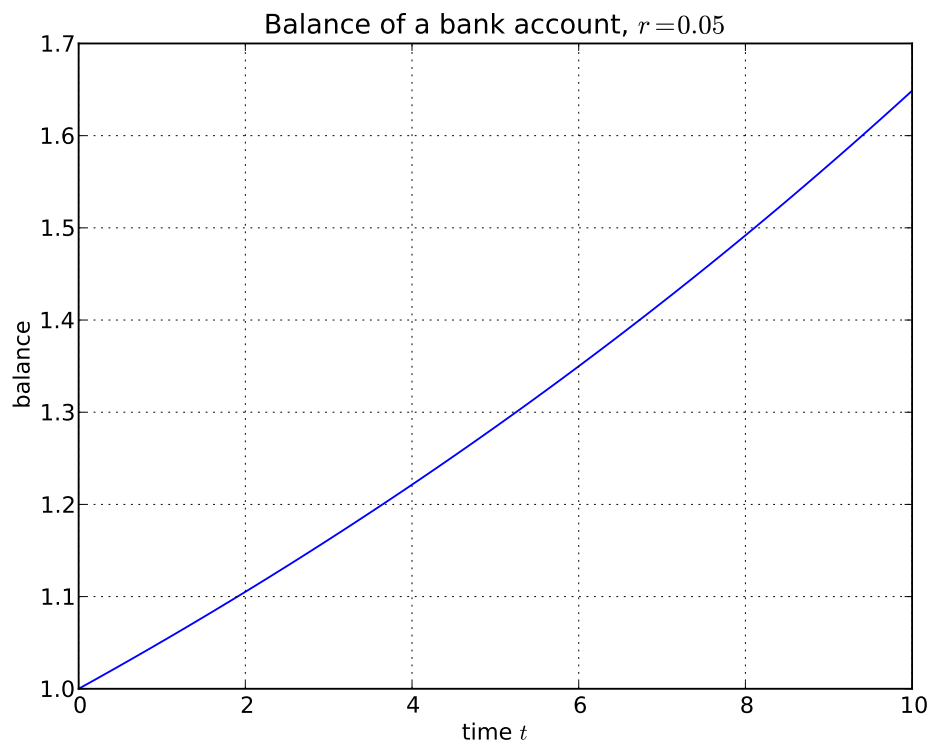


Fig. 1.2 Balance of a bank account with fixed continuously compounded interest rate $r = 0.05$ and initial capital of 1 euro. The scale of the time axis is measured in years.

A typical value of the parameter r is 0.05 which corresponds to an interest rate of 5% per year. Figure 1.2 illustrates the development of the balance $B(t)$ over ten years.

A consequence of our assumption is that two offers of

- a) 100 euros at time $t = 0$, or
- b) e^{rt} 100 euros at time $t > 0$

can be considered to be equal. In fact, by borrowing or investing the money both offers can be transformed into the other.

In a similar way, 100 euros at time $t > 0$ are worth $100e^{-rt}$ euros at time zero. This concept is called *discounting for interest* or *discounting for inflation*.

The following assumption provides a structural framework for our financial market.

Assumption 1.3. a) In addition to the bank account our financial market only consists of one risky asset. By $S(t)$ we denote the nonnegative value of one unit of the risky asset at time $0 \leq t \leq T$.

- b) It is possible to buy and sell any real number of units of the asset at the market price $S(t)$ at any time $0 \leq t \leq T$.
- c) There are no transaction costs and the asset is paying no dividends.
- d) Short selling is allowed, that is, it is possible to hold a negative amount of the asset.

The theory of Black and Scholes also builds on the following fundamental assumption which states the absence of *arbitrages*. Although there exists a rigorous definition of an arbitrage we follow [14] and only provide a verbal formulation of the idea.

Assumption 1.4. There is never an opportunity to make a risk-free profit that gives a greater return than that provided by the interest from the bank account.

For a moment, assume that Assumption 1.4 is violated and there exists an arbitrage opportunity, then investors would simply borrow cash from the bank and take advantage of the arbitrage on a large scale. By the forces of supply and demand this would affect the market prices or interest rates until the opportunity has vanished. Therefore, in practice, if arbitrage opportunities exist in an efficient and liquid financial market then they should be short lived.

After setting a framework of our financial market we introduce a model for the asset price movements. Which properties should we demand from this model? First of all, since $S(t)$ describes the evolution of a price process it is reasonable to force $S(t)$ to be nonnegative. Further, since the asset S represents the erratic dynamics of stock prices, S should statistically behave in a similar way as real stock market data. This leads to the idea to model S as a stochastic process.

As it was proposed by Black and Scholes in [6], we assume that $S(t)$, $0 \leq t \leq T$, is given as the solution to the *stochastic ordinary differential equation* (SODE)

$$dS(t) = \mu S(t) dt + \sigma S(t) dW(t), \quad S(0) = S_0, \quad (1.3)$$

where S_0 denotes the initial price at time $t = 0$. The parameter μ is called the average rate of growth of the asset price or simply the *drift parameter* and the number $\sigma > 0$ denotes the *volatility*, which measures the standard deviation of the returns. By $W : [0, T] \times \Omega \rightarrow \mathbb{R}$ we denote a real-valued standard *Wiener process* which is defined on the time interval $[0, T]$ and on a probability space (Ω, \mathcal{F}, P) . As it is common in stochastic analysis we omit the $\omega \in \Omega$ in the argument of W , that is, $W(t, \omega) = W(t)$.

The SODE (1.3) is not a differential equation in the usual sense. It is short hand for the integral equation

$$S(t) = S_0 + \int_0^t \mu S(\tau) d\tau + \int_0^t \sigma S(\tau) dW(\tau) \quad \text{for all } t \in [0, T]. \quad (1.4)$$

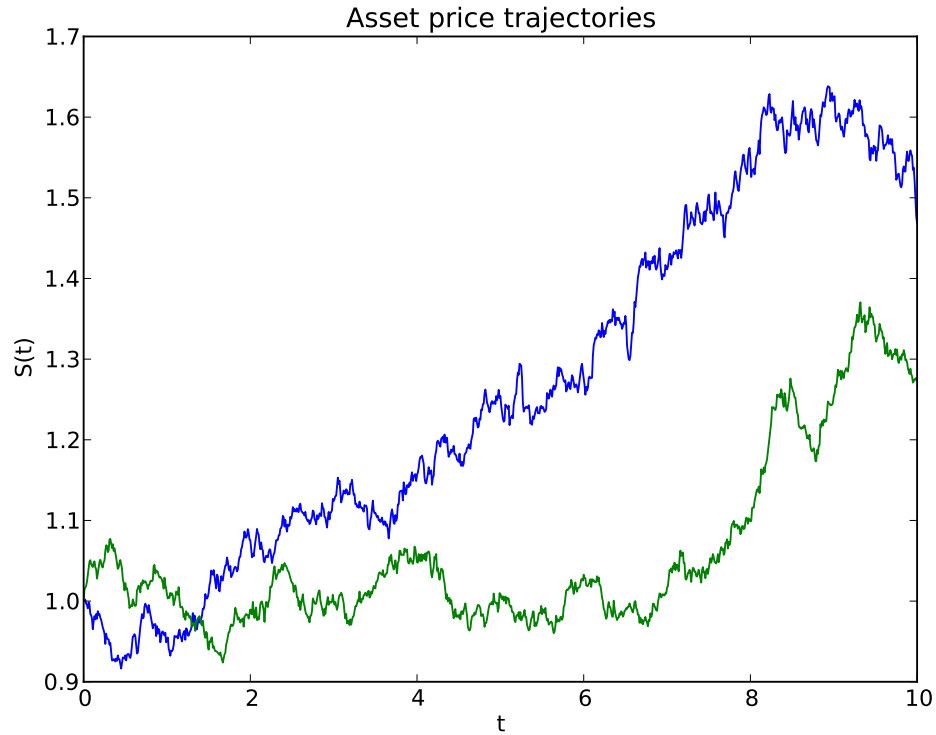


Fig. 1.3 Two typical realizations of the stochastic asset price process (1.5) with parameter values $T = 10$, $S_0 = 1$, $\mu = 0.05$, and $\sigma = 0.2$.

While the first integral is a well-known Lebesgue or Riemann-integral, the second integral is a so called stochastic Itô-integral, named after Kiyoshi Itô. The solution process $S: [0, T] \times \Omega \rightarrow \mathbb{R}$ to (1.3) is explicitly given by

$$S(t) = S_0 e^{(\mu - \frac{1}{2}\sigma^2)t + \sigma W(t)}. \quad (1.5)$$

Without using any knowledge on the Wiener process $W(t)$ we derive from (1.5) that $S(t)$ is in fact nonnegative.

For $\omega \in \Omega$ we call the mapping $t \mapsto S(t, \omega)$ a *sample path* of the stochastic process S . Figure 1.3 shows two typical sample paths of the asset price model (1.5).

For the definition of S we used terminology from stochastic analysis which we did not explain so far. This will be done in full detail in Chapter 6. For the rest of this section we aim to facilitate an intuitive understanding of the solution to (1.3).

Let $0 \leq t_0 < t_1 \leq T$ be arbitrary. Then by (1.4) we have

$$S(t_1) - S(t_0) = \int_{t_0}^{t_1} \mu S(\tau) d\tau + \int_{t_0}^{t_1} \sigma S(\tau) dW(\tau).$$

Next, under the assumption that the time length $t_1 - t_0$ is sufficiently small, we approximate both integrals by

$$\int_{t_0}^{t_1} \mu S(\tau) d\tau + \int_{t_0}^{t_1} \sigma S(\tau) dW(\tau) \approx \mu S(t_0)(t_1 - t_0) + \sigma S(t_0)(W(t_1) - W(t_0)).$$

By using this we get

$$\frac{S(t_1) - S(t_0)}{S(t_0)} \approx \mu(t_1 - t_0) + \sigma(W(t_1) - W(t_0)). \quad (1.6)$$

The left hand side of this equation is the *return* of the asset from time t_0 to t_1 . By the second summand in (1.6) we model the risk and uncertainty of the asset price movements. This is established by the Wiener increment $W(t_1) - W(t_0)$, which is a Gaussian random variable with mean zero and variance $t_1 - t_0$.

Doing the same steps for the balance of the bank account (1.1) yields

$$\frac{B(t_1) - B(t_0)}{B(t_0)} \approx r(t_1 - t_0). \quad (1.7)$$

Taking the expectation in (1.6) and comparing the result with (1.7) shows that the drift parameter μ in fact plays the same role for the asset price as the interest rate r for the bank account and can be interpreted as the average growth rate.

From (1.6) we also see that σ controls how strong the observed returns fluctuate around the average growth rate. So, if σ is very small the asset price movement will be dominated by the drift term and S behaves in a similar way as a bank account. On the other hand, a larger σ signals a larger influence of the stochastic disturbance and the asset price movements will be more erratic.

Finally, we want to indicate two reasons for the choice of a standard Wiener process W in the asset price model. First, we interpret the asset price movements as the sum of several independent decisions of a large number of market participants. Thus, by the central limit theorem, the observed price fluctuations should statistically behave like a normal random variable.

The second reason is that another property of Wiener processes fits well together with the *efficient market hypothesis*. This hypothesis states that the asset price responds immediately to any new information. Thus, a prediction of the future asset price cannot be improved if one also uses historical prices in addition to the present price. This corresponds very well to the independence of Wiener increments.

Remark 1.5. In practice, the asset price model (1.3) only gives useful approximations of the reality on very short time scales. Since the publication of the Black-Scholes formula in [6] more advanced asset price models have been developed. As a starting point we refer to [11].

1.3 Black-Scholes Formula

In this section we determine a mapping $C: [0, T] \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ such that $C(t, s)$ denotes the fair value of a European call at time t , if the call option expires at time T and the price of the underlying asset at time t is $S(t) = s$.

Further, by $E > 0$ we denote the exercise price. Then, by Definition 1.1, we already know that

$$C(T, s) = \max(0, s - E). \quad (1.8)$$

Under the asset price model and assumptions from Section 1.2 Black and Scholes [6] proved that the mapping C is the solution to the partial differential equation (PDE)

$$\frac{\partial C(t, s)}{\partial t} + \frac{1}{2} \sigma^2 s^2 \frac{\partial^2 C(t, s)}{\partial s^2} + rs \frac{\partial C(t, s)}{\partial s} - rC(t, s) = 0. \quad (1.9)$$

This is the so called *Black-Scholes PDE*. In Chapter 6 we will present a derivation of this equation when we have the Itô formula at our disposal.

Together with the final time condition (1.8) there exists a unique solution to the Black-Scholes PDE. Moreover, Black and Scholes also presented an explicit representation of the solution, the famous *Black-Scholes formula* for European call options.

Theorem 1.6 (Black-Scholes formula for European call options). *There exists a unique solution $C: [0, T] \times \mathbb{R}_+ \rightarrow \mathbb{R}_+$ to the Black-Scholes PDE (1.9) which satisfies the final time condition (1.8). The solution is explicitly given by*

$$C(t, s) = sF_{N(0,1)}(d_1) - Ee^{-r(T-t)}F_{N(0,1)}(d_2), \quad (1.10)$$

where $F_{N(0,1)}$ is the cumulative probability distribution of the standard normal distribution, that is

$$F_{N(0,1)}(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{1}{2}z^2} dz.$$

In addition, d_1 and d_2 are given by

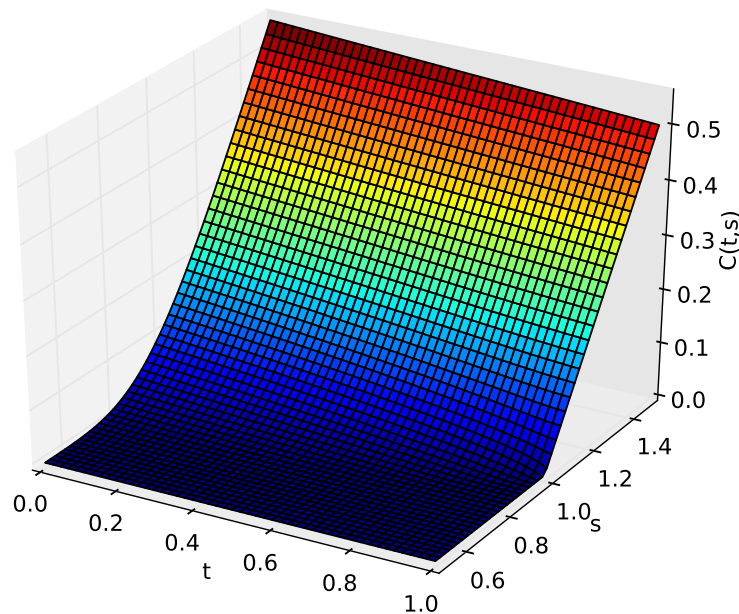


Fig. 1.4 Surface plot of the Black-Scholes formula for a European call option with fixed parameters values $T = 1$, $r = 0.05$, $\sigma = 0.05$, and $E = 1$.

$$d_1 = d_1(t, s) = \frac{\ln(s/E) + (r + \frac{1}{2}\sigma^2)(T - t)}{\sigma\sqrt{T - t}}$$

and

$$d_2 = d_2(t, s) = d_1 - \sigma\sqrt{T - t}.$$

A proof is given in, for instance, [22, Ch. 9]. Figure 1.4 shows a surface plot of (1.10) with fixed parameter values $T = 1$, $r = 0.05$, $\sigma = 0.05$, and $E = 1$.

Remarks 1.7. a) In order to compute the value of a European call we need to know the present asset price S_0 , the expiry time T , the exercise price E , the interest rate r , and the volatility σ . But the mapping C and the Black-Scholes PDE are independent of the value of the drift parameter μ .

b) In addition to the explicit representation of C , a further reason for the popularity of the Black-Scholes theory is that it also provides a portfolio trading strategy which can be used to replicate the payoff of a European call option. By simply

following this *hedging* strategy a bank can sell options without taking risks. For further reading we refer to [14].

1.4 Monte Carlo Methods for Financial Option Valuation

The Black-Scholes formula (1.10) gives an analytic solution to the problem of determining the fair value of a European call option. Therefore, we could consider the problem as being solved. However, as it is pointed out by D. Higham [13, 14], there are many variations of the option valuation problem, where a simple analytic solution does not exist.

For instance, this is true for so called exotic options [14], where the payoff not only depends on the final time asset price, but also on its behaviour during the time interval $[0, T]$. The same problem occurs when we invoke a different asset price model. For example, we refer to [22, Ch. 9.2] for several variations.

The aim of this section and, in fact, of the whole lecture is to present and analyze numerical methods for applications where analytic solutions are not available. Here, we focus on *Monte Carlo methods*. Roughly speaking, a Monte Carlo method tries to approximate the mean of a random variable X by the average over a large number N of independent outcomes of X , that is

$$\mathbf{E}[X] \approx \frac{1}{N} \sum_{i=1}^N X(\omega_i) \quad \text{for } N \text{ large.}$$

In the context of the option valuation problem, other numerical approaches contain numerical approximation of the solution to the Black-Scholes PDE (1.9) or the use of simplified asset price models, for example, the binomial method. For these approaches, we refer to [14, 12].

In the following we will again discuss the problem of the valuation of a European call option under the conditions of Section 1.2. The basic idea, how to apply a Monte Carlo method, is to use the *discounted expected payoff* as the option price, that is, the value is given by

$$V(S_0) = V(S_0; \mu, \sigma, T, E, r) = e^{-rT} \mathbf{E}[\max(0, S(T) - E)], \quad (1.11)$$

where, as above, S_0 is the present price of the asset, T and E denote the expiry date and exercise price of the option, and r is the riskless interest rate. The asset price $S(T)$ is given by the SODE (1.3) with drift parameter μ and volatility $\sigma > 0$.

This approach leads to two questions:

- a) What is the relationship between the discounted expected payoff V and the Black-Scholes formula (1.10)?

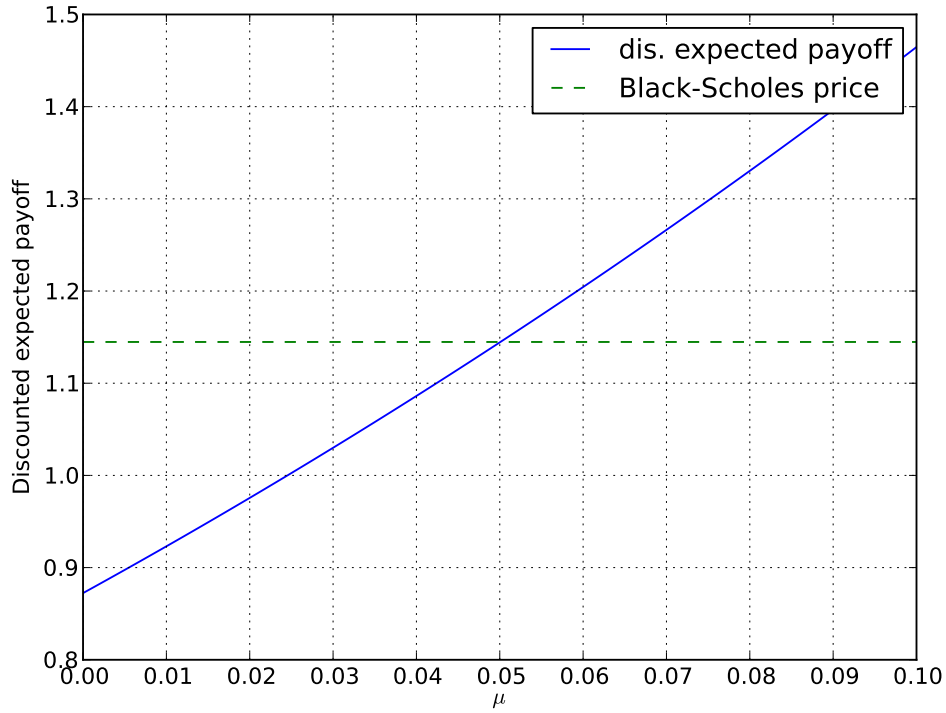


Fig. 1.5 Discounted expected payoff (1.11) for a European call option with varying parameter values for μ and with $S_0 = 2$, $T = 3$, $E = 1$, $r = 0.05$, and $\sigma = 0.25$. The dashed line is the corresponding option value of the Black-Scholes formula (1.10).

b) How do we compute the expected value in (1.11)?

Concerning a) we note that, in contrast to $C(0, S_0)$, the expected value $V(S_0)$ does depend on the parameter value of μ . Further, the value of V is derived without any reference to Assumption 1.4.

As it is illustrated in Figure 1.5 the value of V does indeed vary for different values of μ . Thus, by following [14, Ch. 12], it seems that two investors have to agree on the drift μ in order to use the discounted expected payoff for determining the value of a European call option, which contradicts the results from Section 1.3. Therefore, the ansatz of the discounted expected payoff seems to be of little use.

However, this can be resolved by the following observation: It is not a coincidence that the discounted expected payoff and the Black-Scholes value agree for $\mu = r$. In [14] this setting is known as the *risk neutrality assumption*.

A risk neutral investor is a person who is indifferent to two investments where the first offers a guaranteed rate of return r , and the second is a risky investment

with same expected rate of return r (Usually, investors are assumed to be risk-averse, which means that they will prefer the first investment).

In the case $\mu = r$ one can show that $\mathbf{E}[S(t)] = S_0 e^{rt}$ which coincides with the balance process of the bank account (1.1). Thus, risk neutral investors have no preferences between investing in the bank account and in the risky asset S .

Without going into details, it is possible to transform the underlying measure \mathbf{P} into a measure $\tilde{\mathbf{P}}$ in such a way that investors behave risk neutral with respect to $\tilde{\mathbf{P}}$. This is achieved by an application of Girsanov's Theorem [22, Ch. 8, Th. 2.2] and leads to the so called martingale approach to option valuation. For further reading we refer to [15].

Next, we come to question b). Following the considerations in [12, Ch. 5.1] the discounted expected payoff (1.11) can be approximated by processing the next two steps:

- 1) Compute N independent outcomes of the random variable $(S(T, \omega_i))_{i=1}^N$.
- 2) For each outcome, determine the payoff $\max(0, S(T, \omega_i) - E)$. The Monte-Carlo estimator of the discounted expected payoff is given by

$$e^{-rT} \mathbf{E}[\max(0, S(T) - E)] \approx e^{-rT} \frac{1}{N} \sum_{i=1}^N \max(0, S(T, \omega_i) - E).$$

While the implementation of step 2) is elementary, the theoretical background of the Monte Carlo estimator is provided in Chapter 5. There, we will discuss in which sense the right hand side is an approximation of the left hand side. We will discuss the order of convergence which usually is of the form $\mathcal{O}(\sqrt{N^{-1}})$ and show some techniques to accelerate convergence.

A larger focus lies on step 1). Under the risk neutrality assumption $\mu = r$ the asset price process $S(t)$ is given as the solution to the SODE

$$dS(t) = rS(t) dt + \sigma S(t) dW(t), \quad S(0) = S_0. \quad (1.12)$$

By (1.5) an analytic representation of the solution is given by

$$S(t) = S_0 e^{(r - \frac{1}{2}\sigma^2)t + \sigma W(t)}. \quad (1.13)$$

Thus, simulating $S(t)$ is easy if we know how to simulate the Wiener process $W(t)$. But, by the definition of a Wiener process (see Chapter 6), $W(t)$ is a $N(0, t)$ distributed random variable. Therefore, we can simulate the random variable $S(T)$ on a computer by

$$S(T, \omega_i) = S_0 \exp\left(\left(r - \frac{1}{2}\sigma^2\right)T + \sigma\sqrt{T}Z_i\right),$$

where the $(Z_i)_{i=1}^N$ are a sequence of $N(0, 1)$ distributed random numbers, which are produced by a random generator such as `randn` in Matlab. Chapters 3 and 4 will provide more details on pseudo-random number generators.

But we note, that this method still relies on the explicitly known analytic representation (1.13). Chapters 7 and 8 are concerned with numerical methods which are used to approximate the random variable $S(T)$ if a simple analytic solution is not available.

The easiest numerical method to approximate the solution to a SODE is the *Euler-Maruyama method* which for the SODE (1.12) is given by the recursion

$$\begin{aligned} S^j &= S^{j-1} + h\mu S^{j-1} + \sigma S^{j-1} (W(t_j) - W(t_{j-1})), \quad \text{for } j = 1, \dots, N_h, \\ S^0 &= S_0, \end{aligned}$$

where $0 = t_0 < t_1 < \dots < t_{N_h} = T$ is an equidistant partition of the time interval $[0, T]$ with step size $h = \frac{T}{N_h}$, $N_h \in \mathbb{N}$.

In order to simulate the increments $(W(t_j) - W(t_{j-1}))_{j=1}^{N_h}$ we again make use of the definition of the Wiener process which states that the increments are mutually independent $N(0, t_j - t_{j-1})$ -distributed random variables. Thus, an increment $W(t_j) - W(t_{j-1})$ can be simulated by $\sqrt{t_j - t_{j-1}}Z_j$, where Z_j is a $N(0, 1)$ -random number.

There exist two different concepts to analyze the error of the Euler-Maruyama method. The first concept, the so called *strong convergence*, compares the analytic solution $S(t_j)$ and the approximation S^j in an ω -wise fashion. That is we analyze the strong error

$$\left(\mathbf{E} \left[\max_{j=0, \dots, N_h} |S(t_j) - S^j|^2 \right] \right)^{\frac{1}{2}}.$$

As we will see in Chapter 7, the Euler-Maruyama method converges with strong order $\frac{1}{2}$.

The second concept is called *weak convergence* of the numerical method. Here, we analyze the error

$$\left| \mathbf{E}[\varphi(S(T))] - \mathbf{E}[\varphi(S^{N_h})] \right|,$$

where the real-valued function φ varies over a sufficiently large set of test functions.

The concept of weak convergence is closer related to our application of computing the discounted expected payoff. One result of Chapter 8 is that the Euler-Maruyama method converges with weak order 1.

Exercises

Problem 1.8. Show that the Black-Scholes solution $C(t, s)$, $0 \leq t < T$, satisfies the final time condition

$$\lim_{t \nearrow T} C(t, s) = \max(0, s - E), \quad \text{for all } s, E \geq 0.$$

Problem 1.9. The hedging strategy which enables banks to sell a European call without risks in the Black-Scholes model is known under the term *Delta-hedging*. It works as follows: Consider a European call with exercise price $E > 0$ and expiry date $T > 0$. At any time $t \in [0, T]$, one needs to have $\Delta(t)$ units of the underlying asset in order to replicate the payoff function of the call option. Here, $\Delta(t)$ is given by

$$\Delta(t) := \frac{\partial C}{\partial s}(t, S(t)),$$

where C denotes the Black-Scholes formula (1.10), $S(t)$ the asset price at time t . For fixed volatility $\sigma > 0$ and riskless interest rate $r > 0$ show that

$$\Delta(t) = F_{N(0,1)}\left(\frac{\log(S(t)/E) + (r + \sigma^2/2)(T - t)}{\sigma\sqrt{T - t}}\right).$$

Write a program which produces a surface plot of the function

$$(t, s) \mapsto \frac{\partial C}{\partial s}(t, s)$$

for $t \in [0, 1]$, $s \in [0, 2]$, $E = 1$, $r = 0.05$ and $\sigma = 0.2$.

Chapter 2

Preliminaries from Probability Theory

In this chapter we recall some basic results from measure and probability theory that will be used in the sequel. Readers familiar with the topic can safely skip this section or briefly read it for adapting to the notations used.

For measure and probability theory we use classical references such as [4, Kap. I-III], [3, Ch.II], [5] but there are many more books that cover the basic theory. In some instances the notions in [3],[4] differ from what has become standard . For example, distribution functions in [3] are taken to be left continuous, while the common use nowadays is to take them right continuous, see e.g. [28],[21].

The summary in this section will mainly follow the presentation in [28, Ch. 1 - 2],[21, Ch. 1],[22, Ch. 1].

After recalling the main notions of random variables, distribution functions and densities we turn to the concept of independent and identically distributed (i.i.d.) sequences of random variables. Then we discuss the transformation theorem for integrals which is one of the main tools for doing explicit calculations with distribution functions. Conditional expectations also play a dominant role in the theory of stochastic differential equations as well as in generating nonuniformly distributed random numbers. Finally, we write down the most important limit theorems in probability theory: the Law of Large Numbers (LLN) and the Central Limit Theorem (CLT). Both theorems lie at the heart of any convergence result for Monte Carlo simulation.

2.1 Probability Spaces and Random Variables

In this section we follow [28, Ch. 1] and [22, Ch. 1.2].

Probability theory provides mathematical models to analyze random phenomena. By Ω we denote the set of all possible *outcomes* and the typical elements $\omega \in \Omega$ are called *elementary events*. Usually, we are interested in combinations

of elementary events. If it is possible to determine if a given combination has occurred we call it an *event*.

More formally, we combine the set Ω with a family \mathcal{F} of subsets of Ω which satisfies

- (i) $\emptyset \in \mathcal{F}$, where \emptyset denotes the empty set,
- (ii) $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$, where $A^c = \Omega \setminus A$ is the complement of $A \subset \Omega$,
- (iii) $\{A_i\}_{i=1}^{\infty} \subset \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$.

A family \mathcal{F} of subsets of Ω with these properties is called a σ -*algebra*. The pair (Ω, \mathcal{F}) is called a *measurable space*, and an element of $A \in \mathcal{F}$ is called a *measurable set* or, simply, an *event*.

Frequently, we will encounter the measurable space $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$, where $\mathcal{B}(\mathbb{R}^d)$ denotes the *Borel- σ -algebra* on \mathbb{R}^d which is generated by all open sets in \mathbb{R}^d .

A mapping $\mathbf{P}: \mathcal{F} \rightarrow [0, 1]$ which satisfies

- (i) $\mathbf{P}(\Omega) = 1$,
- (ii) for any disjoint sequence $\{A_i\}_{i=1}^{\infty} \subset \mathcal{F}$, that is $A_i \cap A_j = \emptyset$ if $i \neq j$, we have

$$\mathbf{P}\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} \mathbf{P}(A_i),$$

is called a *probability measure* on (Ω, \mathcal{F}) and the triple $(\Omega, \mathcal{F}, \mathbf{P})$ is named a *probability space*.

Further, if $(\Omega, \mathcal{F}, \mathbf{P})$ is a probability space, we consider the family of subsets

$$\overline{\mathcal{F}} = \{A \subset \Omega : \exists B, C \in \mathcal{F} \text{ with } B \subset A \subset C, \mathbf{P}(B) = \mathbf{P}(C)\}. \quad (2.1)$$

Then $\overline{\mathcal{F}}$ is a σ -algebra and called the *completion* of \mathcal{F} . It is clear that \mathcal{F} is a sub- σ -algebra of $\overline{\mathcal{F}}$ and if $\mathcal{F} \neq \overline{\mathcal{F}}$ one can extend \mathbf{P} to $\overline{\mathcal{F}}$ by setting $\mathbf{P}(A) = \mathbf{P}(B) = \mathbf{P}(C)$ for all subsets $A \subset \Omega$ with $B \subset A \subset C$ for $B, C \in \mathcal{F}$. Then, $(\Omega, \overline{\mathcal{F}}, \mathbf{P})$ is called a *complete probability space*.

A function $X: \Omega \rightarrow \mathbb{R}$ is said to be \mathcal{F} - $\mathcal{B}(\mathbb{R})$ -*measurable* if

$$X^{-1}(A) \in \mathcal{F} \quad \text{for all } A \in \mathcal{B}(\mathbb{R}),$$

or, equivalently,

$$X^{-1}((-\infty, a]) = \{\omega \in \Omega : X(\omega) \leq a\} \in \mathcal{F} \quad \text{for all } a \in \mathbb{R}.$$

In this case, we say that X is a real-valued *random variable* on (Ω, \mathcal{F}) . A function $X: \Omega \rightarrow \mathbb{R}^d$ is called an \mathbb{R}^d -valued random variable on (Ω, \mathcal{F}) if it is \mathcal{F} - $\mathcal{B}(\mathbb{R}^d)$ -measurable. We remark that an \mathbb{R}^d -valued function is a random variable if and only if all components are real-valued random variables (see [3, §22, Rem. 2]).

For example, the indicator function $\mathbb{1}_A: \Omega \rightarrow \mathbb{R}$ which is given by

$$\mathbb{1}_A(\omega) = \begin{cases} 1, & \text{for } \omega \in A, \\ 0, & \text{for } \omega \in A^c, \end{cases} \quad (2.2)$$

is a random variable if and only if $A \in \mathcal{F}$.

Next, as in [3, §9], we compactify \mathbb{R} to $\overline{\mathbb{R}}$ by adding the points $-\infty$ and ∞ . The Borel- σ -algebra $\mathcal{B}(\overline{\mathbb{R}})$ consists of all sets of the form B , $B \cup \{-\infty\}$, $B \cup \{+\infty\}$, $B \cup \{-\infty, \infty\}$ with $B \in \mathcal{B}(\mathbb{R})$. A \mathcal{F} - $\mathcal{B}(\overline{\mathbb{R}})$ -measurable function $X: \Omega \rightarrow \overline{\mathbb{R}}$ is called a *numerical function*.

For a function $X: \Omega \rightarrow \mathbb{R}^d$ we define $\sigma(X)$ to be the smallest σ -algebra on Ω which contains all sets $X^{-1}(A) = \{\omega \in \Omega : X(\omega) \in A\}$ with $A \in \mathcal{B}(\mathbb{R}^d)$. Consequently, X is $\sigma(X)$ - $\mathcal{B}(\mathbb{R}^d)$ -measurable and we say $\sigma(X)$ is the σ -algebra generated by X .

Finally, let $(X(t))_{t \in \mathbb{T}}$ with $\mathbb{T} \subset \mathbb{R}$ be a family of \mathbb{R}^d -valued random variables, that is for all $t \in \mathbb{T}$ the mapping

$$\Omega \ni \omega \mapsto X(t, \omega) \in \mathbb{R}^d$$

is a random variable. The family $(X(t))_{t \in \mathbb{T}}$ is called a *stochastic process* on \mathbb{T} and for a fixed $\omega \in \Omega$ the mapping

$$\mathbb{T} \ni t \mapsto X(t, \omega) \in \mathbb{R}^d$$

is called a *sample path* of the process. As usual in stochastic analysis, we often suppress ω as an argument of a random variable.

We will come back to the theory of stochastic processes in Chapter 6.

2.2 Independence and Distributions of Random Variables

This section is based on [22, Ch. 1.2] and [4, §3].

We fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a measurable space (Ω', \mathcal{F}') . Consider a measurable function $X: \Omega \rightarrow \Omega'$. Then, for any $A' \in \mathcal{F}'$ the number

$$\mathbf{P}_X(A') := \mathbf{P}(X^{-1}(A')) = \mathbf{P}(\{\omega \in \Omega : X(\omega) \in A'\}) \quad (2.3)$$

is well-defined. In fact, the mapping $\mathcal{F}' \ni A' \mapsto \mathbf{P}(X^{-1}(A'))$ is a probability measure on (Ω', \mathcal{F}') , the *image measure induced by X* (see [3, §7]). In probability theory the image measure \mathbf{P}_X is also denoted by $X \circ \mathbf{P}$ and called the *distribution of the random variable X* .

If $(\Omega', \mathcal{F}') = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ the distribution of X is completely characterized by the *cumulative distribution function* $F: \mathbb{R} \rightarrow [0, 1]$ (*c.d.f.* for short), which is given by

$$F(x) = \mathbf{P}(X \leq x) = \mathbf{P}_X((-\infty, x]) \quad \text{for all } x \in \mathbb{R}. \quad (2.4)$$

The cumulative distribution function is increasing and right-continuous and satisfies $\lim_{x \rightarrow \infty} F(x) = 1$, $\lim_{x \rightarrow -\infty} F(x) = 0$. Conversely, any function with these properties generates a probability measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ (see [31] and note that the reference [3, Theorem 6.6] uses left-continuous distribution-functions with $(-\infty, x)$ instead of $(-\infty, x]$ in (2.4)).

Now, we come to the very important concept of independent random variables. For a formal definition we refer to [4, §7]. Here we will work with the following characterization [4, Th. 7.2].

Theorem 2.1. *Consider a finite family of measurable spaces $(\Omega_i, \mathcal{F}_i)$, $i = 1, \dots, n$. A finite family $(X_i)_{i=1}^n$ of random variables $X_i: \Omega \rightarrow \Omega_i$ is independent if and only if*

$$\mathbf{P}(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n \mathbf{P}(X_i \in A_i) \quad (2.5)$$

for all $A_i \in \mathcal{F}_i$, $i = 1, \dots, n$.

As it is proved in [4, §7] it is enough to show (2.5) for all sets A_i from a \cap -stable generator of the σ -algebra \mathcal{F}_i . For example, let $(\Omega_i, \mathcal{F}_i) = (\mathbb{R}, \mathcal{B}(\mathbb{R}))$ for all $i = 1, \dots, n$. Then a finite family of real-valued random variables $(X_i)_{i=1}^n$ is independent if and only if (2.5) holds for all half-open intervals $A_i = [a_i, b_i)$.

Further, in the situation of Theorem 2.1 consider the mapping $Y: \Omega \rightarrow \Omega_1 \times \dots \times \Omega_n$ which is given by

$$Y(\omega) = (X_1(\omega), \dots, X_n(\omega)).$$

Then Y is a \mathcal{F} - $\bigotimes_{i=1}^n \mathcal{F}_i$ -measurable random variable, where $\bigotimes_{i=1}^n \mathcal{F}_i$ denotes the product- σ -algebra. The distribution \mathbf{P}_Y of Y is called the *joint distribution of $(X_i)_{i=1}^n$* . The following characterization is a consequence of Theorem 2.1 and proved in [4, Th. 7.5].

Theorem 2.2. *A finite family of random variables $(X_i)_{i=1}^n$ is independent if and only if*

$$\mathbf{P}_Y = \mathbf{P}_{X_1} \otimes \dots \otimes \mathbf{P}_{X_n}.$$

Now we turn to the independence of infinitely many random variables $(X_i)_{i=1}^{\infty}$ which map into measurable spaces $(\Omega_i, \mathcal{F}_i)$. As in [4, §7] we say that the family $(X_i)_{i=1}^{\infty}$ is independent if and only if every choice of finitely many random variables $(X_i)_{i \in I}$ with $I \subset \mathbb{N}$ is independent. By [4, Th. 9.4] the statement of Theorem 2.2 stays valid for the case $n = \infty$.

We also recall the following useful result.

Theorem 2.3. *Let $(X_i)_{i \in I}$ be a finite or infinite family of random variables with values in measurable spaces $(\Omega_i, \mathcal{F}_i)$. Consider measurable mappings $Y_i: \Omega_i \rightarrow \Omega'_i$. Then the family of random variables $(Z_i)_{i \in I}$ given by $Z_i := Y_i \circ X_i$ is independent.*

For the proof we refer to [4, Th. 7.4].

We call a family of random variables $(X_i)_{i \in I}$ with $I \subset \mathbb{N}$ *independent and identically distributed*, for short *i.i.d.*, if the family is independent and $\mathbf{P}_{X_i} = \mathbf{P}_{X_j}$ for all $i, j \in I$.

In the following we often assume that an i.i.d. sequence $(X_i)_{i=1}^{\infty}$ of random variables is given. In probability theory this is a common assumption and it is easy to construct a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a sequence $(X_i)_{i=1}^{\infty}$ of measurable mappings such that the $(X_i)_{i=1}^{\infty}$ are i.i.d. with an arbitrary target distribution \mathbf{P}_X . For this realization problem we refer to [4, §9].

2.3 Integrability and Moments of Random Variables

As above, we fix a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Consider a real-valued random variable $X: \Omega \rightarrow \mathbb{R}$. We say that X is *integrable* with respect to the probability measure \mathbf{P} if the integral

$$\mathbf{E}[X] := \mathbf{E}_{\mathbf{P}}[X] := \int_{\Omega} X(\omega) d\mathbf{P}(\omega)$$

exists. In this case we call $\mathbf{E}[X]$ the *expectation* or *mean value* of X .

Following [3, Th. 12.2] we have that X is integrable if and only if $\mathbf{E}[|X|] < \infty$. By $\mathcal{L}^1(\Omega) := \mathcal{L}^1(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R})$ we denote the set of all integrable real-valued random variables on $(\Omega, \mathcal{F}, \mathbf{P})$.

Since for $\alpha \in \mathbb{R}$ and $X, Y \in \mathcal{L}^1(\Omega)$ it also holds that $(\alpha X), (X + Y) \in \mathcal{L}^1(\Omega)$ with

$$\mathbf{E}[\alpha X] = \alpha \mathbf{E}[X], \quad \text{and} \quad \mathbf{E}[(X + Y)] = \mathbf{E}[X] + \mathbf{E}[Y],$$

we obtain that $\mathcal{L}^1(\Omega)$ is a vector space.

Further, we have the inequalities (see [3, Th. 12.4])

$$|\mathbf{E}[X]| \leq \mathbf{E}[|X|]$$

and for all $X, Y \in \mathcal{L}^1(\Omega)$ with $X \leq Y$ it holds that

$$\mathbf{E}[X] \leq \mathbf{E}[Y].$$

Therefore, the mapping $X \mapsto \mathbf{E}[X]$ is an isotone linear form on $\mathcal{L}^1(\Omega)$.

Given $p \geq 1$ we say that X is *p-integrable* if $|X|^p$ is integrable, that is

$$\mathbf{E}[|X|^p] = \int_{\Omega} |X(\omega)|^p d\mathbf{P}(\omega) < \infty.$$

The value $\mathbf{E}[|X|^p]$ is called the *p-th moment* of X . Note that the set $\mathcal{L}^p(\Omega)$ of all random variables X with existing *p-th moment* forms a subspace of $\mathcal{L}^1(\Omega)$. In the case $p = 2$ we also say that X is *square-integrable* with respect to \mathbf{P} .

As in [3, §14] we assign the seminorms

$$N_p(X) := \left(\int_{\Omega} |X(\omega)|^p d\mathbf{P}(\omega) \right)^{\frac{1}{p}}$$

to the spaces $\mathcal{L}^p(\Omega)$ for $p \geq 1$. In particular, the seminorm satisfies

$$N_p(\alpha X) = |\alpha| N_p(X) \quad \text{for all } \alpha \in \mathbb{R}, X \in \mathcal{L}^p(\Omega),$$

and *Minkowski's inequality*

$$N_p(X + Y) \leq N_p(X) + N_p(Y) \quad \text{for all } X, Y \in \mathcal{L}^p(\Omega).$$

A further important inequality is *Hölder's inequality*

$$N_1(XY) \leq N_p(X) N_q(Y) \quad \text{for all } X \in \mathcal{L}^p(\Omega), Y \in \mathcal{L}^q(\Omega),$$

with $p, q > 1$, $\frac{1}{p} + \frac{1}{q} = 1$. A generalized version of Hölder's inequality is found in [1, Lem. 1.16]: For $n \in \mathbb{N}$ and $X_i \in L^{p_i}(\Omega)$, $i = 1, \dots, n$, with $p_i \in [1, \infty)$ and $r \in [1, \infty)$ which satisfy

$$\frac{1}{r} = \sum_{i=1}^n \frac{1}{p_i}$$

we have

$$N_r\left(\prod_{i=1}^n X_i\right) \leq \prod_{i=1}^n N_{p_i}(X_i).$$

The seminorm N_p turns into a norm if we identify random variables which coincide almost surely. To be more precise, we say that two random variables X and Y are equal P -almost surely or with probability 1, if there exists a measurable set $N \in \mathcal{F}$ with $\mathbf{P}(N) = 0$ such that

$$X(\omega) = Y(\omega) \quad \text{for all } \omega \in \Omega \setminus N.$$

For short we write $X = Y$ a.s.

Since it holds that

$$N_p(X) = 0 \quad \Leftrightarrow \quad |X|^p = 0 \text{ a.s.} \quad \Leftrightarrow \quad |X| = 0 \text{ a.s.} \quad \Leftrightarrow \quad X = 0 \text{ a.s.},$$

the set

$$\mathcal{N} = N_p^{-1}(0)$$

is a linear subspace of $L^p(\Omega)$. As noted in [3, §15] the subspace \mathcal{N} is independent of p and the quotient vector space

$$L^p(\Omega) := \mathcal{L}^p(\Omega) / \mathcal{N}$$

is well-defined. By defining

$$\|\tilde{X}\|_{L^p(\Omega)} := N_p(X)$$

for an equivalence class $\tilde{X} \in L^p(\Omega)$ and an arbitrary element X of \tilde{X} we obtain a norm on $L^p(\Omega)$. In fact, $(L^p(\Omega), \|\cdot\|_{L^p(\Omega)})$ is a Banach space [3, Th. 15.7]. For probability spaces it holds that $L^p(\Omega) \subset L^q(\Omega) \subset L^1(\Omega)$ whenever $p \geq q \geq 1$.

Note that in our notation we usually make no difference between random variables $X \in \mathcal{L}^p(\Omega)$ and their equivalence classes $\tilde{X} \in L^p(\Omega)$.

Of special interest is the case $p = 2$. By setting

$$(X, Y) = \mathbf{E}(XY) = \int_{\Omega} X(\omega)Y(\omega) d\mathbf{P}(\omega) \quad \text{for } X, Y \in L^2(\Omega)$$

we obtain an inner product and the space $L^2(\Omega)$ becomes a Hilbert space.

Another useful inequality is *Jensen's inequality*. Let $J \subset \mathbb{R}$ denote an interval which contains the range of a random variable $X \in L^1(\Omega)$ and consider a convex function $\varphi: J \rightarrow \mathbb{R}$ such that $\varphi(X) \in L^1(\Omega)$. Then Jensen's inequality

$$\varphi(\mathbf{E}[X]) \leq \mathbf{E}[\varphi(X)]$$

holds. For a proof we refer to [5, p. 276] (see also Problem 2.22).

Next, we introduce the *variance* of a random variable. The variance of a real-valued random variable X is defined by

$$\text{var}(X) = \mathbf{E}[(X - \mathbf{E}[X])^2].$$

If $X \in L^2(\Omega)$ then $\text{var}(X) < \infty$. The square root of the variance is called the *standard deviation* of X and often used in statistics to measure the spread of X around its mean. A simple calculation shows

$$\text{var}(X) = \mathbf{E}[X^2] - (\mathbf{E}[X])^2. \quad (2.6)$$

If Y is another real-valued random variable, we call

$$\text{cov}(X, Y) = \mathbf{E}[(X - \mathbf{E}[X])(Y - \mathbf{E}[Y])]$$

the *covariance* of X and Y . If $\text{cov}(X, Y) = 0$ we say that X and Y are *uncorrelated*. In particular, if X and Y are independent then it follows that $\text{cov}(X, Y) = 0$. Moreover, if $(X_i)_{i=1}^n$ are pairwise uncorrelated random variables then it holds that (see [4, Th. 8.3])

$$\text{var}(X_1 + \dots + X_n) = \text{var}(X_1) + \dots + \text{var}(X_n). \quad (2.7)$$

This is due to the fact that

$$\text{var}(X + Y) = \text{var}(X) + 2\text{cov}(X, Y) + \text{var}(Y). \quad (2.8)$$

We conclude this section with a brief look at \mathbb{R}^d -valued random variables $X(\omega) = (X_1(\omega), \dots, X_d(\omega))$. The mean of X is given by

$$\mathbf{E}[X] = (\mathbf{E}[X_1], \dots, \mathbf{E}[X_d]).$$

Following [4, §30] one defines the covariance of X by

$$\text{cov}(X) = \mathbf{E}[(X - \mathbf{E}[X])(X - \mathbf{E}[X])^T] \in \mathbb{R}^{d,d},$$

that is, $\text{cov}(X)$ is a symmetric matrix with entries $\text{cov}(X_i, X_j)$. In fact, one can show that C is always negative semidefinite.

In the same way as above, vector valued random variables give rise to a scalar of Banach spaces $L^p(\Omega; \mathbb{R}^d)$ with norm

$$\|X\|_{L^p(\Omega; \mathbb{R}^d)} := \left(\int_{\Omega} \|X(\omega)\|^p dP(\omega) \right)^{\frac{1}{p}},$$

where $\|\cdot\|$ denotes the Euclidean norm in \mathbb{R}^d .

2.4 The Transformation Theorem for Integrals

In this section we present some useful versions of the well-known substitution rule which we first state for general measure spaces.

The first version is concerned with integration with respect to an image measure. For this consider a measure space $(\Omega, \mathcal{F}, \mu)$, a measurable space (Ω', \mathcal{F}') and an \mathcal{F} - \mathcal{F}' measurable mapping $T: \Omega \rightarrow \Omega'$. The mapping T induces an image measure (see (2.3) or [3, Th. 7.5])

$$\mu' := \mu_T$$

on the measurable space (Ω', \mathcal{F}') .

Theorem 2.4. *Let f' be a numerical function on Ω' . Then the μ_T -integrability of f' is equivalent to the μ -integrability of $f' \circ T$. In case of integrability it holds that*

$$\int_{\Omega'} f' d\mu_T = \int_{\Omega} f' \circ T d\mu.$$

For the proof we refer to [3, §19]. In the case of a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ and a random variable X which takes values in a measurable space (Ω', \mathcal{F}') a probabilistic version of Theorem 2.4 reads as follows:

Theorem 2.5. *Let f' be a numerical function on Ω' . Then the \mathbf{P}_X -integrability of f' is equivalent to the \mathbf{P} -integrability of $f' \circ X$. In case of integrability it holds that*

$$\mathbf{E}_{\mathbf{P}_X}[f'] = \mathbf{E}_{\mathbf{P}}[f' \circ X].$$

In particular, for a real-valued random variable X the \mathbf{P}_X -integrability of the mapping $x \mapsto x$ is equivalent to the \mathbf{P} -integrability of X and we have

$$\mathbf{E}[X] = \int_{\mathbb{R}} x d\mathbf{P}_X(x).$$

The next theorem turns to Lebesgue integrals and is often called the *general transformation theorem for integrals*.

Theorem 2.6. *Let G, G' be open subsets of \mathbb{R}^d , and $\Phi: G \rightarrow G'$ a C^1 -diffeomorphism of G onto G' . A numerical function f' on G' is λ^d -integrable if and only if the function $(f' \circ \Phi)|\det D\Phi|$ is λ^d -integrable over G , and in this case*

$$\int_{G'} f' d\lambda^d = \int_G (f' \circ \Phi)|\det D\Phi| d\lambda^d.$$

A proof can be found in [3, §19] or [2, Th. 8.4].

2.5 Some Standard Distributions

In this section we focus on random variables which take values in \mathbb{R}^d . In particular, we are interested in random variables X whose distribution \mathbf{P}_X is *absolutely continuous* with respect to the Lebesgue measure λ^d on \mathbb{R}^d , that is $\mathbf{P}_X(A) = 0$ for all $A \in \mathcal{B}(\mathbb{R}^d)$ with $\lambda^d(A) = 0$.

In this case, by the Radon-Nikodym theorem [3, §17], there exists a measurable function $f: \mathbb{R}^d \rightarrow [0, \infty)$ with

$$\mathbf{P}_X(A) = \int_A f(x) d\lambda^d(x) \quad \text{for all } A \in \mathcal{B}(\mathbb{R}^d). \quad (2.9)$$

The function f is called the *probability density function* (*p.d.f.* for short) of X and uniquely determined λ^d -almost surely. In addition to (2.9), a $\mathcal{B}(\mathbb{R}^d)$ -measurable function $h: \mathbb{R}^d \rightarrow \mathbb{R}$ is integrable with respect to \mathbf{P}_X if and only if hf is Lebesgue integrable, and in this case we have

$$\int_{\mathbb{R}^d} h(x) d\mathbf{P}_X(x) = \int_{\mathbb{R}^d} h(x)f(x) dx. \quad (2.10)$$

For $d = 1$, the cumulative distribution function $F: \mathbb{R} \rightarrow [0, 1]$ of X satisfies

$$F(x) = \mathbf{P}(X \in (-\infty, x]) = \int_{-\infty}^x f(y) dy.$$

The following examples are well-known standard distributions.

Example 2.7 (Uniform distribution). For $a, b \in \mathbb{R}$ with $a < b$ we say that a random variable X is *uniformly distributed* on the interval $[a, b]$ if it has the probability density function

$$f(x) = \begin{cases} \frac{1}{b-a}, & x \in [a, b], \\ 0, & \text{otherwise.} \end{cases}$$

For short, we write $X \sim U(a, b)$.

Example 2.8 (Normal distribution). A *normal* or *Gaussian* random variable X with mean $\mu \in \mathbb{R}$ and variance σ^2 , $\sigma > 0$, has the probability density function

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right), \quad \text{for all } x \in \mathbb{R}.$$

For short, we write $X \sim N(\mu, \sigma^2)$. In fact, we have that $\mathbf{E}(X) = \mu$ and $\text{var}(X) = \sigma^2$.

We say that X is *standard normally distributed* if $\mu = 0$ and $\sigma = 1$. It holds that $\sigma X + \mu \sim N(\mu, \sigma)$ if $X \sim N(0, 1)$.

Example 2.9 (Normal distribution in \mathbb{R}^d). As in [4, §30] we say that an \mathbb{R}^d -valued random variable X is normally distributed, if for all linear forms $\ell: \mathbb{R}^d \rightarrow \mathbb{R}$, $\ell \neq 0$ there exist values $\mu_\ell \in \mathbb{R}$ and $\sigma_\ell > 0$ such that

$$\ell(X) \sim N(\mu_\ell, \sigma_\ell).$$

Set

$$\mu := \mathbf{E}[X] \in \mathbb{R}^d \quad \text{and} \quad C := \text{cov}(X) \in \mathbb{R}^{d,d}. \quad (2.11)$$

If X is normally distributed, then its distribution is uniquely determined by μ and C . For short, we write $X \sim N(\mu, C)$.

If C is invertible, then the density of X is given by

$$f(x; \mu, C) = (2\pi)^{-\frac{d}{2}} (\det(C))^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^t C^{-1}(x - \mu)\right) \quad (2.12)$$

for all $x \in \mathbb{R}^d$. While the proof that X has this distribution is somewhat advanced and uses Fourier transform [4, Satz 30.2] it is easy to show that a random variable with p.d.f. (2.12) satisfies (2.11), see Exercise 2.19.

Example 2.10 (Chi-square distribution). If Z_1, \dots, Z_r , $r \geq 1$, are independent and $N(0, 1)$ -distributed random variables, then the sum of their squares are *chi-square* (χ^2) distributed with r degrees of freedom, that is

$$X = \sum_{i=1}^r Z_i^2 \sim \chi_r^2.$$

For a proof we refer to Problem 2.21. The probability density function of the chi-square distribution is given by

$$f(x; r) = \begin{cases} 2^{-\frac{r}{2}} \Gamma(\frac{r}{2})^{-1} x^{\frac{r}{2}-1} e^{-\frac{1}{2}x}, & \text{if } x \geq 0, \\ 0, & \text{if } x < 0, \end{cases}$$

where Γ denotes the *Gamma function*. The chi-square distribution is a special case of the *gamma distribution* and often arises in statistical tests.

2.6 Conditional Expectations

In this section we briefly review the concept of *conditional expectations*. For the reader who is unfamiliar with this topic the probabilistic name is somewhat confusing since unlike the expectation of a random variable, the conditional expectation is in general not a real number but a random variable. Here we present two different approaches to define the conditional expectation.

As usual a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ is given. The first thing one may associate with the word “conditional” is perhaps the *conditional probability of $A \in \mathcal{F}$ under condition $B \in \mathcal{F}$ with $\mathbf{P}(B) > 0$* which is given by

$$\mathbf{P}(A|B) = \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(B)}.$$

This is the probability of the event A if we already know that the event B will occur. The concept of conditional expectations aims to generalize conditional probabilities to a family of conditions.

First, we follow [4, §15] and present the definition which makes use of the Radon-Nikodym theorem [3, §17]. Let $X \in L^1(\Omega)$ be given and consider a sub- σ -algebra $\mathcal{G} \subset \mathcal{F}$. In general \mathcal{G} contains significantly less events than \mathcal{F} and we cannot expect X to be \mathcal{G} -measurable. Therefore the question arises: How does X behave if we assume that only events from \mathcal{G} occur? To answer this question we look for a random variable Y which is measurable with respect to \mathcal{G} and satisfies

$$\mathbf{E}[\mathbb{1}_A X] = \int_A X \, d\mathbf{P} = \int_A Y \, d\mathbf{P} = \mathbf{E}[\mathbb{1}_A Y] \quad \text{for all } A \in \mathcal{G}. \quad (2.13)$$

We find Y by noting that the map $\mathcal{G} \ni A \mapsto \mathbf{E}[\mathbb{1}_A X]$ defines a signed measure on the measure space $(\Omega, \mathcal{G}, \mathbf{P}|_{\mathcal{G}})$ that is absolutely continuous with respect to $\mathbf{P}|_{\mathcal{G}}$. Then Y is its Radon-Nikodym density function which is unique \mathbf{P} -almost surely. We have

$$\int_A Y \, d\mathbf{P}|_{\mathcal{G}} = \int_A X \, d\mathbf{P} \quad \text{for all } A \in \mathcal{G}.$$

But for $A \in \mathcal{G}$ we have $\int_A Y \, d\mathbf{P}|_{\mathcal{G}} = \int_A Y \, d\mathbf{P}$ so that (2.13) follows. In the following we will always use the symbol \mathbf{P} when we integrate \mathcal{G} -measurable functions with respect to the restriction $\mathbf{P}|_{\mathcal{G}}$.

We say that Y is the *conditional expectation of X under the condition \mathcal{G}* and we use the notation

$$\mathbf{E}[X|\mathcal{G}] := Y.$$

If X is \mathcal{G} -measurable then X and Y coincide.

Before we discuss the properties of $\mathbf{E}[X|\mathcal{G}]$ we present an alternative way to define the conditional expectation. For this note that $L^2(\Omega, \mathcal{G}, \mathbf{P}; \mathbb{R})$ is a closed subspace of $L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R})$. Therefore, since $L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R})$ is a Hilbert space there exists the orthogonal projector $Q_{\mathcal{G}}$ onto $L^2(\Omega, \mathcal{G}, \mathbf{P}; \mathbb{R})$ which satisfies

$$\mathbf{E}[XZ] = (X, Z)_{L^2(\Omega)} = (Q_{\mathcal{G}}(X), Z)_{L^2(\Omega)} = \mathbf{E}[Q_{\mathcal{G}}(X)Z] \quad (2.14)$$

for all $X \in L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R})$, $Z \in L^2(\Omega, \mathcal{G}, \mathbf{P}; \mathbb{R})$. In particular, (2.14) holds for all $Z = \mathbb{1}_A$ with $A \in \mathcal{G}$ and, thus, $Q_{\mathcal{G}}(X)$ satisfies (2.13). By the uniqueness of the conditional expectation it follows that $\mathbf{E}[X|\mathcal{G}] = Q_{\mathcal{G}}(X)$.

Without going into details, by using the density of $L^2(\Omega)$ -functions in $L^1(\Omega)$ it is possible to extend the projector $Q_{\mathcal{G}}$ to functions in $L^1(\Omega)$ such that $Q_{\mathcal{G}}(X) = \mathbf{E}[X|\mathcal{G}]$ for all $X \in L^1(\Omega)$.

We conclude this section by stating useful properties of the conditional expectation. This list can be found in [22, Ch. 1.3]. For proofs we refer to [4, §15] and [5, Sec. 34].

$$\begin{aligned}
\mathcal{G} = \{\emptyset, \Omega\} &\Rightarrow \mathbf{E}[X|\mathcal{G}] = \mathbf{E}[X]\mathbb{1}_{\Omega}, \\
X \geq 0 &\Rightarrow \mathbf{E}[X|\mathcal{G}] \geq 0, \\
X \text{ is } \mathcal{G}\text{-measurable} &\Rightarrow \mathbf{E}[X|\mathcal{G}] = X, \\
X \equiv c &\Rightarrow \mathbf{E}[X|\mathcal{G}] = c, \\
a, b \in \mathbb{R} &\Rightarrow \mathbf{E}[aX + bY|\mathcal{G}] = a\mathbf{E}[X|\mathcal{G}] + b\mathbf{E}[Y|\mathcal{G}], \\
X \text{ is } \mathcal{G}\text{-measurable} &\Rightarrow \mathbf{E}[XY|\mathcal{G}] = X\mathbf{E}[Y|\mathcal{G}], \\
\mathcal{G}_1 \subset \mathcal{G}_2 \subset \mathcal{F} &\Rightarrow \mathbf{E}[\mathbf{E}[X|\mathcal{G}_2]|\mathcal{G}_1] = \mathbf{E}[X|\mathcal{G}_1].
\end{aligned}$$

Also useful is a *conditional version of Jensen's inequality*.

Lemma 2.11. *Let $J \subset \mathbb{R}$ denote an interval containing the range of $X \in L^1(\Omega)$. If $\varphi(X) \in L^1(\Omega)$ for a convex function $\varphi: J \rightarrow \mathbb{R}$ then it holds that*

$$\varphi(\mathbf{E}[X|\mathcal{G}]) \leq \mathbf{E}[\varphi(X)|\mathcal{G}].$$

For the proof we refer to [5, p. 449] (see also Problem 2.22).

Problem 2.23 asks the reader to investigate a link between the conditional probability and the conditional expectation.

2.7 Limit Theorems

Several fundamental theorems in probability theory describe the limit behavior of averages of independent and identically distributed random variables. We begin with the strong Law of Large Numbers (LLN). A proof of the following theorem is found in [4, Satz 12.1].

Theorem 2.12. *Let $(X_i)_{i \in \mathbb{N}}$ be a sequence of pairwise independent real valued random variables identically distributed with $\mathbf{E}(X_i) = \eta, i \in \mathbb{N}$. Then the following convergence holds almost surely,*

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \eta \quad \text{as } n \rightarrow \infty, \quad (2.15)$$

that is there exists a set $A \in \mathcal{F}$ with $\mathbf{P}(A) = 0$ such that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i(\omega) = \eta \quad \text{for } \omega \notin A. \quad (2.16)$$

Sometimes almost sure convergence is also denoted as convergence almost everywhere (a.e. for short). There are other notions than almost sure convergence that will play a role in the following. We list them in the following definition.

Definition 2.13. Let $Y_n, n \in \mathbb{N}$ and Y be random variables on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$. Then one says that Y_n converges to Y

- in L^p or in p -th mean with $p \geq 1$, if

$$\mathbf{E}(|Y_n - Y|^p) = \|Y_n - Y\|_{L^p}^p \rightarrow 0 \quad \text{as } n \rightarrow \infty,$$

- in probability, if for all $\varepsilon > 0$,

$$\mathbf{P}(|Y_n - Y| \geq \varepsilon) \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

- weakly (or in distribution), if for all continuous bounded functions $\varphi \in \mathbb{C}_b(\mathbb{R})$

$$\int_{\Omega} \varphi d\mathbf{P}_{Y_n} \rightarrow \int_{\Omega} \varphi d\mathbf{P}_Y \quad \text{as } n \rightarrow \infty.$$

The notion of convergence in distribution comes from the fact (see [3, Satz 30.13, §30 Aufgabe 7]) that weak convergence is equivalent to the statement that the c.d.f.'s F_n, F of Y_n, Y satisfy for all $x \in \mathbb{R}$ where F is continuous,

$$F_n(x) \rightarrow F(x) \quad \text{as } n \rightarrow \infty.$$

The relation between these various notions of convergence is illustrated by the following implications (cf. [4, §5])

$$\left. \begin{array}{l} L^p\text{-convergence} \implies L^1\text{-convergence} \implies \\ \text{almost sure convergence} \implies \end{array} \right\} \text{convergence in probability,} \\ \text{convergence in probability} \implies \text{weak convergence.} \quad (2.17)$$

The Central Limit Theorem (CLT) gives more information about the error of convergence in (2.15). The simplest version assumes an i.i.d. sequence of random variables with identical variance (see [5, Th. 27.1]).

Theorem 2.14. (De Moivre-Laplace, CLT) *Let $(X_i)_{i \in \mathbb{N}}$ be an i.i.d. sequence of square integrable random variables with expectation $\mathbf{E}(X_i) = \eta$ and variance $\text{var}(X_i) = \sigma^2$ for $i \in \mathbb{N}$. Then*

$$S_n = \frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n (X_i - \eta) \longrightarrow N(0, 1) \quad \text{as } n \rightarrow \infty \quad \text{in distribution.} \quad (2.18)$$

Since $N(0, 1)$ has continuous c.d.f. $F_{N(0,1)}$ the convergence of the c.d.f.'s F_{S_n} is uniform (cf. [3, Th.30.13]), i.e.

$$\|F_{S_n} - F_{N(0,1)}\|_\infty = \sup_{x \in \mathbb{R}} |F_{S_n}(x) - F_{N(0,1)}(x)| \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

Moreover, by the theorem of Berry and Esséen (see [4, eq. (28.23)]) one has an estimate of the form

$$\|F_{S_n} - F_{N(0,1)}\|_\infty \leq \frac{6}{\sigma^3\sqrt{n}} \mathbf{E}(|X_1 - \eta|^3), \quad (2.19)$$

provided the random variables have finite third moments. The order of convergence $\mathcal{O}(n^{-\frac{1}{2}})$ in this estimate cannot be improved in general.

Theorem 2.14 holds under much weaker assumptions on the random variables X_i than stated above, see the Lindeberg conditions in [4, §28]. It is only assumed that the X_i are independent with expectation $\eta_i = \mathbf{E}(X_i)$ and variance $\sigma_i^2 = \text{var}(X_i)$. The sum S_n in (2.18) is then replaced by

$$S_n = \frac{1}{s_n} \sum_{i=1}^n (X_i - \eta_i), \quad s_n^2 = \text{var}\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n \sigma_i^2,$$

and the Berry and Esséen estimate (2.19) generalizes to

$$\|F_{S_n} - F_{N(0,1)}\|_\infty \leq \frac{6}{s_n^3} \sum_{i=1}^n \mathbf{E}(|X_i - \eta_i|^3). \quad (2.20)$$

Finally, we also note the following multidimensional version of the central limit theorem (see [5, Th. 29.5]).

Theorem 2.15. *Let $X_i = (X_{i,1}, \dots, X_{i,d})$ be an i.i.d. sequence of square integrable random vectors with values in \mathbb{R}^d . Set $\mu := \mathbf{E}[X_i] \in \mathbb{R}^d$ and $C := \text{cov}(X_i) \in \mathbb{R}^{d,d}$ for $i \in \mathbb{N}$. Then*

$$S_n = \frac{1}{\sqrt{n}} \sum_{i=1}^n (X_i - \mu) \longrightarrow N(0, C) \quad \text{as } n \rightarrow \infty \quad \text{in distribution.}$$

Exercises

Problem 2.16. Prove the following statement: If X, Y, Z are independent random variables, then so are

- (i) $(X + Y)$ and Z ,
- (ii) XY and Z .

Problem 2.17. Given a random variable $X: \Omega \rightarrow \mathbb{R}$ with probability density function $f: \mathbb{R} \rightarrow [0, \infty)$, determine the probability density function of

- (i) $X + a$, for $a \in \mathbb{R}$,
- (ii) bX , for $b \neq 0$,
- (iii) $\exp(X)$,
- (iv) X^2 .

Problem 2.18. (i) Show that $\text{var}(X) = \mathbf{E}[X^2] - (\mathbf{E}[X])^2$.

(ii) Calculate the first and second moments and the variance of a random variable $X: \Omega \rightarrow \mathbb{N}_0$ with a Poisson distribution, i.e. for some $\lambda > 0$,

$$p_n = \mathbf{P}(X = n) = \frac{\lambda^n}{n!} \exp(-\lambda) \quad \text{for } n = 0, 1, \dots$$

Problem 2.19. Let $C \in \mathbb{R}^{d,d}$ be a symmetric, positive definite matrix and let X be an \mathbb{R}^d -valued random variable with density function

$$f(x; \mu, C) = (2\pi)^{-\frac{d}{2}} (\det(C))^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(x - \mu)^T C^{-1}(x - \mu)\right)$$

where $\mu \in \mathbb{R}^d$. Show that $\mathbf{E}[X] = \mu$ and $\text{cov}(X) = C$.

Hint: Substitute $y = C^{-\frac{1}{2}}(x - \mu)$.

Problem 2.20. Let $(X_i)_{i=1}^n$ be a finite family of i.i.d. $N(0, 1)$ random variables. For an orthogonal matrix $V \in \mathbb{R}^{n,n}$ consider the random vector $Y := VX$, where $X := (X_1, \dots, X_n)^T$. Show that the components Y_i , $i = 1, \dots, n$, of Y are also a finite family of i.i.d. $N(0, 1)$ random variables.

Hint: Use without a proof that a finite family of $N(0, 1)$ -distributed random variables is independent if and only if they are pairwise uncorrelated.

Problem 2.21. The *Gamma distribution* $\Gamma(a, b)$ with parameters $a, b > 0$ has the density function

$$f(x; a, b) = \begin{cases} b^a \frac{1}{\Gamma(a)} x^{a-1} e^{-bx}, & x > 0, \\ 0, & x \leq 0, \end{cases}$$

with $\Gamma(a) = \int_0^\infty t^{a-1} e^{-t} dt$. Prove the following statements:

(i) If $(X_i)_{i=1, \dots, n}$ are independent with $X_i \sim \Gamma(a_i, b)$, then

$$\sum_{i=1}^n X_i \sim \Gamma\left(\sum_{i=1}^n a_i, b\right).$$

Hint: Use that the density of a sum of two independent random variable is the convolution of their respective densities.

(ii) Let Z be a real-valued random variable with $Z \sim N(0, 1)$. Then $Z^2 \sim \Gamma\left(\frac{1}{2}, \frac{1}{2}\right)$.

(iii) Let $(Z_i)_{i=1}^n$ be independent and $N(0, 1)$ -distributed random variables. Then their sum is chi-square distributed, that is

$$\sum_{i=1}^n Z_i^2 \sim \chi_n^2.$$

Problem 2.22. Let $\varphi: \mathbb{R}^n \rightarrow \mathbb{R}$ be a two-times differentiable function such that the Hessian matrix $\text{Hess}(\varphi)(x)$ is nonnegative definite for all $x \in \mathbb{R}^n$.

(i) Show that

$$\varphi(x) \geq \varphi(y) + D\varphi(x)(y-x) \text{ and } \varphi\left(\frac{x+y}{2}\right) \leq \frac{1}{2}(\varphi(x) + \varphi(y))$$

for all $x, y \in \mathbb{R}^n$.

(ii) Prove *Jensen's inequality*, that is

$$\varphi(\mathbf{E}[X]) \leq \mathbf{E}[\varphi(X)]$$

for an arbitrary random variable $X: \Omega \rightarrow \mathbb{R}^n$.

(iii) For a given sub- σ -algebra $\mathcal{G} \subset \mathcal{F}$ prove the conditional version of Jensen's inequality

$$\varphi(\mathbf{E}[X|\mathcal{G}]) \leq \mathbf{E}[\varphi(X)|\mathcal{G}]$$

for an arbitrary random variable $X: \Omega \rightarrow \mathbb{R}^n$.

Problem 2.23. Let $(B_j)_{j=1}^n$ be a finite partition of Ω , that is

$$\bigcup_{j=1}^n B_j = \Omega, \quad B_j \in \mathcal{F}, \quad \mathbf{P}(B_j) > 0, \quad B_j \cap B_k = \emptyset \text{ for } j \neq k.$$

Set $\mathcal{G} = \sigma((B_j)_{j=1}^n)$. For $X \in L^1(\Omega)$ show that

$$\mathbf{E}[X|\mathcal{G}] = \sum_{j=1}^n \frac{\mathbf{E}[\mathbb{1}_{B_j} X]}{\mathbf{P}(B_j)} \mathbb{1}_{B_j}.$$

By using this, derive for $X = \mathbb{1}_A$ with $A \in \mathcal{F}$ that

$$\mathbf{E}[X|\mathcal{G}](\omega) = \mathbf{P}(A|B_j) \quad \text{for all } \omega \in B_j.$$

Chapter 3

Generating Random Numbers

In this section we describe different approaches to generate random numbers. In particular, we discuss some algorithms which produce *pseudo-random numbers*. The goodness of these algorithms is analysed through a set of statistical *tests*. At the end of this section we have a look at the Mersenne Twister, which is a widely used pseudo-random number generator.

3.1 Motivation

As it was pointed out in the introduction there exists a growing interest in modelling real world phenomenas which appear to be random. But before we can use a computer to get any insights from one of these models we are facing the very elementary problem that we work with a completely deterministic machine to model random behaviour. So we are in need of some source of randomness.

The first approaches to overcome this problem were built-in physical devices in computers which generated random numbers by atomic decay or cosmic ray counters. Another possibility are large databases of random numbers which were generated by real random phenomena.

But in practice, these solutions have several shortcomings. For example, in modern applications in finance it is important to recompute the value of stock options very fast after a significant change of one of the model parameters. Thus, it is too slow if our physical device gives only one random number every ten seconds or too expensive to buy a million of these devices to get sufficiently many random numbers in the desired time horizon.

In that regard databases are better suited. But here we have the problem that such a database may be too small or may take too much memory. A database with one billion random numbers takes already one gigabyte of memory if every random number takes exactly one byte.

Another issue is that one can think of several physical or chemical processes that could be used to generate random numbers. But they are only useful for statistical purposes if one exactly knows their distribution, which may vary over time.

In this section we mainly consider another approach, the so called pseudo-random number generator (PRNG). These algorithms produce numbers U_1, U_2, \dots which are completely deterministic but mimic a certain random behaviour. In particular, we will focus on generators whose outputs look like an independent and identically distributed sequence of $U(0, 1)$ random numbers. In today's practice, PRNGs are most often used in statistical applications.

Which random source one should choose in practice depends on the importance of the following criteria in the given application:

- Statistical properties,
- Speed and efficiency,
- Number of available random numbers,
- Reproducibility.

In any case one should always be careful about using results which are derived with the help of random number generators. As N. Madras points out in [21, Ch. 2.3] it is only recommended to use random number generators which have been tested thoroughly. In critical applications, it is worth to run simulations twice using different random number generators.

One may think of many more sources of randomness which are not mentioned here. As a starting point to this very active research field we refer to the discussion in [10, Ch. 1]. Let us finally mention that, despite considerable progress in actual computations, the question of a proper notion of a *random sequence* remains one of the fundamental problems in Mathematics, see the enlightening discussion in [16, Ch.3].

3.2 Pseudo-Random Number Generators

In this paragraph we loosely follow [9, Ch. 3.2.2] and [21, Ch. 2.2]. We give a definition of the class of generators which produce independent $U(0, 1)$ -distributed pseudo-random numbers and introduce some terminology. We conclude with three examples.

Since computers can only store values of finite accuracy it is natural to consider generators which produce random integers in a finite set $\{0, 1, \dots, M - 1\}$. Then the output X is transformed into a random number $U \in (0, 1)$, for example, by an auxiliary function which divides X by M .

More formally, we have the following

Definition 3.1. A *pseudo-random number generator* (PRNG) is given by a choice of a positive integer M , a function $h: \{0, 1, \dots, M-1\}^k \rightarrow \{0, 1, \dots, M-1\}$ and a deterministic recurrence relation

$$X_i = h(X_{i-k}, X_{i-k+1}, \dots, X_{i-1}) \quad i \geq k$$

for some fixed integer $k \geq 1$. The required initial vector (X_0, \dots, X_{k-1}) is called the *seed*.

So far, the definition does not require any statistical properties of the generated sequence $(X_i)_{i \geq k}$. This connection is established by the next definition.

Definition 3.2. A $U(0, 1)$ -PRNG is a pseudo-random number generator together with an auxiliary function $g: \{0, 1, \dots, M-1\} \rightarrow (0, 1)$ such that the sequence $(U_i)_{i \geq k} := (g(X_i))_{i \geq k}$ passes a set of tests which verify that the sequence $(U_i)_{i \geq k}$ has the same statistical properties as an independent and identically distributed sequence of $U(0, 1)$ random variables.

At this moment we stay somewhat vague about the set of statistical tests but we will be more specific about this point in Section 3.4. Let us first prove the rather obvious property that all PRNG cycle if we let them run long enough. The next lemma is taken from [9, Lemma 3.1].

Lemma 3.3. For every PRNG and every seed (X_0, \dots, X_{k-1}) there exist positive integers $\alpha, \beta \in \mathbb{N}$ such that $X_{i+\beta} = X_i$ for all $i \geq \alpha$.

Proof. Since there exist at most M different outcomes there are at most M^k different k -tuples. Thus, after M^k iterations some k -tuple (X_i, \dots, X_{i+k-1}) must have occurred before. Let $\alpha \in \mathbb{N}$ be such that the tuple $(X_\alpha, \dots, X_{\alpha+k-1})$ is the first to reappear and let $\beta \geq 1$ be the smallest number such that $(X_{\alpha+\beta}, \dots, X_{\alpha+\beta+k-1}) = (X_\alpha, \dots, X_{\alpha+k-1})$. Since the iterates are determined by the deterministic function h we find $X_{i+\beta} = X_i$ for all $i \geq \alpha$ by induction. \square

Note that we take minimal numbers α, β in the proof and that these satisfy $\alpha \in \{0, \dots, M^k - 1\}$, $\beta \in \{1, \dots, M^k\}$. We call the number $\beta = \beta(X_0, \dots, X_{k-1})$ the *period* of the PRNG given the seed (X_0, \dots, X_{k-1}) . The number

$$p := \inf_{(X_0, \dots, X_{k-1}) \in \mathcal{S}} \beta(X_0, \dots, X_{k-1})$$

is called the *period* of the PRNG, which is now independent of the seed. Here the infimum is taken over the set $\mathcal{S} \subset \{0, \dots, M-1\}^k$ of all valid seeds. The number M^k is the trivial upper bound of the period.

The following three examples are taken from [21, Ch. 2].

Example 3.4 (RANDU). The RANDU generator is intended to be an $U(0,1)$ -PRNG. This algorithm fits into our definition by setting $k = 1$, $M = 2^{31}$ and the function h is given by the recursion

$$X_{i+1} = 65539X_i \pmod{2^{31}}.$$

Here, mod is understood in the usual way: If l is an integer and m is a positive integer, then $l \pmod{m}$ is the unique $r \in \{0, \dots, m-1\}$ such that $l = nm + r$ for some integer n .

The auxiliary function is given by $g(x) = \frac{x}{M}$. Note that for the seed $X_0 = 0$ we have $X_i = 0$ for all i . Therefore, we set $\mathcal{S} = \{1, \dots, M-1\}$ as the set of all valid seeds.

The RANDU generator was used in the IBM 360/370 library for many years, although it has rather poor statistical properties. Nevertheless, it is an important historical example of a linear congruential generator (LCG). We will study LCGs in more detail in Section 3.3, where we also have a closer look at RANDU.

Example 3.5 (The Middle-Square Method). This generator was proposed by John von Neumann in 1949. Set $k = 1$, $M = 10^5$. The sequence $(X_i)_{i \geq 0}$ is defined as follows: Given is an integer $0 \leq X_i < 10^5$, that is, a number with up to four digits. Take the square of X_i and, if necessary, add leading zeros to obtain a number with exactly eight digits. Then X_{i+1} is the integer consisting of the middle four digits. For example, let $X_i = 6553$, then $X_i^2 = 42941809$ and $X_{i+1} = 9418$.

However, in practice the middle-square method is not a good pseudo-random number generator since it possesses a relatively short period, some fixed points (0, 100, 2500, 3792 and 7600) and short cycles (for example $540 \rightarrow 2916 \rightarrow 5030 \rightarrow 3009$) and, more seriously, many initial seeds converge to a fixed point or a short cycle (for example, if $X_i < 100$ then the resulting sequence will converge to the fixed point zero).

Example 3.6 (Fibonacci generator). For this generator we set $k = 2$ and M a large positive integer. The recursion is given by

$$X_{i+1} = (X_i + X_{i-1}) \pmod{M}.$$

As the other two examples, the Fibonacci generator does not give rise to a good $U(0,1)$ -PRNG. If we use the auxiliary function $g(x) = \frac{x}{M}$ and set $U_i = g(X_i)$, then $\mathbf{P}(U_i < U_{i+1} < U_{i-1}) = 0$ compared to $\frac{1}{6}$ for i.i.d. $U(0,1)$ -random variables. Compare further with Problem 3.13.

3.3 Linear Congruential Generators

In the following we focus on a specific and historically important class of PRNGs. The presented material is taken from [9, Ch. 3.2], [21, Ch. 2.2].

Definition 3.7. A *linear congruential generator* (LCG) is a pseudo-random number generator with $k = 1$ and the function $h : \{0, \dots, M - 1\} \rightarrow \{0, \dots, M - 1\}$ is given by

$$h(x) = (ax + c) \pmod{M},$$

where $a, c \in \mathbb{N}$ with $a > 0$. In this case we call the generator the (a, c, M) -LCG.

If $c = 0$ we say that the LCG is *multiplicative*.

In this form LCGs were first proposed by D. H. Lehmer in 1951 [20]. Since LCGs are relatively simple their statistical properties can be analyzed mathematically. Therefore, the decision if a given choice of parameters (a, c, M) produces a good PRNG does not only rely on passing a set of tests.

We already know an example: The RANDU generator in Example 3.4 is the multiplicative $(65539, 0, 2^{31})$ -LCG.

Typically, we want to have LCGs with huge periods and, hence, we have to consider M very large. The following theorem characterizes LCGs with maximal periods. For a proof we refer to [16, Ch. 3.2].

Theorem 3.8. *The period of the (a, c, M) -LCG is M if and only if the following conditions are satisfied*

- (i) c and M are relatively prime¹,
- (ii) every prime factor of M divides $a - 1$, and
- (iii) if 4 divides M , then 4 divides $a - 1$.

For example, if $M = 2^{31}$ then the (a, c, M) -LCG has period M if and only if c is odd and $a = 4n + 1$ for some $n \geq 1$. In particular, there exists no multiplicative LCG with period M . In case M is prime, Theorem 3.8 gives the conditions $c \neq 0 \pmod{M}, a = 1 \pmod{M}$, see Problem 3.14.

Although the last theorem is a very useful tool for finding parameter values (a, c, M) such that the resulting LCG possesses a long period there exists another issue which causes LCGs to have poor statistical properties. As in [21, Ch. 2.2] we demonstrate this for the RANDU generator for which d -tuples of consecutive outputs $(X_i, X_{i+1}, \dots, X_{i+d-1}), i \geq 0$, always lie on hyperplanes in \mathbb{R}^d -space.

For the RANDU generator first note that $a = 65539 = 2^{16} + 3$. Now we have

¹ Two nonnegative integers a and b are relatively prime if 1 is the greatest common divisor.

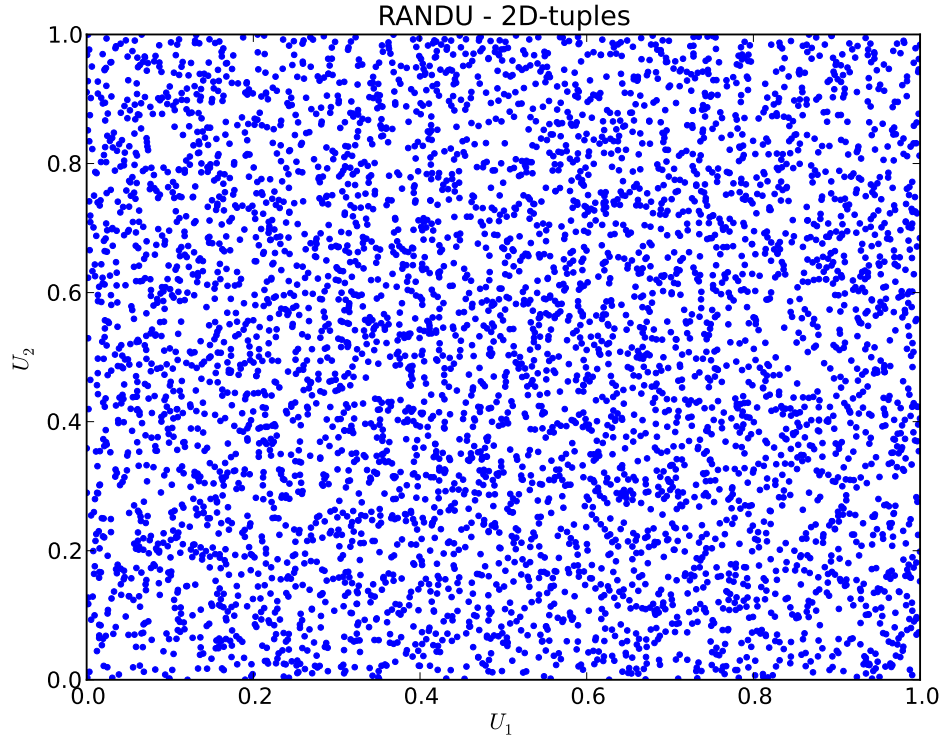


Fig. 3.1 Plot of (U_i, U_{i+1}) with $i = 0, \dots, 4999$ for the RANDU generator and seed $X_0 = 1$.

$$\begin{aligned}
 X_{i+2} &= (2^{16} + 3)X_{i+1} \pmod{2^{31}} \\
 &= (2^{16} + 3)^2 X_i \pmod{2^{31}} \\
 &= (0 + 6 \cdot 2^{16} + 9)X_i \pmod{2^{31}} \\
 &= (6(2^{16} + 3) - 9)X_i \pmod{2^{31}} \\
 &= (6X_{i+1} - 9X_i) \pmod{2^{31}}.
 \end{aligned}$$

This shows that $X_{i+2} - 6X_{i+1} + 9X_i$ is always divisible by 2^{31} and hence $X_{i+2} - 6X_{i+1} + 9X_i = n2^{31}$ for some $n \in \mathbb{Z}$. Since $1 \leq X_j < 2^{31}$ for all $j \geq 0$ it is also clear that $n \in \{-5, -4, \dots, 9\}$. Therefore, the triples (X_i, X_{i+1}, X_{i+2}) lie on one of at most 15 parallel hyperplanes in \mathbb{R}^3 .

Thus, if one knows the values of X_i and X_{i+1} , then the value of X_{i+2} is contained in a set of up to 15 integer values. This strongly contradicts the aim of generating pseudo-random numbers which behave as if they were independent. In that case the knowledge of X_i and X_{i+1} does not help to restrict the set of possible values of X_{i+2} .

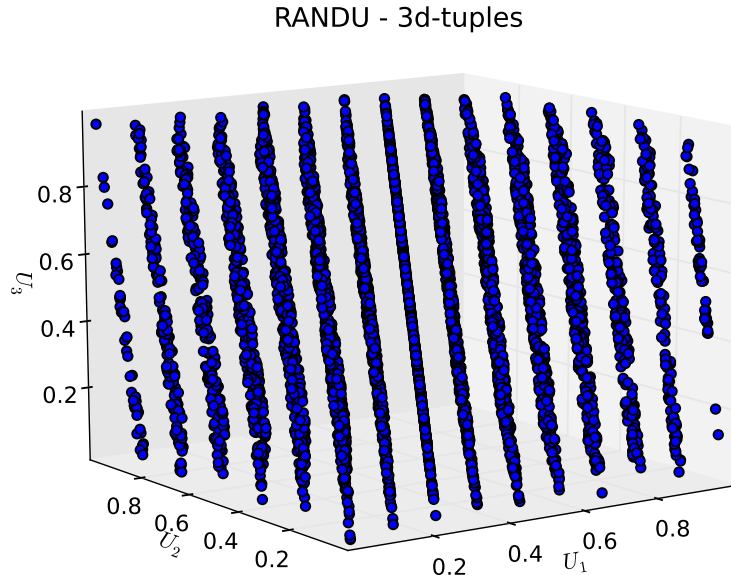


Fig. 3.2 Plot of (U_i, U_{i+1}, U_{i+2}) with $i = 0, \dots, 4999$ for the RANDU generator and seed $X_0 = 1$.

We illustrate this behaviour in Figures 3.1 and 3.2, where we drew the tuples (U_i, U_{i+1}) in Figure 3.1 and (U_i, U_{i+1}, U_{i+2}) in Figure 3.2 for $i = 0, \dots, 4999$ and $U_i = \frac{X_i}{M}$. For the naked eye Figure 3.1 looks as expected, that is, the tuples seem to be uniformly distributed over the unit-square. But in Figure 3.2 we see that the triples lie in parallel hyperplanes.

In fact, the same phenomenon can also be observed in the unit square if one uses sufficiently many tuples and zooms into a small subregion. We will further investigate this in Problem 3.16.

The following theorem, due to G. Marsaglia [23], shows that this behaviour is typical for all multiplicative LCGs. The proof of the first part is deferred to Problem 3.15.

Theorem 3.9. *Let the sequence $(X_i)_{i \geq 0}$ be generated by the $(a, 0, M)$ -LCG with seed $X_0 \in \{1, \dots, M - 1\}$. If $k_1, k_2, \dots, k_d \in \mathbb{Z}$ is any choice of integers such that*

$$k_1 + k_2 a + k_3 a^2 + \dots + k_d a^{d-1} = 0 \pmod{M},$$

then all points $\pi_i = (\frac{X_i}{M}, \dots, \frac{X_{i+d-1}}{M}) \in [0, 1)^d$, $i \geq 0$, lie in one of the hyperplanes defined by the equations

$$k_1x_1 + k_2x_2 + \dots + k_dx_d = 0, \pm 1, \pm 2, \dots$$

There are at most

$$|k_1| + |k_2| + \dots + |k_d|$$

of these hyperplanes which intersect the unit d -cube $[0, 1)^d$. Moreover, there is always a choice of k_1, k_2, \dots, k_d such that all of the points $(\pi_i)_{i \geq 0}$ fall in fewer than $(d!M)^{\frac{1}{d}}$ hyperplanes.

For the RANDU example with $a = 2^{16} + 3$ we have

$$a^2 - 6a + 9 = 0 \pmod{2^{31}}.$$

Therefore, Theorem 3.9 guarantees that triples lie in at most 16 hyperplanes which is a good upper bound of the actual number 15 found above.

Because of this behaviour most LCGs are usually not recommended to be used in Monte-Carlo simulations. However, since they are easily implemented and analyzed LCGs are still considered in practice.

For further reading we refer to the discussions in [10, Ch. 1.2] and [11, Ch. 2.1]. In [9, Ch. 3.1] the authors give more details on how to choose the parameter value a to get a period of $M - 1$ for multiplicative LCGs.

3.4 Empirical Tests

Definition 3.2 incorporated the condition that a $U(0, 1)$ -PRNG should pass a set of statistical tests to verify its statistical properties. In this paragraph we describe how these tests are designed. Here we follow [21, Ch. 2.2] and [10, Ch. 2].

Goodness-of-fit Test

In a goodness-of-fit test we want to test the hypotheses that d -tuples formed from the output of a $U(0, 1)$ -PRNG is uniformly distributed in $(0, 1)^d$.

More formally, the test works in the following way:

- (i) Choose $n \geq 2$ and a partition of $[0, 1]^d$ into n disjoint subsets S_1, \dots, S_n with corresponding “volumes” $p_1, \dots, p_n > 0$.
- (ii) For a given number $N > 0$ of test vectors use the PRNG to generate N d -tuples $(U_1, \dots, U_d), \dots, (U_{(N-1)d+1}, \dots, U_{Nd})$.

(iii) Then compute

$$X_n^2(N) := \sum_{i=1}^n \frac{(\sigma_i - Np_i)^2}{Np_i},$$

where σ_i denotes the observed number of d -tuples which lie in S_i .

The test generalizes the idea of making a histogram to visualize the density of the random numbers and is based on a theorem, which is due to K. Pearson. In our situation the theorem reads as follows:

Theorem 3.10. *If the d -tuples are generated from independent and $U(0, 1)$ -distributed random variables, then it holds that for all $x \in \mathbb{R}$*

$$\mathbf{P}(X_n^2(N) \leq x) \rightarrow F_{n-1}(x) \quad \text{for } N \rightarrow \infty,$$

where F_{n-1} denotes the cumulative distribution function of the chi-square distribution with $n - 1$ degrees of freedom from Example 2.10.

For a proof we refer to [17, Th. 14.5].

Hence, if the d -tuples generated with data from a PRNG behave as if they are uniformly distributed in $(0, 1)^d$, then for large N the computed value $X_n^2(N)$ should also behave in the same way as a random variable that has an approximate chi-square distribution.

Therefore, large values of $X_n^2(N)$ indicate that the observed counts differ by large amounts from the expected counts. A possible decision rule is to reject the tested PRNG as a $U(0, 1)$ -PRNG if the value of $X_n^2(N)$ is in the upper 5% of the tail of the χ_{n-1}^2 distribution.

The test should be repeated for several choices of d , n , $(S_i)_{i=1, \dots, n}$, and N .

Quantile-Quantile Plot

A quantile-quantile plot is a graphical test which indicates if the random numbers generated by a PRNG follow the law of a given distribution. To some extent this can be seen as a variant of the goodness-of-fit test.

To be more precise, let \mathbf{P} be a probability measure on $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$ and consider the cumulative distribution-function $F_{\mathbf{P}} : \mathbb{R} \rightarrow [0, 1]$ given by

$$F_{\mathbf{P}}(x) = \mathbf{P}((-\infty, x])$$

(compare further with (2.4)).

Definition 3.11. For any $p \in (0, 1)$ the value

$$z(p) := \inf \{z \in \mathbb{R} \mid F_{\mathbf{P}}(z) \geq p\} > -\infty \quad (3.1)$$

exists and is called the p -th quantile of the measure \mathbf{P} .

The function $z(\cdot)$ is a partial inverse of $F_{\mathbf{P}}$ and also denoted by $F_{\mathbf{P}}^{-1}$. By the right continuity of $F_{\mathbf{P}}$ we always have

$$F_{\mathbf{P}}(z(p)) = F_{\mathbf{P}}(F_{\mathbf{P}}^{-1}(p)) \geq p, \quad (3.2)$$

see Section 4.1 for an application. If \mathbf{P} is given by a probability density function (*p.d.f.*) f , that is

$$F_{\mathbf{P}}(x) = \mathbf{P}((-\infty, x]) = \int_{-\infty}^x f(y) dy,$$

then $F_{\mathbf{P}}$ is indeed continuous and we have $F_{\mathbf{P}}(z(p)) = p$. For further properties of the quantile function see Exercise 3.17.

The basic idea of a quantile-quantile plot is as follows: We are given a sequence of pseudo-random numbers $(U_i)_{i=1, \dots, N}$ and we want to test the hypothesis that the $(U_i)_{i=1, \dots, N}$ follow the law \mathbf{P} . Thus, if this is true, for large N , approximately one quarter of the U_i should fall in each of the intervals $(-\infty, z(1/4)]$, $(z(1/4), z(1/2)]$, $(z(1/2), z(3/4)]$, $(z(3/4), \infty)$. More generally, partition $(0, 1]$ into $(p_{i-1}, p_i]$, $i = 1, \dots, N$ where $p_i = \frac{i}{N}$ and observe that a random variable U with continuous c.d.f. $F_{\mathbf{P}}$ satisfies

$$\mathbf{P}(z(p_{i-1}) < U \leq z(p_i)) = F_{\mathbf{P}}(z(p_i)) - F_{\mathbf{P}}(z(p_{i-1})) = \frac{1}{N}.$$

Hence, if we generate kN values U_j , $j = 1 \dots Nk$ we expect that, on average, k values lie in each interval $(z(p_{i-1}), z(p_i)]$. In the extreme case $k = 1$ every interval $(z(p_{i-1}), z(p_i)]$ is hit, on average, by one random number U_j , $j = 1, \dots, N$. Therefore, if $(\hat{U}_j)_{j=1, \dots, N}$ denotes the given sequence of pseudo-random numbers $(U_j)_{j=1, \dots, N}$ in increasing order then we should see something close to a straight line if we draw the tuples $(\hat{U}_i, z(p_i))$ for $i = 1, \dots, N$.

We demonstrate this in Figure 3.3. There we generate a sequence of pseudo-random numbers with the RANDU generator (see Example 3.4) and another with the Fibonacci generator (see Example 3.6).

Since both generators are supposed to be $U(0, 1)$ generators we have to determine the distribution-function $F_{U(0,1)}$. We get

$$F_{U(0,1)}(x) = \int_{-\infty}^x \mathbb{1}_{(0,1)}(y) dy = x \quad \text{for } x \in (0, 1).$$

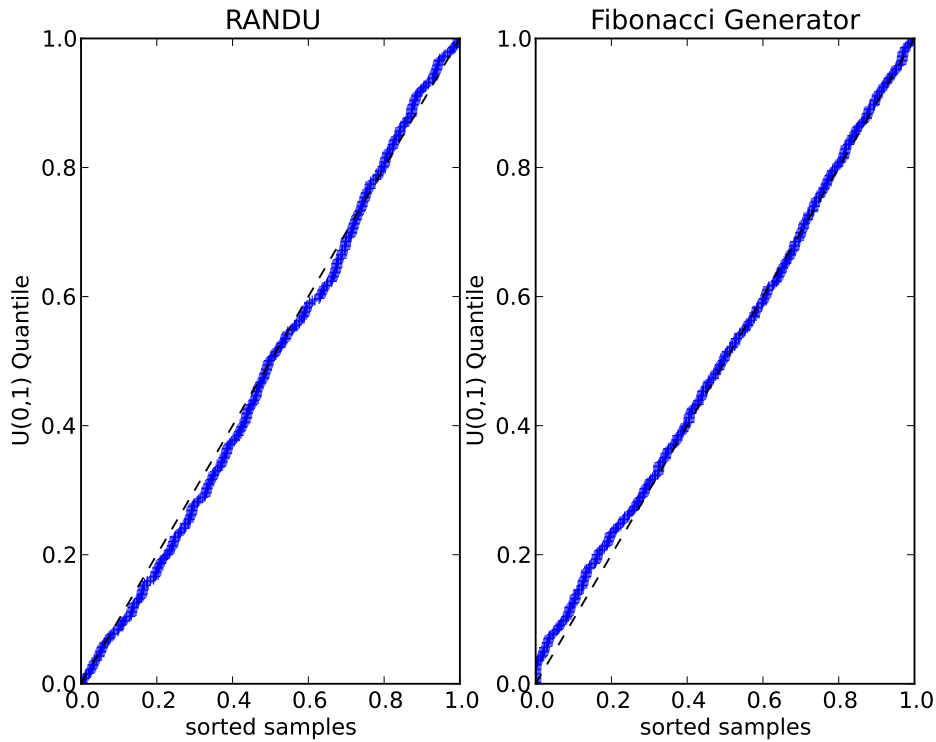


Fig. 3.3 Quantile-Quantile Plots of sorted samples generated by RANDU (left hand picture, seed $X_0 = 1000$) and the Fibonacci generator (right hand picture, seeds $X_0 = 1000$, $X_1 = 3456$) plotted against the quantiles of the $U(0,1)$ distribution. Number of samples $N = 500$.

Hence, we obtain $z(p) = p$. Thus, in order to draw Figure 3.3 we simply sorted the generated pseudo-random numbers and plotted the tuples $(\hat{U}_i, i/N)$ for $i = 1, \dots, N$.

As we can see, for both generators the plotted tuples lie close to the straight dashed line $x \mapsto x$. Thus, both sequences behave in this test as if they are uniformly distributed.

Runs Test

While a goodness-of-fit test and a quantile-quantile plot are static tests which do not take into account in which sequence the output is generated, a runs test is a so-called dynamic test. It is designed to test if a given output sequence U_1, U_2, \dots behaves in the same way as a sequence which is generated by independent and

identically distributed random variables. But note that this test gives no information about the distribution of the U_i .

In more detail, fix $N \geq 2$ and let R_N denote the statistic which records the number of increasing or decreasing “runs” in the output U_1, \dots, U_N .

For example, for $N = 10$ consider the sequence

$$(0.55, 0.57, 0.54, 0.08, 0.29, 0.24, 0.20, 0.80, 0.67, 0.91).$$

Then $(0.55, 0.57)$, $(0.08, 0.29)$, $(0.20, 0.80)$ and $(0.67, 0.91)$ are “up runs” and $(0.57, 0.54, 0.08)$, $(0.29, 0.24, 0.20)$ and $(0.80, 0.67)$ are “down runs”. Therefore, we have $R_{10} = 7$. More formally, we have

$$R_N = \#\{i \in \{2, \dots, N-1\} : (U_{i+1} - U_i)(U_i - U_{i-1}) < 0\} + 1.$$

If the output is i.i.d. then one can show that $\mathbf{E}(R_N) = (2N-1)/3$ and $\text{var}(R_N) = (16N-29)/30$ (compare Problem 3.18). Moreover, one can invoke the central limit theorem to show that $(R_N - \mathbf{E}(R_N))/\sqrt{\text{var}(R_N)}$ converges in distribution to $N(0, 1)$ as $N \rightarrow \infty$. The latter fact may then be used to formulate a decision rule.

Test Suites

There exist collections of standardized statistical tests for $U(0, 1)$ -PRNG. We introduce some of the important suites.

DIEHARD is a battery of eighteen statistical tests which consists of several variants of the goodness-of-fit tests and runs tests. DIEHARD is due to G. Marsaglia and based on [24]. The source code of the tests is available at

<http://www.stat.fsu.edu/pub/diehard/>

A detailed presentation of the tests is found in [10, Ch. 2].

A similar test suite is maintained by the National Institute of Standards and Technology (NIST) of the USA. The software and its documentation is found on the webpage

<http://csrc.nist.gov/rng/>

A very extensive test suite is TestU01. It includes the tests from DIEHARD and NIST and many more. It is written in C and maintained by P. L’Ecuyer and R. Simard [19]. The source code and a detailed documentation is found on

www.iro.umontreal.ca/~simardr/testu01/tu01.html

3.5 The Mersenne Twister

A nowadays widely used PRNG is the Mersenne Twister (MT) . For example, it is the built-in random number generator in the Python module Numpy and in Matlab®², where it is used since version 7.4. The Mersenne Twister [27] was proposed by M. Matsumoto and T. Nishimura in 1998 and passed several statistical tests including DIEHARD.

The MT-PRNG has a very huge period of $M = 2^{19937} - 1$, a Mersenne prime. Moreover, it is 623-equidistributed to 32-bit accuracy which is defined as follows.

Definition 3.12. A sequence X_i of w -bit integers of period M is called k -equidistributed to b -bit accuracy if the following property holds. In each string of length kb

$$(X_{i,b}, X_{i+1,b}, \dots, X_{i+k-1,b}), \quad i = 0, \dots, M-1,$$

where $X_{i,b}$ is the string of the first b bits of X_i , each of the possible 2^{kb} combinations of bits occurs the same number of times, except for the all-zero combination which occurs once less often. The largest k such that k -equidistribution to b -bit accuracy holds is denoted by $k(b)$.

Note that this is a finite version of the equidistribution of k -tuples chosen from an i.i.d. $U(0,1)$ -sequence. For the parameters in (3.4) below it is shown in [27] that $k(b) \geq 623$ holds for $b = 1, \dots, 32$, in particular $k(1) = 19937$, $k(32) = 623$.

For the rest of this section we sketch the idea behind the Mersenne Twister. For more details, in particular the choice of seeds, we refer to the material on the authors' webpage². The content of this section is based on [27].

The Mersenne Twister is a variant of the so called generalized feedback shift register (GFSR) generator. This generator makes use of very fast bit operations which correspond to the polynomial algebra over the two-element field \mathbb{F}_2 .

Let X_i , $i \geq 0$, denote a *word vector* of dimension w , where we usually have $w = 32$. More formally, X_i is a row vector with entries over the field \mathbb{F}_2 . Thus, we can identify X_i with a machine word consisting of w bits. The MT algorithm generates a sequence of word vectors $(X_i)_{i \geq 0}$ which are then considered to be uniform pseudo-random integers between 0 and $M = 2^w - 1$.

We denote by $X \oplus Y$ the addition of two row vectors $X, Y \in \mathbb{F}_2^w$. Moreover, for a fixed integer $0 \leq r \leq w-1$ we define the upper and lower bits of $X = (x_{w-1}, \dots, x_0) \in \mathbb{F}_2^w$ by

$$X^u = (x_{w-1}, \dots, x_r), \quad X^l = (x_{r-1}, \dots, x_0).$$

The following linear recursion builds the core element of the algorithm

² <http://www.math.sci.hiroshima-u.ac.jp/~m-mat/MT/emt.html>

Exercises

Problem 3.13. Let $(X_i)_{i \geq 1}$ be a sequence of natural numbers which is generated by the Fibonacci-PRNG. Show that $\mathbf{P}(X_i < X_{i+1} < X_{i-1}) = 0$, but $\mathbf{P}(U_i < U_{i+1} < U_{i-1}) = \frac{1}{6}$ for an i.i.d. $U(0, 1)$ -sequence $(U_i)_{i \geq 1}$.

Problem 3.14. Prove Theorem 3.8 for the special case when M is a prime.
Hint: Use Fermat's theorem ($a^M = a \pmod{M}$ for all primes M and integers a).

Problem 3.15. Prove the first part of Theorem 3.9. That is, given $k \in \mathbb{Z}^d$ with $\sum_{i=1}^d k_i a^{i-1} = 0 \pmod{M}$, then all points $\pi_i = \frac{1}{M}(X_i, \dots, X_{i+d-1}) \in [0, 1)^d, i \geq 0$ generated by the $(a, 0, M)$ -LCG with seed $X_0 \in \{1, \dots, M-1\}$ lie in one of the hyperplanes

$$H_n(k) = \{x \in \mathbb{R}^d : k^T x = n\}, \quad n \in \mathbb{Z}.$$

There are at most $\sum_{i=1}^d |k_i|$ such hyperplanes that intersect the d -cube $[0, 1)^d$.

Problem 3.16. Implement the RANDU-LCG in Matlab. For an arbitrary seed $X_0 \neq 0$ generate the first $N = 10^7$ appropriately normalized pseudo random numbers $(X_i)_{i=1}^N \subset [0, 1]$, group them into pairs $(X_i, X_{i+1}), i = 1, \dots, N-1$, and plot only those pairs which lie in the subregion $[0, 0.001] \times [0, 1]$. Interpret the result in view of Theorem 3.9 and Figure 3.1.

Problem 3.17. Show that the quantile function $F_{\mathbf{P}}^{-1}$ defined in (3.1) is monotone increasing, left continuous and satisfies

$$F_{\mathbf{P}}^{-1}(p) = \sup\{z \in \mathbb{R} : F_{\mathbf{P}}(z) < p\}.$$

Problem 3.18. For a runs test we are interested in the random variable

$$R_N = \#\{i \in \{2, \dots, N-1\} : (U_{i+1} - U_i)(U_i - U_{i-1}) < 0\} + 1, \quad N \geq 2.$$

- (i) Show that $\mathbf{E}[R_N] = (2N-1)/3$ if $(U_i)_{i=1}^N$ is a family of i.i.d. $U(0, 1)$ random variables.
- (ii) Show that $\mathbf{E}[R_N] = (2N-1)/3$ remains true if $(U_i)_{i=1}^N$ is a family of i.i.d. random variables whose distribution is given by a continuous probability density function $f : \mathbb{R} \rightarrow [0, \infty)$.

Hint: For (ii) use integration by parts.

Chapter 4

Generating Random Variables with Non-Uniform Distribution

In this section we work under the general assumption that we have available an endless stream of independent and identically distributed random numbers $(U_i)_{i \geq 0}$, with known distribution. Our aim is to generate a second sequence of i.i.d. random variables which follow a different distribution.

In practice, the $(U_i)_{i \geq 0}$ will often be generated by a $U(0,1)$ -PRNG from the previous section. We will discuss general methods like the inverse transformation and the accept-reject-algorithm to produce an arbitrary target distribution. Then, we will focus on the Box-Muller method and Marsaglia's Ziggurat method which generate normally distributed random variables.

Literature: [30, Ch. 2], [10, Ch. 4 - 5], [11, Ch. 2.3] [21, Ch. 2.3]

4.1 Inversion Method

Suppose we are given a target cumulative distribution-function $F: \mathbb{R} \rightarrow [0, 1]$ and a $U(0,1)$ distributed random variable U . Then the *inversion method* uses the quantile function (3.1) to generate a random variable with distribution F .

Proposition 4.1. *Let $F: \mathbb{R} \rightarrow [0, 1]$ be a distribution-function and let U be a random variable with $U \sim U(0,1)$. Then $F^{-1}(U) \sim F$.*

Proof. Recall from (3.1) that

$$F^{-1}(U(\omega)) = \inf \{ \xi \in \mathbb{R} : U(\omega) \leq F(\xi) \}.$$

Using (3.2) and the isotony of F we find that $F^{-1}(U(\omega)) \leq x$ is equivalent to $U(\omega) \leq F(x)$ for every $x \in \mathbb{R}$. Note that this also holds in case $U(\omega) = 0$ by setting $F^{-1}(0) = -\infty$. Therefore, $\mathbf{P}(F^{-1}(U) \leq x) = \mathbf{P}(U \leq F(x)) = F(x)$ and the assertion follows. \square

Of course, generating an F -distributed sequence $F^{-1}(U_i)$ from the given sequence U_i needs an efficient algorithm for evaluating F^{-1} . In some cases this can be done explicitly.

Example 4.2. The exponential distribution $\text{Exp}(\lambda)$ with parameter $\lambda > 0$ has the distribution-function $F(x) = 1 - \exp(-\lambda x)$ for $x \geq 0$ and $F(x) = 0$ for $x < 0$. By Proposition 4.1 the random variable $-\frac{1}{\lambda} \log(1 - U)$ then is exponentially distributed.

One way of inverting a continuous distribution-function F numerically is to compute and store in advance the quantiles $z_1 < z_2 < \dots < z_{k-1}$ such that $F(z_i) = \frac{i}{k}$, $i = 1, \dots, k-1$, where k is sufficiently large. If F is smooth then a Newton type method is suitable for doing this. Moreover, let z_0 be the largest value such that F can be regarded as 0 to the left of z_0 , and let z_k be the smallest value such that F can be regarded as 1 to the right of z_k .

Then generate a pseudo random number U and let $i = \lfloor kU \rfloor$ so that $F(z_i) \leq U < F(z_{i+1})$. The value of X is then determined approximately by one step of regula falsi (linear inverse interpolation) from

$$X = z_i + (z_{i+1} - z_i) \frac{U - F(z_i)}{F(z_{i+1}) - F(z_i)} = z_i + (z_{i+1} - z_i)(kU - i). \quad (4.1)$$

4.2 Rejection Method

The rejection method (also called *acceptance/rejection method*) is due to John von Neumann. It generates a random variable with a given p.d.f. f by using two independent $U(0, 1)$ generators. The idea is to construct a box $[a, b] \times [0, c]$ that contains the graph (of the positive part) of f , i.e. $\{(x, f(x)) : x \in \mathbb{R}, f(x) > 0\}$. Then generate a random point (U, Y) in the box and accept U if (U, Y) is below the graph (see Figure 4.1).

More precisely, choose $a < b$, $0 < c$ such that $f(x) = 0$ for $x \notin [a, b]$ and $f(x) \leq c$ for all $x \in [a, b]$. Then perform the following

Simple rejection algorithm

1. Generate $U \sim U(a, b)$ from $U = a + (b - a)U_1$ with $U_1 \sim U(0, 1)$.
2. Generate $Y \sim U(0, c)$ from $Y = cU_2$ with $U_2 \sim U(0, 1)$.
3. If $Y \leq f(U)$ then *accept* $X = U$, else *reject* U and return to step 1.

This algorithm can be generalized for a d -dimensional p.d.f. $f : \mathbb{R}^d \rightarrow \mathbb{R}$ under the assumption that we can generate another random variable with p.d.f. $g : \mathbb{R}^d \rightarrow \mathbb{R}$ where for some $c > 0$,

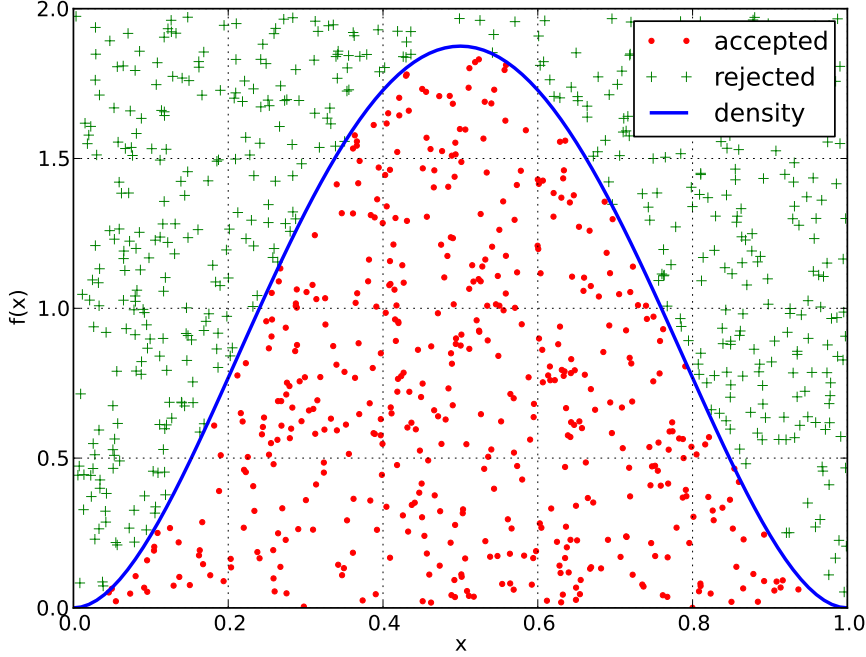


Fig. 4.1 Illustration of the rejection method. All x -coordinates of points (U, Y) are accepted if the condition $Y \leq f(U)$ is satisfied for a given p.d.f. f .

$$f(x) \leq cg(x), \quad \text{for all } x \in \mathbb{R}^d. \quad (4.2)$$

In case $d = 1$, $g = \frac{1}{b-a} \mathbb{1}_{[a,b]}$ the method reduces to the simple rejection algorithm above.

General rejection algorithm

1. Generate $U \in \mathbb{R}^d$ with $U \sim g$.
2. Generate $Y \sim U(0, c)$ from $Y = cU_2$ with $U_2 \sim U(0, 1)$.
3. If $Yg(U) \leq f(U)$ then *accept* $X = U$, else *reject* U and return to step 1.

Let us first note that the rejection algorithm produces a random variable X which is the restriction of a uniformly distributed random variable to a smaller probability space. The probability of $X \in A'$ for some $A' \in \mathcal{B}(\mathbb{R}^d)$ is the probability of $U \in A'$ given that the condition $Yg(U) \leq f(U)$ holds. More formally, consider the event

$$B = \{\omega \in \Omega : Y(\omega)g(U(\omega)) \leq f(U(\omega))\} \quad (4.3)$$

and assume $\mathbf{P}(B) > 0$. Then the random variable X generated by the algorithm, agrees with the restriction $U|_B$ considered as a random variable on the probability

space $(B, \mathcal{F}_B, \mathbf{P}_B)$ with $\mathcal{F}_B = \mathcal{F} \cap B$ and

$$\mathbf{P}_B(A \cap B) = \mathbf{P}(A|B) = \frac{\mathbf{P}(A \cap B)}{\mathbf{P}(B)}, \quad A \in \mathcal{F}. \quad (4.4)$$

Theorem 4.3 (Fundamental Theorem of Simulation). *The random variable X generated by the general rejection algorithm has the probability density function f .*

Proof. Using (4.4) with $A = \{U \in A'\}$, $A' \in \mathcal{B}(\mathbb{R}^d)$, we obtain

$$\mathbf{P}_B(X \in A') = \frac{\mathbf{P}(\{U \in A'\} \cap B)}{\mathbf{P}(B)}. \quad (4.5)$$

Since U and Y are independent the joint distribution of $T = (U, Y)$ is the product of the single distributions (Theorem 2.2)

$$\mathbf{P}_T = \mathbf{P}_U \otimes \mathbf{P}_Y = g(u)\lambda^d(u) \otimes \frac{1}{c}\mathbb{1}_{[0,c]}(y)\lambda^1(y).$$

Consider $B' = \{(u, y) \in \mathbb{R}^d \times [0, c] : yg(u) \leq f(u)\}$. Since $\mathbb{1}_B = \mathbb{1}_{B'} \circ T$ the transformation theorem (Theorem 2.5) yields

$$\begin{aligned} \mathbf{P}(B) &= \int \mathbb{1}_{B'} \circ T \, d\mathbf{P} = \int \mathbb{1}_{B'} \, d\mathbf{P}_T \\ &= \int_{\mathbb{R}^d} \frac{1}{c} \int_{\{y \in [0, c] : yg(u) \leq f(u)\}} dy \, g(u) \, du \\ &= \int_{\mathbb{R}^d} \frac{1}{c} f(u) \, du = \frac{1}{c}. \end{aligned} \quad (4.6)$$

The inner integral in (4.6) equals $\frac{f(u)}{g(u)} \leq c$ if $g(u) > 0$, and its value is c in case $g(u) = 0$ (and hence $f(u) = 0$ by (4.2)). In the same way as above, replace \mathbb{R}^d by A' and find

$$\mathbf{P}(\{U \in A'\} \cap B) = \int_{A'} \frac{1}{c} f(u) \, du.$$

Thus (4.5) yields $\mathbf{P}_B(X \in A') = \int_{A'} f(u) \, du$ for all $A' \in \mathcal{B}(\mathbb{R}^d)$. \square

We discuss the numerical effort for the general rejection algorithm. By (4.6) the probability that the algorithm accepts a value U is $\frac{1}{c}$. Hence c is the expected number of proposals needed to generate one value of X . We would like to minimize this value while satisfying (4.2), i.e. the optimal value is

$$c = \sup_{g(x) > 0} \frac{f(x)}{g(x)}. \quad (4.7)$$

Example 4.4. Let us generate a random variable $X \sim N(0, 1)$ by using a generator for a random variable U with p.d.f. $g(x) = \frac{1}{2}e^{-|x|}, x \in \mathbb{R}$. In this case (4.7) gives

$$c = \sup_{x \in \mathbb{R}} \frac{2 \exp(-\frac{x^2}{2})}{\sqrt{2\pi} \exp(-|x|)} = \sqrt{\frac{2}{\pi}} \sup_{x \geq 0} \exp(x - \frac{x^2}{2}) = \sqrt{\frac{2e}{\pi}} \approx 1.3155. \quad (4.8)$$

The variable U can be generated by taking logarithms of a $U(0, 1)$ -distributed variable with random signs (for a justification and an alternative via the inversion method see Problem 4.8)

1. Generate independent $U_1, U_2 \sim U(0, 1)$.
2. Set $U = -\log(U_2)$ if $U_1 \leq \frac{1}{2}$ and $U = \log U_2$ otherwise.
3. Let $Y = cg(U)U_3$ where $U_3 \sim U(0, 1)$.
4. If $\sqrt{2\pi}Y \leq \exp(-\frac{U^2}{2})$ then accept $X = U$.

As for efficiency, we require an average of $3c \approx 3.9465$ random numbers in order to generate one X -value. By (4.8) the test in step 4 can be simplified to $U_3 \leq \exp(-\frac{1}{2}(1 - |U|)^2)$. For a more detailed efficiency analysis one should also count the additional number of flops and function evaluations of log and exp.

4.3 The Box-Muller Method

This method uses polar coordinates in order to generate two i.i.d. $N(0, 1)$ variables from two i.i.d. $U(0, 1)$ variables:

Box-Muller Method

1. Generate two independent $U_1 \sim U(0, 1), U_2 \sim U(0, 1)$.
2. Set $X_1 = \sqrt{-2 \log U_1} \cos(2\pi U_2), X_2 = \sqrt{-2 \log U_1} \sin(2\pi U_2)$.

Theorem 4.5 (Box-Muller Method). *The Box-Muller Method generates two independent $N(0, 1)$ distributed random variables X_1, X_2 .*

Proof. It is sufficient to show that for every $A \in \mathcal{B}(\mathbb{R}^2)$

$$\mathbf{P}_{(X_1, X_2)}(A) = \mathbf{P}((X_1, X_2) \in A) = \int_A \frac{1}{2\pi} \exp\left(-\frac{1}{2}(x_1^2 + x_2^2)\right) d(x_1, x_2). \quad (4.9)$$

Following the proof in [12, Ch. 5.2] we consider the open sets $G = \mathbb{R}^2 \setminus \{(y, 0) \in \mathbb{R}^2 : y \geq 0\}$ and $G' = (0, 1)^2$ and define a function $h: G' \rightarrow G$ by

$$h(x) = \left(\frac{\sqrt{-2 \log x_1} \cos(2\pi x_2)}{\sqrt{-2 \log x_1} \sin(2\pi x_2)} \right), \quad x = (x_1, x_2) \in G'.$$

Note that h is bijective with inverse

$$h^{-1}(y) = \begin{pmatrix} \exp\left(-\frac{1}{2}|y|^2\right) \\ \frac{1}{2\pi}\text{atan2}(y) \end{pmatrix}, \quad y = (y_1, y_2) \in G.$$

Here, $\text{atan2}: G \rightarrow (0, 2\pi)$ is given by

$$\text{atan2}(y) = \text{atan2}(y_1, y_2) = \begin{cases} \arctan\left(\frac{y_2}{y_1}\right), & \text{for } y_1, y_2 > 0, \\ \frac{1}{2}\pi, & \text{for } y_1 = 0, y_2 > 0, \\ \arctan\left(\frac{y_2}{y_1}\right) + \pi, & \text{for } y_1 < 0, y_2 \in \mathbb{R}, \\ \frac{3}{2}\pi, & \text{for } y_1 = 0, y_2 < 0, \\ \arctan\left(\frac{y_2}{y_1}\right) + 2\pi, & \text{for } y_1 > 0, y_2 < 0. \end{cases}$$

By this definition h^{-1} is continuously differentiable with

$$Dh^{-1}(y) = \begin{pmatrix} -y_1 x_1 & -y_2 x_1 \\ -\frac{1}{2\pi} \frac{y_2}{y_1^2 + y_2^2} & \frac{1}{2\pi} \frac{y_1}{y_1^2 + y_2^2} \end{pmatrix},$$

where $x_1 = \exp(-\frac{1}{2}|y|^2)$. Therefore, h^{-1} is a C^1 -diffeomorphism of G onto G' .

Now, for all $A \in \mathcal{B}(G)$ it holds that

$$\mathbf{P}_{(X_1, X_2)}(A) = \mathbf{P}((X_1, X_2) \in A) = \mathbf{P}(h(U_1, U_2) \in A) = \mathbf{P}((U_1, U_2) \in h^{-1}(A)).$$

Since $h^{-1}(A) \subset G'$ is a measurable set, the function $f': G' \rightarrow \mathbb{R}$, which is given by the indicator function $f'(u) = \mathbb{1}_{h^{-1}(A)}(u)$, is λ^2 -integrable on G' . Hence, we are able to apply the general transformation theorem for integrals (see Theorem 2.6) and get

$$\mathbf{P}((U_1, U_2) \in h^{-1}(A)) = \int_{G'} \mathbb{1}_{h^{-1}(A)} d\lambda^2 = \int_G (\mathbb{1}_{h^{-1}(A)} \circ h^{-1}) |\det Dh^{-1}| d\lambda^2.$$

Next, we use

$$(\mathbb{1}_{h^{-1}(A)} \circ h^{-1})(y) = \mathbb{1}_{h^{-1}(A)}(h^{-1}(y)) = \mathbb{1}_A(y),$$

and

$$|\det Dh^{-1}(y)| = \frac{x_1}{2\pi} = \frac{1}{2\pi} \exp\left(-\frac{1}{2}|y|^2\right).$$

Altogether, this proves the desired result for all measurable sets $A \in \mathcal{B}(G)$, since

$$\begin{aligned}\mathbf{P}_{(X_1, X_2)}(A) &= \int_G \mathbb{1}_A(y) \frac{1}{2\pi} \exp\left(-\frac{1}{2}|y|^2\right) dy \\ &= \frac{1}{2\pi} \int_A \exp\left(-\frac{1}{2}(y_1^2 + y_2^2)\right) d(y_1, y_2).\end{aligned}$$

For general $A \in \mathcal{B}(\mathbb{R}^2)$ consider the decomposition in disjoint sets

$$A = (A \cap G) \cup (A \setminus G).$$

Then, we have $A \cap G \in \mathcal{B}(G)$ and since $\lambda^2(A \setminus G) = 0$ we get

$$\begin{aligned}\mathbf{P}_{(X_1, X_2)}(A \cap G) &= \frac{1}{2\pi} \int_{A \cap G} \exp\left(-\frac{1}{2}(y_1^2 + y_2^2)\right) d(y_1, y_2) \\ &= \frac{1}{2\pi} \int_A \exp\left(-\frac{1}{2}(y_1^2 + y_2^2)\right) d(y_1, y_2).\end{aligned}$$

Thus, the proof is complete if it holds that

$$\mathbf{P}_{(X_1, X_2)}(A \setminus G) = 0.$$

But this is true since

$$\mathbf{P}_{(X_1, X_2)}(A \setminus G) \leq \mathbf{P}(\sin(2\pi U_2) = 0) = \mathbf{P}(U_2 \in \{0, 1\}) = 0.$$

□

4.4 Marsaglia's Ziggurat Method

Because of the function evaluations involved, both the Box-Muller and the inversion method are considered to be quite expensive. A fast alternative method that is nowadays used in MATLAB's `randn` (algorithm `mcg16807` since MATLAB 5) is the *Ziggurat Method* proposed by Marsaglia and Tsang [25, 26]. The name *ziggurat* (in German: *die Zikkurat*) refers to the stepped pyramids with a temple on top that were built in ancient Mesopotamia (~ 2100 B.C.) and that appeared as terraced pyramids about 3000 years later in Central America with the Aztecs and other peoples.

The ziggurat method is an improvement of the acceptance/rejection method in order to increase the efficiency $\frac{1}{c}$ and to reduce the number of function evaluations. Rather than covering the graph of a given p.d.f. f by one box $[a, b] \times [0, c]$ one uses a 'ziggurat of rectangles', a cap and a tail which all have equal area. Then one chooses one of the slices at random and applies the rejection method to the chosen slice (see Figure 4.2 for an illustration).

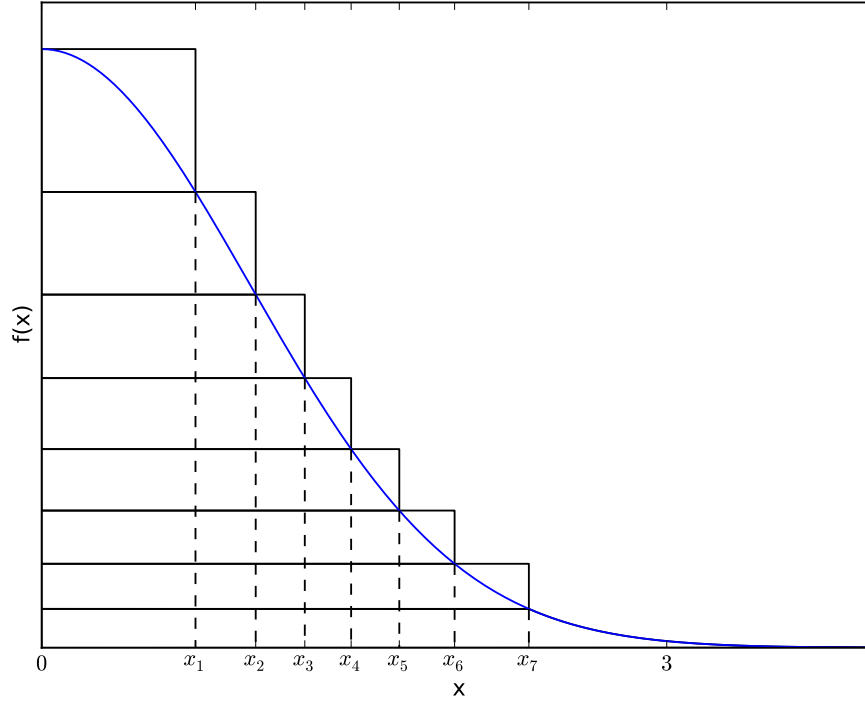


Fig. 4.2 The ziggurat method with $M = 7$ rectangles and a bottom base strip for the p.d.f. of the Gaussian distribution.

More precisely, let us assume that f has support in $[0, \infty)$ and is monotone decreasing. For a given number M we ask for nodes

$$0 = x_0 < x_1 < \dots < x_M < \infty,$$

such that the ziggurat rectangles

$$Z_i = [0, x_i) \times (f(x_i), f(x_{i-1})], \quad i = 1, \dots, M$$

and the tail

$$Z_0 = \{(x, y) \in [0, \infty) \times [0, \infty) : y \leq \min(f(x), f(x_M))\}$$

have the same area, i.e. find the nodes $(x_i)_{i=1}^M$ and $v > 0$ such that

$$v = \lambda^2(Z_0) = x_M f(x_M) + \int_{x_M}^{\infty} f(x) dx, \quad (4.10)$$

$$v = (f(x_{i-1}) - f(x_i))x_i, \quad i = 1, \dots, M. \quad (4.11)$$

In Problem 4.10 we compute the nodes $(x_i)_{i=1}^M$ and $v > 0$ by using Newton's method.

Then the **Ziggurat Method** proceeds as follows:

1. Choose $i \in \{0, \dots, M\}$ at random, uniformly distributed.
2. Generate $U \sim U(0, 1)$.
3. If $i \geq 1$ then
 - Let $X = Ux_i$.
 - If $X < x_{i-1}$ then return X ,
 - else generate $Y \sim U(0, 1)$ independent of U .
If $f(x_i) + Y(f(x_{i-1}) - f(x_i)) < f(X)$ then accept X , otherwise reject.
4. If $i = 0$ then
 - Set $X = \frac{vU}{f(x_M)}$.
 - If $X < x_M$, accept X ,
 - else generate a random value $X \in [x_M, \infty)$ from the tail.

Note that in step 3 the condition $X < x_{i-1}$ holds in most cases. Then neither the second variable Y is generated nor $f(X)$ is evaluated. The same remark applies to the case $i = 0$. In that case a rectangle of area v and height $f(x_M)$ has length $\frac{v}{f(x_M)}$, hence $X = \frac{vU}{f(x_M)}$ generates a random point on the base line of this rectangle that is accepted if $X < x_M$. Otherwise one generates an X -value from the tail for which different methods are available (see below).

The efficiency of the method is determined by the quotient of the integral of f and the area of the ziggurat

$$\text{eff} = \left(\lambda^2 \left(\bigcup_{i=0}^M Z_i \right) \right)^{-1} = \frac{1}{(M+1)v}. \quad (4.12)$$

For the exponential p.d.f. $f(x) = \max(x, 0)e^{-x}$ the paper [26] gives the following data:

$$\begin{aligned} M &= 255, & x_M &= 7.697117470131104972, \\ \text{eff} &= 0.989, & v &= 0.00394965598225815571993 \quad . \end{aligned}$$

For the tail one generates $U_1 \sim U(0, 1)$ and takes $X = x_M - \log(U_1)$ since then (compare with Example 4.2)

$$\mathbf{P}(X \in [a, b]) = \mathbf{P}(e^{x_M-b} \leq U_1 \leq e^{x_M-a}) = \int_a^b e^{x_M-x} dx.$$

For the normal distribution one generates an additional random sign \pm for the value of X which is determined from the ziggurat method applied to the distribution in $[0, \infty)$. For the nonnormalized density function $f(x) = \exp(-\frac{x^2}{2})$ Marsaglia and Tsang [26] give the data

$$M = 255, \quad x_M = 3.6541528853610088, \\ \text{eff} = 0.9933, \quad v = 0.00492867323399 \quad .$$

As a tail method they propose (see Problem 4.11)

1. Generate $U_1 \sim U(0, 1), U_2 \sim U(0, 1)$.
2. Set $X_1 = -\frac{\log(U_1)}{x_M}, X_2 = -\log(U_2)$.
3. If $X_1^2 < 2X_2$ accept $X = x_M + X_1$, otherwise reject.

For the implementation it is convenient to take $M = 2^N - 1$ and to generate the index i in step 1 and the value U in step 2 from a random integer $j \in \{0, \dots, 2^K - 1\}$ where $K > N$. Typical values are $N = 8, K = 32$. In addition, one saves time by putting all operations with K -bit integers into the preparatory phase

$$k_i, w_i = \begin{cases} \left\lfloor 2^K \frac{x_{i-1}}{x_i} \right\rfloor, & 2^{-K} x_i, \quad i = 1, \dots, 2^N - 1, \\ \left\lfloor 2^K \frac{x_M f(x_M)}{v} \right\rfloor, & 2^{-K} \frac{v}{f(x_M)}, \quad i = 0. \end{cases} \quad (4.13)$$

With these settings the algorithm becomes rather short:

The Ziggurat Algorithm

1. Generate a random integer $j \in \{0, \dots, 2^K - 1\}$ and determine $i \in \{0, \dots, 2^N - 1\}$ from the rightmost N bits of j .
2. Set $X = jw_i$.
3. If $j < k_i$ return X , else
 - If $i = 0$ return an X from the tail, else
 - if $(f(x_{i-1}) - f(x_i))U < f(X) - f(x_i)$, return X .
4. Go to step 1.

Exercises

Problem 4.6. The standardized logistic distribution has the density function

$$f(x) = \frac{\exp(-x)}{(1 + \exp(-x))^2}, \quad x \in \mathbb{R}.$$

Show how to use the inversion method to generate a random variable having this distribution.

Problem 4.7. Write a Matlab program with the following properties:

- (i) Implement an $N(0, 1)$ -pseudo-random number generator which is based on the inversion method and uses the Matlab functions `erfinv` and `rand`.

- (ii) Modify the idea of a quantile-quantile plot in such a way that it can be used to visualize if two pseudo-random number generators draw numbers from the same distribution. Write a function which generates a quantile-quantile plot for two given sequences of random data $U = (U_1, \dots, U_N)$ and $V = (V_1, \dots, V_N)$.
- (iii) Produce a quantile-quantile plot which compares the distribution of the Matlab function `randn` versus your PRNG from (i).

Problem 4.8. Show that the first two steps of the algorithm in Example 4.4 generate a random variable with p.d.f. $g(x) = \frac{1}{2}e^{-|x|}$. Set up an alternative algorithm which uses the inversion method instead. Compare the numerical efficiency of both methods.

Problem 4.9. In order to implement Marsaglia's Ziggurat method one first needs to determine the nodes $0 = x_0 < x_1 < \dots < x_M < \infty$ for a given number M . Write a Matlab program which approximates the nodes in the following way: Solve the nonlinear equation $g(x_M) = 0$ by using the bisection method, where the function $g: (0, \infty) \rightarrow \mathbb{R}$ is given by

$$g(x_M) = v - x_1(f(x_0) - f(x_1)),$$

with

$$v = x_M f(x_M) + \int_{x_M}^{\infty} f(x) dx,$$

$$x_{i-1} = f^{-1}\left(\frac{v}{x_i} + f(x_i)\right), \quad i = M, \dots, 2.$$

Compute the nodes for $M = 7$, $f(x) = e^{-\frac{1}{2}x^2}$ and initial values $x_7 = 2.3$ and 3 .

Problem 4.10. Write a MATLAB program that uses Newton's method to determine the solutions $(x_i)_{i=1}^M$ and $v > 0$ to the nonlinear system of equations (4.10) and (4.11) for $M = 7$ and the normal distribution density function, that is

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

Determine the efficiency of the ziggurat algorithm with $M = 7$.

Hint: Reformulate (4.10) and (4.11) into the problem of finding the root $\tilde{x} \in \mathbb{R}^{M+1}$ of a mapping $F: \mathbb{R}_{>0}^{M+1} \rightarrow \mathbb{R}^{M+1}$, where

$$\tilde{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_M \\ v \end{pmatrix}.$$

The implementation of F may use the built-in function `erf`.

Problem 4.11 (Tail method for the Ziggurat algorithm). For $x_M > 0$ consider the algorithm

1. Generate independent random variables $U_1 \sim U(0, 1)$, $U_2 \sim U(0, 1)$,
2. Set $X_1 = -\frac{\log(U_1)}{x_M}$, $X_2 = -\log(U_2)$,
3. If $X_1^2 < 2X_2$ accept $X = x_M + X_1$.

Show that this algorithm generates a random variable X whose distribution is given by the probability density function

$$f(x) = \begin{cases} 0, & x < x_M \\ c^{-1} \exp(-\frac{1}{2}x^2), & x \geq x_M, \end{cases}$$

with the constant $c = \int_{x_M}^{\infty} e^{-\frac{1}{2}y^2} dy$.

Problem 4.12. Show that the Ziggurat algorithm is obtained from the general Ziggurat method via the preparatory steps (4.13).

Chapter 5

Monte Carlo Methods

The general term *Monte Carlo Methods* refers to numerical methods that use random numbers for conducting experiments on a computer. Roughly speaking one may categorize Monte Carlo Methods into two different subclasses:

- Direct simulation of a random system;
- Addition of artificial randomness into a given (deterministic) problem followed by a simulation of the new system.

While our main goal, the numerical solution of stochastic differential equations belongs to the first category, we consider in this chapter a typical example from the second category: approximating definite integrals from a sequence of i.i.d. random variables. For example, it is quite natural to use the acceptance/rejection method from Section 4.2 in order to estimate the area under the graph of the function f in Figure 4.1 (assuming we do not know that f is a p.d.f.).

5.1 Statistical Analysis of Simulation Output

Let us first summarize basic results from statistical analysis that follow from the law of large numbers and the central limit theorem in Section 2. Given a sequence of i.i.d. random variables $X_j, j = 1, 2, \dots$ we want to estimate its expectation, the variance, and a confidence interval. To be specific, let $\mu = \mathbf{E}[X_j]$ and $\text{var}(X_j) = \sigma^2$ for $j \in \mathbb{N}$. We estimate the expectation μ by the *sample mean*

$$\bar{X}_N = \frac{1}{N} \sum_{j=1}^N X_j. \quad (5.1)$$

From the Strong Law of Large Numbers (Theorem 2.12) we infer

$$\bar{X}_N \rightarrow \mu \quad \text{almost surely} \quad \text{as} \quad N \rightarrow \infty.$$

By the linearity of the expectation,

$$\mathbf{E}[\bar{X}_N] = \frac{1}{N} \mathbf{E}\left[\sum_{j=1}^N X_j\right] = \mu,$$

and by the independence of the X_j (see (2.7)),

$$\text{var}(\bar{X}_N) = \frac{1}{N^2} \sum_{j=1}^N \text{var}(X_j) = \frac{1}{N} \sigma^2.$$

Because of these properties \bar{X}_N is called a *consistent (or unbiased) estimator*. Thus the standard deviation of \bar{X}_N is $\frac{\sigma}{\sqrt{N}}$. This decays at the rate $\frac{1}{\sqrt{N}}$ which dominates the asymptotic behavior of practically all Monte Carlo simulations.

The Central Limit Theorem 2.14 tells us that

$$(\bar{X}_N - \mu) \frac{\sqrt{N}}{\sigma} \rightarrow N(0, 1) \quad \text{in distribution.} \quad (5.2)$$

In order to use this for estimating how far our sample mean is from the expected value μ , we compute an approximation of σ^2 from the *sample variance*

$$S_N^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X}_N)^2. \quad (5.3)$$

The value S_N is called the *sample standard deviation*. The following Lemma shows that S_N^2 is an unbiased estimator of the true variance.

Lemma 5.1. *Under the assumptions above the sample variance S_N^2 has expectation $\mathbf{E}[S_N^2] = \sigma^2$.*

Proof. First note that

$$\begin{aligned} S_N^2 &= \frac{1}{N-1} \sum_{i=1}^N \left(\frac{N-1}{N} X_i - \frac{1}{N} \sum_{j \neq i} X_j \right)^2 \\ &= \frac{1}{(N-1)N^2} \sum_{i=1}^N \left((N-1)X_i - \sum_{j \neq i} X_j \right)^2, \end{aligned}$$

where $(N-1)X_i - \sum_{j \neq i} X_j$ is a sum of independent random variables which has expectation 0. Thus by (2.7) and (2.6),

$$\begin{aligned}
\mathbf{E}[S_N^2] &= \frac{1}{(N-1)N^2} \sum_{i=1}^N \mathbf{E} \left[\left((N-1)X_i - \sum_{j \neq i} X_j \right)^2 \right] \\
&= \frac{1}{(N-1)N^2} \sum_{i=1}^N \text{var} \left((N-1)X_i - \sum_{j \neq i} X_j \right) \\
&= \frac{1}{(N-1)N^2} \sum_{i=1}^N \left((N-1)^2 \sigma^2 + (N-1) \sigma^2 \right) = \sigma^2.
\end{aligned}$$

□

Definition 5.2. Let Z be a random variable with $Z \sim N(0, 1)$ and let $0 < \alpha < 1$ be given. Then the interval $[-z_\alpha, z_\alpha]$ is called a $1 - \alpha$ confidence interval (or a $100(1 - \alpha)\%$ confidence interval) provided $\mathbf{P}(|Z| \leq z_\alpha) = 1 - \alpha$. The value of α is referred to as the *significance level*.

Usually we think of α being a small number, a common value in Statistics is $\alpha = 0.05$. If f is the p.d.f. for $N(0, 1)$ then z_α can be computed from

$$1 - \alpha = \int_{-z_\alpha}^{z_\alpha} f(\xi) d\xi = 2 \int_{-\infty}^{z_\alpha} f(\xi) d\xi - 1,$$

i.e. by using a quantile from the $N(0, 1)$ -distribution. We evaluate z_α numerically with the help of the Matlab functions `erf` and `erfinv`, where `erf` is given by

$$\text{erf}(x) := \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt, \quad \text{for all } x \geq 0, \quad (5.4)$$

and, consequently,

$$z_\alpha = \sqrt{2} \text{erfinv}(1 - \alpha), \quad (5.5)$$

since

$$1 - \alpha = 2 \int_0^{z_\alpha} f(\xi) d\xi = \frac{2}{\sqrt{2\pi}} \int_0^{z_\alpha} e^{-\frac{1}{2}\xi^2} d\xi = \text{erf}\left(\frac{z_\alpha}{\sqrt{2}}\right).$$

For example, $\alpha = 0.05$ leads to $z_\alpha \approx 1.96$.

We apply this to $Z_N = (\bar{X}_N - \mu) \frac{\sqrt{N}}{S_N}$ where we replace the unknown variance σ^2 in (5.2) by the sample variance S_N^2 from (5.3). Then we obtain for N large

$$\begin{aligned}
1 - \alpha &\approx \mathbf{P}\left(|\bar{X}_N - \mu| \frac{\sqrt{N}}{\sigma} \leq z_\alpha\right) \\
&\approx \mathbf{P}\left(|\bar{X}_N - \mu| \frac{\sqrt{N}}{S_N} \leq z_\alpha\right) \\
&= \mathbf{P}\left(\mu \in \left[\bar{X}_N - z_\alpha \frac{S_N}{\sqrt{N}}, \bar{X}_N + z_\alpha \frac{S_N}{\sqrt{N}}\right]\right).
\end{aligned} \quad (5.6)$$

Therefore, we take

$$\left[\bar{X}_N - z_\alpha \frac{S_N}{\sqrt{N}}, \bar{X}_N + z_\alpha \frac{S_N}{\sqrt{N}} \right] \quad (5.7)$$

as our $100(1 - \alpha)\%$ confidence interval. Note that the width of this interval is proportional to the sample standard deviation and inverse proportional to the square root of the number of samples.

Note that (5.6) involves two approximations: N should be so large that the Central Limit Theorem applies and such that our sample variance is sufficiently accurate.

Example 5.3 ([14], Section 15.2). Consider the i.i.d. sequence $X_j = \exp(Z_j)$ obtained from an i.i.d. sequence $Z_j \sim N(0, 1)$. The expectation is

$$\begin{aligned} \mathbf{E}[X_j] &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \exp(x - \frac{1}{2}x^2) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \exp(-\frac{1}{2}(x-1)^2 + \frac{1}{2}) dx = \sqrt{e}. \end{aligned}$$

We compute a Monte Carlo approximation of $\mu = \mathbf{E}[X_j]$ by using $N = 2^k, k = 5, \dots, 17$ samples. In Figure 5.1 the resulting values \bar{X}_N are plotted versus N (logarithmic scale for both axes), the dashed line is at height $\mu = \sqrt{e}$, and the 95% confidence intervals (5.7) are indicated by *error bars*.

For the second example we return to the computation of the expected payoff from (1.11).

Example 5.4 ([14], Section 15.3). Given a time horizon $T > 0$, volatility $\sigma > 0$, interest rate $r > 0$, and initial value S_0 we compute the discounted expected payoff

$$\begin{aligned} V(S_0) &= \exp(-rT) \mathbf{E}[\max(0, S(T) - E)], \quad \text{where} \\ S(T) &= S_0 \exp\left(\left(r - \frac{1}{2}\sigma^2\right)T + \sigma\sqrt{T}Z\right) \end{aligned}$$

and $Z \sim N(0, 1)$. Here we again apply the risk neutrality assumption $\mu = r$.

As in the previous example Figure 5.2 shows the values of the sample mean \bar{V}_N together with the 95% confidence intervals computed from (5.7) for the data

$$T = 1, \quad S_0 = 10, \quad E = 9, \quad \sigma = 0.1, \quad r = 0.06, \quad N = 2^k \quad (k = 5, \dots, 17).$$

For comparison the dashed line gives the value of the Black-Scholes formula (1.10).

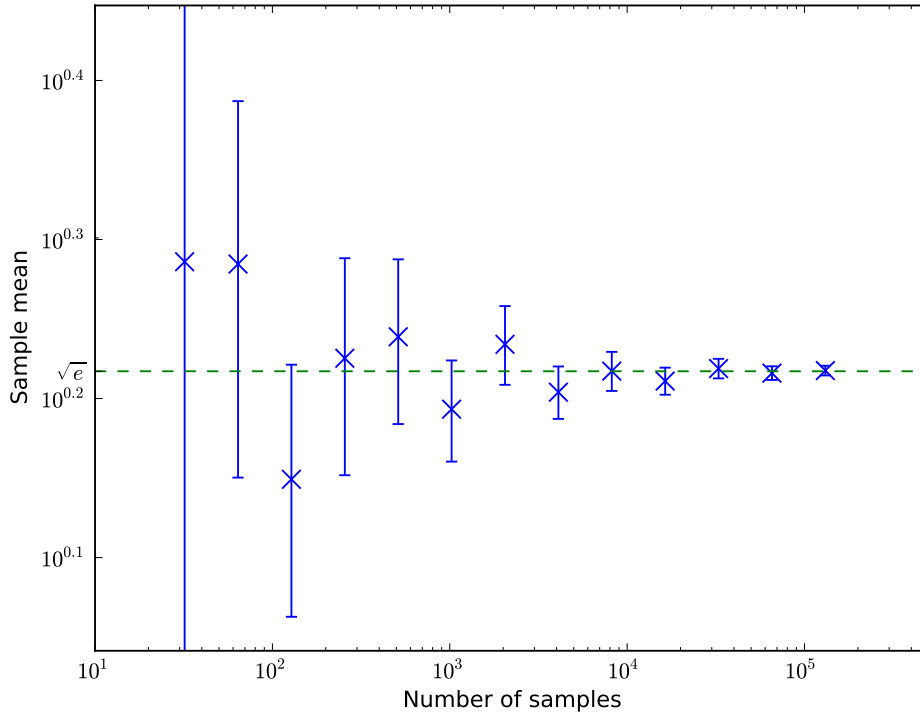


Fig. 5.1 Monte-Carlo approximations to $\mathbf{E}[\exp(Z)]$ with $Z \sim N(0, 1)$. The sample means \bar{X}_N for a given number of samples N are marked by crosses, which lie in the middle of their respective confidence intervals. The dashed line is at height $\sqrt{e} = \mathbf{E}[\exp(Z)]$.

5.2 Monte Carlo Integration

The principles from the previous section can be easily generalized to approximate integrals like

$$I_h = \int_{\mathbb{R}^d} h(x)f(x) dx, \quad (5.8)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a given p.d.f. and $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is a \mathcal{B}^d -measurable function with $hf \in L^1(\mathbb{R}^d)$. An example of (5.8) is the one-dimensional integral

$$I_h = \int_0^1 h(x)dx, \quad h \in L^1[0, 1], \quad (5.9)$$

where $d = 1$ and $f = \mathbb{1}_{[0,1]}$.

Suppose we can simulate a d -dimensional random variable $X : \Omega \rightarrow \mathbb{R}^d$ with p.d.f. f . Then we can express (5.8) as

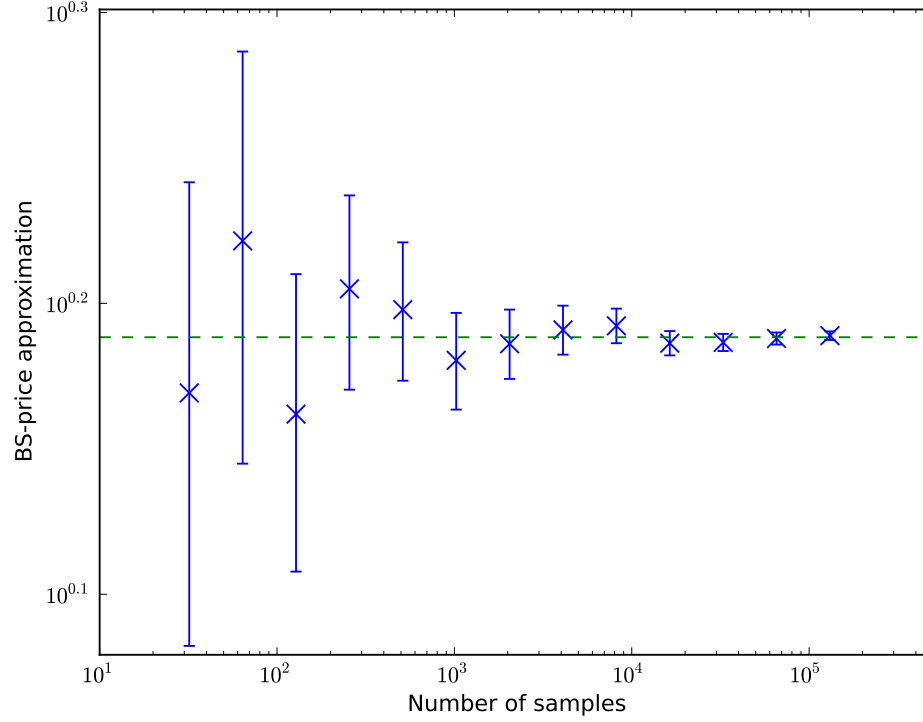


Fig. 5.2 Monte-Carlo approximations to discounted expected payoff $V(S_0)$ of a European call with parameters $T = 1$, $S_0 = 10$, $E = 9$, $\sigma = 0.1$, and $r = 0.06$. The dashed line indicates the option value computed by the Black-Scholes formula (1.10).

$$\mathbf{E}[h \circ X] = \int_{\mathbb{R}^d} h(x)f(x) \, dx = I_h. \quad (5.10)$$

Recall from [4, Satz 3.6] (cf. Problem 2.18) that $h \circ X$ has variance

$$\sigma^2 := \text{var}(h \circ X) = \int_{\mathbb{R}^d} h^2(x)f(x) \, dx - \left(\int_{\mathbb{R}^d} h(x)f(x) \, dx \right)^2. \quad (5.11)$$

In view of Section 5.1 this suggests to use the sample mean

$$\bar{h}_N = \frac{1}{N} \sum_{j=1}^N h(X_j), \quad (5.12)$$

where $X_j, j \in \mathbb{N}$, form an \mathbb{R}^d -valued i.i.d. sequence which all have the p.d.f. f . Since X_j are independent, so are the $h \circ X_j$ (see Theorem 2.3) and hence,

$$\text{var}(\bar{h}_N) = \frac{1}{N} \sigma^2.$$

By the Law of Large Numbers, we have convergence $\bar{h}_N \rightarrow I_h$ almost surely. Moreover, the Central Limit Theorem (Theorem 2.14) asserts

$$(\bar{h}_N - I_h) \frac{\sqrt{N}}{\sigma} \rightarrow N(0, 1) \quad \text{in distribution.}$$

As in Section 5.1 we estimate the variance σ^2 by

$$S_N^2 = \frac{1}{N-1} \sum_{i=1}^N (h(X_i) - \bar{h}_N)^2.$$

We then obtain a $100(1 - \alpha)\%$ confidence interval from

$$\mathbf{P}\left(I_h \in \left[\bar{h}_N - z_\alpha \frac{S_N}{\sqrt{N}}, \bar{h}_N + z_\alpha \frac{S_N}{\sqrt{N}}\right]\right) \approx 1 - \alpha \quad (5.13)$$

where z_α is determined by (5.5).

There are two important observations that we will discuss in the following subsections:

- While the convergence rate $\frac{1}{\sqrt{N}}$ is unavoidable we can try to speed up the process by transforming the integral (5.8) such that its variance σ^2 becomes smaller;
- The rate of convergence $\frac{1}{\sqrt{N}}$ is independent of the dimension d of the phase space and of the smoothness of the integrand. This is in contrast to standard quadrature methods (trapezoidal sum, Romberg integration, Gauß-quadrature) and suggests that Monte Carlo methods are an alternative to classical methods for higher dimensions, see Section 5.4.

5.3 Variance Reduction Techniques

This section is based on [14, Ch.21] and [21, Ch.3]. A rather extensive treatment under the heading of *importance sampling* can be found in [9, Ch. 6].

Variance reduction by antithetic variates

Consider the sample mean from (5.12)

$$\bar{h}_N = \frac{1}{N} \sum_{j=1}^N h(X_j), \quad X_j \sim U(0, 1) \quad (5.14)$$

for approximating $\mathbf{E}[h \circ X] = \int_0^1 h(x) dx = I_h$. We replace \bar{h}_N by the new antithetic sampler

$$\hat{h}_N = \frac{1}{N} \sum_{j=1}^N \frac{1}{2} (h(X_j) + h(1 - X_j)). \quad (5.15)$$

Since the summands are independent we obtain from (2.7)

$$\text{var}(\hat{h}_N) = \frac{1}{N^2} \sum_{j=1}^N \text{var}\left(\frac{1}{2}(h(X_j) + h(1 - X_j))\right).$$

Using $\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y)$ ((2.8)) we can write

$$\text{var}\left(\frac{1}{2}(h(X_j) + h(1 - X_j))\right) = \frac{1}{2} (\text{var}(h \circ X_j) + \text{cov}(h(X_j), h(1 - X_j))). \quad (5.16)$$

Thus we have a reduction of variance if $\text{cov}(h(X_j), h(1 - X_j))$ is negative, the smaller the better. Assuming $\text{cov}(h(X_j), h(1 - X_j)) = -c^2$ we find

$$\text{var}(\hat{h}_N) = \frac{1}{N^2} N \frac{1}{2} (\sigma^2 - c^2) = \frac{1}{2N} (\sigma^2 - c^2). \quad (5.17)$$

According to (5.13) the width of the confidence interval improves by a factor of

$$\rho = \left(\frac{\text{var}(\bar{h}_N)}{\text{var}(\hat{h}_N)} \right)^{1/2} = \sqrt{2} \left(\frac{\sigma^2}{\sigma^2 - c^2} \right)^{1/2}. \quad (5.18)$$

The factor $\sqrt{2}$ is natural since the antithetic sampler needs $2N$ function evaluations, so the improvement factor is $\frac{\sigma}{\sqrt{\sigma^2 - c^2}}$.

The following theorem gives a sufficient criterion for two random variables to have positive covariance.

Theorem 5.5. *Let X be a random variable on $(\Omega, \mathcal{F}, \mathbf{P})$ with values in a Borel set $A \subset \mathbb{R}$ and let $f, g \in L^2(A, \mathbf{P}_X)$ be real-valued functions that are either both increasing or both decreasing. Then*

$$\text{cov}(g(X), h(X)) \geq 0. \quad (5.19)$$

Proof. Let Y be another random variable which is independent of X and has the same distribution (for the existence see [4, Kor.9.5]). By the monotonicity assumption

$$(g(x) - g(y))(h(x) - h(y)) \geq 0 \quad \text{for all } x, y \in A,$$

and hence by the independence of $h(X)$ and $g(X)$ (compare Theorem 2.3),

$$\begin{aligned} 0 &\leq \mathbf{E}[(g(X) - g(Y))(h(X) - h(Y))] \\ &= \mathbf{E}[g(X)h(X)] - \mathbf{E}[g(Y)h(X)] - \mathbf{E}[g(X)h(Y)] + \mathbf{E}[g(Y)h(Y)] \\ &= 2(\mathbf{E}[g(X)h(X)] - \mathbf{E}[g(X)]\mathbf{E}[h(X)]) \\ &= 2\text{cov}(g(X), h(X)). \end{aligned}$$

□

Remark 5.6. There is a generalization of Theorem 5.5 to d variables that can be proved by induction on d , see [21, Theorem 3.17]. Consider an \mathbb{R}^d -valued random variable $X = (X_1, \dots, X_d)$ with independent components X_j that have values in Borel sets $A_j \subset \mathbb{R}$. Two square integrable functions $g, h : \prod_{j=1}^d A_j \rightarrow \mathbb{R}$ that are monotone increasing with respect to each argument then satisfy (5.19).

The application of Theorem 5.5 to the antithetic sampler (5.15) is quite obvious: if $h(x)$ is monotone increasing or decreasing on $[0, 1]$, then so is $-h(1-x)$, hence

$$\text{cov}(h(X_j), h(1-X_j)) = -\text{cov}(h(X_j), -h(1-X_j)) \leq 0.$$

Example 5.7. ([14, Ch. 21.4])

Let $h(x) = e^{\sqrt{x}}$, $x \in [0, 1]$ and $X \sim U(0, 1)$. Then

$$\begin{aligned} \mathbf{E}[h(X)] &= \int_0^1 e^{\sqrt{x}} dx = 2 \\ \text{var}(h(X)) &= \int_0^1 e^{2\sqrt{x}} dx - 4 = \frac{1}{2}(e^2 - 7), \end{aligned}$$

and

$$\begin{aligned} \text{cov}(e^{\sqrt{X}}, e^{\sqrt{1-X}}) &= \mathbf{E}[e^{\sqrt{X}\sqrt{1-X}}] - \mathbf{E}[e^{\sqrt{X}}]\mathbf{E}[e^{\sqrt{1-X}}] \\ &= \int_0^1 e^{\sqrt{x}+\sqrt{1-x}} dx - 4. \end{aligned}$$

From (5.18) we obtain the factor $\rho = 13.463$ by which the confidence interval should shrink. Table 5.1 shows the result of a Monte Carlo simulation with values $N = 10^i$, $i = 2, 3, 4, 5$. The last column gives the ratio of widths of confidence intervals for the standard and the antithetic sampler which is close to the predicted value.

Table 5.1 95%-Confidence Intervals for standard and antithetic MC-approximations of $\mathbf{E}[e^{\sqrt{U}}] = 2$ with $U \sim U(0, 1)$.

N	Standard	Antithetic	Ratio of widths
10^2	[1.85131, 2.02127]	[1.99413, 2.00732]	12.887
10^3	[1.96229, 2.01589]	[1.99945, 2.00340]	13.563
10^4	[1.99145, 2.00886]	[1.99899, 2.00028]	13.535
10^5	[1.99961, 2.00508]	[1.99966, 2.00007]	13.433

In case of normal distributions $X_j \sim N(0, 1)$ the antithetic sampler (5.15) is modified as follows

$$\hat{h}_N = \frac{1}{N} \sum_{j=1}^N \frac{1}{2} (h(X_j) + h(-X_j)). \quad (5.20)$$

If h is monotone increasing or decreasing on \mathbb{R} then Theorem 5.5 yields

$$\begin{aligned} \text{var}\left(\frac{1}{2}(h(X_j) - h(-X_j))\right) &= \frac{1}{2} (\text{var}(h(X_j)) - \text{cov}(h(X_j), -h(-X_j))) \\ &\leq \frac{1}{2} \text{var}(h(X_j)). \end{aligned}$$

Hence the reduction factor is (5.18) with $-c^2 = \text{cov}(h(X_j), -h(-X_j))$.

Example 5.8. ([14, Ch. 21.6])

For the function $h(x) = \exp(x - \frac{1}{2})$ and $X \sim N(0, 1)$ we find

$$\mathbf{E}[h \circ X] = \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2}(x-1)^2) dx = 1,$$

as well as

$$\text{var}(h(X)) = e - 1, \quad \text{cov}(h(X), -h(-X)) = -\frac{1}{e},$$

leading to a factor of $\rho = 1.595$ in (5.18).

Table 5.2 95%-Confidence intervals for standard and antithetic MC-approximations of $\mathbf{E}[e^U/\sqrt{e}] = 1$ with $U \sim N(0, 1)$.

N	Standard	Antithetic	Ratio of widths
10^2	[0.78713, 1.20907]	[0.87478, 1.07900]	2.066
10^3	[0.91377, 1.05352]	[0.94943, 1.02464]	1.858
10^4	[0.96705, 1.01504]	[0.97921, 1.00661]	1.751
10^5	[0.99783, 1.01426]	[0.99708, 1.00628]	1.786

In view of Remark 5.6 it is clear that variance reduction by the antithetic sampler (5.15) also occurs if h is a function of several variables that is monotone in each of its arguments (either increasing or decreasing). As an example we treat **Antithetic variates in option valuation.**

Consider an exotic option where the payoff function does not only depend on the final asset price $S(T)$ but also on the intermediate values

$$S(t_j), \quad t_j = j\Delta t, \quad \Delta t = \frac{T}{M}, \quad j = 0, \dots, M.$$

For example, a barrier $B > E$ is introduced and the payoff function (1.8) of the European call is modified to

$$C((S(t_j)_{j=0}^M)) = \begin{cases} L & \text{if } \max_{j=0, \dots, M} S(t_j) > B, \\ \max(0, S(T) - E) & \text{otherwise.} \end{cases}$$

If the threshold value L equals zero then we have an *up-and-out call* whereas $L = B - E$ leads to a *lock-in call*. If the asset price passes the barrier B during the time interval $[0, T]$ then the payoff is set to zero or to the limit value $B - E$. We assume that asset prices form a stochastic process generated by the solution (1.13) of the SODE (1.12)

$$S(t_{j+1}) = S(t_j) \exp\left(\left(r - \frac{1}{2}\sigma^2\right)\Delta t + \sigma\sqrt{\Delta t}Z_j\right), \quad Z_j \sim N(0, 1), \quad j = 0, \dots, M-1. \quad (5.21)$$

Figure 5.3 shows a process that crosses the barrier and another one that doesn't. Using the same Z_j from (5.21) we generate antithetic variates through

$$\bar{S}(t_{j+1}) = \bar{S}(t_j) \exp\left(\left(r - \frac{1}{2}\sigma^2\right)\Delta t - \sigma\sqrt{\Delta t}Z_j\right), \quad Z_j \sim N(0, 1), \quad j = 0, \dots, M-1. \quad (5.22)$$

Figure 5.4 shows the corresponding zigzag and zagzig curves. Performing N runs, indexed by $i = 1, \dots, N$, we finally evaluate the discounted expected payoffs from

$$V_i = e^{-rT} \begin{cases} L & \text{if } \max_{j=0, \dots, M} S_i(t_j) > B, \\ \max(0, S_i(T) - E) & \text{otherwise,} \end{cases}$$

$$\bar{V}_i = e^{-rT} \begin{cases} L & \text{if } \max_{j=0, \dots, M} \bar{S}_i(t_j) > B, \\ \max(0, \bar{S}_i(T) - E) & \text{otherwise.} \end{cases}$$

Note that the value function satisfies the monotonicity condition for the lock-in call $L = B - E$ but not for the up-and-out call $L = 0$. Table 5.3 confirms that antithetic variates give only a slight improvement over the factor $\sqrt{2} \approx 1.414$ for the up-and-out call while the improvement for the lock-in-call is substantial, see Table 5.4.

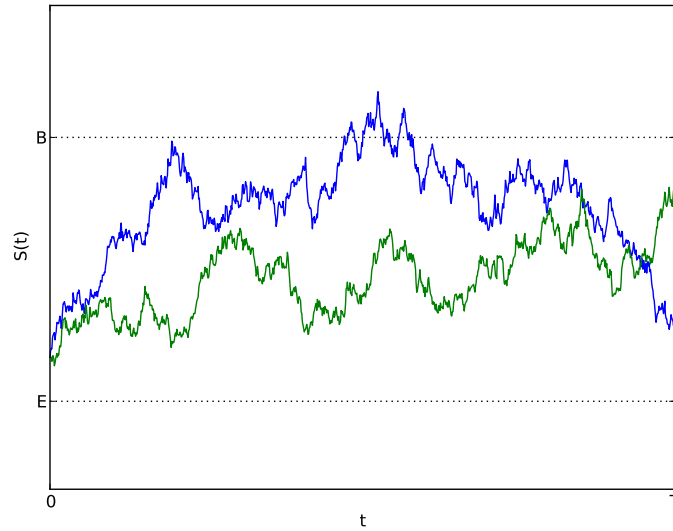


Fig. 5.3 Two runs of asset prizes that give different payoffs for an exotic option with barrier B .

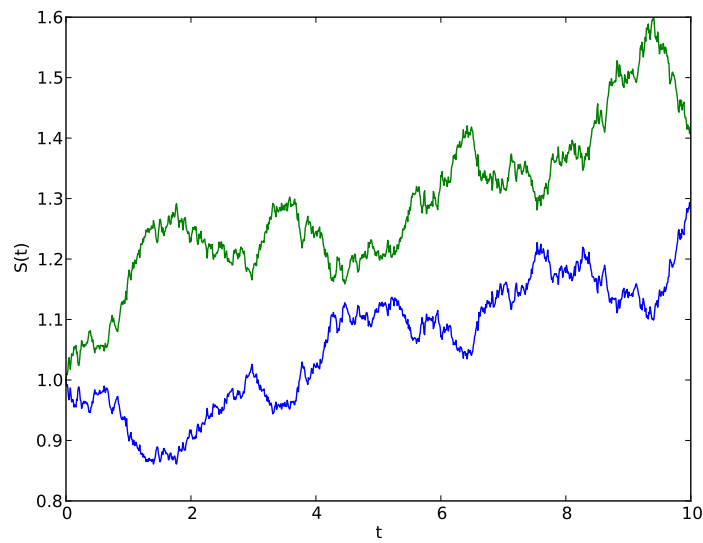


Fig. 5.4 Simulaton of asset prizes (zigzag) and their antithetic variates (zagzig) via equations (5.21) and (5.22).

Table 5.3 95%-Confidence Intervals for MC-approximations of the discounted expected payoff of an Up-and-Out Call with $S_0 = 5$, $r = 0.05$, $\sigma = 0.3$, $E = 4$, $B = 8$, $T = 10$, $M = 500$

N	Standard	Antithetic	Ratio of widths
10^2	[0.48009, 0.71627]	[0.60471, 0.73706]	1.785
10^3	[0.54037, 0.61606]	[0.55460, 0.60152]	1.613
10^4	[0.58772, 0.61197]	[0.58889, 0.60352]	1.658
10^5	[0.59321, 0.60091]	[0.59450, 0.59913]	1.661

Table 5.4 95%-Confidence Intervals for MC-approximations of the discounted expected payoff of an “Lock-in Call” with $S_0 = 5$, $r = 0.05$, $\sigma = 0.3$, $E = 4$, $B = 8$, $T = 10$, $M = 500$

N	Standard	Antithetic	Ratio of widths
10^2	[1.09702, 1.67199]	[1.14481, 1.36578]	2.602
10^3	[1.06575, 1.22454]	[1.13721, 1.20512]	2.339
10^4	[1.13987, 1.19037]	[1.15421, 1.17543]	2.379
10^5	[1.15785, 1.17386]	[1.16575, 1.17251]	2.371

A note on importance sampling:

The Monte Carlo methods of Section 5.2 assume that the p.d.f. $f \in L^1(\mathbb{R}^d)$ which occurs in the integral (5.8) is given. Suppose that this is not the case. We want to approximate an integral $\int_{\mathbb{R}^d} \tilde{h} dx$ with $\tilde{h} \in L^1(\mathbb{R}^d)$ but still have a choice of the density function. Then we write

$$I_{\tilde{h}} = \int_{\mathbb{R}^d} \tilde{h} dx = \int_{\mathbb{R}^d} \frac{\tilde{h}(x)}{f(x)} f(x) dx, \quad (5.23)$$

where $f \in L^1(\mathbb{R}^d)$ is a p.d.f. to be determined such that

$$f(x) = 0, x \in \mathbb{R}^d \implies \tilde{h}(x) = 0. \quad (5.24)$$

Let us define $h(x) = \frac{\tilde{h}(x)}{f(x)}$ if $f(x) \neq 0$ and $h(x) = 0$ otherwise, and let us assume that we can generate random variables $X \sim f$. The *importance sampler* then reads

$$\bar{h}_{N,f} = \frac{1}{N} \sum_{i=1}^N h(X_i) = \frac{1}{N} \sum_{i=1}^N \frac{\tilde{h}(X_i)}{f(X_i)}. \quad (5.25)$$

It is natural to ask for a good choice of the density function f that minimizes the variance of this sampler.

Theorem 5.9. *The variance of the importance sampler (5.25) among all p.d.f.s f satisfying (5.24) becomes minimal at*

$$f_{\min}(x) = \frac{|\tilde{h}(x)|}{\|\tilde{h}\|_{L^1}}, \quad x \in \mathbb{R}^d, \quad (5.26)$$

and the minimal value is

$$\sigma^2(f_{\min}) = \left(\int_{\mathbb{R}^d} |\tilde{h}(x)| dx \right)^2 - \left(\int_{\mathbb{R}^d} \tilde{h}(x) dx \right)^2.$$

Remark 5.10. The theorem shows that it is optimal to choose as density the absolute value of the function itself (suitably normalized). Of course, it is unrealistic to assume that one can generate random variables with this density. However, the result suggests to look for densities that have some weight where the integrand has. An example of this type of importance sampler will be discussed below.

Proof. Recall from (5.11) that $I_{\tilde{h}}$ has variance

$$\sigma^2(f) = \int_{\mathbb{R}^d} \frac{\tilde{h}(x)^2}{f(x)} dx - \int_{\mathbb{R}^d} \tilde{h}(x)^2 dx. \quad (5.27)$$

Then our assertion follows from the equality

$$\sigma^2(f) = \|\tilde{h}\|_{L^1}^2 \mathbf{E} \left[\left(\frac{f_{\min} - f}{f} \right)^2 \right] + \|\tilde{h}\|_{L^1}^2 - \left(\int_{\mathbb{R}^d} \tilde{h} dx \right)^2.$$

The latter relation is obtained from (5.26) and (5.27) by a short calculation

$$\begin{aligned} & \|\tilde{h}\|_{L^1}^2 \mathbf{E} \left[\frac{f_{\min}^2}{f^2} - 2 \frac{f_{\min}}{f} + 1 \right] \\ &= \|\tilde{h}\|_{L^1}^2 \int_{\mathbb{R}^d} \left[\frac{|\tilde{h}|^2}{\|\tilde{h}\|_{L^1}^2 f^2} - 2 \frac{|\tilde{h}|}{\|\tilde{h}\|_{L^1} f} + 1 \right] f dx \\ &= \int_{\mathbb{R}^d} \frac{\tilde{h}^2}{f} dx - 2 \|\tilde{h}\|_{L^1}^2 + \|\tilde{h}\|_{L^1}^2. \end{aligned}$$

□

Finally we discuss a method of variance reduction called **Stratified sampling**.

Consider again the integral (5.10) and the sample mean (5.12). Suppose that we can decompose \mathbb{R}^d into disjoint subsets $S_i, i = 1, \dots, M$ (the *strata*) for which we know conditional probabilities

$$a_i = \int_{S_i} f(x) dx, \quad i = 1, \dots, M, \quad \mathbb{R}^d = \bigcup_{i=1}^M S_i, \quad S_i \cap S_j = \emptyset \quad (i \neq j). \quad (5.28)$$

Note that $\sum_{i=1}^M a_i = 1$. With this information at hand we ask for an estimator of I_h which has smaller variance than \bar{h}_N . Note that (5.28) implies a decomposition of the density function f as follows

$$f = \sum_{i=1}^M a_i f_i, \quad \text{where} \quad f_i(x) = \mathbb{1}_{S_i}(x) \frac{f(x)}{a_i}, \quad x \in \mathbb{R}^d.$$

We decompose the number of samples $N = \sum_{i=1}^M n_i$ and assume that we can generate i.i.d. random variables $X_{i,1}, \dots, X_{i,n_i} \in S_i$ such that $X_{i,j} \sim f_i$. Then the sample mean in stratum S_i is given by

$$T_i = \frac{1}{n_i} \sum_{j=1}^{n_i} h(X_{i,j}),$$

and our *stratified sampler* is defined by

$$T = \sum_{i=1}^M a_i T_i. \quad (5.29)$$

A simple computation shows

$$\mathbf{E}[T_i] = \frac{I_i}{a_i}, \quad \text{where} \quad I_i = \int_{S_i} h(x) f(x) dx, \quad (5.30)$$

$$\mathbf{E}[T] = \sum_{i=1}^M a_i \mathbf{E}[T_i] = I_h. \quad (5.31)$$

For the variance we find

$$\text{var}(T) = \sum_{i=1}^M a_i^2 \text{var}(T_i) = \sum_{i=1}^M \frac{a_i^2}{n_i} \left(\int_{S_i} h^2(x) \frac{f(x)}{a_i} dx - \left(\frac{I_i}{a_i} \right)^2 \right). \quad (5.32)$$

Theorem 5.11. *In the setting above let $n_i = Na_i, i = 1, \dots, M$ be integers, then the variance of the sample mean (5.12) and the stratified mean (5.29) satisfy*

$$\text{var}(\bar{h}_N) = \text{var}(T) + \frac{1}{N} \sum_{i=1}^M a_i \left(\frac{I_i}{a_i} - I_h \right)^2. \quad (5.33)$$

Remark 5.12. The choice $n_i = Na_i$ is called *proportional allocation*. In general, one cannot assume that $n_i = Na_i$ are integers. Then one takes nearby values such that $N = \sum_{i=1}^M n_i$ still holds. Equation (5.33) shows that the variance is reduced by proportional allocation and that the reduction is better the further apart the values of $\mathbf{E}(T_i)$ and I_h are.

Proof. By (5.32) and the choice of n_i the right-hand side of (5.33) equals

$$\begin{aligned} & \sum_{i=1}^M a_i^2 \frac{1}{a_i N} \left(\int_{S_i} h^2 \frac{f}{a_i} dx - \left(\frac{I_i}{a_i} \right)^2 \right) + \frac{1}{N} \sum_{i=1}^M \left(\frac{I_i^2}{a_i} - 2I_i I_h + a_i I_h^2 \right) \\ &= \frac{1}{N} \sum_{i=1}^M \int_{S_i} h^2(x) f(x) dx - 2 \frac{I_h^2}{N} + \frac{I_h^2}{N} = \text{var}(\bar{h}_N). \end{aligned}$$

□

Example 5.13. ([21, Ch.3.2])

Suppose you ask N Bielefeld residents whether their city mayor is doing a good job. For simplicity think of each resident to represent an interval of length 1 and let $R(x)$ denote the resident of the interval to which x belongs. You then consider the sample mean $\bar{h}_N = \frac{1}{N} \sum_{i=1}^N h(x_i)$ of the function

$$h(x) = \begin{cases} 1 & \text{if } R(x) \text{ says yes,} \\ 0 & \text{if } R(x) \text{ says no.} \end{cases}$$

We stratify the set of residents by

$$S_i = \{x : R(x) \text{ votes for party } i\}, i = 1, \dots, 6,$$

where i stands for SPD, CDU, GRÜNE, FDP, LINKE, OTHER. From the last election we know the proportion $a_i = \frac{1}{B} \int_{S_i} 1 dx$ of residents (B total number of voters) that voted for party i . Suppose you also ask the N residents how they voted during the last election and you obtain numbers n_i for the i -th party. Suppose 44.2% of them voted SPD while you know from the newspapers that we had only $a_1 = 0.396$ during the last election. You may then worry about the result of your sample mean because you probably have too many SPD-voters in your sample and the SPD is the mayor's party. Therefore you decide to extract from your N -sample a proportional allocation by taking the answers of $\tilde{n}_i = \tilde{N} a_i \leq n_i$ voters $x_{i,j}$ of party i such that $\tilde{N} = \sum_{i=1}^6 \tilde{n}_i \leq N$ (note that $\tilde{N} = \min \frac{n_i}{a_i}$ is the maximum number for doing this). Then the stratified estimate of the mayor's supporters is

$$T = \sum_{i=1}^6 a_i T_i, \quad T_i = \frac{1}{\tilde{n}_i} \sum_{j=1}^{\tilde{n}_i} h(x_{i,j}).$$

The following theorem shows that proportional allocation is not the best choice for minimizing the variance of the stratified sampler.

Theorem 5.14. (*Tschuprow-Neyman allocation*)

The variance of the stratified sampler (5.29) is minimal for

$$n_i^{\text{opt}} = N \frac{a_i \sigma_i}{\sum_{j=1}^M a_j \sigma_j}, \quad \sigma_i = \text{var}(h(X_{i,j})), j = 1, \dots, n_i, \quad (5.34)$$

and the minimal value is

$$\text{var}(T^{\text{opt}}) = \frac{1}{N} \left(\sum_{j=1}^M a_j \sigma_j \right)^2.$$

Proof. From (5.32) we infer for a general decomposition $N = \sum_{i=1}^M n_i$

$$V(n_1, \dots, n_M) := \text{var}(T) = \sum_{i=1}^M a_i^2 \frac{\sigma_i^2}{n_i}.$$

An application of the Cauchy Schwarz inequality yields the assertion

$$\begin{aligned} V(n_1^{\text{opt}}, \dots, n_M^{\text{opt}}) &= \frac{1}{N} \left(\sum_{j=1}^M a_j \sigma_j \right)^2 \\ &= \frac{1}{N} \left(\sum_{j=1}^M \sqrt{n_j} \frac{a_j \sigma_j}{\sqrt{n_j}} \right)^2 \leq \frac{1}{N} \left(\sum_{j=1}^M n_j \right) \left(\sum_{j=1}^M \frac{a_j^2 \sigma_j^2}{n_j} \right) \\ &= V(n_1, \dots, n_M). \end{aligned}$$

□

It is known that the Cauchy Schwarz inequality is strict unless the vectors are linearly dependent, i.e. $\sqrt{n_j} = c \frac{a_j \sigma_j}{\sqrt{n_j}}$ for some $c \in \mathbb{R}$. Since $\sum_{j=1}^M n_j = N$ this leads to $n_j = n_j^{\text{opt}}$, hence (5.34) is the unique allocation where the minimum is attained. In general it is hard to achieve this minimum since one needs to know the strata variances σ_i^2 . Moreover, n_i^{opt} need not be integers in general. A practical way out is to start sampling with proportional allocation and then estimate σ_i^2 from the sample variance

$$S_{N,i}^2 = \frac{1}{N a_i} \sum_{j=1}^{n_i} (h(X_{i,j}) - T_i)^2.$$

One may then use the allocation (5.34) with $S_{N,i}$ instead of σ_i and continue with increasing the sample size.

5.4 Approximation of Multiple Integrals

This section is based on [9], [10, Ch. 7.1], [30, Ch. 4.3]. We first review quadrature rules for one-dimensional integrals, with an emphasis on Gaussian quadrature for weight functions that are probability density functions. Then we show how to generalize to multiple integrals via product rules. Finally we compare the efficiency of these quadrature rules with Monte-Carlo integration.

The general goal is to approximate an integral

$$I(h) = \int_A h(x)f(x) dx, \quad A \subset \mathbb{R}^d \quad \text{measurable} \quad (5.35)$$

with a p.d.f. $f \in L^1(A)$ by a *quadrature rule* of the form

$$Q_M(h) = \sum_{i=1}^M w_i h(x_i), \quad \text{where } x_1, \dots, x_M \in A. \quad (5.36)$$

We call $x_i \in A$ the *nodes* and $w_i \in \mathbb{R}$ the *weights* of the quadrature rule.

In the following we assume $A \subset \mathbb{R}$ and the nodes to be ordered by $x_1 < \dots < x_M$. Let \mathcal{P}_{M-1} denote the space of polynomials of degree $\leq M-1$. The following result about interpolatory quadrature is well known (the proof is an exercise).

Theorem 5.15. *The quadrature rule (5.36) is exact for polynomials in \mathcal{P}_{M-1} (i.e. $Q_M(h) = I(h)$ for all $h \in \mathcal{P}_{M-1}$) if and only if the weights are given by*

$$w_i = \int_A \ell_i(x)f(x) dx, \quad i = 1, \dots, M, \quad (5.37)$$

with the Lagrange polynomials defined by

$$\ell_i(x) = \prod_{j=1, j \neq i}^M \frac{x - x_j}{x_i - x_j}.$$

For equidistant nodes $x_j = a + (j-1)\Delta x$, $j = 1, \dots, M$ with $A = [a, b]$, $\Delta x = \frac{b-a}{M-1}$ we obtain the Newton-Cotes formulae. The most familiar ones are

$$Q_2(h) = \frac{b-a}{2} (h(a) + h(b)), \quad \text{trapezoidal rule,} \quad (5.38)$$

$$Q_3(h) = \frac{b-a}{6} (h(a) + 4h(\frac{a+b}{2}) + h(b)), \quad \text{Simpson's rule.} \quad (5.39)$$

For Gauß quadrature rules consider an interval $A \subset \mathbb{R}$ (not necessarily compact) and the function space

$$L_f^2(A) = \{h : A \mapsto \mathbb{R} \text{ measurable and } h^2 f \in L^1(A)\}.$$

Assuming $f > 0$ a.e. in A , the space $L_f^2(A)$ is a Hilbert space (cf. Section 2.3) with respect to

$$\langle h, g \rangle = \int_A h(x)g(x)f(x) dx, \quad \|h\|^2 = \langle h, h \rangle.$$

From the monomials $m_j(x) = x^j$, $j = 0, \dots, M-1$ one generates orthogonal polynomials $P_i \in \mathcal{P}_i$ by the Gram-Schmidt orthogonalization process

$$P_i = m_i - \sum_{j=0}^{i-1} \frac{\langle m_i, P_j \rangle}{\langle P_j, P_j \rangle} P_j, \quad i = 0, \dots, M-1. \quad (5.40)$$

Rather than normalizing $\langle P_i, P_i \rangle = 1$ the recursion (5.40) enforces the leading coefficient to be 1, i.e. $P_i(x) = x^i + l.o.t..$ We call $P_i, i \in \mathbb{N}$ the *orthogonal polynomials associated with the weight function f and the interval A* . Using this property it is easy to verify that orthogonal polynomials satisfy a three term recursion.

Theorem 5.16. *The orthogonal polynomials associated with the weight function f and the interval A satisfy the recursion*

$$\begin{aligned} P_{i+1}(x) &= (x - a_{i+1})P_i(x) - b_{i+1}P_{i-1}(x), \quad i \geq 0, \\ a_{i+1} &= \frac{\langle m_1 P_i, P_i \rangle}{\langle P_i, P_i \rangle}, \quad b_{i+1} = \frac{\langle m_1 P_i, P_{i-1} \rangle}{\langle P_{i-1}, P_{i-1} \rangle}, \end{aligned}$$

where we set $P_{-1} \equiv 0, P_0 \equiv 1, b_1 = 0$.

From this we conclude important information about the zeroes of the orthogonal polynomials.

Theorem 5.17. *Let the orthogonal polynomials be defined by (5.40). Then P_i has exactly i distinct zeroes in the interior of A .*

Proof. Let z_1, \dots, z_ν be the zeroes of P_i in the interior of A repeated according to their multiplicity. We have $\nu \leq i$ since $P_i \in \mathcal{P}_i$. Then the polynomial $Q_i(x) = P_i(x) \prod_{l=1}^\nu (x - z_l)$ is of one sign in the interior of A and hence

$$0 \neq \int_A Q_i(x)f(x) dx = \int_A P_i(x) \prod_{l=1}^\nu (x - z_l)f(x) dx.$$

If $\nu < i$ then this integral vanishes due to the orthogonality $\langle P_i, P_l \rangle = 0$ for $l < i$. Hence we conclude $\nu = i$. Now assume that there is a zero of multiplicity $k \geq 2$, say z_1 . Then we can repeat the above argument with $Q_i(x) = P_i(x)(x - z_1)^{k-2} \prod_{l=2}^i (x - z_l)$ and arrive at a contradiction. \square

Definition 5.18. Let

$$x_{1,M} < \dots < x_{M,M} \quad (5.41)$$

denote the distinct zeroes of P_M in A and let

$$w_{i,M} = \int_A \ell_{i,M}(x) f(x) dx, \quad \ell_{i,M}(x) = \prod_{j=1, j \neq i}^M \frac{x - x_{j,M}}{x_{i,M} - x_{j,M}}, \quad i = 1, \dots, M \quad (5.42)$$

be the corresponding Lagrangian weights (cf. (5.37)). Then the quadrature rule

$$Q_M(h) = \sum_{i=1}^M w_{i,M} h(x_{i,M}) \quad (5.43)$$

is called the M -th Gaussian rule associated with domain A and density f .

Theorem 5.19. The M -th Gaussian rule associated with domain A and density f is exact for polynomials in \mathcal{P}_{2M-1} , i.e.

$$Q_M(h) = \int_A h(x) f(x) dx \quad \text{for all } h \in \mathcal{P}_{2M-1}.$$

Proof. For $h \in \mathcal{P}_{2M-1}$ let $P_{[h]} = \sum_{i=1}^M h(x_{i,M}) \ell_{i,M} \in \mathcal{P}_{M-1}$ be its interpolating polynomial. By Theorem 5.17 there exists some $r \in \mathcal{P}_{M-1}$ such that

$$h(x) = P_{[h]}(x) + r(x)P_M(x), \quad x \in \mathbb{R}.$$

From this factorization and the orthogonality we obtain

$$\begin{aligned} I(h) &= \int_A P_{[h]}(x) f(x) dx + \langle r, P_M \rangle \\ &= I(P_{[h]}) = Q_M(P_{[h]}) = \sum_{i=1}^M w_{i,M} P_{[h]}(x_{i,M}) \\ &= \sum_{i=1}^M w_{i,M} (P_{[h]}(x_{i,M}) + r(x_{i,M})P_M(x_{i,M})) = Q_M(h). \end{aligned}$$

□

The favorable properties of Gaussian quadrature rules allow for a rather simple convergence proof as $M \rightarrow \infty$.

Theorem 5.20. The weights of the M -th Gaussian rule satisfy

$$w_{i,M} > 0 \quad (i = 1, \dots, M), \quad \sum_{i=1}^M w_{i,M} = 1. \quad (5.44)$$

If A is a compact interval then one has convergence

$$\lim_{M \rightarrow \infty} Q_M(h) = \int_A h(x) f(x) dx \quad \text{for all } h \in C(A). \quad (5.45)$$

Proof. The squares $h_i = (\ell_{i,M})^2$ of the Lagrange polynomials (5.42) satisfy $h_i \in \mathcal{P}_{2M-2}$ and hence by Theorem 5.19

$$0 < \int_A h_i(x) f(x) \, dx = \sum_{j=1}^M w_{j,M} h_i(x_{j,M}) = w_{i,M}.$$

Moreover, $\sum_{i=1}^M w_{i,M} = Q_M(\mathbb{1}_A) = \int_A f(x) \, dx = 1$.

Given $\varepsilon > 0$ and $h \in C(A)$ we can select by the Weierstraß approximation theorem some $M \in \mathbb{N}$ and $g \in \mathcal{P}_{2M-1}$ such that $\|h - g\|_\infty \leq \frac{\varepsilon}{2}$. Using (5.44) and Theorem 5.19 leads to

$$\begin{aligned} & \left| \int_A h(x) f(x) \, dx - \sum_{i=1}^M w_{i,M} h(x_{i,M}) \right| \\ & \leq \left| \int_A (h(x) - g(x)) f(x) \, dx \right| + \left| \sum_{i=1}^M w_{i,M} (g(x_{i,M}) - h(x_{i,M})) \right| \\ & \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} \sum_{i=1}^M |w_{i,M}| = \varepsilon. \end{aligned}$$

□

The following table collects some special orthogonal polynomials with their domains and (nonnormalized) density functions.

Table 5.5 Table of orthogonal polynomials with domain A and density function f

name	domain A	density $f(x)$	restrictions
Legendre	$[-1, 1]$	1	
Jacobi	$[-1, 1]$	$(1-x)^\alpha(1+x)^\beta$	$\alpha, \beta > -1$
Laguerre	$[0, \infty)$	$x^\alpha e^{-x}$	$\alpha > -1$
Hermite	$(-\infty, \infty)$	$\exp(-x^2)$	

Well known particular cases are

$$P_j(x) = \frac{1}{2^j j!} \frac{d^j}{dx^j} ((x^2 - 1)^j) \quad \text{Legendre,}$$

$$T_j(x) = \cos(j \arccos(x)) \quad \text{Chebyshev (Jacobi with } \alpha = \beta = -\frac{1}{2}),$$

$$H_j(x) = (-1)^j e^{x^2} \frac{d^j}{dx^j} (e^{-x^2}) \quad \text{Hermite.}$$

Error estimates

As noted above classical integration rules aim at accurate integration of sufficiently smooth functions. Approximating smooth functions by a Taylor expansion then suggests to set up quadrature rules that integrate polynomials exactly to the highest degree possible. For an extensive treatment of error estimates as $M \rightarrow \infty$ we refer to the classical monographs of Davis, Rabinowitz [7] and Stroud [33]. Here we draw some simple conclusions that follow from the results above.

Theorem 5.21. *Let $A \subset \mathbb{R}$ be a compact interval and consider a quadrature rule Q_M as in (5.36) that is exact for polynomials in \mathcal{P}_{M-1} . Then for all $h \in C^M(A)$,*

$$|I(h) - Q_M(h)| \leq C_{1,M} \|h^{(M)}\|_{L^\infty}, \quad C_{1,M} = \frac{1}{M!} \int_A \prod_{i=1}^M |x - x_i| f(x) dx. \quad (5.46)$$

Proof. Let $P_{[h]} \in \mathcal{P}_{M-1}$ interpolate the data $(x_i, h(x_i)), i = 1, \dots, M$. For every $x \in A$ there exists some $\xi = \xi(x) \in (x_1, x_M)$ such that

$$h(x) - P_{[h]}(x) = \frac{h^{(M)}(\xi)}{M!} \omega_M(x), \quad \omega_M(x) = \prod_{i=1}^M (x - x_i). \quad (5.47)$$

With this representation we find from Theorem 5.15

$$\begin{aligned} |I(h) - Q_M(h)| &= \left| \int_A (h(x) - P_{[h]}(x)) f(x) dx \right| \\ &\leq \frac{\|h^{(M)}\|_{L^\infty}}{M!} \int_A |\omega_M(x)| f(x) dx. \end{aligned}$$

□

For the Gauss rules with respect to domain A and density f the corresponding result is the following.

Theorem 5.22. *Let $A \subset \mathbb{R}$ be a compact interval and consider a Gaussian quadrature rule (5.43) associated with domain A and density f . Then for all $h \in C^{2M}(A)$,*

$$|I(h) - Q_M(h)| \leq C_{2,M} \|h^{(2M)}\|_{L^\infty}, \quad C_{2,M} = \frac{1}{(2M)!} \int_A \prod_{i=1}^M (x - x_i)^2 f(x) dx. \quad (5.48)$$

Proof. We modify the proof of Theorem 5.21 by letting $P_{[h]} \in \mathcal{P}_{2M-1}$ solve the Hermite interpolation problem

$$P_{[h]}(x_{i,M}) = h(x_{i,M}), \quad P'_{[h]}(x_{i,M}) = h'(x_{i,M}), \quad i = 1, \dots, M.$$

Then the error formula (5.47) holds with $2M$ and ω_M^2 instead of M and ω_M , and the estimate (5.48) follows as in the previous proof. □

Compound rules

One way of increasing the accuracy of quadrature rules is to partition the domain of integration into smaller subdomains and to use a simple rule on each subdomain. In this way one obtains so called *compound rules*. In one dimension the resulting formulae integrate exactly continuous and piecewise polynomial functions (splines). Note however, that this approach usually requires the density to be constant, for otherwise we will need a new quadrature rule on each subdomain.

Let $A = [a, b]$ and $f(x) = 1, x \in [a, b]$ and partition as follows

$$[a, b] = \bigcup_{j=0}^{K-1} [y_j, y_{j+1}], \quad y_j = a + j\Delta x, j = 0, \dots, K, \quad \Delta x = \frac{b-a}{K}.$$

We apply a formula (5.36) for $\int_0^1 h dx$ to the functions $g_j(y) = h(y_j + y\Delta x), y \in [0, 1], j = 0, \dots, K-1$,

$$\begin{aligned} \int_a^b h(x) dx &= \sum_{j=0}^{K-1} \int_{y_j}^{y_{j+1}} h(x) dx = \Delta x \sum_{j=0}^{K-1} \int_0^1 g_j(y) dy \\ &\approx \Delta x \sum_{j=0}^{K-1} Q_M(g_j) = \Delta x \sum_{j=0}^{K-1} \sum_{i=1}^M w_i h(y_j + x_i \Delta x) =: Q_{M,K}(h). \end{aligned} \quad (5.49)$$

Theorem 5.23. *Under the assumptions of Theorem 5.21 the compound quadrature rule $Q_{M,K}$ defined in (5.49) satisfies for all $h \in C^M(A)$*

$$|I(h) - Q_{M,K}(h)| \leq C_{1,M}(b-a)(\Delta x)^M \|h^{(M)}\|_{L^\infty}. \quad (5.50)$$

For the compound Gaussian rules we have under the assumptions of Theorem 5.22

$$|I(h) - Q_{M,K}(h)| \leq C_{2,M}(b-a)(\Delta x)^{2M} \|h^{(2M)}\|_{L^\infty}. \quad (5.51)$$

Proof. We consider only the first case, the second case is treated analogously. From (5.46) we obtain

$$\begin{aligned} |I(h) - Q_{M,K}(h)| &= \left| \Delta x \sum_{j=0}^{K-1} \left(\int_0^1 g_j(y) dy - Q_M(g_j) \right) \right| \\ &\leq \Delta x \sum_{j=0}^{K-1} C_{1,M} \|g_j^{(M)}\|_{L^\infty} \\ &\leq \Delta x C_{1,M} \sum_{j=0}^{K-1} (\Delta x)^M \|h^{(M)}\|_{L^\infty} \\ &= (b-a) C_{1,M} (\Delta x)^M \|h^{(M)}\|_{L^\infty}. \end{aligned}$$

□

A simple compound rule is the trapezium sum that derives from the trapezoidal rule (5.38)

$$Q_{2,K}(h) = \Delta x \left(\frac{1}{2} h(a) + \sum_{j=1}^{K-1} h(a + j\Delta x) + \frac{1}{2} h(b) \right). \quad (5.52)$$

This is an $\mathcal{O}((\Delta x)^2)$ approximation of the integral if $h \in C^2[a, b]$.

Product rules

Exercises

Problem 5.24. Let X be a random variable with $X \sim U(0, 1)$. Compute the expectation and the variance of

$$h_\alpha(X) = (X - \alpha)^2$$

for $0 \leq \alpha \leq 1$. For the values $\alpha \in \{0, \frac{1}{2}, 1\}$ discuss whether antithetic sampling of $h_\alpha(X)$ reduces the variance.

Problem 5.25. The acceptance-rejection algorithm of Section 4.2 can be modified to produce a Monte-Carlo integrator for

$$I_h = \int_{\mathbb{R}^d} h(x) dx, \quad h \in L^1(\mathbb{R}^d).$$

Assume $|h(x)| \leq cf(x)$, $x \in \mathbb{R}^d$, for some p.d.f. $f \in L^1(\mathbb{R}^d)$. For $i = 1, \dots, N$ generate independent random variables $Y_i \sim U(0, c)$ and $U_i \sim f$ and set

$$\bar{h}_N = \frac{c}{N} \sum_{i=1}^N X_i, \quad \text{where } X_i = \begin{cases} 1, & \text{if } Y_i f(U_i) < h(U_i), \\ -1, & \text{if } Y_i f(U_i) < -h(U_i), \\ 0, & \text{otherwise.} \end{cases}$$

Show that \bar{h}_N has expectation I_h and variance $\frac{1}{N} (c \|h\|_{L^1} - I_h^2)$. Show that the importance sampler

$$\hat{h}_N = \frac{1}{N} \sum_{i=1}^N \frac{h(Z_i)}{f(Z_i)}, \quad \text{where } Z_i \sim f$$

has the same expectation but smaller variance.

Problem 5.26. Implement the following Monte-Carlo algorithm which approximates the Delta $\Delta(t)$ of a European call option with exercise price $E > 0$ and expiry date $T > 0$ at time $t = 0$:

1. For $i \in \{1, \dots, N\}$ generate $Z_i \sim N(0, 1)$ i.i.d.
2. For $i \in \{1, \dots, N\}$ compute

$$S_i := S_0 \exp\left(\left(r - \frac{1}{2}\sigma^2\right)T + \sigma\sqrt{T}Z_i\right)$$

and for a given step size $\eta > 0$

$$S_i^\eta := (S_0 + \eta) \exp\left(\left(r - \frac{1}{2}\sigma^2\right)T + \sigma\sqrt{T}Z_i\right).$$

3. Compute $\Delta_i = \frac{1}{\eta}e^{-rT} (\max(0, S_i^\eta - E) - \max(0, S_i - E))$
4. Compute the sample mean $\Delta_N = \frac{1}{N} \sum_{i=1}^N \Delta_i$.

As parameter values use $T = 3$, $E = 5$, $S_0 = 6$, $r = 0.05$, $\sigma = 0.3$. For the number of samples choose $N = 2^5, \dots, 2^{17}$ and for the step sizes set $\eta = 10^{-2}, \dots, 10^{-4}$. For each η make a separate error bar plot which shows the exact value of $\Delta(0)$ and the approximations Δ_N for the given number of samples N together with their 95%-confidence intervals.

Problem 5.27. Let X be a real-valued random variable with unknown expectation and let Y be a random variable with known expectation. Determine the expectation and variance of

$$X_\theta = X + \theta(\mathbf{E}[Y] - Y), \quad \text{for } \theta \in \mathbb{R}.$$

How would you choose θ in order to minimize the variance when sampling X_θ and what is its minimal value? Generalize this to the case of k random variables Y_j , $j = 1, \dots, k$, by minimizing the variance of

$$X_\theta = X + \sum_{j=1}^k \theta_j(\mathbf{E}[Y_j] - Y_j)$$

with $\theta = (\theta_1, \dots, \theta_k) \in \mathbb{R}^k$. Give a criterion that guarantees a unique minimizer.

Problem 5.28. Determine the interpolatory quadrature rule with $M \in \{2, 3\}$ equidistant nodes for

$$I(h) = \int_0^1 h(x)f(x) dx,$$

where $f(x) = 6x(1-x)$. Determine nodes and weights for the corresponding Gauß-rules with $M \in \{1, 2\}$ nodes.

Problem 5.29. A *fixed strike lookback call* option has the payoff function

$$\Lambda((S(t_j))_{j=0}^M) = \max\left(\max_{j=0, \dots, M} S(t_j) - E, 0\right), \quad t_j = j\Delta t, \Delta = \frac{T}{M}, j = 0, \dots, M.$$

The stochastic process $S(t_j)$ is determined in the usual way through

$$S(t_{j+1}) = S(t_j) \exp\left(\left(r - \frac{1}{2}\sigma^2\right)\Delta t - \sigma\sqrt{\Delta t}Z_j\right), \quad Z_j \sim N(0, 1), \quad j = 0, \dots, M-1.$$

Write a program which calculates the 95%-confidence intervals for the standard and the antithetic Monte-Carlo approximations of the discounted expected payoff of the option. Present your results in a table together with the ratio of widths.

Take parameter values $T = 10$, $M = 500$, $r = 0.05$, $\sigma = 0.3$, $S(0) = S_0 = 5$, $E = 4$ and generate $N = 10^2, \dots, 10^5$ samples.

Chapter 6

Theory of Continuous Time Stochastic Processes and Itô-Integrals

6.1 Continuous Time Stochastic Processes

In this section we introduce basic notions of the theory of continuous time stochastic processes. In particular, we deal with filtrations, adaptiveness of stochastic processes and stopping times. For the material presented in this section our main references are [22, Sec. 1.3] and [34, Ch. 2].

Throughout this section let $(\Omega, \mathcal{F}, \mathbf{P})$ be the underlying probability space.

Definition 6.1. A family $(X(t))_{t \in \mathbb{T}}$, $\mathbb{T} \subset [0, \infty)$, of random variables with values in \mathbb{R}^d is called a *stochastic process* with *index set* \mathbb{T} and *state space* \mathbb{R}^d . We say that a stochastic process is *integrable* (*square-integrable*) if $X(t) \in L^1(\Omega; \mathbb{R}^d)$ ($X(t) \in L^2(\Omega; \mathbb{R}^d)$) for all $t \in \mathbb{T}$.

For $\omega \in \Omega$ the mapping $\mathbb{T} \ni t \mapsto X(t, \omega) \in \mathbb{R}^d$ is called a *sample path* or *trajectory* of the stochastic process.

Although most definitions and results also make sense for discrete index sets \mathbb{T} , for the rest of this chapter we take \mathbb{T} to be equal to a bounded or unbounded subinterval of $[0, \infty)$. In this case we say that we consider a stochastic process in *continuous time* – in contrast to a *discrete time* stochastic process.

We say that a stochastic process $(X(t))_{t \in \mathbb{T}}$ is *continuous* if the mappings $\mathbb{T} \ni t \mapsto X(t, \omega) \in \mathbb{R}^d$ are continuous for almost every $\omega \in \Omega$. In a similar way, we define a stochastic process $(X(t))_{t \in \mathbb{T}}$ to be *left continuous* or *right continuous*. Further, we say that a stochastic process is *cadlag* (French: “continue à droite, limite à gauche”) if it is right continuous and the left limits $\lim_{s \nearrow t} X(s, \omega)$ exists for almost every $\omega \in \Omega$.

Next, we introduce a very important concept which is used to control the evolution of information in a stochastic system.

Definition 6.2. A family of sub- σ -algebras $(\mathcal{F}_t)_{t \in \mathbb{T}} \subset \mathcal{F}$ is called a *filtration* if $\mathcal{F}_s \subset \mathcal{F}_t$ for all $s \leq t$.

The σ -algebra \mathcal{F}_t should be regarded as the set of all events for which we can decide at time t whether they have occurred or not. An important example of a filtration is the so called *natural filtration* of a stochastic process $(X(t))_{t \in \mathbb{T}}$ which is defined by

$$\mathcal{F}_t := \sigma(X(s) : s \in \mathbb{T}, s \leq t), \quad t \in T,$$

that is, \mathcal{F}_t is the smallest σ -algebra for which the random variables $X(s)$, $s \leq t$, are measurable.

Definition 6.3. A stochastic process $(X(t))_{t \in \mathbb{T}}$ is called *adapted* to a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ if the random variables $X(t)$ are \mathcal{F}_t -measurable for every $t \in \mathbb{T}$.

Clearly, the smallest filtration to which a stochastic process X is adapted is its natural filtration. A useful enlargement of the natural filtration is given by the next definition.

Definition 6.4. For a given filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ consider

$$\mathcal{F}_t^+ := \bigcap_{\varepsilon > 0} \mathcal{F}_{t+\varepsilon}$$

for every $t \in \mathbb{T}$. We say that the filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ is *right-continuous* if $\mathcal{F}_t = \mathcal{F}_t^+$ for all $t \in \mathbb{T}$.

From the definition it follows at once that $\mathcal{F}_t \subset \mathcal{F}_t^+$ and that $(\mathcal{F}_t^+)_{t \in \mathbb{T}}$ is a right-continuous filtration.

We are also interested in stopping a stochastic process if given random events occur. For example, we may sell a certain stock if the stock price hits a prescribed barrier. For this it is important that the decision to stop at time t only depends on those information which are available at time t . This is formalised in the following definition.

Definition 6.5. Consider a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ and a random variable $\tau: \Omega \rightarrow \mathbb{T} \cup \{\infty\}$. We say that τ is an $(\mathcal{F}_t)_{t \in \mathbb{T}}$ -*stopping time* if $\{\tau \leq t\} \in \mathcal{F}_t$ for all $t \in \mathbb{T}$.

Obviously, a constant random variable is a stopping time. For right-continuous filtrations the following characterization holds true (c.f. [34, Lemma 2.2.2]).

Lemma 6.6. *Let $(\mathcal{F}_t)_{t \in \mathbb{T}}$ be a right-continuous filtration. A random variable $\tau: \Omega \rightarrow \mathbb{T} \cup \{\infty\}$ is a $(\mathcal{F}_t)_{t \in \mathbb{T}}$ -stopping time if and only if $\{\tau < t\} \in \mathcal{F}_t$ for all $t \in \mathbb{T}$.*

Proof. For the proof we follow [34, Lemma 2.2.2]. First we assume that $\{\tau < t\} \in \mathcal{F}_t$ for all $t \in \mathbb{T}$. Then it follows that

$$\{\tau \leq t\} = \bigcap_{n=1}^{\infty} \left\{ \tau < t + \frac{1}{n} \right\} \in \bigcap_{n=1}^{\infty} \mathcal{F}_{t+\frac{1}{n}} = \mathcal{F}_t^+ = \mathcal{F}_t \quad \text{for all } t \in \mathbb{T}.$$

Conversely, let τ be a stopping time. Thus

$$\{\tau < t\} = \bigcup_{n=1}^{\infty} \left\{ \tau \leq t - \frac{1}{n} \right\} \in \mathcal{F}_t,$$

since $\mathcal{F}_{t-\frac{1}{n}} \subset \mathcal{F}_t$ for all $n \geq 1$ with $t - \frac{1}{n} \in \mathbb{T}$. In addition, if $t_0 := \inf(\mathbb{T}) \in \mathbb{T}$, then $\{\tau < t_0\} = \emptyset \in \mathcal{F}_{t_0}$. \square

Next, we come to an important example of a stopping time.

Definition 6.7. Let $(X(t))_{t \in \mathbb{T}}$ be an \mathbb{R}^d -valued stochastic process and let D be a subset of \mathbb{R}^d . Then the random variable $\tau^D: \Omega \rightarrow \mathbb{T} \cup \{\infty\}$ which is given by

$$\tau^D(\omega) := \inf \{t \in \mathbb{T} : X(t, \omega) \notin D\},$$

with the usual convention $\inf \emptyset = \infty$, is called the (*first*) *exit time* of X from D .

We have the following result (c.f. [22, Th. 1.3.2] and [34, Prop. 2.3.4]).

Lemma 6.8. *Let $(\mathcal{F}_t)_{t \in \mathbb{T}}$ be a right-continuous filtration and $(X(t))_{t \in \mathbb{T}}$ an \mathbb{R}^d -valued, $(\mathcal{F}_t)_{t \in \mathbb{T}}$ -adapted stochastic process. If $(X(t))_{t \in \mathbb{T}}$ is left-continuous or right-continuous then the exit time τ^D of X from a closed subset $D \subset \mathbb{R}^d$ is an $(\mathcal{F}_t)_{t \in \mathbb{T}}$ -stopping time.*

Proof. Following the proof of [34, Prop. 2.3.4], we have

$$\{\tau^D < t\} = \bigcup_{s < t} \{X(s) \in D^c\} = \bigcup_{s < t, s \in \mathbb{Q}} \{X(s) \in D^c\},$$

which belongs to \mathcal{F}_t by the adaptiveness of X . The last equality holds since D^c is open and X is left-continuous or right-continuous. Therefore, if $X(s) \notin D$ for some $s < t$ then there also exists $\tilde{s} \in \mathbb{Q}$, $\tilde{s} < t$, with $X(\tilde{s}) \in D^c$.

By applying Lemma 6.6 the proof is complete. \square

We close this section with what is known as the “*usual conditions*” in the literature (c.f. [34, Sec. 2.1])

Assumption 6.9. Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space together with a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ which satisfies the following *usual conditions*:

- (i) $(\Omega, \mathcal{F}, \mathbf{P})$ is a complete probability space, that is, every subset $A \subset N$ of a \mathbf{P} -nullset N is \mathcal{F} -measurable.
- (ii) Every \mathbf{P} -nullset $N \in \mathcal{F}$ also belongs to \mathcal{F}_t for all $t \in \mathbb{T}$.
- (iii) The filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ is right-continuous.

Throughout the rest of these lecture notes we usually assume that Assumption 6.9 is satisfied if not stated otherwise.

6.2 Martingales

In this section we consider an important class of stochastic processes, the so called martingales. Historically, martingales have been introduced by Paul Pierre Lévy, while many important results of this theory were contributed by Joseph Leo Doob.

The name is related to a betting system with the simplest variant working in the following way: Imagine a game of flipping a fair coin. The gambler's stake is doubled if the coin comes up heads, otherwise the stake is completely lost. The betting strategy asks the gambler to double his stake after every loss. Then the reward of the first win after a streak of several losses equals the sum of all previous losses plus a gain of the initial stake.

However, due to the exponential growth of the stakes during a long sequence of consecutive losses the gambler might hit its credit limit and the strategy results in the gambler's bankruptcy. This potentially catastrophic loss, which occurs with a small but strictly positive probability, results in the fact that the overall expected return is equal to zero.

In probability theory a martingale can be seen as a fair game. A consequence of Doob's optional stopping theorem is that there exist no betting system which provides a positive expected return in a fair game. For a more detailed discussion in this direction we refer to [4, §17] and [5, Sec. 35].

In the following we review some aspects of the martingale theory which will prove useful later on. The presented material is standard and is found in, for example, [34, Ch. 3]. For the following definition we recall the notion of conditional expectation from Section 2.6.

Definition 6.10. Let $(\mathcal{F}_t)_{t \in \mathbb{T}}$ be a filtration. An adapted, real-valued, integrable stochastic process $(X(t))_{t \in \mathbb{T}}$ is called a *martingale* with respect to $(\mathcal{F}_t)_{t \in \mathbb{T}}$ if

$$\mathbf{E}[X(t)|\mathcal{F}_s] = X(s) \quad \text{for all } s \leq t. \quad (6.1)$$

If (6.1) only holds with ' \geq ' (' \leq ') then $(X(t))_{t \in \mathbb{T}}$ is called a *submartingale* (*supermartingale*).

In the following section we will introduce Brownian motion which is an important stochastic process. The first part of Example 6.11 gives the arguments why Brownian motion is a martingale. The second part is a more artificial example.

Example 6.11. (i) Let $(X(t))_{t \in \mathbb{T}}$ be an adapted, real-valued, integrable stochastic process such that the increments $X(t) - X(s)$ are independent of \mathcal{F}_s with $\mathbf{E}[X(s)] = \mathbf{E}[X(t)]$ for all $s < t$.

Here, we say that a random variable Y is independent of a σ -algebra $\mathcal{G} \subset \mathcal{F}$ if every \mathcal{G} -measurable random variable is independent of Y . In particular, for $0 \leq t_1 < t_2 \leq t_3 < t_4 < \infty$ it follows that the increments $X(t_2) - X(t_1)$ and $X(t_4) - X(t_3)$ are independent.

In order to show the martingale property the following fact on conditional expectation is needed

$$\mathbf{E}[Y|\mathcal{G}] \equiv \mathbf{E}[Y], \quad (6.2)$$

which holds for all integrable random variables Y which are independent of a sub- σ -algebra \mathcal{G} . Using this, we compute

$$\mathbf{E}[X(t)|\mathcal{F}_s] = \mathbf{E}[X(t) - X(s)|\mathcal{F}_s] + \mathbf{E}[X(s)|\mathcal{F}_s] = \mathbf{E}[X(t) - X(s)] + X(s) = X(s).$$

It remains to show (6.2), that is

$$\int_A Y \, d\mathbf{P} = \int_A \mathbf{E}[Y] \, d\mathbf{P}|_{\mathcal{G}} \quad \text{for all } A \in \mathcal{G}.$$

Indeed, we have

$$\int_A Y \, d\mathbf{P} = \mathbf{E}[\mathbb{1}_A Y] = \mathbf{E}[\mathbb{1}_A] \mathbf{E}[Y] = \mathbf{P}(A) \mathbf{E}[Y] = \int_A \mathbf{E}[Y] \, d\mathbf{P}|_{\mathcal{G}},$$

since $\mathbb{1}_A$ is a \mathcal{G} -measurable random variable and thus independent of Y . This proves (6.2).

(ii) Consider a real-valued random variable $X \in L^1(\Omega)$ and a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$. Then the process $(X(t))_{t \in \mathbb{T}}$ defined by

$$X(t) := \mathbf{E}[X|\mathcal{F}_t] \quad \text{for all } t \in \mathbb{T}$$

is a martingale. In fact, by the properties of the conditional expectation it holds that $(X(t))_{t \in \mathbb{T}}$ is an adapted integrable process and

$$\mathbf{E}[X(t)|\mathcal{F}_s] = \mathbf{E}[\mathbf{E}[X|\mathcal{F}_t]|\mathcal{F}_s] = \mathbf{E}[X|\mathcal{F}_s] = X(s) \quad \text{for all } s \leq t.$$

□

As our first lemma shows we often encounter submartingales as the concatenation of a martingale and a convex function. An important application of this result is $\varphi(x) = |x|$.

Lemma 6.12. *Consider a convex function $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ and a real-valued martingale $(X(t))_{t \in \mathbb{T}}$ such that $\varphi(X(t))$ is integrable for all $t \in \mathbb{T}$. Then the process $(\varphi(X(t)))_{t \in \mathbb{T}}$ is a submartingale.*

Proof. As in the proof of [34, Lemma 3.0.16] we apply Jensen's inequality for conditional expectations which yields

$$\varphi(X(s)) = \varphi(\mathbf{E}[X(t)|\mathcal{F}_s]) \leq \mathbf{E}[\varphi(X(t))|\mathcal{F}_s]$$

for all $s \leq t$. □

Next, we are concerned with the regularity of typical sample paths of a martingale. For this, we say that a stochastic process $(Y(t))_{t \in \mathbb{T}}$ is a *modification* of a stochastic process $(X(t))_{t \in \mathbb{T}}$ if $\mathbf{P}(X(t) = Y(t)) = 1$ for all $t \in \mathbb{T}$.

Theorem 6.13. *Under the usual conditions on $(\mathcal{F}_t)_{t \in \mathbb{T}}$ every martingale has a cadlag modification.*

A proof is found in [34, Th. 3.2.6 (ii)].

A consequence of Doob's martingale stopping theorem ([34, Th. 3.2.7]) is that the martingale property remains valid after stopping the process by a stopping time.

Theorem 6.14. *Consider a right-continuous (sub-)martingale $(X_t)_{t \in \mathbb{T}}$ and a stopping time τ . Then the stopped process $(X(t \wedge \tau))_{t \in \mathbb{T}}$ is also a (sub-)martingale.*

For a proof we refer to [34, Cor. 3.2.8].

The next result is the continuous version of Doob's submartingale inequality. The presented formulation is taken from [22, Th. 1.3.7]. A proof can be found in [34, Th. 3.2.10].

Lemma 6.15. *Let $p > 1$ and consider a real-valued nonnegative and right-continuous submartingale $(X(t))_{t \in \mathbb{T}}$ such that $X(t) \in L^p(\Omega)$. Then it holds that*

$$\mathbf{E}\left[\sup_{s \leq t} X(s)^p\right] \leq \left(\frac{p}{p-1}\right)^p \mathbf{E}[X(t)^p] \quad \text{for all } t \in \mathbb{T}.$$

Finally, we say that an \mathbb{R}^d -valued stochastic process is a martingale if all components are real-valued martingales. Note that all presented results also hold in the multidimensional case. In particular, we have the following generalization of Lemma 6.12:

Lemma 6.16. *Let $(X(t))_{t \in \mathbb{T}}$ be an \mathbb{R}^d -valued martingale with respect to a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$. If $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\varphi(X(t))$ is integrable for all $t \in \mathbb{T}$, then $(\varphi(X(t)))_{t \in \mathbb{T}}$ is a real-valued submartingale with respect to $(\mathcal{F}_t)_{t \in \mathbb{T}}$.*

The proof is similar to the proof of Lemma 6.12 but uses a suitable generalized version of Jensen's inequality for convex functions in \mathbb{R}^d .

6.3 Brownian Motion

In 1826–27 Robert Brown discovered what today is known as the *Brownian motion*: Under the microscope he noticed that pollen particles suspended in water move erratically. In particular, he saw that the path of a single particle seems to be nowhere differentiable and that the paths of two particles appear to be independent of each other. As it was explained later, the Brownian motion is caused by a large number of collisions between the pollen particles and the molecules of water.

In his PhD-thesis from 1900 Louis Bachelier introduced a first mathematical model to describe fluctuations in stock prices using stochastic analysis. In 1905 Albert Einstein studied the Brownian motion and used his results to give an indirect proof of the existence of atoms and molecules.

These findings resulted in a rigorous mathematical theory for the Brownian motion which was first developed by Norbert Wiener during the 1920's. To his honour the stochastic process which describes the physical phenomena known as Brownian motion is also often called a *Wiener process*.

Definition 6.17. A (standard) *Wiener process* or a (standard) *Brownian motion* is a real-valued continuous stochastic process $(W(t))_{t \in \mathbb{T}}$ on $\mathbb{T} = [0, \infty)$, adapted to a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$, which satisfies the following properties:

- (i) $W(0) = 0$ a.s.;
- (ii) the increments $W(t) - W(s)$ for $0 \leq s < t < \infty$ are normally distributed with mean zero and variance $t - s$, that is $W(t) - W(s) \sim N(0, t - s)$;
- (iii) the increments $W(t) - W(s)$ are independent of \mathcal{F}_s for all $0 \leq s < t < \infty$.

For a discussion of property (iii) we refer to Example 6.11 and Figure 6.1 shows a typical sample path of a Wiener process. But before we continue to study some important properties we present a motivation of the definition from [8, Ch. 3] which follows the way Albert Einstein studied Brownian motion.

If we inject a unit amount of ink at position $x = 0$ and at time $t = 0$ into a long thin tube filled with water then let $u(x, t)$ denote the density of ink at time $t \geq 0$ and at position $x \in \mathbb{R}$. Thus, at time $t = 0$ the density is given by

$$u(x, 0) = \delta_0(x), \quad \text{the } \delta\text{-distribution centered at 0.}$$

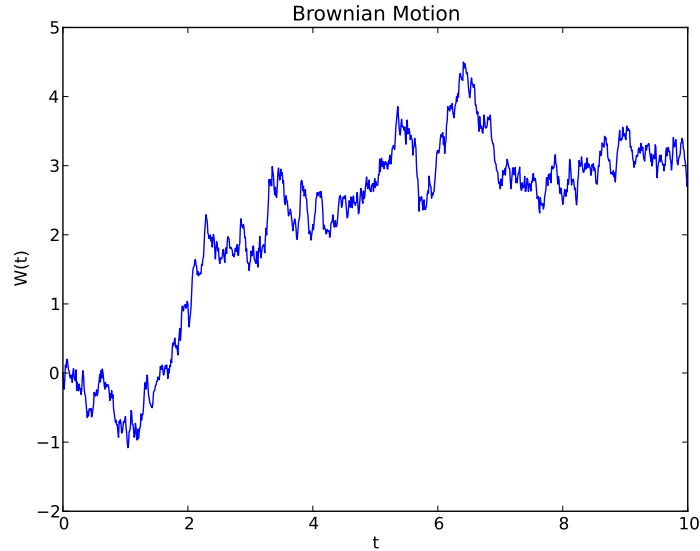


Fig. 6.1 Sample path of a Wiener process on the time interval $[0, 10]$.

Further, let us denote by $\rho(y, \tau)$ the probability density that an ink particle changes its position from x to $x + y$ in (small) time $\tau > 0$. In particular, it holds that

$$\int_{-\infty}^{\infty} \rho(y, \tau) dy = 1.$$

By an application of Taylor's theorem we get

$$\begin{aligned} u(x, t + \tau) &= \int_{-\infty}^{\infty} u(x - y, t) \rho(y, \tau) dy \\ &= \int_{-\infty}^{\infty} \left(u(x, t) - y \frac{\partial u}{\partial x}(x, t) + \frac{1}{2} y^2 \frac{\partial^2 u}{\partial x^2}(x, t) + \dots \right) \rho(y, \tau) dy. \end{aligned}$$

Further, we impose the following additional assumptions on ρ

$$\rho(y, \tau) = \rho(-y, \tau) \quad \text{and} \quad \int_{-\infty}^{\infty} y^2 \rho(y, \tau) dy = D\tau, \quad \text{for } D > 0,$$

that is, we assume that it is equally likely that a particle moves to the left or to the right and that the variance of the movements is linear in τ . The former yields

$$\int_{-\infty}^{\infty} y \rho(y, \tau) dy = 0$$

and, consequently,

$$\frac{1}{\tau} (u(x, t + \tau) - u(x, t)) = \frac{1}{2} \frac{\partial^2 u}{\partial x^2}(x, t) + \dots$$

Finally, the limit $\tau \rightarrow 0$ leads to the partial differential equation

$$\begin{aligned} \frac{\partial u}{\partial t} &= \frac{D}{2} \frac{\partial^2 u}{\partial x^2}, \\ u(x, 0) &= \delta_0(x), \end{aligned} \tag{6.3}$$

which is usually referred to as the *heat equation* or *diffusion equation*. The (fundamental) solution to (6.3) is given by

$$u(x, t) = \frac{1}{\sqrt{2\pi Dt}} \exp\left(-\frac{x^2}{2Dt}\right), \tag{6.4}$$

which is the p.d.f. of $N(0, Dt)$. Since we started with a unit amount of ink we may interpret u as the probability density of the distribution of the ink particles.

Concerning the constant D , Albert Einstein computed

$$D = \frac{RT}{N_A f},$$

where R is the gas constant, T the absolute temperature, f the friction coefficient and N_A denotes Avogadro's number ($\equiv 6 \times 10^{23}$ = number of molecules in a mole).

To sum up, in an idealized situation with $D = 1$ the new position at time $t > 0$ of an ink particle which started at $x = 0$ is $N(0, t)$ -distributed.

A further motivation of this result is given by the Central Limit Theorem. If we follow the assumption that changes to the position of an ink particle are due to a large number of independent collisions with water molecules this will result in normally distributed position changes. A more rigorous derivation of this result, which also uses the Central Limit Theorem can also be found in [8, Sec. 3].

Before we discuss some useful properties of Wiener processes we first note that from Definition 6.17 it is not immediately clear if there even exists such a stochastic process. The first existence result was presented by Norbert Wiener in 1923. The following formulation of an existence theorem is found in [8, Ch. 3].

Theorem 6.18. *Let $(\Omega, \mathcal{F}, \mathbf{P})$ be a probability space on which a sequence $(Z_i)_{i=1}^{\infty}$ of countably many independent and $N(0, 1)$ -distributed random variables exists. Then there exists a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ with $\mathbb{T} = [0, \infty)$ and an adapted stochastic process $(W(t))_{t \in \mathbb{T}}$ such that W is a standard Wiener process.*

For an easy-to-access proof of this theorem which relies on the Lévy-Ciesielski construction of a Wiener process we also refer to [8, Ch. 3] (which in turn gives [18] as a reference).

As it is discussed in [22, Sec. 15] the filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ is part of the definition of a Wiener process. The filtration which is given by Theorem 6.18 does usually not satisfy the usual conditions. But since this is often a convenient property we remark that if $(W(t))_{t \in \mathbb{T}}$ is a Wiener process on a probability space $(\Omega, \mathcal{F}, \mathbf{P})$ then it is also a Wiener process on $(\Omega, \overline{\mathcal{F}}, \mathbf{P})$, where $\overline{\mathcal{F}}$ denotes the completion of \mathcal{F} (see (2.1)).

Further, if \mathcal{N} denotes the set of all \mathbf{P} -null sets, that is $\mathcal{N} = \{A \in \overline{\mathcal{F}} : \mathbf{P}(A) = 0\}$, then we define

$$\overline{\mathcal{F}}_t := \sigma(\mathcal{F}_t \cup \mathcal{N}).$$

Clearly, $(\overline{\mathcal{F}}_t)_{t \in \mathbb{T}}$ is again a filtration to which $(W(t))_{t \in \mathbb{T}}$ is adapted. This filtration is called the *augmentation under \mathbf{P} of $(\mathcal{F}_t)_{t \in \mathbb{T}}$* and, by a possibly further enlargement to $(\overline{\mathcal{F}}_t^+)_{t \in \mathbb{T}}$ we may also assume that it is right-continuous. Therefore, without loss of generality we may assume that $(W(t))_{t \in \mathbb{T}}$ is a Wiener process with respect to a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ which satisfies the usual conditions.

Our first lemmas are direct consequences of Definition 6.17. By property (iii) in the definition of Wiener processes and the arguments given in Example 6.11 Brownian motion is a martingale.

Lemma 6.19. *Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be a Wiener process with filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$. Then $(W(t))_{t \in \mathbb{T}}$ is a square-integrable $(\mathcal{F}_t)_{t \in \mathbb{T}}$ -martingale.*

For weak convergence of numerical schemes it will be important that we can calculate the mean of observables that depend on the Brownian motion.

Lemma 6.20. *Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$ be a Wiener process. Consider a measurable function $h: \mathbb{R}^n \rightarrow \mathbb{R}$ with $n \geq 0$ and a choice of time points $0 = t_0 < t_1 < \dots < t_n$. Then we have*

$$\mathbf{E}[h(W(t_1), \dots, W(t_n))] = \int_{\mathbb{R}} \dots \int_{\mathbb{R}} h(x_1, \dots, x_n) f(x_1; 0, t_1) \times \\ f(x_2; x_1, t_2 - t_1) \dots f(x_n; x_{n-1}, t_n - t_{n-1}) dx_1 \dots dx_n,$$

where

$$f(x; y, t) = \frac{1}{\sqrt{2\pi t}} e^{-\frac{(x-y)^2}{2t}}$$

denotes the density of a Gaussian random variable with mean $y \in \mathbb{R}$ and variance $t > 0$.

Proof. Here we follow [8, Ch. 3, pp. 37–38]. First we define

$$\tilde{h}(y_1, \dots, y_n) = h(y_1, y_1 + y_2, \dots, y_1 + \dots + y_n),$$

then it holds

$$h(W(t_1), \dots, W(t_n)) = \tilde{h}(W(t_1), W(t_2) - W(t_1), \dots, W(t_n) - W(t_{n-1})).$$

Since the increments $W(t_i) - W(t_{i-1})$, $i = 1, \dots, n$ are independent it follows

$$\begin{aligned} & \mathbf{E}[\tilde{h}(W(t_1), W(t_2) - W(t_1), \dots, W(t_n) - W(t_{n-1})))] \\ &= \int_{\mathbb{R}^n} \tilde{h}(y_1, \dots, y_n) f(y_1; 0, t_1) f(y_2; 0, t_2 - t_1) \cdots f(y_n; 0, t_n - t_{n-1}) dy_1 \cdots dy_n \\ &= \int_{\mathbb{R}^n} h(x_1, \dots, x_n) f(x_1; 0, t_1) f(x_2; x_1, t_2 - t_1) \cdots f(x_n; x_{n-1}, t_n - t_{n-1}) dx_1 \cdots dx_n \end{aligned}$$

where we used the transformation $y_i = x_i - x_{i-1}$, $i = 1, \dots, n$, with $x_0 = 0$ for the last equality. Note that the determinant of the Jacobian of this transformation equals 1. \square

Further, the temporal regularity of Wiener processes plays an important role since this determines the order of convergence of the numerical schemes which we study later.

Lemma 6.21. *Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be a Wiener process. Then we have*

$$\mathbf{E}[|W(t) - W(s)|^p] = C|t - s|^{\frac{p}{2}}$$

for all real numbers $p \geq 1$, $t, s \in \mathbb{T}$ and a constant $C > 0$ only depending on p .

Proof. Assume that $t \geq s$. Since $W(t) - W(s) \sim N(0, t - s)$ we get for the p -th moment

$$\mathbf{E}[|W(t) - W(s)|^p] = \frac{2^{\frac{p}{2}} \Gamma(\frac{p+1}{2})}{\sqrt{\pi}} (t - s)^{\frac{p}{2}}.$$

Thus, the assertion follows by setting

$$C := \frac{2^{\frac{p}{2}} \Gamma(\frac{p+1}{2})}{\sqrt{\pi}}.$$

\square

The previous lemma provides information on the temporal regularity of $(W(t))_{t \in \mathbb{T}}$ with respect to the L^p -norm. Together with the following famous theorem we can use this result to analyze the temporal continuity of typical sample paths.

Theorem 6.22 (Kolmogorov-Chentsov). *Let $(X(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be an \mathbb{R}^d -valued stochastic process such that there exists $\alpha, \beta > 0$ and $C > 0$ with*

$$\mathbf{E}[|X(t) - X(s)|^\alpha] \leq C|t - s|^{1+\beta} \quad \text{for all } s, t \in \mathbb{T}.$$

Then there exists a modification $(\tilde{X}(t))_{t \in \mathbb{T}}$ of $(X(t))_{t \in \mathbb{T}}$ such that every sample path $t \mapsto \tilde{X}(t, \omega)$, $\omega \in \Omega$, is locally Hölder continuous for each exponent $0 < \gamma < \frac{\beta}{\alpha}$, that is, for every $T > 0$, $\omega \in \Omega$ there exists a constant $K(T, \gamma, \omega)$ such that

$$|\tilde{X}(t, \omega) - \tilde{X}(s, \omega)| \leq K|t - s|^\gamma \quad \text{for all } 0 \leq s, t \leq T.$$

A proof of this theorem is found in [4, Kor. 39.5] or in [8, p. 47].

In the following we often assume the existence of a Wiener process with continuous sample paths. This is justified by the next corollary.

Corollary 6.23. *Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be a Wiener process. Then there exists a modification of $(\tilde{W}(t))_{t \in \mathbb{T}}$ of $(W(t))_{t \in \mathbb{T}}$ such that every sample path $t \mapsto \tilde{W}(t, \omega)$, $\omega \in \Omega$, is locally Hölder continuous for each exponent $0 < \gamma < \frac{1}{2}$.*

Proof. For $p \geq 1$ set $\alpha = 2p$ and $\beta = p - 1$. Then by Lemma 6.21 we get

$$\mathbf{E}[|W(t) - W(s)|^\alpha] = C|t - s|^{1+\beta}.$$

and Theorem 6.22 yields the existence of a modification $(\tilde{W}(t))_{t \in \mathbb{T}}$ of $(W(t))_{t \in \mathbb{T}}$ such that every sample path of $(\tilde{W}(t))_{t \in \mathbb{T}}$ is locally Hölder continuous for each exponent $0 < \gamma < \frac{\beta}{\alpha}$. Since $p \geq 1$ is arbitrary and

$$\frac{\beta}{\alpha} = \frac{p-1}{2p} = \frac{1}{2} - \frac{1}{2p}$$

the assertion follows for each exponent $0 < \gamma < \frac{1}{2}$. □

The regularity result from Corollary 6.23 can not be improved as our next theorem shows. In fact, a typical sample path of a Wiener process is nowhere differentiable. Following an idea of [5, p. 504] we provide an indication for this.

Let $(W(t))_{t \in \mathbb{T}}$ be a Wiener process. Then Problem 6.36 (ii) (with $c = N \in \mathbb{N}$) asks the reader to confirm that also the transformed process

$$W_N(t) := N^{-1}W(N^2t)$$

is a Wiener process. For every $\varepsilon > 0$ there exist $N \gg 1$ large enough such that

$$\mathbf{P}\left(W(n) - W(n-1) > 1 \text{ for some } n = 1, \dots, N\right) > 1 - \varepsilon.$$

Now fix a sample path $\omega \in \Omega$ and $n = 1, \dots, N$ such that $W(n, \omega) - W(n-1, \omega) > 1$ and note that an possible enlargement of N does not affect our choice of n and ω . Then it holds that

$$W(n) - W(n-1) = N\left(W_N\left(\frac{n}{N^2}\right) - W_N\left(\frac{n-1}{N^2}\right)\right) = \frac{1}{N} \frac{W_N\left(\frac{n}{N^2}\right) - W_N\left(\frac{n-1}{N^2}\right)}{\frac{n}{N^2} - \frac{n-1}{N^2}},$$

or equivalently,

$$\frac{W_N\left(\frac{n}{N^2}, \omega\right) - W_N\left(\frac{n-1}{N^2}, \omega\right)}{\frac{n}{N^2} - \frac{n-1}{N^2}} = N(W(n, \omega) - W(n-1, \omega)).$$

Therefore, the above difference quotient of the Wiener process W_N exceeds N on the short time interval $[0, \frac{1}{N}]$, where we can choose ω from a set with probability $1 - \varepsilon$. In particular, by letting $N \rightarrow \infty$ and since W and W_N are identically distributed for all N , it is equally likely to observe an arbitrarily large difference quotient in a sample path of W on every arbitrarily small time interval. This strongly contradicts with the behaviour of differentiable functions.

In order to formulate the next result we recall that a mapping $Y: [a, b] \rightarrow \mathbb{R}$ is said to have *bounded variation* on the interval $[a, b]$, if

$$V_a^b(Y) := \sup_{\pi_n} \sum_{t_i^n \in \pi_n} |Y(t_i^n) - Y(t_{i-1}^n)| < \infty, \quad (6.5)$$

where the supremum is taken over the set of all *partitions* $\pi_n = \{a = t_0^n < t_1^n < \dots < t_{M_n}^n = b\}$. The *mesh size* of π_n is denoted by $|\pi_n| = \max_{1 \leq i \leq M_n} |t_i^n - t_{i-1}^n|$.

Theorem 6.24. *Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be a Wiener process with continuous sample paths. For every $\frac{1}{2} < \gamma \leq 1$ the sample paths $t \mapsto W(t, \omega)$ are nowhere locally Hölder continuous with exponent γ .*

On the other hand, let $(\pi_n)_{n \in \mathbb{N}}$ be a sequence of partitions of $[a, b]$ with $|\pi_n| \rightarrow 0$ for $n \rightarrow \infty$ and $\pi_n \subset \pi_{n+1}$ for all $n \in \mathbb{N}$. Then the quadratic variation of a typical sample path is given by

$$\langle W \rangle_{[a, b]} := \lim_{n \rightarrow \infty} \sum_{i=1}^{M_n} (W(t_i^n) - W(t_{i-1}^n))^2 = b - a \quad \text{in } L^2(\Omega). \quad (6.6)$$

In particular, the sample paths are nowhere differentiable and have unbounded variation on any finite interval $[a, b] \subset \mathbb{R}$.

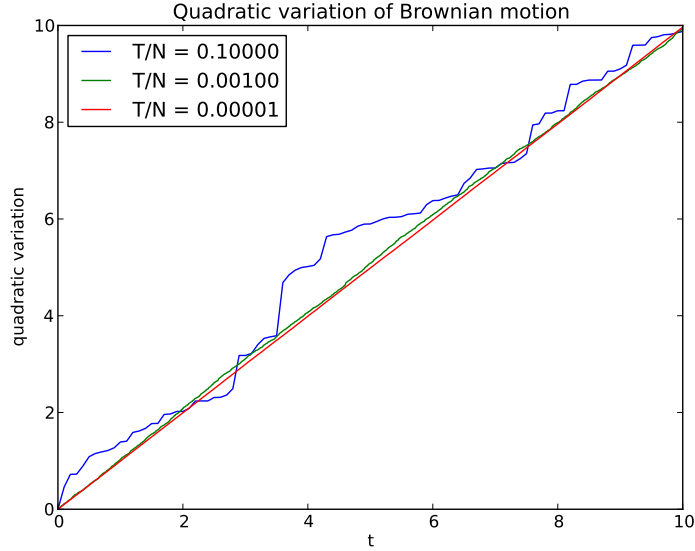


Fig. 6.2 Quadratic variation of a sample path of a Wiener process on the time interval $[0, 10]$.

Remark 6.25. The limit in (6.6) also converges \mathbf{P} -a.s. This is demonstrated in Figure 6.2 which shows for one sample path of a Wiener process on $[0, T]$ with $T = 10$ the time evolution of the quadratic variation. More precisely, for the figure we used the equidistant partitions $\pi_N = \{t_i^N = i\frac{T}{N} : i = 0, \dots, N\}$ with $N = 10^2, 10^4, 10^6$. Then the figure shows a sample path of the mappings

$$t \mapsto \sum_{\substack{t_i^N \in \pi_N \\ t_i^N < t}} (W(t_i^N) - W(t_{i-1}^N))^2 \rightarrow t \quad \text{as } N \rightarrow \infty.$$

Proof (of Theorem 6.24). We only prove the last part of the theorem. For a proof of the first part we refer to [8, pp. 50–51].

In order to prove the assertion on the quadratic variation we follow [8, p. 57] and for the partition π_n we set

$$Q_n := \sum_{i=1}^{M_n} (W(t_i^n) - W(t_{i-1}^n))^2.$$

Then we get

$$Q_n - (b - a) = \sum_{i=1}^{M_n} [(W(t_i^n) - W(t_{i-1}^n))^2 - (t_i^n - t_{i-1}^n)] =: \sum_{i=1}^{M_n} Y_i,$$

where $\mathbf{E}[Y_i] = 0$. It follows

$$\mathbf{E}[(Q_n - (b-a))^2] = \sum_{i=1}^{M_n} \mathbf{E}[Y_i^2] + 2 \sum_{\substack{i,j=1 \\ i < j}}^{M_n} \mathbf{E}[Y_i Y_j].$$

Note that Y_j and Y_i are independent for $i \neq j$ and, consequently, $\mathbf{E}[Y_i Y_j] = \mathbf{E}[Y_i] \mathbf{E}[Y_j] = 0$. Further, we have

$$\begin{aligned} \mathbf{E}[Y_i^2] &= \mathbf{E}\left[\left(\left(\frac{W(t_i^n) - W(t_{i-1}^n)}{\sqrt{t_i^n - t_{i-1}^n}}\right)^2 - 1\right)^2\right] (t_i^n - t_{i-1}^n)^2 \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} (x^2 - 1)^2 e^{-\frac{x^2}{2}} dx (t_i^n - t_{i-1}^n)^2 = 2(t_i^n - t_{i-1}^n)^2, \end{aligned}$$

since

$$\frac{W(t_i^n) - W(t_{i-1}^n)}{\sqrt{t_i^n - t_{i-1}^n}} \sim N(0, 1).$$

Therefore,

$$\mathbf{E}[(Q_n - (b-a))^2] = 2 \sum_{i=1}^{M_n} (t_i^n - t_{i-1}^n)^2 \leq 2|\pi_n|(b-a),$$

which goes to zero as $n \rightarrow \infty$. This proves the L^2 -convergence of the quadratic variation to $(b-a)$. In particular, by going to a subsequence $\mathbb{N}' \subset \mathbb{N}$ we also get for almost every $\omega \in \Omega$ that

$$\sum_{i=1}^{M_n} (W(t_i^n, \omega) - W(t_{i-1}^n, \omega))^2 \rightarrow b-a \quad \text{as } n \rightarrow \infty, n \in \mathbb{N}'.$$

Then for $0 < \gamma < \frac{1}{2}$

$$\begin{aligned} 0 < b-a &= \limsup_{n \rightarrow \infty} \sum_{i=1}^{M_n} |W(t_i^n, \omega) - W(t_{i-1}^n, \omega)|^2 \\ &\leq \limsup_{n \rightarrow \infty} C \sum_{i=1}^{M_n} |t_i^n - t_{i-1}^n|^\gamma |W(t_i^n, \omega) - W(t_{i-1}^n, \omega)| \\ &\leq \limsup_{n \rightarrow \infty} C |\pi_n|^\gamma \sum_{i=1}^{M_n} |W(t_i^n, \omega) - W(t_{i-1}^n, \omega)|. \end{aligned}$$

Thus, if a typical sample path of the Wiener process were of bounded variation then this implies that the right hand side goes to zero. This contradiction proves that sample paths of Wiener processes are of unbounded variation and, consequently, nowhere differentiable. \square

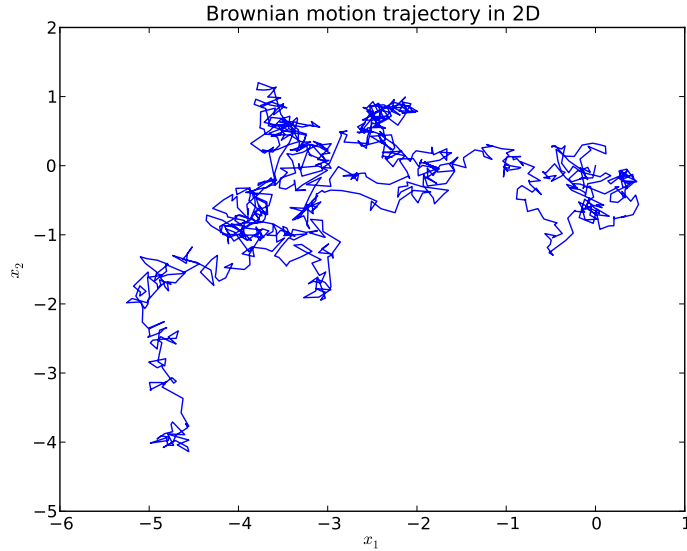


Fig. 6.3 Trajectory of the sample path of a two-dimensional Wiener process on the time interval $[0, 10]$.

So far, we considered Brownian motion only as a real-valued process. We conclude this section with a generalization to \mathbb{R}^d .

Definition 6.26. An \mathbb{R}^d -valued stochastic process

$$(W(t))_{t \in \mathbb{T}} = (W^1(t), \dots, W^d(t))_{t \in \mathbb{T}}$$

is a d -dimensional Wiener process on $\mathbb{T} = [0, \infty)$ if it holds that

- (i) each component W^i , $i = 1, \dots, d$, is a real-valued Wiener process on \mathbb{T} and
- (ii) the components are independent.

Figure 6.3 shows the trajectory of a sample path of a two-dimensional Wiener process. All results on the time regularity remain valid in the d -dimensional case. For the quadratic variation we get

$$\langle W \rangle_{[a,b]} = \lim_{n \rightarrow \infty} \sum_{i=1}^{M_n} |W(t_i^n) - W(t_{i-1}^n)|^2 = d(b-a), \quad (6.7)$$

where $|\cdot|$ denotes the Euclidean norm in \mathbb{R}^d .

Remark 6.27. Here is a remark on “white noise” missing...

6.4 The Itô-Integral

In this section we mainly follow [29, Sec. 2.3] and [22, Sec. 1.5]. Our aim is to define the stochastic integral

$$\int_0^t h(s) dW(s), \quad t \geq 0, \quad (6.8)$$

with respect to an \mathbb{R}^d -valued Wiener process $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$. We also provide a characterization of the class of admissible integrands which consists of $m \times d$ -matrix valued stochastic processes $h: \mathbb{T} \times \Omega \rightarrow \mathbb{R}^{m,d}$.

As we learned in the previous section (see Theorem 6.24) typical sample paths of a Wiener process are of unbounded variation. Therefore, the integral in (6.8) cannot be defined in the usual sense of a Lebesgue-Stieltjes-integral.

In 1949 Kiyoshi Itô was the first who gave a meaningful definition of (6.8) using certain Riemann sums. Since then it is usually called the *stochastic Itô-integral*.

In the whole section let $(\Omega, \mathcal{F}, \mathbf{P})$ denote a complete probability space. Further, let $(W(t))_{t \in \mathbb{T}}$ be an \mathbb{R}^d -valued Wiener process with continuous sample paths and adapted to a filtration $(\mathcal{F}_t)_{t \in \mathbb{T}}$ which satisfies the usual conditions. We also fix a time $T > 0$.

Following [29, Sec. 2.3] the construction of the stochastic integral is done in several steps:

Step 1: First we specify a class \mathcal{E} of elementary stochastic processes which take values in the set $\mathbb{R}^{m,d}$. Then we define the integral for this class of integrands, that is, we define a mapping

$$\begin{aligned} \text{Int}: \mathcal{E} &\rightarrow \mathcal{M}^2([0, T]) \\ h &\mapsto \int_0^t h(s) dW(s), \quad t \in [0, T]. \end{aligned}$$

Here $\mathcal{M}^2([0, T]) := \mathcal{M}^2([0, T]; \mathbb{R}^m)$ denotes the set of continuous, square-integrable, \mathbb{R}^m -valued (\mathcal{F}_t) -martingales. Below we show that this space is a Banach space.

Step 2: There exists a norm on \mathcal{E} such that $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0, T])$ is an isometry.

Therefore, by using an extension theorem we can uniquely extend the domain of Int to the completion $\bar{\mathcal{E}}$.

Step 3: It remains to provide an explicit representation of $\bar{\mathcal{E}}$.

Concerning the first step we begin with the definition of the elementary processes.

Definition 6.28. An $\mathbb{R}^{m,d}$ -valued stochastic process $(h(t))_{t \in [0, T]}$ is called *elementary* if there exists a partition $0 = t_0 < \dots < t_n = T$ of $[0, T]$, $n \in \mathbb{N}$, such that

$$h(t) = \sum_{j=0}^{n-1} h_j \mathbb{1}_{(t_j, t_{j+1}]}(t), \quad t \in [0, T], \quad (6.9)$$

where $h_j: \Omega \rightarrow \mathbb{R}^{m,d}$ is \mathcal{F}_{t_j} -measurable for $0 \leq j \leq n-1$ and only takes a finite number of values in $\mathbb{R}^{m,d}$.

By \mathcal{E} we denote the set of all elementary processes.

Note that from (6.9) it follows directly that every elementary process is left-continuous. However, the representation (6.9) is not unique in general.

In the same way as in [29, Sec. 2.3.2] we define the mapping $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0, T])$ by

$$\text{Int}(h)(t) := \int_0^t h(s) dW(s) := \sum_{j=0}^{n-1} h_j (W(t_{j+1} \wedge t) - W(t_j \wedge t)) \quad (6.10)$$

for every $t \in [0, T]$, where h is an elementary process with representation 6.9. As the next lemma shows by this setting $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0, T])$ is well-defined.

Lemma 6.29. *Let $h \in \mathcal{E}$. Then the stochastic process $(\text{Int}(h)(t))_{t \in [0, T]}$ given by (6.10) is a continuous, square-integrable $(\mathcal{F}_t)_{t \in [0, T]}$ -martingale which takes values in \mathbb{R}^m . The mapping $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0, T])$ is well-defined in the sense that its definition does not depend on the particular representation of h .*

Proof. In Problem 6.39 we ask the reader to confirm that the definition (6.10) is independent of the representation of $h \in \mathcal{E}$. For the rest of the proof we follow [29, Prop. 2.3.2].

For every representation of $h \in \mathcal{E}$ it directly follows that the mapping

$$t \mapsto \int_0^t h(s) dW(s) = \sum_{j=0}^{n-1} h_j (W(t_{j+1} \wedge t) - W(t_j \wedge t))$$

is continuous for almost every $\omega \in \Omega$ by the continuity of the Wiener process and since the multiplication of the matrix $h_j(\omega) \in \mathbb{R}^{m,d}$ with the increment vector of the Wiener process is also continuous for all $\omega \in \Omega$, $0 \leq j \leq n-1$.

Next, since each $h_j: \Omega \rightarrow \mathbb{R}^{m,d}$ only takes a finite number of values it is clear that $\sup_{\omega \in \Omega} |h_j(\omega)| < \infty$ for every $0 \leq j \leq n-1$, where in this case $|\cdot|$ denotes the matrix norm in $\mathbb{R}^{m,d}$ which is induced by the Euclidean norms in \mathbb{R}^m and \mathbb{R}^d . Hence,

$$|h_j(W(t_{j+1} \wedge t) - W(t_j \wedge t))| \leq |h_j| |W(t_{j+1} \wedge t) - W(t_j \wedge t)|$$

is square-integrable for every $0 \leq j \leq n-1$, and, consequently, $\text{Int}(h)(t)$ is square-integrable for every $t \in [0, T]$.

Now, for $t = 0$ we have $h(0) \equiv 0$ which is obviously \mathcal{F}_0 -measurable. For a fixed $t \in (0, T]$ consider the index $i \in \{0, \dots, n-1\}$ with $t_i < t \leq t_{i+1}$. Then we get

$$\text{Int}(h)(t) = \sum_{j=0}^{i-1} h_j(W(t_{j+1}) - W(t_j)) + h_i(W(t) - W(t_i)),$$

which is a \mathcal{F}_t -measurable random variable and, therefore, $(\text{Int}(h)(t))_{t \in [0, T]}$ is adapted.

It remains to prove the martingale property. For this we take $0 \leq s < t \leq T$ and a set $A \in \mathcal{F}_s$ arbitrarily. Then we have

$$\mathbf{E}[\mathbb{1}_A \text{Int}(h)(t)] = \sum_{j=0}^{n-1} \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge t) - W(t_j \wedge t))]. \quad (6.11)$$

If we consider all summands with index $0 \leq j \leq n-1$ such that $t_{j+1} < s$ then we get

$$\mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge t) - W(t_j \wedge t))] = \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge s) - W(t_j \wedge s))].$$

Further, for all summands with index $0 \leq j \leq n-1$ such that $s \leq t_j$ the increment $(W(t_{j+1} \wedge t) - W(t_j \wedge t))$ is independent of $\mathbb{1}_A h_j \in \mathcal{F}_{t_j}$ since $A \in \mathcal{F}_s \subset \mathcal{F}_{t_j}$. Therefore, we obtain

$$\mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge t) - W(t_j \wedge t))] = \mathbf{E}[\mathbb{1}_A h_j] \mathbf{E}[W(t_{j+1} \wedge t) - W(t_j \wedge t)] = 0.$$

For the remaining summand with index $j \in \{0, \dots, n-1\}$ such that $t_j < s \leq t_{j+1}$ we have

$$\begin{aligned} \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge t) - W(t_j \wedge t))] &= \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge t) - W(t_{j+1} \wedge s))] \\ &\quad + \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge s) - W(t_j \wedge s))] \end{aligned}$$

and the first term again vanishes by the independence of $W(t_{j+1} \wedge t) - W(t_{j+1} \wedge s)$ and $\mathbb{1}_A h_j \in \mathcal{F}_s = \mathcal{F}_{t_{j+1} \wedge s}$. Altogether, we get for (6.11)

$$\mathbf{E}[\mathbb{1}_A \text{Int}(h)(t)] = \sum_{j=0}^{n-1} \mathbf{E}[\mathbb{1}_A h_j(W(t_{j+1} \wedge s) - W(t_j \wedge s))] = \mathbf{E}[\mathbb{1}_A \text{Int}(h)(s)].$$

Since $A \in \mathcal{F}_s$ is arbitrary this proves $\mathbf{E}[\text{Int}(h)(t) | \mathcal{F}_s] = \text{Int}(h)(s)$. \square

Having this established we can proceed with the second step of the construction of the Itô-integral. For this we first prove that the set $\mathcal{M}^2([0, T])$ is a Banach space.

Lemma 6.30. *The set $\mathcal{M}^2([0, T])$ of all \mathbb{R}^m -valued continuous, square-integrable martingales endowed with the norm*

$$\|X\|_{\mathcal{M}^2([0, T])} := \sup_{t \in [0, T]} (\mathbf{E}[|X(t)|^2])^{\frac{1}{2}} = (\mathbf{E}[|X(T)|^2])^{\frac{1}{2}}, \quad (6.12)$$

for $X \in \mathcal{M}^2([0, T])$, is a Banach space.

Proof. The norm $\|\cdot\|_{\mathcal{M}^2([0, T])}$ is the standard norm on $C([0, T]; L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R}^m))$, which is a Banach space (see [1, Sec. 1.2]). We show that $\mathcal{M}^2([0, T])$ is a closed subspace of $C([0, T]; L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R}^m))$.

For $X \in \mathcal{M}^2([0, T])$ we have by Lemma 6.16 that $t \mapsto |X(t)|^2$ is a square-integrable submartingale. In particular, we have

$$\|X\|_{\mathcal{M}^2([0, T])} = \sup_{t \in [0, T]} (\mathbf{E}[|X(t)|^2])^{\frac{1}{2}} = (\mathbf{E}[|X(T)|^2])^{\frac{1}{2}} < \infty.$$

Further, we need to show that the mapping $t \mapsto X(t) \in L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R}^m)$ is continuous. For this fix $t \in [0, T]$ and a sequence $(t_n)_{n \in \mathbb{N}}$ with $t_n \rightarrow t$ as $n \rightarrow \infty$. Since X is a continuous martingale we get

$$\lim_{n \rightarrow \infty} |X(t, \omega) - X(t_n, \omega)| \rightarrow 0 \quad \text{for almost every } \omega \in \Omega.$$

Now we want to apply the following fact (see [5, Th. 16.14]): Consider a sequence $(Y_n)_{n \in \mathbb{N}}$ of uniformly integrable random variables such that $\lim_{n \rightarrow \infty} Y_n = Y$ almost surely for some random variable Y . Then $Y \in L^1(\Omega)$ and the sequence $(Y_n)_{n \in \mathbb{N}}$ also converges to Y with respect to the norm in $L^1(\Omega)$. Here $(Y_n)_{n \in \mathbb{N}}$ is a uniformly integrable sequence if

$$\lim_{R \rightarrow \infty} \sup_{n \in \mathbb{N}} \int_{\{|Y_n| \geq R\}} |Y_n| d\mathbf{P} = 0.$$

Thus, if we can show the uniform integrability of $Y_n := |X(t) - X(t_n)|^2$ we get the convergence of $(X(t_n))_{n \in \mathbb{N}}$ to $X(t)$ with respect to the L^2 -norm. In the following we achieve this by the same arguments as in the proof of [34, Lem. 3.2.2].

For every $\varepsilon > 0$ we can choose a $\delta > 0$ such that

$$\int_A |X(T)|^2 d\mathbf{P} < \varepsilon \quad \text{for all } A \in \mathcal{F} \text{ with } \mathbf{P}(A) < \delta.$$

This is possible since $|X(T)|^2$ is integrable. Further, the conditional version of Jensen's inequality (see Lemma 2.11) yields

$$\begin{aligned}
\int_A |X(s)|^2 d\mathbf{P} &= \int_A |\mathbf{E}[X(T)|\mathcal{F}_s]|^2 d\mathbf{P} \\
&\leq \int_A \mathbf{E}[|X(T)|^2|\mathcal{F}_s] d\mathbf{P} \\
&= \int_A |X(T)|^2 d\mathbf{P} \quad \text{for all } A \in \mathcal{F}_s, s \in [0, T].
\end{aligned} \tag{6.13}$$

Finally, by Markov's inequality it holds the estimate

$$\mathbf{P}(|X(s)|^2 > R) \leq \frac{1}{R} \mathbf{E}[|X(s)|^2] \leq \frac{1}{R} \mathbf{E}[|X(T)|^2] \quad \text{for all } s \in [0, T],$$

where the last inequality is due to the submartingale property of $t \mapsto |X(t)|^2$. Thus, for every $R > \delta^{-1} \mathbf{E}[|X(T)|^2]$ we derive

$$\mathbf{P}(|X(s)|^2 > R) \leq \frac{1}{R} \mathbf{E}[|X(T)|^2] < \delta.$$

Altogether, since $\{|X(t_n)|^2 > R\} \in \mathcal{F}_{t_n}$ we get the uniform integrability of the family $(X(t_n))_{n \in \mathbb{N}}$ since

$$\int_{\{|X(t_n)|^2 > R\}} |X(t_n)|^2 d\mathbf{P} \leq \int_{\{|X(t_n)|^2 > R\}} |X(T)|^2 d\mathbf{P} < \varepsilon,$$

where we used that $\{|X(t_n)|^2 > R\} \in \mathcal{F}_{t_n}$. Therefore, since

$$Y_n = |X(t) - X(t_n)|^2 \leq 2(|X(t)|^2 + |X(t_n)|^2)$$

we also have that the family $(Y_n)_{n \in \mathbb{N}}$ is uniformly integrable. By applying the above mentioned result from [5, Th. 16.14] this proves the continuity of the mapping $t \mapsto X(t)$ and we conclude

$$\mathcal{M}^2([0, T]) \subset C([0, T]; L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R}^m)).$$

It remains to show that $\mathcal{M}^2([0, T])$ is closed in $C([0, T]; L^2(\Omega, \mathcal{F}, \mathbf{P}; \mathbb{R}^m))$. Let $(X_n)_{n \in \mathbb{N}} \subset \mathcal{M}^2([0, T])$ denote a Cauchy-sequence with respect to $\|\cdot\|_{\mathcal{M}^2([0, T])}$. By Doob's inequality (see Lemma 6.15) we get

$$\begin{aligned}
\|X_n - X_k\|_{\mathcal{M}^2([0, T])} &= \sup_{t \in [0, T]} \|X_n(t) - X_k(t)\|_{L^2(\Omega; \mathbb{R}^m)} \\
&\leq \left(\mathbf{E} \left[\sup_{t \in [0, T]} |X_n(t) - X_k(t)|^2 \right] \right)^{\frac{1}{2}} \\
&\leq 2 \sup_{t \in [0, T]} \left(\mathbf{E} [|X_n(t) - X_k(t)|^2] \right)^{\frac{1}{2}} = 2 \|X_n - X_k\|_{\mathcal{M}^2([0, T])}.
\end{aligned}$$

That is, the norms $\|\cdot\|_{\mathcal{M}^2([0,T])}$ and $\|\cdot\|_{2,\infty}$ are equivalent on $\mathcal{M}^2([0,T])$, where

$$\|X\|_{2,\infty} := \left(\mathbf{E} \left[\sup_{t \in [0,T]} |X(t)|^2 \right] \right)^{\frac{1}{2}} \quad \text{for } X \in \mathcal{M}^2([0,T]).$$

Here, it is important to note that $\|\cdot\|_{2,\infty}$ is the norm on $L^2(\Omega, \mathcal{F}, \mathbf{P}; C([0,T], \mathbb{R}^m))$, which is a Banach space [1, Satz 1.19]. Therefore, $(X_n)_{n \in \mathbb{N}}$ is also a Cauchy-sequence in this space and there exists $X \in L^2(\Omega, \mathcal{F}, \mathbf{P}; C([0,T], \mathbb{R}^m))$ such that $\lim_{n \rightarrow \infty} \|X - X_n\|_{2,\infty} = 0$. In particular, X has \mathbf{P} -almost surely continuous sample paths. In the following we show that $X \in \mathcal{M}^2([0,T])$.

Since $X(t)$ is an \mathbb{R}^m -valued, square integrable random variable for every $t \in [0,T]$ it is clear that X is a square-integrable stochastic process. Further, as the limit of \mathcal{F}_t -measurable random variables $X(t)$ is itself \mathcal{F}_t -measurable and, therefore, X is adapted to $(\mathcal{F}_t)_{t \in [0,T]}$.

Further, X is a martingale. For this, fix $0 \leq s < t \leq T$ arbitrary. By the martingale property of the X_n we get for all $n \in \mathbb{N}$

$$\begin{aligned} \|X(s) - \mathbf{E}[X(t)|\mathcal{F}_s]\|_{L^2(\Omega; \mathbb{R}^m)} &\leq \|X(s) - X_n(s)\|_{L^2(\Omega; \mathbb{R}^m)} \\ &\quad + \|\mathbf{E}[X_n(t) - X(t)|\mathcal{F}_s]\|_{L^2(\Omega; \mathbb{R}^m)} \end{aligned}$$

Thus, in the limit $n \rightarrow \infty$ the first summand vanishes. For the second summand we apply Jensen's inequality (Lemma 2.11) and get

$$\begin{aligned} \|\mathbf{E}[X_n(t) - X(t)|\mathcal{F}_s]\|_{L^2(\Omega; \mathbb{R}^m)} &= \left(\mathbf{E} \left[\|\mathbf{E}[X_n(t) - X(t)|\mathcal{F}_s]\|^2 \right] \right)^{\frac{1}{2}} \\ &\leq \left(\mathbf{E} \left[\mathbf{E} \left[\|X_n(t) - X(t)\|^2 | \mathcal{F}_s \right] \right] \right)^{\frac{1}{2}} \\ &= \|X_n(t) - X(t)\|_{L^2(\Omega; \mathbb{R}^m)} \end{aligned}$$

which also goes to zero as $n \rightarrow \infty$. Therefore, $X(s) = \mathbf{E}[X(t)|\mathcal{F}_s]$ almost everywhere.

Altogether, this proves that $X \in \mathcal{M}^2([0,T]) \subset C([0,T]; L^2(\Omega; \mathbb{R}^m))$. Since the norms $\|\cdot\|_{2,\infty}$ and $\|\cdot\|_{\mathcal{M}^2([0,T])}$ are equivalent we also get that (X_n) converges to X with respect to the norm in $\mathcal{M}^2([0,T])$. \square

Next, we prove that $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0,T])$ is an isometry between these spaces. Equation (6.15) is usually called the *Itô-isometry*. It will play an important role for the estimates of the strong error of convergence in Chapter 8.

But before we proceed we need to introduce a norm such that \mathcal{E} becomes a normed space. For this first note that \mathcal{E} clearly is a linear space, since for $h, g \in \mathcal{E}$ one can always find representations with respect to the same partition $0 = t_0 <$

$\dots < t_n = T$ of $[0, T]$. But then it is obvious that also $\alpha h + \beta g$ is an elementary stochastic process for all $\alpha, \beta \in \mathbb{R}$.

Since elementary processes are matrix-valued, we have to use a matrix norm. It turns out that the *Hilbert-Schmidt norm* or *Frobenius norm* is suitable for our purposes. We recall that the Hilbert-Schmidt norm of a matrix $A = (A_{ij})_{i \leq m, j \leq d} \in \mathbb{R}^{m,d}$ is given by

$$|A|_{HS}^2 := \sum_{j=1}^d \sum_{i=1}^m A_{ij}^2 = \sum_{j=1}^d (A^T A)_{jj} = \text{tr}(A^T A). \quad (6.14)$$

Note, that the Hilbert-Schmidt norm is consistent with the Euclidean norm, that is

$$|Ax|^2 = \sum_{i=1}^m (Ax)_i^2 = \sum_{i=1}^m \left(\sum_{j=1}^d A_{ij} x_j \right)^2 \leq \sum_{i=1}^m \left(\sum_{j=1}^d A_{ij}^2 \right) \sum_{j=1}^d x_j^2 = |A|_{HS}^2 |x|^2.$$

Let $h \in \mathcal{E}$, then we define a norm on \mathcal{E} by

$$\|h\|_T^2 := \mathbf{E} \left[\int_0^T |h(s)|_{HS}^2 ds \right] = \mathbf{E} \left[\sum_{j=0}^{n-1} |h_j|_{HS}^2 (t_{j+1} - t_j) \right],$$

which is the norm of the Hilbert space $L^2([0, T] \times \Omega, \mathcal{B}([0, T]) \otimes \mathcal{F}, \lambda^1 \otimes \mathbf{P}; \mathbb{R}^{m,d})$, where $\mathbb{R}^{m,d}$ is equipped with the Hilbert-Schmidt norm.

As in [29, p. 27] we remark that by using the norm $\|\cdot\|_T$ we have to identify those elementary processes in \mathcal{E} which are equal $\lambda^1 \otimes \mathbf{P}$ -almost everywhere.

Proposition 6.31. *The mapping $\text{Int}: \mathcal{E} \rightarrow \mathcal{M}^2([0, T])$ is a linear isometry between \mathcal{E} equipped with the norm $\|\cdot\|_T$ and the Banach space $(\mathcal{M}^2([0, T]), \|\cdot\|_{\mathcal{M}^2([0, T])})$. In particular, the Itô-isometry*

$$\|\text{Int}(h)\|_{\mathcal{M}^2([0, T])}^2 = \mathbf{E} \left[\int_0^T |h(s)|_{HS}^2 ds \right] \quad (6.15)$$

is satisfied for all $h \in \mathcal{E}$.

Proof. By choosing representations of $g, h \in \mathcal{E}$ with respect to the same partition $0 = t_0 < \dots < t_n = T$ it is clear that $\text{Int}(\alpha h + \beta g) = \alpha \text{Int}(h) + \beta \text{Int}(g)$. In the light of Lemma 6.29 it only remains to show (6.15).

Consider $h = \sum_{j=0}^{n-1} h_j \mathbb{1}_{(t_j, t_{j+1}]} \in \mathcal{E}$ and set $\Delta W_m := W(t_{m+1}) - W(t_m)$. Then

$$\begin{aligned} \|\text{Int}(h)\|_{\mathcal{M}^2([0, T])}^2 &= \mathbf{E} \left[\left| \int_0^T h(s) dW(s) \right|^2 \right] = \mathbf{E} \left[\left| \sum_{j=0}^{n-1} h_j \Delta W_j \right|^2 \right] \\ &= \mathbf{E} \left[\sum_{i,j=0}^{n-1} (h_i \Delta W_i)^T h_j \Delta W_j \right] \\ &= \sum_{j=0}^{n-1} \mathbf{E} [(h_j \Delta W_j)^T h_j \Delta W_j] + 2 \sum_{0 \leq i < j \leq n-1} \mathbf{E} [(h_i \Delta W_i)^T h_j \Delta W_j]. \end{aligned}$$

For the second sum we use the fact that $(h_j^T h_i \Delta W_i)^T$ is \mathcal{F}_{t_j} -measurable for $i < j$ and, therefore, is independent of $\Delta W_j = W(t_{j+1}) - W(t_j)$. Thus,

$$\mathbf{E}[(h_i \Delta W_i)^T h_j \Delta W_j] = \mathbf{E}[(h_j^T h_i \Delta W_i)^T] \mathbf{E}[\Delta W_j] = 0.$$

For the first sum we have

$$\begin{aligned} \mathbf{E}[\Delta W_j^T h_j^T h_j \Delta W_j] &= \sum_{k,l=1}^d \mathbf{E}[\Delta W_{j,k} (h_j^T h_j)_{k,l} \Delta W_{j,l}] \\ &= \sum_{k=1}^d \mathbf{E}[(h_j^T h_j)_{k,k} \Delta W_{j,k}^2] + 2 \sum_{1 \leq k < l \leq d} \mathbf{E}[\Delta W_{j,k} (h_j^T h_j)_{k,l} \Delta W_{j,l}], \end{aligned}$$

where $\Delta W_{j,k}$ denotes the k -th component of ΔW_j . Now, it holds by the independence of $(h_j^T h_j)_{k,k}$ and $\Delta W_{j,k}^2$

$$\begin{aligned} \sum_{k=1}^d \mathbf{E}[(h_j^T h_j)_{k,k} \Delta W_{j,k}^2] &= \sum_{k=1}^d \mathbf{E}[(h_j^T h_j)_{k,k}] \mathbf{E}[\Delta W_{j,k}^2] \\ &= (t_{j+1} - t_j) \sum_{k=1}^d \mathbf{E}[(h_j^T h_j)_{k,k}] \\ &= (t_{j+1} - t_j) \mathbf{E}[|h_j|_{HS}^2]. \end{aligned}$$

Further, we recall that the components of the Wiener process are also independent of each other. Thus we get

$$\mathbf{E}[\Delta W_{j,k} (h_j^T h_j)_{k,l} \Delta W_{j,l}] = \mathbf{E}[\Delta W_{j,k}] \mathbf{E}[(h_j^T h_j)_{k,l}] \mathbf{E}[\Delta W_{j,l}] = 0.$$

Altogether, this proves

$$\|\text{Int}(h)\|_{\mathcal{M}^2([0,T])}^2 = \sum_{j=0}^{n-1} (t_{j+1} - t_j) \mathbf{E}[|h_j|_{HS}^2] = \mathbf{E} \left[\int_0^T |h(s)|_{HS}^2 ds \right],$$

which is Itô's isometry for elementary stochastic processes. \square

At this point we have all ingredients to apply the following extension theorem, which is a slight modification of [35, Th. II.1.5]. Here we take $\mathcal{D} = \mathcal{E}$ and the normed space N to be equal to the abstract completion of \mathcal{E} . Thus we obtain the extension of the stochastic integral

$$\overline{\text{Int}}: \overline{\mathcal{E}} \rightarrow \mathcal{M}^2([0,T]).$$

The remark after [35, Def. I.1.2] describes how to construct the completion of a normed space.

Theorem 6.32. *Let \mathcal{D} denote a dense subspace of a normed space N and let B be a Banach space. Consider a bounded linear operator $L: \mathcal{D} \rightarrow B$. Then there exists a unique extension $\bar{L}: N \rightarrow B$ of L , that is $\bar{L}|_{\mathcal{D}} = L$ and \bar{L} is bounded with operator norm $\|\bar{L}\| = \|L\|$.*

In addition, if L is an isometry, then so is \bar{L} .

Proof. Since \mathcal{D} is dense in N , for every $x \in N$ there exists a sequence $(x_n)_{n \in \mathbb{N}} \subset \mathcal{D}$ such that $x_n \rightarrow x$ as $n \rightarrow \infty$. In particular, $(x_n)_{n \in \mathbb{N}}$ is a Cauchy sequence in N . The same is true for the sequence $(L(x_n))_{n \in \mathbb{N}} \subset B$ since

$$\|L(x_n) - L(x_k)\|_B \leq \|L\| \|x_n - x_k\|_N.$$

Therefore, by the completeness of the Banach space B there exists a unique element, which we denote by $\bar{L}(x) \in B$, such that $L(x_n) \rightarrow \bar{L}(x)$ as $n \rightarrow \infty$.

The definition of $\bar{L}: N \rightarrow B$ is independent of the choice of the particular sequence $(x_n)_{n \in \mathbb{N}}$. For this, consider another sequence $(\tilde{x}_n)_{n \in \mathbb{N}} \subset \mathcal{D}$ which converges to $x \in N$. If this sequence gives rise to another element $\tilde{L}(x) \in B$ such that $L(\tilde{x}_n) \rightarrow \tilde{L}(x)$ as $n \rightarrow \infty$ then

$$\|\bar{L}(x) - \tilde{L}(x)\|_B = \lim_{n \rightarrow \infty} \|L(x_n) - L(\tilde{x}_n)\|_B \leq \lim_{n \rightarrow \infty} \|L\| \|x_n - \tilde{x}_n\|_N = 0,$$

which shows $\bar{L}(x) = \tilde{L}(x)$ for all $x \in N$.

That \bar{L} is an extension of L is easy to see by using the constant sequence $(x_n)_{n \in \mathbb{N}} \equiv x \in \mathcal{D}$. Also, the operator \bar{L} is linear. Further, we have

$$\|\bar{L}(x)\|_B = \lim_{n \rightarrow \infty} \|L(x_n)\|_B \leq \lim_{n \rightarrow \infty} \|L\| \|x_n\|_N = \|L\| \|x\|_N$$

for every $x \in N$ and $(x_n)_{n \in \mathbb{N}} \subset \mathcal{D}$ with $x_n \rightarrow x$ as $n \rightarrow \infty$. Hence, $\|\bar{L}\| = \|L\|$.

In addition, if L is an isometry, then

$$\|\bar{L}(x)\|_B = \lim_{n \rightarrow \infty} \|L(x_n)\|_B = \lim_{n \rightarrow \infty} \|x_n\|_N = \|x\|_N.$$

This completes the proof. \square

In the last step of the construction we give a characterization of the completion $\bar{\mathcal{E}}$ of \mathcal{E} with respect to the norm $\|\cdot\|_T$.

In the definition of elementary processes the following system of sets plays an important role:

$$\mathcal{G}_T := \{(s, t] \times F_s \mid 0 \leq s < t \leq T, F_s \in \mathcal{F}_s\} \cup \{\{0\} \times F_0 \mid F_0 \in \mathcal{F}_0\}.$$

\mathcal{G}_T is called the system of *predictable sets*.

Definition 6.33. Consider a stochastic process $h: [0, T] \times \Omega \rightarrow \mathbb{R}^{m,d}$. We say that h is *predictable* if it is \mathcal{P}_T - $\mathcal{B}(\mathbb{R}^{m,d})$ -measurable as a mapping from $[0, T] \times \Omega$ to $\mathbb{R}^{m,d}$, where

$$\mathcal{P}_T := \sigma(\mathcal{G}_T)$$

is the σ -algebra generated by all predictable sets \mathcal{G}_T .

Following the discussion in [34, Ch. 6.1] the σ -algebra \mathcal{P}_T is the smallest σ -algebra such that all real-valued, left-continuous and adapted stochastic processes $Y: [0, T] \times \Omega \rightarrow \mathbb{R}$ are measurable as a mapping from $[0, T] \times \Omega$ to \mathbb{R} . It is clear that $\mathcal{P}_T \subset \mathcal{B}([0, T]) \otimes \mathcal{F}$.

By definition, every elementary stochastic process is predictable. Therefore, it holds that

$$\bar{\mathcal{E}} \subset L^2([0, T] \times \Omega, \mathcal{P}_T, \lambda^1 \otimes \mathbf{P}; \mathbb{R}^{m,d}), \quad (6.16)$$

where as above the space $\mathbb{R}^{m,d}$ is equipped with the Hilbert-Schmidt norm. Our aim is to prove that we have equality in (6.16) and we introduce the notation [29, Sec. 2.3]

$$\mathcal{N}_W^2 := \mathcal{N}_W^2([0, T]; \mathbb{R}^{m,d}) := L^2([0, T] \times \Omega, \mathcal{P}_T, \lambda^1 \otimes \mathbf{P}; \mathbb{R}^{m,d}),$$

that is, \mathcal{N}_W^2 consists of all predictable stochastic processes $h: [0, T] \times \Omega \rightarrow \mathbb{R}^{m,d}$ with $\|h\|_T < \infty$. We refer to \mathcal{N}_W^2 as the set of all *stochastically integrable processes*.

Lemma 6.34. *The set \mathcal{E} is a dense subset of the Hilbert space \mathcal{N}_W^2 . In particular,*

$$\bar{\mathcal{E}} = L^2([0, T] \times \Omega, \mathcal{P}_T, \lambda^1 \otimes \mathbf{P}; \mathbb{R}^{m,d}) = \mathcal{N}_W^2. \quad (6.17)$$

For a proof we refer to [29, Prop. 2.3.8].

We summarize the result of this section in the following theorem.

Theorem 6.35. *For every $h \in \mathcal{N}_W^2$ the stochastic process*

$$t \mapsto \overline{\text{Int}}(h)(t) = \int_0^t h(s) dW(s)$$

is well-defined and a square-integrable, continuous martingale with values in \mathbb{R}^m . It satisfies the following properties:

(i) *The mapping*

$$\mathcal{N}_W^2 \ni h \mapsto \overline{\text{Int}}(h)$$

is linear.

(ii) *The Itô-isometry*

$$\|\overline{\text{Int}}(h)\|_{\mathcal{M}^2([0,T])}^2 = \mathbf{E} \left[\left| \int_0^T h(s) dW(s) \right|^2 \right] = \int_0^T \mathbf{E} [h(s)|_{\mathcal{H}_s}]^2 ds = \|h\|_T^2$$

holds for every $h \in \mathcal{N}_W^2$.

(iii) For every $t \in [0, T]$ and $h \in \mathcal{N}_W^2$ the \mathbb{R}^m -valued random variable $\overline{\text{Int}}(h)(t)$ has expectation zero, and covariance matrix

$$\text{cov}(\overline{\text{Int}}(h)(t)) = \int_0^t \mathbf{E} [h(s)h(s)^T] ds.$$

Proof. Parts (i) and (ii) are clear by the construction of the stochastic integral. It remains to prove part (iii). Let us first note that for every $h \in \mathcal{N}_W^2$ we have

$$\mathbf{E} \left[\int_0^T h(s) dW(s) \right] = 0. \quad (6.18)$$

This follows from the martingale property of processes in $\mathcal{M}^2([0, T])$ and the tower property of conditional expectations

$$\mathbf{E}[\overline{\text{Int}}(h)(t)] = \mathbf{E}[\mathbf{E}[\overline{\text{Int}}(h)(t)|\mathcal{F}_0]] = \mathbf{E}[\overline{\text{Int}}(h)(0)] = 0.$$

A generalization of (6.18) for two scalar functions $h_1, h_2 \in \mathcal{N}_W^2([0, T]; \mathbb{R})$ and two independent Wiener processes $W_1(t), W_2(t)$ is

$$\mathbf{E} \left[\int_0^t h_1(s) dW_1(s) \int_0^t h_2(s) dW_2(s) \right] = 0. \quad (6.19)$$

For this take two elementary functions

$$h_\nu = \sum_{j=1}^{M-1} h_{\nu,j} \mathbb{1}_{(t_j, t_{j+1}]}, \quad \nu = 1, 2.$$

which we assume to have the same partition without loss of generality. Then we obtain

$$\begin{aligned} & \mathbf{E} \left[\int_0^t h_1(s) dW_1(s) \int_0^t h_2(s) dW_2(s) \right] \\ &= \sum_{j,k=0}^{M-1} \mathbf{E} [h_{1,j} h_{2,k} (W_1(t_{j+1} \wedge t) - W_1(t_j \wedge t)) (W_2(t_{k+1} \wedge t) - W_2(t_k \wedge t))]. \end{aligned}$$

For $j \leq k$ the variable $h_{1,j} h_{2,k} (W_1(t_{j+1} \wedge t) - W_1(t_j \wedge t))$ is \mathcal{F}_{t_k} -measurable and hence independent of $W_2(t_{k+1} \wedge t) - W_2(t_k \wedge t)$. A similar argument applies in case

$j > k$ and, therefore, all summands above vanish. Again the general formula (6.19) follows by approximation.

Using (6.18) and (6.19) we calculate the covariance

$$\begin{aligned}
\text{cov}(\overline{\text{Int}}(h)(t))_{i,j} &= \mathbf{E} \left[\int_0^t h(s) dW(s) \left(\int_0^t h(s) dW(s) \right)^T \right]_{i,j} \\
&= \mathbf{E} \left[\left(\sum_{k=1}^d \int_0^t h_{ik}(s) dW_k(s) \right) \left(\sum_{l=1}^d \int_0^t h_{jl}(s) dW_l(s) \right) \right] \\
&= \sum_{k,l=1}^d \mathbf{E} \left[\int_0^t h_{ik}(s) dW_k(s) \int_0^t h_{jl}(s) dW_l(s) \right] \\
&= \sum_{k=1}^d \mathbf{E} \left[\int_0^t h_{ik}(s) dW_k(s) \int_0^t h_{jk}(s) dW_k(s) \right] \\
&= \sum_{k=1}^d \mathbf{E} \left[\int_0^t h_{ik}(s) h_{jk}(s) ds \right] \\
&= \int_0^t (\mathbf{E}[h(s)h(s)^T])_{i,j} ds.
\end{aligned}$$

In the penultimate step we used the Itô isometry and the simple polarization identity $4ab = (a+b)^2 - (a-b)^2$ to obtain

$$\begin{aligned}
4\mathbf{E} \left[\int_0^t h_{ik}(s) dW_k(s) \int_0^t h_{jk}(s) dW_k(s) \right] \\
&= \mathbf{E} \left[\left(\int_0^t (h_{ik}(s) + h_{jk}(s)) dW_k(s) \right)^2 + \left(\int_0^t (h_{ik}(s) - h_{jk}(s)) dW_k(s) \right)^2 \right] \\
&= \int_0^t \mathbf{E} [(h_{ik}(s) + h_{jk}(s))^2] ds + \int_0^t \mathbf{E} [(h_{ik}(s) - h_{jk}(s))^2] ds \\
&= \mathbf{E} \left[\int_0^t (h_{ik}(s) + h_{jk}(s))^2 ds + \int_0^t (h_{ik}(s) - h_{jk}(s))^2 ds \right] \\
&= 4\mathbf{E} \left[\int_0^t h_{ik}(s) h_{jk}(s) ds \right].
\end{aligned}$$

□

6.5 Itô's Formula

Exercises

Problem 6.36. Let $(W(t))_{t \in \mathbb{T}}$, $\mathbb{T} = [0, \infty)$, be a real-valued Wiener process. Show that the following processes W_1 and W_2 are also Wiener processes on \mathbb{T} :

- (i) $W_1(t) := W(t+s) - W(s)$ for a fixed $s \geq 0$,
(ii) $W_2(t) := cW(\frac{t}{c^2})$ for a constant $c > 0$.

In both cases define an appropriate filtration.

Problem 6.37. Let $(W(t))_{t \in [0, \infty)}$ be a real-valued, continuous Wiener process. Show directly that the process

$$I(t) := W(t)^2 - t$$

is a continuous martingale.

Hint: $W(t)^2 = (W(t) - W(s))^2 - W(s)^2 + 2W(t)W(s)$.

Problem 6.38. Let $(W(t))_{t \in [a, b]}$ be a real-valued Wiener process on the probability space $(\Omega, \mathcal{F}, \mathbf{P})$ under the usual conditions. For $0 \leq \lambda \leq 1$ consider the modified quadratic variation

$$Q_n(\lambda) = \sum_{j=0}^{M_n-1} (W(t_j^n + \lambda(t_{j+1}^n - t_j^n)) - W(t_j^n))^2,$$

where $\pi_n = \{a = t_0^n < \dots < t_{M_n}^n = b\}$ is a partition of $[a, b]$. Show that $Q_n(\lambda)$ converges to $\lambda(b-a)$ in $L^2(\Omega)$ as $|\pi_n| \rightarrow 0$.

Problem 6.39. Show that the definition (6.10) of the Itô-integral for elementary stochastic processes $h \in \mathcal{E}$ does not depend on the representation of h . For this, consider two representations of h , that is

$$h(t) = \sum_{j=0}^{n-1} h_j \mathbb{1}_{(t_j, t_{j+1}]}(t) = \sum_{i=0}^{k-1} \bar{h}_i \mathbb{1}_{(s_i, s_{i+1}]}(t),$$

for partitions $0 = t_0 < \dots < t_n = T$ and $0 = s_0 < \dots < s_k = T$ and random variables $h_j \in \mathcal{F}_{t_j}$, $j = 0, \dots, n-1$, and $\bar{h}_i \in \mathcal{F}_{s_i}$, $i = 1, \dots, k-1$, which only take finitely many values in $\mathbb{R}^{m,d}$. Show that

$$\sum_{j=0}^{n-1} h_j (W(t_{j+1} \wedge t) - W(t_j \wedge t)) = \sum_{i=0}^{k-1} \bar{h}_i (W(s_{i+1} \wedge t) - W(s_i \wedge t)),$$

where $(W(t))_{t \in [0, T]}$ denotes an \mathbb{R}^d -valued Wiener process.

Problem 6.40. Let $(N, \|\cdot\|)$ denote a normed space.

(i) Consider the set $N^{\mathbb{N}}$ of all sequences in N and define

$$\hat{N} := \{ \hat{x} = (x_j)_{j \in \mathbb{N}} \in N^{\mathbb{N}} \mid (x_j)_{j \in \mathbb{N}} \text{ is a Cauchy-sequence in } N \}.$$

Equip \hat{N} with the seminorm $\| \cdot \|_{\hat{N}}$, which is given by

$$\| \hat{x} \|_{\hat{N}} = \lim_{j \rightarrow \infty} \| x_j \|.$$

Show that

$$(x_j)_{j \in \mathbb{N}} \sim (y_j)_{j \in \mathbb{N}} \quad :\Leftrightarrow \quad \lim_{j \rightarrow \infty} \| x_j - y_j \| = 0$$

establishes an equivalence relation on \hat{N} . Then, after identifying elements in \hat{N} , which are equivalent with respect to \sim , prove that we obtain a Banach space $(\bar{N}, \| \cdot \|_{\bar{N}})$, where $\bar{N} = \hat{N} / \sim$.

(ii) Consider the mapping $J: N \rightarrow \bar{N}$ defined by $J(x) := (x)_{j \in \mathbb{N}}$, that is J maps $x \in N$ to the constant sequence with only value x . Show that J is one-to-one and isometric. Further, show that for every $\hat{x} = (x_j)_{j \in \mathbb{N}} \in \bar{N}$ there exists a sequence $(y_j)_{j \in \mathbb{N}}$ in N such that

$$\lim_{j \rightarrow \infty} \| \hat{x} - J(y_j) \|_{\bar{N}} = 0.$$

Problem 6.41. Let $(W(t))_{t \in [0, T]}$ be a real-valued Wiener process and $(X_1(t))_{t \in [0, T]}$, $(X_2(t))_{t \in [0, T]}$ be two real-valued Itô-processes with

$$\begin{aligned} dX_1(t) &= f_1(t) dt + g_1(t) dW(t), \\ dX_2(t) &= f_2(t) dt + g_2(t) dW(t). \end{aligned}$$

Show that

$$d(X_1(t)X_2(t)) = X_1(t) dX_2(t) + X_2(t) dX_1(t) + g_1(t)g_2(t) dt,$$

which is often called *the integration by parts formula for stochastic integrals*.

Hint: Apply Itô's formula to $V(x, y) = xy$.

Chapter 7
Stochastic Ordinary Differential Equations

Chapter 8
Numerical Solution of SODEs

Chapter 9
Weak Approximation of SODEs

Chapter 10
Monte Carlo Methods for SODEs

References

1. H. W. Alt. *Lineare Funktionalanalysis*. Springer, Berlin, 5., improv. edition, 2006.
2. H. Amann and J. Escher. *Analysis. III. Grundstudium Mathematik*. [Basic Study of Mathematics]. Birkhäuser Verlag, Basel, 2001.
3. H. Bauer. *Measure and Integration Theory*. De Gruyter studies in mathematics ; 26. de Gruyter, 2001.
4. H. Bauer. *Wahrscheinlichkeitstheorie*. de Gruyter, Berlin [u.a.], 2002.
5. P. Billingsley. *Probability and Measure*. Wiley Series in Probability and Mathematical Statistics. John Wiley & Sons Inc., New York, third edition, 1995. A Wiley-Interscience Publication.
6. F. Black and M. Scholes. The pricing of options and corporate liabilities. *Journal of Political Economy*, 81:637–654, 1973.
7. P. J. Davis and P. Rabinowitz. *Methods of Numerical Integration*. Dover Publications Inc., Mineola, NY, 2007. Corrected reprint of the second (1984) edition.
8. L. C. Evans. An introduction to stochastic differential equations. Lecture Notes, Version 1.2.
9. M. Evans and T. Swartz. *Approximating Integrals via Monte Carlo and Deterministic Methods*. Oxford Statistical Science Series. Oxford University Press, Oxford, 2000.
10. J. E. Gentle. *Random Number Generation and Monte Carlo Methods*. Statistics and Computing. Springer, New York, second edition, 2003.
11. P. Glasserman. *Monte Carlo Methods in Financial Engineering*, volume 53 of *Applications of Mathematics (New York)*. Springer, New York, 2004. Stochastic Modelling and Applied Probability.
12. M. Günther and A. Jüngel. *Finanzderivate mit MATLAB, Mathematische Modellierung und numerische Simulation*. Vieweg+Teubner Verlag / Springer Fachmedien Wiesbaden GmbH, Wiesbaden, 2010.

13. D. J. Higham. Black-Scholes for scientific computing students. *Computing in Science and Engineering (Education Section)*, 6:72–79, 2004.
14. D. J. Higham. *An Introduction to Financial Option Valuation*. Cambridge University Press, Cambridge, 2004. Mathematics, Stochastics and Computation.
15. I. Karatzas and S. E. Shreve. *Methods of mathematical finance*. Applications of mathematics ; 39. Springer, 1999.
16. D. E. Knuth. *The Art of Computer Programming. Vol. 2*. Addison-Wesley Publishing Co., Reading, Mass., second edition, 1981. Seminumerical algorithms, Addison-Wesley Series in Computer Science and Information Processing.
17. U. Krengel. *Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Vieweg Studium: Aufbaukurs Mathematik [Vieweg Studies: Mathematics Course]. Friedr. Vieweg & Sohn, Braunschweig, 2003.
18. J. Lamperti. *Probability*. Benjamin, New York, 1966.
19. P. L'Ecuyer and R. Simard. TestU01: a C library for empirical testing of random number generators. *ACM Trans. Math. Software*, 33(4):Art. 22, 40, 2007.
20. D. H. Lehmer. Mathematical methods in large-scale computing units. In *Proceedings of a Second Symposium on Large-Scale Digital Calculating Machinery, 1949*, pages 141–146, Cambridge, Mass., 1951. Harvard University Press.
21. N. Madras. *Lectures on Monte Carlo methods*, volume 16 of *Fields Institute Monographs*. American Mathematical Society, Providence, RI, 2002.
22. X. Mao. *Stochastic Differential Equations and their Applications*. Horwood Publishing Limited, Chichester, second edition, 2008.
23. G. Marsaglia. Random numbers fall mainly in the planes. *Proc. Nat. Acad. Sci. U.S.A.*, 61:25–28, 1968.
24. G. Marsaglia. A current view of random number generators. In *Computer Science and Statistics: 16th Symposium on the Interface*, pages 3–10, North-Holland, Amsterdam, 1985. (edited by L. Billard).
25. G. Marsaglia and W. W. Tsang. A fast, easily implemented method for sampling from decreasing or symmetric unimodal density functions. *SIAM J. Sci. Statist. Comput.*, 5(2):349–359, 1984.
26. G. Marsaglia and W. W. Tsang. The ziggurat method for generating random variables. *Journal of Stat. Software*, 5:1–7, 2000.
27. M. Matsumoto and T. Nishimura. Mersenne twister: A 623-dimensionally equidistributed uniform pseudo-random generator. *ACM Transactions on Modeling and Computer Simulation*, 8:3–30, 1998.

28. E. Platen and D. Heath. *A Benchmark Approach to Quantitative Finance*. Springer Finance. Springer-Verlag, Berlin, 2006.
29. C. Prévôt and M. Röckner. *A Concise Course on Stochastic Partial Differential Equations*, volume 1905 of *Lecture Notes in Mathematics*. Springer, Berlin, 2007.
30. C. P. Robert and G. Casella. *Monte Carlo Statistical Methods*. Springer Texts in Statistics. Springer-Verlag, New York, second edition, 2004.
31. M. Röckner. Wahrscheinlichkeitstheorie I, II, 2006. Vorlesungsskript.
32. R. Seydel. *Tools for computational finance*. Universitext. Springer-Verlag, Berlin, third edition, 2006.
33. A. H. Stroud. *Approximate calculation of multiple integrals*. Prentice-Hall Inc., Englewood Cliffs, N.J., 1971. Prentice-Hall Series in Automatic Computation.
34. H. v. Weizsäcker and G. Winkler. *Stochastic Integrals*. Vieweg, Braunschweig [u.a.], 1990.
35. D. Werner. *Funktionalanalysis*. Springer, Berlin, 5. Aufl., extended edition, 2005.

Index

- k -equidistributed, 45
- a.e., 28
- a.s., 21
- absolutely continuous, 24
- almost surely, 21
- antithetic variables, 67
- antithetic variate, 67
- arbitrage, 5
- asset, 1
 - return, 7
- atan2, 54
- bank account, 3
- Black-Scholes formula, 8
- Black-Scholes PDE, 8
- Borel- σ -algebra, 16
- bounded variation, 99
- Box-Muller method, 53
- Brownian motion, 93
 - existence, 95
- c.d.f., 18
- Central Limit Theorem, 28, 64, 67
 - in \mathbb{R}^d , 29
- chi-square distribution, 25, 31, 41
- CLT, 15, 28
- completion of a σ -algebra, 16
- compound rule, 83
- conditional expectation, 25, 31
 - independence of random variable and σ -algebra, 91
- conditional probability, 26
- confidence interval, 63, 64, 67
- consistent estimator, 62
- convergence
 - almost everywhere, 28
 - almost sure, 28
 - in L^p , 28
 - in p -th mean, 28
 - in distribution, 28, 29
 - in probability, 28
 - weak, 28
- covariance, 22
 - matrix, 22, 25
- cumulative distribution function, 18, 24
- Delta-hedging, 14
- DIEHARD, 44
- diffusion equation, 95
- discounted expected payoff, 10, 64
- discounting, 4
- distribution, 17
- drift, 5
- efficient market hypothesis, 7
- erf, 63
- erfinv, 58, 63
- error bars, 64
- error function, 63
- Euler-Maruyama method, 13
- European call, 1
 - exercise price, 1

- expiry date, 1
- holder, 1
- underlying, 1
- writer, 1
- event, 16
 - elementary, 15
- exit time, 89
- expectation, 19
- Fibonacci generator, 36, 42, 47
- filtration, 87
 - augmentation under \mathbf{P} , 96
 - natural \sim , 88
 - right-continuous \sim , 88
- Frobenius norm, 109
- Fundamental Theorem of Simulation, 52
- Gamma distribution, 30
- goodness-of-fit test, 40
- Hölder continuity, 98
- Hölder's inequality, 20
- heat equation, 95
- Hilbert-Schmidt norm, 109
- i.i.d., 19
- image measure, 23
- importance sampling, 73
- independent and identically distributed, 19
- independent of a σ -algebra, 91
- independent random variables, 18
- index set, 87
- indicator function, 17
- interest rate, 3
- inversion method, 49, 58
- Itô integral, 6
- Itô-isometry, 109
- Jensen's inequality, 21, 31
 - conditional version, 27
- joint distribution, 18
- Kolmogorov-Chentsov theorem, 98
- Law of large numbers
 - strong, 61, 67
- LCG, *see* linear congruential generator
- linear congruential generator, 37
- LLN, 15, 27
- logistic distribution, 58
- martingale
 - Doob's stopping theorem, 92
- mean value, 19
- measurable function, 16
- measurable set, 16
- measurable space, 16
- Mersenne Twister, 45
- middle-square method, 36
- Minkowski's inequality, 20
- modification, 92
- moment of a random variable, 20
- Monte Carlo method, 10
- NIST Test Suite, 44
- normal distribution, 24
 - \sim in \mathbb{R}^d , 25, 30
- numerical function, 17
- numerical method, 13
 - strong convergence, 13
 - weak convergence, 13
- option, *see* European call
 - fixed strike lookback call, 86
 - exotic, 71
 - lock-in call, 71
 - up-and-out call, 71
- p -integrable, 20
- p.d.f., 24
- Poisson distribution, 30
- polynomial
 - Lagrange, 78
 - orthogonal, 79
- predictable set, 111
- PRNG, *see* pseudo-random number generator
- probability density function, 24
- probability measure, 16

- probability space, 16
 - complete, 16
- proportional allocation, 76
- pseudo-random number generator, 35
 - period, 35
 - seed, 35
 - $U(0, 1)$ -PRNG, 35
- quadrature rule, 78
 - Gaussian, 80
 - nodes, 78
 - weights, 78
- quantile, 42
- quantile-quantile plot, 41, 59
- Radon-Nikodym theorem, 24, 26
- rand, 58
- randn, 55
- random variable, 16
 - distribution, 17
 - expectation, 19
 - Gaussian \sim , 24
 - independence, 18
 - integrable, 19
 - p -th moment, 20
 - variance, 22
- RANDU, 36, 37, 42, 47
- rejection method, 50
- risk neutrality, 11
- runs test, 43, 47
- σ -algebra, 16
 - generated by a function, 17
- sample mean, 61, 62, 66, 67
- sample path, 6, 17, 87
- sample standard deviation, 62, 64
- sample variance, 62, 64
- sampler
 - stratified, 75
- short selling, 5
- significance level, 63
- SODE, *see* stochastic ordinary differential equation
- square integrable, 20
- standard deviation, 22
- state space, 87
- stochastic ordinary differential equation, 5
- stochastic process, 17, 87
 - adapted, 88
 - continuous, 87
 - continuous time, 87
 - discrete time, 87
 - integrable, 87
 - predictable, 112
 - square-integrable, 87
- stopping time, 88
- stratified sampling, 74
- strong convergence, *see* numerical method 13
- Strong Law of Large Numbers, 27
- submartingale
 - Doob's inequality for \sim , 92
- TestU01, 44
- three term recursion, 79
- trajectory, 87
- transformation theorem, 23
- unbiased estimator, 62
- uncorrelated, 22
- uniform distribution, 24
- usual conditions, 89
- variance, 22
- volatility, 5
- weak convergence, *see* numerical method 13
- white noise, 102
- Wiener process, 5, 12, 93
 - existence, 95
 - Hölder continuous sample path, 98
 - in \mathbb{R}^d , 102
 - with probability 1, 21
- Ziggurat algorithm, 55, 58, 59
 - tail method, 60